

Research

Deep learning of the regulatory grammar of yeast 5' untranslated regions from 500,000 random sequences

Josh T. Cuperus,^{1,2,7} Benjamin Groves,^{3,7} Anna Kuchina,^{3,7} Alexander B. Rosenberg,^{3,7} Nebojsa Jojic,⁴ Stanley Fields,^{1,2,5} and Georg Seelig^{3,6}

¹Department of Genome Sciences, University of Washington, Seattle, Washington 98195, USA; ²Howard Hughes Medical Institute, University of Washington, Seattle, Washington 98195, USA; ³Department of Electrical Engineering, University of Washington, Seattle, Washington 98195, USA; ⁴Microsoft Research, Seattle, Washington 98195, USA; ⁵Department of Medicine, University of Washington, Seattle, Washington 98195, USA; ⁶Department of Computer Science & Engineering, University of Washington, Seattle, Washington 98195, USA

Our ability to predict protein expression from DNA sequence alone remains poor, reflecting our limited understanding of *cis*-regulatory grammar and hampering the design of engineered genes for synthetic biology applications. Here, we generate a model that predicts the protein expression of the 5' untranslated region (UTR) of mRNAs in the yeast *Saccharomyces cerevisiae*. We constructed a library of half a million 50-nucleotide-long random 5' UTRs and assayed their activity in a massively parallel growth selection experiment. The resulting data allow us to quantify the impact on protein expression of Kozak sequence composition, upstream open reading frames (uORFs), and secondary structure. We trained a convolutional neural network (CNN) on the random library and showed that it performs well at predicting the protein expression of both a held-out set of the random 5' UTRs as well as native *S. cerevisiae* 5' UTRs. The model additionally was used to computationally evolve highly active 5' UTRs. We confirmed experimentally that the great majority of the evolved sequences led to higher protein expression rates than the starting sequences, demonstrating the predictive power of this model.

[Supplemental material is available for this article.]

Precise control of protein expression is critical for cellular homeostasis and growth. One major layer of this control is exerted via the activity of the 5' untranslated region (UTR). In *Saccharomyces cerevisiae*, the effects of 5' UTRs on protein expression, and in particular on translation, have been characterized in detail for a few genes, pointing to the role of such features as upstream open reading frames (uORFs) (Thireos et al. 1984; Werner et al. 1987), hairpins and other secondary structures (Yoon et al. 1992; Linz et al. 1997; Ringnér et al. 2005), and the Kozak sequence, i.e., the nucleotides (nt) immediately surrounding the AUG start codon (Hamilton et al. 1987). More recent studies have analyzed the functional consequences of polymorphisms and short sequence motifs (≤ 10 nt) in thousands or even tens of thousands of yeast (Dvir et al. 2013) and mammalian (Noderer et al. 2014) 5' UTRs. However, this variation was targeted to nucleotides near the start codon, such that we are still unable to predict from sequence alone how the many distinct sequence and structural features of an entire 5' UTR combine to regulate protein production.

A predictive model relating 5' UTR sequence to protein production would not only provide novel insights into the grammar of biological *cis*-regulation, but it would also enable forward engineering of 5' UTRs with tailor-made properties. Designing sequences with quantitatively predictable properties is a long-standing goal of synthetic biology and a prerequisite for accelerating the design-build-test cycle in metabolic engineering. Models have

been designed, for example, to predict for *Escherichia coli* the impact of a ribosome binding site on translation (Salis et al. 2009) or to understand how combinations of promoters and ribosome binding sites affect RNA and protein expression (Kosuri et al. 2013). However, so far, no generally applicable model has been generated that captures the effect of 5' UTR sequence variation on protein production, primarily due to the lack of a data set large and diverse enough to train such a model. Here, we overcome this limitation by using a library with 489,348 5' UTR variants to generate a predictive model using a convolutional neural network (CNN). While many different types of machine learning models have been applied successfully to biological data, CNNs in particular offer a combination of model power and interpretability, as evidenced by their recent use to predict and visualize transcription factor binding, DNase I hypersensitivity sites, enhancers, and sites of DNA methylation (Alipanahi et al. 2015; Kleftogiannis et al. 2015; Zhou and Troyanskaya 2015; Kelley et al. 2016; Lanchantin et al. 2016; Liu et al. 2016; Quang and Xie 2016; Wang et al. 2016). However, with yeast possessing only about 5000 genes, measurement of the protein expression of this number of 5' UTRs yields far too limited a data set for accurate model building using CNNs.

To generate data on a vastly larger scale, we designed a 5' UTR library composed of completely random 50-nt-long sequences.

⁷These authors contributed equally to this work.

Corresponding authors: gseelig@uw.edu, fields@uw.edu

Article published online before print. Article, supplemental material, and publication date are at <http://www.genome.org/cgi/doi/10.1101/gr.224964.117>.

© 2017 Cuperus et al. This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <http://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

With 4^{50} possible 5' UTR sequences, the size of a resulting data set of protein expression levels is limited only by experimental considerations and measurement capacity. By quantifying our library using a growth assay dependent upon the expression of a functional protein, we capture the effect of variation in the sequence adjacent to the coding region at every step of protein production, including transcription, mRNA processing and stability, translation, and protein stability. The small number of nucleotides typically involved in the binding of proteins to DNA and RNA (4–8 nt) or in forming secondary structures (Weirauch et al. 2013) suggests that functional biological motifs will occur often and in a wide range of contexts within these random 5' UTRs. Our study of alternative splicing corroborates the idea that highly predictive biological models can be learned from fully degenerate sequences (Rosenberg et al. 2015).

Results

5' UTR library and assay

Previous analyses of protein expression resulting from variants in a large library employed fluorescence-activated cell sorting (FACS) for measurement (Kinney et al. 2010; Sharon et al. 2012; Dvir et al. 2013; Kosuri et al. 2013; Noderer et al. 2014; Oikonomou et al. 2014; Lubliner et al. 2015; Shalem et al. 2015), wherein cells are separated into bins of differing fluorescence and the variants within each bin are sequenced. However, the FACS step limits the number of cells that can be assayed, thus also limiting the number of sequence variants that can be tested. To increase the number of 5' UTRs that we could test simultaneously and to improve the resolution in measuring activity, we instead used a competitive growth assay based on the accumulation of the yeast His3 protein; the growth rate of cells in media lacking histidine is proportional to the level of their expressed His3 protein, a selection on continuous fitness values that is not reliant on arbitrary bins. In this selection, yeast are transformed with a library of plasmids carrying a *HIS3* reporter gene, each containing one of the random variants of the 5' UTR sited immediately upstream of the start codon. The number of cells harboring each variant before and after growth in selection media is determined through sequencing, with the relative enrichment or depletion of a variant over time correlating with its *HIS3* expression. Since the number of cells in a selection is not limiting, a single culture can, in principle, be used to assay up to millions of variants. Similar growth-based selections have proven to be accurate in measuring activity differences (Hietpas et al. 2011; Starita et al. 2015; Rich et al. 2016).

We constructed a library of more than half a million 5' UTR variants (of which 489,348 were detected) (Fig. 1A;

Methods). With transcriptional regulation under the well-characterized low expression *CYC1* promoter and the *CYC1* terminator (Chen et al. 1994; Guo et al. 1995; Yagil et al. 1998; Martens et al. 2001; Watanabe et al. 2015) and the use of a low copy number plasmid, the growth of each cell should reflect His3 protein accumulation. We performed a large-batch selection in media lacking histidine and supplemented with 1.5 mM 3-Amino-1,2,4 triazole (3-AT) (Supplemental Fig. S1A; see Methods), a competitive inhibitor of His3, collecting cells after ~6.2 doublings. Using massively parallel sequencing, we quantified the relative change in abundance of each variant before and after selection. These relative changes in abundance are presented as \log_2 enrichments. Because enrichment scores are not normalized to any specific sequence, they can differ between experiments for a single 5' UTR depending on the size of the library undergoing selection and the strength of selection (Supplemental Table S1).

To determine the accuracy of these pooled, competitive enrichment measurements, we chose 13 individual variants from the library with a range of enrichments and individually tested them. The relative growth rates of these 13 were similar to those measured in the bulk assay ($R^2=0.84$) (Supplemental Fig. S1B,C). To further test the validity of our approach, we individually cloned 12 5' UTRs from the library into a yellow fluorescent protein reporter and measured fluorescence levels for these constructs using flow cytometry. We found good correlation between the data from the growth selection and flow cytometry ($R^2=0.61$) (Supplemental Fig. S1D), suggesting that results from the *HIS3* assay generalize to other gene contexts.

