

# Deep generative optimization of mRNA codon sequences for enhanced mRNA translation and therapeutic efficacy

Received: 18 October 2024

Accepted: 30 September 2025

Published online: 12 November 2025



Yupeng Li<sup>1,10</sup>, Fan Wang<sup>1,10</sup>, Jiaqi Yang<sup>1</sup>, Zirong Han<sup>2,3</sup>, Linfeng Chen<sup>1</sup>, Wenbing Jiang<sup>1</sup>, Hao Zhou<sup>1</sup>, Tong Li<sup>1</sup>, Zehua Tang<sup>1</sup>, Jianxiang Deng<sup>1</sup>, Xin He<sup>1</sup>, Gaofeng Zha<sup>4</sup>, Zhaoyu Hu<sup>5</sup>, Yong Hu<sup>5</sup>, Linping Wu<sup>6</sup>, Changyou Zhan<sup>7</sup>, Caijun Sun<sup>2,3,8,9</sup>, Yao He<sup>1</sup>✉ & Zhi Xie<sup>1</sup>✉

Messenger RNA (mRNA) therapeutics show immense promise, but their efficacy is limited by suboptimal protein expression. Here, we present RiboDecode, a deep learning framework that generates mRNA codon sequences for enhanced mRNA translation. RiboDecode introduces several advances, including direct learning from large-scale ribosome profiling data and generative exploration of a large sequence space. In silico analysis demonstrates RiboDecode's robust predictive accuracy for unseen genes and cellular environments. In vitro experiments showed substantial improvements in protein expression, significantly outperforming past methods. In addition, RiboDecode enables mRNA design with consideration of cellular context and demonstrates robust performance across different mRNA formats, including m<sup>1</sup>Ψ-modified and circular mRNAs, an important feature for mRNA therapeutics. In vivo mouse studies showed that optimized influenza hemagglutinin mRNAs induce ten times stronger neutralizing antibody responses against influenza virus compared to the unoptimized sequence. In an optic nerve crush model, optimized nerve growth factor mRNAs achieve equivalent neuroprotection of retinal ganglion cells at one-fifth the dose of the unoptimized sequence. Collectively, RiboDecode represents a paradigm shift from rule-based to a data-driven, context-aware approach for mRNA therapeutic applications, enabling the development of more potent and dose-efficient treatments.

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<sup>1–3</sup>, leading to the application of mRNA therapies in various fields, such as vaccine development and protein replacement therapy<sup>4</sup>. The successful development and deployment of mRNA vaccines during the COVID-19 pandemic have further highlighted the transformative potential of this technology<sup>5</sup>.

Despite the remarkable progress in mRNA vaccines, achieving efficient and consistent protein translation from delivered mRNA molecules remains a key challenge, particularly critical for protein replacement therapy, where sustained, precise, and often higher levels of protein expression are required in specific cellular contexts. However, the biological instability of mRNA and the complex regulatory mechanisms governing mRNA translation in cells can lead to suboptimal protein expression<sup>6–8</sup>. Therefore, improving the expression of

mRNA is a key challenge for enhancing the therapeutic efficacy and reducing the required dose of mRNA-based treatments.

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 expression by changing the synonymous codon of an mRNA molecule while maintaining the encoded amino acid sequence. The choice of synonymous codons can largely impact the efficiency of mRNA translation and the stability of the mRNA molecule<sup>6,7</sup>. For example, it has been shown that optimal codon usage can enhance ribosome engagement and increase translation elongation rates, ultimately leading to higher protein production<sup>8</sup>. Additionally, codon choice can influence mRNA structure. Previous studies have demonstrated that mRNA structure critically influences its stability *in vivo*<sup>9</sup>, *in solution*<sup>10</sup>, and translation<sup>11</sup>. Therefore, codon optimization is a critical step in the design of mRNA-based therapies to achieve maximal protein production, leading to better therapeutic efficacy.

Computational tools have been developed for codon optimization, most of which were designed for DNA, employing various strategies to select optimal codons. Past methods rely on codon usage bias derived from highly expressed genes in a given species, such as codon adaptation index (CAI)<sup>12</sup>. These methods aim to mimic the codon usage patterns of efficiently translated endogenous mRNAs. More recently, LinearDesign<sup>6</sup> has been developed for mRNA optimization, aiming to jointly optimize translation and mRNA stability by increasing CAI and reducing minimum free energy (MFE)<sup>6</sup>, which is a computational metric for evaluating mRNA secondary structure. LinearDesign uses a linear programming approach to explore a wider space of sequence variants compared to previous methods and showed superior performance over the previous codon optimization methods. Additionally, other indices have been used to guide sequence optimization. For instance, higher GC content (GC%) has been associated with enhanced gene expression<sup>13</sup>.

Despite the development of the previous methods, several limitations hinder their effectiveness in consistently improving the protein expression of mRNA molecules. 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 experimentally measured protein expression levels<sup>14,15</sup>, indicating that they do not accurately capture the complex factors governing mRNA translation. Secondly, the existing methods do not adequately account for the activity of translational regulators that influence mRNA translation, such as translation factors and RNA-binding proteins<sup>16,17</sup>. This lack of context-aware optimization may reduce the effectiveness of the optimized mRNA sequences in specific cellular environments. Furthermore, the existing methods explore a limited space of codon sequences due to computational constraints and the reliance on predefined rules. This restricted search space may prevent the discovery of previously unexplored and highly optimized sequences that could 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<sup>18,19</sup>. In the context of mRNA codon optimization, a deep learning approach may enable the model to capture the complex interplay between codon usage and cellular context, without relying on predefined rules. Moreover, deep learning models can explore a vast sequence space and discover novel patterns that may not be apparent to human experts or accessible through traditional optimization methods<sup>20</sup>. This ability has been exemplified in the field of protein engineering, where deep learning has been used to design novel protein sequences with improved stability, binding affinity, and catalytic activity<sup>21–23</sup>. Recent advances in codon optimization research have seen the emergence of deep learning-based algorithms, particularly large

language models trained on cross-species nucleotide sequences. These models have been implemented for predictive modeling of mRNA translation efficiency and degradation kinetics<sup>24,25</sup>. Nevertheless, there persists an urgent demand for developing a rigorously validated optimization framework to improve mRNA-encoded protein expression specifically tailored for therapeutic applications.

Massive parallel reporter assays (MPRA) are commonly used to study the effects of regulatory sequences on gene expression<sup>26</sup>. However, it is not suitable for optimizing coding sequences due to the short sequence limitation, which is generally less than 300 base pairs, for high-throughput DNA synthesis. Additionally, MPRA experiments often rely on artificial reporter constructs and may not fully recapitulate the complex regulatory landscape of endogenous mRNA molecules. Ribosome profiling sequencing (Ribo-seq) is a powerful experimental technique that provides a snapshot of actively translating ribosomes on mRNA molecules<sup>27,28</sup>, where the translation level of an mRNA can be derived from the reads per kilobase per million (RPKM) of Ribo-seq. Recent studies have leveraged Ribo-seq to develop translation-focused deep learning models<sup>29–31</sup>. For instance, RiboNN predicts mRNA translation efficiency in mammalian cells by integrating mRNA sequences with ribosome profiling data, revealing translation-stability regulatory mechanisms<sup>32</sup>. However, the field critically requires a rational design framework that translates data-derived translational signatures into codon optimization strategies, enabling high-throughput exploration of sequence space and generation of optimized mRNA constructs for therapeutic development.

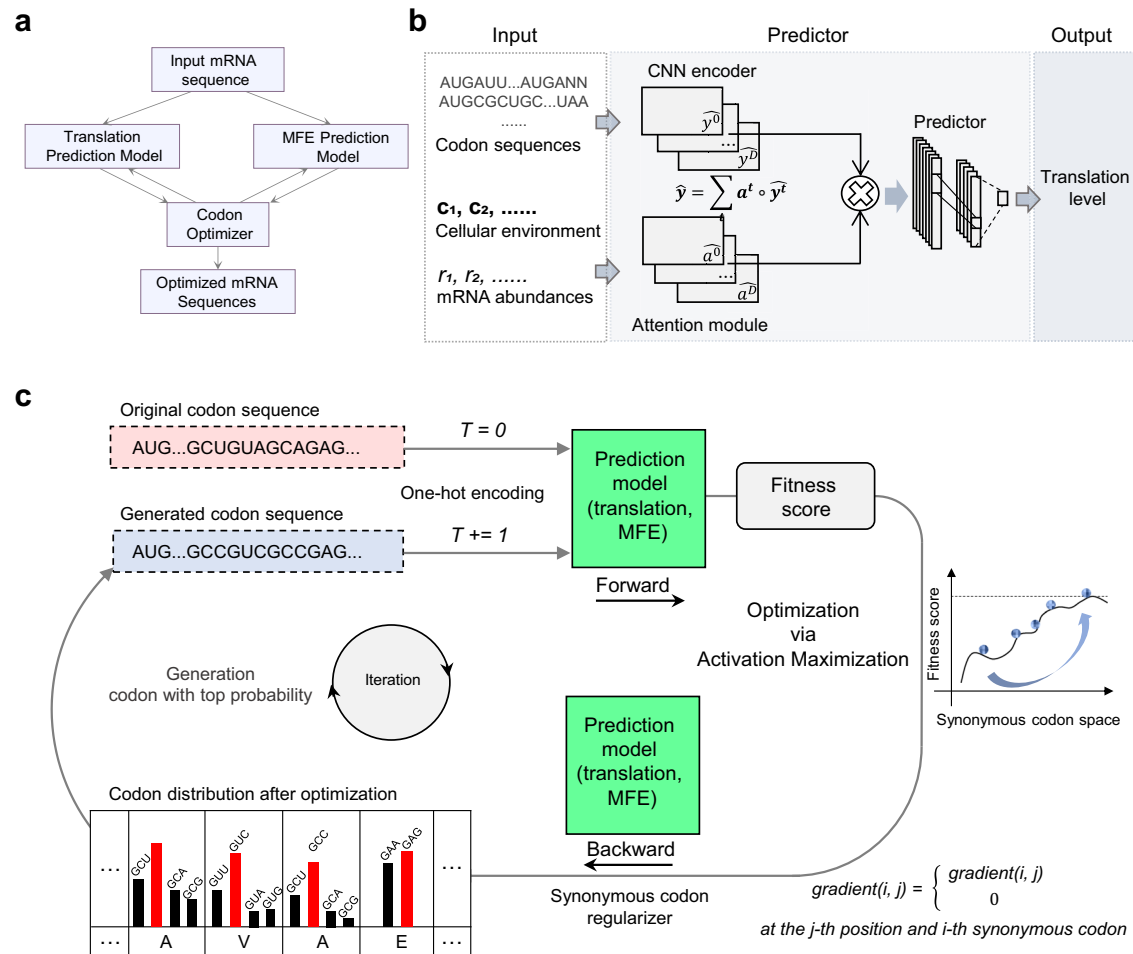
In this study, we present RiboDecode, a deep learning model for mRNA codon optimization that enhances mRNA translation by directly learning complex relationship of mRNA codon sequences to their translation level from large-scale Ribo-seq data. Our prediction model demonstrated robust performance, while analysis of RiboDecode's optimization strategies revealed a complex interplay between sequence characteristics and translation. *In vitro* experiments showed significantly increased in mRNA translation and protein expression, outperforming past methods. RiboDecode also considered cellular context, and maintained robust performance across unmodified, m<sup>1</sup>Ψ-modified, and circular mRNA formats. *In vivo*, optimized influenza virus hemagglutinin (HA) mRNA induced approximately ten times stronger neutralizing antibody responses in mice, while optimized nerve growth factor (NGF) mRNA achieved equivalent neuroprotection of retinal ganglion cells at one-fifth the dose in an optic nerve crush mouse model. This data-driven approach to codon optimization advances our understanding of mRNA translation and facilitates the development of more effective mRNA therapeutics.

## Results

### RiboDecode is a deep learning framework for mRNA codon optimization

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

The translation prediction model estimates the translation level of a given codon sequence by learning the translational expression of diverse mRNA sequences from Ribo-seq experiments (Figs. 1b and S1, “Methods”). In contrast to previous tools that rely on optimizing predefined features such as CAI, our deep learning model automatically extracts relevant features by training on 320 paired Ribo-seq and RNA sequencing (RNA-seq) datasets from 24 different human tissues and cell lines, encompassing translation measurements of over 10,000 mRNAs per dataset (Supplementary Data 1 and “Methods”) <sup>33,34</sup>. In addition, the model incorporates not only codon sequences but also mRNA abundances and cellular context that is presented by gene expression profiles from RNA-seq (“Methods”). This approach enables



**Fig. 1 | Predictive and generative optimization of RiboDecode.** **a** RiboDecode contains three main components, a codon optimizer, a translation prediction model and an 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. Multiple convolutional neural network (CNN) branches encode these inputs, and an attention module assigns weights to branch outputs, indicating their relative importance. The fused

representation is obtained as  $\hat{y} = \sum_t a^t \circ y^t$ , where  $a^t$  are attention weights and  $y^t$  are CNN features. From these representations, a CNN extracts features and outputs the predicted translation level of mRNA (see Methods). **c** Iterative optimization of codon sequences. RiboDecode predicts fitness of an original sequence ( $T = 0$ ), then uses activation maximization to generate optimized synonymous variants ( $T = +1$ ). A synonymous regularizer maintains amino acid sequence. This process iterates until peak fitness is achieved.

