# RNAINFORMER: GENERATIVE RNA DESIGN WITH TERTIARY INTERACTIONS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

The function of an RNA molecule depends on its structure and a strong structureto-function relationship is already achieved on the secondary structure level of RNA. Therefore, the secondary structure based design of RNAs is one of the major challenges in computational biology. A common approach to RNA design is inverse RNA folding. However, existing RNA design methods cannot invert all folding algorithms because they cannot represent all types of base interactions. In this work, we propose RNAinformer, a novel generative transformer based approach to the inverse RNA folding problem. Leveraging axial-attention, we directly model the secondary structure input represented as an adjacency matrix in a 2D latent space, which allows us to invert all existing secondary structure prediction algorithms. Consequently, RNAinformer is the first model capable of designing RNAs from secondary structures with all base interactions, including non-canonical base pairs and tertiary interactions like pseudoknots and base multiplets. We demonstrate RNAinformer's state-of-the-art performance across different RNA design benchmarks and showcase its novelty by inverting different RNA secondary structure prediction algorithms.

025 026 027

024

004

010 011

012

013

014

015

016

017

018

019

021

#### 028 1 INTRODUCTION

029

Ribonucleic acid (RNA) is one of the major regulatory molecules inside the cells of living organisms with key roles during differentiation and development (Morris & Mattick, 2014). RNAs fold hierar-031 chically (Tinoco Jr & Bustamante, 1999) and the structure is key to their function: Base interactions via hydrogen bonds result in a fast formation of a secondary structure, with tertiary interactions sta-033 bilizing the formation of the final 3D shape (Vicens & Kieft, 2022). A strong structure-to-function 034 relationship is already achieved on a secondary structure level (Hammer et al., 2019), and therefore, RNA secondary structure prediction recently got into the focus of the deep learning community, achieving state-of-the-art results (Singh et al., 2019; Fu et al., 2022; Chen et al., 2022; Franke et al., 037 2022; 2024). Compared to more traditional methods, these algorithms predict an adjacency matrix 038 representation of the secondary structure instead of the commonly used but less expressive dotbracket string notation (Hofacker et al., 1994). This has the advantage that they are not limited to the prediction of specific kinds of base pairs but can model non-Watson-Crick interactions (Olson 040 et al., 2019), pseudoknots (Staple & Butcher, 2005), as well as base multiplets (nucleotides that pair 041 with more than one other nucleotide) (Bhattacharya et al., 2019; Singh et al., 2019), which all play 042 significant roles for RNA structures and functions (Reyes et al., 2009; Vicens & Kieft, 2022). 043

Structure-based RNA design considers the inverse problem: Given a target structure, find an RNA
 primary sequence that folds into the desired structure. It is thus intricately tied to RNA folding.
 However, there is currently no structure-based RNA design algorithm available that can invert state of-the-art deep learning-based secondary structure prediction algorithms, which could offer substan tially improved designs, crucial for synthetic biology and the development of RNA-based therapeutics.

In this work, we propose RNAinformer, the first inverse RNA folding algorithm that is capable of de signing RNAs while considering all kinds of base interactions. Inspired by the RNAformer (Franke
 et al., 2024), we show that a transformer architecture enhanced with axial attention can reliably
 design RNAs in different settings including RNA design with non-canonical interactions, pseudo knots, and base multiplets. Figure 1 shows example designs of the RNAinformer solving tasks that



Figure 1: Example designs for experimentally validated RNA structures that include non-canonical base pairs, pseudoknots, and base multiplets. We show predictions of the RNAinformer that solve the respective structures.

contain non-canonical base pairs, pseudoknots, and base multiplets. We see our main contributions as follows:

- We propose RNAinformer, a novel generative transformer model for the inverse RNA folding problem. Using axial attention, our model is the first RNA design algorithm that can design RNAs from secondary structures with all types of base interactions (Section 3).
- We present a data pipeline for creating synthetic datasets for RNA design, with data splits based on RNA families (Section 4).
- We show that our model outperforms existing algorithms on nested and pseudoknot structures, while further being capable of designing sequences that form base multiplets (Section 5).

Our source code, data, and trained models are publicly available<sup>1</sup>.

084 085

087

070

071

072

073 074

075

076

077

078

079

081

082

### 2 RELATED WORK

Traditional Methods The problem of computational RNA design was first introduced as the inverse RNA folding problem by Hofacker et al. (1994). Since then, different methods were proposed for solving the problem using approaches like local search (Hofacker et al., 1994; Andronescu et al., 2004), constraint programming (Garcia-Martin et al., 2013; 2015; Minuesa et al., 2021), evolutionary methods (Esmaili-Taheri et al., 2014; Esmaili-Taheri & Ganjtabesh, 2015), or multifrontier search (Zhou et al., 2023). However, in contrast to our approach, these methods are limited to the design of nested structures, typically considering canonical base pairs only.

094

Learning Based Approaches More recently, RNA design was also approached with learning 096 based methods. One line of research use human priors to design RNAs based on player strategies obtained from the online gaming platform Eterna (Shi et al., 2018; Koodli et al., 2019). However, 098 these models incorporate human strategies that might not be available for all designs and consider nested structures only. The other, more general approach seeks to learn RNA design purely from 099 data. Eastman et al. (2018) propose to use reinforcement learning (RL) to adjust an initial input 100 sequence by replacing nucleotides based on structural information. In contrast, Runge et al. (2019) 101 and Riley et al. (2023) use a generative approach to the problem. Runge et al. (2019) employs a joint 102 architecture and hyperparameter search approach (Bansal et al., 2022) via automated reinforcement 103 learning (AutoRL) (Parker-Holder et al., 2022) to derive an RL system that is capable of generatively 104 designing RNAs that fold into a desired target structure. Riley et al. (2023) uses a GAN (Goodfellow 105 et al., 2020) approach specifically for the design of toehold switches (Green et al., 2014). However,

<sup>106</sup> 107

<sup>&</sup>lt;sup>1</sup>An anonymized repository is available at https://anonymous.4open.science/r/ RNA-design-7204/

108 all learning-based approaches so far consider RNA design for nested structures only, ignoring pseu-109 doknots and base multiplets, while often being limited to the design of canonical base interactions. 110

111 **Pseudoknotted Structures** Pseudoknots are an important type of base pairs that influence the 112 function of an RNA (Staple & Butcher, 2005). Therefore, some approaches tried to design RNAs 113 from pseudoknotted structures (Taneda, 2012; Kleinkauf et al., 2015; Merleau & Smerlak, 2022). However, these algorithms work on a string notation in dot-bracket format (Hofacker et al., 1994), 114 and thus, they cannot express base multiplets. 115

116 Overall none of the existing algorithms can design RNAs including non-canonical base pairs, pseu-117 doknots, and base multiplets. 118

119 120

**RNA Design from 3D Structures** Besides the described approaches to design RNA based on 121 secondary structure information, recently different methods also tackled the design of RNA se-122 quences based on 3D structure information. The current state-of-the-art physics-based toolkit for 123 biomolecular modeling and design is Rosetta (Leman et al., 2020). However, recently deep learning 124 approaches challenged Rosetta's performance. Joshi & Liò (2024) developed gRNAde, a geometric 125 deep learning-based RNA design pipeline that can be conditioned on RNA 3D backbone structures. Similarly, RDesign (Tan et al., 2024), a hierarchical framework that leverages a contrastive learn-126 127 ing approach and incorporates secondary structure information, and RiboDiffusion (Huang et al., 2024), a diffusion model based on a graph neural network (GNN) (Zhou et al., 2020) structure- and 128 a transformer-based (Vaswani et al., 2017) sequence module, showed remarkable results. 129

130 However, in contrast to RNAinformer, these methods leverage additional 3D information for their 131 designs and, therefore, tackle RNA Design from a different perspective. Often, 3D information 132 of RNA structures is not available and RNA 3D structure prediction is still challenging (Das et al., 2023). Therefore, strong secondary structure based RNA design approaches are highly thought after, 133 also, but not limited to, to achieve better 3D predictions. 134

135 136

#### THE RNAINFORMER 3

137

138 RNA secondary structures can be represented in multiple ways, including the common dot-bracket 139 string notation (Hofacker et al., 1994) or adjacency matrices. We show different representations in Figure 8. One advantage of an adjacency matrix representation is that it can model all types of base 140 interactions, especially if a nucleotide interacts with more than one other, a situation prevalent for 141 most experimentally solved structures (Singh et al., 2019). In the following, we detail our generative 142 approach to designing RNAs from secondary structures using matrix representations. 143

144 Model Our model is a modified auto-regressive encoder-decoder transformer model (Vaswani 145 et al., 2017) with a next token prediction objective. The encoder embeds the structure informa-146 tion, while the decoder auto-regressively generates RNA nucleotide sequences by sampling from 147 the softmax distribution (see Figure 5 in Appendix A). For RNAinformer we use axial attention in 148 the first encoder block to process the adjacency matrix input similar to the RNA former (Franke et al., 149 2024) (Figure 6 in Appendix A). To reduce the memory footprint of the 2D latent operations, we use 150 flash-attention-2 (Dao, 2023) in the axial attention modules. For computational efficiency, we use 151 pooling to reduce the 2D latent representation to a 1D vector that is then passed through the encoder and the decoder to generate candidate sequences. During constrained generation, we pass an addi-152 tional input of the masked RNA sequence to the encoder. The masked sequence is embedded into a 153 2D representation and concatenated to the structure embedding. Similarly, for property condition-154 ing, we embed the target GC-Content using a linear layer and add it to the structure embedding (see 155 Figure 7 in Appendix A). For more details about the RNAinformer architecture and the formulation 156 of the loss, please see Appendix A and B. 157

158 159

160

#### A HOMOLOGY AWARE SYNTHETIC DATA PIPELINE FOR RNA DESIGN 4

While secondary structure information obtained from experimentally validated RNA 3D structure 161 data is considered the gold standard, this data is scarce; only roughly 3% of all available 3D struc162 ture data contains RNAs (Schneider et al., 2023). Therefore, most of the available training data 163 is typically derived from comparative sequence analysis (Choudhary et al., 2017) and, thus, is less 164 reliable. Further, the diversity within existing training sets is limited to only a few RNA families, 165 making it challenging for a folding algorithm to generalize to less represented RNA types (Flamm 166 et al., 2021). This recently raised skepticism in the RNA community that learning based models might not be able to generalize to unseen families (Flamm et al., 2021; Szikszai et al., 2022), which 167 could be a serious concern when using a folding algorithm to validate a given design. To evade this 168 problem, we train the RNAinformer exclusively on synthetic data. This allows us to generate large amounts of training data while enabling us to create a family-based split of the data to avoid learning 170 homologies; a known problem in the RNA folding community (Rivas et al., 2012). In the follow-171 ing, we detail our approach to generate clean family-based synthetic train-/test splits for training the 172 RNAinformer in different settings.

