

Thousands of novel translated open reading frames in humans inferred by ribosome footprint profiling

Anil Raj^{1*†}, Sidney H Wang^{2*†}, Heejung Shim^{2†‡}, Arbel Harpak³, Yang I Li¹, Brett Engelmann², Matthew Stephens^{2,4}, Yoav Gilad^{2*}, Jonathan K Pritchard^{1,3,5*}

¹Department of Genetics, Stanford University, Stanford, United States; ²Department of Human Genetics, University of Chicago, Chicago, United States; ³Department of Biology, Stanford University, Stanford, United States; ⁴Department of Statistics, University of Chicago, Chicago, United States; ⁵Howard Hughes Medical Institute, Stanford University, Stanford, United States

Abstract Accurate annotation of protein coding regions is essential for understanding how genetic information is translated into function. We describe riboHMM, a new method that uses ribosome footprint data to accurately infer translated sequences. Applying riboHMM to human lymphoblastoid cell lines, we identified 7273 novel coding sequences, including 2442 translated upstream open reading frames. We observed an enrichment of footprints at inferred initiation sites after drug-induced arrest of translation initiation, validating many of the novel coding sequences. The novel proteins exhibit significant selective constraint in the inferred reading frames, suggesting that many are functional. Moreover, ~40% of bicistronic transcripts showed negative correlation in the translation levels of their two coding sequences, suggesting a potential regulatory role for these novel regions. Despite known limitations of mass spectrometry to detect protein expressed at low level, we estimated a 14% validation rate. Our work significantly expands the set of known coding regions in humans.

DOI: [10.7554/eLife.13328.001](https://doi.org/10.7554/eLife.13328.001)

*For correspondence: rajanil@stanford.edu (AR); siddisis@uchicago.edu (SHW); gilad@uchicago.edu (YG); pritch@stanford.edu (JKP)

†These authors contributed equally to this work

Present address: ‡Department of Statistics, Purdue University, West Lafayette, United States

Competing interests: The authors declare that no competing interests exist.

Funding: See page 21

Received: 24 December 2015

Accepted: 26 May 2016

Published: 27 May 2016

Reviewing editor: Nicholas T Ingolia, University of California, Berkeley, United States

© Copyright Raj et al. This article is distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use and redistribution provided that the original author and source are credited.

Introduction

Annotations for coding sequences (CDSs) are fundamental to genomic research. The GENCODE Consortium (*Harrow et al., 2012*) has played an important role in annotating the set of protein coding sequences in the human genome, predominantly relying on manual annotation from the Human and Vertebrate Analysis and Annotation (HAVANA) group (*Wilming et al., 2008*). Their annotation pipeline identifies coding sequences using homology with sequences in large cDNA/EST databases and the SWISS-PROT protein sequence database (*Bairoch and Apweiler, 2000*), and validates them using sequence homology across vertebrates and using tandem mass spectrometry. Despite being the most comprehensive database of CDSs available, the current set is conservative and does not include several classes of CDSs, including translated upstream open reading frames (ORFs), dually coded transcripts, and N-terminal extensions and truncations.

Recent work has made it increasingly clear that much of the human genome is transcribed in at least one tissue during some stage of development (*Hangauer et al., 2013; Djebali et al., 2012; Birney et al., 2007; Clark et al., 2011; Kapranov et al., 2007; van Bakel et al., 2010*). However, the functional implication for most of these transcripts remains unclear; in particular, the set of sequences translated from these transcripts are not yet completely characterized. For example, there are several recent studies in which RNA transcripts that were previously annotated as noncoding were shown to encode short functional peptides. One well characterized example is the *polished*

rice (pri) / tarsal-less (tal) locus in flies, a polycistronic mRNA encoding four short peptides active during embryogenesis (Kondo et al., 2007, 2010; Galindo et al., 2007). While short peptides are known to play critical roles in multiple biological processes (Laressergues et al., 2015; Oelkers et al., 2008; Le Mercier et al., 2006; Jung et al., 2009), annotating genomic regions that encode them remains challenging.

Direct proteogenomic mass spectrometry has the potential to fill this gap but suffers from variable accuracy in assignment of peptide sequences to spectra and assignment of identified peptides to proteins (for peptides shared across database entries). Moreover, these approaches suffer from a “needle in a haystack” problem when searching all six translational frames over the transcribed portion of the genome (Nesvizhskii, 2014; Le Mercier et al., 2006; Ma, 2015). Alternative approaches that utilize empirically-derived phylogenetic codon models to distinguish coding transcripts from non-coding transcripts are promising (Lin and Kellis, 2011). However, the success of such approaches is contingent on the duration, strength and stability of purifying selection and these methods may be underpowered for short coding sequences or for newly evolved coding sequences.

Ribosome profiling utilizes high throughput sequencing of ribosome-protected RNA fragments (RPFs) to quantify levels of translation (Ingolia et al., 2009). Briefly, the technique consists of isolating monosomes from RNase-digested cell lysates and extracting and sequencing short mRNA fragments protected by ribosomes. Early studies of ribosome profiling have shown that RPFs are substantially more abundant within the CDS of annotated transcripts compared to the 5′ or 3′ untranslated regions (UTRs) (Ingolia et al., 2009; Weinberg et al., 2016). More importantly, when aggregated across annotated coding transcripts, the RPF abundance within the CDS has a clear three base-pair periodicity while the RPF abundance in the UTRs lacks this periodic pattern.

Recently, using ribosome profiling data, several studies reported conflicting results on the coding potential of long intergenic noncoding RNA (Ingolia et al., 2011; Guttman et al., 2013; Ingolia et al., 2014). These studies assessed coding potential using either i) the enrichment of RPFs within the translated CDS relative to background, or ii) the difference in length of RPFs within the translated CDS compared to background. However, these approaches may lack power for several reasons. First, they make little distinction between ribosomes scanning through the transcript and ribosomes decoding the message. Second, the enrichment signal can be severely diminished if the transcript is significantly longer than the coding region within it. Third, there is often substantial variance in RPF abundance within the CDSs, which can decrease power to detect translated sequences when using a simple RPF enrichment statistic alone. Other studies have used the periodicity structure in RPF counts to identify dual coding sequences and short translated CDSs (Michel et al., 2012; Bazzini et al., 2014), but the methods reported high false positive rates and could only identify a few hundred CDSs.

In this work, we developed riboHMM; a model to identify translated CDSs by leveraging both the total abundance and the codon periodicity structure in RPFs. We used this model to identify thousands of novel CDSs in the transcriptome of human lymphoblastoid cell lines (LCLs).

Probabilistic model to infer translated coding sequences

Ribosome footprint profiling data, when aggregated across annotated coding transcripts centered at their translation initiation (or termination) sites (Figure 1A), show two distinct features that distinguish the CDS from untranslated regions (UTRs).

- **Higher abundance within the CDS.** RPF counts are highly enriched within the CDS overall. Moreover, base positions within the CDS close to the translation initiation and termination sites have substantially higher RPF counts compared to base positions in the rest of the CDS. Untranslated regions have very low RPF counts, with the 5′UTR having a slightly higher RPF count compared to the 3′UTR. Furthermore, base positions in the 5′UTR immediately flanking the initiation site have a slightly higher RPF count compared to the rest of the 5′UTR; a similar pattern is observed in the 3′UTR.
- **Three-base periodicity within the CDS.** RPF counts typically peak at the first position of each codon. The RPF count over the initiation and termination codons tend to have a stronger peak (thus, a slightly different periodic pattern) compared to the rest of the CDS. The RPF counts in the UTRs lack this periodic pattern with similar aggregate counts among base positions in the 5′UTR and 3′UTR.

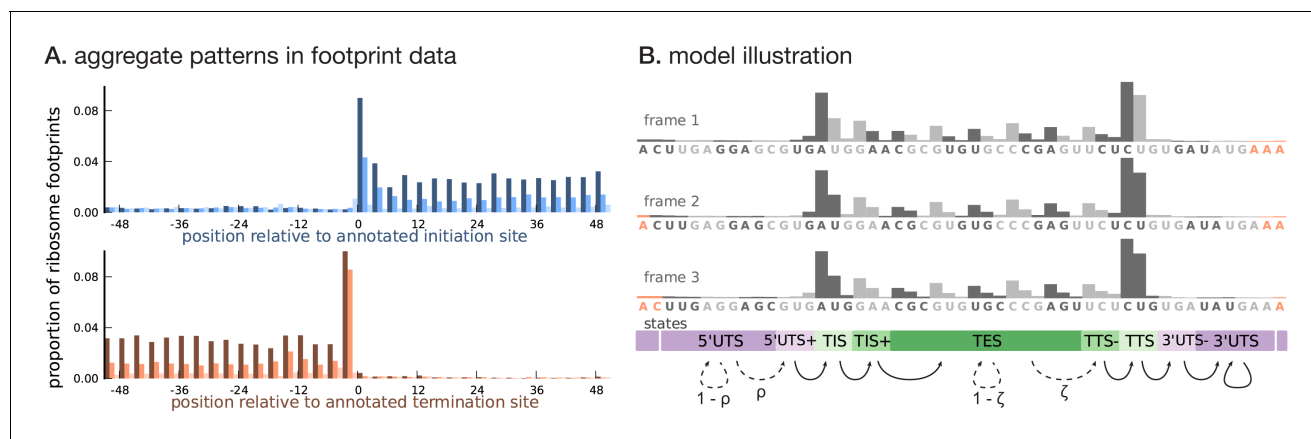


Figure 1. Illustrating the model. (A) Proportion of footprint counts aggregated across 1000 highly expressed annotated coding transcripts, centered at their translation initiation (blue) and termination (orange) sites. In aggregate, RPF count data have higher abundance within the CDS than the UTRs and exhibit a 3-base periodicity within the CDS. (B) Each transcript belongs to one of three unobserved reading frames, and is represented as a sequence of base triplets (highlighted by differing shades of gray) that depends on the reading frame. Each triplet belongs to one of nine unobserved states. The state sequence shown corresponds to frame 3 and varying shades from purple to green highlight the different states. Base positions marked in orange are modeled independently and always belong to the relevant UTS state. Transitions with non-zero probabilities are indicated by arrows, with solid arrows denoting a probability of 1 and dotted arrows denoting probabilities that are a function of the underlying sequence.

DOI: [10.7554/eLife.13328.002](https://doi.org/10.7554/eLife.13328.002)

The following figure supplements are available for figure 1:

Figure supplement 1. Robustness of periodicity parameter estimates.

DOI: [10.7554/eLife.13328.003](https://doi.org/10.7554/eLife.13328.003)

Figure supplement 2. Robustness of occupancy parameter estimates.

DOI: [10.7554/eLife.13328.004](https://doi.org/10.7554/eLife.13328.004)

Figure supplement 3. Decision rules to identify matches and mismatches of mCDS to annotation.

DOI: [10.7554/eLife.13328.005](https://doi.org/10.7554/eLife.13328.005)

Figure supplement 4. Model accuracy.

DOI: [10.7554/eLife.13328.006](https://doi.org/10.7554/eLife.13328.006)

Figure supplement 5. Comparing footprint abundance and gene expression.

DOI: [10.7554/eLife.13328.007](https://doi.org/10.7554/eLife.13328.007)

Figure supplement 6. Comparing the periodicity in ribosome footprint counts for footprints of different lengths.

DOI: [10.7554/eLife.13328.008](https://doi.org/10.7554/eLife.13328.008)

Figure supplement 7. Robustness of parameters for start codon usage to choice of learning set.

DOI: [10.7554/eLife.13328.009](https://doi.org/10.7554/eLife.13328.009)

We developed a framework to infer the translated CDS in a transcript using a model that 1) captures these distinct features of ribosome profiling data and 2) integrates RNA sequence information and transcript expression. As illustrated in **Figure 1B**, to capture the three-base structure in the RPF count data within the CDS, we represented a transcript as a sequence of non-overlapping base triplets. The CDS of the transcript is required to belong to one of three reading frames. To account for all three reading frames, each transcript has a latent frame variable that specifies at which base position of the transcript we begin enumerating the triplets.

Conditional on the frame, we modeled the data for each transcript, represented as a sequence of base triplets, using a hidden Markov model (HMM). Each triplet belongs to one of nine latent states — 5'UTS (5' Untranslated State), 5'UTS+ (the last untranslated triplet prior to the initiation site), TIS (Translation Initiation State), TIS+ (the triplet immediately following the initiation site), TES (Translation Elongation State), TTS- (the translated triplet prior to the termination site), TTS (Translation Termination State), 3'UTS- (the first untranslated triplet immediately following the termination site), and 3'UTS (3' Untranslated State). The states {TIS, TIS+, TES, TTS-, TTS} denote translated triplets and {5'UTS, 5'UTS+, 3'UTS-, 3'UTS} denote untranslated triplets. The probability distribution over the possible sequence of latent states is a function of the underlying RNA sequence. **Figure 1B** illustrates these states, and how they relate to each other, in conjunction with the transcript representation. The groups

of states {5'UTS+, TIS, TIS+} and {TTS-, TTS, 3'UTS-} help model the distinct structure of the RPF counts around the translation initiation and termination sites, respectively.

Assuming each transcript has either 0 or 1 CDS, we restricted the possible transitions between latent states as shown in **Figure 1B**: transitions from 5'UTS to 5'UTS+ occur with probability ρ , transitions from TES to TTS- occur with probability ζ , and all other allowed transitions have probability 1. The transition probabilities ρ and ζ are estimated from the data, and are allowed to depend on the base sequence of the triplet; in addition, the probability ρ also depends on the base sequence context around the triplet (Kozak, 1987). In this work, we assume that translation termination occurs at the first in-frame stop codon (Equation 8), i.e., we do not consider stop codon readthrough.

Conditional on the state assignments, we modeled 1) the total RPF abundance within a triplet, to account for the observation that translated base positions have a higher average RPF count compared to untranslated base positions, and 2) the distribution of RPF counts among the base positions in a triplet, to account for the periodicity in RPF counts within translated triplets. We explicitly accounted for differences in RPF abundance due to differences in transcript expression levels by using transcript-level RNA-seq data as a normalization factor. The short lengths of ribosome footprints mean that many base positions are unmappable; we treated triplets with unmappable positions by modifying the emission probabilities accordingly. Finally, we accounted for the additional variation in RPF counts across triplets assigned to the same state by modeling overdispersion in the triplet RPF abundance (see Materials and methods for details).

