

RNAGYM: BENCHMARKS FOR RNA FITNESS AND STRUCTURE PREDICTION

Rohit Arora *
Systems Biology,
Harvard Medical School

Murphy Angelo *
Systems Biology,
Harvard Medical School

Christian Andrew Choe *
Department of Bioengineering,
Stanford University

Aaron W. Kollasch
Systems Biology,
Harvard Medical School

Fiona Qu
Systems Biology,
Harvard Medical School

Courtney A. Shearer
Systems Biology,
Harvard Medical School

Ruben Weitzman
Systems Biology,
Harvard Medical School

Artem Gazizov
Systems Biology,
Harvard Medical School

Sarah Gurev
Department of EECS,
MIT

Erik Xie
Department of EECS,
MIT

Debora S. Marks †
Systems Biology,
Harvard Medical School

Pascal Notin †
Systems Biology,
Harvard Medical School

ABSTRACT

Predicting the structure and the effects of mutations in RNA are pivotal for numerous biological and medical applications. However, the evaluation of machine learning-based RNA models has been hampered by disparate and limited experimental datasets, along with inconsistent model performances across different RNA types. To address these limitations, we introduce RNAGym, a comprehensive and large-scale benchmark specifically tailored for RNA fitness and structure prediction. This benchmark suite includes nearly 60 standardized deep mutational scanning assays, covering hundreds of thousands of mutations, and curated RNA structure datasets. We have developed a robust evaluation framework that integrates multiple metrics suitable for both predictive tasks while accounting for the inherent limitations of experimental methods. RNAGym is designed to facilitate a systematic comparison of RNA models, offering an essential resource to enhance the development and understanding of these models within the computational biology community.

1 INTRODUCTION

RNA, once considered a passive intermediate between DNA and protein, is now recognized as a dynamic and crucial agent in cellular process regulation. The complexity of RNA lies in its structure — the base-pairing patterns that form the backbone as well as its two- and three-dimensional architecture — and in its functional versatility, which ranges from catalytic activities to gene regulation. Predicting RNA structure and assessing the functional impact of sequence variations, ie. RNA fitness, are key challenges in computational biology and machine learning. These interrelated tasks are critical for advancing our understanding of RNA biology and its applications in fields such as drug discovery and synthetic biology.

The prediction of RNA structure remains a significant challenge. While computational methods have made substantial progress, they still face considerable hurdles, especially for larger RNAs (> 100 nt) with complex features like multi-branched loops and pseudoknots. Experimental methods like

*Equal contributions. †Senior authorship. Correspondence: rohitarora@g.harvard.edu, man-gelo@fas.harvard.edu, cachoe@stanford.edu, debbie@hms.harvard.edu, pascal_notin@hms.harvard.edu

nuclear magnetic resonance, Cryo-EM, and X-ray crystallography can determine RNA 3D structures, but they face technical limitations when applied to RNA, resulting in RNA structures comprising less than 1% of entries in the Protein Data Bank (PDB) (Burley et al., 2017). This scarcity of experimental data further complicates the development and validation of computational prediction methods.

An equally critical task is predicting RNA fitness – the functional capacity of RNA sequences when subjected to mutations. Understanding the impact of these mutations on RNA function is crucial for advancing our knowledge of RNA evolution and its role in cellular processes. This task is also vital for the development of RNA-based therapeutics and the expansion of synthetic biology applications, such as designing riboswitches for gene regulation or engineering RNA sensors for metabolite detection. Despite its importance, accurately predicting the functional consequences of RNA mutations, especially from sequence data alone, remains a significant challenge in the field. Both structure and fitness prediction can benefit from evolutionary information. For structure prediction, approaches such as maximum entropy models leverage sequence co-variation to infer evolutionary constraints (Weinreb et al., 2016; Hopf et al., 2017; Frazer et al., 2021). Similarly, fitness prediction methods can utilize evolutionary data to identify functionally crucial sequence features. However, robust methodologies for integrating this information and accurately predicting both structure and fitness, particularly in zero-shot scenarios, remain elusive.

To address these challenges and support progress in the field, we present RNAGym, a comprehensive benchmarking framework designed to evaluate and compare computational methods for RNA structure and fitness prediction. RNAGym provides a diverse collection of curated RNA mutational scanning assays and chemical mapping data for structure prediction, multiple metrics for model evaluation, and assesses the relative performance of diverse baselines across both tasks.

RNAGym aims to accelerate progress in computational RNA biology by offering a common platform for assessing different approaches. By providing a systematic way to evaluate model performance across various RNA types and prediction tasks, RNAGym can help identify strengths and weaknesses of current methods, guide the development of more accurate algorithms, and ultimately contribute to advancing our understanding of RNA biology and its applications in areas such as personalized medicine, RNA-based drug design, and engineered RNA devices for synthetic biology.

2 RNAGYM BENCHMARKS

2.1 OVERVIEW

RNAGym is a comprehensive benchmark suite advancing the development and analysis of machine learning RNA models. It comprises three integrated layers: datasets, models, and analytics (Fig. 1), supporting three core RNA tasks:

- **Fitness prediction:** Prediction of RNA functionality across diverse RNA types, leveraging a broad set of deep mutational scanning assays.
- **Secondary Structure prediction:** Prediction of RNA secondary structure, focusing on identifying nucleotide contacts, which is crucial for understanding RNA function.
- **Tertiary Structure prediction:** Prediction of RNA tertiary structure, extending beyond secondary contacts to evaluate critical higher-order interactions.

These tasks, evaluated in a zero-shot setting, challenge models to generalize across varied RNA contexts without task-specific fine-tuning. Our data layer includes curated datasets that are specifically structured for these three tasks. These datasets are enriched with detailed annotations for a variety of RNA types and are classified by mutation depth, enhancing the granularity of the data available for analysis. Across all tasks, RNAGym integrates a diverse array of 15 predictive models, each tailored to address the nuances of the specific tasks at hand—whether predicting RNA fitness or determining RNA structure. The analytics layer of RNAGym is designed to provide a deep and comprehensive evaluation of model performance. It utilizes eleven distinct performance metrics to assess the effectiveness of each model in a clear and quantifiable manner. Further, the framework allows for detailed exploration of model performance across different RNA types and mutation depths, with the goal to understand model strengths and limitations in varied biological contexts.

2.2 DATASETS

2.2.1 FITNESS PREDICTION ASSAYS

Screening methodology Coding mRNA datasets were sourced from ProteinGym when nucleotide-level deep mutational scanning information was available (Notin et al., 2023). For noncoding and splicing studies we conducted a broad PubMed search for RNA mutational studies that yielded over 11,000 results, which we then screened using a LLM with carefully designed prompts adapted from systematic review methods. After narrowing down to 52 studies through the LLM screening, we conducted expert manual review using specific inclusion/exclusion criteria to ensure data quality and relevance. All details regarding search terms, prompts and inclusions/exclusions criteria are provided in Appendix B.

Selected assays RNAGym includes 59 Deep Mutational Scanning assays containing 998,296 variants measuring. Beyond the sheer size of RNAGym, assays span a broad range of experimental contexts from in vivo cellular fitness measurements to in vitro biochemical assays across various tRNA, aptamers, ribozymes, and both coding and splicing disrupting mRNAs (Table A1). This effort represents a *fourfold increase* in size over the largest prior RNA benchmarks for non coding fitness prediction Brixi et al. (2025). Unlike previous efforts, these assays are integrated into a standardized, reusable public resource, making RNAGym a more accessible and broadly applicable tool for RNA fitness prediction. Our assays all follow a standard format (see Appendix C).

2.2.2 2° STRUCTURE PREDICTION DATA

In preparing the benchmark for our research paper, we utilized the dataset from the Stanford Ribonanza Challenge, which contains chemical mapping data for many RNA sequences. Chemical mapping is a way to chemically probe an RNA structure by reacting the RNA with various reagent, commonly DMS (dimethyl sulfate) and 2A3 (2-aminopyridine-3-carboxylic acid imidazolide). Once an RNA is reacted with the chemical mapping reagent, the reactivity of each nucleotide can be read by reverse transcription and high throughput sequencing. The reactivity profile gives a one dimensional view of the RNA structure where high reactivity is associated with single-stranded regions (unpaired RNA) and low reactivity is associated with double-stranded regions (paired RNA) (Spitale & Incarnato, 2023). This type of data is invaluable for validating computational models of RNA secondary structure prediction, as it offers direct evidence of the RNA’s physical structure (Table A2). For completeness, we included data on ligand binding, cotranscriptional folding, and RNA degradation. We compiled the data from RMDB (Cordero et al., 2012) and filtered for the sequences with signal-to-noise ratio ≥ 1.0 . For our analysis we focused on only in vitro chemical mapping data involving DMS and 2A3 for ease of interpretation. This resulted in a high quality dataset with 901k reactivity profiles for 583k unique sequences and 100M nucleotides (Table A10). The test dataset provides a score for each nucleotide for all RNA sequences (1 row per nucleotide), reflecting the propensity of that nucleotide to be single-stranded in the RNA structure.

2.2.3 3° STRUCTURE PREDICTION DATA

RNAGym also introduces a 3D RNA structure dataset that emphasizes evaluation across a wide range of lengths, functional classes, and degree of structure and sequence homology to training data (Table A3). Starting with all 21,000+ RNA structures in the PDB, we filtered for RNAs published after January 12, 2023, ensuring no overlap with any baseline’s training set. The remaining RNAs were annotated with $\mathbf{TM}_{\text{train}}$ and $\mathbf{ID}_{\text{train}}$, representing the maximum structure and sequence homology, respectively, of each evaluation candidate to any RNA published prior to the aforementioned date cutoff. Finally, the top three structures of each RNA family—ranked by $\mathbf{TM}_{\text{train}}$, $\mathbf{ID}_{\text{train}}$, length (up to 2,000 nucleotides), and resolution (up to 5.0Å)—were included in the overall dataset. These metrics were paired with unsupervised evolutionary information from EVCouplings to contextualize model performance. This framework allows us to assess not only how well models predict a variety of RNA tertiary structures but also how they generalize to novel folds and how well they extract coevolutionary signals.

2.3 BASELINES

We benchmarked several RNA models including RiNALMo (650M parameters) (Penić et al., 2024), EVO 1 (7B parameters) (Nguyen et al., 2024) EVO 1.5 (Merchant et al., 2024), EVO 2 (Brix et al., 2025), RNA-FM (Chen et al., 2022), GenSLM (Zvyagin et al., 2023a), RNAErnie (Wang et al., 2024a), and Nucleotide Transformer (Dalla-Torre et al., 2023) for fitness prediction, as well as EternaFold (Wayment-Steele et al., 2022), CONTRAfold (Do et al., 2006), Vienna (Gruber et al., 2008), an RNAstructure (Reuter & Mathews, 2010) for 2° structure prediction. Models for RNA fitness prediction are generally large transformer based models with hundreds of millions to billions of parameters. For example, EVO (7B) and nucleotide transformer (2.5B) are foundation models trained on genome-scale data, while smaller BERT-like models like RiNALMo (650M) and RNAErnie (80M) target non-coding RNAs. On the other hand, RNA secondary structure prediction models include classical thermodynamic and statistical folding algorithms (Vienna RNAfold, RNAstructure) and earlier ML models (CONTRAfold) alongside newer specialized methods (EternaFold). Our tertiary structure prediction benchmarks currently include AlphaFold3 (Abramson et al.), NuFold (Kagaya et al.), RhoFold+ (Shen et al.), RosettaFold2NA (Baek et al.), and trRosettaRNA (Wang et al.), selected to encompass a range of methodological approaches in RNA 3D structure prediction. These models represent the forefront of RNA structure prediction, each bringing unique strengths to address the complexities of RNA folding and interactions. All details about baselines are provided in Appendix D.

