

1 **Thousands of novel translated open reading frames in**
2 **humans inferred by ribosome footprint profiling.**

3

4 Anil Raj^{1*†}, Sidney H. Wang^{2*†}, Heejung Shim^{2,6†}, Arbel Harpak⁴,
5 Yang I. Li¹, Brett Engelmann², Matthew Stephens^{2,3}, Yoav Gilad^{2*},
6 Jonathan K. Pritchard^{1,4,5*}

7

8 ¹Department of Genetics, Stanford University, Stanford, CA

9 ²Department of Human Genetics, University of Chicago, Chicago, IL

10 ³Department of Statistics, University of Chicago, Chicago, IL

11 ⁴Department of Biology, Stanford University, Stanford, CA

12 ⁵Howard Hughes Medical Institute

13 ⁶Current address: Department of Statistics, Purdue University, West Lafayette, IN

14 [†]Equal contribution

15 * To whom correspondence should be addressed: rajanil@stanford.edu,
16 siddisis@uchicago.edu, gilad@uchicago.edu, pritch@stanford.edu

1 **Abstract**

2 Accurate annotation of protein coding regions is essential for understanding how
3 genetic information is translated into biological functions. Here we describe
4 riboHMM, a new method that uses ribosome footprint data along with gene
5 expression and sequence information to accurately infer translated sequences.
6 We applied our method to human lymphoblastoid cell lines and identified 7,273
7 previously unannotated coding sequences, including 2,442 translated upstream
8 open reading frames. We observed an enrichment of harringtonine-treated
9 ribosome footprints at the inferred initiation sites, validating many of the novel
10 coding sequences. The novel sequences exhibit significant signatures of
11 selective constraint in the reading frames of the inferred proteins, suggesting that
12 many of these are functional. Nearly 40% of bicistronic transcripts showed
13 significant negative correlation in the levels of translation of their two coding
14 sequences, suggesting a key regulatory role for these novel translated
15 sequences. Our work significantly expands the set of known coding regions in
16 humans.

1 Introduction

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

13

14 Recent work has made it increasingly clear that much of the human genome is
15 transcribed in at least one tissue during some stage of development (Hangauer et al.
16 2013; Djebali et al. 2012; ENCODE Project Consortium et al. 2007; Clark et al. 2011;
17 Kapranov et al. 2007; van Bakel et al. 2010). However, the functional implication for
18 most of these transcripts remains unclear; in particular, the set of sequences translated
19 from these transcripts are not yet completely characterized. For example, there are
20 several recent studies in which RNA transcripts that were previously annotated as
21 noncoding were shown to encode short functional peptides. One well characterized
22 example is the *polished rice (pri) / tarsal-less (tal)* locus in flies, a polycistronic mRNA
23 encoding four short peptides active during embryogenesis (Kondo et al. 2007, 2010;
24 Galindo et al. 2007). While short peptides are known to play critical roles in multiple

1 biological processes (Lauregger et al. 2015; Oelkers et al. 2008; Camby et al. 2006;
2 Jung et al. 2009), annotating genomic regions that encode them remains challenging.

3
4 Direct proteogenomic mass spectrometry has the potential to fill this gap but suffers from
5 variable accuracy in assignment of peptide sequences to spectra and assignment of
6 identified peptides to proteins (for peptides shared across database entries). Moreover,
7 these approaches suffer from a “needle in a haystack” problem when searching all six
8 translational frames over the transcribed portion of the genome (Nesvizhskii 2014;
9 Pevtsov et al. 2006; Ma 2015). Alternative approaches that utilize empirically-derived
10 phylogenetic codon models to distinguish coding transcripts from non-coding transcripts
11 are promising (Lin et al. 2011). However, the success of such approaches is contingent
12 on the duration, strength and stability of purifying selection and these methods may be
13 underpowered for short coding sequences or for newly evolved coding sequences.

14
15 Ribosome profiling utilizes high throughput sequencing of ribosome-protected RNA
16 fragments (RPFs) to quantify levels of translation (Ingolia et al. 2009). Briefly, the
17 technique consists of isolating monosomes from RNase-digested cell lysates and
18 extracting and sequencing short mRNA fragments protected by ribosomes. Early studies
19 of ribosome profiling have shown that RPFs are substantially more abundant within the
20 CDS of annotated transcripts compared to the 5' or 3' untranslated regions (UTRs)
21 (Ingolia et al. 2009; Weinberg et al. 2015). More importantly, when aggregated across
22 annotated coding transcripts, the RPF abundance within the CDS has a clear three
23 base-pair periodicity while the RPF abundance in the UTRs lacks this periodic pattern.