To quantify the accuracy of our model, we designed a simulation scheme to estimate what fraction of our inferred translated sequences are false discoveries. We first estimated the Type 1 error rate – i.e., the probability of inferring a translated region when no such region exists – using a set of simulated transcripts that had no signal of translation (null transcripts). The simulated transcripts were constructed by permuting the observed footprint counts in annotated coding transcripts. We then used this estimate to quantify the false discovery rate for each class of translated CDSs identified by riboHMM. Independently, using a simulated set of transcripts containing some signal of translation, we quantified the proportion of transcripts where our model incorrectly identified the precise translation initiation site conditional on having identified a translated sequence (see Materials and methods for details on the simulations).

Results

Application to human lymphoblastoid cell lines

We applied riboHMM to infer translated CDSs in human lymphoblastoid cell lines (LCLs) for which gene expression phenotypes were measured genome-wide: mRNA in 86 individuals, ribosome occupancy in 72 individuals and protein levels in 60 individuals (Lappalainen et al., 2013; Battle et al., 2015). We first assembled over 2.8 billion RNA sequencing reads into transcripts using StringTie (Pertea et al., 2015). This assembly gives us annotated transcripts that are expressed in LCLs, along with novel transcripts that do not overlap any GENCODE annotated gene. (We do not consider novel isoforms of annotated genes in our analyses.) Restricting to transcripts with at least five footprints mapped to each exon, we used riboHMM to identify high-confidence translated CDS. We learned the maximum likelihood estimates of the model parameters using the top five thousand highly expressed genes. The estimated parameters are robust to the choice of the learning set (Figure 1—figure supplements 1 and 2). Using these parameters, we then inferred the maximum *a posteriori* (MAP) frame and latent state sequence for each of the assembled transcripts. We retained transcripts whose MAP frame and state sequence corresponded to a pair of translation initiation and termination sites and had a joint posterior probability greater than 0.8. Using a set of simulated null transcripts, we estimated that this posterior cutoff corresponded to a Type 1 error rate of 4.5% per transcript. The MAP frame and state sequence directly give us the nucleotide sequence with the strongest signal of translation; we refer to these as main coding sequences or mCDS.

Detection of novel CDSs in LCLs

Among 7801 GENCODE annotated coding genes for which we could infer a high posterior mCDS, we recovered the annotated reading frame for at least one transcript isoform in 7491 genes (96%); of these, we recovered the exact annotated CDS in 4500 genes. In the remaining 310 genes, among

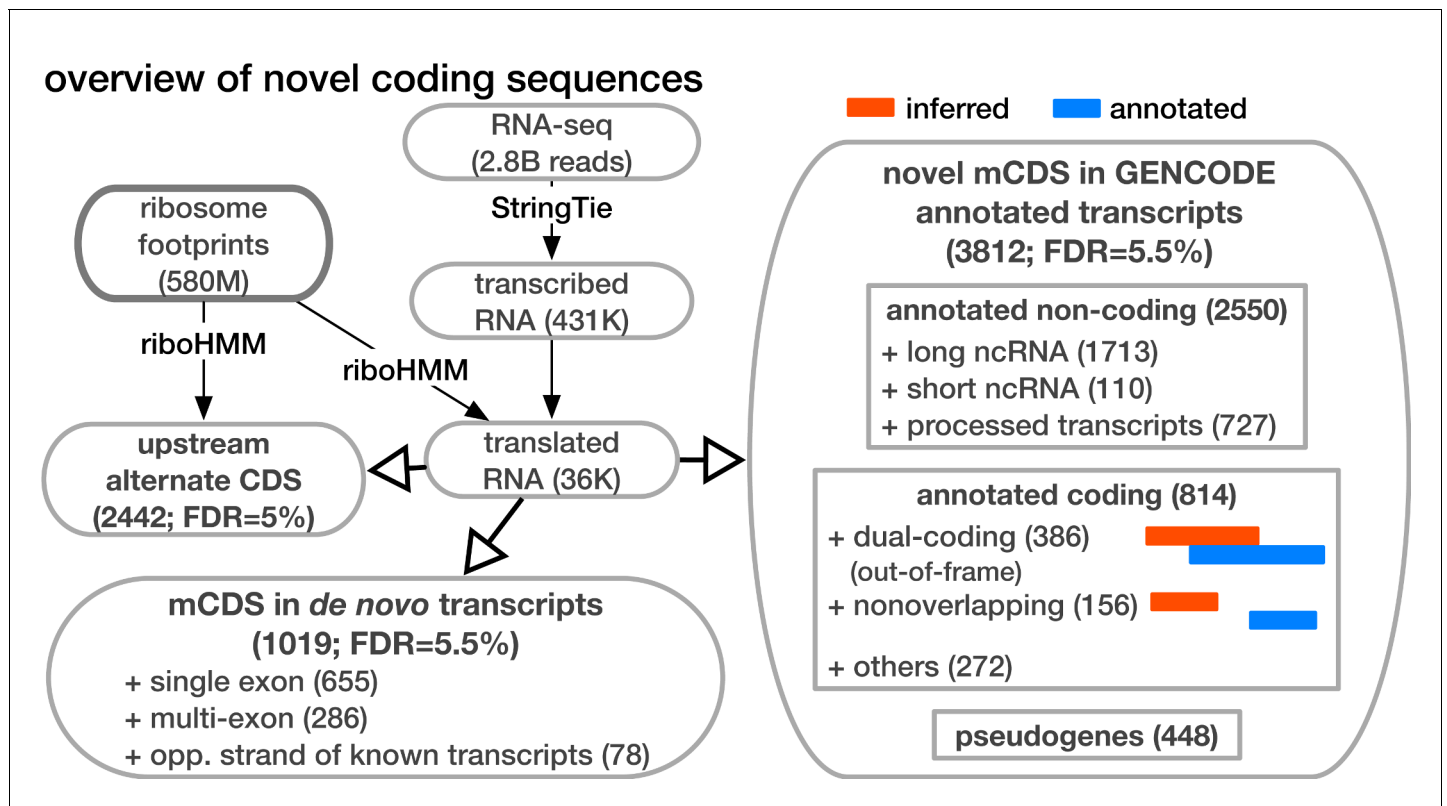


Figure 2. Overview of novel coding sequences. The analysis workflow indicates the size of the data (in read/footprint depth, or number of transcripts) at each step and the numbers and classes of transcript within which novel translated sequences were identified. Transcripts assembled by StringTie that do not overlap any annotated gene are called ‘novel transcripts’. Long non-coding RNA includes lincRNAs, antisense transcripts and transcripts with retained introns, short non-coding RNA includes snRNA, snoRNA and miRNA, processed transcripts are transcripts without a long, canonical ORF, and pseudogenes include all subclasses of such genes annotated by GENCODE.

DOI: [10.7554/eLife.13328.010](https://doi.org/10.7554/eLife.13328.010)

The following figure supplement is available for figure 2:

Figure supplement 1. Decision rules to identify novel mCDS.

DOI: [10.7554/eLife.13328.011](https://doi.org/10.7554/eLife.13328.011)

all isoforms where we inferred an mCDS, the mCDS had a distinct reading frame from the annotated CDS (**Figure 1—figure supplement 3** details the rules that decide how our inference agrees with GENCODE). Of all GENCODE coding genes, we identified 814 GENCODE isoforms where our method identified an mCDS with a distinct reading frame from the annotated CDS. This set of 814 includes both isoforms within the 310 genes and additional isoforms within the 7491 genes (i.e., excluding the isoforms where the mCDS matched the frame of the annotated CDS).

We used simulations to estimate the accuracy of our inferences. For transcripts that do contain a translated sequence, we find that riboHMM inaccurately identifies an overlapping, translated sequence in a different frame at extremely low rates (Type 1 error rate = 0.31%). In contrast, riboHMM has a relatively high error rate for identifying the precise translation initiation site (false discovery proportion = 38%; see Materials and methods for details). Among transcripts where riboHMM correctly identified the reading frame, the concordance between the inferred and annotated translation initiation site does not correlate with the length of CDS (Mann-Whitney test; p-value = 0.12). Amongst these, when riboHMM did not identify the annotated initiation site, the inferred initiation site was equally likely to be upstream or downstream of the annotated initiation site (Mann-Whitney test; p-value = 0.41). Our analysis is robust to sequencing depth; **Figure 1—figure supplement 4** illustrates that nearly 60% of annotated coding sequences identified with the full data set (580 million footprints) could be accurately recovered even when the sequencing depth was reduced by almost two orders of magnitude.

Thus, in summary, it is likely that most of the 814 mCDS that were identified within GENCODE annotated protein-coding transcripts and have a distinct reading frame compared to GENCODE annotations are indeed novel alternate translated sequences. To ensure that an mCDS is truly novel, we verified that it does not overlap any known CDS annotated by GENCODE, UCSC (Rosenbloom et al., 2015), or CCDS (Farrell et al., 2014) in the same frame. (See Figure 2 for the different classes of LCL transcripts that contain a novel mCDS; Figure 2—figure supplement 1 illustrates the decision rules used to identify a novel mCDS). Among these 814 novel mCDS, 386 mCDS overlap an annotated CDS but have a different reading frame (labeled 'dual-coding') and 156 do not overlap the annotated CDS. An example of a novel dual-coding region – an mRNA sequence that codes for proteins in two different frames – inferred in the POLR2M gene is illustrated in Figure 3A. Using tandem mass-spectrometry data (Battle et al., 2015), we found four unique spectra matching peptides in the mCDS and no spectra matching peptides in the annotated CDS (protein level posterior error probability = 3×10^{-35}). However, our simulations suggest that most, or all, of the 39% of genes where riboHMM infers the annotated reading frame but disagrees with the annotated start site are false discoveries, and these are not considered further here.

In addition, we identified 2550 novel mCDS in annotated non-coding transcripts and 1019 mCDS within novel transcripts assembled *de novo* by StringTie (FDR = 5.6%). Using simulations, we estimated that given a transcript has no translated sequence; riboHMM inaccurately identifies a translated sequence at fairly low error rates (Type I error rate = 4.5%). Over 60% of the mCDS in novel transcripts were identified in single-exon transcripts transcribed from regions containing no annotated genes, while about 8% were identified in novel transcripts transcribed from the strand opposite to an annotated transcript. Finally, we inferred mCDS in 448 expressed pseudogenes, among 14,065 pseudogenes annotated in humans (Pei et al., 2012); nearly 90% of these mCDS were identified in processed pseudogenes. An mCDS in pseudogene GAPDHP72 is shown in Figure 3—figure supplement 1, comparing the ribosome abundance and peptide matches to the pseudogene mCDS with those of its parent gene GAPDH.

Unlike current CDS annotations, which almost exclusively start at the methionine codon AUG, these novel mCDS taken together have a substantially higher usage of non-canonical codons, particularly CUG (Figure 3B), consistent with recent observations in mouse embryonic stem cells (Ingolia et al., 2011) and human embryonic kidney cells (Lee et al., 2012). This is despite the fact that we inferred the initiation site by assuming shared properties between novel and annotated CDS. Although riboHMM has a high error rate when identifying translation initiation sites, our use of a hierarchical model for the initiation sites suggests that the errors in our inferred start codons are likely to be unbiased. These novel mCDS are also significantly shorter than annotated CDSs (median lengths 23 vs. 339 amino acids, Mann-Whitney test p-value < 2.2×10^{-16} ; Figure 3C). The overall amino acid content within novel mCDS is comparable to that within annotated CDS, with a slight enrichment for arginine, alanine, cysteine, glycine, proline, and tryptophan residues (binomial test, p-value < 1.1×10^{-16} ; Figure 3—figure supplement 2).

Below, using an alternative measure of ribosome occupancy, we first assess independent evidence for translation initiation at many of these novel mCDS. Then, we test whether these mCDS are functional both using human polymorphism data and using substitution patterns across vertebrates. Finally, we characterize those mCDS whose peptide products were identified in mass-spectrometry data.

