2

Deep Generative Optimization of mRNA Codon Sequences for Enhanced Protein Production and Therapeutic Efficacy

3

| 4 | Yupeng Li ^{1,#} , Fan | Wang ^{1,#} , Jiaqi Yang | ¹ , Zirong Han ^{2,3} , Linfen | ng Chen ¹ , Wenbing Jian | g ¹ , Hao Zhou ¹ , |
|---|--------------------------------|----------------------------------|---|-------------------------------------|--|
|---|--------------------------------|----------------------------------|---|-------------------------------------|--|

5 Tong Li¹, Zehua Tang¹, Jianxiang Deng¹, Xin He¹, Gaofeng Zha⁴, Jiekai Hu⁵, Yong Hu⁵, Linping Wu⁶,

6 Changyou Zhan⁷, Caijun Sun^{2,3,8,9}, Yao He^{1,*}, Zhi Xie^{1,*}

- 7 ¹State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-sen University,
- 8 Guangzhou, 510060, China.
- 9 ²School of Public Health (Shenzhen), Sun Yat-sen University, Shenzhen, 518107, China.
- 10 ³Shenzhen Key Laboratory of Pathogenic Microbes and Biosafety, Shenzhen Campus of Sun Yat-sen
- 11 University, Shenzhen, 518107, China.
- ⁴Scientific Research Center, The Seventh Affiliated Hospital. Sun Yat-sen University, Shenzhen, 518107,
 China.
- 14 ⁵ENO Bio mRNA Innovation Institute, Rhegen Biotechnology Co., Ltd. Shenzhen, China
- 15 ⁶Center for Chemical Biology and Drug Discovery, China -New Zealand Joint Laboratory of Biomedicine
- 16 and Health, Guangdong Provincial Key Laboratory of Biocomputing, Guangzhou Institutes of
- 17 Biomedicine and Health, Chinese Academy of Sciences, Guangzhou, 510530, China.
- 18 ⁷Department of Pharmacology, School of Basic Medical Sciences, Fudan University, Shanghai, 200032,
- 19 China.
- 20 ⁸Key Laboratory of Tropical Disease Control (Sun Yat-sen University), Ministry of Education;
- Cuangzhou, 514400, China.
- 22 ⁹State Key Laboratory of Anti-Infective Drug Discovery and Development, School of Pharmaceutical
- 23 Sciences, Sun Yat-sen University, Guangzhou, 510006, China.
- 24 [#]Equally contributed.
- 25 *Corresponding author. Email: xiezhi@gmail.com (Z.X.), scheyao@hotmail.com (Y. He)
- 26

27 ABSTRACT

Messenger RNA (mRNA) therapeutics show immense promise, but their efficacy is 28 29 limited by suboptimal protein expression. Here, we present RiboCode, a deep learning 30 framework that generates mRNA codon sequences for enhanced protein production. 31 RiboCode introduces several advances, including direct learning from large-scale 32 ribosome profiling data, context-aware mRNA optimization and generative exploration 33 of a large sequence space. In silico analysis demonstrate RiboCode's robust predictive 34 accuracy for unseen genes and cellular environments. In vitro experiments show 35 substantial improvements in protein expression, with up to a 72-fold increase, 36 significantly outperforming past methods. In addition, RiboCode achieves cell-type specific expression and demonstrates robust performance across different mRNA 37 38 formats, including m¹Ψ-modified and circular mRNAs, an important feature for mRNA therapeutics. In vivo mouse studies show that optimized influenza hemagglutinin 39 40 mRNAs induce ten times stronger neutralizing antibody responses against influenza 41 virus compared to the unoptimized sequence. In an optic nerve crush model, optimized 42 nerve growth factor mRNAs achieve equivalent neuroprotection of retinal ganglion 43 cells at one-fifth the dose of the unoptimized sequence. Collectively, RiboCode represents a paradigm shift from rule-based to data-driven, context-sensitive approach 44 45 for mRNA therapeutic applications, enabling the development of more potent and dose-46 efficient treatments.

48 INTRODUCTION

Messenger RNA (mRNA) therapy has emerged as a promising approach for treating diseases. This innovative therapeutic strategy harnesses the cell's protein synthesis machinery to produce desired proteins encoded by the delivered mRNA^{1–3}, leading to the application of mRNA therapies in various fields, such as vaccine development and protein replacement therapy⁴. The successful development and deployment of mRNA vaccines during the COVID-19 pandemic have further highlighted the transformative potential of this technology⁵.

Despite the remarkable progress in mRNA vaccines, achieving efficient and 56 57 consistent protein translation from delivered mRNA molecules remains a key challenge, particularly critical for protein replacement therapy where sustained, precise, and often 58 59 higher levels of protein expression are required in specific cellular contexts. However, the biological instability of mRNA and the complex regulatory mechanisms governing 60 61 mRNA translation in cells can lead to suboptimal protein expression^{6–8}. Therefore, improving the expression of mRNA is a key challenge for enhancing the therapeutic 62 63 efficacy and reducing the required dose of mRNA-based treatments.

64 An amino acid can be encoded by multiple synonymous codons, ranging from one to six codons per amino acid. Codon optimization is a strategy to improve protein 65 66 expression by changing the synonymous codon of an mRNA molecule while 67 maintaining the encoded amino acid sequence. The choice of synonymous codons can 68 largely impact the efficiency of mRNA translation and the stability of the mRNA molecule^{6,7}. For example, it has been shown that optimal codon usage can enhance 69 70 ribosome engagement and increase translation elongation rates, ultimately leading to higher protein production⁸. Additionally, codon choice can influence mRNA stability, 71 72 which is important as mRNAs are prone to degradation. The minimum free energy 73 (MFE) of an mRNA, a computational indicator of its secondary structure stability, is often used to assess this aspect. A lower (more negative) MFE indicates a more stable 74 secondary structure, as more energy would be required to disrupt the base pairing and 75

unfold the RNA⁶. Therefore, codon optimization is a critical step in the design of
mRNA-based therapies to achieve maximal protein production, leading to better
therapeutic efficacy.

79 Computational tools have been developed for codon optimization, most for DNAs, 80 employing various strategies to select optimal codons. Past methods rely on codon 81 usage bias derived from highly expressed genes in a given species, such as codon 82 adaptation index (CAI)⁹. These methods aim to mimic the codon usage patterns of 83 efficiently translated endogenous mRNAs. More recently, LinearDesign has been developed for mRNA optimization, by increasing CAI and reducing MFE, to jointly 84 optimize translation and mRNA stability⁶. LinearDesign uses a linear programming 85 86 approach to explore a wider space of sequence variants compared to previous methods 87 and showed superior performance over the previous codon optimization methods.

88 Despite the development of the previous methods, several limitations hinder their effectiveness in consistently improving the protein expression of mRNA molecules. 89 90 Firstly, the existing methods primarily rely on predefined sequence features, such as CAI, to guide codon selection. However, these metrics often fail to correlate with the 91 experimentally measured protein expression levels^{10,11}, indicating that they do not 92 93 accurately capture the complex factors governing mRNA translation. Secondly, the 94 existing methods do not adequately account for the activity of translational regulators that influence mRNA translation¹². This lack of context-aware optimization may reduce 95 96 the effectiveness of the optimized mRNA sequences in specific cellular environments. 97 Furthermore, the existing methods explore a limited space of codon sequences due to 98 computational constraints and the reliance on predefined rules. This restricted search 99 space may prevent the discovery of novel and highly optimized sequences that could 100 potentially yield significant improvements in protein expression.