2.4 EVALUATION

For **fitness prediction**, the evaluation was primarily based on the Spearman’s rank correlation between the model predictions and experimental measurements, the Area Under the Curve (AUC) and the Matthews Correlation Coefficient (MCC). These metrics are complementary and were chosen to provide a comprehensive evaluation: Spearman correlation assesses the overall ranking of predictions, AUC measures the model’s ability to distinguish between functional and non-functional mutations, while MCC offers a balanced measure for potentially imbalanced datasets. To mitigate biases associated with uneven assay distributions across different RNA types, we calculated an average performance for each RNA type separately and then computed the overall performance as the mean of these RNA-type-level averages. This approach ensures that our results are robust and reflective of true model capabilities across varied biological categories. For the **2° structure prediction task**, we employed three standard metrics: F1-score, Area Under the Curve (AUC), and Mean Absolute Error (MAE). F1-score provides a balanced measure of precision and recall in identifying nucleotide pairings. AUC assesses the model’s ability to distinguish between paired and unpaired nucleotides. MAE offers a direct measure of prediction accuracy by quantifying the average magnitude of errors. For the **3° structure prediction task**, we employed TM score, Δ TM score, Watson-Crick intra-network fidelity (INF_{WC}), and non-Watson-Crick INF (INF_{NWC}), ρ_{TM} , and ρ_{EC} . TM-score captures overall fold correctness, while Δ TM-score captures predictive performance relative to the best scoring template structure. INF_{WC} and INF_{NWC} assess whether models correctly infer canonical and non-canonical base pairing. ρ_{TM} , and ρ_{EC} represent Spearman’s rank correlation with TM_{train} and with the number of correct evolutionary couplings in the input alignment, respectively. In combination, these metrics gauge both global RNA fold and local base-pair fidelity, while ensuring equitable assessment despite distinct baseline training sets.

3 RESULTS

3.1 FITNESS PREDICTION PERFORMANCE

The overall fitness prediction benchmark results (see Table A6) show Evo (2.0) and RNAErnie as the leading performers, with Spearman correlations of 0.286 and 0.221 respectively. These results indicate a leading capability for Evo 2.0 in predicting RNA fitness outcomes based on experimental data (statistical significance analysis is included in Appendix E.1). The relatively low scores across all models, particularly when compared to the stronger correlations reported for protein language models Notin et al. (2023), suggest substantial room for improvement and warrant deeper investigation. Several factors may contribute to this performance gap. A primary consideration is the limited availability of large-scale, diverse RNA datasets for model training compared to the

abundance of protein sequence data. Additionally, there may be a potential misalignment between the training data used for these models and the taxonomical and functional distribution of our fitness landscapes. Lastly, differences in evolutionary conservation patterns between RNAs (especially non-coding RNAs) and proteins could also play a role, potentially affecting the models’ ability to capture fitness-relevant features. When examining performance by RNA type (Table A7), several models show specialized strengths across different RNA categories. RNA-FM achieves the highest correlations across all non coding RNAs including tRNA (0.464), Aptamer (0.190), and Ribozymes (0.201). On the other hand, the Evo family of models, and particularly Evo 2.0 have the best performances on mRNA-splicing (0.431) and mRNA coding (0.323) assays. These performance variations likely reflect the diverse training data of each model. RNA-FM’s particular strength with tRNAs, aptamers, and ribozymes aligns with its training on non-coding RNAs from RNACentral, a database rich in these RNA types. The consistently strong performance of Evo models across different RNA types, and particularly in mRNAs, suggests their training approach may capture broader sequence-function relationships. These observations underscore the importance of targeted model training and selection based on the specific RNA type being studied. They also suggest that performance could potentially be improved by more tailored training data selection or by developing ensemble methods that leverage the strengths of different models for specific RNA types.

3.2 2° STRUCTURE PREDICTION PERFORMANCE

The RNAGym structure prediction benchmark (Table A11) reveals interesting performance patterns across different RNA structure prediction methods. When considering overall performance, EternaFold demonstrates the strongest results among unsupervised methods, achieving an F1-score of 0.658, AUC of 0.714, and MAE of 0.337. CONTRAfold follows closely with an F1-score of 0.654, then RNAstructure (0.652) and Vienna (0.645). This relatively tight performance distribution suggests that these unsupervised approaches, while effective, may be approaching a methodological ceiling in their current framework.

3.3 3° STRUCTURE PREDICTION PERFORMANCE

The overall tertiary structure benchmark results (see Table A13) show that NuFold leads for monomeric RNAs, with a top TM score of 0.393, while AlphaFold3 emerges as the best performer on complexes (TM = 0.381). NuFold’s correlations with training-set similarity ($\rho_{TM} = 0.74$) and evolutionary couplings ($\rho_{EC} = 0.67$) are also the highest for monomers, implying that it leverages both memorized structural priors and covariation signals effectively. Notably, the top monomeric models exhibit modest ΔTM values ($\Delta TM \approx -0.15$), suggesting that, despite relying on different training sets they learn comparable structural models with room for improvement when generalizing beyond memorized templates. AlphaFold3 tops total Watson-Crick ($INF_{WC}=0.83$) and non-Watson-Crick interactions ($INF_{NWC}=0.26$), followed closely by NuFold, RoseTTAFold2NA, and trRosettaRNA. The notably lower scores for non-Watson-Crick interactions underscores a significant gap in RNA structure prediction, highlighting an area where future advances are urgently needed. For RNA complexes, AlphaFold3 outperforms RoseTTAFold2NA, achieving a TM score of 0.380 vs. 0.167. Interestingly, AlphaFold3’s performance on complexes matches or exceeds its performance on monomers despite less homology to training ($\Delta TM = -0.13$), highlighting its robust handling of multi-chain interfaces. Moreover, it shows higher correlation with evolutionary couplings ($\rho_{EC} = 0.77$) than with training-set similarity ($\rho_{TM} = 0.55$), suggesting that it can leverage covariation signals effectively even in these more challenging contexts. The pronounced drop-off in performances for RoseTTAFold2NA indicates that complexes remain difficult for certain architectures.

4 CONCLUSION

RNAGym addresses the significant gap in large-scale benchmarks for the robust evaluation of models tailored for RNA structure prediction and fitness assessment. It enables the direct comparison of methods across several dimensions of interest (e.g., RNA type, mutation type). We anticipate that the RNAGym benchmarks and the accompanying data assets we release to the public will serve as invaluable resources for the Machine Learning and Computational Biology communities. We plan to continually update the benchmarks as new data and baseline models become available.

REFERENCES

- 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, Alexander I. Cowen-Rivers, Andrew Cowie, Michael Figurnov, Fabian B. Fuchs, Hannah Gladman, Rishub Jain, Yousuf A. Khan, Caroline M. R. Low, Kuba Perlin, Anna Potapenko, Pascal Savy, Sukhdeep Singh, Adrian Stecula, Ashok Thillaisundaram, Catherine Tong, Sergei Yakneen, Ellen D. Zhong, Michal Zielinski, Augustin Židek, Victor Bapst, Pushmeet Kohli, Max Jaderberg, Demis Hassabis, and John M. Jumper. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature*, 630(8016):493–500. ISSN 1476-4687. doi: 10.1038/s41586-024-07487-w. URL <https://doi.org/10.1038/s41586-024-07487-w>.
- Bharat V. Adkar, Arti Tripathi, Anusmita Sahoo, Kanika Bajaj, Devrishi Goswami, Purbani Chakrabarti, Mohit K. Swarnkar, Rajesh S. Gokhale, and Raghavan Varadarajan. Protein model discrimination using mutational sensitivity derived from deep sequencing. *Structure*, 20(2):371–381, Feb 2012. ISSN 0969-2126. doi: 10.1016/j.str.2011.11.021. URL <https://doi.org/10.1016/j.str.2011.11.021>.
- Johan O. L. Andreasson, Andrew Savinov, Steven M. Block, and William J. Greenleaf. Comprehensive sequence-to-function mapping of cofactor-dependent rna catalysis in the glms ribozyme. *Nature Communications*, 2020. doi: <https://doi.org/10.1038/s41467-020-15540-1>. URL <https://www.nature.com/articles/s41467-020-15540-1#citeas>.
- Minkyung Baek, Ryan McHugh, Ivan Anishchenko, Hanlun Jiang, David Baker, and Frank DiMaio. Accurate prediction of protein–nucleic acid complexes using RoseTTAFoldNA. *Nature Methods*, 21(1): 117–121. ISSN 1548-7105. doi: 10.1038/s41592-023-02086-5. URL <https://doi.org/10.1038/s41592-023-02086-5>.
- David Baker and George Church. Protein design meets biosecurity. *Science*, 383:349 – 349, 2024. URL <https://api.semanticscholar.org/CorpusID:267212249>.
- James D Beck, Jessica M Roberts, Joey M Kitzhaber, Ashlyn Trapp, Edoardo Serra, Francesca Spezzano, and Eric J Hayden. Predicting higher-order mutational effects in an rna enzyme by machine learning of high-throughput experimental data. *Frontiers Mol. Biosci.*, 2022. doi: 10.3389/fmolb.2022.893864. URL <https://www.frontiersin.org/articles/10.3389/fmolb.2022.893864/full>.
- Garyk Brixi, Matthew G. Durrant, Jerome Ku, Michael Poli, Greg Brockman, Daniel Chang, Gabriel A. Gonzalez, Samuel H. King, David B. Li, Aditi T. Merchant, Mohsen Naghipourfar, Eric Nguyen, Chiara Ricci-Tam, David W. Romero, Gwanggyu Sun, Ali Taghibakshi, Anton Vorontsov, Brandon Yang, Myra Deng, Liv Gorton, Nam Nguyen, Nicholas K. Wang, Etowah Adams, Stephen A. Baccus, Steven Dillmann, Stefano Ermon, Daniel Guo, Rajesh Ilango, Ken Janik, Amy X. Lu, Reshma Mehta, Mohammad R.K. Mofrad, Madelena Y. Ng, Jaspreet Pannu, Christopher Ré, Jonathan C. Schmok, John St. John, Jeremy Sullivan, Kevin Zhu, Greg Zynda, Daniel Balsam, Patrick Collison, Anthony B. Costa, Tina Hernandez-Boussard, Eric Ho, Ming-Yu Liu, Thomas McGrath, Kimberly Powell, Dave P. Burke, Hani Goodarzi, Patrick D. Hsu, and Brian L. Hie. Genome modeling and design across all domains of life with evo 2. *bioRxiv*, 2025. doi: 10.1101/2025.02.18.638918. URL <https://www.biorxiv.org/content/early/2025/02/21/2025.02.18.638918>.
- Stephen K Burley, Helen M Berman, Gerard J Kleywegt, John L Markley, Haruki Nakamura, and Sameer Velankar. Protein data bank (pdb): the single global macromolecular structure archive. *Protein Crystallography*, pp. 627–641, 2017.
- Christian Cao, Jason Sang, Rohit Arora, Robert Kloosterman, Matthew Cecere, Jaswanth Gorla, Richard Saleh, David Chen, Ian Drennan, Bijan Teja, Michael Fehlings, Paul Ronksley, Alexander A Leung, Dany Weisz, Harriet Ware, Mairead Whelan, David B Emerson, Rahul Krishan Arora, and Niklas Bobrovitz. Prompting is all you need: Lms for systematic review screening. *medRxiv*, 2024. doi: 10.1101/2024.06.01.24308323.
- Jiayang Chen, Zhihang Hu, Siqi Sun, Qingxiong Tan, Yixuan Wang, Qinze Yu, Licheng Zong, Liang Hong, Jin Xiao, Tao Shen, Irwin King, and Yu Li. Interpretable rna foundation model from unannotated data for highly accurate rna structure and function predictions, 2022.
- Pablo Cordero, Julius B. Lucks, and Rhiju Das. An rna mapping database for curating rna structure mapping experiments. *Bioinformatics*, 28(22):3006–3008, 09 2012. ISSN 1367-4803. doi: 10.1093/bioinformatics/bts554. URL <https://doi.org/10.1093/bioinformatics/bts554>.
- José Almeida Cruz, Marc-Frédéric Blanchet, Michal J. Boniecki, Janusz M. Bujnicki, Shi-Jie Chen, Song Cao, Rhiju Das, Feng Ding, Nikolay V. Dokholyan, Samuel Coulbourn Flores, Lili Huang, Christopher A. Lavender, Véronique Lisi, François Major, Katarzyna Mikolajczak, Dinshaw J. Patel, Anna Philips, Tomasz