Translation at novel mCDS validated using harringtonine-treated ribosome footprints

We next sought to provide independent experimental validation for the novel mCDS. A direct approach to validate translation initiation sites is to assay ribosome occupancy in cells treated with harringtonine (Ingolia et al., 2011). Harringtonine interacts with and arrests the initiation complex while leaving the elongation complex to continue translating and run off the transcript. Harringtonine-treated ribosome footprint profiling data therefore show a specific enrichment pattern at the translation initiation site; this pattern has previously been used to identify translation initiation sites in mouse embryonic stem cells (Ingolia et al., 2011). We measured harringtonine-treated ribosome footprints in two LCLs and aggregated the counts of footprints across all novel mCDS. We observed an enrichment of footprints at the inferred initiation site of the novel mCDS (binomial test, p-value = 9.5×10^{-79} ; Figure 4), similar to the enrichment of aggregate ribosome occupancy at the

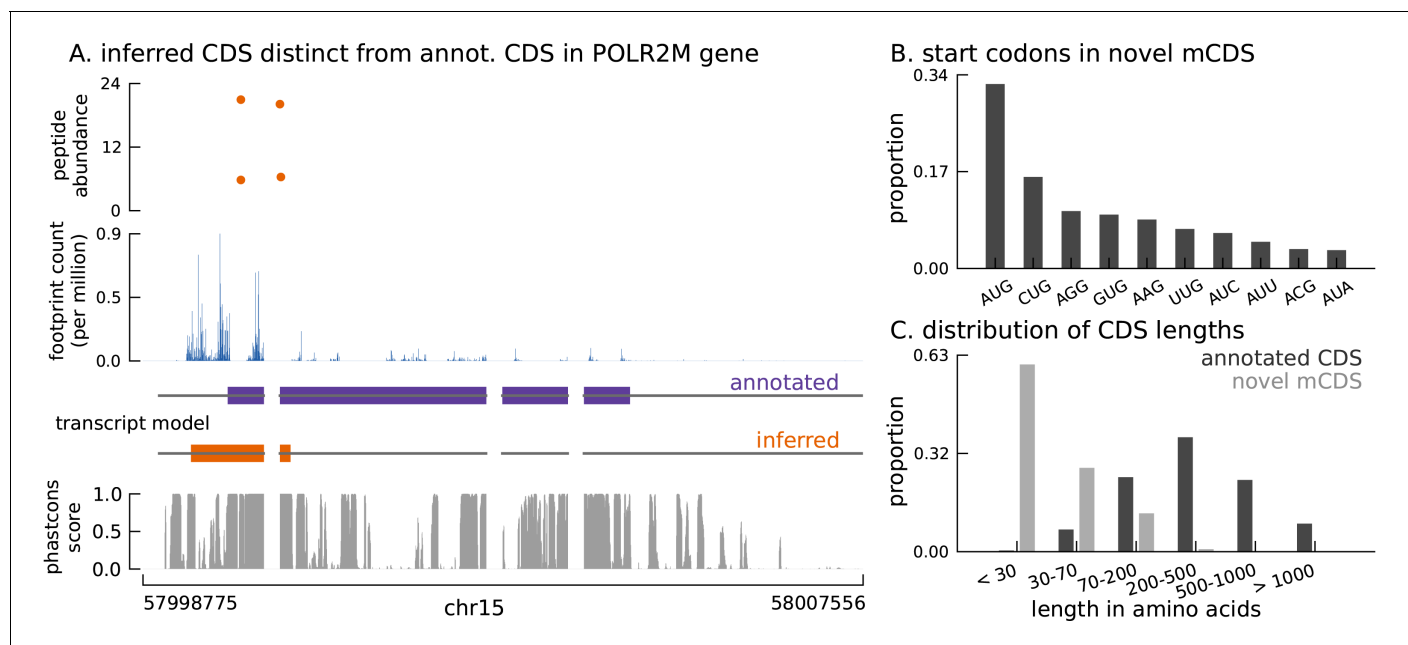


Figure 3. Thousands of novel translated sequences identified in annotated and novel transcript isoforms. (A) The inferred CDS for an isoform of the POLR2M gene overlaps its annotated CDS and is in a different frame. All four distinct peptides uniquely mapping to this gene match the inferred CDS (protein-level posterior error probability = 3×10^{-35}). (The introns and the last exon have been shortened for better visualization.) (B) Distribution of start codon usage across all novel mCDS. (C) Distribution of the lengths of the novel mCDS (gray) compared with the lengths of GENCODE annotated CDSs (black).

DOI: [10.7554/eLife.13328.012](https://doi.org/10.7554/eLife.13328.012)

The following figure supplements are available for figure 3:

Figure supplement 1. Translated coding sequences identified in hundreds of pseudogenes.

DOI: [10.7554/eLife.13328.013](https://doi.org/10.7554/eLife.13328.013)

Figure supplement 2. Comparing the amino acid content between annotated and novel CDS.

DOI: [10.7554/eLife.13328.014](https://doi.org/10.7554/eLife.13328.014)

Figure supplement 3. Characteristics of peptides matched to novel CDS.

DOI: [10.7554/eLife.13328.015](https://doi.org/10.7554/eLife.13328.015)

Figure supplement 4. Annotated genes with peptide hits tend to be longer, have higher expression and a distinct amino acid composition.

DOI: [10.7554/eLife.13328.016](https://doi.org/10.7554/eLife.13328.016)

initiation sites of a matched number of mCDS that agreed exactly with the annotated CDS (see **Figure 4—figure supplement 1** for mCDS in pseudogenes). We observed a significant enrichment at both AUG ($p\text{-value} = 5.2 \times 10^{-79}$) and non-AUG ($p\text{-value} = 9.4 \times 10^{-25}$) initiation sites. The reduced enrichment for the novel mCDS compared to annotated CDSs is likely due to the lower levels of translation of these novel mCDS and the high error rate in identifying the precise base at which translation is initiated. Accounting for these limitations, our observation of enrichment suggests that ribosomes do initiate the translation of many of the novel mCDS identified by riboHMM.

Selective constraint on coding function in novel mCDS

We next ascertained the functional importance of these novel mCDS based on the selective constraint imposed on random mutations that occur within them. A bi-allelic single nucleotide polymorphism (SNP) that falls within an mCDS can be inferred as synonymous or nonsynonymous depending on whether switching between the two alleles of the SNP changes the amino acid sequence of the mCDS. If the mCDS do not produce proteins that are functionally important, we expect the two classes of variants to have similar selection pressures on average, and thus to segregate at similar frequencies. Only if the novel mCDS produce functionally important peptides do we expect inferred nonsynonymous SNPs to segregate at lower frequencies than inferred synonymous SNPs.

Starting with biallelic SNPs identified using whole genome sequences of 2504 individuals (**Auton et al., 2015**), we examined the set of SNPs falling within all novel mCDS (13,907 variants

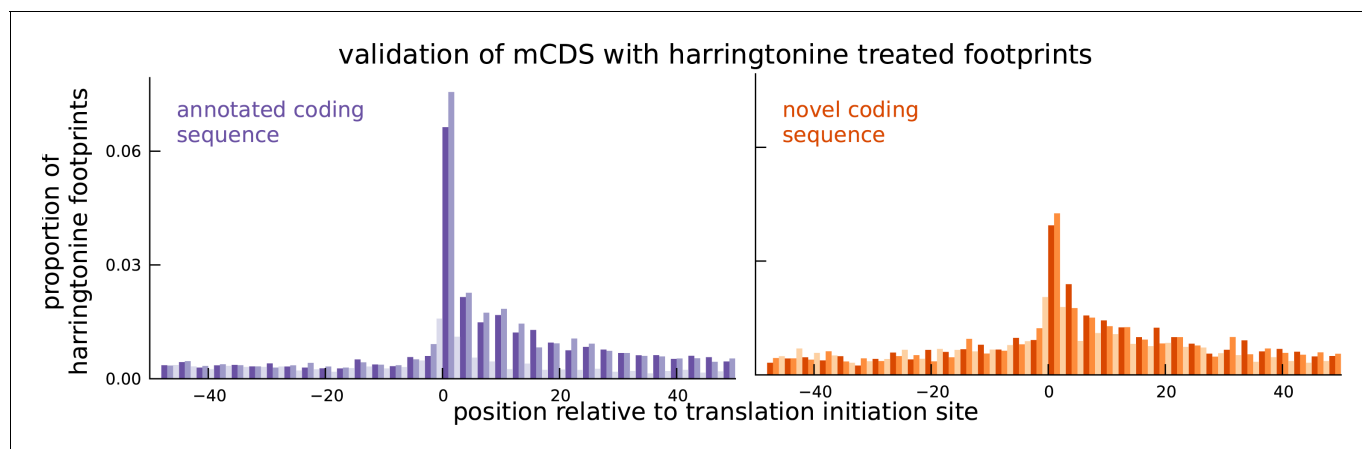


Figure 4. Validation of novel mCDS using harringtonine-treated ribosome profiling data. Harringtonine-treated ribosome footprints show enrichment at the inferred translation initiation sites, when aggregated across novel mCDS (orange), similar to the enrichment at the initiation sites of a matched number of mCDS that agreed exactly with the annotated CDS (purple), suggesting that ribosomes do initiate translation of the novel mCDS.

DOI: [10.7554/eLife.13328.017](https://doi.org/10.7554/eLife.13328.017)

The following figure supplement is available for figure 4:

Figure supplement 1. Validation of translated sequences identified in pseudogenes.

DOI: [10.7554/eLife.13328.018](https://doi.org/10.7554/eLife.13328.018)

within 3096 novel mCDS). We labeled each SNP as synonymous or nonsynonymous with respect to the inferred CDS and show the cumulative distribution of minor allele frequencies (MAF) of each SNP class (**Figure 5A**). We observed that nonsynonymous SNPs have an excess of rare variants compared with synonymous SNPs (Mann-Whitney test; p -value = 1.08×10^{-4}), implying a difference in the intensity of purifying selection (**Nielsen, 2005**). This observed excess suggests that the novel mCDS are under significant constraint, consistent with functional peptides, albeit weaker than at annotated CDS. The mCDS identified within pseudogenes alone also showed a similar excess of rare variants among nonsynonymous SNPs (Mann-Whitney test; p -value = 5.6×10^{-3}). Such an excess was not observed for pseudogenes that had detectable ribosome occupancy but lacked a high-confidence inferred coding sequence (**Figure 5—figure supplement 1**); for these pseudogenes, the SNPs were labeled based on the reading frame of the parent gene. This highlights that ribosome occupancy alone is insufficient to identify translated sequences, and our method is able to leverage finer scale structure in ribosome footprint data to detect functional coding sequences.

While the allele frequency spectra provide evidence that some of the novel mCDS are functional in present-day human populations, they are less informative about the long-term selective constraint on these sequences. To identify whether the novel mCDS have been under long-term functional constraint, we compared the substitution rates at synonymous and nonsynonymous sites within the novel mCDS using whole-genome multiple sequence alignments across 100 vertebrates. (We excluded mCDS identified in pseudogenes from this analysis due to difficulties in assigning orthology.) In **Figure 5B**, 232 novel mCDS have a significantly lower nonsynonymous substitution rate (dN) compared to their synonymous substitution rate (dS) after Bonferroni correction (p -value < 2.91×10^{-5}), suggesting that these mCDS have been under long-term purifying selection. Since the power to detect significantly low values of dN/dS depends on the length of the CDS and the qualities of the genome assemblies and the multiple sequence alignments across distant species at these sequences, the number of functional novel CDSs identified is a conservative lower bound.

Detection of novel proteins by mass spectrometry

We next tested whether we could detect the novel mCDS predictions using mass spectrometry data. We used a large data set of SILAC-labeled tandem mass-spectra generated by trypsin-cleavage of large, stable proteins in many of the same LCLs (**Battle et al., 2015**). Running MaxQuant (**Cox and Mann, 2008**) against the sequence database of 4831 novel mCDS, at 10% FDR, we identified 161 novel mCDS sequences that have at least one unique peptide hit – a tryptic peptide that matches a

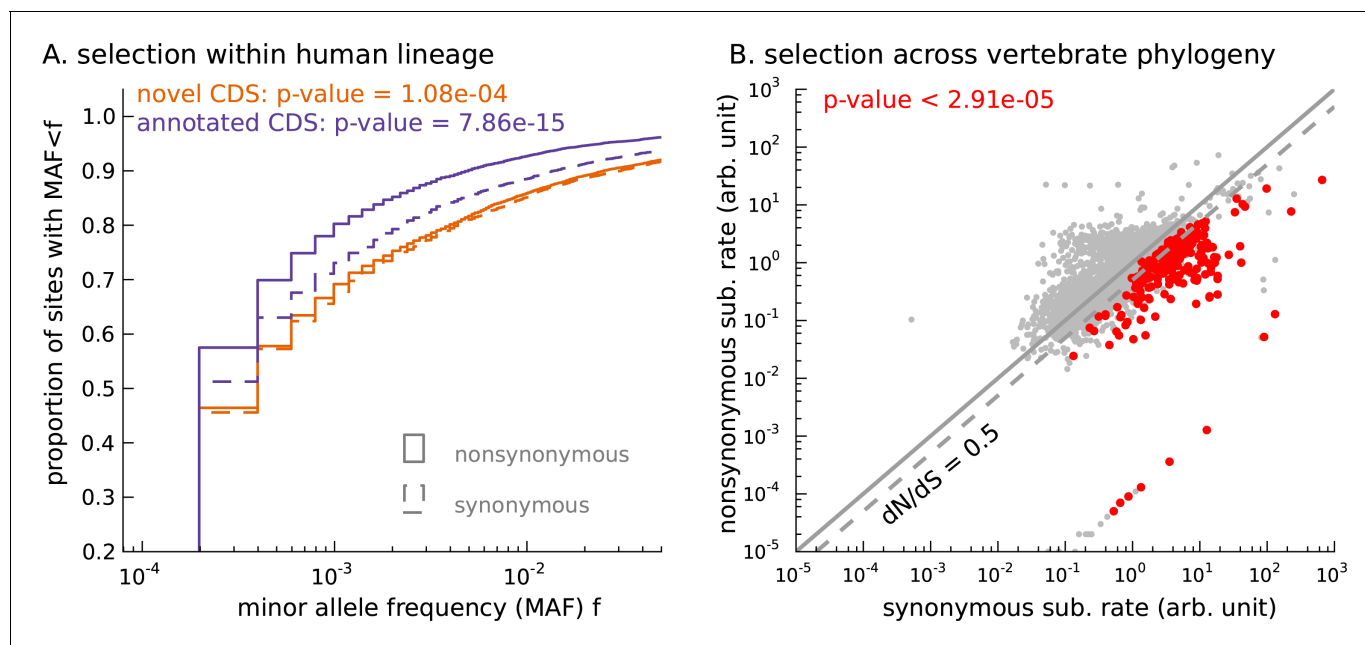


Figure 5. Novel translated sequences show significant signatures of coding function. (A) Genetic variants that are nonsynonymous with respect to the inferred mCDS segregate at significantly lower frequencies in human populations than synonymous variants. The novel regions are under weaker selective constraint compared to known CDS. (The numbers of variants in each class are matched between novel and annotated CDS.) (B) Scatter plot comparing the substitution rate at inferred synonymous variants versus inferred nonsynonymous variants for each novel mCDS, computed using multiple sequence alignments across 100 vertebrate species. Highlighted in red are 232 novel mCDS identified to be under significant long-term purifying selection after Bonferroni correction (testing for $dN/dS < 1$), indicating conserved coding function for these sequences.

DOI: [10.7554/eLife.13328.019](https://doi.org/10.7554/eLife.13328.019)

The following figure supplement is available for figure 5:

Figure supplement 1. Signature of coding function in translated sequences identified in pseudogenes.

DOI: [10.7554/eLife.13328.020](https://doi.org/10.7554/eLife.13328.020)

mass-spectrum (**Supplementary file 1**). More than 70% of novel mCDS with a peptide hit have at least 2 distinct peptides matched to it and, in almost all cases, the unique peptides were independently identified in two or more LCLs (**Figure 3—figure supplement 3**).