173

174 **Initial Training Data Pool** We generate an initial training data pool using the families of the 175 Rfam database (Version 14.10) (Kalvari et al., 2020). We select all families with covariance models 176 having a maximum CLEN (the number of columns from a sequence alignment defined as consensus 177 (match) columns) of 500 and sample 1,000 sequences for each family from the covariance models 178 using Infernal (Nawrocki & Eddy, 2013). However, while our initial length cutoff is set to 500, 179 roughly 80% of the samples had a max length below 200 nucleotides. Since we use the provided test 180 sets from Singh et al. (2021), which all have a maximum length below 200 nucleotides (see below), for our evaluations on known RNAs obtained from PDB, we decide to use a length cutoff at 200 181 nucleotides to decrease computational costs. The sampled sequences are then annotated using the 182 Rfam covariance models and Infernal. Sequences that hit multiple families, families other than the 183 family they were sampled from, or did not hit any of the families were removed. To reduce intra-clan 184 and intra-family sequence similarity we use CD-HIT (Fu et al., 2012) with a 0.8 threshold to cluster 185 the sequences within a clan and if there is no clan information then within the family. Families or clans with less than 50 clusters are removed and we keep a maximum of 300 representatives of the 187 clusters for each family/clan.

188

200

201

202

203

204

205

189 **Data Splits** We split the data into training, validation, and test sets based on the clan information. 190 All families without a clan annotation are put into the training set. We randomly sample 30 and 25 191 clans for the test and validation sets. For each test clan, we sample 100 sequences to form a test set and 50 sequences from each validation clan to form a validation set. The samples from all other clans 192 are used for training. We then apply CD-HIT with a similarity cutoff of 80% to remove sequence 193 similarity between the training, validation, and test sets, followed by a BLAST-search (Altschul 194 et al., 1997) to further remove training and validation samples that are hit by BLAST for any of the 195 test samples at a high e-value of 10 similar to Singh et al. (2021) but for all test data. 196

Obtaining Structure Information We fold all sequences using different folding algorithms to create three datasets with different structural complexity:

- 1. SynNested: Folded using RNAfold (Lorenz et al., 2011) containing only nested structures.
- 2. **SynPseudoknot**: Folded using HotKnots2.0 (Andronescu et al., 2010) containing both pseudoknotted and nested structures.
- 3. **SynMultiplet**: Folded using RNAformer (Franke et al., 2024) containing structures with all base interactions, including base multiplets and non-canonical base pairs.

Further filtering is done to remove structures with no base pairs and structure duplicates.

206 **Experimental Structures from PDB** To evaluate RNAinformer on known RNAs, we use addi-207 tional test sets, TS1, TS2, TS3, and TS\_hard from Singh et al. (2021), derived from experimental 208 structures of the Protein Data Bank (PDB) (Berman et al., 2000). All the test sets contain struc-209 tures with non-canonical base pairs, pseudoknots and base multiplets. To ensure non-homologous 210 data, we apply an additional homology pipeline to remove homologous RNAs from SynMultiplet 211 that share any sequence or structure similarity to any test sample. Specifically, we build covariance 212 models from multiple sequence alignments for every test RNA employing LocARNA-P (Will et al., 213 2012) and remove any sequences from the training and validation sets that have a hit with any of the resulting covariance models using Infernal as previously described (Runge et al., 2024a). This 214 ensures that there is no data homology between the training and test sets based on structure and 215 sequence similarity.

- All datasets used for our experiments are detailed in Appendix C.
- 218

## 5 EXPERIMENTS

219 220

We evaluate the RNA informer on three RNA design paradigms, inverse RNA folding, constrained design, and RNA design with desired properties. We show that the RNAinformer can approach 222 inverse RNA folding and RNA design with desired properties for tasks with increasing structural 223 diversity by first evaluating RNA design for nested structures in Section 5.1, before tackling RNA 224 design for pseudoknotted structures in Section 5.2. We then demonstrate RNAinformer's capa-225 bility to conditionally generate RNA sequences for the real-world task of designing theophylline 226 riboswitches following Runge et al. (2024a) in Section 5.3. We conclude with an assessment of 227 RNAinformer's ability to design RNAs from secondary structures that contain all kinds of base in-228 teraction in Section 5.4. Here, we also compare our strategy to train on synthetic data with a more 229 commonly used fine-tuning strategy using two versions of the RNAinformer.

During evaluation, we generate 20 candidate sequences with the RNAinformer for each task except for the PDB structures where we instead generate 100 candidate sequences. The first sequence is generated using a greedy strategy and the rest are generated using multinomial sampling. We set the threshold for satisfying the property constraint as  $\epsilon = 0.01$ .

234

240

Metrics The ultimate goal of structure-based RNA design is to generate sequences that fold back
 into the target structure. Following the common convention in the field of RNA design, we report
 the number of solved tasks for a given benchmark dataset. However, we provide a more com prehensive analysis of all experiments with different performance measures, described in detail in
 Appendix D.2.

241 Training Details We train our model with 6 encoder blocks and 6 decoder blocks with an em-242 bedding dimension of 256. The model is trained using cosine annealing learning rate schedule with warm-up and AdamW (Loshchilov & Hutter, 2019). We train separate models for each training 243 dataset and for each training dataset we also train a separate GC-content conditioned model. The 244 constrained design models were trained with a maximum length of 100 while the rest were trained 245 with length 200. The longer models were trained across 2 A40 GPUs with an effective batch size 246 of 128 for 50,000 steps having a runtime of  $\sim$ 18 hours. The GC-content of the original sequences 247 is used as the target GC-content. The hyperparameters used for training our model are described in 248 Table 2 in Appendix B.

- 249 250 251
- 5.1 RNA DESIGN FOR NESTED STRUCTURES

252 We first evaluate the RNA informer's ability to design RNA sequences for nested structures with only 253 canonical base pair interactions on the SynNested test set (see Table 3 in Appendix C). Additionally, 254 we also evaluate its ability to design sequences with desired GC-content. We compare the perfor-255 mance of the RNAinformer for inverse folding against one of the currently best-performing set of algorithms, LEARNA, Meta-LEARNA, Meta-LEARNA-Adapt (Runge et al., 2019), libLEARNA 256 (Runge et al., 2024b) and SAMFEO (Zhou et al., 2023). For design with desired GC-content we 257 compare it against libLEARNA. The LEARNA suite algorithms and libLEARNA were run with a 258 timeout of 30 seconds per sample and SAMFEO was run for 1,000 iterations for each task. We 259 then select the best 20 candidates for evaluation. Corresponding to the data generation, the designed 260 candidates are folded using RNAfold Lorenz et al. (2011) for evaluation; note that all competitors 261 also used RNAfold for training and/or evaluation in their respective original publications. 262

263 264 RESULTS

Inverse Folding From Table 1 we observe that RNAinformer outperforms most of its competitors
 except SAMFEO, solving 91.8% of the tasks. Furthermore, RNAinformer generates multiple, highly
 diverse solutions for each task, indicated by a high diversity score of 0.699 as shown in Table 10
 in Appendix E.1. Remarkably, this performance is achieved by sampling only 20 sequences from
 the RNAinformer without any post-processing strategies as e.g. implemented in the local search
 strategy of the LEARNA suite of algorithms. However, while SAMFEO is capable of solving more

270 Table 1: Performance on the nested and pseudoknotted structures of the SynNested and SynPseudo-271 knot datasets, respectively, for Inverse Folding (IF) and Desired GC-content design (GC). We report 272 the % tasks solved in 20 designed sequences.

Model	Synl	Nested	SynPseudoknot	
Model	Solved (IF) [%]	Solved (GC) [%]	Solved PK (IF) [%]	Solved PK (GC) [%]
RNAinformer	91.8	69.6	68.5	33.7
LEARNA	63.9	X	X	X
Meta-LEARNA	36.8	X	X	X
Meta-LEARNA-Adapt	37.2	X	X	X
libLEARNA	77.2	59.0	X	X
SAMFEO	99.6	X	X	X
antaRNA	X	×	15.6	1.2

284 285

291

294

273

tasks, it generates solutions with much lower diversity (0.106) compared to the RNA informer. One 286 reason is that SAMFEO uses an initialization strategy that itself could already solve 78% of the tasks, 287 leveraging biases in the internal scoring functions of RNAfold by placing low energy GC-pairs at 288 paired positions and single A nucleotides at unpaired positions that cannot pair with G or C in the 289 limited model of RNAfold. Despite increasing the ability to solve the design tasks, this approach 290 has the disadvantage that the resulting candidates rarely contain U nucleotides, typically resulting in high GC nucleotide ratios (GC-content), which can drastically impact the function of the resulting 292 RNAs (Isaacs et al., 2006). In contrast, the designs of the RNAinformer do not show similar bias as 293 indicated by the high sequence diversity.