Deep learning has achieved remarkable success in tasks such as image recognition,
 natural language processing, and protein structure prediction, where it has
 outperformed conventional algorithms by learning complex patterns and relationships

from vast amounts of data^{13,14}. In the context of mRNA codon optimization, a deep 104 105 learning approach may enable the model to capture the complex interplay between 106 codon usage and cellular context, without relying on predefined rules. Moreover, deep 107 learning models can explore a vast sequence space and discover novel patterns that may 108 not be apparent to human experts or accessible through traditional optimization methods¹⁵. This ability has been exemplified in the field of protein engineering, where 109 110 deep learning has been used to design novel protein sequences with improved stability, binding affinity, and catalytic activity¹⁶⁻¹⁸. For the codon optimization problem, if we 111 assume an average of three synonymous codons per amino acid, the number of possible 112 codon sequences for a 500 amino acid protein would 3^{500} , which is more than 10^{238} . 113 This number is much larger than the estimated number of atoms ($\sim 10^{80}$) in the 114 observable universe¹⁹. A deep learning-based codon optimization method could 115 116 potentially generate novel and highly optimized mRNA sequences by exploring the 117 immense sequence space.

118 Massive parallel reporter assays (MPRA) are commonly used to study the effects of regulatory sequences on gene expression²⁰. However, it is not suitable for optimizing 119 120 coding sequences due to the short sequence limitation, which is generally less than 300 121 base pairs, for high-throughput DNA synthesis. Additionally, MPRA experiments often 122 rely on artificial reporter constructs and may not fully recapitulate the complex 123 regulatory landscape of endogenous mRNA molecules. Ribosome profiling (Ribo-seq) 124 is a powerful experimental technique that provides a snapshot of actively translating 125 ribosomes on mRNA molecules, offering a direct and quantitative measurement of 126 translation at a transcriptome-wide scale^{21,22}. By training a deep learning model on Ribo-seq data from diverse codon sequences in diverse cell types and conditions, we 127 128 may capture the complex relationship between codon usage and translation which 129 occurs in the natural cellular context. However, such attempts have been lacking so far. In this study, we present RiboCode, a deep learning model for mRNA codon 130 131 optimization that enhances protein expression by directly learning complex relationship

132 of mRNA codon sequences to their translation level from large-scale Ribo-seq data. Our prediction model demonstrated robust performance, while analysis of RiboCode's 133 134 optimization strategies revealed a complex interplay between sequence characteristics 135 and translation. In vitro experiments showed up to a 72-fold increase in protein expression, significantly outperforming past methods. RiboCode also achieved cell-136 type specific expression, and maintained robust performance across unmodified, $m^{1}\Psi$ -137 modified, and circular mRNA formats. In vivo, optimized influenza virus 138 139 hemagglutinin (HA) mRNA induced approximately ten times stronger neutralizing antibody responses in mice, while optimized nerve growth factor (NGF) mRNA 140 achieved equivalent neuroprotection of retinal ganglion cells at one-fifth the dose in an 141 142 optic nerve crush mouse model. This data-driven approach to codon optimization 143 advances our understanding of mRNA translation and facilitates the development of 144 more effective mRNA therapeutics.

145 **RESULTS**

146 **RiboCode is a Deep Learning Framework for mRNA Codon Optimization**

RiboCode is a deep learning-based framework for optimizing mRNA codon sequences
to enhance protein production. It integrates three key components: a translation
prediction model, an MFE prediction model, and a codon optimizer that explores and
optimizes codon choices guided by the prediction models (Figure 1a).

151 The translation prediction model estimates the translation amount of a given codon 152 sequence by learning the translational expression of diverse mRNA sequences from Ribo-seq experiments (Figures 1b and S1). In contrast to previous tools that rely on 153 154 optimizing predefined features such as CAI, our deep learning model automatically 155 extracts relevant features by training on 320 paired Ribo-seq and RNA sequencing 156 (RNA-seq) datasets from 24 different human tissues and cell lines, encompassing translation measurements of over 10,000 mRNAs per dataset^{23,24}. In addition, the model 157 158 incorporates not only codon sequences but also mRNA abundances and cellular context that is presented by gene expression profiles from RNA-seq, via a multi-head attention
mechanism. This integrative approach enables the prediction of mRNA translation by
jointly considering these important factors influencing translation.

To address mRNA stability, we developed an MFE prediction model. Current MFE prediction tools, such as RNA fold²⁵ and Linearfold²⁶, use dynamic programming, which is non-differentiable and incompatible with our codon optimizer described below. Our MFE model employs a deep neural network architecture and undergoes an iterative optimization process, to simultaneously improve its predictive capability and optimize sequences for lower MFE values (Figure S2).

168 The codon optimizer of RiboCode begins with the original codon sequence of a 169 given protein (Figure 1c). The prediction models then predict a fitness score for this 170 sequence. Using a gradient ascent optimization approach based on activation maximization (AM)²⁷, the optimizer adjusts the codon distribution to maximize the 171 fitness score. A synonymous codon regularizer ensures that only synonymous codons 172 173 encoding the same amino acids as the original sequence are considered, preserving the protein's amino acid sequence. Through iterative cycles of sequence generation, 174 175 prediction, and optimization, the system produces codon sequences with improved 176 properties. RiboCode can optimizes mRNA translation, stability or both, by interfacing 177 with both the translation and MFE models, which uses a parameter w of 0 to optimize translation only, w of 1 to optimize MFE only and, a value between 0 to 1 to jointly 178 179 optimize both.

By combining data-driven predictions with high-throughput sequence generation, RiboCode overcomes limitations of conventional heuristic approaches. It enables the exploration of a vast mRNA codon space, potentially uncovering optimized sequences.

183 Evaluation of Translation Prediction Model

184 We first evaluated the RiboCode's performance and generalizability using three cross185 validation datasets: "new gene", "new environment", and "new gene in new

environment", which represented unseen genes, unseen cell types and unseen genes in unseen cell types during training (Figure S3). The model achieved a coefficient of determination (R^2) of 0.81, 0.89, and 0.81 for the three datasets, respectively (Figures 2a), indicating its robustness and ability to generalize.

To understand the relative importance of the three model inputs, we performed 190 191 ablation analysis, revealing that mRNA abundances were the most important 192 contributor to the prediction of translation (Figure 2b, Table S2), in agreement with an 193 early study of yeast translation found that the most predictive variable for translation was the mRNA expression of the gene²². The incorporation of codon sequences lifted 194 195 the R^2 by 0.15, and further inclusion of cellular environment improved the R^2 by 0.06. 196 The ablation analysis demonstrated that all the inputs contributed to predicting mRNA 197 translation.

We next investigated whether our model captured complex sequence features beyond common translation-related metrics. While our model learned relevant sequence features directly from the raw codon sequences, we tried to include common translation-related sequence metrics including CAI, MFE, and codon frequencies as additional model inputs and found these metrics did not improve prediction accuracy (Table S1). This suggested that the model could capture the sequence patterns that were predictive of translation, beyond these sequence metrics.

We explored alternative approaches to incorporating cellular context information. We directly incorporated the meta information of Ribo-seq datasets into the model, including cell types and experimental conditions and found it did not improve the performance (Table S1). This indicated that the gene expression profiles used in the model were an effective proxy to capture the relevant cellular environment influencing mRNA translation.

To validate the model's ability to capture cellular environment information, we leveraged a proteomics dataset from GTEx²⁸, which measured the relative protein expression of genes across 32 human tissues. We predicted the mRNA translation levels of 9,582 genes in these tissues. Notably, the predicted translation levels of 6,805 out of
9,582 genes across these tissues (71%) were significantly correlated with their
measured protein levels (p-value < 0.05 after False Discovery Rate (FDR) adjustment,
Pearson's correlation), with a median correlation value of 0.56 (Figure 2c).