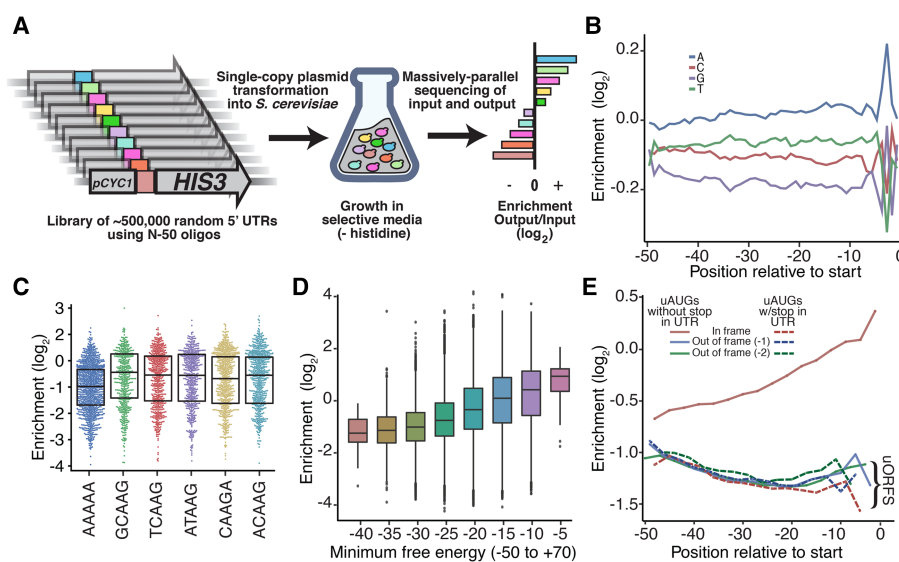


Figure 1. Experimental design and biological discovery. (A) Experimental design of a liquid-based growth assay of 489,348 5' UTR variants. Random 50 nt were introduced directly upstream of the *HIS3* coding sequence, replacing the 56 nt of the 5' UTR of the *CYC1* promoter. These constructs were introduced into a low copy number plasmid, transformed into yeast without a native copy of *HIS3*, and competed in media lacking histidine. The enrichment of each UTR after growth was measured by using massively parallel sequencing before and after selection. (B) 5' UTR enrichment scores per nucleotide were averaged at each position. (C) The Kozak sequences (–5 to –1 position) leading to the highest His3 protein expression compared to the most abundant yeast Kozak sequence (AAAAA). (D) The enrichment of 5' UTRs based on the predicted minimum free energy of the –50 to +70 sequences. (E) The enrichment of 5' UTRs based on the presence of an upstream AUG (uAUG) and a stop codon within the UTR. Upstream open reading frames (uORFs) are characterized by an in-frame uAUG followed by a termination codon before the primary ORF start codon, or an out-of-frame uAUG followed by a stop codon before or after the primary ORF start codon.

Effects of 5' UTR features

The size of our library allowed us to observe particular subsequences many more times than would be possible using genomic sequences alone. By simply comparing the enrichment of 5' UTRs with the subsequence to those without, we could determine whether a feature was, on average, favorable or detrimental (Methods). With this approach, we analyzed the effects of nucleotide identity at each position, focusing in particular on the Kozak sequence—defined here as positions -5 to -1 relative to the start codon. Consistent with prior work (Baim and Sherman 1988; Looman and Kuivenhoven 1993; Dvir et al. 2013), the single nucleotide effects at positions -3 to -1 relative to the start codon were the most important, with an adenine in the -3 position the most beneficial to protein expression (Fig. 1B). This -3 preference for adenine is shared across many eukaryotes, including fungi, mammals, and plants (Nakagawa et al. 2008). We examined the effect on protein expression of all possible Kozak sequences, as the library encompassed the 1024 possible 5-mers at positions -5 to -1 , with each 5-mer occurring, on average, in 478 different 5' UTR contexts (Supplemental Table S2). Although the most common Kozak sequence for yeast is all adenine (Hamilton et al. 1987; Cavener and Ray 1991), we found that this sequence did not lead to the highest protein expression. In fact, 5' UTRs containing 154 other Kozak sequences (122 of which contain an adenine at position -3) led to higher average protein expression than all adenine (the top five are plotted in Fig. 1C), contrary to the widely held belief that the most efficient Kozak is all adenine (Supplemental Table S2). These highly efficient Kozak sequences are also present in the yeast genome in substantial numbers. Each of these top five Kozak sequences from our assay led to higher average protein expression than an all adenine sequence when assessed by ribosome profiling of native yeast genes (Supplemental Fig. S2A; Pop et al. 2014). Of note, genes associated with cytoplasmic translation were enriched for the top five -5 to -1 Kozak sequences (P -value = 9.48×10^{-4}).

We assessed the effect of secondary structure, which can influence ribosome initiation, scanning, and elongation (Rouskin et al. 2013; Pop et al. 2014). We first examined the correlation between the predicted minimum free energy (MFE) of the 5' UTRs and protein expression. To calculate the predicted MFE, we used RNAfold (Gruber et al. 2008, 2015; Lorenz et al. 2011) to fold each 5' UTR sequence along with the first 70 nt of the *HIS3* coding region. Binning the 5' UTRs by their predicted MFE score, we found that lower MFE bins corresponded to decreased protein expression (Fig. 1D). Since the MFE provides only an aggregate measure of structure, we next looked at the effect of structure at each position in the 5' UTR. We found that secondary structure had the largest effect on His3 expression when it occurs either near the 5' end of the UTR or near the start codon (Supplemental Fig. S2B). Access to the 5' cap by the 5' cap binding complex may be reduced by 5' secondary structure, although only highly stable 5' UTR secondary structures (<30 kcal/mol) markedly decrease eukaryotic translation rates (Babendure et al. 2006). Finally, as an alternative, simpler measure of secondary structure, we looked at the presence of hairpins with varying stem (5, 6, or 7 base pairs [bp]) and loop (0–25 nt) lengths within the UTRs. We found that hairpins with longer stems and relatively short loops had the most negative impact on protein expression, perhaps because hairpins with longer loops form more slowly and are therefore scanned more readily by the translation machinery (Supplemental Fig. S2C). In spite of these correlations, secondary structure alone can explain only a small

fraction of overall protein expression (MFE; enrichment correlation of $R^2 = 0.078$) (Supplemental Fig. S2D).

We analyzed the effects of upstream open reading frames, characterized by an in-frame uAUG followed by a termination codon before the primary ORF start codon, or an out-of-frame upstream AUG (uAUG) followed by a stop codon before or after the primary ORF start codon (Morris and Geballe 2000; Wang and Rothnagel 2004; Dvir et al. 2013). uORFs compete with the primary ORF, often producing nonsensical polypeptides and requiring translation to restart at the primary ORF start codon. Consistent with this competition, we found that the presence of a uORF led to greatly reduced protein expression (Fig. 1E; Supplemental Fig. S2E). On the other hand, a uAUG in-frame with the primary ORF—which results only in additional amino acids at the N terminus of the translated protein—caused a minor reduction in expression. The effects of these in-frame uAUGs became more severe as the uAUG was located further toward the 5' end of the UTR (Fig. 1E; Supplemental Fig. S2E), consistent with other reports (Wang and Rothnagel 2004; Dvir et al. 2013; Rich et al. 2016). The additional amino acids might cause a cumulative effect on translation, protein function, or protein stability. Enrichment of 5' UTRs with in-frame uAUGs correlated with the frequency with which the codons added upstream of the true AUG are used in *S. cerevisiae* ($R^2 = 0.75$) (Supplemental Fig. S2F), generally considered a measure of translational efficiency (Sharp and Cowe 1991; Akashi 2003; Pop et al. 2014).

Predicting protein expression with a convolutional neural network

To better understand and engineer UTR sequences, we sought to create a predictive model of protein expression from 5' UTR sequence alone. A comparison between different modeling approaches revealed several trade-offs. For example, a linear regression model with position-dependent 3-mer features ($4^3 \times 48 = 3072$ distinct features; $R^2 = 0.42$) outperformed models with more complex but position-independent features (e.g., 6-mer model; $4^6 = 4096$ features; $R^2 = 0.33$) (Supplemental Fig. S3A,B). Given that many key features of protein expression in yeast have a position dependence—e.g., the identity of the nucleotide at position -3 or the frame of an upstream start codon—it is not surprising that a model that captures such position dependence can outperform a model that does not, even at the expense of using relatively simple features. However, the position-sensitive linear regression model was still unable to capture more complex features, such as uORFs or secondary structure. When features capturing this information were added to the model, the performance was further improved ($R^2 = 0.52$) (Supplemental Fig. S3C; Methods; Dvir et al. 2013). On the other hand, CNNs can capture not only position dependence but also nonlinear interactions between features. Since they do so in an unsupervised fashion, they can also potentially draw attention to unappreciated elements.