To assess how many hits we would expect to the novel mCDS if their properties were like those of annotated CDSs, we developed a model that predicts whether an annotated protein has at least one mass-spectrum match, using features based on expression and sequence composition of the protein (see Materials and methods for more details). The mass-spectrometry data are highly biased towards detection of larger and more highly expressed proteins. Furthermore, the trypsin cleavage step of the experimental protocol imposes strong constraints on the set of unique peptide sequences that can be observed in an experiment. Assuming that the distributions of these predictive features estimated from annotated CDSs can be applied to the novel mCDS, we computed the expected number of novel mCDS with a peptide hit to be 603.

We thus find many fewer mass spectrometry hits to the novel mCDS than expected from a model calibrated on previously annotated mCDS (161 vs. 603). Since our model accounts for translation levels of the mCDS, the low validation rate is unlikely to be due to low rates of protein production. One possible explanation for the low validation rate may be that a large number of the inferred novel mCDS are false discoveries. However, our simulations highlight that our method has a low false positive rate and the Harringtonine data argue that many of the novel mCDS are correct predictions, thus we suggest that some other property of the mCDS may explain their low detection rate. In particular, it is possible that the novel proteins may have higher turnover rates than annotated proteins. For example it is possible that the proteins translated from novel mCDS may have substantially lower half-life than annotated proteins, or may be secreted, and thus have too low concentrations within the cell to be detectable by mass spectrometry assays.

Translation of short alternate coding sequences in addition to the mCDS

Protein-coding transcripts in eukaryotes are typically annotated to have only one CDS (i.e., they are monocistronic). However, a number of studies have demonstrated that ribosomes can initiate translation at alternative start codons (Xu et al., 2010; Ingolia et al., 2011; Lee et al., 2012) and many others have identified transcripts with alternative CDSs encoding functional peptides (Vanderperre et al., 2013; Kochetov, 2008; Barbosa and Romão, 2013). Furthermore, anecdotal evidence has suggested that translation of the alternate CDS serves as a mechanism to suppress translation of the main CDS (Lee et al., 2002; Hernández-Sánchez et al., 2003; Lammich et al., 2004). However, assessing such a mechanism genome-wide has been challenging, mainly due to a lack of appropriate annotations (Calvo et al., 2009).

To this end, we adapted our approach to identify additional coding sequences within transcripts that are translated in LCLs. Assuming that the sub-codon structure of footprint abundance is similar between the main and alternate CDS, we identified 2442 novel CDSs upstream of the mCDS inferred by our method (FDR = 5%); we call them upstream alternate coding sequences or uaCDS (see Materials and methods for details; see also Figure 6—figure supplement 1). Figure 6A illustrates the ribosome footprint density within the uaCDS of the transmembrane gene TM7SF2, and its

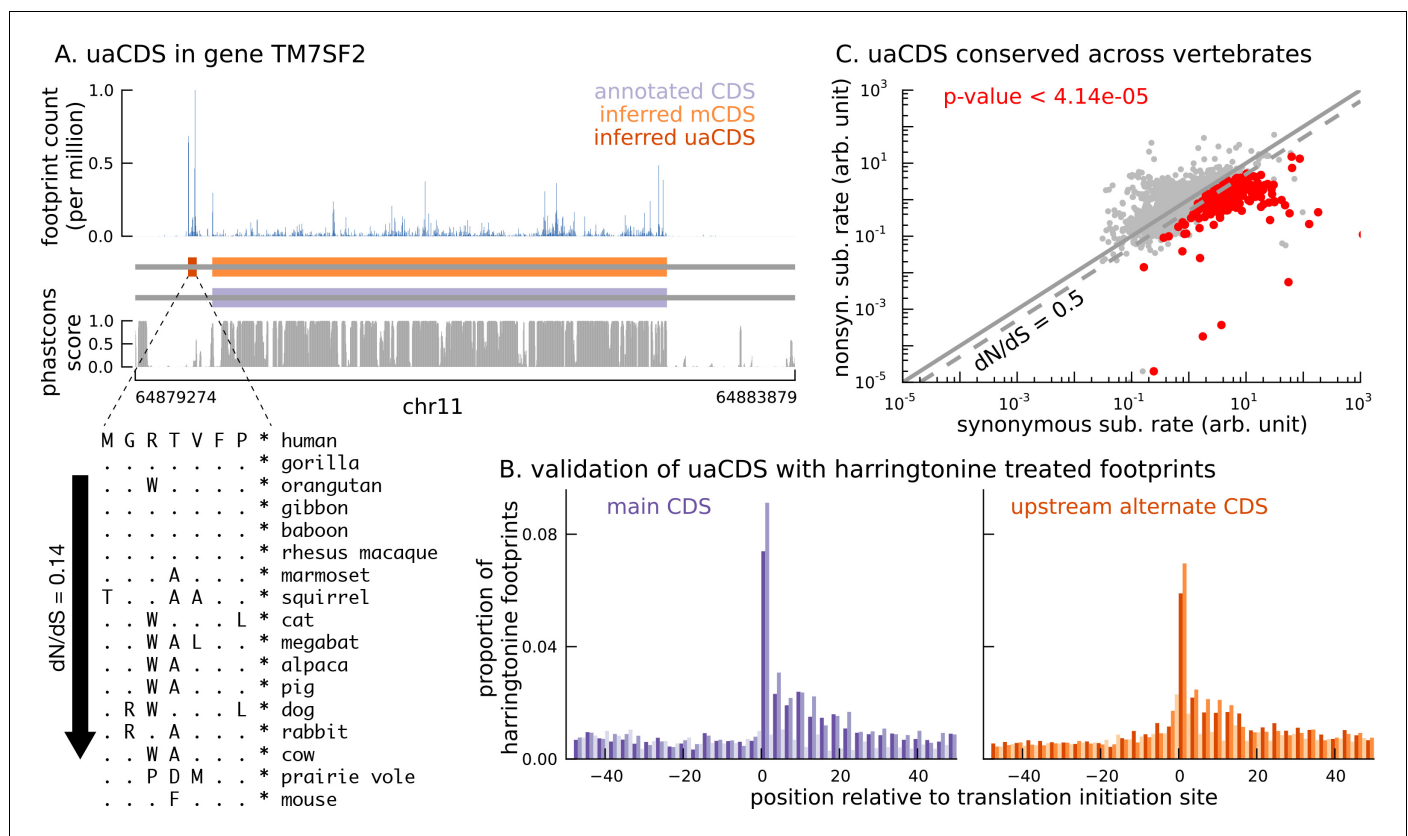


Figure 6. Short translated sequences identified upstream of thousands of translated main coding sequences. (A) An alternate, novel CDS was identified upstream of the inferred main CDS in gene TM7SF2. As shown in its protein sequence alignment across mammals, the uaCDS (in particular, the start and stop codons) is highly conserved with $dN/dS = 0.14$. (B) Harringtonine-treated ribosome footprints show strong enrichment at the inferred initiation sites of uaCDS, comparable to the enrichment at the initiation sites of the corresponding mCDS, suggesting that ribosomes do initiate translation of these uaCDS. (C) Using multiple sequence alignment across 100 vertebrate species, 317 uaCDS were identified to have strong, significant long-term conservation.

DOI: 10.7554/eLife.13328.021

The following figure supplement is available for figure 6:

Figure supplement 1. Characteristics of novel uaCDS.

DOI: 10.7554/eLife.13328.022

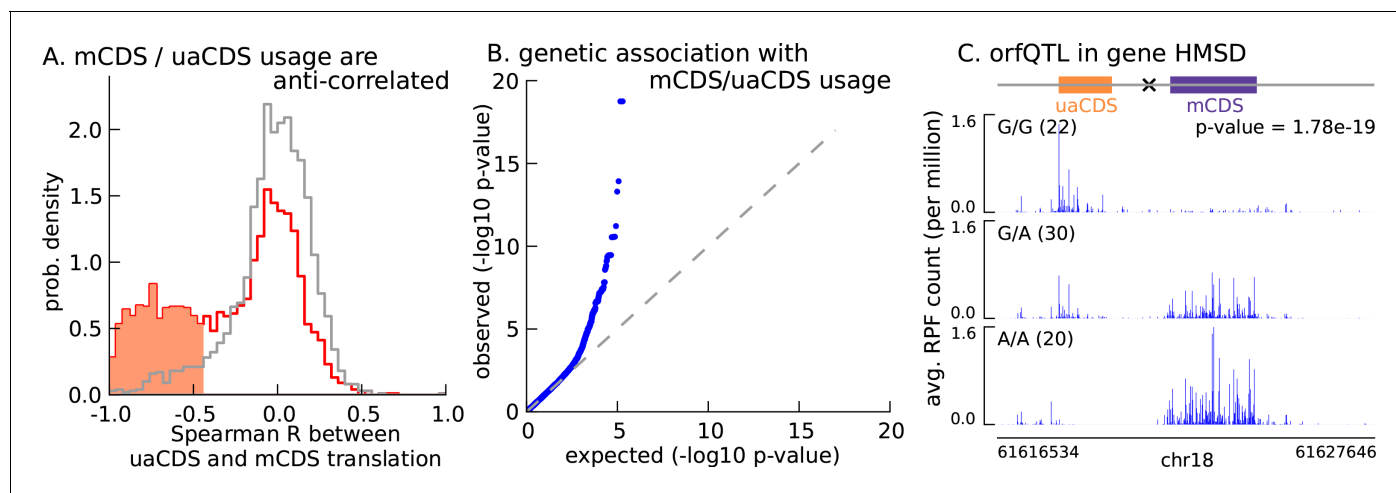


Figure 7. Translation of uaCDS regulates translation of mCDS. (A) Spearman correlation, across LCLs, between mCDS translation and uaCDS translation (red histogram). Using random (mCDS, uaCDS) pairs, matched for length and pairwise distance, we computed an empirical null distribution of Spearman correlations (gray histogram). At 10% FDR, 917 inferred (uaCDS, mCDS) pairs have significant negative correlation (shaded red region). (B) 365 orfQTLs (genetic variants associated with ORF usage; i.e., whether the mCDS or uaCDS of a transcript is translated) were identified at 10% FDR (41 pairs of mCDS/uaCDS). (C) Illustrating an example of an orfQTL in the histocompatibility minor serpin domain-containing (HMSD) gene (introns removed for better visualization). The most significant variant (marked x) lies within an intron between the mCDS and uaCDS of the transcript. DOI: 10.7554/eLife.13328.023

conservation across mammals. We find strong enrichment of harringtonine-treated ribosome footprints at the initiation sites of uaCDS similar to the initiation sites of mCDS in the same transcripts (**Figure 6B**). Using mass-spectrometry data, we identified 46 uaCDS that have at least one peptide hit, substantially lower than the expectation of 891 hits predicted by our model. Finally, comparing the substitution rates at inferred synonymous and nonsynonymous sites, we identified 317 uaCDS with highly constrained coding function (**Figure 6C**). Those uaCDS with a peptide match and those having evidence of constrained coding function are not concordant (Fisher's test; p -value = 0.56), consistent with the low sensitivity of standard mass-spectrometry protocols to identify very short proteins.

Translation of uaCDS negatively correlates with translation of mCDS

With 2442 uaCDS identified as translated in LCLs, we next tested the hypothesis that uaCDS expression negatively correlates with mCDS for each pair. We observed that, at 10% FDR, 917 pairs of uaCDS and mCDS had significant negative correlations across individuals between the proportion of footprints assigned to them (**Figure 7A**). Our observation that nearly 40% of pairs of uaCDS and mCDS are significantly anti-correlated, despite incomplete power due to low sample size, suggests that a potential role of alternate CDSs in a transcript is to regulate the translation of the main CDS.

Variation in ORF usage can be driven by a number of factors including *cis* genetic effects and *trans* effects like variation in expression of RNA binding proteins. To identify *cis* variants that affect ORF usage in a bicistronic transcript, we tested for association of the proportion of RPFs assigned to the mCDS (or uaCDS) with variants in a 10-kilobase window around the transcript; this phenotype effectively controls for variation in gene expression across the LCLs. We identified 365 *cis* orfQTLs (genetic variants associated with ORF usage) across 41 pairs of mCDS and uaCDS at 10% FDR (**Figure 7B**). In **Figure 7C**, we illustrate an example of an orfQTL in a bicistronic transcript of the HMSD gene (histocompatibility minor serpin domain-containing); this gene is also known to have a distinct genetic variant associated with alternative usage of two coding isoforms (*Kawase et al., 2007*). Our observation of orfQTLs in a number of genes distinguishes ORF usage as an additional layer of post-transcriptional regulation of protein expression.

Discussion

We developed riboHMM, a mixture of hidden Markov models to accurately resolve the precise set of mRNA sequences that are being translated in a given cell type, using sequenced RPFs from a ribosome profiling assay, sequenced reads from an RNA-seq assay and the RNA sequence. When applied to human LCLs, this method was able to accurately identify the translated frame in 96% of annotated coding genes that had a high posterior mCDS. In addition, a key advantage of our framework is the ability to infer novel translated sequences that may be missed by annotation pipelines that focus on long CDSs (>100 amino acids), conservation based approaches that require long-term purifying selection, or direct proteomics measurements that are biased toward highly expressed, stable proteins. We used riboHMM to identify 7273 novel CDSs, including 448 of novel translated sequences in pseudogenes and 2442 bicistronic transcripts that contain an upstream CDS in addition to a main CDS. We observed enrichment in harringtonine-arrested ribosome occupancy at the inferred translation initiation sites, suggesting that many of the novel mCDS are real. These novel sequences showed significant differences in the amount of purifying selection acting on inferred non-synonymous versus synonymous sites, suggesting that many of these sequences are conserved as functional peptides, including those mCDS identified in lncRNAs, pseudogenes and novel transcripts.