To exclude random correlations, we conducted a control experiment to calculate the correlations with shuffled cellular environment profiles for 1,000 times. With randomized data, only 0.02% of genes showed significant correlations, compared to 71% with the correct environment. The median correlation coefficients with shuffled environments were near zero and significantly lower than those with the correct environment (p < 0.0005, one-sample t-test, Figure 2d). These results demonstrated that our model effectively captured cellular context information.

Finally, we investigated the positional importance of coding sequences in translation prediction. We analyzed the importance of each nucleotide position for the model's prediction. The results showed that the coding sequences close to the translation start site (TSS) were more important (Figure 2e). This is consistent with a general knowledge that codons near TSS had a greater impact on protein synthesis, by influencing translation initiation⁸.

Overall, the data-driven approach of RiboCode enabled robust predictive capabilities with biological relevance by learning important sequence patterns directly from the Ribo-seq data.

234 **RiboCode's Optimization Strategies for Enhanced mRNA Translation**

Having established the efficacy of our translation prediction model, we next explored how this model could be leveraged to generate sequences with enhanced translation potential. We first generated codon sequences of Gaussia luciferase (Gluc) (Figure S4). T-distributed stochastic neighbor embedding (t-SNE) indicated that the model established an association between the sequence space and translation levels. The red area in the upper right showed that a wide space of high translation sequences was explored (Figure 3a). We next explored how RiboCode optimized translation and 242 stability independently or jointly. A widely used Gluc sequence (MF882921.1) was used as a reference for comparison, which had a predicted translation level of 5.9 and 243 244 an MFE value of -216 (Figure 3b). By optimizing the sequence for translation (w=0), 245 the predicted translation level increased to around 25. On the other hand, codon 246 sequences optimized for MFE (w=1) reduced the MFE from -150 kcal/mol to around -247 350 kcal/mol, with a similar translation level to the reference. With joint optimization 248 $(0 \le w \le 1)$, RiboCode explored a wider sequence space, achieving both enhanced 249 translation and reduced MFE.

To understand RiboCode's optimization strategy, we analyzed codon usage patterns 250 251 between generated sequences with enhanced and reduced translation, as well as 252 between high- and low-translated endogenous sequences. We found that codons 253 preferentially used in highly translated endogenous sequences were also favored in 254 RiboCode-generated sequences with enhanced translation. Notably, the differences in 255 codon usage between RiboCode's enhanced and reduced translation sequences were 256 more pronounced than those observed between high- and low-translated endogenous 257 sequences (Figures 3c). To assess the generalizability of these findings, we extended 258 our analysis to multiple genes across various cell types. Consistently, we observed the 259 same pattern of biased codon usage in all the cases (Figure S5). This suggests that 260 RiboCode not only mimics but amplifies the codon usage patterns of efficiently translated endogenous mRNAs, potentially leading to even greater improvements in 261 262 translation.

We next examined how RiboCode utilized sequence features during generation and optimization. Analysis of sequence features across different mRNAs revealed complex and variable relationships with translation (Figures 3d and S6). Notably, highly translated mRNAs generally showed an increase in uridine content (U%), which may reduce secondary structure formation and facilitate smoother ribosome movement during translation⁷. Additionally, these mRNAs mostly exhibited a decrease in Effective Number of Codons (ENC), suggesting a selection against rare or inefficient codon pairs to enhance translation²⁹. Variations in CAI, Codon Pair Bias (CPB), GC content (GC%)
and MFE across different mRNAs suggested that while these features could influence
translation, their impact might be more mRNA- or context-dependent.

In short, these findings highlight RiboCode's ability to capture complex sequencetranslation relationships, offering a sophisticated approach to mRNA optimization that goes beyond traditional codon optimization metrics.

276 Experimental Validation Demonstrates RiboCode's Versatility and Efficacy

277 While our in silico analyses demonstrated the potential of RiboCode to optimize codon 278 sequences, we next sought to validate these findings experimentally. We first validated 279 RiboCode's ability to optimize codon sequences for enhanced protein expression. For 280 Gluc, protein expression levels of all RiboCode-optimized sequences showed a dramatic enhancement in protein production compared to the reference, with up to a 281 72-fold increase (RD2) (Figure 4a, Table S3), which also significantly outperformed 282 the LinearDesign-optimized sequences (p-value=0.019, one-sided Mann-Whitney U 283 284 test). The predicted translational levels showed a strong positive correlation with the experimentally measured protein levels, although the p-value was slightly above 0.05 285 (Correlation coefficient=0.71, p-value=0.077, Pearson's correlation, Figures 4b and S7). 286 In contrast, the CAI showed negative correlation with the experimental measurements 287 288 (Correlation coefficient=-0.15), indicating that CAI is not a reliable predictor of protein 289 expression levels in this context and may even lead to counterproductive optimization 290 strategies. We additionally optimized another commonly used reporter gene, firefly 291 luciferase (Fluc). The experimental validation showed a 17-fold increase in protein expression compared to the WT Fluc (Figure S8). 292

We next experimentally evaluated RiboCode's ability to design mRNAs with cell type specificity. We designed Gluc mRNA variants optimized for preferential expression in HEK293T cells and compared expression in HEK293T cells to both A549 and ARPE19 cells. For HEK293T vs A549, RiboCode predicted expression ratios ranging from 1.65 to 1.80, which closely matched experimental ratios of 1.35 to 1.73 (Figure 4c, Table S4). In the HEK293T vs ARPE19 comparison, while predicted ratios
were around 1.41 for all variants, experimental results showed higher ratios between
2.57 and 3.13 (Figure 4d, Table S5). All designed mRNAs consistently demonstrated
preferential expression in HEK293T cells across both comparisons. These results
showed RiboCode's capacity to capture cellular environment and design mRNAs with
enhanced expression in desired cell types.

304 Modified mRNAs, such as those with 1-methylpseudouridine $(m^{1}\Psi)$ modifications, 305 and circular RNAs are used in mRNA therapy instead of unmodified mRNAs due to their improved stability and reduced immunogenicity^{3,7,30}. We therefore assessed the 306 307 effectiveness of RiboCode in enhancing translation in these alternative mRNA forms. 308 Among the four codon variants, $m^{1}\Psi$ -modified RD2 and RD4 showed higher protein 309 expression levels compared to the reference, with up to a 4.6-fold higher expression at 48 hours post-transfection (Figure 4e, Table S3). Moreover, all four RiboCode-310 generated codon variants in the circular form outperformed the reference (Figure 4f, 311 312 Table S3). These results demonstrate that RiboCode optimization enhances protein production in both m¹Ψ-modified and circular mRNAs, illustrating its reliability and 313 314 versatility.

These experimental validations demonstrate RiboCode's ability to significantly enhance protein expression, optimize for specific cell types, and improve translation across various mRNA forms, highlighting its potential as a powerful tool for mRNA therapeutic development.

RiboCode Enhances Immunogenicity of mRNA-based Influenza Vaccines

Having established the robustness of our optimization approach, we next aimed to demonstrate its practical application in the development of mRNA-based vaccines. Influenza A viruses are responsible for causing respiratory infections, leading to annual epidemics that result in millions of human infections worldwide³¹. HA, a glycoprotein found on the surface of influenza A viruses, plays a crucial role in the viral infection process and is the primary target for the development of influenza vaccines. Although most of the vaccines were developed using inactivated influenza viruses, mRNA-based
 influenza vaccines are currently actively developed³².