295 **Desired GC-Content Design** The GC-content conditioned RNAinformer model solves  $\sim 10\%$ 296 more tasks than libLEARNA, the only competitor capable of also generating RNAs with desired 297 GC-contents, as shown in Table 1. Further, our results in Table 11 in Appendix E.1 demonstrate 298 that even with the GC-content constraints, RNAinformer can still generate multiple highly diverse 299 solutions for each task. We also note that the average GC-content error of the candidate sequences generated by RNA informer is very low (0.01), indicating that the model actively generates sequences 300 with GC-content close to the desired target value. 301

302 303

#### 5.2 RNA DESIGN WITH PSEUDOKNOTS

304 In this section, we assess the performance of RNAinformer when designing RNAs for pseudoknotted 305 input structures. Pseudoknots are tertiary interactions that typically connect local geometries of the 306 RNA secondary structure by establishing long-range interactions between nucleotides. We compare 307 the RNAinformer against antaRNA (Kleinkauf et al., 2015) with HotKnots2.0 (Andronescu et al., 308 2010) as the folding algorithm; one out of three folding engines available for antaRNA which was 309 also used for data generation. Again, we provide results for both inverse-folding and for the design 310 with desired GC-content. We evaluate the RNAinformer on both the pseudoknotted structures(pK) and nested structures(pK-free) of the SynPseudoknot Dataset (see Table 4 in Appendix C). However, 311 due to the long runtime of antaRNA and its internal ant-colony optimization strategy, we only evalu-312 ate one design candidate, supposed to solve the task, and limit the comparison to the pseudoknotted 313 structures (pK). However, there are many more intermediate sequences designed by antaRNA before 314 outputting the final design. To make a fair comparison, we additionally evaluate the first design of 315 RNAinformer for completeness. 316

317 RESULTS 318

319 **Inverse Folding** RNAinformer significantly outperforms antaRNA, solving  $\sim$ 50% more pseudo-320 knot tasks, as shown in Table 1 (right). Remarkably, RNAinformer also achieves high performance 321 on the nested structures solving more than 90% of the tasks as shown in Table 12 in Appendix E.2. Generally, we observe that RNAinformer is able to generate multiple solutions with high diversity 322 for both the nested and the pseudoknotted structures. The high F1 and MCC scores of RNAinformer 323 further indicate that the designs are close to a solution even for the tasks that could not be solved with 20 candidates (see Table 12 in Appendix E.2). Notably, the RNAinformer also solves 39.1%
 of the tasks with the first sequence generated, still outperforming antaRNA by solving twice the number of tasks (see Table 13 in Appendix E.2).

328 **Desired GC-Content Design** Similar to the unconditional generation, the RNAinformer also outperforms antaRNA for the conditional design of RNAs with desired GC-contents by a large margin as shown in Table 1. The RNAinformer roughly solves one third of the pseudoknotted structures 330 (33.7%) compared to only 1.2% solved tasks by antaRNA. Although the number of solutions gen-331 erated by the conditioned RNAinformer model is less compared to the unconditioned model, the 332 diversity of the solutions is maintained as shown in Table 13 in Appendix E.2. Moreover, the 333 RNAinformer is capable of solving almost two-thirds of the nested structures (65.8%) while again 334 generating sequences with low GC-content error, indicating closeness to the target value (see Ta-335 ble 13 in Appendix E.2). Again we evaluate the first prediction of RNAinformer and observe that it 336 solves nearly 15 times the number of tasks compared to antaRNA (solving 17.7% of the tasks with 337 a single shot; see Table 13 in Appendix E.2).

338 339 340

345

346

347

348

367

#### 5.3 AUTOMATED DESIGN OF THEOPHYLLINE RIBOSWITCHES

We evaluate the RNAinformer's ability to do constrained design by tackling the design of synthetic
theophylline riboswitches. We use the design space formulation from Runge et al. (2024b), which
was created by combining the shared sequence and structure motifs of the proposed constructs by
Wachsmuth et al. (2012), defined as,

349 where ? represent masked out positions and ? represent positions for extensions. The different sections of the construct are highlighted, (i) Aptamer (Red), (ii) Spacer (Green), (iii) the Comple-350 mentary Sequence (Blue), and (iv) the 8-U Stretch (Black). We use the above formulation to sample 351 tasks for our evaluation. Since we do not have any ground truth sequences to get target GC-contents 352 we test RNAinformer's ability for conditional generation on a range of GC-content values for each 353 task. We compare our models against libLEARNA for both inverse folding and design with desired 354 GC-content. For each task, we again generate 20 candidates and use RNAfold to fold them, in line 355 with the original procedure described at Runge et al. (2024b). 356

357 Training Data For training, we use the Training 358 Short and Validation datasets provided by Runge 359 et al. (2024b) (see Table 7 in Appendix C). The 360 datasets were generated by sampling from the Rfam 361 database version 14.1 and folding all the sequences 362 with RNAfold. The structures and sequences were then randomly masked to create the final datasets 363 for constrained design. During training, target GC-364 content values were obtained from the unmasked sequences. 366

**Riboswitch Tasks** We generate an exhaustive set 368 of riboswitch design tasks for evaluation using the 369 design formulation 1. We sample masked sequences 370 for each of the extension positions while consider-371 ing the length constraints for each part. We filter the 372 tasks for the seven GC-content targets, listed in Ta-373 ble 8 in Appendix C, based on the possible range of 374 GC-contents for each task as calculated from their 375 length and masked sequence. Few of the generated tasks have no valid sequences possible when evalu-376



Figure 2: Comparison between RNAinformer and libLEARNA for Riboswitch design with desired GC-content.

ating using RNAfold. As it is not feasible to determine all the un-designable tasks we do not filter them out.





#### RESULTS

**Inverse Folding** Both RNAinformer and libLEARNA are capable of solving almost all of the riboswitch tasks (solving >90% of the tasks) as shown in 14 in Appendix E.3. However, the RNAinformer slightly outperforms libLEARNA. Furthermore, the RNAinformer generates more solutions (valid sequences) compared to libLEARNA. Due to the fixed sequence constraints, the observed diversity of the solution sequences is rather low.

Desired GC-Content Design The results for the desired GC-content design are shown in Figure 2. We observe that libLEARNA outperforms RNAinformer for the smaller GC-content targets (0.3,0.35), whereas RNAinformer outperforms libLEARNA for the larger target values (0.55,0.6) and the performance is almost identical for the mid-range targets. However, the performance of the RNAinformer is generally significantly more consistent across all the target GC-content values. Additional results are shown in Table 15 in Appendix E.3.

408 409 410

393

396 397

399

400

401

402

#### 5.4 RNA DESIGN WITH ALL KINDS OF BASE INTERACTIONS

411 In this section, we evaluate RNA informer for designing RNAs from structure data that contains all 412 kinds of base pairs including pseudoknots and base multiplets on experimentally validated structures 413 from the PDB (Berman et al., 2000) and the SynMultiplet Dataset (see Table 5 in Appendix C). 414 To account for the difficulty of the task, we design 100 candidate sequences instead of only 20 415 sequences. The RNAinformer is the only method capable of tackling the task of designing RNAs 416 for structures that contain base multiplets; consequently, we cannot compare our designs with other 417 methods from the field. Instead, we compare against a simple baseline that uniformly samples RNA sequences for both inverse folding and desired GC-content design and a GNN baseline where we 418 use the implementation of structTransformer provided by Ingraham et al. (2019) and run it with 419 the same batch size and steps as the RNAinformer. According to the data generation, the designed 420 candidate sequences are folded using RNA former Franke et al. (2024) for evaluation. 421

422 423 RESULTS

424 **Inverse Folding** RNAinformer significantly outperforms the randomly designed sequences and 425 the GNN baseline as shown in Table 16. From Figure 3a we observe that RNAinformer further is 426 the first method capable of solving experimentally determined structures across the PDB test sets 427 including structures with base multiplets. However, the overall solved rate is rather low, indicating 428 that designing sequences for structures with all kinds of base pairs seems to be much more challenging than for nested structures or structures with pseudoknots only. Examples of solved experimental 429 structures with base multiplets are shown in Figure 1. Despite relatively low rates of solved tasks, 430 the RNAinformer still achieves high F1 and MCC scores, indicating that the designed candidates 431 have high structural accuracy. Notably, the RNAinformer achieves similar performance on the SynMultiplet Test set. In addition, the RNAinformer is able to generate multiple solutions with high
diversity for all the test sets except TS3 and can generate solutions with non-canonical base pairs,
which are typically ignored by other design algorithms.

Desired GC-Content Design As shown in Figure 3b, the GC-conditioned RNAinformer can solve tasks across the PDB test sets including structures with base multiplets, even with the additional GC-content constraint. While the number of solutions generated drops significantly, the diversity of the generated sequences is maintained. Similar to the unconditioned model, we observe high F1 and MCC scores indicating structural similarity to the ground truth, and that the RNAinformer generates solutions with non-canonical base pairs and candidate sequences with a low GC-content error across the test sets.

Detailed results are shown Table 17 and Table 18 in Appendix E.4.