- Puton, John SantaLucia, Fredrick Sijenyi, Thomas Hermann, Kristian Rother, Magdalena Rother, Alexander Serganov, Marcin Skorupski, Tomasz Soltysinski, Parin Sripakdeevong, Irina Tuszynska, Kevin M. Weeks, Christina Waldsich, Michael Wildauer, Neocles B. Leontis, and Eric Westhof. Rna-puzzles: a casp-like evaluation of rna three-dimensional structure prediction. *RNA*, 18 4:610–25, 2012. URL <https://api.semanticscholar.org/CorpusID:263498187>.
- Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza-Revilla, Nicolas Lopez Carranza, Adam Henryk Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Bernardo P. de Almeida, Hassan Sirelkhatim, Guillaume Richard, Marcin Skwark, Karim Beguir, Marie Lopez, and Thomas Pierrot. The nucleotide transformer: Building and evaluating robust foundation models for human genomics. *bioRxiv*, 2023. doi: 10.1101/2023.01.11.523679. URL <https://www.biorxiv.org/content/early/2023/09/19/2023.01.11.523679>.
- Padideh Danaee, Mason Rouches, Michelle Wiley, Dezhong Deng, Liang Huang, and David Hendrix. bprna: large-scale automated annotation and analysis of rna secondary structure. *Nucleic Acids Research*, 46(11): 5381–5394, 05 2018. ISSN 0305-1048. doi: 10.1093/nar/gky285. URL <https://doi.org/10.1093/nar/gky285>.
- David Ding, Ada Y. Shaw, Sam Sinai, Nathan Rollins, Noam Prywes, David F. Savage, Michael T. Laub, and Debora S. Marks. Protein design using structure-based residue preferences. *Nature Communications*, 15(1): 1639, Feb 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-45621-4. URL <https://doi.org/10.1038/s41467-024-45621-4>.
- Chuong B. Do, Daniel A. Woods, and Serafim Batzoglou. Contrafold: Rna secondary structure prediction without physics-based models. *Bioinformatics*, 22 14:e90–8, 2006. URL <https://api.semanticscholar.org/CorpusID:1646946>.
- Júlia Domingo, Guillaume Diss, and Ben Lehner. Pairwise and higher-order genetic interactions during the evolution of a trna. *Nature*, 2018. doi: 10.1038/s41586-018-0170-7. URL <https://api.semanticscholar.org/CorpusID:240071819>.
- Gregory M. Findlay, Riza M. Daza, Beth Martin, Melissa D. Zhang, Anh P. Leith, Molly Gasperini, Joseph D. Janizek, Xingfan Huang, Lea M. Starita, and Jay Shendure. Accurate classification of brca1 variants with saturation genome editing. *Nature*, 562(7726):217–222, Oct 2018. ISSN 1476-4687. doi: 10.1038/s41586-018-0461-z. URL <https://doi.org/10.1038/s41586-018-0461-z>.
- Elad Firnberg, Jason W. Labonte, Jeffrey J. Gray, and Marc Ostermeier. A comprehensive, high-resolution map of a gene’s fitness landscape. *Molecular Biology and Evolution*, 31(6):1581–1592, 02 2014. ISSN 0737-4038. doi: 10.1093/molbev/msu081. URL <https://doi.org/10.1093/molbev/msu081>.
- Jonathan Frazer, Pascal Notin, Mafalda Dias, Aidan N. Gomez, Joseph K Min, Kelly P. Brock, Yarin Gal, and Debora S. Marks. Disease variant prediction with deep generative models of evolutionary data. *Nature*, 599: 91–95, 2021. URL <https://api.semanticscholar.org/CorpusID:240071819>.
- Andreas R. Gruber, Ronny Lorenz, Stephan H. Bernhart, Richard Neuböck, and Ivo L. Hofacker. The vienna rna websuite. *Nucleic Acids Research*, 36:W70–W74, 2008. URL <https://api.semanticscholar.org/CorpusID:6481000>.
- Michael P. Guy, David L. Young, Matthew J. Payea, Xiaoju Zhang, Yoshiko Kon, Kimberly M. Dean, Elizabeth J. Grayhack, David H. Mathews, Stanley Fields, and Eric M. Phizicky. Identification of the determinants of trna function and susceptibility to rapid trna decay by high-throughput in vivo analysis. *Genes and Development*, 2014. doi: 10.1101/gad.245936.114. URL <https://genesdev.cshlp.org/content/28/15/1721.long>.
- Thomas A Hopf, John B Ingraham, Frank J Poelwijk, Charlotta PI Schärfe, Michael Springer, Chris Sander, and Debora S Marks. Mutation effects predicted from sequence co-variation. *Nature biotechnology*, 35(2): 128–135, 2017.
- Hervé Jacquier, André Birgy, Hervé Le Nagard, Yves Mechulam, Emmanuelle Schmitt, Jérémy Glodt, Beatrice Bercot, Emmanuelle Petit, Julie Poulain, Guilène Barnaud, Pierre-Alexis Gros, and Olivier Tenailon. Capturing the mutational landscape of the beta-lactamase tem-1. *Proceedings of the National Academy of Sciences*, 110(32):13067–13072, 2013. doi: 10.1073/pnas.1215206110. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1215206110>.
- Evan Janzen, Yuning Shen, Alberto Vázquez-Salazar, Ziwei Liu, Celia Blanco, Josh Kenchel, and Irene A. Chen. Emergent properties as by-products of prebiotic evolution of aminoacylation ribozymes. *Nature Communications*, 2022. doi: 10.1038/s41467-022-31387-0. URL <https://www.nature.com/articles/s41467-022-31387-0>.

- Philippe Julien, Belén Miñana, Pablo Baeza-Centurion, Juan Valcárcel, and Ben Lehner. The complete local genotype–phenotype landscape for the alternative splicing of a human exon. *Nature communications*, 2016. doi: <https://doi.org/10.1038/ncomms11558>. URL <https://www.nature.com/articles/ncomms11558>.
- Yuki Kagaya, Zicong Zhang, Nabil Ibtchaz, Xiao Wang, Tsukasa Nakamura, Pranav Deep Punuru, and Daisuke Kihara. Nufold: End-to-end approach for RNA tertiary structure prediction with flexible nucleobase center representation. *Nature Communications*, 16(1):881. ISSN 2041-1723. doi: 10.1038/s41467-025-56261-7. URL <https://doi.org/10.1038/s41467-025-56261-7>.
- Shengdong Ke, Vincent Anquetil, Jorge Rojas Zamalloa, Alisha Maity, Anthony Yang, Mauricio A. Arias, Sergey Kalachikov, James J. Russo, and Jingyue Juand Lawrence A. Chasin. Saturation mutagenesis reveals manifold determinants of exon definition. *Genome Research*, 2017. doi: 10.1101/gr.219683.116. URL <https://genome.cshlp.org/content/28/1/11.long>.
- Eric D. Kelsic, Hattie Chung, Niv Cohen, Jimin Park, Harris H. Wang, and Roy Kishony. Rna structural determinants of optimal codons revealed by mage-seq. *Cell Systems*, 3(6):563–571.e6, Dec 2016. ISSN 2405-4712. doi: 10.1016/j.cels.2016.11.004. URL <https://doi.org/10.1016/j.cels.2016.11.004>.
- Shungo Kobori and Yohei Yokobayashi. High-throughput mutational analysis of a twister ribozyme. *Angewandte Chemie International Edition*, 55(35):10354–10357, 2016. doi: <https://doi.org/10.1002/anie.201605470>. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/anie.201605470>.
- Shungo Kobori, Yoko Nomura, Anh Miu, and Yohei Yokobayashi. High-throughput assay and engineering of self-cleaving ribozymes by sequencing. *Nucleic Acids Research*, 43(13):e85–e85, 03 2015. ISSN 0305-1048. doi: 10.1093/nar/gkv265. URL <https://doi.org/10.1093/nar/gkv265>.
- Shungo Kobori, Kei Takahashi, and Yohei Yokobayashi. Deep sequencing analysis of aptazyme variants based on a pistol ribozyme. *ACS Synthetic Biology*, 6(7):1283–1288, 2017. doi: 10.1021/acssynbio.7b00057. URL <https://doi.org/10.1021/acssynbio.7b00057>. PMID: 28398719.
- Eran Kotler, Odem Shani, Guy Goldfeld, Maya Lotan-Pompan, Ohad Tarcic, Anat Gershoni, Thomas A. Hopf, Debora S. Marks, Moshe Oren, and Eran Segal. A systematic p53 mutation library links differential functional impact to cancer mutation pattern and evolutionary conservation. *Molecular Cell*, 71(1):178–190.e8, Jul 2018. ISSN 1097-2765. doi: 10.1016/j.molcel.2018.06.012. URL <https://doi.org/10.1016/j.molcel.2018.06.012>.
- Andriy Kryshchak, Maciej Antczak, Marta Szachniuk, Tomasz Zok, Rachael C. Kretsch, Ramya Rangan, Phillip Pham, Rhiju Das, Xavier Robin, Gabriel Studer, Janani Durairaj, Jerome Eberhardt, Aaron Sweeney, Maya Topf, Torsten Schwede, Krzysztof Fidelis, and John Moulton. New prediction categories in casp15. *Proteins*, 91:1550–1557, 2023. URL <https://api.semanticscholar.org/CorpusID:259138147>.
- Juhye M. Lee, John Huddleston, Michael B. Doud, Kathryn A. Hooper, Nicholas C. Wu, Trevor Bedford, and Jesse D. Bloom. Deep mutational scanning of hemagglutinin helps predict evolutionary fates of human h3n2 influenza variants. *Proceedings of the National Academy of Sciences*, 115(35):E8276–E8285, 2018. doi: 10.1073/pnas.1806133115. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1806133115>.
- Chuan Li, Wenfeng Qian, Calum J. Maclean, and Jianzhi Zhang. The fitness landscape of a trna gene. *Science*, 352(6287):837–840, 2016. doi: 10.1126/science.aae0568. URL <https://www.science.org/doi/abs/10.1126/science.aae0568>.
- Mark R. MacRae, Dhenesh Puvanendran, Max A.B. Haase, Nicolas Coudray, Ljivica Kolich, Cherry Lam, Minkyung Baek, Gira Bhabha, and Damian C. Ekiert. Protein–protein interactions in the mla lipid transport system probed by computational structure prediction and deep mutational scanning. *Journal of Biological Chemistry*, 299(6), Jun 2023. ISSN 0021-9258. doi: 10.1016/j.jbc.2023.104744. URL <https://doi.org/10.1016/j.jbc.2023.104744>.
- Richard N. McLaughlin Jr, Frank J. Poelwijk, Arjun Raman, Walraj S. Gosal, and Rama Ranganathan. The spatial architecture of protein function and adaptation. *Nature*, 491(7422):138–142, Nov 2012. ISSN 1476-4687. doi: 10.1038/nature11500. URL <https://doi.org/10.1038/nature11500>.
- Ewan K. S. MacRae, Christopher J. K. Wan, Emil L. Kristoffersen, Kalinka Hansen, Edoardo Gianni, Isaac Gallego, Joseph F. Curran, James Attwater, Philipp Holliger, and Ebbe S. Andersen. Cryo-em structure and functional landscape of an rna polymerase ribozyme. *Proceedings of the National Academy of Sciences*, 121(3):e2313332121, 2024. doi: 10.1073/pnas.2313332121. URL <https://www.pnas.org/doi/abs/10.1073/pnas.2313332121>.