328 To enhance the expression of HA and potentially improve the efficacy of HA-based 329 vaccines, we optimized the HA coding sequence. Two RiboCode-optimized HA 330 sequences showed enhanced in vitro protein expression compared to the WT (Figure 331 5a). Particularly, RD1 showed substantial improvement compared to the WT and 332 LinearDesign-optimized sequences. In addition, RD1 exhibited considerably higher 333 expression levels compared to the WT sequence in both $m^{1}\Psi$ -modified and circular mRNA forms (Figures 5b and 5c). These results again highlight the robustness and 334 335 versatility of the RiboCode-optimized sequence.

336 We further assessed the *in vivo* immunogenicity induced by the optimized sequence, 337 for both the prime and boost responses, where split virus influenza vaccine (SV) was 338 served as the positive control (Figure 5d). The RD1 sequence induced significantly stronger neutralizing antibody responses, measured by the micro-neutralization (MN) 339 340 titers, compared to the WT sequence and SV. For the prime response, RD1 elicited significantly higher MN titers compared to WT, with approximately 4.4-fold increase 341 342 (Figure 5e, mean MN titers: RD1=2,560, WT=580; p-value=0.008, one-sided Wilcoxon 343 test). The difference was more pronounced for the boost response, with RD1 inducing 344 a 9.6-fold increase in MN titers compared to WT (Figure 5e, mean MN titers: RD1=83,200, WT=8,640; p-value=0.002, one-sided Wilcoxon test). These results 345 346 demonstrated that the RiboCode-optimized sequence significantly enhanced both the initial and boosted immune responses. This dramatic improvement in immunogenicity 347 348 underscores RiboCode's potential to enable more effective vaccines with lower doses.

349 Enhanced Protein Expression and Therapeutic Efficacy with optimized NGF

350 **mRNA**

Having demonstrated the efficacy of RiboCode in optimizing mRNA for vaccinedevelopment, we next explored its potential in protein replacement therapy. We focused

353 on NGF as a promising candidate for treating glaucoma, which is a leading cause of 354 irreversible blindness³³ and causes death of retinal ganglion cells (RGCs). Our recent 355 study demonstrated that mRNA-based NGF therapy provided robust neuroprotection 356 for RGCs in an optic nerve crush (ONC) mouse model³⁴.

To improve the neuroprotection efficacy, we optimized the codon sequences of 357 358 human NGF mRNA. The protein expression levels of three RiboCode-designed 359 sequences were more than 2-fold higher compared to that of the WT while the 360 LinearDesign sequences did not show improvement (Figure 6a). We further assessed the best performed sequence (RD3) in both m¹ Ψ -modified mRNA and circular mRNA 361 forms. Notably, with $m^{1}\Psi$ -modification, RD3 achieved 8.4- and 9.8-fold higher protein 362 363 levels compared to the WT at 24h and 48h, respectively (Figure 6b). With mRNA 364 circulation, RD3 also achieved a more than 2-fold higher expression than the WT at 365 both 24h and 48h (Figure 6c). These results again demonstrated the robustness of 366 RiboCode-optimized sequences across different mRNA forms.

Based on its superior performance in initial *in vitro* tests, we selected RD3 for further *in vivo* studies. To evaluate the *in vivo* expression of optimized NGF mRNA, we intravitreally administered both the RD3 and WT sequences. Each mRNA was $m1\Psi$ -modified and encapsulated within LNP and administered at two doses: 100 ng/µl and 500 ng/µl. The RD3 sequence demonstrated significantly higher NGF protein expression compared to the WT sequence at both doses. Remarkably, RD3 at 100 ng/µl achieved even slightly higher expression than WT at 500 ng/µl (Figure 6d).

We then investigated the therapeutic potential of optimized NGF mRNA using an ONC mouse model, which mimics RGC injury and resulted in significant RGC loss (Figures 6e-g). Treatment with NGF mRNA showed clear neuroprotective effects, preserving more RGCs after injury. Notably, mice treated with 100 ng/µl RD3 showed significantly higher RGC counts than those treated with the same dose of WT mRNA (Figures 6h and 6i, p-value=0.0002, one-sided Wilcoxon test). Moreover, these counts were comparable to those in mice treated with 500 ng/µl WT mRNA. To sum, the optimized sequence exhibited superior protein expression both *in vitro* and *in vivo*, while maintaining therapeutic efficacy at one-fifth the dose of the unoptimized sequence. These results demonstrated the effectiveness of RiboCode in optimizing NGF mRNA for the treatment of RGC injury.

385 **DISCUSSION**

386 In this study, we present RiboCode, a novel deep learning-based framework for mRNA 387 codon optimization. The generative optimization framework, guided by the deep 388 learning prediction model, enables the efficient exploration of the immense space of 389 possible codon sequences. This allows RiboCode to discover novel, highly optimized sequences that may not be accessible to traditional optimization methods. RiboCode-390 optimized sequences demonstrate superior performance in various mRNA formats, 391 392 including unmodified, $m^{1}\Psi$ modified, and circular mRNAs, highlighting its broad 393 applicability in the rapidly evolving field of mRNA therapeutics. In vitro and in vivo 394 experiments using the optimized sequences of therapeutically relevant proteins show 395 substantial enhancements in protein expression compared to the unoptimized sequences. These improvements further translate into increased therapeutic efficacy, as 396 demonstrated by significantly enhanced immune responses to an optimized influenza 397 398 vaccine and markedly improved RGC protection in mice with optic nerve injury.

399 The superior performance of RiboCode can be attributed to several factor. Firstly, 400 RiboCode's deep learning model learns directly from diverse nature sequences with 401 translation measurements, enabling it to capture complex patterns of codon sequences for mRNAs with high translation level. In contrast, previous approaches, such as 402 403 LinearDesign, relies on limited sequence features, which may not fully capture the 404 intricacies of mRNA translation. Second, our model considered the cellular contexts of 405 mRNA translation whereas the previously codon optimization tools such as CAI-based methods or LinearDesign did not. Third, RiboCode's generative optimization 406 407 framework allows it to explore a large sequence space and to discover novel, highly optimized sequences that may not be accessible to the traditional approaches. 408

409 The findings of our study have important implications for the field of mRNA 410 therapeutics. Firstly, RiboCode can generate and evaluate a vast number of novel codon 411 combinations. This capability allows RiboCode to optimize mRNA sequences beyond 412 the limitations of evolutionary constraints, potentially uncovering more efficient codon usage patterns not found in natural transcripts. By transcending natural sequence 413 414 limitations, RiboCode represents a significant advancement in mRNA optimization, 415 potentially leading to levels of protein expression and therapeutic efficacy that surpass 416 what can be achieved with naturally evolved sequences. Second, by substantially 417 increasing protein production, the optimized sequences can improve the potency and 418 reduce the required dose of mRNA-based treatments, potentially mitigating side effects 419 and enhancing patient outcomes. This is particularly relevant for applications such as 420 protein replacement therapies, where achieving high levels of protein expression is 421 crucial for therapeutic success. Third, RiboCode's robustness and versatility across 422 different mRNA formats, including modified and circular mRNAs, expand the range of 423 therapeutic applications for which it can be employed. As the field of mRNA therapeutics continues to evolve and new mRNA formats are developed to enhance 424 stability, reduce immunogenicity, and improve delivery^{30,35}, RiboCode's ability to 425 426 optimize sequences for these diverse formats will be invaluable.