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 RNAfold<sup>35</sup> and Linearfold<sup>36</sup>, use dynamic programming. However, these methods are non-differentiable and thus 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 (Fig. S2, “Methods”).

The codon optimizer of RiboDecode begins with the original codon sequence of a given protein (Fig. 1c). The prediction models then predict a fitness score for this sequence. Using a gradient ascent optimization approach based on activation maximization (AM)<sup>37</sup>, the optimizer adjusts the codon distribution to maximize the fitness score (Fig. S3). A synonymous codon regularizer ensures that only synonymous codons encoding the same amino acids as the original sequence are considered, preserving the protein’s amino acid sequence. Through iterative cycles of sequence generation, prediction, and optimization, the system produces codon sequences with improved properties. RiboDecode can optimize mRNA translation,

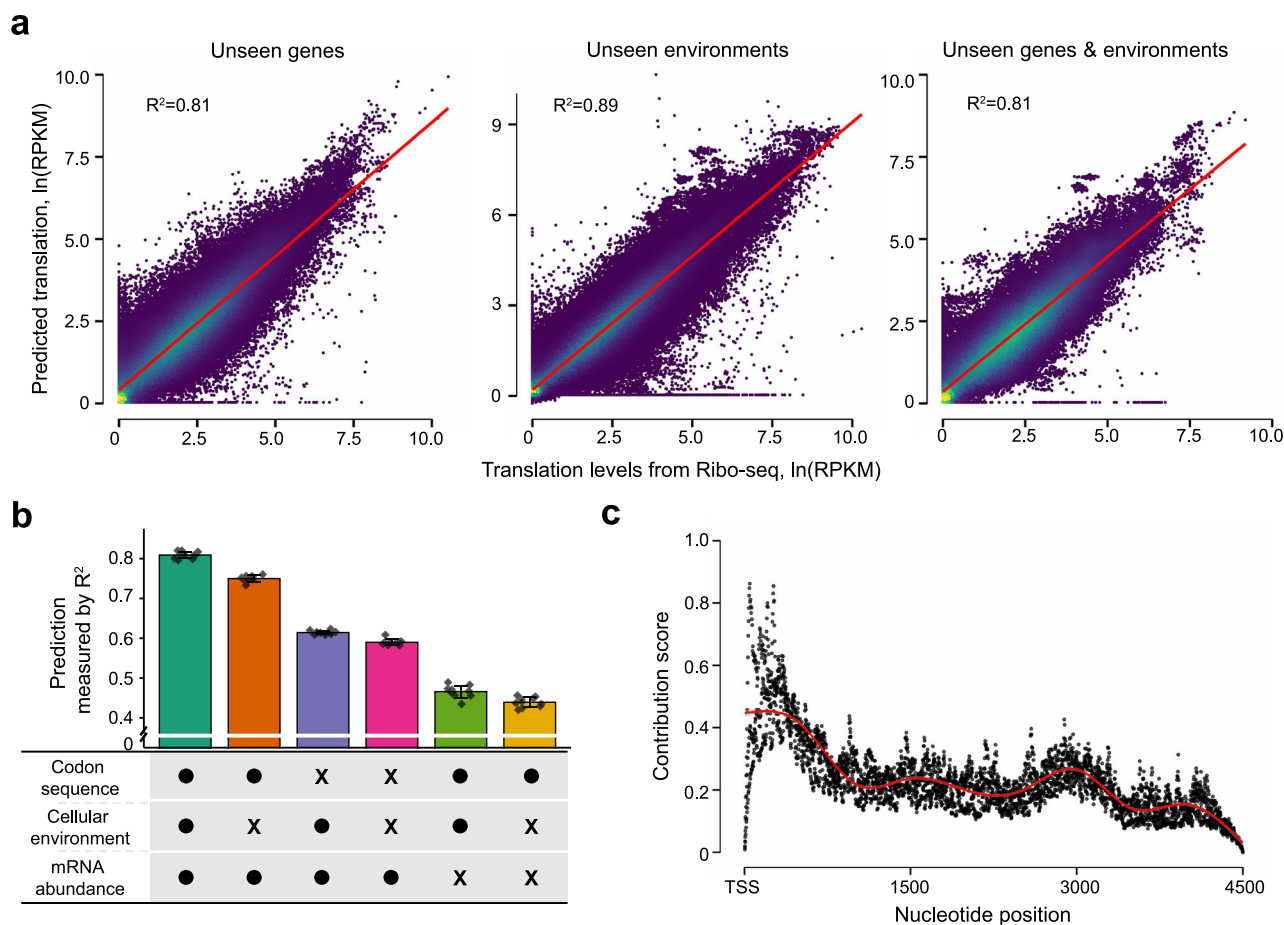
stability or both, by interfacing with both the translation and MFE models. This uses a parameter,  $w$ :  $w = 0$  optimizes translation only,  $w = 1$  optimizes MFE only, and a value  $0 < w < 1$  jointly optimizes both (“Methods”).

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

### Evaluation of translation prediction model

We first evaluated the RiboDecode’s performance and generalizability using three cross-validation datasets: “unseen genes”, “unseen environments”, and “unseen genes and environments”, which represented unseen genes and unseen cell types during training (Fig. S4, “Methods”). The model achieved a coefficient of determination ( $R^2$ ) of 0.81, 0.89, and 0.81 for the three datasets, respectively (Fig. 2a), indicating its robustness and ability to generalize.

To understand the relative importance of the three model inputs, we performed ablation analysis, revealing that mRNA abundances were the most important contributor to the prediction of translation (Fig. 2b and Table S2), in agreement with an early study of yeast



**Fig. 2 | Evaluation of the translation prediction model. a** Experimentally measured translation levels by Ribo-seq versus predicted translation levels in the three validation datasets. The lines denote the linear fit. The translation levels from Ribo-seq were ln-transformed (see “Methods”). **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. The points denote the  $R^2$  values

from the tenfold cross-validation ( $n=10$ ). Data are presented as mean values  $\pm$  SD. **c** The importance of each nucleotide position for the translation prediction. The x-axis represents the nucleotide position from the TSS (translation starting site). Integrated Gradients attribution method was used to obtain the importance score for each nucleotide position ( $n=2000$ ). The line denotes the local polynomial regression fit.

translation that found that the most predictive variable for translation was the mRNA expression of the gene<sup>28</sup>. The incorporation of codon sequences lifted the  $R^2$  by 0.15, and further inclusion of cellular environment improved the  $R^2$  by 0.06. The ablation analysis demonstrated that all the inputs contributed to predicting mRNA 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.

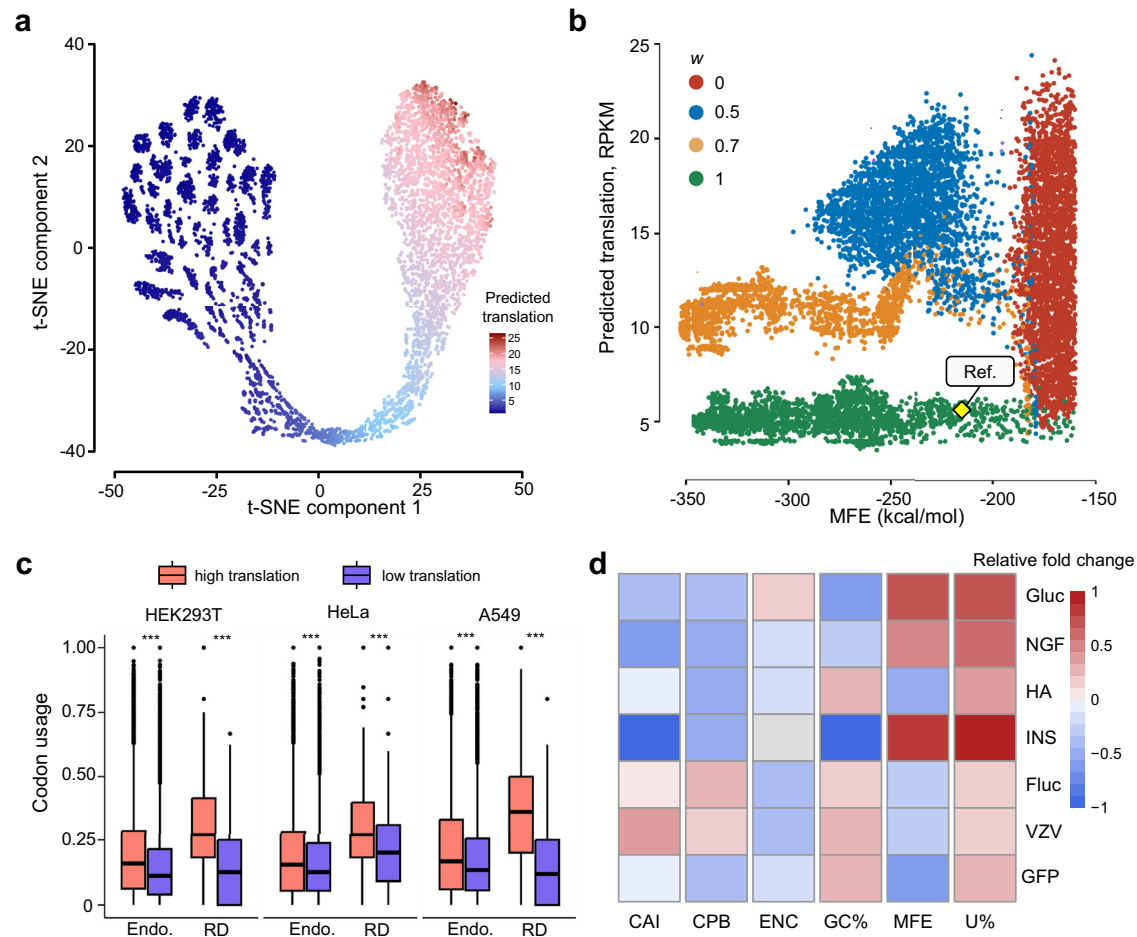
Finally, we investigated the positional importance of coding sequences in translation prediction. We analyzed the importance of each nucleotide position for the model’s prediction (“Methods”). The

results showed that the coding sequences close to the translation start site (TSS) were more important (Fig. 2c). This is consistent with a general knowledge that codons near TSS have a greater impact on protein synthesis, by influencing translation initiation<sup>8</sup>.

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

### RiboDecode’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 *Gussia luciferase* (Gluc) (Fig. S5). T-distributed stochastic neighbor embedding (t-SNE) indicated that the model established an association between the sequence space and translation levels (similar to Ribo-seq-derived RPKM values). The red area in the upper right showed that a wide space of high translation sequences was explored (Fig. 3a). We next explored how RiboDecode-optimized translation and stability independently or jointly. A widely used Gluc sequence was used as a reference for comparison, which had a predicted translation level of 5.9 and an MFE value of  $-216$  (Fig. 3b). By optimizing the sequence for translation ( $w=0$ ), the predicted translation level increased to around 25. On the other hand, codon



**Fig. 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) ( $w = 0$ , Fig. S5). 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). The position of the reference sequence (MF882921.1) is shown. **c** Codons that appeared more frequently in highly translated endogenous sequences were also used more often in highly translated Gluc sequences generated by RiboDecode. In the panels, “RD” denotes RiboDecode-generated sequences, with the dots representing the codon usage frequency of the top 10 most frequent codons from 1000 high- or low-translated

