# ROOT: Rethinking Offline Optimization as Distributional Translation via Probabilistic Bridge

**Manh Cuong Dao**[*1], The Hung Tran[*4], Phi Le Nguyen[2],
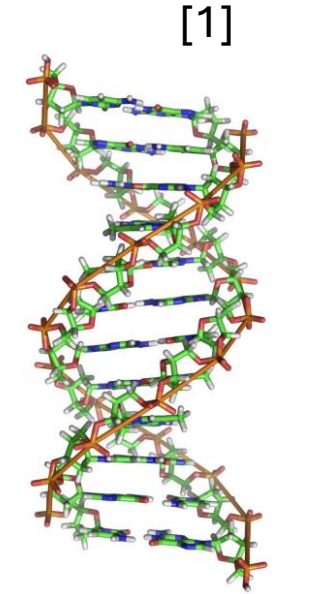Thao Nguyen Truong[3], Trong Nghia Hoang[4]
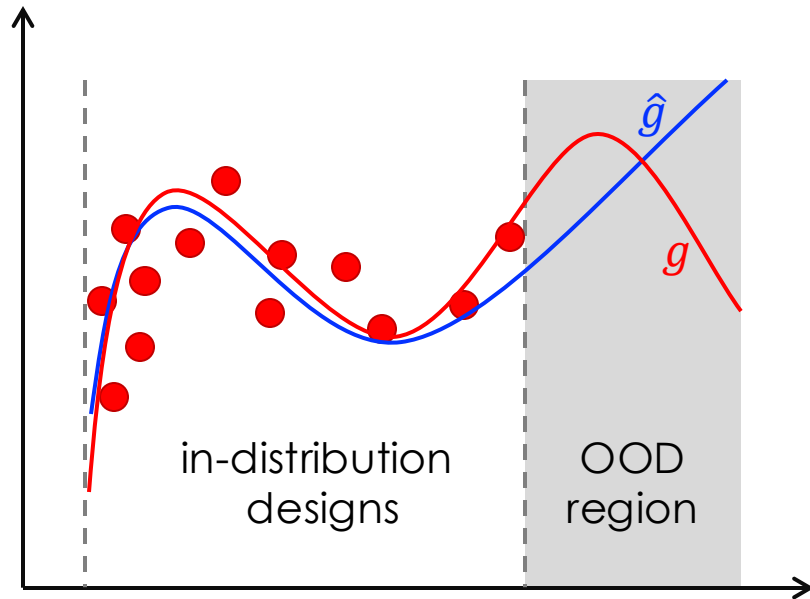
*These authors contributed equally.

# Problem Definition

- Find a design $x$ that maximizes certain desirable properties.
  - For instance:
    - Find a DNA sequence with maximum binding affinity.

- However, evaluation $g(x)$ is prohibitively expensive.
  - For instance:
    - Expensive laboratory experiment to measure binding affinity.

- **Offline Model-based Optimization (MBO):** Given an offline dataset $\mathfrak{D} = (x_i, z_i)_{i=1}^{n}$ where $z_i = g(x_i)$ with $g(.)$ is an **unknown** oracle function, find

$$x_* = \underset{x \in \mathcal{X}}{\operatorname{argmax}} \, g(x)$$

[1] Brandon Trabucco, Xinyang Geng, Aviral Kumar, and Sergey Levine. Design-bench: Benchmarks for data-driven offline model-based optimization. In International Conference on Machine Learning, pages 21658–21676. PMLR, 2022.