427 While our study demonstrates the significant potential of RiboCode in optimizing mRNA codon sequences for enhanced protein expression and therapeutic efficacy, there 428 are several limitations to address and future directions to explore. Firstly, we focused 429 430 exclusively on optimizing the codon sequences of mRNA molecules to enhance protein 431 expression in this study. While this approach yielded significant improvements in 432 protein production and therapeutic efficacy, our future work will expand upon this 433 foundation by jointly optimizing both the untranslated regions and codon sequences. 434 By addressing both coding and non-coding regions, we will develop a more comprehensive optimization strategy that maximizes the potential of mRNA-based 435 436 therapeutics across a wider range of applications and cellular contexts. Second,

437 although we attempted to incorporate mRNA structure-related features such as MFE into our model, we found no significant improvement in prediction accuracy. However, 438 439 this does not conclusively rule out the importance of mRNA structure in translation. 440 Future directions of RiboCode could benefit from more sophisticated integration of 441 structural information, such as local secondary structures or ribosome pause sites. 442 Finally, future studies could expand the validation to more advanced preclinical models, 443 optimizing a broader range of genes and proteins, incorporating additional optimization 444 objectives, and elucidating the underlying mechanisms.

In conclusion, RiboCode represents a paradigm shift from rule-based to data-driven mRNA optimization, potentially uncovering entirely new principles of efficient translation that were previously inaccessible. RiboCode will provide a versatile tool for researchers to maximize the potential of mRNA-based therapeutics, paving the way for more effective treatments in various medical applications.

450 DATA AVAILABILITY

451 The translation model and optimization framework of RiboCode are available on

452 GitHub (<u>https://github.com/wangfanfff/RiboCode</u>).

453 **ACKNOWLEDGMENTS**

We are grateful to the researchers who shared Ribo-seq & RNA-seq data and to the authors of RPFdb, whose contributions made this research possible. We are also immensely grateful to the Precision Medicine Center of Sun Yat-sen University for its long-term support.

458 AUTHORS' CONTRIBUTIONS

Z.X. conceived, designed and supervised the project. YP.L. collected and preprocessed
the datasets, evaluated the model performance and analyzed the data. Y. He, F.W. and
ZH.T. developed the translation deep model. H.Z. developed the MFE deep model. F.W.
and H.Z. developed the AM framework. JX.D. tested the model. JQ.Y. and LF.C.
prepared mRNAs and conducted *in vitro* experiments. WB.J. conducted *in vivo* NGF

464 experiments. ZR.H. and CJ.S. conducted the *in vivo* HA experiments. GF.Z
465 encapsulated mRNA with LNP. X.H., H.J., Y. Hu, LP.W., and CY.Z. contributed to
466 project discussion. YP.L., Z.X., Y. He and F.W. wrote the manuscript. All authors read
467 and approved the manuscript.

469 **REFERENCES**

- 470 1. Baden, L. R. *et al.* Efficacy and Safety of the mRNA-1273 SARS-CoV-2
 471 Vaccine. *N Engl J Med* **384**, 403–416 (2021).
- 472 2. Gebre, M. S. *et al.* Optimization of non-coding regions for a non-modified
 473 mRNA COVID-19 vaccine. *Nature* 601, 410–414 (2022).
- 474 3. Pardi, N., Hogan, M. J., Porter, F. W. & Weissman, D. mRNA vaccines a new
 475 era in vaccinology. *Nat Rev Drug Discov* **17**, 261–279 (2018).
- 476 4. Qin, S. *et al.* mRNA-based therapeutics: powerful and versatile tools to
- 477 combat diseases. *Signal Transduct Target Ther* **7**, 166 (2022).
- 478 5. Fang, E. *et al.* Advances in COVID-19 mRNA vaccine development. *Sig*479 *Transduct Target Ther* 7, 94 (2022).

480 6. Zhang, H. *et al.* Algorithm for optimized mRNA design improves stability and
481 immunogenicity. *Nature* 621, 396–403 (2023).

- 482 7. Leppek, K. et al. Combinatorial optimization of mRNA structure, stability,
- 483 and translation for RNA-based therapeutics. *Nat Commun* **13**, 1536 (2022).
- 484 8. Hanson, G. & Coller, J. Codon optimality, bias and usage in translation and
 485 mRNA decay. *Nat Rev Mol Cell Biol* **19**, 20–30 (2018).
- 486 9. Sharp, P. M. & Li, W. H. The codon Adaptation Index--a measure of
- 487 directional synonymous codon usage bias, and its potential applications.
- 488 *Nucleic Acids Res* **15**, 1281–1295 (1987).
- 489 10. Lu, P., Vogel, C., Wang, R., Yao, X. & Marcotte, E. M. Absolute protein
- 490 expression profiling estimates the relative contributions of transcriptional and
 491 translational regulation. *Nat Biotechnol* 25, 117–124 (2007).
- 492 11. Li, G.-W., Burkhardt, D., Gross, C. & Weissman, J. S. Quantifying absolute
- 493 protein synthesis rates reveals principles underlying allocation of cellular
 494 resources. *Cell* **157**, 624–635 (2014).
- 495 12. Waldman, Y. Y., Tuller, T., Shlomi, T., Sharan, R. & Ruppin, E. Translation
- 496 efficiency in humans: tissue specificity, global optimization and differences
- 497 between developmental stages. *Nucleic Acids Research* **38**, 2964–2974 (2010).
- 498 13. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold.
 499 **596**, 583–589 (2021).

- 500 14. Cao, C. *et al.* Deep Learning and Its Applications in Biomedicine. *Genomics*,
 501 *Proteomics & Bioinformatics* 16, 17–32 (2018).
- 502 15. Zrimec, J. et al. Deep learning suggests that gene expression is encoded in all
- 503 parts of a co-evolving interacting gene regulatory structure. *Nat Commun* **11**,
- 504 6141 (2020).
- 505 16. Dauparas, J. *et al.* Robust deep learning–based protein sequence design
 506 using ProteinMPNN. *Science* 378, 49–56 (2022).
- 507 17. Sumida, K. H. *et al*. Improving Protein Expression, Stability, and Function
- 508 with Protein MPNN. J. Am. Chem. Soc. **146**, 2054–2061 (2024).
- 18. Bennett, N. R. *et al.* Improving de novo protein binder design with deep
- 510 learning. *Nat Commun* **14**, 2625 (2023).
- 511 19. observable universe. https://en.wikipedia.org/wiki/Observable_universe.
- 512 20. Melnikov, A. et al. Systematic dissection and optimization of inducible
- 513 enhancers in human cells using a massively parallel reporter assay. *Nat*
- 514 Biotechnol **30**, 271–277 (2012).
- 515 21. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S.
- 516 Genome-wide analysis in vivo of translation with nucleotide resolution using
- 517 ribosome profiling. *Science* **324**, 218–223 (2009).
- 518 22. Weinberg, D. E. et al. Improved Ribosome-Footprint and mRNA
- 519 Measurements Provide Insights into Dynamics and Regulation of Yeast
- 520 Translation. *Cell Rep* **14**, 1787–1799 (2016).
- 521 23. Xie, S.-Q. *et al.* RPFdb: a database for genome wide information of translated
- 522 mRNA generated from ribosome profiling. *Nucleic Acids Res* 44, D254-258523 (2016).
- 524 24. Wang, H. *et al*. RPFdb v2.0: an updated database for genome-wide
- 525 information of translated mRNA generated from ribosome profiling. *Nucleic*
- 526 Acids Research **47**, D230–D234 (2019).
- 527 25. Lorenz, R. et al. ViennaRNA Package 2.0. Algorithms Mol Biol **6**, 26 (2011).
- 528 26. Huang, L. *et al*. LinearFold: linear-time approximate RNA folding by 5'-to-3'
- 529 dynamic programming and beam search. *Bioinformatics* **35**, i295–i304 (2019).