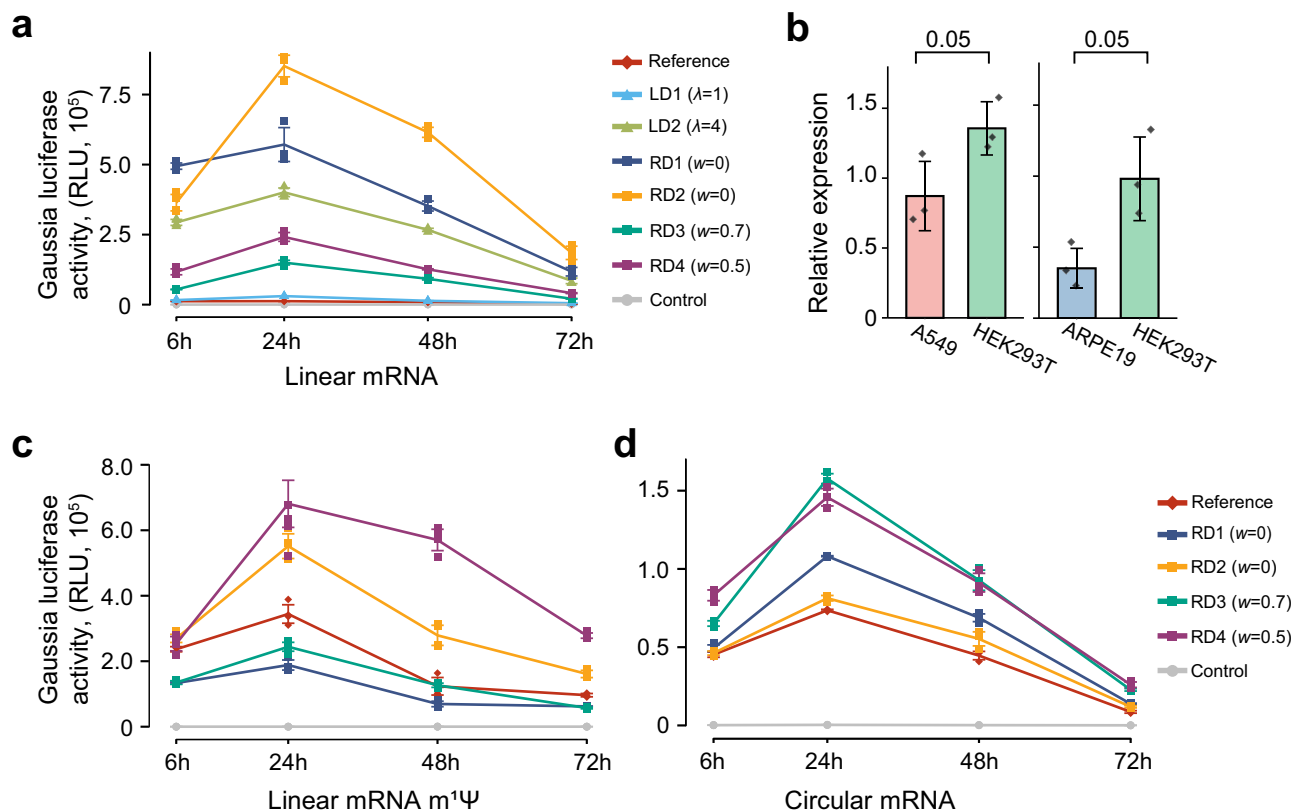
sequences. “Endo.” denotes endogenous genes, with the dots representing the codon usage frequency from endogenous sequences with top or bottom 10% translation level (see “Methods”). Boxes denote interquartile (IQR) ranges, centers mark medians and whiskers extend to 1.5 IQR from the quartiles. (for all the statistical test:  $p < 2.2 \times 10^{-16}$ , two-sided  $t$ -test). **d** Changes of sequence features of optimized sequences compared to the unoptimized, for different mRNAs. For each column (feature), a positive fold-change value indicates that the feature is more abundant (or higher) in optimized sequences compared to unoptimized ones, whereas a negative value signifies the opposite. The difference of ENC for INS shows no significance (the cell in gray). GFP green fluorescent protein, Fluc firefly luciferase, INS insulin, VZV varicella zoster virus glycoprotein E, and HA influenza A hemagglutinin.

sequences optimized for MFE ( $w = 1$ ) reduced the MFE from  $-150$  kcal/mol to around  $-350$  kcal/mol, with a similar translation level to the reference. With joint optimization ( $0 < w < 1$ ), RiboDecode explored a wider sequence space, achieving both enhanced translation and reduced MFE. Moreover, RiboDecode-generated sequences spanned a broader embedding space compared to those produced by Ribotree<sup>10</sup>, CDSfold<sup>38</sup>, and LinearDesign (Fig. S6, see “Methods”), suggesting enhanced sequence diversity. Furthermore, designing Gluc codons for different cell lines showed that the generated codons had distinct sequence patterns for different cellular contexts, reflecting differences in cellular environment (Fig. S7).

To understand RiboDecode’s optimization strategy, we analyzed codon usage patterns between generated sequences with enhanced and reduced translation, as well as between high- and low-translated endogenous sequences. We found that codons preferentially used in highly translated endogenous sequences were also favored in RiboDecode-generated sequences with enhanced translation. Notably,

the differences in codon usage between RiboDecode’s enhanced and reduced translation sequences were more pronounced than the differences observed in endogenous sequences (Fig. 3c, “Methods”). To assess the generalizability of these findings, we extended our analysis to multiple genes across various cell types. Consistently, we observed the same pattern of biased codon usage in all the cases (Fig. S8). This suggests that RiboDecode not only mimics but amplifies the codon usage patterns of efficiently translated endogenous mRNAs, potentially leading to even greater improvements in translation.

We next examined how RiboDecode utilized sequence features during generation and optimization. Analysis of sequence features across different mRNAs revealed complex and variable relationships with translation (Figs. 3d and S9, “Methods”). Notably, highly translated mRNAs generally showed an increase in uridine content (U%), which may reduce secondary structure formation and facilitate smoother ribosome movement during translation<sup>7</sup>. Additionally, these mRNAs mostly exhibited a decrease in Effective Number of Codons



**Fig. 4 | Robustness of optimization across unmodified, modified, and circular mRNAs.** **a** Protein expression of Gluc was measured by fluorescence intensity. RD sequences were designed by RiboDecode, with  $w$  parameter indicated in parentheses. LD sequences were designed by LinearDesign, with  $\lambda$  parameter indicated in parentheses. “RLU”: relative light units. Biological replicates were repeated three to four times (see Supplementary Data 10). Data are presented as mean values  $\pm$  SD. **b** Relative expression of experimentally measured protein expression values of mRNA variants designed by RiboDecode at 24 h in different

cells. One-sided Wilcoxon test was used to calculate  $p$ -values shown in the figure. Each dot shows the expression measurement of a designed mRNA. The expression of each mRNA was measured with four biological replicates. Data are presented as mean values  $\pm$  SD. **c, d** Protein expression of generated Gluc codon variants in **c** the linear mRNA form with  $m^1\Psi$  modification and **d** the circular mRNA form. Biological replicates were repeated three to four times (see Supplementary Data 10). Data are presented as mean values  $\pm$  SD.

(ENC), suggesting a selection against rare or inefficient codon pairs to enhance translation<sup>39</sup>. 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, due to complex sequence or structure feature changes.

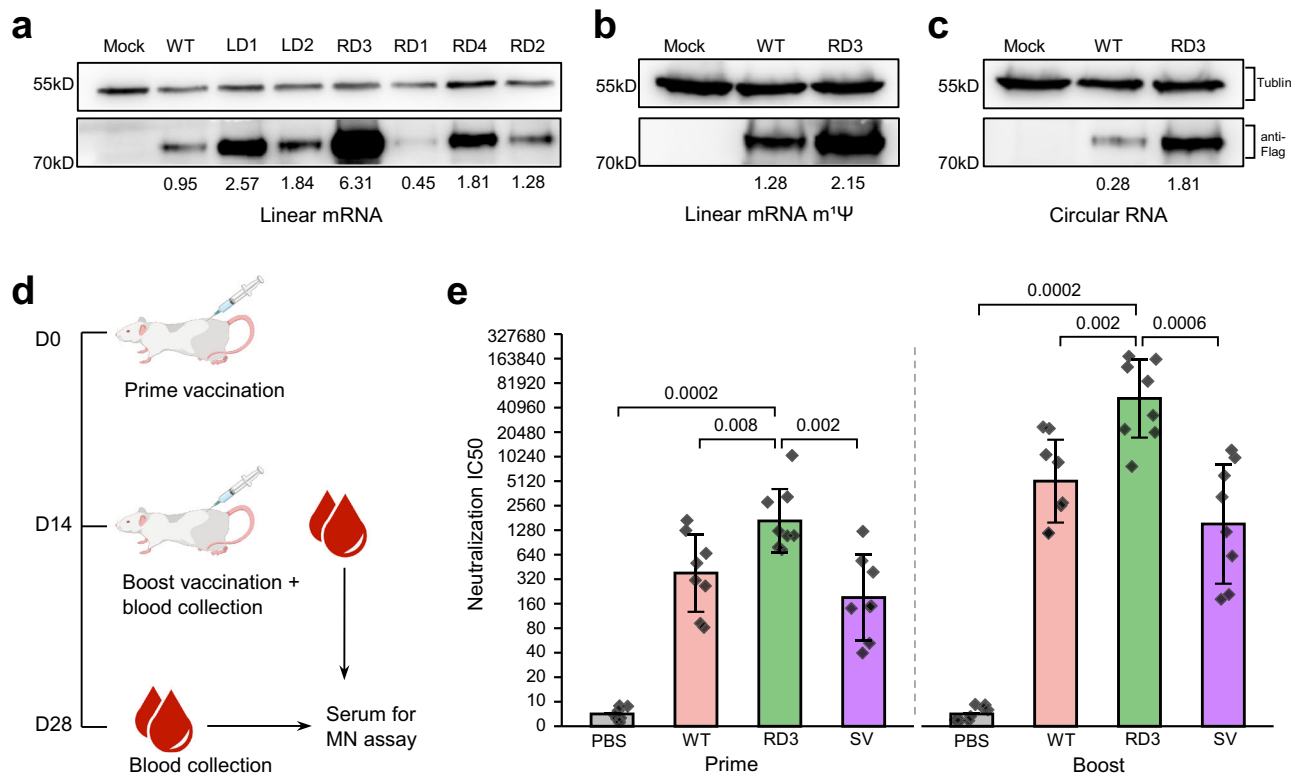
In short, these findings highlight RiboDecode’s ability to capture complex sequence-translation relationships, offering a sophisticated approach to mRNA optimization that goes beyond traditional codon optimization metrics.

#### Experimental validation demonstrates Ribodecode’s versatility and efficacy

While our *in silico* analyses demonstrated the potential of RiboDecode to optimize codon sequences, we next sought to validate these findings experimentally. We first validated RiboDecode’s ability to optimize codon sequences for enhanced protein expression. For Gluc, protein expression levels of the RiboDecode-optimized sequences surpassed the reference and were more than twice as high as that of the best-performing sequences designed by LinearDesign ( $p$ -value = 0.019, one-sided Mann–Whitney  $U$  test, Figs. 4a and S17a and Table S3). The predicted translational levels showed a positive correlation with the experimentally measured protein levels (Correlation coefficient = 0.71,  $p$ -value = 0.077, Pearson’s correlation, Fig. S10). However, the  $p$ -value is marginal, potentially due to the small number of constructs tested. Further validation with larger datasets is needed.

In contrast, the CAI showed negative correlation with the experimental measurements (Correlation coefficient =  $-0.15$ ), indicating that CAI is not a reliable predictor of protein expression levels and ineffective optimization strategy in this context. We noticed that the two RiboDecode-designed mRNAs (RD1 and RD2) had the best performance, which also had the highest MFE values of around  $-200$ . In contrast, the LinearDesign sequences had MFE values of  $-350$  and  $-300$ . To rule out the increased protein production was associated with higher MFE value, we tested additional LinearDesign sequences with higher MFEs ( $-266.8$  and  $-246.70$  kcal/mol, respectively). Compared to four LinearDesign sequences, RiboDecode-designed sequences outperformed (Fig. S11).

We additionally optimized another commonly used reporter gene, firefly luciferase (Fluc). Seven sequences, including four RiboDecode-optimized (RD1–RD4), two LinearDesign-optimized (LD1, LD2), and a WT, were transfected into HEK293T cells, and activity was measured over 72 h (Fig. S12 and Table S6). All optimized sequences significantly outperformed the WT. The LD1 (CAI of 0.766) and LD2 (CAI = 0.952) yielded the increased expression with about 7- and 41-fold changes, respectively. RiboDecode sequences (CAI of  $-0.71$ – $0.73$ ) also achieved substantial improvements with 6- to 16-fold over the WT. The superior performance of the LD2 sequence with a high CAI value here—contrasting with results for Gluc—underscores that the optimal codon strategy is gene-specific. This highlights the need for context-dependent approaches like RiboDecode that capture complex features beyond single metrics like CAI.



**Fig. 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 h after transfection. RD sequences were designed by RiboDecode (*w* parameter were set to 0, 0, 0.7, and 0.5 for RD1-4). LD sequences were designed by LinearDesign (*l* parameter were set to 0 and 4 for LD1 and LD2). The expression values were quantified using GelAnalyzer. Three times of the experiment was repeated independently with similar results. **b**, **c** Protein expression RD3 in **b** the linear mRNA form with m<sup>1</sup>Ψ modification and **c** the circular mRNA form. Three times each experiment was repeated independently with similar results. **d** HA

mRNA immunization and analysis: BALB/c mice were intramuscularly inoculated with two doses (10 μg mRNA for each dose) with an interval of two weeks. The mouse serum was collected at 14 days and 28 days for MN assay. Created with BioGDP.com<sup>51</sup>. **e** Levels of neutralizing antibodies against influenza viruses after prime and boost vaccination. IC<sub>50</sub>, half-maximal inhibitory concentration. PBS and split virus influenza vaccine (SV) were used as the negative and positive controls, respectively. One-sided Wilcoxon test was used to calculate *p*-values shown in the figure. Biological replicates were repeated eight times. Data are presented as mean values ± SD.