CNNs typically consist of several layers of convolution that eventually feed into a classic feed-forward neural network. The first convolutional layer consists of many “filters” that essentially each learn a positional weight matrix (PWM). The output of this layer then feeds into further convolutional layers that can learn interactions between the different motifs recognized by each filter in the first layer. To choose the architecture of the model (such as the size of filters, number of filters, and number of layers), we performed a hyperparameter search using cross-validation on our training set. This search led us to choose a model with three layers of

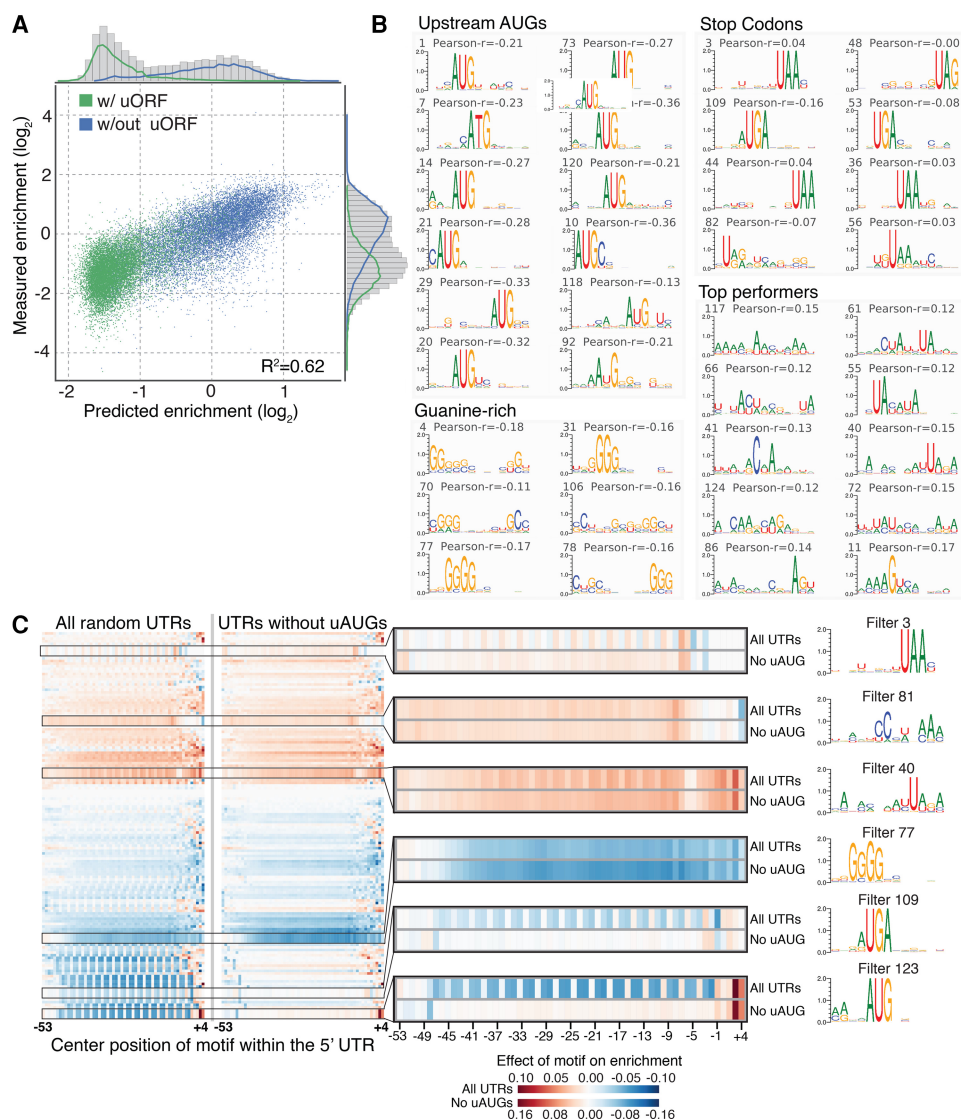


Figure 2. A convolutional neural network (CNN) approach to model random 5' UTR sequences. (A) A three-layer convolutional neural network model trained on random 5' UTRs was tested on a held-out test set of the top 5% based on input read depth. Tested 5' UTRs are specified by color for those with or without an upstream open reading frame. (B) Four hundred eighty-eight thousand random 13-mers were scored for each filter in layer 1 of the CNN. The top 1000 13-mers were used to create a positional weight matrix (PWM) for each filter. These PWMs include motifs of start codons, stop codons, and guanine quadruplexes. Positive Pearson correlations indicate a positive effect on enrichment, while negative correlations indicate a negative effect on enrichment. (C) The effect of each motif per position was measured by assessing the Pearson correlation of motif score and enrichment at each position. Heat maps of all 5' UTRs (left) and those lacking upstream AUGs (right), including specific examples highlighting filters with different positional patterns are shown.

convolution, each with 128 filters of length 13. The convolutional layers then feed into a fully connected layer and finally a linear output layer. The output of our model is the predicted fitness score for each 5' UTR, which should be proportional to protein expression.

Our model accounted for 62% of the observed variability in a held-out test set ($R^2=0.62$) (Fig. 2A), outperforming any of the other models that we tested. When determining the accuracy of the models above, we sought to minimize the impact of experimental noise due to sequencing depth. To do so, we chose those 5' UTRs with the top 5% of input read counts as our test set and used the remaining 95% (464,880 sequences) to train the CNN. Choosing a test set by input read counts allowed us to focus on the higher quality data, while retaining the same distribution of

growth rates as the training set (Supplemental Fig. S4A). As expected, a similarly trained CNN model tested on a held-out 5% that was randomly chosen, and presumably having greater noise due to its lower sequencing depth, had reduced accuracy ($R^2=0.47$) (Supplemental Fig. S4B).

We also wanted to understand whether our approach was making use of the size and other characteristics of our library. To determine whether the scale of our library was an important factor in improving the accuracy of the model, we made learning curves from models trained on different sized subsets of the data. We saw a corresponding decrease in the predictive power of our models as the training size decreased (Supplemental Fig. S4C). We also found that inclusion of the entire 50 nt of sequence was necessary for the high predictive capacity of the model, since a CNN trained

using only the 10 nt adjacent to the start codon performed poorly ($R^2 = 0.097$) (Supplemental Fig. S4D).

To understand the filters presented in the first layer, we scored 488,000 random 13-mers and created a PWM out of the top 1000 scoring sequences for each filter (Methods). Twelve of the 128 filters in the first layer of the model learned uAUG motif variants, while eight learned motifs with stop codons (UAG/UGA/UAA) (Fig. 2B; Supplemental Fig. S5). Additional filters resemble motifs involved in a G-quadruplex, an important motif in RNA secondary structure (Capra et al. 2010). Several other filters have no obvious biological significance at the first layer of convolution as expected; however, in combination, these may correspond to meaningful motifs. Some of these filters might explain the binding sites of RNA-binding proteins, as few binding sites for such proteins have been characterized in *S. cerevisiae* (Ray et al. 2013). Because the model should be learning not just translational efficiency but also features like RNA stability and changes in the transcriptional start site, filters could include several types of potential motifs (Fig. 2B; Supplemental Fig. S5).

Visualizing the positional dependencies of the first-layer motifs resulted in interpretable maps of the 5' UTR sequence-function relationship. Some motifs had positional effects (Fig. 2C), such that they influenced protein expression differentially depending on their location within the 5' UTR. Others showed a striking 3-nt periodicity, reflecting their position relative to the reading frame of uAUGs. This periodicity was not present when 5' UTRs lacking uAUGs were analyzed alone.

The second and third layer of the CNN can learn information about the interplay of lower-level filters. For example, some of the higher layers combine uAUG and stop codon filters to learn the concept of a uORF, as evidenced by the model predicting much lower protein expression for 5' UTR sequences containing a uORF (see Fig. 2A). The model predicts that a 5' UTR with only an in-frame upstream AUG will have a higher enrichment than one with an in-frame uAUG as well as an in-frame stop codon (Supplemental Fig. S3D). The model also predicts that a 5' UTR with an in-frame uAUG as well as an out-of-frame stop codon will have only a small effect on expression (Supplemental Fig. S3D).

Native 5' UTRs may contain a higher density of motifs or higher order motifs not captured using a random library. We therefore asked whether the model could predict protein expression from native *S. cerevisiae* 5' UTRs (Park et al. 2014). We constructed

a library composed of 50-nt segments from all known native 5' UTR sequences in the context of the *HIS3* reporter (Fig. 3A). Our model performed well on the task of predicting the impact of the native sequences ($R^2 = 0.60$) (Fig. 3B; Supplemental Table S3), giving us confidence that it captured the sequence information important for 5' UTR function. On the other hand, a CNN trained only on this native library performed worse on both the native and random library data sets, most likely due to the limited size of the training set (training set = 9492 sequences; $R^2 = 0.47$ native test set; $R^2 = 0.30$ random test set) (Supplemental Fig. S6). As in the case of our model trained on random sequences, a CNN trained on the native sequences using only the 10 nt proximal to the start codon also performed poorly ($R^2 = 0.14$) (Supplemental Fig. S6).