- 530 27. Linder, J. & Seelig, G. Fast activation maximization for molecular sequence
 531 design. *BMC Bioinformatics* 22, 510 (2021).
- 532 28. Jiang, L. *et al.* A Quantitative Proteome Map of the Human Body. *Cell* 183,
 533 269-283.e19 (2020).
- 534 29. Wright, F. The 'effective number of codons' used in a gene. *Gene* 87, 23–29535 (1990).
- 536 30. Chen, R. *et al.* Engineering circular RNA for enhanced protein production.
 537 Nat Biotechnol **41**, 262–272 (2023).
- 538 31. Molinari, N.-A. M. *et al.* The annual impact of seasonal influenza in the US:
 539 measuring disease burden and costs. *Vaccine* 25, 5086–5096 (2007).
- 540 32. Myers, M. L. et al. Commercial influenza vaccines vary in HA-complex
- structure and in induction of cross-reactive HA antibodies. *Nat Commun* 14,1763 (2023).
- 543 33. Lambiase, A. *et al.* Experimental and clinical evidence of neuroprotection by
- 544 nerve growth factor eye drops: Implications for glaucoma. *Proc Natl Acad Sci U S*545 *A* **106**, 13469–13474 (2009).
- 546 34. Jiang, W. *et al.* Circular RNA-based therapy provides sustained and robust
 547 neuroprotection for retinal ganglion cells. *Molecular Therapy Nucleic Acids* 35,
 548 102258 (2024).
- 549 35. Chaudhary, N., Weissman, D. & Whitehead, K. A. mRNA vaccines for
- 550 infectious diseases: principles, delivery and clinical translation. *Nat Rev Drug*
- 551 *Discov* **20**, 817–838 (2021).
- 36. Joshi NA & Fass JN. Sickle: A sliding-window, adaptive, quality-basedtrimming tool for FastQ files. (2011).
- 554 37. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2.
 555 Nat Methods 9, 357–359 (2012).
- 38. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 29,
 15–21 (2013).
- 558 39. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose
- program for assigning sequence reads to genomic features. *Bioinformatics* 30,
 923–930 (2014).

- 40. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq
 data with or without a reference genome. *BMC Bioinformatics* 12, 323 (2011).
- 41. Jia, L. & Qian, S.-B. Therapeutic mRNA Engineering from Head to Tail. Acc
 Chem Res 54, 4272–4282 (2021).
- 42. Ho, J. J. D. *et al*. A network of RNA-binding proteins controls translation
- 566 efficiency to activate anaerobic metabolism. *Nat Commun* **11**, 2677 (2020).
- 43. Luan, Y. *et al.* Deficiency of ribosomal proteins reshapes the transcriptional
 and translational landscape in human cells. *Nucleic Acids Research* 50, 6601–
 6617 (2022).
- 44. Liu, Y., Beyer, A. & Aebersold, R. On the Dependency of Cellular Protein
- 571 Levels on mRNA Abundance. *Cell* **165**, 535–550 (2016).
- 572 45. Schwanhäusser, B. *et al.* Global quantification of mammalian gene 573 expression control. *Nature* **473**, 337–342 (2011).
- 46. Zrimec, J., Buric, F., Kokina, M., Garcia, V. & Zelezniak, A. Learning the
- 575 Regulatory Code of Gene Expression. *Front. Mol. Biosci.* **8**, 673363 (2021).
- 47. Kelley, D. R., Snoek, J. & Rinn, J. L. Basset: learning the regulatory code of the
- 577 accessible genome with deep convolutional neural networks. *Genome Res.* **26**,
- 578 990–999 (2016).
- 48. loffe, S. & Szegedy, C. Batch Normalization: Accelerating Deep Network
- 580 Training by Reducing Internal Covariate Shift. in *International conference on* 581 *machine learning* (2015).
- 49. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R.
- 583 Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *The journal*
- 584 of machine learning research **15**, 1929–1958 (2014).
- 585 50. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep 586 convolutional neural networks. *Commun. ACM* **60**, 84–90 (2017).
- 587 51. Loshchilov, I. & Hutter, F. Decoupled Weight Decay Regularization. Preprint 588 at http://arxiv.org/abs/1711.05101 (2019).
- 589 52. Nair, V. & Hinton, G. E. Rectified Linear Units Improve Restricted Boltzmann
- 590 Machines. in *Proceedings of the 27th international conference on machine* 591 *learning (ICML-10)* (2010).
 - 22

- 53. He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image
- 593 Recognition. in 2016 IEEE Conference on Computer Vision and Pattern
- 594 *Recognition (CVPR)* 770–778 (2016).
- 595 54. Miyato, T., Dai, A. M. & Goodfellow, I. Adversarial Training Methods for Semi-
- 596 Supervised Text Classification. Preprint at http://arxiv.org/abs/1605.07725 597 (2021).
- 598 55. Maas, A. L., Hannun, A. Y. & Ng, A. Y. Rectifier nonlinearities improve neural 599 network acoustic models. in *Proc. icml* (2013).
- 56. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic Attribution for Deep Networks.
- 601 in International conference on machine learning 3319–3328 (2017).
- 602 57. Bach, S. *et al*. On Pixel-Wise Explanations for Non-Linear Classifier
- 603 Decisions by Layer-Wise Relevance Propagation. *PLoS ONE* **10**, e0130140
- 604 (2015).
- 58. Shrikumar, A., Greenside, P. & Kundaje, A. Learning Important Features
- 606 Through Propagating Activation Differences. in *International conference on*
- 607 machine learning 3145–3153 (2017).
- 59. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray
- expression data using empirical Bayes methods. *Biostatistics* **8**, 118–127 (2007).
- 610 60. Wang, Y. et al. L226Q Mutation on Influenza H7N9 Virus Hemagglutinin
- 611 Increases Receptor-Binding Avidity and Leads to Biased Antigenicity Evaluation. J
- 612 *Virol* **94**, e00667-20 (2020).



616 Figure 1. Predictive and Generative Optimization of RiboCode.

a. RiboCode contains three main components, a codon optimizer, a translation prediction model andan MFE prediction model.

b. The framework of the prediction model for translation. The input includes codon sequences in
one-hot encoding, corresponding mRNA abundances and, the cellular environment which is
presented by vectors of gene expression profiles from RNA-seq. The model integrates these inputs
using a multi-head attention mechanism. From the fused representations, a convolutional neural
network extracts features and outputs the predicted translation level of mRNA.

624 c. Iterative optimization of codon sequences. RiboCode predicts fitness of an original sequence

- 625 (T=0), then uses activation maximization to generate optimized synonymous variants (T+=1). A
- 626 synonymous regularizer maintains amino acid sequence. This process iterates until peak fitness is
- 627 achieved.
- 628



629

630 Figure 2. Evaluation of the Translation Prediction Model.

631 a. Experimentally measured translation levels by Ribo-seq versus predicted translation levels in the632 three validation datasets. Red lines denote the linear fit.

b. Ablation analysis shows the contributions of the three inputs to the prediction model. The table

below shows the ablation status of the inputs, with dots and crosses representing the presence and

- absence of corresponding elements, respectively.
- 636 c. Comparison of protein measurements across 32 human tissues by mass spectrum and predicted
- translation levels from RiboCode. We calculated the Pearson's correlation coefficient for all 9,582
- 638 genes, of which 6,805 (71%) showed significant correlations (p-value < 0.05 after FDR adjustment
- 639 for multiple testing). Red vertical line indicates the median of correlation coefficient. Light blue
- 640 represents a significant correlation. Darker blue represents insignificance.