445 Synthetic data vs Real-World data To val-446 idate our strategy to only use synthetic data 447 during training of the RNAinformer, we addi-448 tionally train a model on known RNA data us-449 ing the inter-family dataset from Runge et al. 450 (2024a). The training data was collected from 451 multiple public sources: bpRNA-1m (Danaee et al., 2018), ArchiveII (Sloma & Mathews, 452 2016) and RNAStrAlign (Tan et al., 2017) from 453 Chen et al. (2020), RNA-Strand (Andronescu 454 et al., 2008) and PDB (Berman et al., 2000). 455 Homologies between the training, validation, 456 and test sequences were removed by filtering 457 using CD-Hit (Fu et al., 2012) and BLAST-458 Search (Altschul et al., 1997). An additional 459 homology reduction based on structure similar-460 ity was applied using covariance models of the 461 PDB test sets (TS1, TS2, TS3 and TS-Hard). 462 We train an RNAinformer model (NAT + FT) on it and further fine-tuned it on the PDB train-463 ing samples. We compare it against an RNAin-464



Figure 4: Difference in F1-Scores between the folded structures for the designed sequences and the PDB test set sequences.

former model (Syn) pre-trained on the SynMultiplet dataset and a second model that was also pretrained on synthetic data but finetuned on the PDB samples (Syn + FT) similar to the NAT + FT
model.

The results are shown in Table 19 in Appendix E.4. Surprisingly, the RNAinformer model pre-468 trained on synthetic data performs significantly better than the model that was pre-trained and fine-469 tuned on known RNAs, having nearly double the F1 score across the test sets. However, the addi-470 tional finetuning appears beneficial, as indicated by slightly higher scores of the Syn + FT model. 471 To further assess this, we evaluated the Syn and Syn + FT models on RNA only samples from 472 the Critical Assessment of Structure Prediction 15 (CASP15) competition. As before, we use the 473 RNAformer for secondary structure predictions but additionally employ AlphaFold 3 (Abramson 474 et al., 2024) for 3D structure prediction of the generated sequences. The results are shown in Ta-475 bles 20 and 24 for 2D and 3D predictions, respectively. In contrast to the results for the PDB test 476 set, we observe that the finetuned model achieves slightly worse performance than the model trained 477 on synthetic data only for both folding engines, RNAinformer and AlphaFold 3. We conclude that 478 training on synthetic data appears beneficial for RNA informer compared to training on known RNAs 479 only, and finetuning on experimentally validated structures can lead to slightly better performance in some cases. 480

481

Improved Foldability of Designed Sequences The RNAinformer's designed sequences on aver age achieve higher structure F1 scores across the PDB test sets compared to RNAformer's original
 predictions on the PDB test set sequences. We take this as an indicator that the RNAinformer learns
 to design sequences that are better foldable by the RNAformer for the experimental test structures.
 However, to ensure that these results are not artifacts resulting from overfitting the RNAformer's

486 distribution, we folded the best-designed sequences by the RNAinformer for each of the PDB test 487 structures using other folding algorithms (SPOT-RNA (Singh et al., 2019), MXFold2 (Sato et al., 488 2021) and UFold (Fu et al., 2022)) and compared the F1 scores against the folding algorithms' 489 predictions on the PDB test set sequences. We observe that designed sequences have better or sim-490 ilar F1 scores for almost all the folding algorithms, showing consensus agreement for the designed sequences (see Figure 4). The F1 scores of the designed sequences and the folding algorithms eval-491 uations on the PDB test sets are reported in Table 21 in Appendix E.4. We again also folded the 492 designed sequences with AlphaFold 3 and analyzed the results regarding TM score and RMSD. The 493 results are shown in Table 22. We find that RNAinformer cannot improve the TM scores compared 494 to TM scores achieved when folding the ground truth sequence with AlphaFold 3. However, Al-495 phaFold 3 is known to struggle with predictions for so-called orphan RNAs (Bernard et al., 2024) 496 - sequences where there exist no homologs in the database - due to its dependence on multiple 497 sequence alignments (MSAs). Therefore, we analyzed the AlphaFold 3 predictions in more detail 498 and find that for all the designed sequences AlphaFold 3 was not able to find an MSA during the 499 search while roughly 59% of the PDB samples had multiple homologs. We show the difference in 500 TM-Score when splitting the data into orphan and non-orphan RNAs in Table 23. We observe that 501 the designed sequences of the RNAinformer slightly improve the TM score for three out of four datasets for orphan RNAs. However, when there is MSA available for the ground truth but not for 502 the designed sequences, the TM score drops drastically for the designed candidates. 503

504 505

506

#### 6 CONCLUSION, LIMITATIONS & FUTURE WORK

507 In this work, we propose RNAinformer, the first RNA design algorithm capable of designing RNA 508 sequences from secondary structures that contain all kinds of base interactions, including non-509 canonical base pairs, pseudoknots, and base multiplets. Using axial-attention, the RNAinformer 510 leverages a 2D latent representation to process adjacency matrix representations of RNA secondary structures to achieve state-of-the-art results in structure based RNA design. We demonstrate the 511 strong performance of RNAinformer on tasks with nested structures only, tasks that contain pseudo-512 knots, as well as on experimentally derived structures with all kinds of base interactions. We observe 513 high diversity across all designs and tasks and improved foldability of the designed sequences com-514 pared to their known counterparts. 515

516

Limitations While showing overall strong performance, there is still room for improvement, particularly for the design for known RNA structures. Further, while we reduce computational complexity using a pooling operation to map the 2D latent representation to a 1D vector, training the RNAinformer is memory intensive. As a result, we only train the RNAinformer with a sequence length cutoff at 200 nucleotides. While this is sufficient for current benchmarks, a higher length cutoff would further increase the usability of our approach.

Future Work We think that RNAinformer is a useful basis for future approaches to RNA design and expect it to be of great value for the RNA design community. Future work could e.g. focus on improving the memory footprint of the RNAinformer.

- 526
- 527
- 528 529
- 530
- 531
- 532
- 533 534
- 535
- 536
- 536
- 538
- 539

## 540 REPRODUCIBILITY STATEMENT

541 542 543

544

546

547

548

549

550

551

To ensure the reproducibility of our results, we have made our source code, model checkpoints, and datasets publicly available in the anonymous repository https://anonymous.4open. science/r/RNA-design-7204/. The repository contains detailed instructions for setting up the environment, including specific Python package versions (see environment.yml). Model checkpoints and predictions for all experiments are provided in the runs.ta.xz file. Our datasets are provided in the data.tar.xz. Links to download both files are in the repository. We provide scripts for evaluating trained models (eval.py) and for reproducing our training procedures (train.py with configs for all trained models also provided. Scripts for running inference on the test sets is also provided (inference.py). Hardware requirements (GPU specifications) for both training and inference are clearly stated. Evaluation on provided predictions can be done without gpus.

552 553 554

555

556

557

558

562

563

564

565

573

586

592

#### References

- Josh Abramson, Jonas Adler, Jack Dunger, Richard Evans, Tim Green, Alexander Pritzel, Olaf Ronneberger, Lindsay Willmore, Andrew J Ballard, Joshua Bambrick, et al. Accurate structure prediction of biomolecular interactions with alphafold 3. *Nature*, pp. 1–3, 2024.
- Stephen F Altschul, Thomas L Madden, Alejandro A Schäffer, Jinghui Zhang, Zheng Zhang, Webb
   Miller, and David J Lipman. Gapped blast and psi-blast: a new generation of protein database
   search programs. *Nucleic acids research*, 25(17):3389–3402, 1997.
  - Mirela Andronescu, Anthony P Fejes, Frank Hutter, Holger H Hoos, and Anne Condon. A new algorithm for rna secondary structure design. *Journal of molecular biology*, 336(3):607–624, 2004.
- Mirela Andronescu, Vera Bereg, Holger H. Hoos, and Anne Condon. Rna strand: The rna secondary
   structure and statistical analysis database. *BMC Bioinformatics*, 9:340 340, 2008.
- Mirela S Andronescu, Cristina Pop, and Anne E Condon. Improved free energy parameters for rna pseudoknotted secondary structure prediction. *Rna*, 16(1):26–42, 2010.
- Viraj Bagal, Rishal Aggarwal, PK Vinod, and U Deva Priyakumar. Molgpt: Molecular generation
   using a transformer-decoder model. *Journal of Chemical Information and Modeling*, 2021.
- Archit Bansal, Danny Stoll, Maciej Janowski, Arber Zela, and Frank Hutter. Jahs-bench-201: A
  foundation for research on joint architecture and hyperparameter search. *Advances in Neural Information Processing Systems*, 35:38788–38802, 2022.
- Helen M Berman, John Westbrook, Zukang Feng, Gary Gilliland, Talapady N Bhat, Helge Weissig,
  Ilya N Shindyalov, and Philip E Bourne. The protein data bank. *Nucleic acids research*, 28(1):
  235–242, 2000.
- Clément Bernard, Guillaume Postic, Sahar Ghannay, and Fariza Tahi. Has alphafold 3 reached its success for rnas? *bioRxiv*, pp. 2024–06, 2024.
- Sohini Bhattacharya, Ayush Jhunjhunwala, Antarip Halder, Dhananjay Bhattacharyya, and Abhijit
   Mitra. Going beyond base-pairs: topology-based characterization of base-multiplets in rna. *RNA*, 25(5):573–589, 2019.
- Jiayang Chen, Zhihang Hu, Siqi Sun, Qingxiong Tan, Yixuan Wang, Qinze Yu, Licheng Zong, Liang
   Hong, Jin Xiao, Irwin King, et al. Interpretable rna foundation model from unannotated data for
   highly accurate rna structure and function predictions. *arXiv preprint arXiv:2204.00300*, 2022.
- Xinshi Chen, Yu Li, Ramzan Umarov, Xin Gao, and Le Song. {RNA} secondary structure prediction
   by learning unrolled algorithms. In *International Conference on Learning Representations*, 2020.
- 593 Davide Chicco and Giuseppe Jurman. The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation. *BMC genomics*, 21:1–13, 2020.