24
25 Recently, using ribosome profiling data, several studies reported conflicting results on
26 the coding potential of long intergenic noncoding RNA (Ingolia et al. 2011; Guttman et al.

1 2013; Ingolia et al. 2014). These studies assessed coding potential using either i) the
2 enrichment of RPFs within the translated CDS relative to background, or ii) the
3 difference in length of RPFs within the translated CDS compared to background.
4 However, these approaches may lack power for several reasons. First, they make little
5 distinction between ribosomes scanning through the transcript and ribosomes decoding
6 the message. Second, the enrichment signal can be severely diminished if the transcript
7 is significantly longer than the coding region within it. Third, there is often substantial
8 variance in RPF abundance within the CDSs, which can decrease power to detect
9 translated sequences when using a simple RPF enrichment statistic alone. Other studies
10 have used the periodicity structure in RPF counts to identify dual coding sequences and
11 short translated CDSs (Michel et al. 2012, Bazzini et al. 2014), but the methods reported
12 high false positive rates and could only identify a few hundred CDSs.

13

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

18

19 **Probabilistic model to infer translated coding sequences**

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

- 23 • **Higher abundance within the CDS.** RPF counts are highly enriched within the
24 CDS overall. Moreover, base positions within the CDS close to the translation
25 initiation and termination sites have substantially higher RPF counts compared to

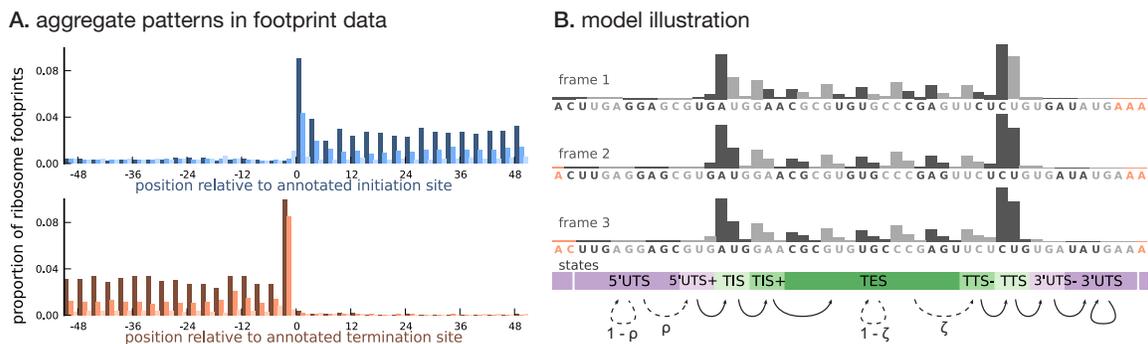
1 base positions in the rest of the CDS. Untranslated regions have very low RPF
2 counts, with the 5'UTR having a slightly higher RPF count compared to the
3 3'UTR. Furthermore, base positions in the 5'UTR immediately flanking the
4 initiation site have a slightly higher RPF count compared to the rest of the 5'UTR;
5 a similar pattern is observed in the 3'UTR.

6 • **Three-base periodicity within the CDS.** RPF counts typically peak at the first
7 position of each codon. The RPF count over the initiation and termination codons
8 tend to have a stronger peak (thus, a slightly different periodic pattern) compared
9 to the rest of the CDS. The RPF counts in the UTRs lack this periodic pattern
10 with similar aggregate counts among base positions in the 5'UTR and 3'UTR.

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

19
20 Conditional on the frame, we modeled the data for each transcript, represented as a
21 sequence of base triplets, using a hidden Markov model (HMM). Each triplet belongs to
22 one of nine latent states — 5'UTS (5' Untranslated State), 5'UTS+ (the last untranslated
23 triplet prior to the initiation site), TIS (Translation Initiation State), TIS+ (the triplet
24 immediately following the initiation site), TES (Translation Elongation State), TTS- (the
25 translated triplet prior to the termination site), TTS (Translation Termination State),

1 3'UTS- (the first untranslated triplet immediately following the termination site), and
 2 3'UTS (3' Untranslated State). The states {TIS, TIS+, TES, TTS-, TTS} denote translated
 3 triplets and {5'UTS, 5'UTS+, 3'UTS-, 3'UTS} denote untranslated triplets. The probability
 4 distribution over the possible sequence of latent states is a function of the underlying
 5 RNA sequence. Figure 1B illustrates these states, and how they relate to each other, in
 6 conjunction with the transcript representation. The groups of states {5'UTS+, TIS, TIS+}
 7 and {TTS-, TTS, 3'UTS-} help model the distinct structure of the RPF counts around the
 8 translation initiation and termination sites, respectively.
 9