# Motivation



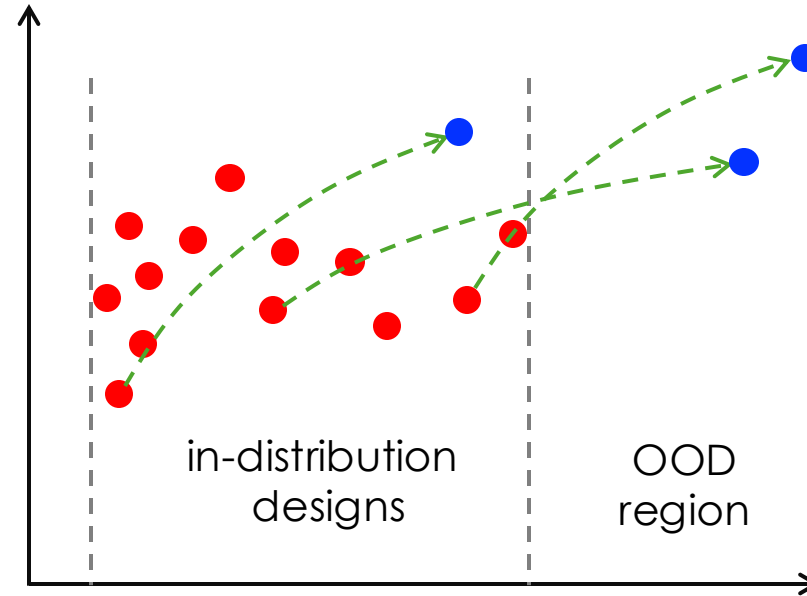**Surrogate model-based approaches**

- Learn a surrogate model $\hat{g}$
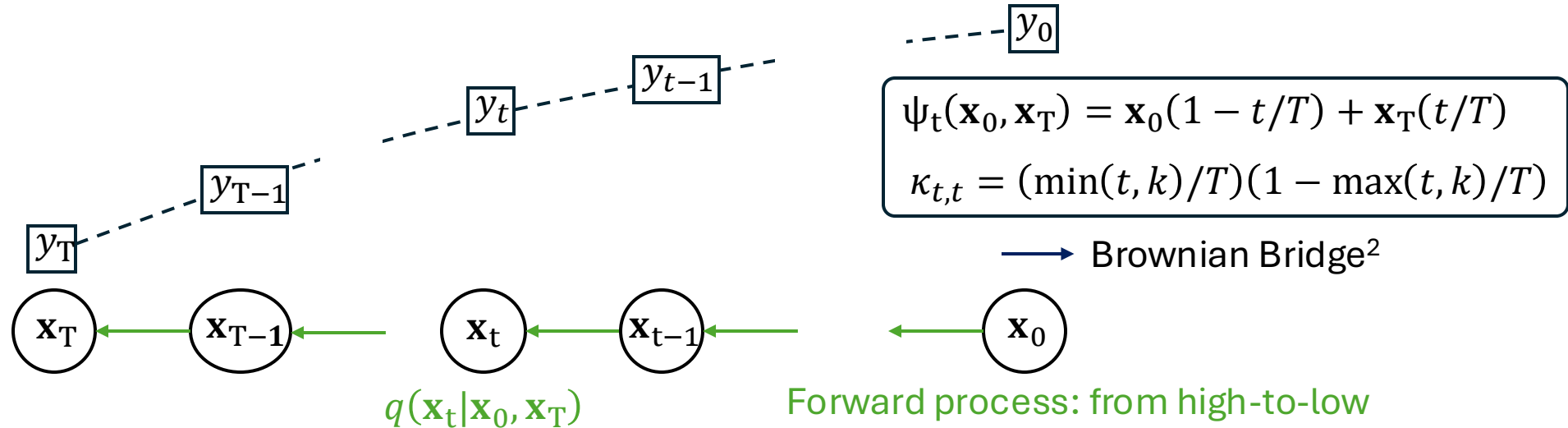- Approximate optimal design
$$x_* = \arg\max_{x} \hat{g}$$

**Challenge**: $\hat{g}$ is **unreliable** in OOD regime

**Our approach: ROOT**

- Frame **MBO** as a **distribution translation task** from low-value to high-value regime
- To enable this translation, we introduce a new concept of **probabilistic bridge.**

# Probabilistic Bridge Construction

$$\psi_t(\mathbf{x}_0, \mathbf{x}_T) = \mathbf{x}_0(1 - t/T) + \mathbf{x}_T(t/T)$$

$$\kappa_{t,t} = (\min(t,k)/T)(1 - \max(t,k)/T)$$

$\longrightarrow$ Brownian Bridge[2]