In silico evolution of 5' UTRs

The design of functional sequences with user-defined properties is a compelling demonstration of the predictive power of a model. As a goal, we sought to use our model to improve the expression of a sample of random and native 5' UTRs. We performed a model-guided in silico evolution of 200 5' UTR sequences, half chosen from our random library and half from the native library, representing UTRs over the entire range of activity. During each step of the evolution, we made all possible mutations and selected the single nucleotide substitution predicted to result in the greatest increase in protein expression. By making sequential single base changes, we were able to track how sequence features changed over the course of evolution. We continued making changes until the predicted expression of each 5' UTR plateaued (Fig. 4A; Supplemental Fig. S7A).

For 98 of the sequences derived from the random library and 93 from the native library, we were able to construct *HIS3* constructs for the starting, midpoint, and endpoint of the evolutions. We then tested these 573 sequences in our growth assay. Our approach yielded improved expression for ~84% and ~84% of the sequences selected from the random and native libraries, respectively (Fig. 4B; Supplemental Fig. S7B; Supplemental Table S4). The relative expression of these 5' UTRs held in different 3-AT conditions ($R^2 > 0.93$) (Supplemental Fig. S8A). For the majority of sequences from both libraries, the largest increases in expression occurred between the starting and midpoint sequences, consistent with prediction from the evolutions. In both data sets, we also

found that the degree to which the expression improved negatively correlated with the starting value, suggesting that it is easier to improve upon low expression 5' UTRs than on high expression ones (random: $R^2 = 0.54$; native: $R^2 = 0.86$) (Supplemental Fig. S8B). We also found that the majority of the endpoint sequences from the evolutions (88 out of 98 of the random library and 75 out of 93 of the native library) performed better than 90% of their corresponding larger library after normalization (Fig. 4C; Methods). Even with this technically simple approach, we found that we could increase protein production starting from almost any native or random starting sequence.

To analyze where and how the CNN made changes, we expanded the number

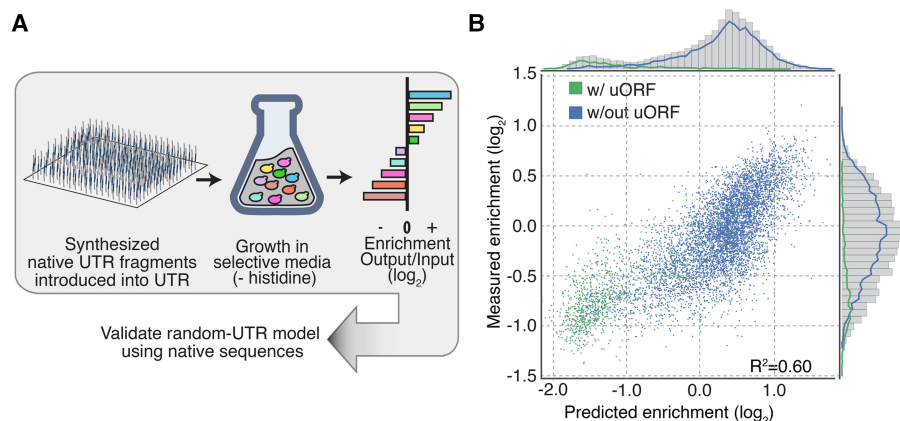


Figure 3. Validation of the CNN model on native 5' UTRs. (A) Native 5' UTR sequences were synthesized in 50-nt fragments and introduced into the *HIS3*-based selection system. (B) Correlation of a native library with the predictions from our convolutional neural network built from random sequences.

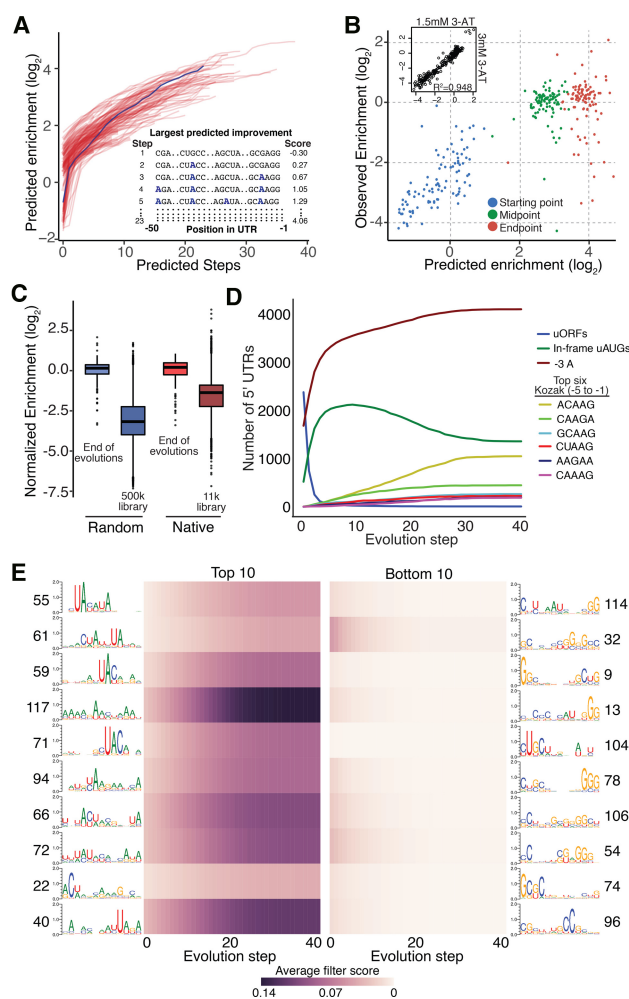


Figure 4. Model-guided optimization of 5000 random sequences. (A) Using our convolutional neural network, we iteratively predicted the optimal single nucleotide change in 100 random 5' UTR sequences until no additional increase in enrichment was predicted. An example of these changes can be seen in the *inset*. (B) The start, midpoint, and endpoints from evolutions in A were tested experimentally. The predicted and observed enrichments are plotted. (C) Experimental data from endpoints of the optimized 5' UTR sequences derived from both the random and native sets of sequence are compared to the enrichment distribution from the original random and native libraries. (D) Five thousand sequences from our random library were evolved over 40 steps and assayed for enrichment and depletion of common nucleotide features. (E) Analysis of the enrichment (*left*) and depletion (*right*) of motifs identified from the first convolutional layer of our model—the same as described in Figure 2.

of random sequences that we computationally evolved to 5000, with each proceeding through 40 steps (Supplemental Tables S7, S8). We looked at the prevalence of simple characteristics, including uORFs, in-frame uAUGs, an A in the -3 position, favorable Kozak sequences, and nucleotide bias (Fig. 4D; Supplemental Fig. S9A). The model selected against uORFs and structure (Fig. 4D; Supplemental Fig. S9B) and selected for an in-frame uAUG, A at the -3 position, and overall A-rich composition except at positions -1 and -49 , where Gs predominated (Supplemental Fig. S7A). Although one Kozak sequence (ACAAG) was the most prevalent, no single 5-nt sequence dominated. During in silico evolutions, the CNN did not collapse sequences to the same sequence. Rather, it maintained a similar and large Hamming distance (~ 34

at the beginning, ~ 32 at the end), suggesting that many of the motifs added are position-independent.

These more predictable changes were accompanied by more complex ones, revealed by analyzing the increase and decrease of the PWMs that we derived from the filters from the first convolutional layer of our model (and the addition and removal of specific 4-mers) (Fig. 4E; Supplemental Fig. S9). Consistent with our initial characterization of the model, the “top performing” filters (Fig. 2B) were selected for enrichment over the course of the rounds of evolution. Moreover, except at either end of the 5' UTR, the spatial distribution of enrichment and depletion of most of the PWMs was largely uniform across the UTRs. Notable exceptions included the negative selection of out-of-frame positions for filters containing strong AUG signals (Supplemental Fig. S10). We also note that the 4-mers most enriched in the evolved sequences often appeared multiple times in a single 5' UTR (Supplemental Fig. S11A). We re-analyzed the experimental data collected from the full random and native libraries and found that additional copies of the enriched 4-mers correlated with continued increases in enrichment (Supplemental Fig. S11B,C, respectively). Similarly, each additional copy of the depleted 4-mers correlated with reduced expression (Supplemental Fig. S11B,C, right side). The most enriched 4-mers and the most enriched Kozak sequence partially overlap with the reverse complement of part of the consensus motif for Nab3 (UCUUGU), a component of the transcription termination Nrd1 complex (Creamer et al. 2011). The three 4-mers (CAAG, ACAAG, and AAGA) that match the reverse complement of the consensus site were highly enriched at the end of the evolutions compared to other motifs with the same nucleotide composition ($P = 4.6 \times 10^{-7}$, *t*-test; 42,906 occurrences in total for the three motifs), while 4-mers found within the motif itself (UCUU, CUUG, and UUGU) occurred only 48 times at step 40.