10
 11 **Figure 1: Illustrating the model. (A)** Proportion of footprint counts aggregated across 1000 highly
 12 expressed annotated coding transcripts, centered at their translation initiation (blue) and termination
 13 (orange) sites. In aggregate, RPF count data have higher abundance within the CDS than the UTRs and
 14 exhibit a 3-base periodicity within the CDS. **(B)** Each transcript belongs to one of three unobserved reading
 15 frames, and is represented as a sequence of base-triplets (highlighted by differing shades of gray) that
 16 depends on the reading frame. Each triplet belongs to one of nine unobserved states. The state sequence
 17 shown corresponds to frame 3 and varying shades from purple to green highlight the different states. Base
 18 positions marked in orange are modeled independently and always belong to the relevant UTS state.
 19 Transitions with nonzero probabilities are indicated by arrows, with solid arrows denoting a probability of 1
 20 and dotted arrows denoting probabilities that are a function of the underlying sequence.
 21

22 Assuming each transcript has either 0 or 1 CDS, we restricted the possible transitions
 23 between latent states as shown in Figure 1B: transitions from 5'UTS to 5'UTS+ occur
 24 with probability ρ , transitions from TES to TTS- occur with probability ζ , and all other
 25 allowed transitions have probability 1. The transition probabilities ρ and ζ are estimated

1 from the data, and are allowed to depend on the base sequence of the triplet; in addition,
2 the probability ρ also depends on the base sequence context around the triplet (Kozak
3 1987). In this work, we assume that translation termination occurs at the first in-frame
4 stop codon (Equation 9), i.e., we do not consider stop codon readthrough.

5

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

17

18 To quantify the accuracy of our model, we designed a simulation scheme to estimate
19 what fraction of our inferred translated sequences are false discoveries. We first
20 estimated the Type 1 error rate – i.e., the probability of inferring a translated region when
21 no such region exists – using a set of simulated transcripts that had no signal of
22 translation (null transcripts). The simulated transcripts were constructed by permuting
23 the observed footprint counts in annotated coding transcripts. We then used this
24 estimate to quantify the false discovery rate for each class of translated CDSs identified
25 by riboHMM. Independently, using a simulated set of transcripts containing some signal
26 of translation, we quantified the proportion of transcripts where our model incorrectly

1 identified the precise translation initiation site conditional on having identified a
2 translated sequence (see **Materials and Methods** for details on the simulations).

3

4 **Results**

5 **Application to human lymphoblastoid cell lines**

6 We applied riboHMM to infer translated CDSs in human lymphoblastoid cell lines (LCLs)
7 for which gene expression phenotypes were measured genome-wide: mRNA in 86
8 individuals, ribosome occupancy in 72 individuals and protein levels in 60 individuals
9 (Lappalainen et al. 2013; Battle et al. 2015). We first assembled over 2.8 billion RNA
10 sequencing reads into transcripts using StringTie (Pertea et al. 2015). This assembly
11 gives us annotated transcripts that are expressed in LCLs, along with novel transcripts
12 that do not overlap any GENCODE annotated gene. (We do not consider novel isoforms
13 of annotated genes in our analyses.) Restricting to transcripts with at least five footprints
14 mapped to each exon, we used riboHMM to identify high-confidence translated CDS. We
15 learned the maximum likelihood estimates of the model parameters using the top five
16 thousand highly expressed genes. The estimated parameters are robust to the choice of
17 the learning set (Figures S1 and S2). Using these parameters, we then inferred the
18 maximum *a posteriori* (MAP) frame and latent state sequence for each of the assembled
19 transcripts. We retained transcripts whose MAP frame and state sequence
20 corresponded to a pair of translation initiation and termination sites and had a joint
21 posterior probability greater than 0.8. Using a set of simulated null transcripts, we
22 estimated that this posterior cutoff corresponded to a Type 1 error rate of 4.5% per
23 transcript. The MAP frame and state sequence directly give us the nucleotide sequence
24 with the strongest signal of translation; we refer to these as main coding sequences or
25 mCDS.