- Gianmarco Meier, Sujani Thavarasah, Kai Ehrenbolger, Cedric A. J. Hutter, Lea M. Hürlimann, Jonas Barandun, and Markus A. Seeger. Deep mutational scan of a drug efflux pump reveals its structure–function landscape. *Nature Chemical Biology*, 19(4):440–450, Apr 2023. ISSN 1552-4469. doi: 10.1038/s41589-022-01205-1. URL <https://doi.org/10.1038/s41589-022-01205-1>.
- Aditi T. Merchant, Samuel H. King, Eric Nguyen, and Brian L. Hie. Semantic mining of functional de novo genes from a genomic language model. *bioRxiv*, 2024. URL <https://api.semanticscholar.org/CorpusID:274893904>.
- Eric Nguyen, Michael Poli, Matthew G Durrant, Armin W Thomas, Brian Kang, Jeremy Sullivan, Madelena Y Ng, Ashley Lewis, Aman Patel, Aaron Lou, Stefano Ermon, Stephen A Baccus, Tina Hernandez-Boussard, Christopher Ré, Patrick D Hsu, and Brian L Hie. Sequence modeling and design from molecular to genome scale with evo. *bioRxiv*, 2024. doi: 10.1101/2024.02.27.582234.
- Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: large-scale benchmarks for protein fitness prediction and design. *Advances in Neural Information Processing Systems*, 36, 2023.
- Christina Nutschel, Alexander Fulton, Olav Zimmermann, Ulrich Schwaneberg, Karl-Erich Jaeger, and Holger Gohlke. Systematically scrutinizing the impact of substitution sites on thermostability and detergent tolerance for bacillus subtilis lipase a. *Journal of Chemical Information and Modeling*, 60(3):1568–1584, Mar 2020. ISSN 1549-9596. doi: 10.1021/acs.jcim.9b00954. URL <https://doi.org/10.1021/acs.jcim.9b00954>.
- Rafael Josip Penić, Tin Vlašić, Roland G. Huber, Yue Wan, and Mile Šikić. Rinalmo: General-purpose rna language models can generalize well on structure prediction tasks, 2024.
- Gianluca Peri, Clémentine Gibard, Nicholas H Shults, Kent Crossin, and Eric J Hayden. Dynamic RNA Fitness Landscapes of a Group I Ribozyme during Changes to the Experimental Environment. *Molecular Biology and Evolution*, 39(3):msab373, 01 2022. ISSN 1537-1719. doi: 10.1093/molbev/msab373. URL <https://doi.org/10.1093/molbev/msab373>.
- Jason N. Pitt and Adrian R. Ferré-D’Amaré. Rapid construction of empirical rna fitness landscapes. *Science*, 330(6002):376–379, 2010. doi: 10.1126/science.1192001. URL <https://www.science.org/doi/abs/10.1126/science.1192001>.
- Yuchen Ren, Zhiyuan Chen, Lifeng Qiao, Hongtai Jing, Yuchen Cai, Sheng Xu, Peng Ye, Xinzhu Ma, Siqi Sun, Hongliang Yan, Dong Yuan, Wanli Ouyang, and Xihui Liu. Beacon: Benchmark for comprehensive rna tasks and language models, 2024. URL <https://arxiv.org/abs/2406.10391>.
- Jessica S. Reuter and David H. Mathews. Rnastructure: software for rna secondary structure prediction and analysis. *BMC Bioinformatics*, 11:129 – 129, 2010. URL <https://api.semanticscholar.org/CorpusID:10356201>.
- Jessica M Roberts, James D Beck, Tanner B Pollock, Devin P Bendixsen, and Eric J Hayden. Rna sequence to structure analysis from comprehensive pairwise mutagenesis of multiple self-cleaving ribozymes. *eLife*, 12:e80360, jan 2023. ISSN 2050-084X. doi: 10.7554/eLife.80360. URL <https://doi.org/10.7554/eLife.80360>.
- Frederic Runge, Karim Farid, Jorg Franke, and Frank Hutter. Rnabench: A comprehensive library for in silico rna modelling, 01 2024.
- Karen S. Sarkisyan, Dmitry A. Bolotin, Margarita V. Meer, Dinara R. Usmanova, Alexander S. Mishin, George V. Sharonov, Dmitry N. Ivankov, Nina G. Bozhanova, Mikhail S. Baranov, Onuralp Soylemez, Natalya S. Bogatyreva, Peter K. Vlasov, Evgeny S. Egorov, Maria D. Logacheva, Alexey S. Kondrashov, Dmitry M. Chudakov, Ekaterina V. Putintseva, Ilgar Z. Mamedov, Dan S. Tawfik, Konstantin A. Lukyanov, and Fyodor A. Kondrashov. Local fitness landscape of the green fluorescent protein. *Nature*, 533(7603):397–401, May 2016. ISSN 1476-4687. doi: 10.1038/nature17995. URL <https://doi.org/10.1038/nature17995>.
- Tao Shen, Zhihang Hu, Siqi Sun, Di Liu, Felix Wong, Jiuming Wang, Jiayang Chen, Yixuan Wang, Liang Hong, Jin Xiao, Liangzhen Zheng, Tejas Krishnamoorthi, Irwin King, Sheng Wang, Peng Yin, James J. Collins, and Yu Li. Accurate RNA 3D structure prediction using a language model-based deep learning approach. *Nature Methods*, 21(12):2287–2298. ISSN 1548-7105. doi: 10.1038/s41592-024-02487-0. URL <https://doi.org/10.1038/s41592-024-02487-0>.
- Jaswinder Singha, Jack Hanson, Kuldip Paliwal, and Yaoqi Zhou. Rna secondary structure prediction using an ensemble of two-dimensional deep neural networks and transfer learning. *Nature communications*, 2019. doi: <https://doi.org/10.1038/s41467-019-13395-9>. URL <https://www.nature.com/articles/s41467-019-13395-9>.

- Valerie W. C. Soo, Jacob B. Swadling, Andre J. Faure, and Tobias Warnecke. Fitness landscape of a dynamic rna structure. *PLOS Genetics*, 17(2):1–21, 02 2021. doi: 10.1371/journal.pgen.1009353. URL <https://doi.org/10.1371/journal.pgen.1009353>.
- Robert C. Spitale and Danny Incarnato. Probing the dynamic rna structurome and its functions. *Nature Reviews Genetics*, 24(3):178–196, Mar 2023. ISSN 1471-0064. doi: 10.1038/s41576-022-00546-w. URL <https://doi.org/10.1038/s41576-022-00546-w>.
- Tyler N. Starr, Allison J. Greaney, Sarah K. Hilton, Daniel Ellis, Katharine H.D. Crawford, Adam S. Dingens, Mary Jane Navarro, John E. Bowen, M. Alejandra Tortorici, Alexandra C. Walls, Neil P. King, David Veessler, and Jesse D. Bloom. Deep mutational scanning of sars-cov-2 receptor binding domain reveals constraints on folding and ace2 binding. *Cell*, 182(5):1295–1310.e20, 2020. ISSN 0092-8674. doi: <https://doi.org/10.1016/j.cell.2020.08.012>. URL <https://www.sciencedirect.com/science/article/pii/S0092867420310035>.
- Shunsuke Sumi, Michiaki Hamada, and Hirohide Saito. Deep generative design of rna family sequences. *Nature Methods*, 21:435–443, 2024.
- Jacob M Tome, Abdullah Ozer, John M Pagano, Dan Gheba, Gary P Schroth, and John T Lis. Comprehensive analysis of rna-protein interactions by high-throughput sequencing–rna affinity profiling. *Nature methods*, 2014. doi: 10.1038/nmeth.2970. URL <https://www.nature.com/articles/nmeth.2970>.
- Kotaro Tsuboyama, Justas Dauparas, Jonathan Chen, Elodie Laine, Yasser Mohseni Behbahani, Jonathan J. Weinstein, Niall M. Mangan, Sergey Ovchinnikov, and Gabriel J. Rocklin. Mega-scale experimental analysis of protein folding stability in biology and design. *Nature*, 620(7973):434–444, Aug 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06328-6. URL <https://doi.org/10.1038/s41586-023-06328-6>.
- Fabio Urbina, Filippa Lentzos, Cédric Invernizzi, and Sean Ekins. Dual use of artificial-intelligence-powered drug discovery. *Nature Machine Intelligence*, 4(3):189–191, 2022.
- F. H. D. van Batenburg, Alexander P. Gulyaev, Cornelis W. A. Pleij, J. Ng, and J. Oliehoek. Pseudobase: a database with rna pseudoknots. *Nucleic acids research*, 28 1:201–4, 2000. URL <https://api.semanticscholar.org/CorpusID:27636094>.
- Rosario Vanella, Christoph Küng, Alexandre A. Schoepfer, Vanni Doffini, Jin Ren, and Michael A. Nash. Understanding activity-stability tradeoffs in biocatalysts by enzyme proximity sequencing. *Nature Communications*, 15(1):1807, Feb 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-45630-3. URL <https://doi.org/10.1038/s41467-024-45630-3>.
- Ning Wang, Jiang Bian, Yuchen Li, Xuhong Li, Shahid Mumtaz, Linghe Kong, and Haoyi Xiong. Multi-purpose rna language modelling with motif-aware pretraining and type-guided fine-tuning. *Nature Machine Intelligence*, 6(5):548–557, May 2024a. ISSN 2522-5839. doi: 10.1038/s42256-024-00836-4. URL <https://doi.org/10.1038/s42256-024-00836-4>.
- Ning Wang, Jiang Bian, Yuchen Li, Xuhong Li, Shahid Mumtaz, Linghe Kong, and Haoyi Xiong. Multi-purpose rna language modelling with motif-aware pretraining and type-guided fine-tuning. *Nature Machine Intelligence*, 6(5):548–557, 2024b. ISSN 2522-5839. doi: 10.1038/s42256-024-00836-4. URL <https://doi.org/10.1038/s42256-024-00836-4>.
- Wenkai Wang, Chenjie Feng, Renmin Han, Ziyi Wang, Lisha Ye, Zongyang Du, Hong Wei, Fa Zhang, Zhenling Peng, and Jianyi Yang. trrosettarna: Automated prediction of RNA 3D structure with transformer network. *Nature Communications*, 14(1):7266. ISSN 2041-1723. doi: 10.1038/s41467-023-42528-4. URL <https://doi.org/10.1038/s41467-023-42528-4>.
- Hannah K. Wayment-Steele, Wipapat Kladwang, Alexandra I. Strom, Jeehyung Lee, Adrien Treuille, Alexander J Becka, and Rhiju Das. Rna secondary structure packages evaluated and improved by high-throughput experiments. *bioRxiv*, 2020. URL <https://api.semanticscholar.org/CorpusID:219311051>.
- Hannah K. Wayment-Steele, Wipapat Kladwang, Alexandra I. Strom, Jeehyung Lee, Adrien Treuille, Alex Becka, Rhiju Das, and Eterna Participants. Rna secondary structure packages evaluated and improved by high-throughput experiments. *Nature Methods*, 19(10):1234–1242, Oct 2022. ISSN 1548-7105. doi: 10.1038/s41592-022-01605-0. URL <https://doi.org/10.1038/s41592-022-01605-0>.
- Ryan Weeks and Marc Ostermeier. Fitness and functional landscapes of the e. coli rna iii gene rnc. *Molecular Biology and Evolution*, 40(3):msad047, 02 2023. ISSN 1537-1719. doi: 10.1093/molbev/msad047. URL <https://doi.org/10.1093/molbev/msad047>.
- Caleb Weinreb, Adam J Riesselman, John B Ingraham, Torsten Gross, Chris Sander, and Debora S Marks. 3D RNA and functional interactions from evolutionary couplings. *Cell*, 165(4):963–975, May 2016.

Zhe Zhang, Peng Xiong, Tongchuan Zhang, Junfeng Wang, Jian Zhan, and Yaoqi Zhou. Accurate inference of the full base-pairing structure of rna by deep mutational scanning and covariation-induced deviation of activity. *Nucleic Acids Research*, 48(3):1451–1465, 12 2019. ISSN 0305-1048. doi: 10.1093/nar/gkz1192. URL <https://doi.org/10.1093/nar/gkz1192>.

Zhe Zhang, Xu Hong, Peng Xiong, Junfeng Wang, Yaoqi Zhou, and Jian Zhan. Minimal twister sister-like self-cleaving ribozymes in the human genome revealed by deep mutational scanning. *eLife*, 12:RP90254, dec 2024. ISSN 2050-084X. doi: 10.7554/eLife.90254. URL <https://doi.org/10.7554/eLife.90254>.