Discussion

Here, we built and analyzed a library of approximately 500,000 random 5' UTRs. We used the resulting data to train a CNN model that can predict the effect that any 5' UTR in yeast will have on protein expression. Although the model was trained only on data from the random library, it performed equally well at predicting the behavior of native 5' UTRs. The high quality of the predictions is a direct result of the large training data set, compared to methods that consider only the limited set of approximately 5000 native yeast 5' UTRs.

Even though the sequences in our library were randomly generated, the size of the library allowed us to confidently quantify the effect of naturally occurring sequence features on protein expression. For example, we identified 154 variants of the -5 to -1 region that outperformed the consensus Kozak sequence of five adenines. While functional roles for these motifs are supported by ribosome profiling data (Supplemental Fig. S2), the vast majority appear in the yeast genome in such low frequency that they could not be uncovered using native genes.

Analysis and visualization of the model features allowed us to identify *cis*-regulatory motifs. These include motifs such as G-quadruplex sequences known to influence expression, as well as several novel motifs with unknown mechanisms. Some of these motifs may represent target sites for RNA-binding proteins, for which only a limited number of recognition sites have been identified to date in yeast 5' UTRs. Similarly, other motifs likely correspond to regulators that act at other levels of protein expression. For example, a handful of the motifs appear to contain the binding

site of components of the Nrd1 antisense transcriptional regulatory complex (Creamer et al. 2011). Since these motifs are enriched in noncoding RNAs and lead to early transcription termination, they may potentially reduce the amount of antisense transcription, which is known to control expression in a subset of yeast genes (Schulz et al. 2013; Huber et al. 2016). Since the genomic *CYC1* gene—whose elements we use in our assay—has an antisense transcript that initiates within its terminator, the plasmid carrying our reporter gene likely shares this same source of antisense transcription.

We also generated a comprehensive spatial map of the impact of *cis*-regulatory motifs on protein production. In doing so, we observed that the majority of position-dependent effects are observed at either end of the 5' UTR, while along much of the length the positional effects are largely uniform. A notable exception is for motifs containing start or stop codons, where the model is sensitive to the reading frame with respect to the primary *HIS3* start codon.

We demonstrated that our model can be used for the forward engineering of sequences with improved properties. Using a simple model-driven evolution approach, we selected for sequences that were enriched for characteristics correlated with higher protein expression (Fig. 4; Supplemental Figs. S9–S11). Such computationally performed evolutions can dramatically reduce experimental overhead in the design of regulatory elements for synthetic pathways.

Any approach that uses protein expression as its readout is potentially limited by its inability to distinguish among transcription, RNA processing and stability, translation, and protein stability. Transcription and posttranscriptional effects could be disentangled by direct measurement of RNA levels, for which RNA-seq-based approaches are well established (Kwasniewski et al. 2012; Melnikov et al. 2012; Patwardhan et al. 2012). Moreover, because our experimental approach relies upon growth selection, it is inherently less sensitive in detecting sequence variants that lead to poor protein expression. This lack of sensitivity is apparent when we compare the individual R^2 that includes only UTRs containing uORFs ($R^2 = 0.245$) to UTRs without uORFs ($R^2 = 0.478$). Because most UTRs with uORFs do poorly in the growth assay, they are sequenced less frequently and are therefore subject to less accurate measurement. However, such variants that drastically reduce protein expression have been of limited interest, at least for engineering applications.

Yeast have been the source of much of our knowledge of the highly conserved process of translation. Thus, we expect that our approach developed here will be similarly useful for understanding aspects of the biology of other organisms, for example, allowing predictions about the impact of human genetic variation on transcription and translation (Dunham and Fowler 2013).

Methods

Library construction

Synthetic 5' UTR library

We replaced a 56-bp *CYC1* 5' UTR fragment upstream of the *HIS3* ATG on a p415-*CYC1* plasmid (Mumberg et al. 1995) with a library of 50-bp synthetic 5' UTR fragments. The *CYC1* promoter is short (298 nt), with well-established TATA-binding protein sites, upstream activating sequences (UASs) for *HAPI* (Pfeifer et al. 1987) and *MIG1* (Olesen et al. 1987; Treitel and Carlson 1995), and transcriptional start sites, and is regularly used as a consistent low-expression promoter. The synthetic 5' UTR fragments were

constructed by annealing primers 126 and 127 containing an overlap region (ggaccttgcagca) and making the sequence double-stranded using the Klenow fragment of DNA polymerase I (NEB). The resulting fragment had a 50-bp random region and 60-bp and 33-bp 5' and 3' overlaps with the *CYC1* promoter and the *HIS3* coding sequence, respectively, including the ATG start codon. We inverse-PCR-amplified the p415-*CYC1* plasmid backbone with primers 132 and 133 using KAPA HiFi polymerase (Kapa Biosystems), excluding the ATG start codon. Including the start codon in the library fragment served to prevent background plasmids not containing a library fragment from resulting in growth in media lacking histidine. The final library (YTLR200) was assembled using Gibson assembly (Gibson et al. 2009) and electroporated into 40 μ L of DH5 α electrocompetent *E. coli* (NEB), to yield approximately 500,000 colonies.

Native 5' UTR library

For the native library, we constructed 11,962 sequences representing native 5' UTRs from the yeast genome (Park et al. 2014) in 50-bp fragments with 25-bp overlap if the UTR exceeded 50 bp in length, and in smaller fragments for UTRs shorter than 50 bp. Twenty-base-pair overhangs were added to both 5' (acattaggaccttgcagca) and 3' (ATGacagagcagaagcct) ends of these sequences, again overlapping the *CYC1* promoter and *HIS3* gene on the p415-*CYC1* plasmid. The library sequences were purchased from CustomArray Inc. as a mixed oligo pool and amplified by qPCR using primers 126 and 142 in 15 cycles. The resulting fragment was assembled with the plasmid backbone via Gibson reaction and electroporated as described above, resulting in about 200,000 colonies (YTLN200).

Yeast transformation

For the library transformation into yeast, we followed the electroporation protocol described (Benatuil et al. 2010). For the large synthetic 5' UTR library (YTLR200), we used an overnight culture of BY4741 (Baker Brachmann et al. 1998) diluted 1:50 into 50 mL of YPAD media (Amberg et al. 2005) and grown to OD 1.6. We prepared 400 μ L of electroporation-competent cells as described and transformed with a mixture of 3.66 μ g library plasmid YTLR200 linearized with EcoRI and 11.2 μ g of DNA fragment PCR-amplified from YTLR200 with primers 134 and 135 to contain regions of overlap both upstream of and downstream from the EcoRI restriction site (Supplemental Table S5). We grew the transformed library in 500 mL of synthetic dextrose media (Amberg et al. 2005) without leucine (SD-Leu) overnight and used colony counts from serial dilutions plated on SD-Leu to estimate library size. Using a longer region of homology (2.3 kb) led to improved transformation, resulting in $\sim 2 \times 10^6$ colonies. For the generation of the native 5' UTR library (YTLN), the same protocol was followed; 6.7 μ g of EcoRI-digested library plasmid YTLN200 and 15.55 μ g of PCR-amplified fragment (primers 134 and 135) were transformed into 800 μ L of electrocompetent BY4741 yeast cells with similar efficiency as the YTLR library described above. For the transformation of individual plasmids into yeast strains, we followed a lithium acetate method (Gietz and Schiestl 2007). Although yeast can experience relatively high cotransformation rates, there is usually only one CEN-containing plasmid per cell after 36 h of outgrowth (Scanlon et al. 2009). Posttransformation, cells in our experiments were grown well beyond 36 h before growth selections were begun.

Growth rates measurements

Yeast cultures were grown overnight at 30°C in 5 mL until saturated. In 96-well plates, cultures were diluted 1:20 in 200 μ L volume

of minimal selective media. The plates were shaken at 30°C in media lacking histidine and leucine and with 3-Amino-1,2,4 triazole (3-AT, Sigma) (Brennan and Struhl 1980) in a Synergy H1 hybrid reader (Biotek). Mean ($n = 6$) maximum doubling rate was determined by measuring the maximum slope of O.D. 660 measurements over six points of measurement \pm standard error and compared to the calculated enrichment from the competition assay (Supplemental Table S6).