$q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_T)$

Forward process: from high-to-low

Given $\mathbf{x}_0$ and $\mathbf{x}_T$, a **probabilistic bridge** is defined as observations $[\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{T-1}]$ of a random function distributed by Gaussian Process $\boldsymbol{GP}\big(\psi_t(\mathbf{x}_0, \mathbf{x}_T), \kappa_{t,k}\boldsymbol{I}\big)$ where $\psi_0(\mathbf{x}_0, \mathbf{x}_T) = \mathbf{x}_0$, and $\psi_T(\mathbf{x}_0, \mathbf{x}_T) = \mathbf{x}_T$
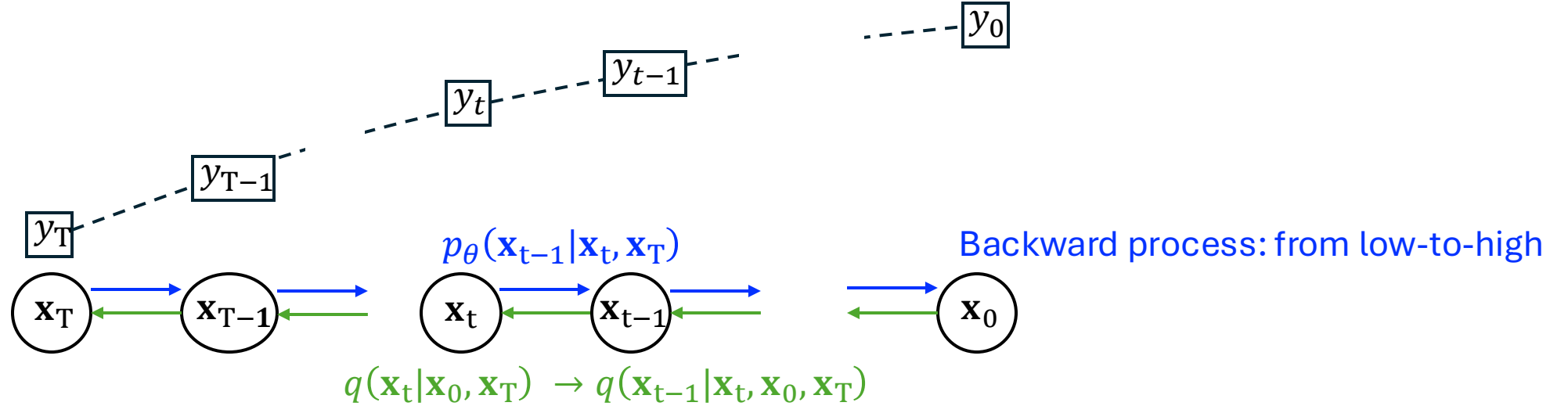
$\rightarrow$This implies a marginal Gaussian:

$$q(\mathbf{x}_t | \mathbf{x}_0, \mathbf{x}_T) = \mathbb{N}\big(\mathbf{x}_t; \psi_t(\mathbf{x}_0, \mathbf{x}_T), \kappa_{t,t}\boldsymbol{I}\big)$$

and reveals the backward transition:

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0, \mathbf{x}_T) = \mathbb{N}(\mathbf{x}_{t-1}; \mu(\mathbf{x}_t, \mathbf{x}_0, \mathbf{x}_T), \tilde{\kappa}_{t-1}\boldsymbol{I})$$

with $\mu(\mathbf{x}_t, \mathbf{x}_0, \mathbf{x}_T) = \psi_{t-1}(\mathbf{x}_0, \mathbf{x}_T) + \kappa_{t-1,t}\kappa_{t,t}^{-1}\big(\mathbf{x}_t - \psi_t(\mathbf{x}_0, \mathbf{x}_T)\big)$ and $\tilde{\kappa}_{t-1} = \kappa_{t-1,t-1} - \kappa_{t-1,t}\kappa_{t,t}^{-1}\kappa_{t,t-1}$

[2] Li, Bo, et al. "Bbdm: Image-to-image translation with brownian bridge diffusion models." *Proceedings of the IEEE/CVF conference on computer vision and pattern Recognition*. 2023.