619

632

394	Krishna Choudhary, Fei Deng, and Sharon Aviran. Comparative and integrative analysis of rna
595	structural profiling data: current practices and emerging questions. <i>Quantitative Biology</i> , 5(1):
596	3-24, 2017.
597	

- 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, 2018.
- Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv* preprint arXiv:2307.08691, 2023.
- Rhiju Das, Rachael C Kretsch, Adam J Simpkin, Thomas Mulvaney, Phillip Pham, Ramya Rangan,
  Fan Bu, Ronan M Keegan, Maya Topf, Daniel J Rigden, et al. Assessment of three-dimensional
  rna structure prediction in casp15. *Proteins: Structure, Function, and Bioinformatics*, 91(12):
  1747–1770, 2023.
- Peter Eastman, Jade Shi, Bharath Ramsundar, and Vijay S Pande. Solving the rna design problem with reinforcement learning. *PLoS computational biology*, 14(6):e1006176, 2018.
- Ali Esmaili-Taheri and Mohammad Ganjtabesh. Erd: a fast and reliable tool for rna design including
   constraints. *BMC bioinformatics*, 16(1):20, 2015.
- Ali Esmaili-Taheri, Mohammad Ganjtabesh, and Morteza Mohammad-Noori. Evolutionary solution for the rna design problem. *Bioinformatics*, 30(9):1250–1258, 2014.
- Christoph Flamm, Julia Wielach, Michael T Wolfinger, Stefan Badelt, Ronny Lorenz, and Ivo L
   Hofacker. Caveats to deep learning approaches to rna secondary structure prediction. *Biorxiv*, pp. 2021–12, 2021.
- Jörg Franke, Frederic Runge, and Frank Hutter. Probabilistic transformer: Modelling ambiguities and distributions for rna folding and molecule design. *Advances in Neural Information Processing Systems*, 35:26856–26873, 2022.
- Jörg K.H. Franke, Frederic Runge, Ryan Köksal, Rolf Backofen, and Frank Hutter. Rnaformer: A
   simple yet effective deep learning model for rna secondary structure prediction. *bioRxiv*, 2024.
   doi: 10.1101/2024.02.12.579881.
- Laiyi Fu, Yingxin Cao, Jie Wu, Qinke Peng, Qing Nie, and Xiaohui Xie. Ufold: fast and accurate rna secondary structure prediction with deep learning. *Nucleic acids research*, 50(3):e14–e14, 2022.
- Limin Fu, Beifang Niu, Zhengwei Zhu, Sitao Wu, and Weizhong Li. Cd-hit: accelerated for cluster ing the next-generation sequencing data. *Bioinformatics*, 28(23):3150–3152, 2012.
- Juan Antonio Garcia-Martin, Peter Clote, and Ivan Dotu. Rnaifold: a constraint programming algorithm for rna inverse folding and molecular design. *Journal of Bioinformatics and Computational Biology*, 11(02):1350001, 2013.
- Juan Antonio Garcia-Martin, Ivan Dotu, and Peter Clote. Rnaifold 2.0: a web server and software to design custom and rfam-based rna molecules. *Nucleic Acids Research*, 43(W1):W513–W521, 2015.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
   Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- Alexander A Green, Pamela A Silver, James J Collins, and Peng Yin. Toehold switches: de-novodesigned regulators of gene expression. *Cell*, 159(4):925–939, 2014.
- Stefan Hammer, Christian Günzel, Mario Mörl, and Sven Findeiß. Evolving methods for rational de novo design of functional rna molecules. *Methods*, 161:54 63, 2019. ISSN 1046-2023. Development and engineering of artificial RNAs.

651

660

665

682

683

684

685

688

694

- Ivo Hofacker, Walter Fontana, Peter Stadler, Sebastian Bonhoeffer, Manfred Tacker, and Peter Schuster. Fast Folding and Comparison of RNA Secondary Structures. *Monatshefte fuer Chemic/Chemical Monthly*, 125:167–188, 02 1994.
- Han Huang, Ziqian Lin, Dongchen He, Liang Hong, and Yu Li. Ribodiffusion: tertiary structurebased rna inverse folding with generative diffusion models. *Bioinformatics*, 40(Supplement\_1):
  i347–i356, 2024.
- John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graphbased protein design. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper\_files/ paper/2019/file/f3a4ff4839c56a5f460c88cce3666a2b-Paper.pdf.
- Farren J Isaacs, Daniel J Dwyer, and James J Collins. Rna synthetic biology. *Nature biotechnology*, 24(5):545–554, 2006.
- Chaitanya K Joshi and Pietro Liò. grnade: A geometric deep learning pipeline for 3d rna inverse design. In *RNA Design: Methods and Protocols*, pp. 121–135. Springer, 2024.
- Ioanna Kalvari, Eric P Nawrocki, Nancy Ontiveros-Palacios, Joanna Argasinska, Kevin
  Lamkiewicz, Manja Marz, Sam Griffiths-Jones, Claire Toffano-Nioche, Daniel Gautheret, Zasha
  Weinberg, Elena Rivas, Sean R Eddy, Robert D Finn, Alex Bateman, and Anton I Petrov. Rfam
  14: expanded coverage of metagenomic, viral and microRNA families. *Nucleic Acids Research*,
  49(D1):D192–D200, 11 2020. ISSN 0305-1048.
- Robert Kleinkauf, Torsten Houwaart, Rolf Backofen, and Martin Mann. antaRNA–Multi-objective
  inverse folding of pseudoknot RNA using ant-colony optimization. *BMC bioinformatics*, 16(1):
  389, 2015.
- Rohan V Koodli, Benjamin Keep, Katherine R Coppess, Fernando Portela, Eterna participants, and
  Rhiju Das. Eternabrain: Automated rna design through move sets and strategies from an internetscale rna videogame. *PLoS computational biology*, 15(6):e1007059, 2019.
- Julia Koehler Leman, Brian D Weitzner, Steven M Lewis, Jared Adolf-Bryfogle, Nawsad Alam, Rebecca F Alford, Melanie Aprahamian, David Baker, Kyle A Barlow, Patrick Barth, et al. Macromolecular modeling and design in rosetta: recent methods and frameworks. *Nature methods*, 17 (7):665–680, 2020.
  - Ronny Lorenz, Stephan H. Bernhart, Christian Höner zu Siederdissen, Hakim Tafer, Christoph Flamm, Peter F. Stadler, and Ivo L. Hofacker. Viennarna package 2.0. Algorithms for Molecular Biology, 6(1):26, Nov 2011. ISSN 1748-7188.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer- ence on Learning Representations*, 2019.
- Nono S. C. Merleau and Matteo Smerlak. arnaque: an evolutionary algorithm for inverse pseudo knotted rna folding inspired by lévy flights. *BMC Bioinformatics*, 23, 2022.
- Gerard Minuesa, Cristina Alsina, Juan Antonio Garcia-Martin, Juan Carlos Oliveros, and Ivan Dotu.
   Moirnaifold: a novel tool for complex in silico rna design. *Nucleic acids research*, 49(9):4934–4943, 2021.
- Kevin V Morris and John S Mattick. The rise of regulatory rna. *Nature Reviews Genetics*, 15(6): 423–437, 2014.
- Eric P. Nawrocki and Sean R. Eddy. Infernal 1.1: 100-fold faster rna homology searches. *Bioinformatics*, 29:2933 2935, 2013.
- Wilma K Olson, Shuxiang Li, Thomas Kaukonen, Andrew V Colasanti, Yurong Xin, and Xiang-Jun
   Lu. Effects of noncanonical base pairing on rna folding: structural context and spatial arrangements of g· a pairs. *Biochemistry*, 58(20):2474–2487, 2019.

702 703 704 705	Jack Parker-Holder, Raghu Rajan, Xingyou Song, André Biedenkapp, Yingjie Miao, Theresa Eimer, Baohe Zhang, Vu Nguyen, Roberto Calandra, Aleksandra Faust, Frank Hutter, and Marius Lin- dauer. Automated reinforcement learning (autorl): A survey and open problems. <i>Journal of</i> <i>Artificial Intelligence Research (JAIR)</i> , 74:517–568, 2022.
706 707 708	Francis E Reyes, Andrew D Garst, and Robert T Batey. Strategies in rna crystallography. <i>Methods in enzymology</i> , 469:119–139, 2009.
709 710 711	Aidan T. Riley, James M. Robson, and Alexander A. Green. Generative and predictive neural net- works for the design of functional rna molecules. <i>bioRxiv</i> , 2023.
712 713 714	Elena Rivas, Raymond Lang, and Sean R Eddy. A range of complex probabilistic models for rna secondary structure prediction that includes the nearest-neighbor model and more. <i>RNA</i> , 18(2): 193–212, 2012.
715 716	Frederic Runge, Danny Stoll, Stefan Falkner, and Frank Hutter. Learning to design RNA. In Inter- national Conference on Learning Representations, 2019.
717 718 719	Frederic Runge, Karim Farid, Jorg K.H. Franke, and Frank Hutter. Rnabench: A comprehensive library for in silico rna modelling. <i>bioRxiv</i> , 2024a.
720 721	Frederic Runge, Jörg Franke, Daniel Fertmann, Rolf Backofen, and Frank Hutter. Partial rna design. <i>Bioinformatics</i> , 40(Supplement_1):i437–i445, 2024b.
723 724	Kengo Sato, Manato Akiyama, and Yasubumi Sakakibara. Rna secondary structure prediction using deep learning with thermodynamic integration. <i>Nature communications</i> , 12(1):1–9, 2021.
725 726 727 728	Bohdan Schneider, Blake Alexander Sweeney, Alex Bateman, Jiri Cerny, Tomasz Zok, and Marta Szachniuk. When will rna get its alphafold moment? <i>Nucleic Acids Research</i> , 51(18):9522–9532, 2023.
729 730	Jade Shi, Rhiju Das, and Vijay S Pande. Sentrna: Improving computational rna design by incorporating a prior of human design strategies. <i>arXiv preprint arXiv:1803.03146</i> , 2018.
731 732 733 734	Jaswinder Singh, Jack Hanson, Kuldip Paliwal, and Yaoqi Zhou. Rna secondary structure predic- tion using an ensemble of two-dimensional deep neural networks and transfer learning. <i>Nature</i> <i>communications</i> , 10(1):1–13, 2019.
735 736 737	Jaswinder Singh, Kuldip Paliwal, Tongchuan Zhang, Jaspreet Singh, Thomas Litfin, and Yaoqi Zhou. Improved rna secondary structure and tertiary base-pairing prediction using evolutionary profile, mutational coupling and two-dimensional transfer learning. <i>Bioinformatics</i> , 37, 2021.
738 739 740	Michael F Sloma and David H Mathews. Exact calculation of loop formation probability identifies folding motifs in rna secondary structures. <i>RNA</i> , 22(12):1808–1818, 2016.
741 742	David W Staple and Samuel E Butcher. Pseudoknots: Rna structures with diverse functions. <i>PLoS biology</i> , 3(6):e213, 2005.
743 744 745 746	Marcell Szikszai, Michael Wise, Amitava Datta, Max Ward, and David H Mathews. Deep learn- ing models for rna secondary structure prediction (probably) do not generalize across families. <i>Bioinformatics</i> , 38(16):3892–3899, 2022.
747 748 749 750	Cheng Tan, Yijie Zhang, Zhangyang Gao, Bozhen Hu, Siyuan Li, Zicheng Liu, and Stan Z. Li. RDesign: Hierarchical data-efficient representation learning for tertiary structure-based RNA design. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=RemfXx7ebP.
751 752 753 754	Zhen Tan, Yinghan Fu, Gaurav Sharma, and David H Mathews. Turbofold ii: Rna structural alignment and secondary structure prediction informed by multiple homologs. <i>Nucleic acids research</i> , 45(20):11570–11581, 2017.
	Akita Tanada, Multi abiastiva ganatia algorithm for psaudoknottad ma saguanga dasign. <i>Frontians</i>