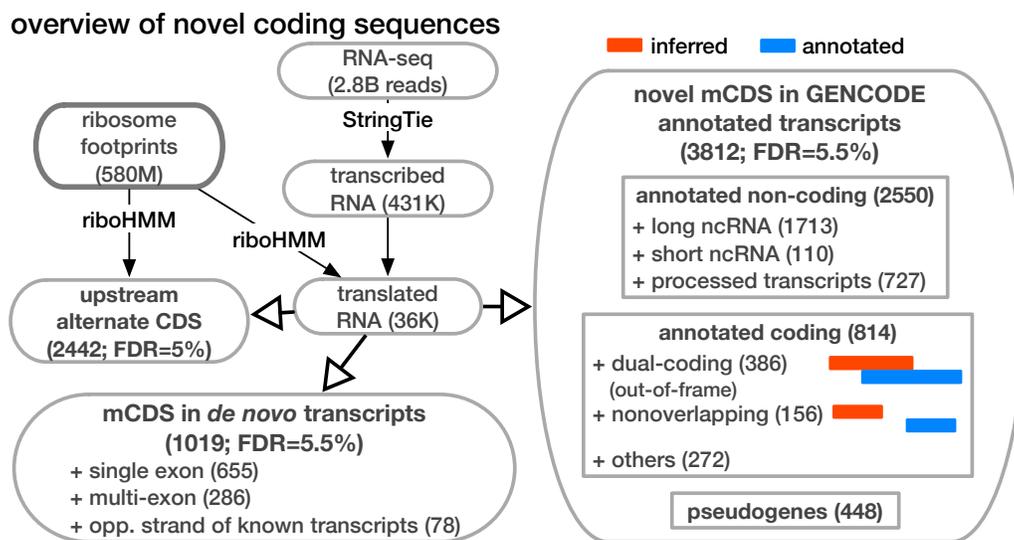
1 **Detection of novel CDSs in LCLs**

2 Among 7,801 GENCODE annotated coding genes for which we could infer a high
3 posterior mCDS, we recovered the annotated reading frame for at least one isoform in
4 7,491 genes (96%); of these, we recovered the exact annotated CDS in 4,500 genes. In
5 310 genes, we inferred mCDS that were entirely distinct from the annotated CDS.
6 (Figure S3 details the rules that decide how our inference agrees with GENCODE.)
7 Thus, for a subset of GENCODE coding genes, our method identified an mCDS distinct
8 from the annotation. Using simulations, we estimated that riboHMM inaccurately detects
9 a completely novel mCDS at a fairly low rate (Type 1 error rate = 4.5%), but has a higher
10 error rate when identifying the precise translation initiation site (false discovery
11 proportion = 38%; see **Materials and Methods** for details). Thus, it is likely that many of
12 the mCDS that have a distinct reading frame compared to annotation are novel alternate
13 translated sequences; however, for those mCDS that only matched the annotated
14 reading frame, the inferred start sites were false discoveries. Our analysis is also robust
15 to sequencing depth; Figure S4 illustrates that nearly 60% of annotated coding
16 sequences identified with the full data set (580 million footprints) could be accurately
17 recovered even when the sequencing depth was reduced by almost two orders of
18 magnitude.

19

20 We identified 4,831 novel mCDS in transcripts expressed in LCLs (FDR = 5.5%). To
21 ensure that these are truly novel, we verified that they do not overlap any known CDS
22 annotated by GENCODE, UCSC (Rosenbloom et al. 2015), or CCDS (Farrell et al.
23 2014) in the same frame. (See Figure 2 for the different classes of LCL transcripts that
24 contain a novel mCDS; Figure S5 illustrates the decision rules used to identify a novel
25 mCDS). Among these, 814 novel mCDS were identified within GENCODE annotated
26 protein-coding transcripts; in these transcripts, the mCDS encodes for a protein distinct

1 from that annotated by GENCODE. Of these, 386 mCDS overlap an annotated CDS but
 2 have a different reading frame (labeled ‘dual-coding’) and 156 do not overlap the
 3 annotated CDS. An example of a novel dual-coding region – an mRNA sequence that
 4 codes for proteins in two different frames – inferred in the POLR2M gene is illustrated in
 5 Figure 3A. Using tandem mass-spectrometry data (Battle et al. 2015), we found four
 6 unique spectra matching peptides in the mCDS and no spectra matching peptides in the
 7 annotated CDS (protein level posterior error probability = 3×10^{-35}).
 8