# Learning Probabilistic Bridge Model



$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_T)$$

Backward process: from low-to-high

$$q(\mathbf{x}_t|\mathbf{x}_0, \mathbf{x}_T) \rightarrow q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{x}_T)$$

To learn the target-agnostic transformation, we parameterize:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_T) = \mathbb{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, \mathbf{x}_T, t), \tilde{\kappa}_{t-1}\boldsymbol{I})$$

Then $\theta$ is learnt via:
$$\theta_{PB} = \underset{\theta}{\text{argmin}} \, \mathbb{E}_{(\boldsymbol{x}_0, \boldsymbol{x}_T, t)}\big[D_{KL}\big(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0, \mathbf{x}_T)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_T)\big)\big] \qquad (1)$$
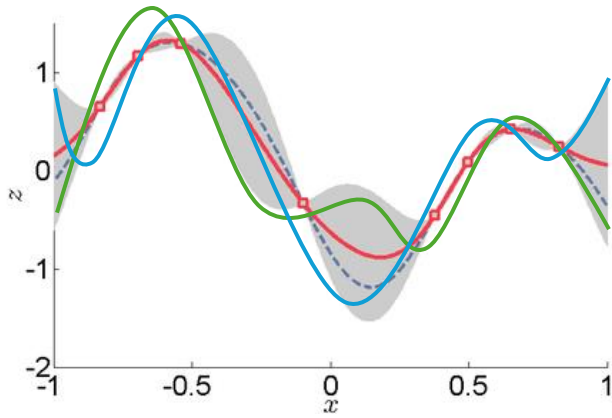
Given $\theta_{PB}$ and low-value design $\mathbf{x}_T$, we can simulate high-value $\mathbf{x}_0$:

$$\mathbf{x}_{t-1} = \mu_{\theta_{PB}}(\mathbf{x}_t, \mathbf{x}_T, t) + \sqrt{\tilde{\kappa}_t}\boldsymbol{\epsilon}$$
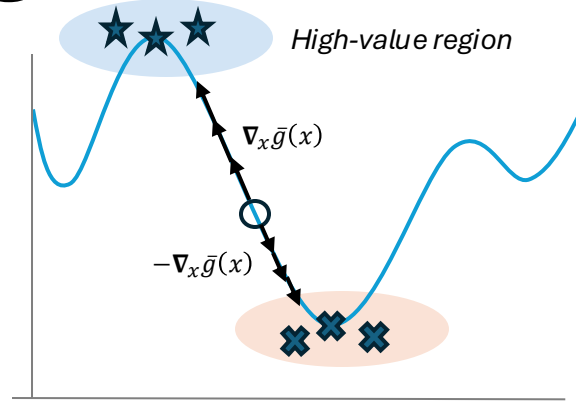
where $\boldsymbol{\epsilon} \sim \mathbb{N}(0, \boldsymbol{I})$ when $t > 1$ and $\boldsymbol{\epsilon} = 0$ otherwise.

**Challenge**: The loss function (1) requires pairs of data $(\mathbf{x}_0, \mathbf{x}_T)$
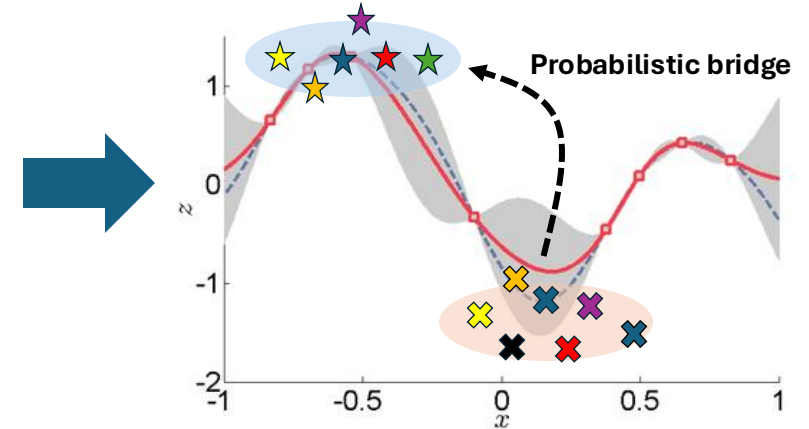
# Synthetic data generation