755 Akito Taneda. Multi-objective genetic algorithm for pseudoknotted rna sequence design. *Frontiers in Genetics*, 3:36, 2012.

756 757 758	Ignacio Tinoco Jr and Carlos Bustamante. How rna folds. <i>Journal of molecular biology</i> , 293(2): 271–281, 1999.
759 760 761 762	<ul> <li>Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017.</li> </ul>
763 764	Quentin Vicens and Jeffrey S Kieft. Thoughts on how to think (and talk) about RNA structure. <i>Proceedings of the National Academy of Sciences</i> , 119(17):e2112677119, 2022.
765 766 767 768	Manja Wachsmuth, Sven Findeiß, Nadine Weissheimer, Peter F. Stadler, and Mario Mörl. De novo design of a synthetic riboswitch that regulates transcription termination . <i>Nucleic Acids Research</i> , 41(4):2541–2551, 12 2012. ISSN 0305-1048.
769 770	Sebastian Will, Tejal Joshi, Ivo L Hofacker, Peter F Stadler, and Rolf Backofen. Locarna-p: accurate boundary prediction and improved detection of structural rnas. <i>Rna</i> , 18(5):900–914, 2012.
771 772 773 774	Chengxin Zhang, Morgan Shine, Anna Marie Pyle, and Yang Zhang. Us-align: universal structure alignments of proteins, nucleic acids, and macromolecular complexes. <i>Nature methods</i> , 19(9): 1109–1115, 2022.
775 776 777	Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. <i>AI open</i> , 1:57–81, 2020.
778 779 780 781 782 783 784	Tianshuo Zhou, Ning Dai, Sizhen Li, Max Ward, David H Mathews, and Liang Huang. RNA de- sign via structure-aware multifrontier ensemble optimization. <i>Bioinformatics</i> , 39(Supplement_1): i563-i571, 06 2023. ISSN 1367-4811. doi: 10.1093/bioinformatics/btad252. URL https: //doi.org/10.1093/bioinformatics/btad252.
785 786 787 788	
789 790 791	
792 793 794 795	
796 797 798	
799 800 801	
802 803 804 805	
806 807 808	
808	

## 810 A MODEL DETAILS





Figure 7: Overview of RNAinformer encoder for constrained design and GC-Content conditioning with an adjacency matrix structure representation.

#### 

## **B** TRAINING DETAILS

**Loss** The problem of RNA design is often addressed by defining a structural loss function  $L_{\omega} = d(\omega, \mathcal{F}(\phi))$  that quantifies the difference between the target structure  $\omega$  and the folding,  $\mathcal{F}(\cdot)$ , of the designed candidate sequence  $\phi$  (Runge et al., 2019). However, the folding process is generally not differentiable making it difficult to use the structural loss for training in deep learning based approaches. Instead, we cast the problem as a conditional language modeling problem and train a conditional transformer language model on RNA sequences.

914The inverse folding problem can then be formulated as conditioning RNA sequences on the target915structures. The conditional probability for an RNA sequence  $\phi$  conditioned on a target structure  $\omega$ 916is,

$$p(\phi) = p(\phi \,|\, \omega) \qquad . \tag{2}$$

To design RNA sequences with a certain set of desired properties C we extend the inverse folding formulation and condition the sequence on both the target structure  $\omega$  and the desired properties C. Equation 2 is extended to, 

$$p(\phi) = p(\phi \mid \omega, C) \qquad . \tag{3}$$

Similar to the task of scaffold-based generation for small molecules (Bagal et al., 2021) but for both the RNA sequence and the structure, we can further extend the above formulation for con-strained design to include constraints on the target structure as well as the designed sequence at certain positions. By imposing these constraints on the sequence  $\phi$  and the structure  $\omega$  we get the masked sequence  $\phi$  and the masked structure  $\hat{\omega}$ , respectively, to condition the sequence generation on. Equation 3 then becomes,

$$p(\phi) = p(\phi \,|\, \hat{\omega}, \phi, C) \qquad . \tag{4}$$

Using auto-regressive modeling allows us to factorize the probability of the whole sequence  $p(\phi)$ into,

> $p(\phi) = \prod_{i=1}^{l} p(\phi_i \,|\, \phi_{< i}, \hat{\omega}, \hat{\phi}, C) \qquad ,$ (5)

where *l* is the length of the sequence  $\phi$ . 

This decomposes the design problem into a next token prediction problem. Now we can train a model with parameters  $\theta$  over a dataset  $\mathcal{D} = \{(\phi, \hat{\omega}, \hat{\phi}, C)\}^n$ , where  $n = |\mathcal{D}|$  using the loss,

$$\mathcal{L}_{\mathcal{D}} = \frac{1}{n} \sum_{k=1}^{n} \frac{1}{l^{k}} \sum_{i=1}^{l^{k}} l_{CE}(\phi_{i}^{k}, \theta(\phi_{$$

where  $L_{CE}(\psi_i, \phi_i)$  is the cross entropy loss between the target sequence and the designed sequence at position *i*. 

#### Table 2: Hyperparmeters for RNAinformer training.

Group	Parameter	Value
Trainer	Batch Size Training Steps	128 50,000
Optimizer	LR Weight Decay Betas	0.0005 0.1 0.9,0.98
LR Schedule	Schedule LR Decay Factor Warmup Steps	Cosine Annealing 0.1 1,000
Model	Model dim Layers Num Head FeedForward factor FeedForward kernel Dropout	256 6 4 3 0.1

#### DATASETS С

Table 3: Overview of the SynNested dataset.

968					
969	Set	#Samples	Avg Length	Pseudoknots	Multiplets
070	Train	444766	100	0(0.00%)	0(0.00%)
970	Valid	1108	100	0(0.00%)	0(0.00%)
971	Test	2722	96	0(0.00%)	0(0.00%)



Figure 8: Representations of RNA secondary structures. (Left) Common graph representation of the RNA. (Middle) Dot-bracket notation in the graph structure. A pair of nucleotides is indicated by a pair of matching brackets, unpaired nucleotides are indicated by a dot. (Right) Matrix representation of the RNA. The matrix is a binary  $L \times L$  square matrix, where L is the sequence length of the RNA. Pairing nucleotides are shown in yellow.

Table 4:	Overview	of the Sy	vnPseudoknot	dataset.
14010 11	0,01,10,0		, m beaution	aucubec

Set	#Samples	Avg Length	Pseudoknots	Multiplets
Train	444768	100	19.82%	0.00%
Valid	1108	100	20.31%	0.00%
Test	2732	96	29.61%	0.00%

#### Table 5: Overview of the SynMultiplet dataset.

Set	#Samples	Avg Length	Pseudoknots	Multiplets	Non Canonical BP
Train	441028	100	44.92%	57.44%	63.85%
Valid	1101	101	46.91%	61.25%	68.14%
Test	2721	95	45.61%	57.00%	62.27%

#### Table 6: Overview of the PDB Test sets.

Set	#Samples	Avg Length	Pseudoknots	Multiplets	Non-Canonical BP
TS1	67	74	83.58%	79.10%	92.54%
TS2	39	52	66.67%	74.36%	97.44%
TS3	19	79	94.74%	94.74%	94.74%
TS-Hard	28	66	71.43%	75.00%	85.71%

Table 7: Overview of the Rfam Constrained Design Dataset.