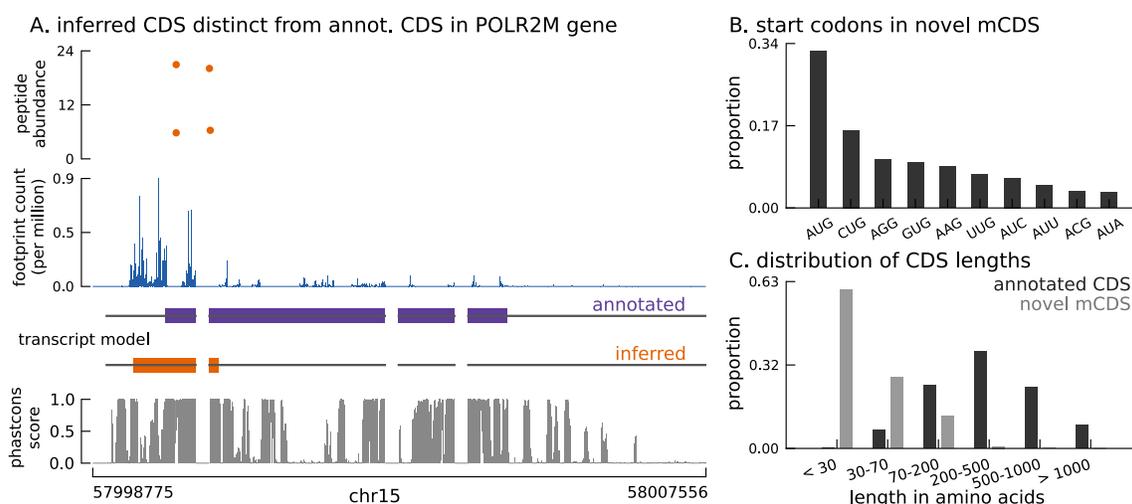
9
 10 **Figure 2: Overview of novel coding sequences.** The analysis workflow indicates the size of the data (in
 11 read/footprint depth, or number of transcripts) at each step and the numbers and classes of transcript within
 12 which novel translated sequences were identified. Transcripts assembled by StringTie that do not overlap
 13 any annotated gene are called “novel transcripts”. Long non-coding RNA includes lincRNAs, antisense
 14 transcripts and transcripts with retained introns, short non-coding RNA includes snRNA, snoRNA and
 15 miRNA, processed transcripts are transcripts without a long, canonical ORF, and pseudogenes include all
 16 subclasses of such genes annotated by GENCODE.

17

18 In addition, we identified 2,550 mCDS in annotated non-coding transcripts and 1,019
 19 mCDS within novel transcripts assembled *de novo* by StringTie. Over 60% of the mCDS
 20 in novel transcripts were identified in single-exon transcripts transcribed from regions

1 containing no annotated genes, while about 8% were identified in novel transcripts
2 transcribed from the strand opposite to an annotated transcript. Finally, we inferred
3 mCDS in 448 expressed pseudogenes, among 14,065 pseudogenes annotated in
4 humans (Pei et al. 2012); nearly 90% of these mCDS were identified in processed
5 pseudogenes. An mCDS in pseudogene GAPDHP72 is shown in Figure S6A, comparing
6 the ribosome abundance and peptide matches to the pseudogene mCDS with those of
7 its parent gene GAPDH.

8



9

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

16

17 Unlike annotated CDS, which almost exclusively start at the methionine codon AUG,
18 these novel mCDS taken together have a substantially higher usage of non-canonical
19 codons, particularly CUG (Figure 3B), consistent with recent observations in mouse
20 embryonic stem cells (Ingolia et al. 2011) and human embryonic kidney cells (Lee et al.
21 2012). Although riboHMM has a high error rate when identifying translation initiation

1 sites, our use of a hierarchical model for the initiation sites suggests that the errors in our
2 inferred start codons are likely to be unbiased. These novel mCDS are also significantly
3 shorter than annotated CDSs (median lengths 23 vs. 339 amino acids, Mann-Whitney
4 test p -value $< 2.2 \times 10^{-16}$; Figure 3C). The overall amino acid content within novel mCDS
5 is comparable to that within annotated CDS, with a slight enrichment for arginine,
6 alanine, cysteine, glycine, proline, and tryptophan residues (binomial test,
7 p -value $< 1.1 \times 10^{-16}$; Figure S7).