We next evaluated RiboDecode's ability to design for cellular context by optimizing Gluc mRNA for preferential expression in HEK293T cells over A549 and ARPE19 cells. The designed variants successfully exhibited the intended higher expression in HEK293T compared to both cell lines (Fig. 4b and Table S4). Predicted expression ratios favoring HEK293T closely matched experimental results in the comparison against A549 (-1.7-fold predicted vs. -1.55-fold experimental). Preferential expression was also achieved against ARPE19 (Fig. 4b and Table S5), although the experimental fold-change (-2.8-fold) doubled the prediction (-1.4-fold). These results confirm RiboDecode's capability for context-aware design, while the discrepancy in the ARPE19 comparison indicates potential for refining the model to better capture quantitative differences across diverse cellular environments.

Modified mRNAs, such as those with 1-methylpseudouridine (m<sup>1</sup>Ψ) modifications, and circular RNAs are used in mRNA therapy instead of unmodified mRNAs due to their improved stability and reduced immunogenicity<sup>3,7,40</sup>. We therefore assessed the effectiveness of RiboDecode in enhancing translation in these alternative mRNA forms. Among the four codon variants, m<sup>1</sup>Ψ-modified RD2 and RD4 showed higher protein expression levels compared to the reference, with up to a 4.6-fold higher expression at 48 h post-transfection (Figs. 4c and S17b). Moreover, all four RiboDecode-generated codon variants in the circular form outperformed the reference (Figs. 4d and S17c). These results demonstrate that RiboDecode optimization enhances protein production in both

m<sup>1</sup>Ψ-modified and circular mRNAs, illustrating its reliability and versatility.

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

### RiboDecode 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 worldwide<sup>41</sup>. 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 developed<sup>42</sup>.

To enhance the expression of HA and potentially improve the efficacy of HA-based vaccines, we optimized the HA coding sequence. Three out of four RiboDecode-optimized HA sequences and two LinearDesign-optimized sequences showed higher in vitro protein expression compared to the WT (Fig. 5a and Table S7). Particularly, RD3 showed approximately sixfold increase compared to the WT and

LinearDesign-optimized sequences. In addition, RD3 exhibited considerably higher expression levels compared to the WT sequence in both m<sup>1</sup>Ψ-modified and circular mRNA forms (Fig. 5b, c). These results again highlight the robustness and versatility of the RiboDecode-optimized sequence.

We further assessed the in vivo immunogenicity induced by the optimized sequence for both the prime and boost responses, where split virus influenza vaccine (SV) was served as the positive control (Fig. 5d, “Methods”). The RD3 sequence induced significantly stronger neutralizing antibody responses, measured by the micro-neutralization (MN) titers, compared to the WT sequence and SV. For the prime response, RD3 elicited significantly higher MN titers compared to WT, with approximately 4.4-fold increase (Fig. 5e, mean MN titers: RD3 = 2560, WT = 580; *p*-value = 0.008, one-sided Wilcoxon test). The difference was more pronounced for the boost response, with RD3 inducing a 9.6-fold increase in MN titers compared to WT (Fig. 5e, mean MN titers: RD3 = 83,200, WT = 8640; *p*-value = 0.002, one-sided Wilcoxon test). These results demonstrated that the RiboDecode-optimized sequence significantly enhanced both the initial and boosted immune responses. This dramatic improvement in immunogenicity underscores RiboDecode’s potential to enable more effective vaccines with lower doses.

### Enhanced protein expression and therapeutic efficacy with optimized NGF mRNA

Having demonstrated the efficacy of RiboDecode in optimizing mRNA for vaccine development, we next explored its potential in protein replacement therapy. We focused on NGF as a promising candidate for treating glaucoma, which is a leading cause of irreversible blindness<sup>43</sup> and causes death of retinal ganglion cells (RGCs). Our recent study demonstrated that mRNA-based NGF therapy provided robust neuroprotection for RGCs in an optic nerve crush (ONC) mouse model<sup>44</sup>.

To improve the neuroprotection efficacy, we optimized the codon sequences of human NGF mRNA. The protein expression levels of three RiboDecode-designed sequences were more than twofold higher compared to that of the WT, whereas the LinearDesign sequences did not show improvement (Figs. 6a and S18a and Table S8). We further assessed the best-performing sequence (RD3) in both m<sup>1</sup>Ψ-modified mRNA and circular mRNA forms. Notably, with m<sup>1</sup>Ψ-modification, RD3 achieved 8.4- and 9.8-fold higher protein levels compared to the WT at 24 h and 48 h, respectively (Figs. 6b and S18b). With mRNA circulation, RD3 also achieved a more than twofold higher expression than the WT at both 24 h and 48 h (Figs. 6c and S18c). These results again demonstrated the robustness of RiboDecode-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 m<sup>1</sup>Ψ-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 (Fig. 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 (Fig. 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 (Fig. 6h, *i*, *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 RiboDecode in optimizing NGF mRNA for the treatment of RGC injury.

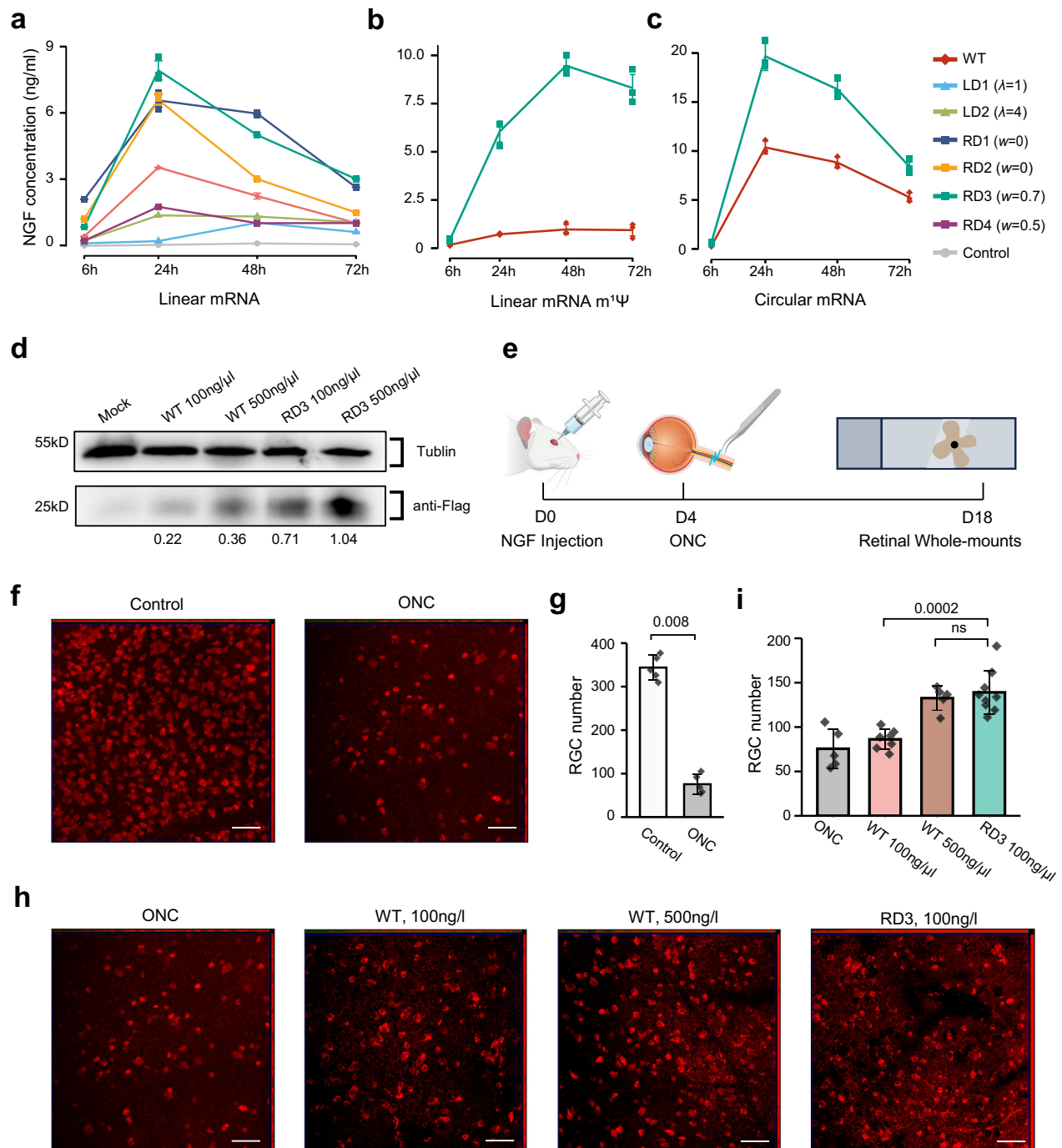
## Discussion

In this study, we present RiboDecode, a data-driven, deep learning-based framework for mRNA codon optimization. The generative optimization framework, guided by the deep learning prediction model, enables the efficient exploration of the immense space of possible codon sequences. This allows RiboDecode to discover previously unexplored, highly optimized sequences that may not be accessible to traditional optimization methods. RiboDecode-optimized sequences demonstrate superior performance in various mRNA formats, including unmodified, m<sup>1</sup>Ψ-modified, and circular mRNAs, highlighting its broad applicability in the rapidly evolving field of mRNA therapeutics. In vitro and in vivo experiments using the optimized sequences of therapeutically relevant proteins show substantial enhancements in protein expression compared to the unoptimized sequences. These improvements further translate into increased therapeutic efficacy, as demonstrated by significantly enhanced immune responses to an optimized influenza vaccine and markedly improved RGC protection in mice with optic nerve injury.

While mRNA abundance serves as both an input feature and intrinsically linked to the Ribo-seq RPKM target variable, potentially explaining its high predictive weight in ablation analysis, the model demonstrated significant improvements in prediction accuracy by integrating codon sequences and cellular context. Importantly, RiboDecode’s in vitro and in vivo validation, where optimized sequences substantially enhanced protein expression and therapeutic efficacy, confirms its ability to optimize biologically meaningful translational signals rather than merely reflecting transcript abundance. The superior performance of RiboDecode may be attributed to several factors. Firstly, RiboDecode’s deep learning model learns directly from diverse nature sequences with translation measurements, enabling it to capture complex patterns of codon sequences for mRNAs with high translation level. Second, the model considered the cellular contexts of mRNA translation. Third, RiboDecode’s generative optimization framework allows it to explore a large sequence space and to discover distinct, highly optimized sequences that may not be accessible to the traditional approaches.

These findings also align with the broader framework of gene expression regulation: while mRNA abundance is a prerequisite for protein synthesis (consistent with its predictive dominance), our approach directly targets translational control, a critical layer highlighted by Schwahnhäuser et al.<sup>45</sup>, by leveraging Ribo-seq-derived codon usage and cellular context to maximize translational efficiency. Although the model excludes post-translational events (as Ribo-seq captures ribosome activity prior to protein maturation), the achieved gains in protein output and in vivo efficacy robustly validate translational optimization as a key determinant of functional protein levels, complementing established regulatory hierarchies.

The findings of our study have important implications for the field of mRNA therapeutics. Firstly, RiboDecode can generate and evaluate a vast number of diverse codon combinations. This capability allows RiboDecode to optimize mRNA sequences beyond the limitations of evolutionary constraints, potentially uncovering more efficient codon usage patterns not found in natural transcripts. Second, by substantially increasing protein production, the optimized sequences can improve the potency and reduce the required dose of mRNA-based treatments, potentially mitigating side effects and enhancing patient outcomes. This is particularly relevant for applications such as protein replacement therapies, where achieving high levels of protein expression is crucial for therapeutic success. Third, RiboDecode’s robustness and versatility across different mRNA formats, including modified and circular mRNAs, expand the range of therapeutic



**Fig. 6 | Enhanced protein expression and therapeutic efficacy with optimized NGF mRNA.** **a** NGF protein expression in HEK293T cells, measured by ELISA. RD sequences were designed by RiboDecode, with  $w$  parameter indicated in parentheses. LD sequences were designed by LinearDesign. Biological replicates were repeated three times. Data are presented as mean values  $\pm$  SD. **b, c** Protein expression levels of RD3 in  $m^1\psi$ -modified (**b**) and circular (**c**) mRNA formats. Biological replicates were repeated three times. Data are presented as mean values  $\pm$  SD. **d** In vivo protein expression of RD3. The  $m^1\psi$ -modified mRNAs were injected into mouse retinas, and protein levels were measured by western blot 48 h post-injection. The expression values were quantified using GelAnalyzer. Three times of the experiment was repeated independently with similar results. **e** Timeline of NGF mRNA therapy in ONC model: At Day 0 (D0), the  $m^1\psi$ -modified NGF mRNAs were injected into mouse retinas. At Day 4 (D4), the optic nerve was

subjected to a physical crush injury. At Day 18 (D18), RGC numbers were quantified using immunofluorescence staining. Created with BioGDP.com<sup>81</sup>. **f, g** RGC numbers measured by immunofluorescence staining in the ONC mouse were significantly reduced compared to that of the control (one-sided Wilcoxon test). Scale bar, 50  $\mu$ m. Biological replicates were repeated five to nine times (see Supplementary Data 10). Data are presented as mean values  $\pm$  SD. **h** RGC numbers in the ONC mouse retina after injection of NGF  $m^1\psi$ -modified mRNA with 100 and 500 ng/ $\mu$ l dosages. Scale bar, 50  $\mu$ m. **i** The RGC number in mice treated with 100 ng/ $\mu$ l RD3 was similar to those treated with 500 ng/ $\mu$ l WT and significantly higher than those treated with 100 ng/ $\mu$ l WT. One-sided Wilcoxon test was used to calculate  $p$ -values shown in the figure; ns:  $p > 0.05$ . Biological replicates were repeated five to nine times (see Supplementary Data 10). Data are presented as mean values  $\pm$  SD.

applications for which it can be employed. As the field of mRNA therapeutics continues to evolve and new mRNA formats are developed to enhance stability, reduce immunogenicity, and improve delivery<sup>40,46</sup>, RiboDecode's ability to optimize sequences for these diverse formats will be invaluable.