**Step 1**: Fit multiple GPs          **Step 2**: Take GA, GD on each GP          Final synthetic dataset

**Step 1:** Multiple GP posterior mean functions are used as synthetic functions:

$$\bar{g}_{\phi_s}(\boldsymbol{x}) = k(\phi_s)^T (K(\phi_s) + \sigma^2 I)^{-1} y \quad \text{where kernel } k \text{ is parameterized by } \phi_s$$

**Step 2:** Find high- and low- value samples of $\bar{g}_{\phi_s}$ via taking $M$ gradient ascent and descent steps:

$$X_s^+ = \left\{ \boldsymbol{x}_M^+ = \boldsymbol{x}_0^+ + \eta \sum_{m=0}^{M} \nabla_{\boldsymbol{x}} \bar{g}_{\phi_s}(\boldsymbol{x}_m^+) \Big|_{\boldsymbol{x}_0^+ \in D} \right\}$$

$$X_s^- = \left\{ \boldsymbol{x}_M^- = \boldsymbol{x}_0^- - \eta \sum_{m=0}^{M} \nabla_{\boldsymbol{x}} \bar{g}_{\phi_s}(\boldsymbol{x}_m^-) \Big|_{\boldsymbol{x}_0^- \in D} \right\}$$

Eventually, a **synthetic dataset $D_s$**:

$$D_s = \{(X_s^-, y_s^-); (X_s^+, y_s^+)\}_{i=1}^{n_g} \text{ where } y_s^- = \bar{g}_{\phi_s}(X_s^-), y_s^+ = \bar{g}_{\phi_s}(X_s^+)$$

# Experiment settings

- **Benchmark Tasks.** 4 real-world tasks from the Design-Bench[3] and 3 RNA-Binding tasks from ViennaRNA[4].

| Dataset | Size | Dimensions | Categories | Type |
|---|---|---|---|---|
| Ant Morphology | 25009 | 60 | N/A | Continuous |
| D'Kitty Morphology | 25009 | 56 | N/A | Continuous |
| TF Bind 8 | 32898 | 8 | 4 | Discrete |
| TF Bind 10 | 50000 | 10 | 4 | Discrete |
| RNA-Binding | 5000 | 14 | 4 | Discrete |

- **Baselines.** 21 widely recognized methods.
- **Evaluation Protocol.** 128 candidates are generated. The performances are recorded at 50th, 80th, and 100th percentiles (8 seeds).

[3] Trabucco et al, "*Design-Bench: Benchmarks for Data-Driven Offline Model-Based Optimization*", ICML 2022
[4] Lorenz et al, "*ViennaRNA Package 2.0*", *Algorithms for Molecular Biology 2011*