Oligonucleotides and DNA sequencing

Oligonucleotides were obtained from Integrated DNA Technologies with standard desalting purification.

Sanger sequence and analysis was performed as described (Sanger et al. 1977). Deep sequencing of plasmid DNA was performed on an Illumina NextSeq after purifying plasmid DNA using the ZymoPrep yeast plasmid prep II (Zymo Research) and PCR amplification for 12 to 20 cycles.

Library selection

Cells from the input population were collected for sequencing and for back dilution into the selection medium (SD-His-Leu + 1.5 mM 3-AT) in triplicate, adding 1×10^8 cells to 1 L medium. Each replicate was cultured for 20 h to logarithmic phase (O.D. A660 = 1.0, 6×10^9 cells), after which 3×10^8 cells were collected for sequencing.

Optimization of the dynamic range of the selection assay

To optimize the dynamic range of our selection assay, we compared the growth of two yeast strains, one harboring the *HIS3* construct with the native *CYC1* 5' UTR and the other with a 5' UTR containing a strong hairpin known to impair translation (Dvir et al. 2013; Lamping et al. 2013). In the presence of various concentrations of 3-AT, we found a maximal separation of growth rate between the two strains at 1.5 mM 3-AT (Supplemental Fig. S1A).

Strains and media

Yeast experiments used the BY4741 strain. p*CYC1-HIS3* was cloned into the pRS series¹⁰ of yeast vectors with the *LEU2* nutrient marker (pRS415). To construct the plasmids harboring the individual synthetic and native 5' UTRs, we designed a set of one forward and two reverse primers, each 30 bases long with a 10-base overlap in the middle of the sequence for each sequence listed above. We added a 5' acattagacacattgcagca overhang to the forward primer (overlapping *CYC1* promoter) and either agggctttctgctgctcat 3' overhang (overlapping *HIS3* gene) or attcttcacatttagacat 3' overhang (overlapping Venus gene) to the reverse primers. We obtained the oligos in a 96-well array (IDT), annealed them, filled them in with the Klenow fragment, and cloned them into either the p415-p*CYC1* backbone or p415-p*CYC1*-Venus backbone as described. The p415-p*CYC1*-Venus plasmid was constructed by replacing the *HIS3* sequence in the p415-p*CYC1* plasmid used in our library construction with Venus via Gibson assembly.

Enrichment analysis

For the random libraries, we first listed all identified UTRs. We then collapsed any sequences with a Hamming distance of less than 3 and removed any with length less than 3. We used STAR (Dobin et al. 2013) to align reads from both our input libraries and selection libraries to this complete list of UTRs. Next, we counted the number of alignments to each UTR. To calculate the enrichment scores, we first added a “pseudocount” of one to the counts of each UTR in both inputs and selections and normalized the adjust-

ed counts of each UTR by the total counts in each time point (input or selection), calculating the log enrichment of each sequence in the selection relative to the input. Native sequences were quantified similarly; however, because we started with known sequences, we were able to simply count the occurrences of each UTR in both the input and selection libraries as described above.

Identifying features of 5' UTRs

Using the enrichment scores derived from deep sequencing, we determined the average per-position score for each base, resulting in the plot in Figure 2B. Ribosome profiling scores of native genes were calculated as the log-ratio of ribosome footprint counts over mRNA fragment counts. To isolate the effect of the -5 to -1 positions comprising the Kozak sequence, we considered each possible 5-mer separately. We first generated a subset of all 5' UTR sequences containing a specific 5-mer in the -5 to -1 positions. We calculated an average enrichment score for this subset and compared it to the enrichment score calculated for all other 5' UTRs. This process averaged out effects of all sequence elements other than the Kozak sequence. We then repeated this process to get scores for all possible Kozak sequences (Supplemental Table S2). Minimum free energy was calculated using a window of -56 (the predicted transcriptional start site) and +70 using RNAfold (Gruber et al. 2015), then binned based on this MFE in increments of 5. Free base probabilities were also calculated using RNAfold. We searched for potential hairpins comprising combinations of hairpin length (5–7 nt) and loop length (0–24 nt), and then searched for perfect complementary pairs of 5–7 nt contained in a UTR. For each type of hairpin, we calculated the average enrichment scores of the subset of UTRs containing that type of hairpin. Plots were generated using Matplotlib (Hunter 2007) or ggplot2 (Wickham 2009).

Convolutional neural network training

All models were trained using the Python package Keras (<https://keras.io/>). The test set was made from the 5' UTRs with the most reads before selection (top 5%), using the rest of the data as a training set. One hot encoding was used to convert the DNA sequence into a binary matrix; each column in the matrix is associated with a position in the DNA sequence and each row with one of the 4 nt. In each column, a single entry is set to a logical “1” (in the row corresponding to the nucleotide at that position), while the other three entries are “0.” All of our models were trained with the Adam optimizer (Kingma and Ba 2014), and early stopping was used to prevent overfitting to the training data. Cross-validation was performed on the training set to choose the model architecture. We tested combinations of the following hyperparameters: convolutional filter width: [9, 13, 17, 25], number of convolutional filters per layer: [32, 64, 128, 256], number of convolutional layers: [2, 3, 4], number of dense layers: [1, 2], dropout probability in convolutional layers: [0, 0.15], dropout probability in dense layers: [0, 0.1, 0.25, 0.5], number of units in each dense layer: [32, 64, 128, 256]. The best combination of hyperparameters proved to be the following model architecture:

Layer 1: Convolutional, 128 filters (4×13), relu activation, 0.15 dropout probability
 Layer 2: Convolutional, 128 filters (1×13), relu activation, 0.15 dropout probability
 Layer 3: Convolutional, 128 filters (1×13), relu activation, 0.15 dropout probability
 Layer 4: Fully connected layer, 64 hidden units, relu activation, no dropout
 Layer 5: Linear output layer, 1 output unit

Hyperparameter searches for the CNN models presented in Supplemental Figures S4 and S6 were performed in the same fashion as described above and with the same parameters—with the exception of the CNN trained using only the 10 nt adjacent to the start codon, which tested convolution filter widths of [2, 3, 5, 10]. All CNN models were trained as described above. Code for the hyperparameter search and model training are available in the Supplemental Code file and on GitHub at <https://github.com/Seeliglab/2017--Deep-learning-yeast-UTRs>.

k-mer models

The same training and test data were used to train linear regression models based on *k*-mer features. We trained models that simply used *k*-mer counts in each 5' UTR as features as well as training models using *k*-mers at each position as features (e.g., for a 3-mer model, there are 64 possible 3-mer sequences and 48 positions, leading to 3072 model weights; $3072 = 43 \times 48$). We also added additional features used in previous work (Dvir et al. 2013): presence of a uORF, MFE (positions -56 to $+70$) by RNAfold (Gruber et al. 2015), a purine at position -3 , adenine at position -1 , number of GG-dinucleotide occurrences, and number of CACC pattern occurrences. We cross-validated to choose the optimal L2 regularization parameter for all *k*-mer models.

Visualization and analysis of convolutional filters

To visualize each filter in the first layer of convolution, we scored 488,000 newly created random 13-mers with each filter. We then used the top 1000 (0.2%) scoring 13-mers as input into the WebLogo 3 (Crooks et al. 2004) program to generate motifs. However, this visualization does not inform us as to whether the given filter/motif has a positive or negative influence on protein expression. To assess which motifs might increase or decrease protein expression, we calculated the correlation between each filter score and the observed enrichment scores by first calculating the maximum score for each filter in each UTR sequence (across all positions). A high filter score simply indicates a strong match to the given motif within a UTR. We then calculated the Pearson correlation between these maximal filter scores and the enrichment scores.

Forward engineering of sequences

Five thousand random sequences and 100 native sequences were analyzed for the single nucleotide change that led to the largest predicted change from our CNN model. This was done iteratively for 40 steps. From these, the start, midpoint, and endpoint of 100 sequences from the random library and the 100 native sequences were chosen for synthesis. Endpoints were chosen based on the step at which no additional predicted enrichment was attained. Sequences were synthesized by oligonucleotide array (CustomArray Inc.), introduced using Gibson assembly, and transformed into yeast. These yeast transformants were grown, collected, and sequenced as before. Deep-sequencing data were analyzed using the Enrich2 package to assess enrichment of sequences (Rubin et al. 2016). To directly compare the evolved sequences with our larger random and native libraries, we determined the differences in enrichment scores of the starting point sequences (present in both libraries). We then normalized the rest of the larger libraries by the slope of these starting point scores to account for the differences in the strength of selection due to the differences in the sizes of the larger libraries compared to the evolutions. Code for the forward engineering of the selected random and native UTRs is available in the Supplemental Code file and on GitHub at <https://github.com/Seeliglab/2017--Deep-learning-yeast-UTRs>.