While our study demonstrates the significant potential of RiboDecode in optimizing mRNA codon sequences for enhanced mRNA translation and therapeutic efficacy, there are several future directions to explore. Firstly, we focused exclusively on optimizing the codon sequences while not explicitly modeling the 5'UTR. Recognizing the critical role of the 5'UTR in regulating translation initiation, our future work aims to expand RiboDecode to jointly optimize both UTRs and the codon sequence. Second, our results indicate that while our primary goal was to enhance translation through codon optimization, incorporating MFE optimization synergistically modulating mRNA secondary structure stability and translation efficiency. However, we observed the best  $w$  value appears context-dependent, suggesting a one-size-fits-all approach may be suboptimal. Consequently, we recommend that experimental designs include the testing of multiple  $w$  values to identify the optimal balance for each mRNA target. Further research into the determinants of the optimal  $w$  value will be critical to refine this strategy and could lead to more systematic approaches for integrating MFE into mRNA design. Third, our MFE model cannot predict MFEs for unseen mRNAs. For an unseen mRNA, the model must first train the sequences with MFEs labelled by RNAfold. A general MFE model should be developed and used in future. Finally, the model was trained exclusively on Ribo-seq data from endogenous, unmodified mRNAs. Although our results showed significant expression enhancements for optimized sequences in both m1 $\Psi$ -modified and circular forms compared to their respective controls, the relative fold-improvement compared to the unmodified format varied between constructs. Future iterations could potentially incorporate data from modified and circular transcripts to further refine optimization rules.

In conclusion, RiboDecode represents a paradigm shift from rule-based to data-driven mRNA optimization, potentially uncovering previously inaccessible principles of efficient translation that were previously inaccessible. RiboDecode 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.

## Methods

### Data collection and processing

**Data preprocessing and filtering.** We downloaded translation counts of Ribo-seq datasets from the RPFdb database<sup>33,34,47</sup>. The following steps were implemented for processing the Ribo-seq data in accordance with RPFdb: First, to prevent adapter interference in downstream analyses, the 3' adapter sequences were manually extracted for each dataset from the original publications or the corresponding MultiQC<sup>48</sup> outputs. Adapter sequences, if present at the ends of sequencing reads, were subsequently removed using Cutadapt (version 1.16)<sup>49</sup>. Next, to minimize rRNA and tRNA contamination, sequences corresponding to rRNA and tRNA were retrieved for each species from ENSEMBL<sup>50</sup> and UCSC<sup>51</sup> databases and removed post-mapping using Bowtie2<sup>52</sup>. Finally, to ensure the retention of high-quality ribosome-protected footprints, which exhibit a characteristic read-length distribution reflecting the size of a translating ribosome on the RNA, only footprints within the 25–34 nucleotide length range were retained after contaminant removal and alignment. The count tables were transformed to reads per kilobase per million (RPKM). Because some of the paired RNA-seq were not available in RPFdb, we reprocessed the RNA-seq datasets. The raw FASTQ files of RNA-seq were trimmed by sickle<sup>53</sup> (v1.33) for adapter removal and quality control. Then, to filter out reads from tRNA or rRNA, we mapped the reads to human tRNA and rRNA reference sequences (hg38) using bowtie2<sup>54</sup>

(v2.3.5.1, -L 20). The unmapped reads were then mapped to the human genome (GRCh38, gencode.v28, <https://www.gencodegenes.org/>) using STAR<sup>55</sup> (v2.7.4a). Finally, read counts for each gene were summarized by featureCounts<sup>56</sup> (v2.0.1, -t exon). The expression counts of Ribo-seq and RNA-seq were  $\ln(\text{RPKM} \times 5 + 1)$  transformed. Genes with low expression (median RPKM < 1) were filtered out. Finally, 11,725 coding genes from 320 samples with 24 cell types were included in this study (Supplementary Data 1).

**Cross-validation dataset preparation.** Out of 11,725 genes in the Ribo-seq data, we randomly selected 1173 genes (1/10 of the total genes) that were not included in the training datasets, as the “unseen genes” dataset. Out of 320 Ribo-seq datasets, we randomly selected 120 datasets whose cell types were not included in the training datasets, as the “unseen environments” dataset (Supplementary Data 1). The “unseen genes and environments” dataset was also defined (Fig. S4).

**mRNA isoform selection and sequence encoding.** To address the complexity of alternative splicing while managing computational feasibility, we utilized the major isoform of each gene as a representative for mRNA codon variants. This approach was necessitated by the inherent limitations of NGS data in accurately quantifying the proportions of individual isoforms. We defined the major isoform as the transcript with the highest expression level, estimated using RSEM<sup>57</sup> (v1.3.3). In total, our dataset contained 60,255 different mRNA sequences.

### Translation model architecture

Translation is influenced by multiple factors, including codon sequences as the pivotal signals modulating translation<sup>58</sup>, trans-acting elements modulating the cellular environment of translation<sup>59,60</sup>, and mRNA abundance, which provides more templates for translation<sup>61</sup>. To capture these interacting variables, a translation model was developed using a deep neural network architecture comprising 2 convolutional layers and 5 fully connected (FC) layers<sup>62,63</sup>. The model inputs are a codon sequence in one-hot encoding, a transcript abundance, and a vector of gene expression profiles from RNA-seq, presenting cellular environment (Fig. S1). The codon sequences were fixed to 4500 base pairs (bp) starting from the translation start site, which covers over 98% of coding sequences, and those shorter than 4500 bp were zero-padded on 3' end. To process these inputs, 3 sequential FC layers extract features from the cellular environment vector, which are then concatenated with the transcript abundance to form an attention vector. In parallel, convolutional neural networks (CNNs) processes the one-hot encoded codon sequence, generating embeddings with four distinct feature sets. Subsequently, the averaged feature is encoded through the other convolutional layers and flattened. Finally, 2 FC layers yield the predicted translation level output. To enhance model robustness and prevent overfitting, batch normalization<sup>64</sup> and dropout<sup>65</sup> techniques are employed after each layer, and max-pooling<sup>66</sup> is utilized following the convolutional layers. For optimization, the AdamW optimizer<sup>67</sup> is adopted in conjunction with the SmoothL1Loss loss function and ReLU activation functions<sup>68</sup>. Additionally, gradient clipping and learning rate decay strategies are implemented to mitigate gradient explosions and instability during training. The model was parameterized by a total of 15 hyper-parameters, as detailed in Table S9. The model was pretrained for 20 epochs before being used by the codon optimizer of RiboDecode. Model training and validation were assessed using the R-squared metric ( $R^2$ ). The best-performing model was selected based on the highest  $R^2$  value achieved on the validation set. The  $R^2$  metric, defined as  $R^2 = 1 - (\text{SS}_{\text{Residual}} / \text{SS}_{\text{Total}})$ , quantifies the goodness of fit by comparing the sum of squared residuals ( $\text{SS}_{\text{Residual}}$ ) to the total sum of squares ( $\text{SS}_{\text{Total}}$ ). The PyTorch (v1.12.0) framework was leveraged for the implementation of the deep model.

The translation model's architecture, which uses a fixed 4500 bp input, results in a high nominal parameter count. This is an artifact of flattening the zero-padded input sequences required for the majority of genes, which have a median length of only 1200 nt. To mitigate the resulting risk of overfitting, we implemented a strong dropout regularization strategy (rate of 0.9) before the final layers. This high dropout rate forces the model to learn robust features rather than fitting to noise in the padded regions. A comparative study (Fig. S15) confirms this strategy is essential: models with no or low dropout overfit immediately, whereas the 0.9 dropout rate eliminates overfitting and ensures stable, improving performance on all validation sets.

### MFE model architecture

The MFE prediction model was formed by a deep neural network architecture comprising 2 convolutional layers, 9 residual blocks derived from ResNet<sup>69</sup>, and 3 FC layers (Fig. S2a). The model input is the one-hot encoded codon sequence, which is the same as the codon sequence used in the translation model. Each residual block consists of 2 convolutional layers with a residual connection. Initially, shallow features are extracted through 2 convolutional layers, followed by 5 residual blocks with max-pooling for deep feature extraction. Subsequently, 4 residual blocks without max-pooling are used to maintain spatial resolution. After the final residual block, features are flattened and fed into 3 FC layers to predict the MFE value. Batch normalization<sup>64</sup> is applied after the first 2 convolutional layers, and dropout<sup>65</sup> is introduced after the subsequent 2 residual blocks. The Fast Gradient Method (FGM)<sup>70</sup> is integrated during training to enhance generalization by applying perturbations to input sequences. We utilized the AdamW optimizer<sup>67</sup>, SmoothL1Loss function, and LeakyReLU activations<sup>71</sup>. The final loss is defined as  $\text{Loss} = \text{Loss}_{\text{mfe}} + \text{Loss}_{\text{fgm}}$ , where  $\text{Loss}_{\text{mfe}}$  is the SmoothL1Loss between predictions and RNAfold<sup>35</sup> MFE values, and  $\text{Loss}_{\text{fgm}}$  is the SmoothL1Loss with added perturbations. The model is trained alongside optimization by the codon optimizer of RiboDecode using generated sequences as the training data, with performance evaluated using the  $R^2$  metric. Details of 15 hyperparameters are provided in Table S10. The PyTorch 1.12.0 framework and RNAfold (v2.4.18) were used through Python interfaces (v3.8.19). We evaluated the MFE values of mRNAs between our model and RNAfold. We found that our MFE value highly agreed with the one from RNAfold (Fig. S13a), showing the reliability of our MFE model. To determine the optimal training set size, we compared our model trained on a compact set of 210,000 sequences against a model trained on an extended dataset of 11.5 million sequences (achieving a data-to-parameter ratio >1). Both models achieved comparable optimization performance and yielded MFE predictions highly correlated with values from RNAfold (Fig. S13). Given the similar performance, we selected the 210 K training set for our framework, as it reduces computational time by approximately 75% without sacrificing model reliability or accuracy.

We evaluated our MFE optimization framework against existing MFE optimization methods, including the general-purpose, differentiable model (JAX-RNAfold)<sup>72</sup>. Our analysis revealed that the high computational complexity of JAX-RNAfold (v2.0.0-beta) severely limits its use to short sequences (under 600 nt on our V100 GPU with 32 GB video RAM), making it unsuitable for our primary goal of optimizing full-length human mRNAs. To perform a direct, head-to-head comparison, we used a compatible sequence (Gaussia luciferase, 558 nt). Our model achieved a lower (more favorable) MFE, completing the optimization faster and with substantially less GPU memory than JAX-RNAfold (Table S12). While LinearDesign is another powerful tool, it is not a differentiable model and is thus incompatible with our gradient-based optimization framework. Given the critical need for scalability to therapeutically relevant sequence lengths and the efficiency observed

in direct comparison, we proceeded with our sequence-specific MFE optimization approach.

### Codon optimizer architecture

**Fitness score.** The fitness score combines mRNA translation level and MFE predicted by above models.

For translation level, the prediction of a mRNA sequence is conducted using the pre-trained translation model with the cellular environment and transcript abundance holding constant during the optimization. The loss function is designed as follows:

$$L_t = \sqrt{(r - S_t)^2} / \alpha \quad (1)$$

Here,  $S_t$  is the predicted translation level,  $r$  represents the desired output value of the translation model, and  $\alpha$  is a constant based on  $S_t$ .

For MFE, the loss function is formulated based on model parameters of the current epoch since the MFE model is trained alongside the optimization:

$$L_m = -\beta / S_m \quad (2)$$

where  $S_m$  is the predicted value of the current epoch and  $\beta$  is a constant set based on it.

The introduction of  $\alpha$  and  $\beta$  is intended to balance the loss in the translation optimization and MFE optimization processes. The output range of the translation model ( $S_t$ ) typically lies between 0 and -100 (e.g., 0 to 25 for Gluc), while RNAfold-derived MFE values generally range between -100 and -1000 (e.g., -350 to -150 for Gluc). To ensure comparable magnitudes and numerical stability during optimization, we introduced scaling factors  $\alpha$  and  $\beta$  to normalize both loss components. For most codon sequences exhibiting translation prediction values <100 and MFE values >-1000 kcal/mol,  $\alpha = \beta = 100$  provide appropriate normalization (Table S11). However, we recommend parameter adjustments under extreme value conditions:  $\alpha$  should be increased to 1000 when translation predictions exceed 100, while  $\beta$  should be elevated to 1000 when MFE values fall below -1000 kcal/mol.