Maxim Zvyagin, Alexander Brace, Kyle Hippe, Yuntian Deng, Bin Zhang, Cindy Orozco Bohorquez, Austin Clyde, Bharat Kale, Danilo Perez-Rivera, Heng Ma, Carla M. Mann, Michael Irvin, Defne G. Ozgulbas, Natalia Vassilieva, James Gregory Pauloski, Logan Ward, Valerie Hayot-Sasson, Murali Emani, Sam Foreman, Zhen Xie, Diangen Lin, Maulik Shukla, Weili Nie, Josh Romero, Christian Dallago, Arash Vahdat, Chaowei Xiao, Thomas Gibbs, Ian Foster, James J. Davis, Michael E. Papka, Thomas Brettin, Rick Stevens, Anima Anandkumar, Venkatram Vishwanath, and Arvind Ramanathan. Genslms: Genome-scale language models reveal sars-cov-2 evolutionary dynamics. *The International Journal of High Performance Computing Applications*, 37(6):683–705, 2023a. doi: 10.1177/10943420231201154. URL <https://doi.org/10.1177/10943420231201154>.

Maxim Zvyagin, Alexander Brace, Kyle Hippe, Yuntian Deng, Bin Zhang, Cindy Orozco Bohorquez, Austin Clyde, Bharat Kale, Danilo Perez-Rivera, Heng Ma, Carla M. Mann, Michael Irvin, Defne G. Ozgulbas, Natalia Vassilieva, James Gregory Pauloski, Logan Ward, Valerie Hayot-Sasson, Murali Emani, Sam Foreman, Zhen Xie, Diangen Lin, Maulik Shukla, Weili Nie, Josh Romero, Christian Dallago, Arash Vahdat, Chaowei Xiao, Thomas Gibbs, Ian Foster, James J. Davis, Michael E. Papka, Thomas Brettin, Rick Stevens, Anima Anandkumar, Venkatram Vishwanath, and Arvind Ramanathan. Genslms: Genome-scale language models reveal sars-cov-2 evolutionary dynamics. *The International Journal of High Performance Computing Applications*, 37(6):683–705, 2023b. doi: 10.1177/10943420231201154.

APPENDIX

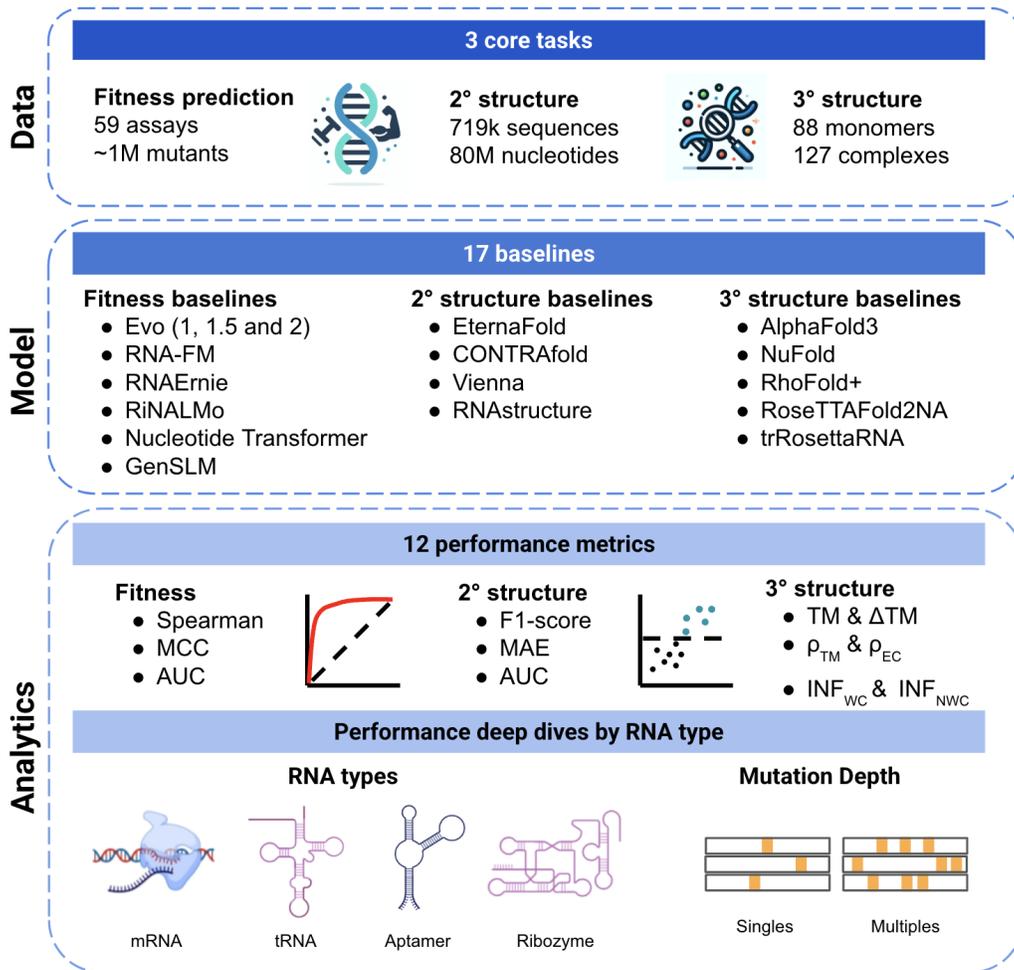


Figure 1: **RNAGym benchmarks.** RNAGym is a comprehensive RNA analysis framework designed specifically for fitness and structure prediction tasks. It evaluates the performance of diverse baselines across these tasks, and offers in-depth assessments by RNA type and mutation depth.

A RELATED WORK AND BACKGROUND

A.1 PRIOR RNA BENCHMARKS

Existing RNA benchmarks for noncoding variant effect prediction have been limited and fragmented, primarily focusing on testing individual models rather than serving as comprehensive benchmarking platforms. For instance, the RfamGen model was evaluated using five assays, including datasets on ribozymes and tRNAs (Sumi et al., 2024). Similarly, the Evo model was assessed using nine assays, incorporating ncRNA mutational scanning datasets (Nguyen et al., 2024). Both studies relied on overlapping but distinct datasets to evaluate their models, making it difficult to compare performance metrics between studies directly. These small benchmark sets restrict the ability to generalize findings and were primarily used to test the respective models’ performance, rather than providing a broad and standardized benchmarking framework spanning the diversity of RNA types.

This limited scope stands in stark contrast to the field of protein research, where platforms like ProteinGym have been established to offer extensive and standardized benchmarking datasets (Notin et al., 2023). RNAGym addresses this gap by introducing a comprehensive benchmarking platform for RNA variant effect prediction that offers more than four times the number of assays compared to previous efforts, across a broader array of RNA classes. Previous RNA benchmarks include BEACON Ren et al. (2024), which covers 13 diverse

tasks spanning secondary structure, functional genomics, and RNA engineering, and RNABench Runge et al. (2024), which provides six benchmarks centered on RNA structure prediction and design. In contrast, RNAGym specifically targets RNA fitness (variant effect) prediction and RNA folding tasks across both secondary and tertiary structures.

With respect to 3D RNA structure prediction, several competitive benchmarks have been developed, including the Critical Assessment of Structure Prediction (CASP) and RNA-Puzzles. CASP15 (Kryshtafovych et al., 2023), the latest iteration of CASP, introduced a dedicated category for RNA structure prediction, reflecting the growing recognition of RNA’s importance and the need for accurate computational models. RNA-Puzzles (Cruz et al., 2012), on the other hand, is a community-driven initiative that presents real-world challenges to participants, who submit their models to be evaluated against experimentally determined RNA structures. Notably, only a few RNA molecules are evaluated at CASP and through RNA-Puzzles, limiting the ability of these high quality datasets to act as benchmarking standards.

Train-test splits are commonly used to evaluate 2D RNA structure prediction models (Penić et al., 2024; Chen et al., 2022). Intra-family splits involve training models on RNA sequences from the same family, with sequences from these families appearing in both the training and test sets (Singha et al., 2019). This approach tests a model’s ability to learn and predict structures within known families. In contrast, inter-family splits ensure that sequences from the same RNA family in the test set are excluded from the training set (Penić et al., 2024). This method assesses whether a model can generalize to entirely new RNA families that were not included in the training data. While existing benchmarks offer valuable insights, they often lack the scale and diversity to comprehensively evaluate model performance across various RNA types and structures. Zero-shot benchmarks for structure prediction models are crucial, as they avoid biases inherent in supervised approaches and potential overfitting, thus providing a more robust assessment of true generalization capabilities.

A.2 BACKGROUND: THE DIVERSITY OF RNA MOLECULES

RNA molecules exhibit a remarkable range of structures and functions, highlighting their essential role in both the fundamental processes of biology and their growing utility in medical and biotechnological applications. From the synthesis and regulation of proteins to the catalysis of key biochemical reactions, RNA types such as mRNA, tRNA, ribozymes, and aptamers demonstrate the diversity of RNA molecular diversity and complexity. **mRNAs (Messenger RNAs)** act as the intermediary transcript that carry genetic information from DNA to the ribosome, where they serve as a template for protein synthesis. The sequence of mRNA dictates the amino acid sequence in a protein, thereby directly influencing gene expression and regulatory mechanisms. While alterations in splicing impact how exons are joined together to form the mature mRNA transcript, missense changes directly alter the codon sequence itself, resulting in the incorporation of a different amino acid in the final protein product. **tRNAs (Transfer RNAs)** play a crucial role in translation, the process of protein synthesis in the ribosome. Each tRNA molecule transports a specific amino acid to the ribosome; its anticodon loop pairs with the corresponding codon in the mRNA, ensuring that the correct amino acid is incorporated into the growing protein chain. **Ribozymes** are catalytic RNA molecules that perform specific biochemical reactions akin to protein enzymes. These include critical activities such as RNA splicing during gene expression, where ribozymes help in the excision of introns from a pre-mRNA. **Aptamers** are short, single-stranded RNA molecules designed to bind with high specificity and affinity to certain targets, including proteins, small molecules, and various cellular components. Their high specificity and adaptability make aptamers highly valuable for therapeutic uses, as well as in diagnostic and biosensing applications.

B DATASET COLLECTION

Prioritization of studies for expert review Coding mRNA datasets were sourced from ProteinGym when nucleotide-level deep mutational scanning information was available (Notin et al., 2023).

The selection process for prioritizing noncoding RNA studies for expert review was structured as follows. We initiated with a targeted PubMed search, utilizing specific queries to ensure a comprehensive capture of relevant literature:

1. **Deep Mutational Scan:** (deep[All Fields] OR comprehensive[All Fields]) AND (mutational[All Fields] OR muta*) AND (scan OR scans OR scanning) AND RNA
2. **Saturation Mutagenesis:** (saturat* muta*) AND RNA
3. **MAVE:** (“Multiplex* assay” AND “variant*”) AND RNA
4. **MPRA:** (“MPRA” OR “Massively parallel reporter assay*”) AND RNA NOT protein
5. **Other:** (“Fitness Landscape” AND muta*) AND RNA

These initial searches proved to be either overly restrictive or too broad, which complicated the manual screening process. Ultimately, this approach resulted in the identification of only 20 primary articles. To enrich our review,

an additional 10 articles were identified by scraping references from other pertinent studies, including those cited in previous research such as (Sumi et al., 2024; Nguyen et al., 2024), and RNA-related datasets from studies like ProteinGym (Notin et al., 2023).

Hypothesizing that our initial search strategy may have missed relevant studies, we conducted a comprehensive PubMed search using the following query:

Broad Search: (deep OR comprehensive OR MPRA OR multiplex assay OR massively parallel OR landscape OR saturation) AND (muta* OR variant OR variants) AND (scan OR scans OR screen OR landscape OR assay) AND (RNA OR ribozyme* OR microRNA OR miRNA OR siRNA OR snoRNA OR tRNA OR lncRNA OR (RNA AND aptamer) OR circRNA)

B.1 LITERATURE PRE-SCREENING WITH LLM

The prior search yielded an overwhelming 11,635 results. To efficiently handle this volume, we utilized a large language model (LLM), specifically GPT4-0125-preview, for secondary screening. We adapted a recent prompting approach designed for systematic review screening (Cao et al., 2024). The LLM was instructed with clear study objectives and specific inclusion/exclusion criteria, effectively narrowing down the pool to fewer than 500 articles, thereby making manual curation manageable. To enhance the sensitivity of this process, the LLM's prompt was refined using an initial set of 30 positively identified articles as a control group. This novel use of LLMs for data extraction markedly improved our capacity to pinpoint relevant studies. Consequently, we were able to incorporate an additional 22 studies into our initial screen, resulting in a total of 52 studies ready for manual expert review.