One caveat of our model is its restriction on one CDS per transcript. In this study, we worked around this limitation using a greedy approach and identified thousands of transcripts with multiple CDSs (either two non-overlapping inferred CDSs or an inferred mCDS distinct from the annotated CDS). Indeed, in some instances where the frame of the mCDS and annotated CDS of a transcript disagreed, we found strong support from mass-spec data for the inferred mCDS frame (**Figure 3A**). These observations highlight the existence of a large number of transcripts in humans that have multiple CDSs and the variation in alternative usage of CDSs across tissues, an area that has largely been overlooked. Additionally, riboHMM does not effectively distinguish footprints arising from different isoforms and, thus, cannot resolve overlapping translated sequences from multiple coding isoforms of a gene. Extending riboHMM to model multiple, possibly overlapping CDSs jointly across multiple isoforms could help uncover this additional layer of complexity in the human genome.

In addition to identifying individual novel coding sequences, our method enables us to observe general properties shared across these coding regions. Interestingly, we found novel coding sequences to have a higher usage of non-AUG start codons than would be expected by considering current translation initiation site annotation (**Figure 3B**). We emphasize that although our model assumes shared properties between novel CDS and annotated CDS, we did not use any information about annotated translation initiation and termination sites when learning the model parameters. We used well-expressed genes as our learning set to ensure that when the footprint data do not provide very strong evidence regarding the initiation site, novel coding sequences identified by our method are as similar as possible to annotated coding sequences in the sequence composition of their initiation sites. While this allows us to be conservative and identify novel CDS that are similar to annotated CDS in their ribosome footprint patterns, our approach will not be able to identify translation events that differ in their footprint patterns from the majority of translation events. In other words, our choice of learning set could potentially bias the inference. Nevertheless, similar start codon usage frequencies were observed when random sets of 5000 genes were used as learning set (**Figure 1—figure supplement 7**) further confirming the robustness of our method.

To improve our ability to identify the translation initiation site, we attempted to incorporate harringtonine treated data in the model by introducing an additional covariate in the transition probabilities, providing independent information on the positions of translation initiation sites. However, the codon usage at the inferred initiation sites showed no significant change (K-S test; p-value = 0.88) and the set of inferred coding sequences showed very little difference when harringtonine data were incorporated into the model. Since the footprint data without treatment show clear enrichment at initiation sites, it is likely that harringtonine treated data do not provide much additional information. Thus, while the harringtonine treated data were useful as independent validation for our inferred initiation sites, the data did not have sufficient additional information to calibrate the confidence in our predicted initiation sites for each transcript.

While the precise function of these novel CDSs remains unclear, we found evidence supporting a regulatory role for novel alternate CDSs identified upstream of the mCDS (uaCDS). Although it is

unclear whether the down regulation of mCDS by uaCDS is dependent on the peptide sequences of uaCDS, our finding is consistent with previous assertions under which translation of upstream ORFs regulates translation of the main CDS in cap-dependent translation initiation (*Morris and Geballe, 2000*).

Our method provides an alternative framework for annotating the coding elements of the genome. Compared to current methods that use sequence information in cDNA and protein databases and those that rely on high-quality genome annotations in closely related species, riboHMM provides a relatively unbiased CDS annotation and opportunities for finding entirely novel CDSs. In particular, one could use riboHMM to identify the set of CDS for a species within a poorly annotated evolutionary clade, using ribosome profiling and RNA seq data immediately after its genome is sequenced and assembled. In addition, given ribosome footprint profiling data from multiple cell types, riboHMM can be used to investigate cell-type-specific translation of coding elements beyond cell-type-specific gene or isoform expression. These features render this tool particularly useful in studying molecular evolution of newly arisen coding genes and linking tissue-specificity of CDS usage to disease.

Materials and methods

Assembling expressed transcripts in LCLs

We mapped paired-end 75 bp RNA-seq reads pooled across 85 Yoruba lymphoblastoid cell lines (*Lappalainen et al., 2013*) to the Genome Reference Consortium Human Reference 37 (GRCh37) assembly using STAR (*Dobin et al., 2013*), with the additional flag `-outSAMstrandField intronMotif` to aid transcript assembly downstream, resulting in 2.8 billion uniquely mapped fragments. Using the mapped reads, we assembled models of transcripts expressed in LCLs using StringTie v1.0.4 (*Pertea et al., 2015*), and used GENCODE v.19 transcript models to guide the assembly. In addition, we required that the lowest expressed isoform of a gene have no less than 1% the expression of the highest expressed isoform ($-f 0.01$), and that each exon-exon junction be supported by at least 2 spliced reads ($-j 2$). Since the RNA-seq protocol did not produce strand-specific reads, we treated the forward strand and reverse strand of a transcript model assembled by StringTie as distinct transcripts. Our final set of 430,754 expressed transcripts included 122,168 GENCODE annotated transcript isoforms and 308,586 novel isoforms. (We did not consider novel isoforms of annotated genes identified by StringTie.)

Ribosome footprint profiling

Cell lines used in this study were ordered from Coriell Institute for Medical Research (<https://www.coriell.org>). To verify the identity of each cell line, we used genotype information derived from the sequencing data. To inspect potential contamination by mycoplasma, we used Universal Mycoplasma Detection Kit from ATCC (ATCC 30-1012K). Ribosome footprint profiling experiments and sequencing data processing were performed as previously described (*Battle et al., 2015*), with the exception of a harringtonine treatment step to arrest ribosomes at the sites of translation initiation. Briefly, lymphoblastoid cell lines, GM19204 and GM19238, were cultured at 37°C with 5% CO₂ in RPMI media with 15% FBS. The media were further supplemented with 2 mM L-glutamate, 100 IU/ml penicillin, and 100 µg/ml streptomycin. Right before cell lysate preparation, each culture was treated with 2 µg/ml harringtonine (final concentration in media) for 2 min followed by 100 µg/ml cycloheximide (final concentration in media). For ribosome profiling experiments, ARTseq Ribosome Profiling kit for mammalian cells (RPHMR12126) was used following vendor's instructions. Sephacryl S400 spin columns (GE; 27-5140-01) were used for monosome isolation. Libraries were sequenced on an Illumina HiSeq 2500. For sequencing data processing and mapping, adaptor sequences were removed from the 3' end of each read using the Clipper tool from the FASTX-Toolkit. In addition, the 5' most nucleotide (commonly resulted from non-templated additions) was removed using the Trimmer tool from the FASTX-Toolkit. To increase mapping efficiency, we filtered out sequence reads that mapped to rRNA, tRNA or snoRNA (FASTA files downloaded from Ensembl on 05/02/13) using Bowtie 2, version 2.0.2 (*Langmead and Salzberg, 2012*). Processed reads were aligned to genome build hg19 (human) using TopHat v2.0.6 (*Trapnell et al., 2009*). The mapping step was guided by transcriptome annotations (downloaded from Ensembl on 01/31/13).

Mixture of HMMs to model translated coding sequences

Consider N transcripts where the n^{th} transcript has length of L_n assumed to be a multiple of three ($L_n = 3M_n$; see *Transcripts with length not a multiple of three* for details on how our model handles the remaining one or two base positions when L_n is not a multiple of three). Our data consist of RPF counts $T = (T^n)_{n=1}^N$, RNA sequence $S = (S^n)_{n=1}^N$, and transcript expression $E = (E^n)_{n=1}^N$ (in units of RNA-seq reads per base position per million sequenced reads) on N transcripts, where T^n and S^n are vector quantities and E^n is a scalar aggregated over the entire length of the transcript. Let $T^n = (T_1^n, \dots, T_{L_n}^n)$ and $S^n = (S_1^n, \dots, S_{L_n}^n)$, where T_b^n and S_b^n denote the RPF counts and the base at the b^{th} position in the n^{th} transcript, respectively. We model the footprint data T using a mixture of HMMs that incorporates S and E . Assuming independence across transcripts, the probability of T given S and E is written as $P(T|\Theta, S, E) = \sum_{n=1}^N P(T^n|\Theta, S^n, E^n)$ where Θ denotes the set of model parameters.

Mixture of three reading frames for a transcript

To capture the three-base structure in RPF data within the CDS, we represent each transcript as a sequence of non-overlapping base triplets, some of which potentially represent codons. Since the CDS of the transcript could belong to one of three reading frames (as illustrated in **Figure 1B**), we introduced a latent frame variable, $F^n \in \{1, 2, 3\}$, that specifies the reading frame for the n^{th} transcript. Then, given $F^n = f$, T^n can be represented as a sequence of $M_n - 1$ triplets and three remaining base positions (see **Figure 1B**). Specifically, $T^n|F^n = f := (X_{f,1}^n, \dots, X_{f,(M_n-1)}^n, R_f^n)$, where $X_{f,m}^n = (T_{3m-3+f}^n, T_{3m-2+f}^n, T_{3m-1+f}^n)$ and

$$R_f^n = \begin{cases} (T_{L_n-2}^n, T_{L_n-1}^n, T_{L_n}^n) & \text{if } f = 1 \\ (T_1^n, T_{L_n-1}^n, T_{L_n}^n) & \text{if } f = 2 \\ (T_1^n, T_2^n, T_{L_n}^n) & \text{if } f = 3 \end{cases} \quad (1)$$

The probability of T^n is then given by

$$\begin{aligned} P(T^n|\Theta, S^n, E^n) &= \sum_{f=1}^3 P(T^n|F^n = f, \Theta, S^n, E^n) P(F^n = f|\Theta, S^n, E^n) \\ &= \sum_{f=1}^3 P(X_{f,1}^n, \dots, X_{f,(M_n-1)}^n, R_f^n|F^n = f, \Theta, S^n, E^n) P(F^n = f|\Theta, S^n, E^n) \end{aligned} \quad (2)$$

We assumed that the probability over F^n is independent of S^n and E^n , and is uniform over all three frames, $P(F^n = f|\Theta, S^n, E^n) = \frac{1}{3}$. In addition, we assumed that the RPF data from the sequence of triplets and the RPF data from the three remaining base positions are independent, leading to

$$P(X_{f,1}^n, \dots, X_{f,(M_n-1)}^n, R_f^n|F^n = f) = P(X_{f,1}^n, \dots, X_{f,(M_n-1)}^n|F^n = f) P(R_f^n|F^n = f). \quad (3)$$

(For notation convenience, we have dropped highlighting the dependence of X^n and R^n on Θ , S_n , and E_n .) We modeled the probability of the data from the sequence of triplets, $P(X_{f,1}^n, \dots, X_{f,(M_n-1)}^n|F^n = f)$, using an HMM, and the probability of the data from the remaining positions, $P(R_f^n|F^n = f)$, using a Poisson-gamma model as described below.

HMM for each frame of a transcript

The pattern of RPF count data in triplets depends on whether the triplet is being translated or not. To model these patterns, we assumed that each triplet belongs to one of nine states (see **Figure 1B**): 5' Untranslated State (5'UTS), last untranslated triplet 5' to the CDS (5'UTS+), Translation Initiation State (TIS), state after TIS (TIS+), Translation Elongation State (TES), state before TTS (TTS-), Translation Termination State (TTS), first untranslated triplet 3' to the CDS (3'UTS-), and 3' Untranslated State (3'UTS). The five states (TIS+, TIS, TES, TTS-, TTS) denote translated triplets and

the remaining four states (5'UTS, 5'UTS+, 3'UTS-, 3'UTS) denote untranslated triplets. In particular, the start codon corresponds to the base triplet assigned to the TIS state and the stop codon corresponds to the base triplet assigned to the 3'UTS- state. The groups of states (5'UTS+, TIS, TIS+) and (TTS-, TTS, 3'UTS-) help model the distinct features of the footprint data around the translation initiation and termination sites, respectively. We introduced a sequence of $M_n - 1$ hidden variables $Z_f^n = (Z_{f,1}^n, \dots, Z_{f,(M_n-1)}^n)$ for each frame of the n^{th} transcript, where $Z_{f,m}^n$ denotes the state for the m^{th} triplet in the f^{th} frame.

For each state, an emission probability for $X_{f,m}^n$ can be modeled as follows. Let $Y_{f,m}^n$ denote the sum of three elements in $X_{f,m}^n$ (i.e., the total RPF count for the m^{th} triplet). Then, $P(X_{f,m}^n | Z_{f,m}^n = z) = P(X_{f,m}^n | Y_{f,m}^n, Z_{f,m}^n = z) P(Y_{f,m}^n | Z_{f,m}^n = z)$ and

$$X_{f,m}^n | Y_{f,m}^n, Z_{f,m}^n = z \sim \text{multinomial}(Y_{f,m}^n, \pi_z), \tag{4}$$

$$Y_{f,m}^n | Z_{f,m}^n = z \sim \text{Poisson}(\mu_{zfm}^n E^n), \tag{5}$$

$$\mu_{zfm}^n \sim \text{gamma}(\alpha_z, \beta_z), \tag{6}$$

where the density of the gamma distribution is $P(\mu) = \frac{\beta^\alpha}{\Gamma(\alpha\beta)} \mu^{\alpha\beta-1} \exp^{-\beta\mu}$ with the mean and variance equal to α and $\frac{\alpha}{\beta}$, respectively.

The periodicity of RPF counts within the CDS is captured by the multinomial distribution with parameters $\pi_z = (\pi_{z,1}, \pi_{z,2}, \pi_{z,3})$, where we assume $\pi_z = (\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$ for $z \in \{5'UTS, 5'UTS+, 3'UTS-, 3'UTS\}$ to capture the lack of periodicity in the RPF data in untranslated regions. Furthermore, we allow the pattern of periodicity to differ across five states (TIS, TIS+, TTS, TTS-, TES).

The Poisson distribution for $Y_{f,m}^n$ captures the difference in RPF abundance between translated and untranslated regions (precisely, difference in abundance between triplets in different states). We corrected for differences in RPF abundance across transcripts due to differences in transcript expression levels by using E^n as a transcript-specific normalization factor (see **Figure 1—figure supplement 5**). To account for additional variation in the RPF counts across triplets in the same state (e.g., due to varying translation rates across transcripts, and translational pausing), we allowed for triplet-specific parameters μ_{zfm}^n in the Poisson intensity and assumed that those parameters follow a gamma distribution. Under this model, $E[Y_{f,m}^n | Z_{f,m}^n = z] = \alpha_z E^n$ and $\text{Var}[Y_{f,m}^n | Z_{f,m}^n = z] = \frac{\alpha_z}{\beta_z} E^{n^2} + \alpha_z E^n$.