The final optimized loss function is defined as follows:

$$\text{Loss} = L_m * w + L_t * (1 - w) \quad (3)$$

Where  $w$  is a constant ranging from 0 to 1. When  $w$  is set to 0, only the translation of mRNA is optimized. When  $w$  set to 1, only the MFE of mRNA is optimized. When  $w$  set to a constant value between 0 and 1, both models are optimized simultaneously, with the magnitude indicating the relative emphasis on optimizing each model.

**Optimization process of the codon optimizer.** The codon optimizer uses a gradient ascent optimization approach based on activation maximization (AM)<sup>30</sup> to generate synonymous codon sequences with optimized fitness score. The optimization process involves the following steps:

1. Initial representation: the codon distribution is initially represented in a one-hot manner, where each position is assigned to the specific codon of the original codon sequence.
2. Optimization: the optimized codon distribution by AM becomes probabilistic, assigning a likelihood to each possible synonymous codon at every position, with the goal of maximizing the fitness score.
3. Regularization: during the optimization, a synonymous codon regularizer is used to ensure that the optimization process only adjusts the selection probabilities within the synonymous codons capable of encoding the same amino acids as in the original sequence. The regularizer applies a constraint on the codon

distribution,  $T$ , given by a synonymous substitution mask matrix  $W$ :

$$W = \begin{bmatrix} w_{11} & \cdots & w_{1l} \\ \vdots & \ddots & \vdots \\ w_{lj} & \cdots & w_{lj} \end{bmatrix} \quad (4)$$

where  $w_{ij}$  signifies the selection of the  $j$ th coding category at the  $i$ th position, with value 1 for a synonymous codon, and 0 otherwise. Here,  $i = 1, 2, \dots$  and  $j = 1, 2, \dots$

Subsequently, regularization is performed on  $T$ :

$$T = \frac{T \odot W}{\sum_{i=1}^L (T \odot W)_i} \quad (5)$$

The maximum index values are then converted to a one-hot representation to obtain  $T_{one-hot}$ :

$$T_{one-hot_{ij}} = \begin{cases} 1, & \text{if } j = \operatorname{argmax}(T_i) \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

4. Sequence generation: after the gradient updates, a new codon sequence is generated by selecting codons with the highest probabilities.
5. Iteration: This sequence then re-enters the codon optimizer for further rounds of optimization.

We employed the Adam optimizer and utilized a learning rate decay strategy at different training stages. The optimization process was regularized with a weight decay of  $1 \times 10^{-4}$ . The learning rate was initialized at  $5 \times 10^{-4}$  and decayed by a factor of 10 when the generated data reached 1/2, 3/4, 7/8, and 16/17 of the total 40,000 sequences for one epoch. The optimization process was conducted over 20 epochs. However, we observed that all genes converged to maximum translation levels before the 7th epoch. Consequently, a total of 280,000 sequences were utilized to plot the progression of the generation.

**Integrative training of MFE model.** To train the MFE model, we adopted an active learning strategy that integrates model training with sequence optimization by the codon optimizer. This approach uses the RNAfold tool<sup>31</sup> for labeling and involves simultaneous training, optimization, and sequence generation. The process consists of four interconnected steps that run concurrently with the sequence optimization:

**Initial sampling:** we begin by generating 20,000 sequences through random synonymous substitutions (up to 10% of codons) of the original sequence.

**Initial training:** these 20,000 sequences are input into both the MFE model for predictions and RNAfold for ground truth MFE labeling.

**Sequence generation:** utilizing the trained model of current epoch, new sequences are generated by the codon optimizer of RiboDecode. This pool is sorted by predicted MFE, with lower values being better. The top 480 sequences are selected as generated sequences. The top sequence undergoes two operations: random replacement (up to 10% of codons) and distributed replacement according to its generating codon distribution. A total of 10,000 sequences are generated through this process.

**Model retraining:** The generated sequences, along with those from distributed and random replacements, are used as input for MFE prediction and RNAfold labeling. The model is then retrained on this new data.

**Model running time.** The computational infrastructure employed for our model training and sequence optimization comprised a single NVIDIA Tesla V100 SXM2 GPU with 32 GB of memory and an Intel Xeon Gold 5218 CPU operating at 2.30 GHz.

Training the translation model required approximately 24 h. The training time for the MFE model varied depending on the input sequence length. For example, for the Gluc, NGF, and HA mRNA codon sequences, the MFE model training took 1.08, 1.5, and 3.13 h, respectively. After the MFE model training was completed, the optimization phase was carried out, which required about 1 h for any sequence.

**Optimized codon sequences screening strategy for experimental validation.** To select optimal codon sequences for experimental validation, we employed a multi-step screening process for each gene:

1. Sequence Generation: We generated three rounds of candidate sequences using different optimization strategies: (a) Translation-only optimization ( $w = 0$ ), (b) Joint optimization with moderate MFE consideration ( $w = 0.5$ ), (c) Joint optimization with stronger MFE consideration ( $w = 0.7$ ).
2. Initial Filtering: For each round, we filtered out potentially over-optimized sequences by retaining only those variants with a predicted translation level below the 90th percentile of all generated variants. This step helps avoid unreliable over-optimization that might not translate to real-world performance.
3. Selection Criteria: From the filtered sequences in each round, we selected candidate sequences based on two criteria: (a) High predicted translation level, (b) Low MFE value.
4. Final Selection: We recommend selecting one or more candidates from each of the three optimization rounds ( $w = 0$ ,  $w = 0.5$ ,  $w = 0.7$ ) for experimental validation. This ensures a diverse set of optimized sequences, balancing pure translation optimization with different levels of MFE consideration.

### Model evaluation and analysis

The following software and packages were used for statistical analysis and figure generation: R (version 4.1.0) with the packages sva (v3.40.0), ggplot2 (v3.4.2), viridis (v0.6.4), ggpointdensity (v0.1.0), and data.table (v1.14.4).

**Translation model evaluation.** To evaluate the importance of each input component, we independently trained multiple translation models on the cellular environment vector, transcript abundance, and codon sequence and their combinations<sup>20</sup>. During training, the hyperparameter settings and dataset partitioning were consistent for each model. We replaced the input to be ablated with a zero tensor of the same shape and independently observed the impact of each input component on the final translation prediction.

When evaluating the performance of the model with CAI, MFE, and cellular information as additional inputs, these new features were concatenated into the attention vector of cellular environment features and transcript abundance for mRNA translation prediction.

Furthermore, potential associations between genes in the unseen dataset and those in the seen dataset may lead to overestimation of model performance. To evaluate this potential issue, we obtained gene family annotations from the HUGO Gene Nomenclature Committee (<https://www.genenames.org/>) and excluded genes in the unseen dataset that belonged to the same gene families as those in the seen dataset (resulting in the removal of 130 genes). The results demonstrated only a marginal decrease in prediction accuracy on the “unseen gene test set” (with  $R^2$  decreasing from 0.81329 to 0.81295). Therefore, although some of the genes in the unseen dataset are related to those in the seen dataset, the performance of the model was not overestimated.

**Nucleotide position contribution for translation prediction.** To evaluate the importance of each nucleotide at different positions in the codon sequence, we explored the attribution of translation model predictions to their input features. Here, we implemented the attribution method of Integrated Gradients<sup>73</sup>, obtaining an importance score for each nucleotide position. This method combines the implementation invariance of gradients with the sensitivity of techniques such as LRP<sup>74</sup> or DeepLift<sup>75</sup>. Firstly, we determined a vector with all features set to zero as the baseline value. Then, we linearly interpolated the input features from the baseline value to the actual input values, with these intermediate values representing different strengths or combinations of features. Subsequently, for each interpolated input, we calculated the gradient of the model output relative to that input and multiplied the gradient at each interpolation point with the difference between the input feature value and the baseline value, obtaining the contribution of that feature at each interpolation point. Finally, we weighted and summed the contribution values at all interpolation points to obtain the Integrated Gradient for that feature. For the final analytical outcomes and visualization, the importance score assigned to each nucleotide position was calculated as the mean of absolute values derived from non-padding regions.

**Cellular environment and mRNA abundance for translation model.** The default mRNA abundance level was set to 4.5 (ln-transformed RPKM  $\times 5 + 1$ , median transcript abundance; Fig. S14). To evaluate robustness, we additionally tested input values of 3.82 and 5.19, corresponding to half and double the original RPKM expression level, respectively. Notably, Gluc codon sequence predictions remained highly consistent across these input variations, demonstrating the stability of our model.

optimizing mRNA sequences using RiboDecode, we recommend using the RNA expression profile of the target or the most similar tissue or cell type through either experimental sequencing or publicly available datasets. To facilitate implementation, we have incorporated a user-defined environment input parameter into the software package, accompanied by comprehensive documentation and step-by-step instructions.

For example, to predict translation in HEK293T, the input of cellular environment was constructed as follows. First, the mRNA expression levels of genes representing environmental factors were obtained from RNA-seq data of untreated HEK293T cells. Then, the batch effect arising from different data sources was corrected with R package sva (v3.40.0) ComBat<sup>76</sup>. Finally, the mean value of the mRNA expression was taken as the cellular environment input.

**Codon usage analysis.** Codon sequence variants with enhanced and reduced predicted translation levels were generated in different cellular contexts, including HEK293T, HeLa, and A549. Codon sequences of endogenous genes with high (top 10%) and low (bottom 10%) translation level were selected from Ribo-seq data. We further evaluated the impact of varying expression thresholds (top/bottom 5, 10, and 20%, Figs. S8 and S16), with consistent results observed across all cutoff values, except Gluc in A549 cells at the 5% threshold. Then, the codon usage of generated and endogenous sequences was calculated by the proportion of each codon among synonym codons. To get the codons that appeared more in high-translated sequences, we performed *t*-test (*p*-value < 0.05, after FDR adjustment for multiple-testing) on high and low-translated sequences for both endogenous and generated sets. Codons with the significant and greatest differences on codon usage in endogenous sequences were chosen as the top-10 codons.

**Sequence feature analysis.** To analyze the sequence features of RiboDecode-optimized sequences compared to non-optimized sequences, we followed these steps:

1. Sequence Generation: (a) We randomly selected 2000 RiboDecode-optimized codon sequences. (b) We generated 2000 non-optimized sequences by performing random synonymous codon substitutions on the unoptimized input sequence.
2. Feature Calculation: We calculated several sequence features for both sets of sequences, including CAI, CPB, ENC, GC content (GC %), MFE, and Uracil content (U%).
3. Fold Change Calculation: For each feature, we calculated the fold change by dividing the median value of the optimized sequences by the median value of the non-optimized sequences.
4. Data Transformation: We ln-transformed the fold change values for each feature to normalize the distribution.
5. Data Scaling: The ln-transformed fold change values were scaled to a range of -1 to 1 for visualization purposes.
6. Visualization: The scaled values were used to create a heatmap representation of the feature changes (Fig. 3d). For detailed distributions of each feature, refer to Fig. S9.

All statistical analyze were performed through R (v4.1.0). R package ggplot2 (v3.4.2) were used to make graphs.

**Parameters used for sequence optimization.** For Gluc, HA, and NGF sequence optimization, four sequences were designed by RiboDecode (RD1, RD2, RD3, and RD4:  $w = 0, 0, 0.7$ , and  $0.5$ , respectively), and two sequences were designed by LinearDesign (LD1 and LD2:  $\lambda = 0$  and  $4$ ). The WT sequence was also used as a reference.

To preliminarily evaluate mRNA optimization in vitro, we first performed codon optimization for HA and NGF for the HEK293T cellular environment and measured their protein expression in HEK293T cells (Figs. 5a–c and 6a–c).

**Comparative analysis of sequence space across design methods.** We generated 1000 full-length Gluc codon sequences using Ribotree (v1.1.8), LinearDesign (v1.0.0), CDSfold<sup>38</sup> (<https://github.com/gterai/CDSfold>), and RiboDecode, respectively. High-dimensional sequence embeddings were extracted using CodonBERT (<https://github.com/Sanofi-Public/CodonBERT>)<sup>24</sup>, followed by t-SNE dimensionality reduction for visualization.

## mRNA preparation

**Plasmid construction.** The 5'homology sequence, IRES sequence, 3' homology sequence, E1/E2 sequence and protein coding sequence, were chemically synthesized, and cloned into the vector pUC57, which contains a T7 RNA polymerase promoter.

For linear mRNA, the plasmids contain the 5'UTR, protein coding region, 3'UTR and 105 nt poly-A elements. A  $3 \times$  flag tag was added after the coding region, in order to detect the protein expression by Western blot. The UTR and coding sequences were listed in Supplementary Data 2.

**Linear mRNA production and modification.** The linear mRNAs were produced using the HiScribe T7 High Yield RNA Synthesis Kit and capped with m7G(5')ppp(5')G RNA Cap Structure Analog (NEB, #S1404). Then, the RNA was column-purified. Primers for mRNA amplification are listed in Supplementary Data 2.