8

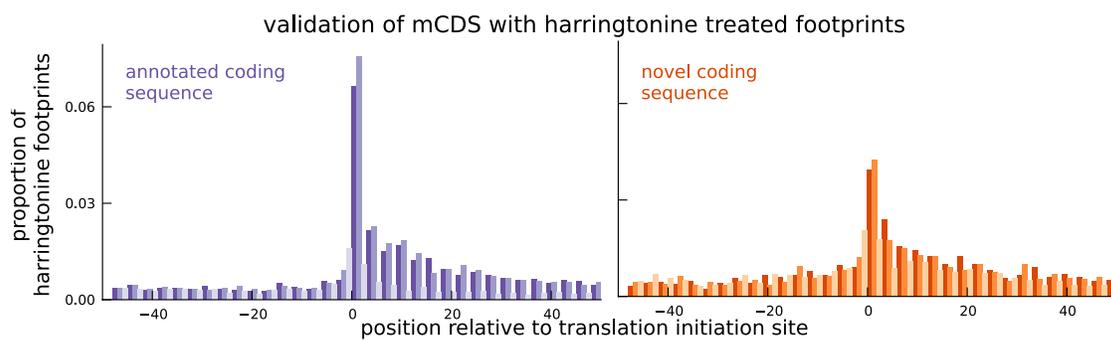
9 Below, using an alternative measure of ribosome occupancy, we first assess
10 independent evidence for translation initiation at many of these novel mCDS. Then, we
11 test whether these mCDS are functional both using human polymorphism data and using
12 substitution patterns across vertebrates. Finally,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

14

15 **Translation at novel mCDS validated using harringtonine-treated ribosome** 16 **footprints**

17 We next sought to provide independent experimental validation for the novel mCDS. A
18 direct approach to validate translation initiation sites is to assay ribosome occupancy in
19 cells treated with harringtonine (Ingolia et al. 2011). Harringtonine interacts with and
20 arrests the initiation complex while leaving the elongation complex to continue
21 translating and run off the transcript. Harringtonine-treated ribosome footprint profiling
22 data therefore show a specific enrichment pattern at the translation initiation site; this
23 pattern has previously been used to identify translation initiation sites in mouse
24 embryonic stem cells (Ingolia et al. 2011). We measured harringtonine-treated ribosome
25 footprints in two LCLs and aggregated the counts of footprints across all novel mCDS.

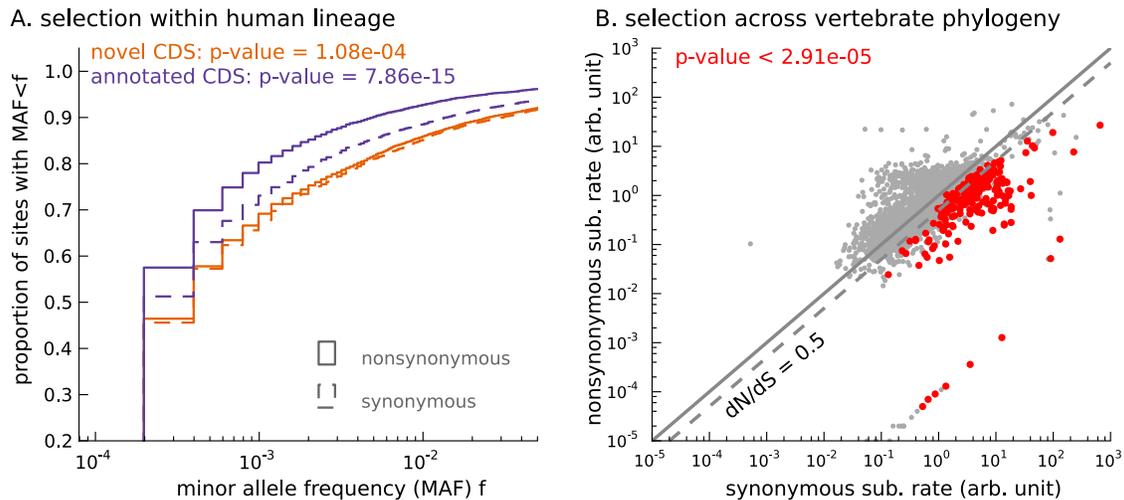
1 We observed an enrichment of footprints at the inferred initiation site of the novel mCDS
2 (binomial test, p -value = 9.5×10^{-79} ; Figure 4), similar to the enrichment of aggregate
3 ribosome occupancy at the initiation sites of a matched number of mCDS that agreed
4 exactly with the annotated CDS (see Figure S6B for mCDS in pseudogenes). We
5 observed a significant enrichment at both AUG (p -value = 5.2×10^{-79}) and non-AUG
6 (p -value = 9.4×10^{-25}) initiation sites. The reduced enrichment for the novel mCDS
7 compared to annotated CDSs is likely due to the lower levels of translation of these
8 novel mCDS and the high error rate in identifying the precise base at which translation is
9 initiated. Accounting for these limitations, our observation of enrichment suggests that
10 ribosomes do initiate the translation of many of the novel mCDS identified by riboHMM.
11