We used the following prompt to pre-screen relevant studies during our extensive literature search:

"The goal of this study is to create a benchmark that contains RNA deep mutational scanning or fitness landscape datasets. We are generating these datasets to benchmark RNA fitness prediction algorithms, and need our datasets we evaluate to have information on RNA mutants/variants and their relative 'fitness'.

The following is an excerpt of two sets of criteria. A study is considered included if it meets all the inclusion criteria. If a study meets any of the exclusion criteria, it should be excluded. Here are the two sets of criteria:

Inclusion Criteria (all must be fulfilled): 1. Studies involve RNA. We are also interested in RNA subclasses such as Ribozyme, lncRNA, tRNA, rRNA, microRNA (miRNA), Aptamer, Riboswitch, mRNA 2. Studies report on fitness prediction. Other terms for fitness prediction can include deep mutational scans, comprehensive multiplex assays, or comprehensive fitness landscapes, among others 3. Studies with greater than 100 experimental measurements 4. Studies that report on mutant fitness through reporter assays, bulk RNA-sequencing, single-cell RNA sequencing assay, fluorescence in-situ hybridization (FISH) assay, flow cytometry assay, imaging mass cytometry assay, evolution of ligands by exponential enrichment assay, single cell imaging, multiplexed fluorescent antibody imaging, binding assays, cell proliferation assay, splicing assays, survival assessment assay selection types, or similar. 5. Studies that report on enzymatic activity, binding affinity, stability, fluorescence, proliferation selection assays, or similar assays. 6. The study must be primary research and generate a novel dataset

Exclusion Criteria (if any met then exclude): 1. Studies only reporting on protein mutational scans, with no relevance or mention of RNA being mutated 2. Studies that do not focus on fitness quantification 3. Review articles (systematic reviews, case reports, case series, etc.) or other non-primary research sources.

Instructions

We now assess whether the paper should be included in the systematic review by evaluating it against each and every predefined inclusion and exclusion criterion. First, we will reflect on how we will decide whether a paper should be included or excluded. Then, we will think step by step for each criteria, giving reasons for why they are met or not met. Studies that may not fully align with the primary focus of our inclusion criteria but provide data or insights potentially relevant to our review deserve thoughtful consideration. Given the nature of abstracts as concise summaries of comprehensive research, some degree of interpretation is necessary. Our aim should be to inclusively screen abstracts, ensuring broad coverage of pertinent studies while filtering out those that are clearly irrelevant. We will conclude by outputting (on the very last line) 'XXX' if the paper warrants exclusion, or 'YYY' if inclusion is advised or uncertainty persists. We must output either 'XXX' or 'YYY'.

Title and Abstract in investigation:

Title: #Insert title of study#

Abstract: #Insert abstract of paper#"

Expert review process The process for accepting a paper involved several steps to ensure the quality and relevance of the data. First, we checked whether the data was openly available and could be integrated into our benchmark. If data was not accessible, study authors were contacted.

Next, we used the following inclusion and exclusion criteria during our through expert review process:

Inclusion Criteria

- Assay must focus on RNA
- Assay must have at least 100 experimental variants tested, with a sufficiently wide dynamic range
- Assay must be relevant to fitness prediction, and report on mutant fitness
- Assay must only focus on substitutions, not insertions or deletions

Exclusion Criteria

- Assays focusing on DNA or Proteins
- Assays that are not primary research
- Assays with mutants of varying lengths

C DATASETS DETAILS

RNA Type	Description	# Assays	# Singles	# Multiples
Aptamer	Target binding ability	2	0.4k	2.6k
Ribozyme	Cleavage and splicing	26	4.1k	755k
Transfer RNA (tRNA)	Stability and growth	3	0.4k	95k
Messenger RNA (mRNA)	Splicing ability	2	0.3k	5.4k
Messenger RNA (mRNA)	Coding mRNA fitness	26	17k	118k
Total		59	22k	976k

Table A1: **RNAGym fitness benchmark summary.** RNAGym includes a large collection of DMS assays about diverse RNA types. The table reports the number of assays and number of single and multiple mutants per RNA type.

Dataset Type	Description	# Sequences	# Filtered Sequences
2A3	Standard RNA	2,136k	471k
DMS	Standard RNA	2,148k	411k
Other modifier	Standard RNA	70k	47k
Ligand binding	RNA with ligands	292	292
In vivo	RNA is folded in vivo	30	29
Cotranscriptional	RNA is folded cotranscriptionally	3k	3k
Degradation	Accelerated RNA degradation	18k	16k
Total		4,375k	947k

Table A2: **2° structure dataset.** The secondary structure dataset consists of RNA reactivity data compiled from RMDB. Standard RNA refers to the experiment done in vitro.

C.1 FITNESS ASSAYS

References An exhaustive list of the publications from which the assays included our fitness benchmark originated from is provided in Table C.1 and Table C.1. Our final processed datasets all have a consistent format with the same 3 fields across: "Mutant" (mutation triplets), "Sequence" (mutated sequence), and "DMS score" (experimental measurement). We also corrected the *directionality* of each measured experimental phenotype, to ensure that higher DMS scores always translate to higher fitness across assays.

Licenses All fitness assays were licensed under CC-BY 4.0 (<https://creativecommons.org/licenses/by/4.0/>), or the ACS AuthorChoice Usage Agreement (https://pubs.acs.org/page/policy/authorchoice_termsofuse.html).

Dataset	# Chains	# Rfams	# Rfam signatures	Avg. resolution	Avg. length
Monomers	88	37	60	2.94Å	130 nt
Complexes	127	32	74	2.98Å	52 nt
RNAgym (total)	215	66	127	2.96Å	84 nt

Table A3: **RNAgym 3° structure benchmark summary.** RNAgym features diverse 3D RNA monomers and complexes. Rfams := unique top-matching RNA families (E-value ≤ 1). Rfam signatures capture structural diversity of unique top-matching Rfam footprints (E-value ≤ 10). Monomers := chains with <33% residues bound to another polymer (within 4.5Å). Complexes := all others.

Cross-validation splits For users interested in supervised RNA fitness prediction, we provide two types of cross-validation splits:

- **Random:** a random 80%-20% train-test split;
- **Minimum similarity:** a 80%-20% split in which we minimized the sequence similarity between training and validation RNA sequences.

C.2 STRUCTURE PREDICTION DATASET

We constructed the data for our structure prediction challenge using the data from the Ribonanza Challenge hosted on Kaggle (<https://www.kaggle.com/competitions/stanford-ribonanza-rna-folding/data>) and RMDB (Cordero et al., 2012). The data contains RNA chemical reactivity at various conditions (Table A2). For ease of use we created a dataset of only in vitro chemical mapping data involving 2A3 and DMS. While we do not include the data such as those collected with ligands in our clean train-test data, it is available. To ensure the integrity of our dataset and prevent data leakage, we clustered the sequences with 40% sequence identity and created a 20-80 train-test split. The final dataset consists of 901k reactivity profiles for 583k unique sequences and 100M nucleotides. This extensive dataset underpins the robustness and comprehensive nature of our RNA structure prediction challenge. The original data from the challenge is made available under a CC-BY 4.0 license.

D BASELINES

D.1 RNA FITNESS PREDICTION MODELS

Our fitness prediction benchmarks currently include the following 5 baselines:

- **RiNALMo** (Penić et al., 2024) is a 650 million parameter RNA language model trained on 36 million non-coding RNA sequences, achieving high performance in RNA structural and functional prediction tasks. Variants were scored using the masked marginal scoring strategy from the ESM sequence modeling framework.
- **Evo** (Nguyen et al., 2024) is a 7 billion parameter model trained on 2.7M prokaryotic and phage genomes to generate DNA sequences using a context length of 8k (further extended to 131k tokens). It is based on StripedHyena, a deep signal processing architecture designed to improve efficiency and quality over the prevailing Transformer architecture. We use the 8k version of the model, referred as Evo 1. The recently released Evo 1.5 builds upon Evo 1 by increasing the pretraining data from 300 billion to 450 billion tokens Merchant et al. (2024). We also include the EVO2 (7B) model Brix et al. (2025), EVO2 is trained on over 9.3T tokens at a single nucleotide resolution.
- **RNA-FM** (Chen et al., 2022) is a RNA foundation model based on the BERT language model architecture. It was trained on 23 million unlabeled non-coding RNA sequences from over 800,000 species, collected from RNA-central. Variants were scored using the wild-type marginal scoring strategy from the ESM sequence modeling framework, which RNA-FM builds upon.
- **GenSLM** (Zvyagin et al., 2023b) includes an autoregressive RNA language model with codon-level tokenization, along with stable diffusion, to model genome-scale interactions and predict SARS-CoV-2 evolution. It was trained on 110 million prokaryotic coding sequences from BV-BRC. Variant sequence likelihoods were scored using the 2.5 billion parameter language model.
- **RNAErnie** (Wang et al., 2024b) is a 12-layer transformer-based RNA language model pre-trained on 23 million ncRNA sequences using masked language modeling. With 105 million trainable parameters,

Title	Year	RNA type	# Assays	Reference	License
Saturation mutagenesis reveals manifold determinants of exon definition	2017	mRNA - splicing	1	(Ke et al., 2017)	CC BY 4.0
Fitness landscape of a dynamic RNA structure	2021	ribozyme	1	(Soo et al., 2021)	CC BY 4.0
Comprehensive sequence-to-function mapping of cofactor-dependent RNA catalysis in the glmS ribozyme	2020	ribozyme	1	(Andreasson et al., 2020)	CC BY 4.0
High-throughput assay and engineering of self-cleaving ribozymes	2015	ribozyme	3	(Kobori et al., 2015)	CC BY 4.0
Identification of the determinants of tRNA function and susceptibility to rapid tRNA decay by high-throughput in vivo analysis	2014	tRNA	1	(Guy et al., 2014)	CC BY 4.0
Deep sequencing analysis of aptazyme variants based on a pistol ribozyme	2018	ribozyme	1	(Kobori et al., 2017)	ACS Author-Choice Usage Agreement
Accurate inference of the full base-pairing structure of RNA by deep mutational scanning and covariation-induced deviation of activity.	2020	ribozyme	1	(Zhang et al., 2019)	CC BY 4.0
Rapid Construction of Empirical RNA Fitness Landscapes	2010	ribozyme	1	(Pitt & Ferré-D'Amaré, 2010)	CC-BY 4.0
Dynamic RNA Fitness Landscapes of a Group I Ribozyme during Changes to the Experimental Environment	2022	ribozyme	1	(Peri et al., 2022)	CC-BY 4.0
RNA sequence to structure analysis from comprehensive pairwise mutagenesis of multiple self-cleaving ribozymes	2023	ribozyme	5	(Roberts et al., 2023)	CC-BY 4.0
Pairwise and higher-order genetic interactions during the evolution of a tRNA	2018	tRNA	1	(Domingo et al., 2018)	CC-BY 4.0
Emergent properties as by-products of prebiotic evolution of aminoacylation ribozymes	2022	ribozyme	5	(Janzen et al., 2022)	CC-BY 4.0
Predicting higher-order mutational effects in an RNA enzyme by machine learning of high-throughput experimental data	2022	ribozyme	1	(Beck et al., 2022)	CC-BY 4.0
The fitness landscape of a tRNA gene	2016	tRNA	1	(Li et al., 2016)	CC-BY 4.0
Comprehensive analysis of RNA-protein interactions by high-throughput sequencing-RNA affinity profiling	2014	aptamer	2	(Tome et al., 2014)	CC-BY 4.0
The complete local genotype-phenotype landscape for the alternative splicing of a human exon	2016	mRNA - splicing	1	(Julien et al., 2016)	CC-BY 4.0
Minimal twister sister-like self-cleaving ribozymes in the human genome revealed by deep mutational scanning.	2024	ribozyme	3	(Zhang et al., 2024)	CC-BY 4.0
Cryo-EM structure and functional landscape of an RNA polymerase ribozyme	2024	ribozyme	2	(McRae et al., 2024)	CC-BY 4.0
High-Throughput Mutational Analysis of a Twister Ribozyme	2016	ribozyme	1	(Kobori & Yokobayashi, 2016)	CC-BY 4.0

Table A4: **RNAGym fitness prediction data.** We developed our noncoding and splicing fitness prediction benchmark by curating and processing 33 assays from 19 publications.

it achieves high performance in RNA sequence classification, RNA–RNA interaction, and RNA secondary structure prediction.