For mRNA m<sup>1</sup>Ψ-modification, NI-Me-Pseudo UTP (Yeasen Biotechnology, #10651ES) was used to replace the unmodified UTP.

**Circular mRNA production and purification.** The in vitro transcription (IVT) of circular mRNA was carried out from linearized circular mRNA plasmid templates with the HiScribe T7 High Yield RNA Synthesis Kit (New England Biolabs, E2040S) following the kit manual. After IVT, circular mRNA was purified using the Monarch RNA Cleanup Kit (New England Biolabs, #T2050L). Then, the RNA precursors were heated to 70 °C for 3 min and immediately placed on ice for 2 min. GTP

was added to a final concentration of 2 mM with a buffer (50 mM Tris-HCl, 10 mM MgCl<sub>2</sub>, 1 mM DTT, pH 7.5) for 8 min at 55 °C to catalyze the cyclization. Then RNA was column-purified.

For circular mRNA purification, we collected RNA fractions through UV absorbance at 260 nm on an Agilent 1260 Series HPLC (Agilent) system with a 4.6 × 300 mm column (Sepax Technologies, #215980P-4630) at a flow rate of 0.3 mL/min. The fractions were concentrated with a 4 mL Ultracel-10 regenerated cellulose membrane (Millipore, #UFC8010) and purified by column chromatography. Then, the RNA was treated with RNase R (Beyotime, R7092L) for further enrichment. Finally, RNase R-digested RNA was column-purified.

### In vitro experiments

The HEK293T (#CRL-1573), A549 (#CCL-185), and ARPE-19 (#CRL-2302) cell lines from the American Type Culture Collection were used in this study.

**mRNA transfection.** HEK293T cells were cultured in Dulbecco's Modified Eagle's Medium (BasalMedia, #L110KJ) containing 10% fetal bovine serum (NATOCOR, #SFBE) and 1% penicillin-streptomycin (GIBCO, #15140122) at the condition of 37 °C and 5% CO<sub>2</sub>. mRNAs were transfected into HEK293T cells using Lipofectamine MessengerMax (Invitrogen, #LMRNA015) according to the manufacturer's instructions. At the appropriate time after transfection, the cell lysate or supernatant was collected for protein detection.

**In vitro protein expression measurements.** For measurement of Gaussia luciferase activities, HEK293T cells were seeded in 96-well plate and transfected with 150 ng mRNA per well. The cells were lysed at 6, 24, 48, and 72 h after transfection using 1× cell lysis buffer from Dual Luciferase Reporter Assay Kit (Vazyme Biotech, DL101-01). Then, the luminescence signal was detected following the provided instructions.

For in vitro quantitative measurement of NGF, HEK293T cells were seeded in 24-well plate and 500 ng of RNA was transfected into the cells per well. The cell culture supernatants were collected 6, 24, 48, and 72 h after transfection. Then, the protein expression level was detected using Enzyme-linked Immunosorbent Assay Kit (Cloud-Clone Corp, #SEA105Mu), following the provided instructions.

For measurement of HA protein level, HEK293T cells were seeded in 12-well plate and transfected with 1.25 µg mRNA per well. After 24 h, cells were harvested and collected in 300 µL RIPA lysis buffer (HANGZHOU DUDE BIOLOGICAL CO.LTD, #FD009) that contained 1% PMSF (HANGZHOU DUDE BIOLOGICAL CO.LTD, #FD0100). Then, protein level was analyzed with western blot using anti-flag primary antibody (Sigma-Aldrich, #F1804, 1:1000).

Each experiment was repeated three or four times from distinct samples.

**Cell type specificity experiments.** We used RiboDecode to design three Gluc mRNA variants optimized for preferential expression in HEK293T cells. The optimization process considered the cellular context of HEK293T cells while maintaining or reducing expression levels in A549 and ARPE19 cells. Gluc protein expression was measured 24 h post-transfection. Each experiment was performed in quadruplicate. Expression levels were normalized to the reference (MF882921.1) Gluc mRNA for each cell type.

### In vivo experiments

**Mouse retina histology and microscopy.** For retinal whole-mounts immunofluorescence, eyes were surgically removed from perfused mice and fixed with 4% PFA at room temperature for 1 h. Retinas were detached and whole mount staining was performed. The retinas were blocked for 1 h in PBS staining buffer containing 5% normal donkey serum (Solarbio, SL050) and 0.1% Triton ×-100 (Sigma, ×100-100). The

retinas were incubated with the primary antibody (Novus, #NBP2-20112, 1:500) overnight at 4 °C and washed 3 times with PBS for 5 min each before incubation with the secondary antibody (CST, #4413S, 1:1000) for 2 h at room temperature. The retinas were washed again with PBS 3 times for 5 min each and then mounted.

Confocal images were obtained using a Zeiss LSM 980 microscope. To count retinal ganglion cells (RGCs), we analyzed 320 × 320 µm samples from the peripheral retina. These samples were taken ~500 µm from the center to the edge in all four quadrants of the retina. We then processed the data using ImageJ and ZEN software.

**mRNA intravitreal injection.** Adult mice were anesthetized by intraperitoneal injection of 1% sodium pentobarbital solution (25 mg/kg). Subsequently, a minor incision was made in the eyelid using a 30-gauge needle to facilitate eye exposure. For intravitreal injections, a micropipette was carefully inserted through the serosal opening, and formulations such as LNP-mRNA or other substances were administered into the vitreous body of the eye. To prevent the backflow of the injected solution, the needle was maintained in position for approximately 10 s after the injection before being gently withdrawn. To protect the cornea post-procedure, tobramycin was applied.

**In vivo NGF protein expression.** The m<sup>1</sup>Ψ-modified NGF mRNAs were injected into mouse retina and protein level were measured after 48 h. The detachment and processing of the mouse retina were performed in the same way as mentioned in the section above. Then, the retinas were collected in 300 µL RIPA lysis buffer (Beyotime, P0013B) that contained 1% PMSF (Sigma, 10837091001). Protein level was analyzed with western blot using anti-NGF Antibody- BSA and Azide-free (Abcam, #ab6199, 1:1000).

**Optic nerve crush mouse model.** Mice were anesthetized by intraperitoneal injection of 1% sodium pentobarbital solution (25 mg/kg). Then, the eye surface was dilated with tropicamide drops and surface anesthesia was provided with proparacaine hydrochloride. The mice were fixed on the animal operating table. The optic nerve was completely exposed by cutting open the bulbar fascia under the surgical microscope and using microforceps to separate the surrounding tissues and hold the optic nerve for 5 s with a 0.07 mm wide reverse forceps at 1 mm posterior to the globe in the vertical direction of the longitudinal axis of the optic nerve. Tobramycin was applied daily to the superior orbital rim incision for 3 days postoperatively. Five to nine biological replicates were performed for each experiment.

**In vivo immunogenicity.** First, m<sup>1</sup>Ψ-modified mRNA with WT and RD1 codon sequences were encapsulated within lipid nanoparticles (LNP). BALB/c mice were received two intramuscular doses of 10 µg mRNA each, with the dose determined by previous studies<sup>77,78</sup>, administered on day 0 (prime) and day 14 (boost). (Fig. 5d). We collected mouse serum and performed micro-neutralization (MN) assays to quantify neutralizing antibodies at two time points: day 14 (for prime response) and day 28 (for boost response) (Fig. 5e, "Methods"). PBS buffer and the inactivated Split-virus (A/Victoria/2570/2019 (H1N1)) Influenza Vaccine was used as negative and positive control, respectively.

**Micro-neutralization (MN) assay.** To measure the titer of anti-influenza virus neutralizing antibodies, we treated mouse serum with receptor-destroying enzyme II (RDE II) (Denka-Seiken) at 37 °C for 16 h, followed by heat-inactivation at 56 °C for 30 min. The experimental procedure of MN was the same as previously reported<sup>79</sup>. In brief, mouse serum samples were incubated with receptor-destroying enzyme (RRE, Denka Seiken) at 37 °C for 16 h, followed by heat inactivation at 56 °C for 30 min. Subsequently, 50 µL of the treated serum was subjected to twofold serial dilution, mixed with an equivalent volume of virus solution containing 100 TCID<sub>50</sub> (50% tissue culture

infectious dose), and then transferred to 96-well plates. Following a 1-h incubation at 37 °C, 100 µl of Madin–Darby canine kidney (MDCK) cell suspension at a density of  $3 \times 10^5$  cells/ml was added to each well. The plates were further incubated for 18 h at 37 °C under 5% CO<sub>2</sub>, after which the cells were washed with phosphate-buffered saline (PBS) and fixed with cold 80% acetone for 10 min. Viral nucleoprotein (NP) was detected using an enzyme-linked immunosorbent assay (ELISA) with a monoclonal antibody specific to influenza A virus NP (Abcam). MN titration was defined as the reciprocal of the maximum serum dilution that neutralizes the 50% influenza H1N1/PR8 virus infections in MDCK cells. The minimum MN titer was set to 10. Eight biological replicates were performed for each experiment.

**Ethics statement.** For the in vivo assay of NGF mRNA, 8-week-old male C57BL/6JGpt mice (RRID: N/A) were purchased from GemPharmatech. All experimental procedures involving these mice were conducted in strict accordance with the animal protocols that received approval from the Institutional Animal Care and Use Committee at the Zhongshan Ophthalmic Center, Sun Yat-Sen University, under the animal ethics approval number Z2021067.

For the in vivo assay of HA mRNA, 6–8-week-old male BALB/cAnNCrI mice were supplied by the Laboratory Animal Center of Sun Yat-Sen University. All experimental procedures were conducted in strict accordance with the animal protocols that received approval from the Animal Care and Use Committee Institutional at the Sun Yat-Sen University, under the animal ethics approval number 2022002794.

All mice were group-housed (5 mice per cage) in specific pathogen-free facilities with a 12-h light/12-h dark cycle, an ambient temperature of 20–26 °C, and a relative humidity of 40–60%. Standard laboratory rodent chow and autoclaved water were provided ad libitum.

## Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

## Data availability

The raw RNA-seq and Ribo-seq data analyzed in this study were sourced from the RPFdb database (<https://sysbio.gzzoc.com/rpfdb/>). The relevant sample accession codes and gene lists are provided in Supplementary Data 1. The processed datasets generated for the translation model—including gene expression counts, mRNA sequences, transcript isoform information, associated sample metadata, filtered genes, and training/test group—are publicly available on Figshare (<https://doi.org/10.6084/m9.figshare.28916288.v3>). The specific mRNA sequences utilized in the experimental validations are included in Supplementary Data 2. All source data underlying the figures and statistical analyses are available within the Supplementary Data (specifically, Supplementary Data 3–10).

## Code availability

The complete translation model and optimization framework of RiboDecode are available on GitHub (<https://github.com/wangfanfff/RiboDecode>). The repository has been archived on Zenodo (<https://doi.org/10.5281/zenodo.17096436>, version v1.0.0)<sup>80</sup> and Figshare (<https://doi.org/10.6084/m9.figshare.28916288.v3>).