- d. Median of Pearson's correlation coefficients between protein level and predicted translation level
- 642 with randomized environment input for 1,000 simulations. Only 0.02% of results were significant
- 643 (p-value < 0.05 after FDR adjustment for multiple-testing).
- e. The importance of each nucleotide position for the translation prediction. The x-axis represents
- 645 the nucleotide position from the TSS (translation starting site). Integrated Gradients attribution
- 646 method was used to obtain the importance score for each nucleotide position (black dots). The red
- 647 line denotes the local polynomial regression fit.
- 648
- 649



650

651 Figure 3. Strategies of Enhanced Translation in Generated Sequences.

a. Generation of Gluc codon sequences with low translation level (the upper left area) and high
translation level (the upper right area) (*w*=0, Figure S4). T-SNE of codon sequences is shown. Each
dot represents one sequence, and the color represents the predicted translation level.

b. Generation and optimization of Gluc codon sequences using different *w* of 0, 0.5, 0.7, and 1. Each

dot represents one sequence, positioned in its predicted translation level (y-axis) and MFE (x-axis).

657 The position of the reference sequence (MF882921.1) is shown.

658 c. Codons that appeared more frequently in highly translated endogenous sequences were also used
659 more often in highly translated Gluc sequences generated by RiboCode. "RD": RiboCode-generated
660 sequences, "Endo.": endogenous genes. (***: p<0.001, t-test).

d. Changes of sequence features of optimized sequences compared to the unoptimized, for different

662 mRNAs. For each column (feature), cells in red represent that the values were higher in optimized

- 663 sequences than unoptimized ones while cells in green are opposite. The difference of ENC for INS
- shows no significance (the cell in gray). Abbreviations: GFP: green fluorescent protein, Fluc: firefly
- 665 luciferase, INS: insulin, VZV: varicella zoster virus glycoprotein E, and HA: influenza A
- 666 hemagglutinin.



668 Figure 4. Robustness of Optimization across Unmodified, Modified, and Circular mRNAs.

669 a. Protein expression of Gluc was measured by fluorescence intensity. RD sequences were designed

by RiboCode. LD sequences were designed by LinearDesign. "RLU": relative light units.

b. Correlation of experimentally measured protein expression values of Gluc at 24 hours versuspredicted values of RiboCode and CAI (see Figure S7 for the details).

673 c. Relative protein expression levels of mRNA variants in HEK293T vs. ARPE19 cells. The ratio of674 larger than 1 indicates high protein expression level in HEK293T than that in ARPE19.

d. Relative protein expression levels of mRNA variants in HEK293T vs. A549 cells. The ratio oflarger than 1 indicates high protein expression level in HEK293T than that in A549.

677 e, f. Protein expression of generated Gluc codon variants in (e) the linear mRNA form with $m^{1}\Psi$ 678 modification and (f) the circular mRNA form.

679



682 Figure 5. More Effective mRNA-based Influenza Vaccines Through Codon Optimization

a. Western blot analysis shows protein expression of the HA variants in HEK293T cells 24 hours
after transfection. RD sequences were designed by RiboCode. LD sequences were designed by
LinearDesign.

b, c. Protein expression RD1 in (b) the linear mRNA form with m¹Ψ modification and (c) the circular
 mRNA form.

d. HA mRNA immunization and analysis: BALB/c mice were intramuscularly inoculated with two
doses (10µg mRNA for each dose) with an interval of two weeks. The mouse serum was collected
at 14 days and 28 days for MN assay.

- 691 e. Levels of neutralizing antibodies against influenza viruses after prime and boost vaccination.
- 692 IC50, half-maximal inhibitory concentration. PBS and split virus influenza vaccine (SV) were used
- as the negative and positive controls, respectively. One-sided Wilcox test was used to calculate p
- 694 values shown in the figure.
- 695
- 696



698 Figure 6. Enhanced Protein Expression and Therapeutic Efficacy with Optimized NGF

- 699 mRNA.
- **a.** NGF protein expression in HEK293T cells, measured by ELISA. RD sequences were designed
- 701 by RiboCode. LD sequences were designed by LinearDesign.
- **702 b, c.** Protein expression levels of RD3 in $m^1\Psi$ -modified (b) and circular (c) mRNA formats.
- **d.** *In vivo* protein expression of RD3. The m¹Ψ-modified mRNAs were injected into mouse retinas
- by intravitreal injection, and protein levels were measured by western blot 48 hours post-injection.
- **e.** Timeline of NGF mRNA therapy in ONC model: At Day 0 (D0), the $m^{1}\Psi$ -modified NGF mRNAs
- 706 were injected into mouse retinas by intravitreal injection. At Day 4 (D4), the optic nerve was
- subjected to a physical crush injury. At Day 18 (D18), RGC numbers were quantified using
- 708 immunofluorescence staining.

- f, g. RGC numbers measured by immunofluorescence staining in the ONC mouse were significantlyreduced compared to that of the control (one-sided Wilcox test).
- h. RGC numbers in the ONC mouse retina after injection of NGF m¹Ψ-modified mRNA with 100
 and 500 ng/µl dosages.
- i. The RGC number in mice treated with 100 ng/µl RD3 was similar to those treated with 500 ng/µl
- 714 WT and significantly higher than those treated with 100 ng/µl WT. One-sided Wilcoxon test was
- 715 used to calculate p-values shown in the figure; ns: p>0.05.
- 716
- 717

| 718 | Supplementary Materials |
|-----|--|
| 719 | |
| 720 | Deep Generative Optimization of mRNA Codon Sequences for |
| 721 | Enhanced Protein Production and Therapeutic Efficacy |
| 722 | |



Figure S1. Schematic diagram of the translation model structure. Multi-head CNN contains
four heads of CNN. Multi-head Attention is constructed by cell environment vector and transcript
abundances through the fully connected layers. CNN Encoder consists of two max pooling layers
and one convolutional layer.



730 Figure S2. MFE Model Architecture and Optimization Process

731 a. MFE Model Architecture: a shallow CNN with two convolutional layers, nine Resblocks, and

three fully connected layers.

- **b.** The MFE optimization contains four steps, including initial sampling, initial training, sequence
- 734 generation and model retraining.
- 735



Figure S3. Training and internal validation sets of the translation model. Schematic
representation of how the dataset was divided into training (90% of genes) and test (10% of genes)
sets. The 320 total datasets were split into 200 for training and 120 for validation, creating three
validation sets: "new gene", "new cellular environment", and "new gene in new cellular
environment".





745Figure S4. Generated Gluc codon variants. Plot showing the predicted translation levels of Gluc746codon variants generated using the translation model (w=0). The x-axis represents generation747iterations, with red and blue lines indicating enhanced and reduced predicted translation,748respectively. The green dashed line shows the predicted translation level of the reference Gluc749sequence.



Figure S5. Codon Usage Similarity Between Endogenous and RiboCode-Generated Sequences.
Comparison of codon usage patterns in endogenous high-translation sequences and low-translation
sequences, as well as the RiboCode-designed sequences with enhanced and reduced translation for
HA, and NGF in different cellular environments (HEK293T, A549, and HeLa). "Endo.":
endogenous. "RD": RiboCode-designed. (t-test, ****: p<0.0001).