12
13 **Figure 4: Validation of novel mCDS using harringtonine-treated ribosome profiling data.**
14 Harringtonine-treated ribosome footprints show enrichment at the inferred translation initiation sites, when
15 aggregated across novel mCDS (orange), similar to the enrichment at the initiation sites of a matched
16 number of mCDS that agreed exactly with the annotated CDS (purple), suggesting that ribosomes do initiate
17 translation of the novel mCDS.

1 Selective constraint on coding function in novel mCDS

2 We next ascertained the functional importance of these novel mCDS based on the
3 selective constraint imposed on random mutations that occur within them. A bi-allelic
4 single nucleotide polymorphism (SNP) that falls within an mCDS can be inferred as
5 synonymous or nonsynonymous depending on whether switching between the two
6 alleles of the SNP changes the amino acid sequence of the mCDS. If the mCDS do not
7 produce proteins that are functionally important, we expect the two classes of variants to
8 have similar selection pressures on average, and thus to segregate at similar
9 frequencies. Only if the novel mCDS produce functionally important peptides do we
10 expect inferred nonsynonymous SNPs to segregate at lower frequencies than inferred
11 synonymous SNPs.
12



13
14 **Figure 5: Novel translated sequences show significant signatures of coding function. (A)** Genetic
15 variants that are nonsynonymous with respect to the inferred mCDS segregate at significantly lower
16 frequencies in human populations than synonymous variants. The novel regions are under weaker
17 selective constraint compared to known CDS. (The numbers of variants in each class are matched between novel and
18 annotated CDS.) **(B)** Scatter plot comparing the substitution rate at inferred synonymous variants versus
19 inferred nonsynonymous variants for each novel mCDS, computed using multiple sequence alignments
20 across 100 vertebrate species. Highlighted in red are 232 novel mCDS identified to be under significant
21 long-term purifying selection after Bonferroni correction (testing for $dN/dS < 1$; p -value < 2.91×10^{-5}),
22 indicating conserved coding function for these sequences.

1 Starting with biallelic SNPs identified using whole genome sequences of 2,504
2 individuals (1000 Genomes Project Consortium et al. 2015), we examined the set of
3 SNPs falling within all novel mCDS (13,907 variants within 3,096 novel mCDS). We
4 labeled each SNP as synonymous or nonsynonymous with respect to the inferred CDS
5 and show the cumulative distribution of minor allele frequencies (MAF) of each SNP
6 class (Figure 5A). We observed that nonsynonymous SNPs have an excess of rare
7 variants compared with synonymous SNPs (Mann-Whitney test; p -value = 1.08×10^{-4}),
8 implying a difference in the intensity of purifying selection (Nielsen 2005). This observed
9 excess suggests that the novel mCDS are under significant constraint, consistent with
10 functional peptides, albeit weaker than at annotated CDS. The mCDS identified within
11 pseudogenes alone also showed a similar excess of rare variants among
12 nonsynonymous SNPs (Mann-Whitney test; p -value = 5.6×10^{-3}). Such an excess was
13 not observed for pseudogenes that had detectable ribosome occupancy but lacked a
14 high-confidence inferred coding sequence (Figure S6C); for these pseudogenes, the
15 SNPs were labeled based on the reading frame of the parent gene. This highlights that
16 ribosome occupancy alone is insufficient to identify translated sequences, and our
17 method is able to leverage finer scale structure in ribosome footprint data to detect
18 functional coding sequences.

19

20 While the allele frequency spectra provide evidence that some of the novel mCDS are
21 functional in present-day human populations, they are less informative about the long-
22 term selective constraint on these sequences. To identify whether the novel mCDS have
23 been under long-term functional constraint, we compared the substitution rates at
24 synonymous and nonsynonymous sites within the novel mCDS using whole-genome
25 multiple sequence alignments across 100 vertebrates. (We excluded mCDS identified in

1 pseudogenes from this analysis due to difficulties in assigning orthology.) In Figure 5B,
2 232 novel mCDS have a significantly lower nonsynonymous substitution rate (dN)
3 compared to their synonymous substitution rate (dS) after Bonferroni correction
4 (p -value $< 2.91 \times 10^{-5}$), suggesting that these mCDS have been under long-term
5 purifying selection. Since the power to detect significantly low values of dN/dS depends
6 on the length of the CDS and the qualities of the genome assemblies and the multiple
7 sequence alignments across distant species at these sequences, the number of
8 functional novel CDSs identified is a conservative lower bound.