We assumed that the sequence of hidden variables Z_f^n follow a Markov chain. The assumption of up to one CDS in each transcript leads to a transition probability shown in **Figure 1B**, where $\rho_{f,m}^n = P(Z_{f,m+1}^n = 5'UTS+ | Z_{f,m}^n = 5'UTS)$ and $\zeta_{f,m}^n = P(Z_{f,m+1}^n = TTS- | Z_{f,m}^n = TES)$ depend on the underlying RNA sequence and are given by

$$\rho_{f,m}^n = \begin{cases} \text{logistic} \left(\psi_\kappa K_{f,m+2}^n + \sum_{c \in \Omega_{\text{start}}} \psi_c \mathbb{I}[M_{f,m+2}^n = c] \right), & \text{if } M_{f,m+2}^n \in \Omega_{\text{start}} \\ 0, & \text{otherwise} \end{cases} \tag{7}$$

$$\zeta_{f,m}^n = \begin{cases} 1, & \text{if } M_{f,m+3}^n \in \Omega_{\text{stop}} \\ 0, & \text{otherwise} \end{cases} \tag{8}$$

where $\mathbb{I}[\cdot]$ is the indicator function, $M_{f,m}^n = (S_{3m-3+f}^n, S_{3m-2+f}^n, S_{3m-1+f}^n)$ denotes the base sequence of the m^{th} triplet, and $K_{f,m}^n$ denotes the log of ratio of likelihood under the Kozak model to likelihood under a background model of the base sequence flanking the m^{th} triplet (see *Kozak model* for details). In our analysis, Ω_{start} contained the canonical start codon and all near-cognates, $\Omega_{\text{start}} = \{\text{AUG, CUG, GUG, UUG, AAG, ACG}\}$ and Ω_{stop} contained the canonical stop codons, $\Omega_{\text{stop}} =$

{UAA, UAG, UGA}. The parameters, ψ_c and ψ_κ , indicate the importance of the triplet base sequence and the flanking base sequence in determining transition from untranslated triplets to translated triplets. The current specification of $\zeta_{f,m}^n$ and Ω_{stop} forces the coding sequence to terminate at the first in-frame occurrence of a stop codon. This model can be extended to account for stop codon read-through by using a logistic function for $\zeta_{f,m}^n$ for the same set Ω_{stop} .

Model for R_f^n

We model R_f^n , the RPF counts at bases before or after the sequence of triplets (see **Equation 1**), using the emission probabilities of the 5'UTS or 3'UTS states. Assuming that the three elements of R_f^n are independent, we have $P(R_f^n | F^n = f) = \prod_{i=1}^3 R_{f,i}^n | F^n = f$. Each element can be modeled as

$$R_{f,i}^n \sim \text{Poisson}\left(\frac{1}{3}\lambda_{fi}^n E^n\right), \tag{9}$$

$$\lambda_{fi}^n \sim \text{gamma}(\alpha_z, \beta_z), \tag{10}$$

where $z = 5'$ UTS if $R_{f,i}^n \in \{T_1^n, T_2^n\}$, and $z = 3'$ UTS if $R_{f,i}^n \in \{T_{L_n-2}^n, T_{L_n-1}^n, T_{L_n}^n\}$.

Parameter estimation and inference

We used an EM algorithm to compute the maximum likelihood estimate for the model parameters $\Theta = \{\pi_z, \alpha_z, \beta_z, \psi_\kappa, \psi_c\}$, that is, $\hat{\Theta} = \text{argmax}_\Theta P(T | \Theta, S, E)$.

To infer the translated CDS for the n^{th} transcript, we identified the frame and state sequence that maximizes the joint posterior probability

$$(z^{n^*}, f^{n^*}) := \text{argmax}_{z,f} P(Z_f^n = z, F^n = f | T^n, S^n, E^n, \hat{\Theta}). \tag{11}$$

We first computed the maximum *a posteriori* (MAP) state sequence for each reading frame using the Viterbi algorithm, $z_f^{n^*} := \text{argmax}_z P(Z_f^n = z | F^n = f, T^n, S^n, E^n, \hat{\Theta})$ for $f = 1, 2, 3$. Then, the MAP state sequence and frame is given as

$$(z^{n^*}, f^{n^*}) := \text{argmax}_{z,f} P(Z_f^n = z_f^{n^*} | F^n = f, T^n, S^n, E^n, \hat{\Theta}) P(F^n = f | T^n, S^n, E^n, \hat{\Theta}), \tag{12}$$

where $z_f^{n^*}$ is a function of f , $P(F^n = f | T^n, S^n, E^n, \hat{\Theta}) \propto P(T^n | F^n = f, S^n, E^n, \hat{\Theta}) P(F^n = f)$ and $P(T^n | F^n = f, S^n, E^n, \hat{\Theta})$ is the probability of the data marginalized over the latent states.

In our analyses, we estimated the model parameters using the top five thousand highly expressed genes. Then, we inferred the translated CDS for those transcripts in which each exon has at least five distinct ribosome footprints mapping to it. We restricted our further analyses to transcripts where (1) $P(Z_f^n = z^{n^*}, F^n = f^{n^*} | T^n, S^n, E^n, \hat{\Theta}) > 0.8$, (2) the MAP state sequence z^{n^*} contains a TIS state and a TTS state (i.e., a pair of initiation and termination sites), (3) more than 50% of base positions within the inferred CDS are mappable, and (4) the coding sequence encodes a peptide more than 6 amino acids long – we call these translated sequences as main coding sequences or mCDS.

Modeling ribosome footprints of different lengths

We observed that ribosome footprints with different lengths, arising due to incomplete nuclease digestion, show slightly different patterns of abundance when aggregated across transcripts (see **Figure 1—figure supplement 6**). To model these differences, we partitioned the footprints into multiple groups based on length, and modeled the data in each group with a separate set of parameters in the emission probability (all groups share the same state sequence along a transcript). Specifically, for G groups of footprints, the data at the m^{th} triplet in f^{th} reading frame $X_{f,m}^n$ can be partitioned into G components, $X_{f,m}^n = \left(X_{g, fm}^n\right)_{g=1}^G$, where $X_{g, fm}^n$ denotes the triplet of RPF counts from g^{th} group. Assuming that the RPF counts from different groups at a given triplet are

independent, conditional on the state of the triplet, the emission probability can be written as $P(X_{f,m}^n | Z_{f,m}^n = z) = \prod_{g=1}^G P(X_{g,fm}^n | Z_{f,m}^n = z)$ and

$$X_{g,fm}^n | Y_{g,fm}^n, Z_{f,m}^n = z \sim \text{multinomial}(Y_{g,fm}^n, \pi_{g,z}), \quad (13)$$

$$Y_{g,fm}^n | Z_{f,m}^n = z \sim \text{Poisson}(\mu_{g,z,fm}^n E^n), \quad (14)$$

$$\mu_{g,z,fm}^n \sim \text{gamma}(\alpha_{g,z}, \beta_{g,z}), \quad (15)$$

where group-specific parameters, $(\pi_{g,z}, \alpha_{g,z}, \beta_{g,z})$, capture the distinct patterns in each group. The RPF data used in our analyses had four groups of footprints of lengths 28, 29, 30, and 31 bases.

Base positions with missing data

Approximately 15% of the transcriptome have unmappable base positions, in part due to the short lengths of ribosome footprints. Consider the m^{th} base triplet in frame f in the n^{th} transcript. If $J_{g,fm}^n$ is the set of positions in this triplet that are unmappable for footprints corresponding to group g , the emission probabilities become

$$X_{g,fm}^n | Y_{g,fm}^n, Z_{f,m}^n = z \sim \text{multinomial}(Y_{g,fm}^n, \tilde{\pi}_{g,z}), \quad (16)$$

$$Y_{g,fm}^n | Z_{f,m}^n = z \sim \text{Poisson}(\psi_{g,z,fm}^n \mu_{g,z,fm}^n E^n), \quad (17)$$

$$\mu_{g,z,fm}^n \sim \text{gamma}(\alpha_{g,z}, \beta_{g,z}), \quad (18)$$

where

$$\psi_{g,z,fm}^n = \sum_{j \notin J_{g,fm}^n} \pi_{g,z,j}, \quad (19)$$

$$\tilde{\pi}_{g,z,j} = \begin{cases} 0 & \text{if } j \in J_{g,fm}^n \\ \frac{\pi_{g,z,j}}{\psi_{g,z,fm}^n} & \text{otherwise.} \end{cases} \quad (20)$$

If all three positions in a triplet are unmappable, then we treat the triplet as having missing data for that footprint group and set $P(X_{g,fm}^n | Z_{f,m}^n) = 1$ for all values of $Z_{f,m}^n$.

Kozak model

Using the annotated initiation sites of GENCODE annotated coding transcripts, we estimated a position weight matrix (PWM) that captures the base composition of the -9 to $+6$ positions flanking known initiation sites. Since the consensus sequence of this PWM is the same as the reported consensus Kozak sequence (Kozak, 1987), we refer to this model as the Kozak model. We estimated a background PWM model using the same set of positions relative to random AUG triplets within the same set of transcripts. For the m^{th} triplet in frame f in the n^{th} transcript, using the base sequence from the -9 to $+6$ positions flanking this triplet, we computed $K_{f,m}^n$, the log of ratio of likelihood of the flanking sequence under the Kozak model to likelihood under the background model.

Transcripts with length not a multiple of three

The length of such a transcript can be written as $L_n = 3M_n + B$, where $B \in \{1, 2\}$. We assumed that the RPF data on the first $3M_n$ bases $(T_{1:3M_n}^n)$ and the data on the remaining B bases $(T_{3M_n+1:L_n}^n)$ are independent. We modeled $T_{1:3M_n}^n$ using a mixture of HMMs as described above, and modeled $T_{3M_n+1:L_n}^n$ using the emission probability of the 3'UTS state as follows.

$$P\left(T_{3M_n+1:L_n}^n | E^n, \alpha_z, \beta_z\right) = \prod_{m=3M_n+1}^{L_n} P\left(T_m^n | E^n, \alpha_z, \beta_z\right), \quad (21)$$

$$T_m^n \sim \text{Poisson}\left(\frac{1}{3}\tau_m^n E^n\right), \quad (22)$$

$$\tau_m^n \sim \text{gamma}(\alpha_z, \beta_z), \quad (23)$$

$$z = 3'UTS$$

A Python implementation of riboHMM can be downloaded from <https://rajanil.github.io/riboHMM/>.

Quantifying false discoveries of riboHMM

We characterize the performance of riboHMM by addressing three scenarios: (1) How often does riboHMM identify an mCDS in transcripts with no signal of translation? (2) How often does riboHMM identify an incorrect reading frame in transcripts with signal for translation? (3) When riboHMM identifies the correct reading frame in transcripts with signal for translation, how often does it identify an incorrect initiation site? To address the first question, we started with the transcripts for which riboHMM was able to identify an mCDS and generated a set of “null transcripts” by permuting the footprint counts among base positions within each transcript. Applying a posterior cutoff of 0.8, riboHMM incorrectly identified an mCDS in 4.5% of these null transcripts. We used this estimate of the Type 1 error rate to compute the false discovery rate for novel mCDS in noncoding transcripts and novel uaCDS identified by riboHMM. To address the other two questions, we started with the set of annotated coding transcripts for which riboHMM was able to recover the precise CDS (i.e., the mCDS matched the annotated CDS exactly). We generated a set of “simulated transcripts” using the following strategy: (1) randomly select a new TIS downstream and in-frame to the annotated TIS, ensuring that the codon underlying the new TIS belonged to the set Ω_{start} , (2) permute the footprint counts among bases upstream of the new TIS. Among the simulated transcripts in which riboHMM could identify an mCDS, the inferred reading frame was completely different from the true translated reading frame in 0.31% transcripts. We used this estimate of the Type 1 error rate to quantify false discoveries among novel mCDS in annotated coding transcripts. In the remaining simulated transcripts, the inferred TIS matched the new TIS exactly in 62% of transcripts; this corresponds to a false discovery proportion of 38%.

Translated mCDS in pseudogenes

Starting with 14,065 pseudogenes that have been identified and categorized in humans (*Pei et al., 2012*), 9,375 pseudogenes were identified by StringTie to be expressed in LCLs. Using a very stringent posterior cutoff of 99.99%, we inferred mCDS in 448 of these expressed pseudogenes. Using pairwise alignment of the pseudogene and parent gene transcript, we observed that although the pseudogene mCDS typically code for shorter protein sequences compared with the parent protein, a large fraction of the pseudogene mCDS share coding-frame with their parent gene (see *Figure 3—figure supplement 1*).

Validation with Harringtonine-treated data

Harringtonine-treated ribosome footprints were measured in LCLs with a total sequencing depth of 21 million reads. In *Figure 4*, we illustrate the aggregate proportion of treated ribosome footprints centered at the inferred start codon for all novel mCDS, and compare it with the aggregate proportion of treated footprints around the start codon of an equal number of annotated CDSs that have a posterior probability greater than 0.8 under our model. In *Figure 4—figure supplement 1*, we illustrate the aggregate proportion of treated footprints for mCDS inferred in pseudogenes alone, and in *Figure 6B*, we compare the aggregate treated footprint proportions at the start codons of inferred uaCDS and their corresponding mCDS.

Identifying translated alternate ORFs

For each transcript that had a mCDS with posterior greater than 0.8 and more than 50 base pairs of RNA sequence in the 5'UTS state, we defined an "upstream-restricted transcript" consisting of the exons within the 5'UTS state. Using a random set of 5000 non-overlapping upstream-restricted transcripts in which more than 80% of base positions were mappable, we computed the maximum likelihood estimates of the transition parameters and occupancy parameters to identify additional translated sequences within these upstream-restricted transcripts. Assuming that the fine-scale structure of footprint counts within these translated sequences would be similar to that within the mCDS, we kept the periodicity parameters fixed to their previously estimated values. With these parameter estimates, we inferred the MAP frame and state sequences with posterior greater than 0.8 and filtered out inferences where less than 50% of the inferred CDS was mappable. These additional translated sequences within the upstream-restricted transcripts were called upstream alternate coding sequences or uaCDS.