Figure S6. Features of generated codon sequences for different genes. The graphs show various sequence features for optimized (red) and unoptimized (green) codon sequences for different genes (w=0). (t-test, ***: p<0.001). The full names of mRNAs are noted in the main text.



Figure S7. Correlation of experimental expression with predicted measures. The graphs show
 correlations between experimental protein expression (in log) and translation predicted by

766 RiboCode, and CAI. Pearson's correlation coefficients are provided.



Figure S8. Fluc protein expression in transfected HEK293T cells. Protein expression was
determined by fluorescence intensity. RD1 and RD2 were generated by RiboCode (*w*=0 and 0.5,
respectively). Compared to the wild-type (WT), RD2 expression increased by 14.5-fold at 24h, and
17.9-fold at 48h. "RLU": relative light unit.



Figure S9. Distribution of mRNA expression across all samples. The graph shows the distribution
of mRNA expression for 11,725 genes across all samples. Expression counts were transformed to
log(y×RPKM+1), where y was set to 5 to maximize the correlation between mRNA abundance and
translation level. RPKM stands for reads per kilobase of transcript per million reads mapped. The
default mRNA count for codon optimization was set to 4.5 based on the median value of this
distribution.

| Input | Training | Test set (R ²) | | | | |
|-------------------|-----------------------|----------------------------|----------|-----------------|--|--|
| Input | set (R ²) | New gene | New Env. | New gene & Env. | | |
| basic | 0.86021 | 0.81329 | 0.88719 | 0.80838 | | |
| - codon frequency | 0.86108 | 0.800528 | 0.883919 | 0.790282 | | |
| + MFE | 0.86086 | 0.809827 | 0.889674 | 0.800357 | | |
| + CAI | 0.86127 | 0.802417 | 0.880232 | 0.791394 | | |
| + cell type | 0.8609 | 0.808396 | 0.890336 | 0.799348 | | |
| + treatment | 0.86174 | 0.803956 | 0.885125 | 0.796533 | | |

| 782 | Table S1. Evaluation of translation prediction model on internal test sets. This table compares |
|-----|---|
| 783 | the performance (R^2) of the basic translation model with models incorporating additional inputs on |
| 784 | various test sets. The basic input includes codon sequence, mRNA abundance, and cellular |
| 785 | environment represented by gene expression profiles from RNA-seq. "Env." stands for cellular |
| 786 | environment. |

| Input | Training set (R ²) | Test set (R ²) |
|--|-----------------------------------|----------------------------|
| CDS + mRNA abundances + cellular environment | 0.85924 | 0.819247 |
| mRNA abundances + cellular environment | 0.7331 | 0.62351 |
| CDS + mRNA abundances | 0.80825 | 0.760404 |
| CDS + cellular environment | 0.58303 | 0.467216 |
| CDS | 0.56135 | 0.429161 |
| mRNA abundances | 0.70318 | 0.608484 |

788 Table S2. Ablation analysis of translation model inputs. This table shows the results of an

ablation analysis investigating the contribution of the three main inputs (codon sequences, mRNA

abundances, cellular environment) to the translation model's performance. "CDS": codon

791 sequence.

| Saguanas ID | Predicted | Linear mRNA | | Linear mRNA (m¹Ψ) | | Circular mRNA | | MEE | CAL |
|-------------|-----------|-------------|----------|-------------------|----------|---------------|----------|--------|-------|
| Sequence ID | level | 24h | 48h | 24h | 48h | 24h | 48h | MIFE | CAI |
| Reference | 5.88 | 12369.33 | 8546.333 | 344260.5 | 123505 | 73567.67 | 44620 | -216.4 | 0.808 |
| RD1 | 19.34 | 571477.7 | 352015.7 | 188385.8 | 69464.5 | 108143 | 68785.67 | -195.7 | 0.713 |
| RD2 | 19.75 | 851809 | 615058.3 | 551472.8 | 279046 | 81295 | 55305 | -195.3 | 0.703 |
| RD3 | 11.28 | 149905 | 92193.33 | 244131.3 | 126622.7 | 157691 | 92775.67 | -342.4 | 0.748 |
| RD4 | 18.85 | 241644 | 125983.3 | 680751.3 | 570459.5 | 145763.7 | 91233.33 | -258.4 | 0.736 |
| LD1 | 4.99 | 30501.67 | 13699.67 | | | | | -346.2 | 0.766 |
| LD2 | 5.29 | 400987.7 | 267094.3 | | | | | -302.5 | 0.952 |

Table S3. Predicted translation levels and experimental protein expression of Gluc mRNA variants. This table presents data for various Gluc mRNA constructs, including the reference sequence, RiboCode-generated variants (RD1-RD4 with different *w* values), and LinearDesigngenerated variants (LD1-LD2). For each construct, the table shows predicted translation levels, experimental protein expression in different mRNA formats (linear, m¹ Ψ -modified, and circular) at 24h and 48h post-transfection, as well as MFE and CAI values. Protein expression was measured by fluorescence intensity in HEK293T cells.

| Sequence ID | Protein expression level (24h) | | Relative expre | ssion level | Expression ratio of | |
|-------------|--------------------------------|----------|----------------|-------------|---------------------|--|
| | In HEK293T | In A549 | In HEK293T | In A549 | HEK239T/A549 | |
| Reference | 24527.5 | 28271.25 | 1 | 1 | 1 | |
| HEK_A549_1 | 38522.25 | 32671 | 1.57 | 1.16 | 1.35 | |
| HEK_A549_2 | 29584 | 19712.75 | 1.21 | 0.70 | 1.73 | |
| HEK_A549_3 | 23215.25 | 21544 | 1.29 | 0.76 | 1.70 | |

| 803 | Table S4. Experimental protein expression of Gluc mRNA variants with cellular differential |
|-----|--|
| 804 | expression. This table presents experimental protein expression for various Gluc mRNA construct, |
| 805 | including the reference sequence and RiboCode-designed variants. The designed variants were |
| 806 | predicted to express more in HEK293T than in A549 (HEK_A549_1 to 3, w=0). For each construct, |
| 807 | the table shows original protein expression level, expression level relative to the reference, and |
| 808 | expression ratio of HEK239T/A549. Protein expression was measured by fluorescence intensity. |
| 809 | |

| Sequence ID | Protein express | otein expression level (24h) | | ression level | Expression ratio of | |
|-------------|-----------------|------------------------------|------------|---------------|---------------------|--|
| | In HEK293T | In ARPE19 | In HEK293T | In AREP19 | HEK2391/ARPE19 | |
| Reference | 24527.5 | 106308.5 | 1 | 1 | 1 | |
| HEK_ARPE_1 | 23215.25 | 34516 | 0.95 | 0.32 | 2.97 | |
| HEK_ARPE_2 | 17301 | 23928.5 | 0.71 | 0.23 | 3.09 | |
| HEK_ARPE_3 | 31671.25 | 53428.25 | 1.29 | 0.50 | 2.58 | |

| 812 | Table S5. Experimental protein expression of Gluc mRNA variants with cellular differential |
|-----|---|
| 813 | expression. This table presents experimental protein expression for various Gluc mRNA construct, |
| 814 | including the reference sequence and RiboCode-designed variants. The designed variants were |
| 815 | predicted to express more in HEK293T than in ARPE19 (HEK_ARPE_1 to 3, w=0). For each |
| 816 | construct, the table shows original protein expression level, expression level relative to the reference, |
| 817 | and expression ratio of HEK239T/ARPE19. Protein expression was measured by fluorescence |
| 818 | intensity. |
| 819 | |