- **Nucleotide Transformer** (Dalla-Torre et al., 2023) is a 2.5B parameter model trained on 850 species. The 2.5b-multi-species version was used.

A Comprehensive, High-Resolution Map of a Gene’s Fitness Landscape	2014	mRNA	1	(Firnberg et al., 2014)	CC-BY 4.0
Capturing the mutational landscape of the beta-lactamase TEM-1	2013	mRNA	1	(Jacquier et al., 2013)	CC-BY 4.0
Protein model discrimination using mutational sensitivity derived from deep sequencing	2012	mRNA	1	(Adkar et al., 2012)	CC-BY 4.0
RNA structural determinants of optimal codons revealed by MAGE-Seq	2016	mRNA	1	(Kelsic et al., 2016)	CC-BY 4.0
Fitness and functional landscapes of the E. coli RNase III gene rnc	2023	mRNA	1	(Weeks & Ostermeier, 2023)	CC-BY 4.0
Protein design using structure-based residue preferences	2023	mRNA	1	(Ding et al., 2024)	CC-BY 4.0
Deep mutational scan of a drug efflux pump reveals its structure–function landscape	2023	mRNA	1	(Meier et al., 2023)	CC-BY 4.0
Protein-protein interactions in the Mla lipid transport system probed by computational structure prediction and deep mutational scanning	2023	mRNA	1	(MacRae et al., 2023)	CC-BY 4.0
Systematically Scrutinizing the Impact of Substitution Sites on Thermostability and Detergent Tolerance for Bacillus subtilis Lipase A	2020	mRNA	1	(Nutschel et al., 2020)	CC-BY 4.0
Local fitness landscape of the green fluorescent protein	2016	mRNA	1	(Sarkisyan et al., 2016)	CC-BY 4.0
The spatial architecture of protein function and adaptation	2012	mRNA	1	(McLaughlin Jr et al., 2012)	CC-BY 4.0
Understanding Activity-Stability Tradeoffs in Biocatalysts by Enzyme Proximity Sequencing	2023	mRNA	1	(Vanella et al., 2024)	CC-BY 4.0
A Systematic p53 Mutation Library Links Differential Functional Impact to Cancer Mutation Pattern and Evolutionary Conservation	2018	mRNA	1	(Kotler et al., 2018)	CC-BY 4.0
Accurate classification of BRCA1 variants with saturation genome editing	2018	mRNA	1	(Findlay et al., 2018)	CC-BY 4.0
Deep Mutational Scanning of SARS-CoV-2 Receptor Binding Domain Reveals Constraints on Folding and ACE2 Binding	2020	mRNA	1	(Starr et al., 2020)	CC-BY 4.0
Deep mutational scanning of hemagglutinin helps predict evolutionary fates of human H3N2 influenza variants	2018	mRNA	1	(Lee et al., 2018)	CC-BY 4.0
Mega-scale experimental analysis of protein folding stability in biology and design	2023	mRNA	10	(Tsuboyama et al., 2023)	CC-BY 4.0

Table A5: **RNAGym coding fitness prediction data.** We developed our coding fitness prediction benchmark by curating and processing 26 assays from 17 publications.

D.2 RNA 2° STRUCTURE PREDICTION MODELS

Our 2° structure prediction benchmarks currently include the following 4 baselines:

- **EternaFold** (Wayment-Steele et al., 2020) is built on the principles derived from the Eterna massive open online game, where players design RNA sequences that fold into target shapes. This model

incorporates crowd-sourced insights from thousands of players to refine its algorithms, significantly enhancing its ability to predict RNA structures under varied environmental conditions and complexities.

- **CONTRAFold** (Do et al., 2006) is a machine learning-based RNA secondary structure prediction model that utilizes conditional log-linear models for structure inference.
- **Vienna** (Gruber et al., 2008) is one of the most widely used RNA secondary structure prediction tools. It employs dynamic programming algorithms based on thermodynamic models to accurately predict RNA secondary structures, including handling pseudoknotted structures as extensions.
- **RNAstructure** (Reuter & Mathews, 2010) is a model designed for the prediction and analysis of RNA secondary structures. It is known for its dual ability to use either thermodynamic or machine learning-based methods to predict RNA folding patterns.

D.3 RNA 3° STRUCTURE PREDICTION MODELS

Our tertiary structure prediction benchmarks currently include the following 5 baselines:

- **AlphaFold3** (Abramson et al.) extends the AlphaFold framework to RNA, predicting full 3D coordinates using a diffusion-based architecture trained on protein and nucleic acid data. It captures interchain interactions, supporting both RNA monomers and complexes.
- **NuFold** (Kagaya et al.) is an end-to-end deep learning model dedicated to RNA, adapting iterative refinement techniques from protein structure prediction. It incorporates base-centered representations to accurately capture local RNA geometry.
- **RhoFold+** (Shen et al.) offers an automated, end-to-end platform for RNA 3D structure prediction that integrates the RNA-FM language model and iterative geometry-aware refinement.
- **RoseTTAFold2NA** (Baek et al.) extends the RoseTTAFold framework to handle RNA and protein–RNA complexes. By jointly refining sequence features, residue-pair distances, and cartesian coordinates, it models RNA folds and interfaces in a single pipeline.
- **trRosettaRNA** (Wang et al.) leverages a transformer-based network to predict intra-nucleotide distances and torsion angles, followed by refinement using Rosetta’s energy minimization. This hybrid approach combines deep-learning potentials and physics-based energy terms to generate accurate RNA structures.

E DETAILED EXPERIMENTAL RESULTS

E.1 DETAILED FITNESS PREDICTION PERFORMANCE

Rank	Model name	Spearman	AUC	MCC
1	Evo 2	0.286	0.641	0.220
2	RNA-FM	0.218	0.606	0.155
3	RNAErnie	0.188	0.598	0.153
4	Evo 1.5	0.182	0.586	0.134
5	Nucl. Transformer	0.166	0.579	0.117
6	RiNALMo	0.155	0.582	0.122
7	Evo 1	0.13	0.564	0.096
8	GenSLM	0.122	0.558	0.083

Table A6: **RNAGym - Fitness prediction overall benchmark.** Average of Spearman’s rank correlation, AUC and MCC between model scores and experimental measurements on the full RNAGym fitness prediction benchmark.

Performance by mutation type. The fitness prediction results segmented by mutation type (Table A8) show that Evo models consistently outperform others across both single and multiple mutations. Evo 2.0 achieves the highest correlations for both single mutations (0.245) and multiple mutations (0.280). While these results establish Evo models as the current leaders in mutation effect prediction, the relatively low correlation values indicate substantial room for improvement in capturing RNA sequence-function relationships. Notably, all models, with the exception of EVO 2.0, show somewhat stronger performance on single mutations compared to multiple mutations, highlighting the increased challenge of predicting fitness effects for more complex genetic variations. This performance gap between single and multiple mutations points to opportunities for improving model architectures and training approaches to better handle combinatorial effects of mutations.

Future Directions. To advance the field of RNA fitness prediction, several promising avenues of investigation emerge. First, developing models specifically trained on diverse RNA fitness landscapes could potentially improve performance by more closely aligning the training data with the prediction task. Additionally, incorporating RNA secondary structure predictions or experimental structure data into fitness prediction models may enhance their accuracy by capturing the important relationship between RNA structure and function. Comparing zero-shot performance with fine-tuned models could provide valuable insights into the generalizability of learned RNA features, potentially guiding future model development strategies. Lastly, exploring new architectural elements or pre-training objectives that better capture RNA-specific properties might lead to more robust and accurate predictions.

We report the assay-level Spearman performance, across all assays in the RNAGym fitness prediction benchmark in Fig E.1, aggregated performance by RNA type in Table A7, and by mutational depth in Table A8.

Rank	Model name	mRNA-splic.	tRNA	Aptamer	Ribozyme	mRNA-cod.	All
1	Evo 2	0.431	0.387	0.119	0.172	0.323	0.286
2	RNA FM	0.103	0.464	0.190	0.201	0.131	0.218
3	RNAErnie	0.230	0.416	0.037	0.142	0.116	0.188
4	Evo 1.5	0.131	0.385	0.044	0.169	0.179	0.182
5	NT	0.121	0.317	0.131	0.147	0.099	0.166
6	RiNALMo	0.348	0.260	0.026	0.072	0.071	0.155
7	Evo 1	0.134	0.155	0.035	0.137	0.190	0.13
8	GenSLM	0.173	0.093	0.126	0.135	0.083	0.122
-	All	0.209	0.310	0.090	0.147	0.149	-

Table A7: **RNAGym - Fitness prediction by RNA type.** Average of Spearman’s rank correlation between model scores and experimental measurements by RNA type and overall.

Rank	Model name	Singles	Multiples
1	Evo 2	0.245	0.280
2	RNA FM	0.163	0.162
3	RNAErnie	0.199	0.196
4	Evo 1.5	0.216	0.163
5	Nucl. Transformer	0.129	0.163
6	RiNALMo	0.140	0.153
7	Evo 1	0.174	0.128
8	GenSLM	0.145	0.117

Table A8: **RNAGym - Fitness prediction by mutation type.** Average of Spearman’s rank correlation between model scores and experimental measurements by mutation type.

We also report the statistical significance for the relative Spearman performance by RNA type in Table A9. We follow the same methodology as in ProteinGym (Notin et al., 2023) and assess statistical significance by computing the non-parametric bootstrap standard error of the difference between the Spearman performance of a given model and that of the best overall model.

E.2 DETAILED SECONDARY STRUCTURE PREDICTION RESULTS

We mapped RNA sequences in our evaluation set to pseudoknot annotations from PseudoBase (van Batenburg et al., 2000) (358 test sequences mapped), and report the corresponding global F1 score and crossed pair F1 score (Table A12). Out of our various structure prediction baselines, only RNAstructure is able to score pseudoknots.

Future Directions. Our evaluation reveals that current zero-shot RNA structure prediction methods perform remarkably similarly, with only small differences separating their effectiveness. To overcome this limit, significant advances will likely come from supervised learning approaches that can leverage experimental data more effectively. To facilitate this progress, we provide a carefully curated nonredundant training dataset that researchers can use to develop and benchmark supervised methods without concerns of sequence redundancy between training and test sets. By establishing this resource, we aim to enable fair evaluation of model generalization capabilities and accelerate the development of next-generation RNA structure prediction methods that can significantly outperform current approaches while maintaining robustness across diverse RNA families.



Figure 2: **RNAgym fitness prediction benchmark - Detailed performance for non-coding and splicing assays.** Spearman’s rank correlation between model predictions and experimental values for non coding and splicing assays in the RNAgym fitness prediction benchmark.