Identifying stable peptides with mass spectrometry data

To identify stable proteins translated from the novel CDSs (mCDS and uaCDS), we analyzed quantitative, high-resolution mass spectrometry data derived from 60 LCLs, with MaxQuant v1.5.0.30 (Cox and Mann, 2008) and the Andromeda (Cox et al., 2011) search engine. Sample labeling, processing and data collection details can be found elsewhere (Battle et al., 2015; Khan et al., 2013). Peptides were identified using a database that contained 63,904 GENCODE annotated protein sequences and 7271 novel CDSs identified by our method. For all searches, up to two missed tryptic cleavages were allowed, carbamidomethylation of cysteine was entered as a fixed modification, and N-terminal acetylation and oxidation of methionine were included as variable modifications for all searches. A 'first search' tolerance of 40 ppm with a score threshold of 70 was used for time-dependent mass recalibration followed by a main search MS1 tolerance of 6 ppm and an MS2 tolerance of 20 ppm. The 're-quantify' option was used to aid the assignment of isotopic patterns to labeling pairs. The 'match between runs' option was enabled to match identifications across samples using a matching time window of 42 s and an alignment time window of 20 min. Peptide and protein false discovery rates were set to 10% using a reverted version of the search database. Protein group quantifications were taken as the median \log_2 (sample/standard) ratio for all groups containing at least two independent unique or 'razor' peptide quantifications (including multiple measurements of the same peptide in different fractions) without a modified peptide counterpart.

Bias correction to compute expected number of peptide hits

Proteins with at least one peptide identified by this high-resolution mass-spectrometry protocol tend to be distinct from proteins with no mass-spectrum matches.

1. The median footprint density of annotated coding genes with at least one peptide match is about 125 fold higher than that of coding genes with no peptide match (see [Figure 3—figure supplement 4A](#)).
2. The median length of coding genes with at least one peptide match is 20% higher than that of coding genes without a peptide match (see [Figure 3—figure supplement 4B](#)).
3. The trypsin cleavage step of the protocol ensures that nearly all observable peptides have a C-terminal lysine or arginine residue, and up to two additional lysine or arginine residues within the peptide sequence (called "tryptic peptides"). This step imposes a strict constraint on the set of unique peptide sequences that can be observed from a protein sequence, and genes with fewer tryptic peptides are less likely to have a mass-spectrum match.
4. All tryptic peptides in an expressed protein are not equally likely to be observed. The probability of detecting a tryptic peptide depends on its electrostatic properties relative to other tryptic peptides from all expressed proteins, which in turn depends on the amino acid composition of the tryptic peptides (see [Figure 3—figure supplement 4C–F](#)).

To account for these biases, we developed a predictive model to estimate the probability that a protein has at least one peptide hit in a mass-spectrometry experiment. The predicted label for a protein is whether the protein has at least one mass-spectrum match ($H_n = 1$) or no mass-spectrum match ($H_n = 0$). The predictive features of a protein used in the model are (1) the ribosome footprint density of the corresponding transcript (D_n), (2) the protein length (S_n), and (3) the counts of amino

acids within each of the K tryptic peptides that can be generated from the protein ($L_n = \{L_{n1}, \dots, L_{nK}\}$). Since the relevant feature of an amino acid is its charge, we partitioned the set of amino acids into four groups – positively charged (R, H, K), negatively charged (D, E), polar uncharged (S, T, N, Q), and others. The amino acid count vector L_{nk} was then collapsed into a vector of the counts of each of these four groups. Conditional on $H_n = 1$, we introduced a latent variable for each tryptic peptide that indicates whether the peptide was matched to a mass-spectrum or not ($Z_{nk} \in \{1, 0\}$); this latent variable accounts for differences between matched and unmatched peptides.

Assuming that the three predictive features are independent conditional on the predicted label H_n , the odds of observing at least one peptide hit is then given as

$$\frac{p(H_n = 1 | D_n, S_n, L_n)}{p(H_n = 0 | D_n, S_n, L_n)} = \frac{p(D_n | H_n = 1) p(S_n | H_n = 1) \left\{ \sum_{Z_n} p(L_n | Z_n, H_n = 1) p(Z_n | H_n = 1) \right\} p(H_n = 1)}{p(D_n | H_n = 0) p(S_n | H_n = 0) p(L_n | Z_n = 0, H_n = 0) p(H_n = 0)}$$

We learn the predictive model using annotated coding genes and partitioning them into those that have at least one peptide hit (“hit genes”) and those that do not have a peptide hit (“no-hit genes”). We computed $p(D_n | H_n)$ using an empirical distribution of footprint density within coding genes, $p(S_n | H_n)$ using an empirical distribution of the lengths of coding genes, $p(L_{nk} | Z_{nk} = 1, H_n = 1)$ using tryptic peptides within hit genes matched to mass-spectra, $p(L_{nk} | Z_{nk} = 0, H_n = 1)$ using unmatched tryptic peptides within hit genes, and $p(L_{nk} | Z_{nk} = 0, H_n = 0)$ using tryptic peptides within no-hit genes. Finally, we set $p(H_n = 1) = p(H_n = 0) = 1/2$ and $p(Z_n | H_n = 1) = 1/(2^K - 1)$. Using peptide hits in annotated proteins, we evaluated the accuracy of this model by holding out some annotated proteins as test data, learning the predictive distributions using the remaining training data and computing the expected number of test proteins that had a mass-spectrum match. We estimated the expected number of held-out annotated proteins with at least one mass-spectrum match to be 1206 (s.d. = 34), while the actual number of held-out proteins with a match was 1387 (s.d. = 36).

Test for long-term purifying selection

In order to quantify whether the novel mCDS are evolutionary conserved in terms of their amino acid sequence, we first extracted DNA sequences orthologous to the mCDS from a 100-way vertebrate whole-genome alignments (UCSC), restricting to genomes aligned with either Syntenic net or Reciprocal best net. We next performed a three-frame translation on each orthologous sequence and a multiple alignment to obtain the correct codon alignments. More specifically, for each orthologous sequence, we kept the frame with the highest amino acid identity compared to the human peptide, requiring at least 60% identity for alignable positions and no more than 50% of the alignment as gaps. Finally, we used codeML/PAML (Yang, 2007) to estimate dN and dS rates across the trees consisting of all remaining peptides, first using a model allowing variable omega and then a model with omega fixed to one. To determine whether a specific peptide is under purifying selection or not, we compared the two models using a likelihood ratio test and reported peptides that satisfied a Bonferroni-corrected p -value threshold.

Correlation between uaCDS and mCDS

We computed the correlation across LCLs between the proportion of footprints mapped to a transcript that fall within its uaCDS and the proportion that fall within its mCDS. We evaluated the statistical significance of these correlations using an empirical null distribution of Spearman correlations computed using random pairs of mCDS and uaCDS. A random pair of mCDS and uaCDS was obtained by randomly shifting the coordinates of an observed pair of mCDS and uaCDS, matching for their respective lengths and the distance between them.

Data release

All novel coding sequences identified in this work, along with the harringtonine-treated ribosome profiling data are deposited in GEO Accession GSE75290.

Acknowledgements

We thank Adam Frankish from the GENCODE group for discussions on the sources of information used for annotating coding genes, Zia Khan for discussions on the mass-spectrometry data analysis, and Audrey Fu for discussions on evaluating the correlation between mCDS and uaCDS. We also thank members of the Pritchard, Gilad and Stephens labs for comments and suggestions, and the two anonymous reviewers, the Reviewing Editor, and Naama Barkai, the Senior Editor for their insightful reviews and comments. This work was funded by grants from the NIH (HG007036 to JKP, MH084703 to YG and JKP, and HG02585 to MS), and by the Howard Hughes Medical Institute. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Additional information

Funding

| Funder | Grant reference number | Author |
|---------------------------------|------------------------|------------------------------------|
| National Institutes of Health | HG007036 | Jonathan K Pritchard |
| National Institutes of Health | MH084703 | Yoav Gilad Jonathan K Pritchard |
| National Institutes of Health | HG02585 | Matthew Stephens |
| Howard Hughes Medical Institute | | Jonathan K Pritchard |

The funders had no role in study design, data collection and interpretation, or the decision to submit the work for publication.

Author contributions

AR, HS, MS, YG, JKP, Conception and design, Analysis and interpretation of data, Drafting or revising the article; SHW, Conception and design, Acquisition of data, Analysis and interpretation of data, Drafting or revising the article; AH, YIL, BE, Analysis and interpretation of data, Drafting or revising the article

Author ORCIDs

Anil Raj, <http://orcid.org/0000-0003-4412-0883>

Yang I Li, <http://orcid.org/0000-0002-0736-251X>

Brett Engelmann, <http://orcid.org/0000-0002-9845-6668>

Additional files

Supplementary files

- Supplementary file 1. Peptide matches to novel CDS detected using mass spectrometry. This table lists the 207 novel CDS (161 mCDS and 46 uaCDS) that have at least one mass-spectrum uniquely matching a peptide in the inferred protein sequence, at protein-level 10% FDR.

DOI: [10.7554/eLife.13328.024](https://doi.org/10.7554/eLife.13328.024)

Major datasets

The following dataset was generated:

| Author(s) | Year | Dataset title | Dataset URL | Database, license, and accessibility information |
|---|------|--|---|--|
| Wang SH, Raj A, Shim H, Gilad Y, Pritchard JK, Stephens M, Engelmann B, Li YI, Harpak A | 2015 | Thousands of novel translated open reading frames in humans inferred by ribosome footprint profiling | http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE75290 | Publicly available at the NCBI Gene Expression Omnibus (accession no: GSE75290). |

The following previously published datasets were used:

| Author(s) | Year | Dataset title | Dataset URL | Database, license, and accessibility information |
|--|------|--|---|---|
| Khan Z | 2015 | Mass-spectrometry measurements in 60 human LCLs | http://proteomecentral.proteomexchange.org/cgi/GetDataset?ID=PXD001406 | Publicly available at ProteomeXchange (accession no: PXD001406) |
| Battle A, Khan Z, Wang SH, Mitrano A, Ford MJ, Pritchard JK, Gilad Y | 2015 | Impact of regulatory variation from RNA to protein | http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE61742 | Publicly available at the NCBI Gene Expression Omnibus (accession no: GSE61742) |
| Lappalainen T, Dermitzakis E | 2013 | RNA-seq measurements in 86 human LCLs | http://www.ebi.ac.uk/arrayexpress/experiments/E-GEUV-1/ | Publicly available at EMBL European Bioinformatics Institute (accession no: E-GEUV-1) |

References

- Auton A**, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR, Consortium GP, 1000 Genomes Project Consortium. 2015. A global reference for human genetic variation. *Nature* **526**:68–74. doi: [10.1038/nature15393](https://doi.org/10.1038/nature15393)
- Bairoch A**, Apweiler R. 2000. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Research* **28**:45–48. doi: [10.1093/nar/28.1.45](https://doi.org/10.1093/nar/28.1.45)
- Barbosa C**, Peixeiro I, Romão L. 2013. Gene expression regulation by upstream open reading frames and human disease. *PLoS Genetics* **9**:e1003529. doi: [10.1371/journal.pgen.1003529](https://doi.org/10.1371/journal.pgen.1003529)
- Battle A**, Khan Z, Wang SH, Mitrano A, Ford MJ, Pritchard JK, Gilad Y. 2015. Genomic variation. Impact of regulatory variation from RNA to protein. *Science* **347**:664–667. doi: [10.1126/science.1260793](https://doi.org/10.1126/science.1260793)
- Bazzini AA**, Johnstone TG, Christiano R, Mackowiak SD, Obermayer B, Fleming ES, Vejnar CE, Lee MT, Rajewsky N, Walther TC, Giraldez AJ. 2014. Identification of small ORFs in vertebrates using ribosome footprinting and evolutionary conservation. *The EMBO Journal* **33**:981–1074. doi: [10.1002/embj.201488411](https://doi.org/10.1002/embj.201488411)
- Birney E**, Stamatoyannopoulos JA, Dutta A, Guigó R, Gingeras TR, Margulies EH, Weng Z, Snyder M, Dermitzakis ET, Thurman RE, Kuehn MS, Taylor CM, Neph S, Koch CM, Asthana S, Malhotra A, Adzhubei I, Greenbaum JA, Andrews RM, Flicek P, et al. 2007. Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**:799–816. doi: [10.1038/nature05874](https://doi.org/10.1038/nature05874)
- Calvo SE**, Pagliarini DJ, Mootha VK. 2009. Upstream open reading frames cause widespread reduction of protein expression and are polymorphic among humans. *Proceedings of the National Academy of Sciences of the United States of America* **106**:7507–7512. doi: [10.1073/pnas.0810916106](https://doi.org/10.1073/pnas.0810916106)
- Camby I**, Le Mercier M, Lefranc F, Kiss R. 2006. Galectin-1: a small protein with major functions. *Glycobiology* **16**:137R–157. doi: [10.1093/glycob/cwl025](https://doi.org/10.1093/glycob/cwl025)
- Clark MB**, Amaral PP, Schlesinger FJ, Dinger ME, Taft RJ, Rinn JL, Ponting CP, Stadler PF, Morris KV, Morillon A, Rozowsky JS, Gerstein MB, Wahlestedt C, Hayashizaki Y, Carninci P, Gingeras TR, Mattick JS. 2011. The reality of pervasive transcription. *PLoS Biology* **9**:e1001102. doi: [10.1371/journal.pbio.1000625](https://doi.org/10.1371/journal.pbio.1000625)
- Cox J**, Mann M. 2008. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nature Biotechnology* **26**:1367–1372. doi: [10.1038/nbt.1511](https://doi.org/10.1038/nbt.1511)
- Cox J**, Neuhauser N, Michalski A, Scheltema RA, Olsen JV, Mann M. 2011. Andromeda: a peptide search engine integrated into the MaxQuant environment. *Journal of Proteome Research* **10**:1794–1805. doi: [10.1021/pr101065j](https://doi.org/10.1021/pr101065j)
- Djebali S**, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W, Schlesinger F, Xue C, Marinov GK, Khatun J, Williams BA, Zaleski C, Rozowsky J, Röder M, Kokocinski F, Abdelhamid RF, Alioto T, et al. 2012. Landscape of transcription in human cells. *Nature* **489**:101–108. doi: [10.1038/nature11233](https://doi.org/10.1038/nature11233)
- 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. doi: [10.1093/bioinformatics/bts635](https://doi.org/10.1093/bioinformatics/bts635)
- Evans SN**, Shvets Y, Slatkin M. 2007. Non-equilibrium theory of the allele frequency spectrum. *Theoretical Population Biology* **71**:109–119. doi: [10.1016/j.tpb.2006.06.005](https://doi.org/10.1016/j.tpb.2006.06.005)
- Farrell CM**, O’Leary NA, Harte RA, Loveland JE, Wilming LG, Wallin C, Diekhans M, Barrell D, Searle SM, Aken B, Hiatt SM, Frankish A, Suner MM, Rajput B, Steward CA, Brown GR, Bennett R, Murphy M, Wu W, Kay MP, et al. 2014. Current status and new features of the Consensus Coding Sequence database. *Nucleic Acids Research* **42**:D865–D872. doi: [10.1093/nar/gkt1059](https://doi.org/10.1093/nar/gkt1059)
- Galindo MI**, Pueyo JI, Fouix S, Bishop SA, Couso JP. 2007. Peptides encoded by short ORFs control development and define a new eukaryotic gene family. *PLoS Biology* **5**:e106. doi: [10.1371/journal.pbio.0050106](https://doi.org/10.1371/journal.pbio.0050106)