## References

- Baden, L. R. et al. Efficacy and safety of the mRNA-1273 SARS-CoV-2 vaccine. *N. Engl. J. Med.* **384**, 403–416 (2021).
- Gebre, M. S. et al. Optimization of non-coding regions for a non-modified mRNA COVID-19 vaccine. *Nature* **601**, 410–414 (2022).
- Pardi, N., Hogan, M. J., Porter, F. W. & Weissman, D. mRNA vaccines — a new era in vaccinology. *Nat. Rev. Drug Discov.* **17**, 261–279 (2018).
- Qin, S. et al. mRNA-based therapeutics: powerful and versatile tools to combat diseases. *Signal Transduct. Target Ther.* **7**, 166 (2022).
- Fang, E. et al. Advances in COVID-19 mRNA vaccine development. *Sig Transduct. Target Ther.* **7**, 94 (2022).
- Zhang, H. et al. Algorithm for optimized mRNA design improves stability and immunogenicity. *Nature* **621**, 396–403 (2023).
- Leppek, K. et al. Combinatorial optimization of mRNA structure, stability, and translation for RNA-based therapeutics. *Nat. Commun.* **13**, 1536 (2022).
- Hanson, G. & Collier, J. Codon optimality, bias and usage in translation and mRNA decay. *Nat. Rev. Mol. Cell Biol.* **19**, 20–30 (2018).
- Mauger, D. M. et al. mRNA structure regulates protein expression through changes in functional half-life. *Proc. Natl. Acad. Sci. USA* **116**, 24075–24083 (2019).
- Wayment-Steele, H. K. et al. Theoretical basis for stabilizing messenger RNA through secondary structure design. *Nucleic Acids Res.* **49**, 10604–10617 (2021).
- Hedaya, O. M. et al. Secondary structures that regulate mRNA translation provide insights for ASO-mediated modulation of cardiac hypertrophy. *Nat. Commun.* **14**, 6166 (2023).
- Sharp, P. M. & Li, W. H. The codon adaptation index—a measure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Res.* **15**, 1281–1295 (1987).
- Kudla, G., Lipinski, L., Caffin, F., Helwak, A. & Zylicz, M. High guanine and cytosine content increases mRNA levels in mammalian cells. *PLoS Biol.* **4**, e180 (2006).
- Lu, P., Vogel, C., Wang, R., Yao, X. & Marcotte, E. M. Absolute protein expression profiling estimates the relative contributions of transcriptional and translational regulation. *Nat. Biotechnol.* **25**, 117–124 (2007).
- Li, G.-W., Burkhardt, D., Gross, C. & Weissman, J. S. Quantifying absolute protein synthesis rates reveals principles underlying allocation of cellular resources. *Cell* **157**, 624–635 (2014).
- Waldman, Y. Y., Tuller, T., Shlomi, T., Sharan, R. & Ruppin, E. Translation efficiency in humans: tissue specificity, global optimization and differences between developmental stages. *Nucleic Acids Res.* **38**, 2964–2974 (2010).
- Fabbri, L., Chakraborty, A., Robert, C. & Vagner, S. The plasticity of mRNA translation during cancer progression and therapy resistance. *Nat. Rev. Cancer* **21**, 558–577 (2021).
- Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
- Cao, C. et al. Deep learning and its applications in biomedicine. *Genomics, Proteom. Bioinforma.* **16**, 17–32 (2018).
- Zrimec, J. et al. Deep learning suggests that gene expression is encoded in all parts of a co-evolving interacting gene regulatory structure. *Nat. Commun.* **11**, 6141 (2020).
- Dauparas, J. et al. Robust deep learning-based protein sequence design using ProteinMPNN. *Science* **378**, 49–56 (2022).
- Sumida, K. H. et al. Improving protein expression, stability, and function with proteinMPNN. *J. Am. Chem. Soc.* **146**, 2054–2061 (2024).
- Bennett, N. R. et al. Improving de novo protein binder design with deep learning. *Nat. Commun.* **14**, 2625 (2023).
- Li, S. et al. CodonBERT large language model for mRNA vaccines. *Genome Res.* **34**, 1027–1035 (2024).
- de Almeida, B. P. et al. A multimodal conversational agent for DNA, RNA and protein tasks. *Nat. Mach. Intell.* **7**, 928–941 (2025).
- Melnikov, A. et al. Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**, 271–277 (2012).
- Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science* **324**, 218–223 (2009).

28. Weinberg, D. E. et al. Improved ribosome-footprint and mRNA measurements provide insights into dynamics and regulation of yeast translation. *Cell Rep.* **14**, 1787–1799 (2016).
29. Liu, T.-Y. & Song, Y. S. Prediction of ribosome footprint profile shapes from transcript sequences. *Bioinformatics* **32**, i183–i191 (2016).
30. Shao, B. et al. Riboformer: a deep learning framework for predicting context-dependent translation dynamics. *Nat. Commun.* **15**, 2011 (2024).
31. Clauwaert, J. et al. Deep learning to decode sites of RNA translation in normal and cancerous tissues. *Nat. Commun.* **16**, 1275 (2025).
32. Linder, J., Srivastava, D., Yuan, H., Agarwal, V. & Kelley, D. R. Predicting RNA-seq coverage from DNA sequence as a unifying model of gene regulation. *Nat. Genet.* <https://doi.org/10.1038/s41588-024-02053-6> (2025).
33. Xie, S.-Q. et al. RPFdb: a database for genome wide information of translated mRNA generated from ribosome profiling. *Nucleic Acids Res.* **44**, D254–D258 (2016).
34. Wang, H. et al. RPFdb v2.0: an updated database for genome-wide information of translated mRNA generated from ribosome profiling. *Nucleic Acids Res.* **47**, D230–D234 (2019).
35. Lorenz, R. et al. ViennaRNA Package 2.0. *Algorithms Mol. Biol.* **6**, 26 (2011).
36. Huang, L. et al. LinearFold: linear-time approximate RNA folding by 5'-to-3' dynamic programming and beam search. *Bioinformatics* **35**, i295–i304 (2019).
37. Linder, J. & Seelig, G. Fast activation maximization for molecular sequence design. *BMC Bioinforma.* **22**, 510 (2021).
38. Terai, G., Kamegai, S. & Asai, K. CDSfold: an algorithm for designing a protein-coding sequence with the most stable secondary structure. *Bioinformatics* **32**, 828–834 (2016).
39. Wright, F. The 'effective number of codons' used in a gene. *Gene* **87**, 23–29 (1990).
40. Chen, R. et al. Engineering circular RNA for enhanced protein production. *Nat. Biotechnol.* **41**, 262–272 (2023).
41. Molinari, N.-A. M. et al. The annual impact of seasonal influenza in the US: measuring disease burden and costs. *Vaccine* **25**, 5086–5096 (2007).
42. 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).
43. Lambiase, A. et al. Experimental and clinical evidence of neuroprotection by nerve growth factor eye drops: Implications for glaucoma. *Proc. Natl. Acad. Sci. USA* **106**, 13469–13474 (2009).
44. Jiang, W. et al. Circular RNA-based therapy provides sustained and robust neuroprotection for retinal ganglion cells. *Mol. Ther. Nucleic Acids* **35**, 102258 (2024).
45. Schwanhäusser, B. et al. Global quantification of mammalian gene expression control. *Nature* **473**, 337–342 (2011).
46. Chaudhary, N., Weissman, D. & Whitehead, K. A. mRNA vaccines for infectious diseases: principles, delivery and clinical translation. *Nat. Rev. Drug Discov.* **20**, 817–838 (2021).
47. Wang, Y., Tang, Y., Xie, Z. & Wang, H. RPFdb v3.0: an enhanced repository for ribosome profiling data and related content. *Nucleic Acids Res.* **53**, D293–D298 (2025).
48. Ewels, P., Magnusson, M., Lundin, S. & Käller, M. MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* **32**, 3047–3048 (2016).
49. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet J.* **17**, 10 (2011).
50. Zerbino, D. R. et al. Ensembl 2018. *Nucleic Acids Res.* **46**, D754–D761 (2018).
51. Casper, J. et al. The UCSC genome browser database: 2018 update. *Nucleic Acids Res.* **46**, D762–D769 (2018).
52. Johnstone, T. G., Bazzini, A. A. & Giraldez, A. J. Upstream ORFs are prevalent translational repressors in vertebrates. *EMBO J.* **35**, 706–723 (2016).
53. Joshi, N. A. & Fass, J. N. Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files (Version 1.33). GitHub <https://github.com/najoshi/sickle> (2011).
54. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
55. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
56. 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).
57. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinforma.* **12**, 323 (2011).
58. Jia, L. & Qian, S.-B. Therapeutic mRNA engineering from head to tail. *Acc. Chem. Res.* **54**, 4272–4282 (2021).
59. Ho, J. J. D. et al. A network of RNA-binding proteins controls translation efficiency to activate anaerobic metabolism. *Nat. Commun.* **11**, 2677 (2020).
60. Luan, Y. et al. Deficiency of ribosomal proteins reshapes the transcriptional and translational landscape in human cells. *Nucleic Acids Res.* **50**, 6601–6617 (2022).
61. Liu, Y., Beyer, A. & Aebersold, R. On the dependency of cellular protein levels on mRNA abundance. *Cell* **165**, 535–550 (2016).
62. Zrimec, J., Buric, F., Kokina, M., Garcia, V. & Zelezniak, A. Learning the regulatory code of gene expression. *Front. Mol. Biosci.* **8**, 673363 (2021).
63. Kelley, D. R., Snoek, J. & Rinn, J. L. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* **26**, 990–999 (2016).
64. Ioffe, S. & Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning* (PMLR, 2015).
65. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* **15**, 1929–1958 (2014).
66. Krizhevsky, A., Sutskever, I. & Hinton, G. E. ImageNet classification with deep convolutional neural networks. *Commun. ACM* **60**, 84–90 (2017).
67. Loshchilov, I. & Hutter, F. Decoupled weight decay regularization. Preprint at <http://arxiv.org/abs/1711.05101> (2019).
68. Nair, V. & Hinton, G. E. Rectified linear units improve restricted Boltzmann machines. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, 807–814 (2010).
69. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770–778 (2016).
70. Miyato, T., Dai, A. M. & Goodfellow, I. Adversarial training methods for semi-supervised text classification. Preprint at <http://arxiv.org/abs/1605.07725> (2021).
71. Maas, A. L., Hannun, A. Y. & Ng, A. Y. Rectifier nonlinearities improve neural network acoustic models. In *Proc. ICML* **30**, 1–6 (2013).
72. Krueger, R. K. & Ward, M. JAX-RNAfold: scalable differentiable folding. *Bioinformatics* **41**, btaf203 (2025).
73. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*, 3319–3328 (PMLR, 2017).
74. Bach, S. et al. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE* **10**, e0130140 (2015).
75. Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation differences. In *International Conference on Machine Learning*, 3145–3153 (PMLR, 2017).

76. Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Bio-statistics* **8**, 118–127 (2007).
77. Pardi, N. et al. Nucleoside-modified mRNA immunization elicits influenza virus hemagglutinin stalk-specific antibodies. *Nat. Commun.* **9**, 3361 (2018).
78. Zhuang, X. et al. mRNA vaccines encoding the HA protein of influenza A H1N1 virus delivered by cationic lipid nanoparticles induce protective immune responses in mice. *Vaccines (Basel)* **8**, 123 (2020).
79. Wang, Y. et al. L226Q mutation on influenza H7N9 virus hemagglutinin increases receptor-binding avidity and leads to biased antigenicity evaluation. *J. Virol.* **94**, e00667–20 (2020).
80. Li, Y. et al. Deep generative optimization of mRNA codon sequences for enhanced mRNA translation and therapeutic efficacy. *Zenodo* <https://doi.org/10.5281/zenodo.17096436> (2025).
81. Jiang, S. et al. Generic diagramming platform (GDP): a comprehensive database of high-quality biomedical graphics. *Nucleic Acids Res.* **53**, D1670–D1676 (2025).

## Acknowledgements

This project was supported by the National Key R&D Program of China (Grant No. 2022YFF1203100, Y.H.); National Natural Science Foundation of China (Grant No. 32470705, Z.X.); Science and Technology Program of Guangzhou, China (Grant No. 2025A03J3990, Z.X.). We gratefully acknowledge the researchers who made their Ribo-seq and RNA-seq data publicly available. We also extend our sincere appreciation to the reviewers for their insightful comments and suggestions, which have significantly improved this manuscript.

## Author contributions

Z.X. conceived, designed, and supervised the project. Y.P.L. collected and preprocessed the datasets, evaluated the model performance, and analyzed the data. Y. He, F.W. and Z.H.T. developed the translation deep model. H.Z. developed the MFE deep model. F.W. and H.Z. developed the AM framework. J.X.D. tested the model. J.Q.Y. and L.F.C. prepared mRNAs and conducted in vitro experiments. W.B.J. conducted in vivo NGF experiments. Z.R.H. and C.J.S. conducted the in vivo HA experiments. G.F.Z. encapsulated mRNA with LNP. T.L., X.H., Z.Y.H., Y. Hu, L.P.W. and C.Y.Z. contributed to project discussion. Y.P.L., Z.X., Y. He, and F.W. wrote the manuscript. All authors read and approved the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-025-64894-x>.

**Correspondence** and requests for materials should be addressed to Yao He or Zhi Xie.

**Peer review information** *Nature Communications* thanks Rhiju Das, who co-reviewed with Hamish Blair, Jigyasa Verma and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2025

<sup>1</sup>State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-Sen University, Guangdong Provincial Key Laboratory of Ophthalmology and Visual Science, Guangzhou, China. <sup>2</sup>School of Public Health (Shenzhen), Sun Yat-Sen University, Shenzhen, China. <sup>3</sup>Shenzhen Key Laboratory of Pathogenic Microbes and Biosafety, Shenzhen Campus of Sun Yat-Sen University, Shenzhen, China. <sup>4</sup>Scientific Research Center, The Seventh Affiliated Hospital, Sun Yat-Sen University, Shenzhen, China. <sup>5</sup>ENO Bio mRNA Innovation Institute, Rhogen Biotechnology Co., Ltd, Shenzhen, China. <sup>6</sup>Center for Chemical Biology and Drug Discovery, China-New Zealand Joint Laboratory of Biomedicine and Health, Guangdong Provincial Key Laboratory of Bio-computing, Guangzhou Institutes of Biomedicine and Health, Chinese Academy of Sciences, Guangzhou, China. <sup>7</sup>Department of Pharmacology, School of Basic Medical Sciences, Fudan University, Shanghai, China. <sup>8</sup>Key Laboratory of Tropical Disease Control (Sun Yat-Sen University), Ministry of Education, Guangzhou, China. <sup>9</sup>State Key Laboratory of Anti-Infective Drug Discovery and Development, School of Pharmaceutical Sciences, Sun Yat-Sen University, Guangzhou, China. <sup>10</sup>These authors contributed equally: Yupeng Li, Fan Wang. ✉ e-mail: [scheyao@hotmail.com](mailto:scheyao@hotmail.com); [xiezhi@gmail.com](mailto:xiezhi@gmail.com)