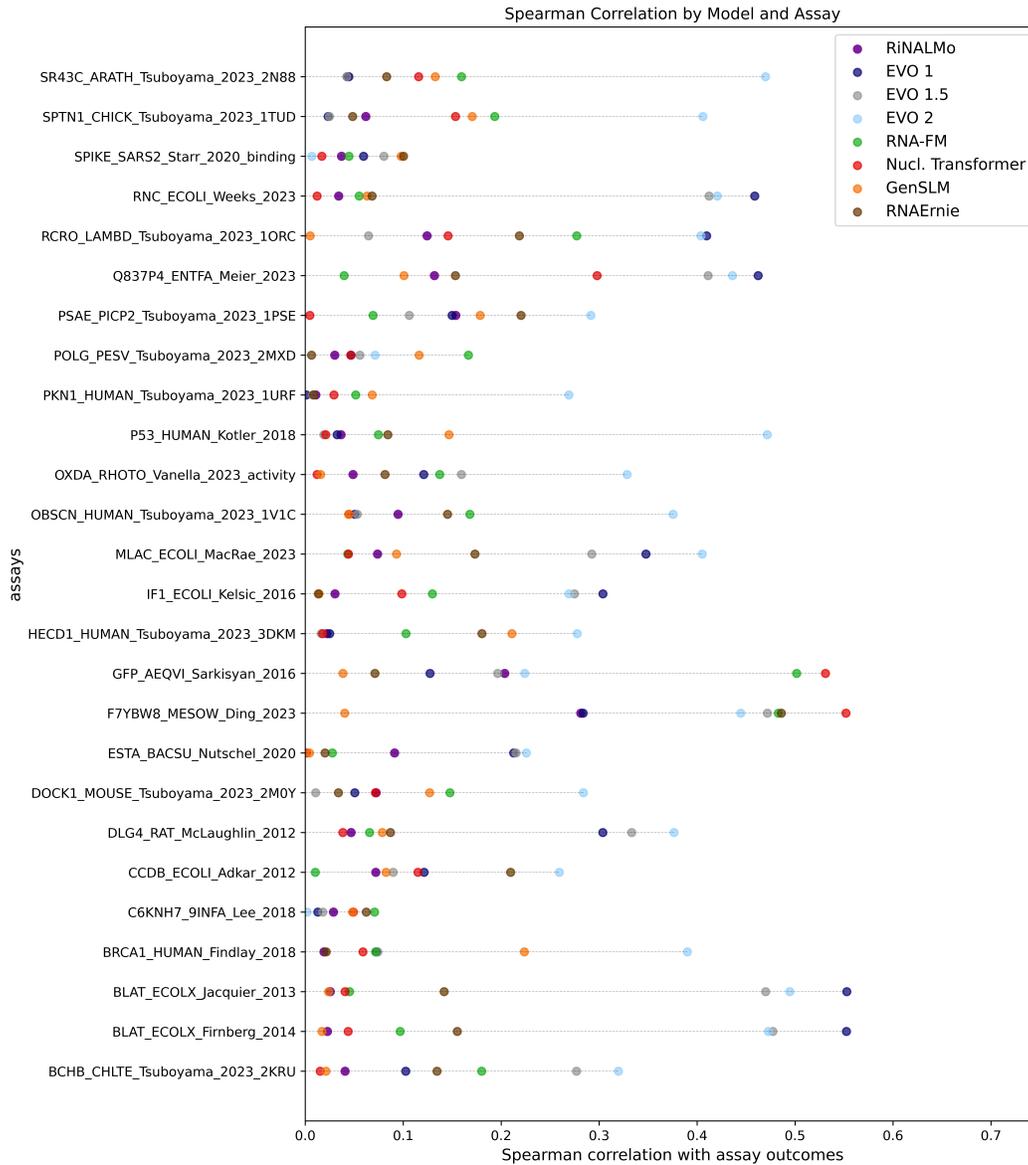


Figure 3: **RNAGym fitness prediction benchmark - Detailed performance for coding assays.** Spearman’s rank correlation between model predictions and experimental values for coding assays in the RNAGym fitness prediction benchmark.

Rank	Model	mRNA-splic.		mRNA-cod.		tRNA		Aptamer		Ribozyme		All	
		Diff	SE	Diff	SE	Diff	SE	Diff	SE	Diff	SE	Diff	SE
1	Evo2	0.000	0.000	0.000	0.000	-0.077	0.050	-0.071	0.013	-0.029	0.038	0.000	0.000
2	RNA-FM	-0.328	0.126	-0.192	0.034	0.000	0.000	0.000	0.000	0.000	0.000	-0.069	0.028
3	RNAErnie	-0.201	0.062	-0.207	0.027	-0.048	0.094	-0.153	0.113	-0.059	0.018	-0.098	0.023
4	Evo1.5	-0.300	0.125	-0.144	0.031	-0.079	0.050	-0.146	0.089	-0.032	0.034	-0.105	0.021
5	NT	-0.311	0.014	-0.224	0.036	-0.147	0.043	-0.044	0.022	-0.054	0.027	-0.121	0.025
6	RiNALMo	-0.083	0.083	-0.252	0.028	-0.204	0.136	-0.164	0.097	-0.129	0.029	-0.131	0.020
7	Evo1	-0.298	0.059	-0.133	0.031	-0.308	0.014	-0.156	0.085	-0.064	0.036	-0.156	0.021
8	GenSLM	-0.258	0.031	-0.240	0.029	-0.370	0.058	-0.065	0.092	-0.065	0.029	-0.164	0.023

Table A9: **Fitness prediction by RNA type - Difference in Spearman to best score by category** Difference in average of Spearman’s rank correlation between model scores and experimental measurements to the best model by category, by RNA type and overall. The standard error reported corresponds to the non-parametric bootstrap standard error of the difference between the Spearman performance of a given model and that of the best overall model for a given category, computed over 10k bootstrap samples from the set of assays in the RNAGym fitness benchmark.

Dataset Type	Train	Test	Total
2A3	97.6k (11.1M)	381.3k (43.1M)	479.0k (54.2M)
DMS	84.6k (9.5M)	337.6k (37.4M)	422.2k (46.9M)
Total	182.3k (20.7M)	718.9k (80.5M)	901.1k (101.2M)

Table A10: **Final 2° structure dataset.** The secondary structure chemical mapping dataset was further cleaned to only consist of 2A3 and DMS reactivity profiles for RNA in standard condition. The numbers represent the count of chemical mapping profiles while the number in parentheses indicate the count of nucleotides.

E.3 DETAILED TERTIARY STRUCTURE PREDICTION RESULTS

Future directions. These results indicate promising yet still limited zero-shot accuracy for RNA tertiary structure prediction. Although TM and Δ TM scores reveal room for improvement, different models demonstrate unique advantages: NuFold excels at template modeling and evolutionary coupling extraction (Δ TM, ρ_{TM} , ρ_{EC}), while AlphaFold3 leads on local basepair network fidelity (INF_{WC} , INF_{NWC}) and complexes. The poor performance of all models on non-Watson-Crick interactions highlights a pressing need for new approaches. Harnessing these complementary strengths in a unified model will be instrumental to achieving highly accurate 3D RNA predictions. The centralized RNAGym benchmark accelerates progress towards this goal by enabling the community to more transparently compare the strengths and weaknesses of new methods as they emerge.

F LIMITATIONS

Experimentally assaying RNA fitness, while resource-intensive, provides critical insights that help advance our understanding of RNA function. However, such experimental assays may not always accurately mimic the cellular environment, which can lead to variations between observed in vitro results and actual in vivo functionality. Chemical mapping experiments using dimethyl sulfate (DMS) offer valuable data on RNA secondary structures by identifying accessible adenine and cytosine bases that interact with DMS. This technique, although powerful for revealing the in-vivo-like structure of RNA in a relatively high-throughput manner, has its limitations. DMS mapping can be affected by incomplete coverage, as it primarily marks only two of the four nucleotide types. Additionally, the resolution of DMS mapping might not always distinguish closely spaced structural features, potentially obscuring important details about RNA folding and interaction sites. The accuracy of predictions from DMS data also heavily depends on the computational tools used to interpret the chemical reactivity patterns, necessitating ongoing improvements in both experimental and analytical methodologies to enhance the precision and utility of RNA structural studies. There is also sampling bias for the RNA families that were assayed using high-throughput fitness assays, with not all families equally represented. The majority of fitness assays were focused on ribozymes while other RNA families such as mRNA and tRNA had far fewer available datasets to evaluate.

Chemical	Model name	F1-score \uparrow	AUC \uparrow	MAE \downarrow
DMS	EternaFold	0.626	0.683	0.352
	CONTRAFold	0.634	0.693	0.363
	Vienna	0.628	0.678	0.308
	RNAstructure	0.631	0.686	0.310
2A3	EternaFold	0.686	0.740	0.324
	CONTRAFold	0.671	0.723	0.342
	Vienna	0.659	0.710	0.308
	RNAstructure	0.673	0.725	0.310
All	EternaFold	0.658	0.714	0.337
	CONTRAFold	0.654	0.709	0.352
	Vienna	0.645	0.695	0.308
	RNAstructure	0.652	0.706	0.310

Table A11: **RNAGym - Structure prediction benchmark.** F1-score, AUC and MAE between model predictions and experimental measurements (DMS and 2A3) on the RNAGym structure prediction benchmark.

Rank	Model name	Global F1 score	Crossed Pair F1 score
1	RNAstructure	0.637	0.204
2	CONTRAFold	0.594	N/A
3	EternaFold	0.585	N/A
4	Vienna	0.580	N/A

Table A12: **Structure prediction - Pseudoknots.** Global and crossed pair F1 score on PseudoBase sequences not related to sequences in the training set (90% sequence identity). 347 sequences were used out of the 358 sequences in the dataset. Crossed pair F1 score is only for the 'crossed' base pairs (base pairs i - j and m - n with $i < m < j < n$).

G SOCIETAL IMPACT

The advancement of RNA models holds transformative potential across a spectrum of applications. By accurately predicting RNA structure and fitness, researchers can unlock new therapies by targeting previously intractable genetic conditions, enhance crop resilience through agricultural biotechnology, and even engineer microbial systems for cleaner energy production. The creation of benchmarks like RNAGym is crucial in this endeavor, as they drive the field forward by setting standards for model performance and fostering innovation through competition and collaboration. However, as it is the case for any approach that facilitates the development of novel biological sequences for good, the potential misuse of these technologies to create harmful biological agents cannot be ignored (Urbina et al., 2022). It is imperative to proceed with a careful framework that promotes secure use, ethical guidelines, and synthesis monitoring (Baker & Church, 2024) to mitigate risks associated with dual-use capabilities. Ultimately, benchmarks like RNAGym not only validate the effectiveness of emerging RNA models but also, by highlighting the methods leading to step-change performance improvements, encourage their integration into real-world applications, ensuring that these innovations contribute positively to society.

H RESOURCES

Codebase. We open source under an MIT License all resources curated for the RNAGym benchmark via our GitHub repository. In particular, we consolidate of the numerous RNA structure and fitness prediction models discussed in Appendix D and make them available via a common interface, which will facilitate the seamless integration and evaluation of new models as they are developed. This resource aims to provide researchers with robust tools, reducing the technical barrier to entry for conducting advanced RNA analysis and enhancing the reproducibility of results across the scientific community.

Processed datasets. We have made available all datasets used in our fitness and structure prediction benchmarks, including both raw and processed versions, as detailed in Section 2.2. Our GitHub repository provides instructions for downloading these resources. To enhance the utility of our benchmarks, we have included several additional components. For the fitness benchmark, where available, we provide tertiary structure PDB files and multiple sequence alignments for the relevant protein families. In the case of the secondary

Rank	Model name	TM	Δ TM	ρ_{TM}	ρ_{EC}	INF _{WC}	INF _{NWC}
1	NuFold	0.393	-0.151	0.74	0.67	0.78	0.19
2	AlphaFold3	0.386	-0.167	0.71	0.60	0.83	0.26
3	trRosettaRNA	0.382	-0.152	0.60	0.58	0.75	0.12
4	RhoFold+	0.367	-0.172	0.72	0.62	0.50	0.09
5	RoseTTAFold2NA	0.365	-0.154	0.67	0.55	0.74	0.21
1	AlphaFold3	0.381	-0.130	0.55	0.77	0.86	0.37
2	RoseTTAFold2NA	0.168	-0.294	0.24	0.55	0.18	0.00

Table A13: **RNA Gym - 3° structure prediction results.** NuFold leads monomers and AlphaFold3 leads complexes with TM scores of about 0.4, reflecting moderate prediction accuracy despite structural diversity. Monomers shown at top, complexes at bottom.

structure prediction benchmark, we have clustered the sequences at a 40% sequence identity cutoff and provide both the representative clusters as well as the clustering results from mmseq2. From the clustering we created a clean non-redundant 20-80 train-test datasplit. We have mapped all sequences in the train and test set to similar RNA sequences found in the RCSB , PseudoBase (van Batenburg et al., 2000), and bpRNA (Danaee et al., 2018) databases, providing easy access to the rich annotations contained in these databases. Furthermore, to support researchers interested in supervised learning approaches, we offer training datasets for both the fitness and secondary structure prediction tasks (Appendix C).