Data access

High-throughput reads of selections from this study have been submitted to the Gene Expression Omnibus repository (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number GSE104252. Individual Sanger sequencing reads from this study have been submitted to the Sequence Read Archive (SRA; <https://www.ncbi.nlm.nih.gov/sra>) under accession number SRP120191.

Acknowledgments

We thank the Cold Spring Harbor Laboratory Sequence-function Journal Club for their helpful comments on our preprint. This work was supported by the Defense Advanced Research Projects Agency (DARPA) under Contract No. W911NF-11-2-0068 to G.S., by the National Science Foundation through grant CCF 1317653 to G.S., and by National Institutes of Health grant P41 GM103533 to S.F. S.F. is an investigator of the Howard Hughes Medical Institute.

References

- Akashi H. 2003. Translational selection and yeast proteome evolution. *Genetics* **164**: 1291–1303.
- Alipanahi B, Delong A, Weirauch MT, Frey BJ. 2015. Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. *Nat Biotechnol* **33**: 831–838.
- Amberg DC, Burke D, Strathern JN. 2005. *Methods in yeast genetics: a Cold Spring Harbor Laboratory course manual*. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, NY.
- Babendure JR, Babendure JL, Ding J-H, Tsien RY. 2006. Control of mammalian translation by mRNA structure near caps. *RNA* **12**: 851–861.
- Baim SB, Sherman F. 1988. mRNA structures influencing translation in the yeast *Saccharomyces cerevisiae*. *Mol Cell Biol* **8**: 1591–1601.
- Baker Brachmann C, Davies A, Cost GJ, Caputo E, Li J, Hieter P, Boeke JD. 1998. Designer deletion strains derived from *Saccharomyces cerevisiae* S288C: a useful set of strains and plasmids for PCR-mediated gene disruption and other applications. *Yeast* **14**: 115–132.
- Benatou L, Perez JM, Belk J, Hsieh C-M. 2010. An improved yeast transformation method for the generation of very large human antibody libraries. *Protein Eng Des Sel* **23**: 155–159.
- Brennan MB, Struhl K. 1980. Mechanisms of increasing expression of a yeast gene in *Escherichia coli*. *J Mol Biol* **136**: 333–338.
- Capra JA, Paeschke K, Singh M, Zakian VA, Pringle T. 2010. G-quadruplex DNA sequences are evolutionarily conserved and associated with distinct genomic features in *Saccharomyces cerevisiae*. *PLoS Comput Biol* **6**: e1000861.
- Cavener DR, Ray SC. 1991. Eukaryotic start and stop translation sites. *Nucleic Acids Res* **19**: 3185–3192.
- Chen J, Ding M, Pederson DS. 1994. Binding of TFIID to the *CYC1* TATA boxes in yeast occurs independently of upstream activating sequences. *Proc Natl Acad Sci* **91**: 11909–11913.
- Creamer TJ, Darby MM, Jamonnak N, Schaughency P, Hao H, Wheelan SJ, Corden JL. 2011. Transcriptome-wide binding sites for components of the *Saccharomyces cerevisiae* non-poly(A) termination pathway: Nrd1, Nab3, and Sen1. *PLoS Genet* **7**: e1002329.
- Crooks GE, Hon G, Chandonia J-M, Brenner SE. 2004. WebLogo: a sequence logo generator. *Genome Res* **14**: 1188–1190.
- Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15–21.
- Dunham MJ, Fowler DM. 2013. Contemporary, yeast-based approaches to understanding human genetic variation. *Curr Opin Genet Dev* **23**: 658–664.
- Dvir S, Velten L, Sharon E, Zeevi D, Carey LB, Weinberger A, Segal E. 2013. Deciphering the rules by which 5'-UTR sequences affect protein expression in yeast. *Proc Natl Acad Sci* **110**: E2792–E2801.
- Gibson DG, Young L, Chuang R-Y, Venter JC, Hutchison CA, Smith HO. 2009. Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nat Methods* **6**: 343–345.
- Gietz RD, Schiestl RH. 2007. High-efficiency yeast transformation using the LiAc/SS carrier DNA/PEG method. *Nat Protoc* **2**: 31–34.
- Gruber AR, Lorenz R, Bernhart SH, Neubock R, Hofacker IL. 2008. The Vienna RNA Website. *Nucleic Acids Res* **36**: W70–W74.
- Gruber AR, Bernhart SH, Lorenz R. 2015. The ViennaRNA Web Services. *Methods Mol Biol* **1269**: 307–326.