Set	#Samples	Avg Length	Pseudoknots	Multiplets
Train	51063	73	0.00%	0.00%
1	49	72	0.00%	0.00%

	Inverse Folding			Targe	ntent			
		0.30	0.35	0.40	0.45	0.50	0.55	0.60
#Tasks	1440	205	1275	1440	1440	1436	1220	364

#### Table 9: Overview of the Inter-family Dataset.

Set	#Samples	Avg Length	Pseudoknots	Non-Canonical BP	Multiplets
Train	19540	73	2047(10.47%)	11114(56.70%)	1330(6.80%)
Valid	494	77	12(2.43%)	287(57.86%)	13(2.63%)
TS1	54	61	43(79.62%)	49(90.74%)	40(74.07%)
TS2	36	45	23(63.88%)	35(97.22%)	26(72.22%)
TS3	16	67	15(93.75%)	15(93.75%)	15(93.75%)
TS-Hard	25	55	17(68.00%)	21(84.0%)	18(72.00%)

1044 D EVALUATION

# 1046 D.1 NOTATION

**Task** We call designing sequences for a particular target structure a task. The task may also have additional constraints for design with desired properties and constrained design.

**Solved Task** If a task has at least one designed sequence that folds back into the target structure and satisfies the other constraints of the task, then the task is considered to be solved.

1054 Candidate Sequence All the designed RNA sequences for a particular task are considered as its candidate sequences.

**Valid Sequence** Candidate sequences that solve a task are considered valid sequences for the task.

Valid Structure If a candidate sequence for a task folds back into the target structure or satisfies the constraints on the structure, then it has a valid structure.

1062 D.2 METRICS

All metrics for the RNAinformer are reported as the mean and standard deviation of three random seed runs.

Solved As the main performance measure of the model we report the percent of solved tasks for a given benchmark dataset.

Valid Sequences (Valid Seq.) We refer any candidate sequence that solves a task as a valid sequence. We measure the efficiency of the generative process by the number of valid sequences that are produced for each task.

$$ValidSequences = \frac{\#ValidSequences}{\#CandidateSequences}$$
(7)

1075 1076 **Diversity (Div.)** To measure the diversity of the valid sequences generated for a target structure, 1077 we use the pairwise Hamming distance. For N valid sequences of length l the diversity is defined 108,

1078  
1079 
$$Diversity = \frac{1}{N} \sum_{i}^{N} \sum_{j}^{N} \frac{1}{l} \sum_{k=1}^{l} H(S_{ik}, S_{jk}) \quad , \quad (8)$$

1027 1028

1029 1030 1031

1032 1033 1034

1035 1036 1037

1050

1053

1061

1072

1073 1074

where  $H(S_{ik}, S_{jk})$  describes the positional Hamming distance:

$$H(S_{ik}, S_{jk}) = \begin{cases} 0 & \text{if } S_{ik} = S_{jk} \\ 1 & \text{else} \end{cases}$$

$$\tag{9}$$

NC To measure the models ability to design with non-canonical base pair interactions we report the number of valid sequences containing non-canonical base pairs.

1088 GC-Content Error (GCE) For design with desired GC-content we report the property constraint
 1089 violation of the candidate sequences with valid structures, given by:

$$GCE = abs(GC_{target} - GC_{Sequence})$$
(10)

where  $GC_{target}$  is the target GC-Content value and  $GC_{Sequence}$  is the GC-Content of a candidate sequence.

F1 Score The F1 Score is a commonly used performance measure to assess the quality of secondary structure prediction algorithms. It is based on the confusion matrix, which describes the number of true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN) when comparing a predicted structure to the ground truth. The F1 score is the harmonic mean of precision and sensitivity, defined as:

F

$$T = \frac{2 \cdot TP}{(2 \cdot TP + FP + FN)} \tag{11}$$

Matthews Correlation Coefficient (MCC) Compared to the F1 score that emphasizes on positives, the MCC is a more balanced measure (Chicco & Jurman, 2020). The MCC can be calculated as follows.

$$MCC = \frac{(TP \cdot TN) - (FP \cdot FN)}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}$$
(12)

For each task in a test set, we take the maximum F1 and MCC scores achieved by a candidate sequence and report the average over these values across three random seeds.

TM-score and RMSD We use US-align (Zhang et al., 2022) to get the TM-scores and RMSD value based on optimal structural alignment between the folded 3D structure of the designed sequences and the ground truth 3D structure from the PDB.

#### Ε ADDITIONAL RESULTS

#### E.1 RNA DESIGN FOR NESTED STRUCTURES

Table 10: Results for the design of RNAs for nested structures of the SynNested Dataset.

Model	Solved	Valid Seq.	Diversity	F1	Seq. Rec.
RNAinformer	$0.918 {\pm} 0.020$	$0.299 {\pm} 0.005$	$0.699 {\pm} 0.001$	$0.992 {\pm} 0.005$	$0.394{\pm}0.002$
LEARNA	0.639	0.372	0.457	0.989	0.371
Meta-LEARNA	0.368	0.244	0.747	0.975	0.366
Meta-LEARNA-Adapt	0.372	0.236	0.746	0.976	0.366
libLEARNA	0.772	0.740	0.733	0.987	0.370
SAMFEO	0.996	0.999	0.106	0.999	0.343

Table 11: Results for the design of RNAs for nested structures of the SynNested dataset with target GC content. 

Model	Solved	Valid Seq.	Diversity	F1	Seq. Rec.	GC Error
RNAinformer(GC)	0.696±0.025	$0.247 {\pm} 0.004$	$0.701 {\pm} 0.003$	0.993±0.002	0.395±0.002	$0.010{\pm}0.001$
libLEARNA	0.590	0.452	0.728	0.923	0.373	0.082

#### E.2 RNA DESIGN WITH PSEUDOKNOTS

Table 12: Results for the design of RNAs including pseudoknots on the SynPseudoknot dataset.

Model	Test Set	Solved	Valid Seq.	Diversity	F1	Seq. Rec.
RNAinformer	pK-free pK	0.934±0.016 <b>0.685±0.064</b>	$0.213 \pm 0.018$ $0.140 \pm 0.028$	$0.692 {\pm} 0.002$ $0.692 {\pm} 0.004$	0.995±0.002 <b>0.947±0.009</b>	0.392±0.001 <b>0.405±0.001</b>
RNAinformer-1	pK-free pK	0.670±0.025 0.391±0.045	-	-	$0.920 \pm 0.007$ $0.734 \pm 0.025$	-
antaRNA	pK	0.156	-	-	0.788	0.253

Table 13: Results for the design of RNAs including pseudoknots on the SynPseudoknot dataset with target GC content. 

Model	Test Set	Solved	Valid Seq.	Diversity	F1	Seq. Rec.	GC-Error
RNAinformer(GC)	pK-free pK	0.658±0.023 0.337±0.021	$0.147 {\pm} 0.009 \\ 0.105 {\pm} 0.003$	$0.691 {\pm} 0.002$ $0.694 {\pm} 0.002$	0.993±0.003 <b>0.937±0.015</b>	0.395±0.001 <b>0.405±0.001</b>	0.010±0.002 0.013±0.001
RNAinformer(GC)-1	pK-free pK	0.447±0.008 <b>0.177±0.016</b>	-	-	0.931±0.016 0.757±0.045	-	0.010±0.001 <b>0.010±0.001</b>
antaRNA	pК	0.012	-	-	0.747	0.257	0.063

E.3 AUTOMATED DESIGN OF THEOPHYLLINE RIBOSWITCHES.

Table 14: Results for the design of Riboswitches

1184					
1185	Model	Set	Solved	Valid Seq.	Diversity
1186	RNAinformer	Riboswitch Tasks	0.919±0.002	0.599±0.031	$0.182 \pm 0.003$
1187	libLEARNA	Riboswitch Tasks	0.915	0.557	0.190

## 

Target	Model	Solved	Valid Seq.	Diversity
0.3	RNAinformer(GC)	0.964±0.020	<b>0.403±0.043</b>	<b>0.126±0.001</b>
	libLEARNA	<b>0.990</b>	0.264	0.125
0.35	RNAinformer(GC)	0.898±0.039	<b>0.511±0.077</b>	0.128±0.004
	libLEARNA	<b>0.941</b>	0.421	0.142
0.4	RNAinformer(GC)	0.893±0.015	<b>0.486±0.048</b>	0.149±0.002
	libLEARNA	<b>0.913</b>	0.398	<b>0.177</b>
0.45	RNAinformer(GC)	0.899±0.023	<b>0.431±0.039</b>	0.145±0.002
	libLEARNA	<b>0.913</b>	0.388	0.189
0.5	RNAinformer(GC)	0.888±0.036	0.399±0.039	0.143±0.002
	libLEARNA	<b>0.894</b>	0.309	0.174
0.55	RNAinformer(GC)	<b>0.871±0.010</b>	<b>0.336±0.033</b>	0.128±0.009
	libLEARNA	0.672	0.205	0.148
0.6	RNAinformer(GC)	<b>0.600±0.083</b>	<b>0.193±0.014</b>	<b>0.128±0.010</b>
	libLEARNA	0.071	0.050	0.000

#### Table 15: Results for the design of Riboswitches with target GC-Content.

#### E.4 RNA DESIGN WITH ALL KINDS OF BASE INTERACTIONS

Table 16: Comparison with different baselines on the experimentally validated structures from PDB.

Test Set	RNAin	former	Ran	dom	GNN		
	F1	Seq. Rec	F1	Seq. Rec	F1	Seq. Rec	
TS1 TS2	$0.832 \pm 0.004$ 0.923 $\pm 0.004$	$0.461 \pm 0.003$ 0.479 ± 0.007	$0.223 \pm 0.006$ 0.349 \pm 0.014	$0.391 \pm 0.001$ 0.416 \pm 0.007	$0.221 \pm 0.017$ 0.315 \pm 0.04	$0.393 \pm 0.007$ 0.405 ± 0.004	
TS3	$0.925 \pm 0.004$ $0.866 \pm 0.010$	$0.479 \pm 0.007$ $0.454 \pm 0.001$	$0.349\pm0.014$ $0.211\pm0.004$	$0.384 \pm 0.009$	$0.313 \pm 0.04$ $0.223 \pm 0.044$	$0.392 \pm 0.004$	
TS-Hard SynMulitplet	$0.834{\pm}0.007$ $0.862{\pm}0.002$	$0.465{\pm}0.006 \\ 0.432{\pm}0.001$	$ \begin{array}{c c} 0.274 \pm 0.023 \\ 0.205 \pm 0.002 \end{array} $	$0.401 \pm 0.003$ $0.371 \pm 0.001$	0.283±0.016 -	0.407±0.012 -	

Table 17: Results for RNA design for SynMultiplet Dataset and experimentally validated structures
 from PDB using RNAinformer.