# Experiment results

## Table 1: Experimental results on Design-bench

| Method | Ant | D'Kitty | TFBind8 | TFBind10 | Mean Rank |
|--------|-----|---------|---------|----------|-----------|
| $D_o$ (best) | 0.565 | 0.884 | 0.439 | 0.467 | - |
| BO-qEI | 0.812 ± 0.000 | 0.896 ± 0.000 | 0.825 ± 0.091 | 0.627 ± 0.033 | 16.75 / 22 |
| CMA-ES | 1.561 ± 0.896 | 0.724 ± 0.001 | 0.939 ± 0.039 | 0.664 ± 0.034 | 8.00 / 22 |
| REINFORCE | 0.263 ± 0.026 | 0.573 ± 0.204 | 0.961 ± 0.034 | 0.618 ± 0.011 | 17.00 / 22 |
| GA | 0.293 ± 0.029 | 0.860 ± 0.021 | 0.985 ± 0.011 | 0.638 ± 0.032 | 12.75 / 22 |
| COMs | 0.882 ± 0.044 | 0.932 ± 0.006 | 0.940 ± 0.027 | 0.621 ± 0.033 | 13.25 / 22 |
| CbAS | 0.846 ± 0.033 | 0.895 ± 0.016 | 0.903 ± 0.028 | 0.649 ± 0.055 | 12.50 / 22 |
| MINs | 0.894 ± 0.022 | 0.939 ± 0.004 | 0.908 ± 0.063 | 0.630 ± 0.019 | 12.50 / 22 |
| GA on GP | 0.948 ± 0.013 | 0.946 ± 0.001 | 0.770 ± 0.087 | 0.654 ± 0.038 | 9.25 / 22 |
| RoMA | 0.593 ± 0.066 | 0.829 ± 0.020 | 0.665 ± 0.000 | 0.553 ± 0.000 | 20.00 / 22 |
| ICT | 0.911 ± 0.030 | 0.945 ± 0.011 | 0.888 ± 0.047 | 0.624 ± 0.033 | 13.50 / 22 |
| Tri-mentoring | 0.944 ± 0.033 | 0.950 ± 0.015 | 0.899 ± 0.045 | 0.647 ± 0.039 | 9.00 / 22 |
| MATCH-OPT | 0.931 ± 0.011 | 0.957 ± 0.014 | 0.977 ± 0.004 | 0.543 ± 0.002 | 9.50 / 22 |
| PGS | 0.949 ± 0.017 | 0.966 ± 0.013 | 0.981 ± 0.015 | 0.532 ± 0.000 | 7.75 / 22 |
| LTR | 0.907 ± 0.032 | 0.960 ± 0.014 | 0.973 ± 0.000 | 0.652 ± 0.039 | 6.25 / 22 |
| DDOM | 0.930 ± 0.029 | 0.925 ± 0.008 | 0.885 ± 0.061 | 0.634 ± 0.015 | 13.75 / 22 |
| GTG | 0.865 ± 0.040 | 0.935 ± 0.010 | 0.901 ± 0.039 | 0.639 ± 0.016 | 12.50 / 22 |
| BDI | 0.964 ± 0.000 | 0.941 ± 0.000 | 0.973 ± 0.000 | 0.636 ± 0.020 | 7.50 / 22 |
| RGD | 0.922 ± 0.020 | 0.883 ± 0.014 | 0.889 ± 0.068 | 0.644 ± 0.048 | 13.00 / 22 |
| BONET | 0.948 ± 0.025 | 0.957 ± 0.008 | 0.894 ± 0.086 | 0.606 ± 0.024 | 10.75 / 22 |
| GABO | 0.224 ± 0.051 | 0.719 ± 0.001 | 0.939 ± 0.038 | 0.639 ± 0.033 | 15.25 / 22 |
| DEMO | 0.948 ± 0.013 | 0.956 ± 0.011 | 0.812 ± 0.054 | 0.648 ± 0.042 | 9.25 / 22 |
| **ROOT (ours)** | 0.965 ± 0.014 | 0.972 ± 0.005 | 0.986 ± 0.007 | 0.685 ± 0.053 | 1.25 / 22 |

## Table 2: Results Biological RNA Tasks

| Method | RNA-A | RNA-B | RNA-C | Mean Rank |
|--------|-------|-------|-------|-----------|
| CbAS | 0.270 ± 0.098 | 0.249 ± 0.088 | 0.261 ± 0.093 | 6.00 / 8 |
| BO-qEI | 0.537 ± 0.106 | 0.517 ± 0.108 | 0.481 ± 0.100 | 3.67 / 8 |
| GA | 0.518 ± 0.120 | 0.499 ± 0.100 | 0.496 ± 0.091 | 4.33 / 8 |
| COMs | 0.187 ± 0.123 | 0.144 ± 0.121 | 0.209 ± 0.100 | 7.67 / 8 |
| REINFORCE | 0.166 ± 0.096 | 0.149 ± 0.081 | 0.225 ± 0.075 | 7.33 / 8 |
| BDI | 0.604 ± 0.000 | 0.505 ± 0.000 | 0.411 ± 0.000 | 4.00 / 8 |
| Boot-Gen | 0.913 ± 0.064 | 0.881 ± 0.024 | 0.786 ± 0.039 | 2.00 / 8 |
| **ROOT (ours)** | 0.956 ± 0.023 | 0.955 ± 0.013 | 0.922 ± 0.013 | 1.00 / 8 |

ROOT establishes a new state-of-the-art performance with the highest rank across benchmarks.