- Guo Z, Russo P, Yun DF, Butler JS, Sherman F. 1995. Redundant 3' end-forming signals for the yeast *CYC1* mRNA. *Proc Natl Acad Sci* **92**: 4211–4214.
- Hamilton R, Watanabe CK, De Boer HA. 1987. Compilation and comparison of the sequence context around the AUG startcodons in *Saccharomyces cerevisiae* mRNAs. *Nucleic Acids Res* **15**: 3581–3593.
- Hietpas RT, Jensen JD, Bolon DNA. 2011. Experimental illumination of a fitness landscape. *Proc Natl Acad Sci* **108**: 7896–7901.
- Huber F, Bunina D, Gupta I, Khamelinskii A, Meurer M, Theer P, Steinmetz LM, Knop M. 2016. Protein abundance control by non-coding antisense transcription. *Cell Rep* **15**: 2625–2636.
- Hunter JD. 2007. Matplotlib: a 2D graphics environment. *Comput Sci Eng* **9**: 90–95.
- Kelley DR, Snoek J, Rinn JL. 2016. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res* **26**: 990–999.
- Kingma DP, Ba JL. 2014. Adam: a method for stochastic optimization. *arXiv:1412.6980* [cs.LG].
- Kinney JB, Murugan A, Callan CG, Cox EC. 2010. Using deep sequencing to characterize the biophysical mechanism of a transcriptional regulatory sequence. *Proc Natl Acad Sci* **107**: 9158–9163.
- Kleftogiannis D, Kalnis P, Bajic VB. 2015. DEEP: a general computational framework for predicting enhancers. *Nucleic Acids Res* **43**: e6.
- Kosuri S, Goodman DB, Cambrey G, Mutalik VK, Gao Y, Arkin AP, Endy D, Church GM. 2013. Composability of regulatory sequences controlling transcription and translation in *Escherichia coli*. *Proc Natl Acad Sci* **110**: 14024–14029.
- Kwasniewski JC, Mogno I, Myers CA, Corbo JC, Cohen BA. 2012. Complex effects of nucleotide variants in a mammalian *cis*-regulatory element. *Proc Natl Acad Sci* **109**: 19498–19503.
- Lamping E, Niimi M, Cannon RD. 2013. Small, synthetic, GC-rich mRNA stem-loop modules 5' proximal to the AUG start-codon predictably tune gene expression in yeast. *Microb Cell Fact* **12**: 74.
- Lanchantin J, Singh R, Lin Z, Qi Y. 2016. Deep Motif: visualizing genomic sequence classifications. *arXiv:1605.01133* [cs.LG].
- Linz B, Koloteva N, Vasilescu S, McCarthy JE. 1997. Disruption of ribosomal scanning on the 5'-untranslated region, and not restriction of translational initiation *per se*, modulates the stability of nonaberrant mRNAs in the yeast *Saccharomyces cerevisiae*. *J Biol Chem* **272**: 9131–9140.
- Liu F, Li H, Ren C, Bo X, Shu W. 2016. PEDLA: predicting enhancers with a deep learning-based algorithmic framework. *Sci Rep* **6**: 28517.
- Looman AC, Kuivenhoven JA. 1993. Influence of the three nucleotides upstream of the initiation codon on expression of the *Escherichia coli lacZ* gene in *Saccharomyces cerevisiae*. *Nucleic Acids Res* **21**: 4268–4271.
- Lorenz R, Bernhart SH, Höner zu Siederdisen C, Tafer H, Flamm C, Stadler PF, Hofacker IL. 2011. ViennaRNA Package 2.0. *Algorithms Mol Biol* **6**: 26.
- Lubliner S, Regev I, Maya L-P, Edelheit S, Weinberger A, Segal E. 2015. Core promoter sequence in yeast is a major determinant of expression level. *Genome Res* **25**: 1008–1017.
- Martens C, Krett B, Laybourn PJ. 2001. RNA polymerase II and TBP occupy the repressed *CYC1* promoter. *Mol Microbiol* **40**: 1009–1019.
- Melnikov A, Murugan A, Zhang X, Tesileanu T, Wang L, Rogov P, Feizi S, Gnirke A, Callan CG, Kinney JB, et al. 2012. Systematic dissection and optimization of inducible enhancers in human cells using a massively parallel reporter assay. *Nat Biotechnol* **30**: 271–277.
- Morris DR, Geballe AP. 2000. Upstream open reading frames as regulators of mRNA translation. *Mol Cell Biol* **20**: 8635–8642.
- Mumberg D, Müller R, Funk M. 1995. Yeast vectors for the controlled expression of heterologous proteins in different genetic backgrounds. *Gene* **156**: 119–122.
- Nakagawa S, Niimura Y, Gojobori T, Tanaka H, Miura K. 2008. Diversity of preferred nucleotide sequences around the translation initiation codon in eukaryote genomes. *Nucleic Acids Res* **36**: 861–871.
- Noderer WL, Flockhart RJ, Bhaduri A, Diaz de Arce AJ, Zhang J, Khavari PA, Wang CL. 2014. Quantitative analysis of mammalian translation initiation sites by FACS-seq. *Mol Syst Biol* **10**: 748.
- Oikonomou P, Goodarzi H, Tavazoie S. 2014. Systematic identification of regulatory elements in conserved 3' UTRs of human transcripts. *Cell Rep* **7**: 281–292.
- Olesen J, Hahn S, Guarente L. 1987. Yeast HAP2 and HAP3 activators both bind to the *CYC1* upstream activation site, UAS2, in an interdependent manner. *Cell* **51**: 953–961.
- Park D, Morris AR, Battenhouse A, Iyer VR. 2014. Simultaneous mapping of transcript ends at single-nucleotide resolution and identification of widespread promoter-associated non-coding RNA governed by TATA elements. *Nucleic Acids Res* **42**: 3736–3749.
- Patwardhan RP, Hiatt JB, Witten DM, Kim MJ, Smith RP, May D, Lee C, Andrie JM, Lee S-I, Cooper GM, et al. 2012. Massively parallel functional dissection of mammalian enhancers *in vivo*. *Nat Biotechnol* **30**: 265–270.
- Pfeifer K, Arcangeli B, Guarente L. 1987. Yeast *HAP1* activator competes with the factor RC2 for binding to the upstream activation site UAS1 of the *CYC1* gene. *Cell* **49**: 9–18.
- Pop C, Rouskin S, Ingolia NT, Han L, Phizicky EM, Weissman JS, Koller D. 2014. Causal signals between codon bias, mRNA structure, and the efficiency of translation and elongation. *Mol Syst Biol* **10**: 770.
- Quang D, Xie X. 2016. DanQ: a hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences. *Nucleic Acids Res* **44**: e107.
- Ray D, Kazan H, Cook KB, Weirauch MT, Najafabadi HS, Li X, Gueroussov S, Albu M, Zheng H, Yang A, et al. 2013. A compendium of RNA-binding motifs for decoding gene regulation. *Nature* **499**: 172–177.
- Rich MS, Payen C, Rubin AF, Ong GT, Sanchez MR, Yachie N, Dunham MJ, Fields S. 2016. Comprehensive analysis of the *SUL1* promoter of *Saccharomyces cerevisiae*. *Genetics* **203**: 191–202.
- Ringnér M, Krogh M, Mignone F, Gissi C, Liuni S, Pesole G, Hurowitz E, Brown P, Jansen R, Bashirullah A, et al. 2005. Folding free energies of 5'-UTRs impact post-transcriptional regulation on a genomic scale in yeast. *PLoS Comput Biol* **1**: e72.
- Rosenberg AB, Patwardhan RP, Shendure J, Seelig G. 2015. Learning the sequence determinants of alternative splicing from millions of random sequences. *Cell* **163**: 698–711.
- Rouskin S, Zubradt M, Washietl S, Kellis M, Weissman JS. 2013. Genome-wide probing of RNA structure reveals active unfolding of mRNA structures *in vivo*. *Nature* **505**: 701–705.
- Rubin AF, Lucas N, Bajjalieh SM, Papenfuss AT, Speed TP, Fowler DM. 2016. Enrich2: a statistical framework for analyzing deep mutational scanning data. *bioRxiv* doi: 10.1101/075150.
- Salis HM, Mirsky EA, Voigt CA. 2009. Automated design of synthetic ribosome binding sites to control protein expression. *Nat Biotechnol* **27**: 946–950.
- Sanger F, Nicklen S, Coulson AR. 1977. DNA sequencing with chain-terminating inhibitors. *Proc Natl Acad Sci* **74**: 5463–5467.
- Scanlon TC, Gray EC, Griswold KE. 2009. Quantifying and resolving multiple vector transformants in *S. cerevisiae* plasmid libraries. *BMC Biotechnol* **9**: 95.
- Schulz D, Schwab B, Kiesel A, Baejen C, Torkler P, Gagneur J, Soeding J, Cramer P. 2013. Transcriptome surveillance by selective termination of noncoding RNA synthesis. *Cell* **155**: 1075–1087.
- Shalem O, Sharon E, Lubliner S, Regev I, Maya L-P, Yakhini Z, Segal E. 2015. Systematic dissection of the sequence determinants of gene 3' end mediated expression control. *PLoS Genet* **11**: e1005147.
- Sharon E, Kalma Y, Sharp A, Raveh-Sadka T, Levo M, Zeevi D, Keren L, Yakhini Z, Weinberger A, Segal E. 2012. Inferring gene regulatory logic from high-throughput measurements of thousands of systematically designed promoters. *Nat Biotechnol* **30**: 521–530.
- Sharp PM, Cowe E. 1991. Synonymous codon usage in *Saccharomyces cerevisiae*. *Yeast* **7**: 657–678.
- Starita LM, Young DL, Islam M, Kitzman JO, Gullingsrud J, Hause RJ, Fowler DM, Parvin JD, Shendure J, Fields S. 2015. Massively parallel functional analysis of BRCA1 RING domain variants. *Genetics* **200**: 413–422.
- Thireos G, Penn MD, Greer H. 1984. 5' untranslated sequences are required for the translational control of a yeast regulatory gene. *Proc Natl Acad Sci* **81**: 5096–5100.
- Treitel MA, Carlson M. 1995. Repression by Ssn6-Tup1 is directed by MIG1, a repressor/activator protein. *Proc Natl Acad Sci* **92**: 3132–3136.
- Wang X-Q, Rothnagel JA. 2004. 5'-Untranslated regions with multiple upstream AUG codons can support low-level translation via leaky scanning and reinitiation. *Nucleic Acids Res* **32**: 1382–1391.
- Wang Y, Liu T, Xu D, Shi H, Zhang C, Mo Y-Y, Wang Z, Gardiner-Garden M, Frommer M, Cedar H, et al. 2016. Predicting DNA methylation state of CpG dinucleotide using genome topological features and deep networks. *Sci Rep* **6**: 19598.
- Watanabe K, Yabe M, Kasahara K, Kokubo T. 2015. A random screen using a novel reporter assay system reveals a set of sequences that are preferred as the TATA or TATA-like elements in the *CYC1* promoter of *Saccharomyces cerevisiae*. *PLoS One* **10**: e0129357.
- Weirauch MT, Cote A, Norel R, Annala M, Zhao Y, Riley TR, Saez-Rodriguez J, Cokelaer T, Vedenko A, Talukder S, et al. 2013. Evaluation of methods for modeling transcription factor sequence specificity. *Nat Biotechnol* **31**: 126–134.
- Werner M, Feller A, Messenguy F, Piérard A. 1987. The leader peptide of yeast gene *CPA1* is essential for the translational repression of its expression. *Cell* **49**: 805–813.
- Wickham H. 2009. *ggplot2: elegant graphics for data analysis*. Springer-Verlag, New York.
- Yagil G, Shimron F, Tal M. 1998. DNA unwinding in the *CYC1* and *DED1* yeast promoters. *Gene* **225**: 153–162.
- Yoon H, Miller SP, Pabich EK, Donahue TF. 1992. *SSL1*, a suppressor of a *HIS4* 5'-UTR stem-loop mutation, is essential for translation initiation and affects UV resistance in yeast. *Genes Dev* **6**: 2463–2477.
- Zhou J, Troyanskaya OG. 2015. Predicting effects of noncoding variants with deep learning-based sequence model. *Nat Methods* **12**: 931–934.

Received May 12, 2017; accepted in revised form October 18, 2017.



Deep learning of the regulatory grammar of yeast 5' untranslated regions from 500,000 random sequences

Josh T. Cuperus, Benjamin Groves, Anna Kuchina, et al.

Genome Res. 2017 27: 2015-2024 originally published online November 2, 2017
Access the most recent version at doi:[10.1101/gr.224964.117](https://doi.org/10.1101/gr.224964.117)

Supplemental Material <http://genome.cshlp.org/content/suppl/2017/11/02/gr.224964.117.DC1>

References This article cites 72 articles, 23 of which can be accessed free at:
<http://genome.cshlp.org/content/27/12/2015.full.html#ref-list-1>

Creative Commons License This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <http://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>