Test Set	Solved	Valid Seq.	Diversity	F1	Seq. Rec.	NC
TS1	$0.119 {\pm} 0.015$	$0.135 {\pm} 0.026$	$0.590 {\pm} 0.021$	$0.832 {\pm} 0.004$	$0.461 {\pm} 0.003$	$0.546 {\pm} 0.038$
TS2	$0.214 \pm 0.015$	$0.150 \pm 0.035$	$0.575 \pm 0.011$	$0.923 \pm 0.004$	$0.479 \pm 0.007$	$0.751 \pm 0.006$
TS3	$0.035 \pm 0.030$	$0.013 \pm 0.015$	$0.191 \pm 0.330$	$0.866 \pm 0.010$	$0.454 {\pm} 0.001$	$0.333 {\pm} 0.577$
TS-Hard	$0.167 \pm 0.041$	$0.081 \pm 0.032$	$0.558 {\pm} 0.037$	$0.834 {\pm} 0.007$	$0.465 {\pm} 0.006$	$0.433 \pm 0.111$
SynMultiplet	$0.089 {\pm} 0.006$	$0.035 {\pm} 0.008$	$0.712 {\pm} 0.001$	$0.862 {\pm} 0.002$	$0.432 {\pm} 0.001$	$0.421 {\pm} 0.038$

Table 18: Results for RNA design for SynMultiplet Dataset and experimentally validated structures from PDB with target GC-content using RNAinformer.

1237	Test Set	Solved	Valid Seq.	Diversity	F1	Seq. Rec	NC	GC-Error
1238	TS1	$0.114 \pm 0.017$	$0.093 {\pm} 0.005$	$0.656 {\pm} 0.032$	$0.834 {\pm} 0.004$	$0.470 {\pm} 0.004$	$0.648 {\pm} 0.070$	$0.005 \pm 0.001$
1239	TS2	$0.188 {\pm} 0.039$	$0.131 {\pm} 0.037$	$0.569 {\pm} 0.035$	$0.922 {\pm} 0.002$	$0.489 {\pm} 0.002$	$0.797 {\pm} 0.056$	$0.007 \pm 0.001$
10/0	TS3	$0.018 {\pm} 0.030$	$0.007 \pm 0.012$	$0.176 {\pm} 0.305$	$0.858 {\pm} 0.002$	$0.462 {\pm} 0.007$	$0.167 {\pm} 0.289$	$0.004 \pm 0.001$
1240	TS-Hard	$0.179 {\pm} 0.036$	$0.059 {\pm} 0.008$	$0.576 {\pm} 0.063$	$0.842 {\pm} 0.007$	$0.478 {\pm} 0.017$	$0.492 {\pm} 0.138$	$0.004 \pm 0.001$
1241	SynMultiplet	$0.037 \pm 0.009$	$0.018 {\pm} 0.002$	$0.711 \pm 0.002$	$0.839 {\pm} 0.008$	$0.433 \pm 0.001$	$0.409 \pm 0.066$	$0.010 \pm 0.001$

Table 19: Comparison between different versions of the RNAinformer pre-trained on synthetic or real-world data with and without finetuning on the experimentally validated structures from PDB.
RNAinformerSyn refers to the RNAinformer model trained on synthetic data; RNAinformerSyn + FT refers to the model that was trained on synthetic data and finetuned with experimentally validated structures from PDB; RNAinformerNAT + FT refers to the model trained on existing (known)
RNA secondary structures from publicly available sources and finetuned on experimentally validated structures from PDB.

Test Set	RNAinfo	RNAinformerSyn			erSyn + FT	RNAinformerNAT + FT		
	F1	Seq. Rec		F1	Seq. Rec	F1	Seq. Rec	
TS1 TS2 TS3 TS-Hard	$\begin{array}{c} 0.832 {\pm} 0.004 \\ 0.923 {\pm} 0.004 \\ 0.866 {\pm} 0.010 \\ 0.834 {\pm} 0.007 \end{array}$	$\begin{array}{c} 0.461 {\pm} 0.003 \\ 0.479 {\pm} 0.007 \\ 0.454 {\pm} 0.001 \\ 0.465 {\pm} 0.006 \end{array}$		$\begin{array}{c} 0.861 {\pm} 0.003 \\ 0.924 {\pm} 0.003 \\ 0.877 {\pm} 0.005 \\ 0.842 {\pm} 0.007 \end{array}$	$\begin{array}{c} 0.520 {\pm} 0.005 \\ 0.537 {\pm} 0.005 \\ 0.500 {\pm} 0.015 \\ 0.503 {\pm} 0.010 \end{array}$	$\begin{array}{c} 0.346 {\pm} 0.022 \\ 0.459 {\pm} 0.021 \\ 0.324 {\pm} 0.004 \\ 0.364 {\pm} 0.036 \end{array}$	$\begin{array}{c} 0.427 {\pm} 0.005 \\ 0.444 {\pm} 0.008 \\ 0.413 {\pm} 0.014 \\ 0.412 {\pm} 0.004 \end{array}$	

Table 20: Comparison of RNAinformer with and without finetuning on experimentally validated structures from PDB evaluated on the CASP15 RNA data.

Test Set	RNAir	RNAinformer		RNAinformer + FT	
	F1	MCC		F1	MCC
CASP15	$0.901 {\pm} 0.005$	0.902±0.005		$0.877 {\pm} 0.021$	$0.878 {\pm} 0.021$

1271Table 21: Comparison between F1 Scores of designed sequences and PDB test set sequences using<br/>different folding algorithms.

Folding Algo.	Sequence	TS1	TS2	TS3	TS-Hard
RNAformer	Designed-Syn	<b>0.832</b>	<b>0.923</b>	<b>0.866</b>	<b>0.834</b>
	Designed-Nat	0.346	0.459	0.324	0.364
	PDB	0.716	0.797	0.709	0.641
SPOT-RNA	Designed-Syn	<b>0.719</b>	<b>0.824</b>	<b>0.731</b>	<b>0.677</b>
	Designed-Nat	0.304	0.436	0.260	0.298
	PDB	0.714	0.800	0.671	0.663
MXFold2	Designed-Syn	<b>0.689</b>	<b>0.792</b>	<b>0.739</b>	0.663
	Designed-Nat	0.269	0.406	0.197	0.242
	PDB	0.663	0.763	0.640	<b>0.667</b>
UFold	Designed-Syn	0.662	0.790	<b>0.682</b>	<b>0.628</b>
	Designed-Nat	0.256	0.410	0.219	0.246
	PDB	<b>0.673</b>	<b>0.892</b>	0.648	0.587

Table 22: Comparison of TM-scores of designed sequences and PDB test set sequences using Al phaFold3 for 3D structure predictions.

12	Folding Algo.	Sequence	TS1	TS2	TS3	TS-Hard	
	AlphaFold3	Designed(Avg) Designed(Best)	0.331	0.289	0.325	0.280	
		PDB	0.570	0.354	0.500	0.410	

\_

\_

\_

Table 23: Difference in TM-scores of designed sequences and PDB test set sequences using Al-phaFold3 for 3D structure predictions. We observe that RNAinformer predictions improve the TM score for orphan RNAs (where there is no MSA available for the ground truth sequence) but become worse for the sequences where MSA is available for the ground truth. Note that for all the designed sequences, AlphaFold did not find any MSA. 

Folding Algo.	MSA	TS1	TS2	TS3	TS-Hard
AlphaFold3	Orphan	0.008	-0.035	0.031	0.006
	Non-Orphan	-0.224	-0.033	-0.188	-0.197

Table 24: Evaluations of the designed sequences for the CASP15 data using AlphaFold 3 as the folding algorithm. We compare a version with and one without finetuning on known sequences (FT). 

Task	Id	FT	TM-score ↑	$\textbf{RMSD}\downarrow$
CPEB3_ribozyme (7QR4_1_B)	R1107	no yes	<b>0.439</b> 0.391	<b>3.103</b> 3.507
CPEB3 Ribozyme (7QR3_1_C)	R1108	no yes	<b>0.373</b> 0.275	<b>3.490</b> 3.533
CPEB3 Ribozyme (7QR3_1_D)	R1108	no yes	<b>0.311</b> 0.307	<b>3.250</b> 3.327
Cloverleaf RNA (8S95_1_C)	R1116	no yes	0.357 <b>0.395</b>	4.797 <b>4.050</b>
SARS-CoV-2 SL5 (8UYS_1_A)	R1149	no yes	0.325 <b>0.329</b>	<b>3.640</b> 4.003
BtCoV-HKU5 SL5 (8UYE_1_A)	R1156	no yes	<b>0.379</b> 0.312	<b>3.803</b> 3.867
BtCoV-HKU5 SL5 (8UYG_1_A)	R1156	no yes	0.326 <b>0.352</b>	<b>3.743</b> 4.293
BtCoV-HKU5 SL5 (8UYJ_1_A)	R1156	no yes	0.313 <b>0.314</b>	<b>3.833</b> 4.090
A-6B (7YR7_1_A)	R1189	no yes	<b>0.259</b> 0.240	<b>4.003</b> 4.167
A-4B (7YR6_1_A)	R1190	no yes	0.251 <b>0.270</b>	3.503 <b>2.873</b>