- Guttman M**, Russell P, Ingolia NT, Weissman JS, Lander ES. 2013. Ribosome profiling provides evidence that large noncoding RNAs do not encode proteins. *Cell* **154**:240–251. doi: [10.1016/j.cell.2013.06.009](https://doi.org/10.1016/j.cell.2013.06.009)
- Hanguer MJ**, Vaughn IW, McManus MT. 2013. Pervasive transcription of the human genome produces thousands of previously unidentified long intergenic noncoding RNAs. *PLoS Genetics* **9**:e1003569. doi: [10.1371/journal.pgen.1003569](https://doi.org/10.1371/journal.pgen.1003569)
- Harrow J**, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D, Zadissa A, Searle S, Barnes I, Bignell A, Boychenko V, Hunt T, Kay M, Mukherjee G, Rajan J, Despacio-Reyes G, Saunders G, Steward C, et al. 2012. GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Research* **22**:1760–1834. doi: [10.1101/gr.135350.111](https://doi.org/10.1101/gr.135350.111)
- Hernández-Sánchez C**, Mansilla A, de la Rosa EJ, Pollerberg GE, Martínez-Salas E, de Pablo F. 2003. Upstream AUGs in embryonic proinsulin mRNA control its low translation level. *The EMBO Journal* **22**:5582–5592. doi: [10.1093/emboj/cdg515](https://doi.org/10.1093/emboj/cdg515)
- Ingolia NT**, Brar GA, Stern-Ginossar N, Harris MS, Talhouarne GJ, Jackson SE, Wills MR, Weissman JS. 2014. Ribosome profiling reveals pervasive translation outside of annotated protein-coding genes. *Cell Reports* **8**:1365–1379. doi: [10.1016/j.celrep.2014.07.045](https://doi.org/10.1016/j.celrep.2014.07.045)
- Ingolia NT**, Ghaemmaghami S, Newman JR, Weissman JS. 2009. Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science* **324**:218–223. doi: [10.1126/science.1168978](https://doi.org/10.1126/science.1168978)
- Ingolia NT**, Lareau LF, Weissman JS. 2011. Ribosome profiling of mouse embryonic stem cells reveals the complexity and dynamics of mammalian proteomes. *Cell* **147**:789–802. doi: [10.1016/j.cell.2011.10.002](https://doi.org/10.1016/j.cell.2011.10.002)
- Jung HW**, Tschaplinski TJ, Wang L, Glazebrook J, Greenberg JT. 2009. Priming in systemic plant immunity. *Science* **324**:89–91. doi: [10.1126/science.1170025](https://doi.org/10.1126/science.1170025)
- Kapranov P**, Cheng J, Dike S, Nix DA, Dutttagupta R, Willingham AT, Stadler PF, Hertel J, Hackermüller J, Hofacker IL, Bell I, Cheung E, Drenkow J, Dumais E, Patel S, Helt G, Ganesh M, Ghosh S, Piccolboni A, Sementchenko V, et al. 2007. RNA maps reveal new RNA classes and a possible function for pervasive transcription. *Science* **316**:1484–1488. doi: [10.1126/science.1138341](https://doi.org/10.1126/science.1138341)
- Kawase T**, Akatsuka Y, Torikai H, Morishima S, Oka A, Tsujimura A, Miyazaki M, Tsujimura K, Miyamura K, Ogawa S, Inoko H, Morishima Y, Kodera Y, Kuzushima K, Takahashi T. 2007. Alternative splicing due to an intronic SNP in HMSD generates a novel minor histocompatibility antigen. *Blood* **110**:1055–1063. doi: [10.1182/blood-2007-02-075911](https://doi.org/10.1182/blood-2007-02-075911)
- Khan Z**, Ford MJ, Cusanovich DA, Mitrano A, Pritchard JK, Gilad Y. 2013. Primate transcript and protein expression levels evolve under compensatory selection pressures. *Science* **342**:1100–1104. doi: [10.1126/science.1242379](https://doi.org/10.1126/science.1242379)
- Kochetov AV**. 2008. Alternative translation start sites and hidden coding potential of eukaryotic mRNAs. *BioEssays* **30**:683–691. doi: [10.1002/bies.20771](https://doi.org/10.1002/bies.20771)
- Kondo T**, Hashimoto Y, Kato K, Inagaki S, Hayashi S, Kageyama Y. 2007. Small peptide regulators of actin-based cell morphogenesis encoded by a polycistronic mRNA. *Nature Cell Biology* **9**:660–665. doi: [10.1038/ncb1595](https://doi.org/10.1038/ncb1595)
- Kondo T**, Plaza S, Zanet J, Benrabah E, Valenti P, Hashimoto Y, Kobayashi S, Payre F, Kageyama Y. 2010. Small peptides switch the transcriptional activity of Shavenbaby during Drosophila embryogenesis. *Science* **329**:336–339. doi: [10.1126/science.1188158](https://doi.org/10.1126/science.1188158)
- Kozak M**. 1987. An analysis of 5′-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic Acids Research* **15**:8125–8148. doi: [10.1093/nar/15.20.8125](https://doi.org/10.1093/nar/15.20.8125)
- Lammich S**, Schöbel S, Zimmer AK, Lichtenthaler SF, Haass C. 2004. Expression of the Alzheimer protease BACE1 is suppressed via its 5′-untranslated region. *EMBO Reports* **5**:620–625. doi: [10.1038/sj.embor.7400166](https://doi.org/10.1038/sj.embor.7400166)
- Langmead B**, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**:357–359. doi: [10.1038/nmeth.1923](https://doi.org/10.1038/nmeth.1923)
- Lappalainen T**, Sammeth M, Friedländer MR, 't Hoen PA, Monlong J, Rivas MA, González-Porta M, Kurbatova N, Griebel T, Ferreira PG, Barann M, Wieland T, Greger L, van Iterson M, Almlöf J, Ribeca P, Pulyakhina I, Esser D, Giger T, Tikhonov A, et al. 2013. Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* **501**:506–511. doi: [10.1038/nature12531](https://doi.org/10.1038/nature12531)
- Laressergues D**, Couzigou JM, Clemente HS, Martinez Y, Dunand C, Bécard G, Combier JP. 2015. Primary transcripts of microRNAs encode regulatory peptides. *Nature* **520**:90–93. doi: [10.1038/nature14346](https://doi.org/10.1038/nature14346)
- Lee J**, Park EH, Couture G, Harvey I, Garneau P, Pelletier J. 2002. An upstream open reading frame impedes translation of the huntingtin gene. *Nucleic Acids Research* **30**:5110–5119. doi: [10.1093/nar/gkf664](https://doi.org/10.1093/nar/gkf664)
- Lee S**, Liu B, Lee S, Huang SX, Shen B, Qian SB. 2012. Global mapping of translation initiation sites in mammalian cells at single-nucleotide resolution. *Proceedings of the National Academy of Sciences of the United States of America* **109**:E2424–E2432. doi: [10.1073/pnas.1207846109](https://doi.org/10.1073/pnas.1207846109)
- Lin MF**, Jungreis I, Kellis M. 2011. PhyloCSF: a comparative genomics method to distinguish protein coding and non-coding regions. *Bioinformatics* **27**:i275–282. doi: [10.1093/bioinformatics/btr209](https://doi.org/10.1093/bioinformatics/btr209)
- Ma B**. 2015. Novor: Real-time peptide de novo sequencing software. *Journal of the American Society for Mass Spectrometry* **26**:1885–1894. doi: [10.1007/s13361-015-1204-0](https://doi.org/10.1007/s13361-015-1204-0)
- Michel AM**, Choudhury KR, Firth AE, Ingolia NT, Atkins JF, Baranov PV. 2012. Observation of dually decoded regions of the human genome using ribosome profiling data. *Genome Research* **22**:2219–2229. doi: [10.1101/gr.133249.111](https://doi.org/10.1101/gr.133249.111)
- Morris DR**, Geballe AP. 2000. Upstream open reading frames as regulators of mRNA translation. *Molecular and Cellular Biology* **20**:8635–8642. doi: [10.1128/MCB.20.23.8635-8642.2000](https://doi.org/10.1128/MCB.20.23.8635-8642.2000)
- Nesvizhskii AI**. 2014. Proteogenomics: concepts, applications and computational strategies. *Nature Methods* **11**:1114–1125. doi: [10.1038/nmeth.3144](https://doi.org/10.1038/nmeth.3144)

- Nielsen R. 2005. Molecular signatures of natural selection. *Annual Review of Genetics* **39**:197–218. doi: [10.1146/annurev.genet.39.073003.112420](https://doi.org/10.1146/annurev.genet.39.073003.112420)
- Oelkers K, Goffard N, Weiller GF, Gresshoff PM, Mathesius U, Frickey T. 2008. Bioinformatic analysis of the CLE signaling peptide family. *BMC Plant Biology* **8**:1. doi: [10.1186/1471-2229-8-1](https://doi.org/10.1186/1471-2229-8-1)
- Pei B, Sisu C, Frankish A, Howald C, Habegger L, Mu XJ, Harte R, Balasubramanian S, Tanzer A, Diekhans M, Reymond A, Hubbard TJ, Harrow J, Gerstein MB. 2012. The GENCODE pseudogene resource. *Genome Biology* **13**:R51. doi: [10.1186/gb-2012-13-9-r51](https://doi.org/10.1186/gb-2012-13-9-r51)
- Perteau M, Perteau GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. 2015. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature Biotechnology* **33**:290–295. doi: [10.1038/nbt.3122](https://doi.org/10.1038/nbt.3122)
- Pevtsov S, Fedulova I, Mirzaei H, Buck C, Zhang X. 2006. Performance evaluation of existing de novo sequencing algorithms. *Journal of Proteome Research* **5**:3018–3028. doi: [10.1021/pr060222h](https://doi.org/10.1021/pr060222h)
- Rosenbloom KR, Armstrong J, Barber GP, Casper J, Clawson H, Diekhans M, Dreszer TR, Fujita PA, Guruvadoo L, Haeussler M, Harte RA, Heitner S, Hickey G, Hinrichs AS, Hubley R, Karolchik D, Learned K, Lee BT, Li CH, Miga KH, et al. 2015. The UCSC Genome Browser database: 2015 update. *Nucleic Acids Research* **43**:D670–681. doi: [10.1093/nar/gku1177](https://doi.org/10.1093/nar/gku1177)
- Trapnell C, Pachter L, Salzberg SL. 2009. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* **25**:1105–1111. doi: [10.1093/bioinformatics/btp120](https://doi.org/10.1093/bioinformatics/btp120)
- van Bakel H, Nislow C, Blencowe BJ, Hughes TR. 2010. Most "dark matter" transcripts are associated with known genes. *PLoS Biology* **8**:e1000371. doi: [10.1371/journal.pbio.1000371](https://doi.org/10.1371/journal.pbio.1000371)
- Vanderperre B, Lucier JF, Bissonnette C, Motard J, Tremblay G, Vanderperre S, Wisztorski M, Salzet M, Boisvert FM, Roucou X. 2013. Direct detection of alternative open reading frames translation products in human significantly expands the proteome. *PLoS One* **8**:e70698. doi: [10.1371/journal.pone.0070698](https://doi.org/10.1371/journal.pone.0070698)
- Weinberg DE, Shah P, Eichhorn SW, Hussmann JA, Plotkin JB, Bartel DP. 2016. Improved ribosome-footprint and mRNA measurements provide insights into dynamics and regulation of yeast translation. *Cell Reports* **14**:1787–99. doi: [10.1016/j.celrep.2016.01.043](https://doi.org/10.1016/j.celrep.2016.01.043)
- Wilming LG, Gilbert JG, Howe K, Trevanion S, Hubbard T, Harrow JL. 2008. The vertebrate genome annotation (Vega) database. *Nucleic Acids Research* **36**:D753–D760. doi: [10.1093/nar/gkm987](https://doi.org/10.1093/nar/gkm987)
- Xu H, Wang P, Fu Y, Zheng Y, Tang Q, Si L, You J, Zhang Z, Zhu Y, Zhou L, Wei Z, Lin B, Hu L, Kong X. 2010. Length of the ORF, position of the first AUG and the Kozak motif are important factors in potential dual-coding transcripts. *Cell Research* **20**:445–457. doi: [10.1038/cr.2010.25](https://doi.org/10.1038/cr.2010.25)
- Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Molecular Biology and Evolution* **24**:1586–1591. doi: [10.1093/molbev/msm088](https://doi.org/10.1093/molbev/msm088)