9

10 **Detection of novel proteins by mass spectrometry**

11 We next tested whether we could detect the novel mCDS predictions using mass
12 spectrometry data. We used a large data set of SILAC-labeled tandem mass-spectra
13 generated by trypsin-cleavage of large, stable proteins in many of the same LCLs (Battle
14 et al. 2015). Running MaxQuant (Cox and Mann 2008) against the sequence database
15 of 4,831 novel mCDS, at 10% FDR, we identified 161 novel mCDS sequences that have
16 at least one unique peptide hit -- a tryptic peptide that matches a mass-spectrum (Table
17 S1). More than 70% of novel mCDS with a peptide hit have at least 2 distinct peptides
18 matched to it and, in almost all cases, the unique peptides were independently identified
19 in two or more LCLs (Figure S8).

20

21 To assess how many hits we would expect to the novel mCDS if their properties were
22 like those of annotated CDSs, we developed a model that predicts whether an annotated
23 protein has at least one mass-spectrum match, using features based on expression and
24 sequence composition of the protein (see **Materials and Methods** for more details). The
25 mass-spectrometry data are highly biased towards detection of larger and more highly

1 expressed proteins. Furthermore, the trypsin cleavage step of the experimental protocol
2 imposes strong constraints on the set of unique peptide sequences that can be observed
3 in an experiment. Assuming that the distributions of these predictive features estimated
4 from annotated CDSs can be applied to the novel mCDS, we computed the expected
5 number of novel mCDS with a peptide hit to be 603.

6

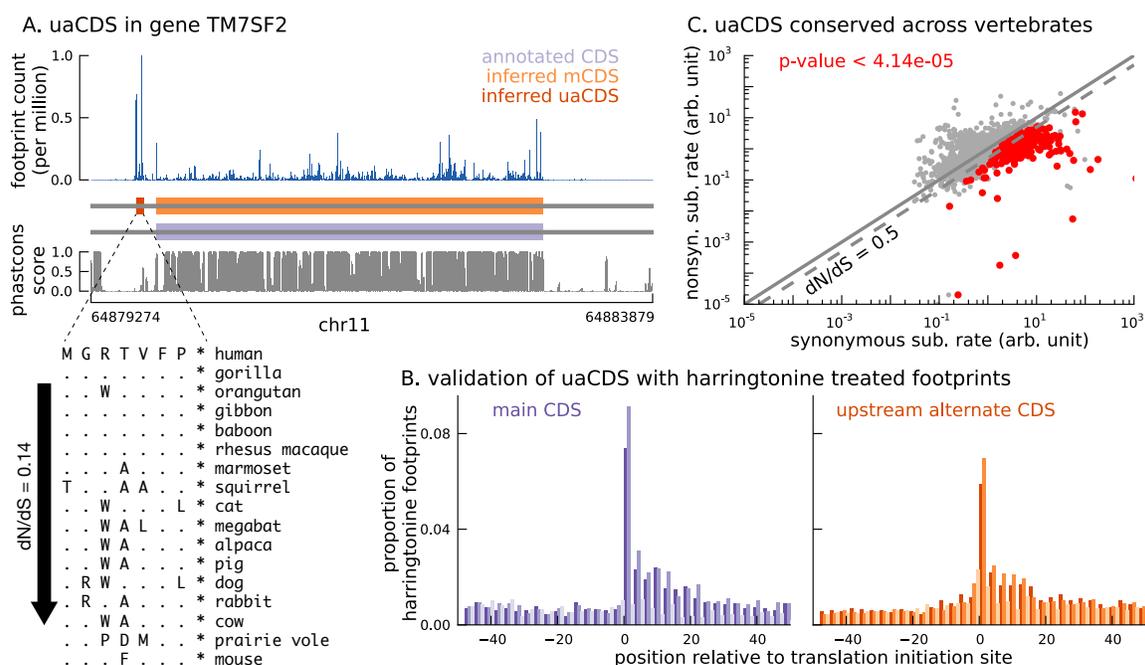
7 We thus find many fewer mass spectrometry hits to the novel mCDS than expected from
8 a model calibrated on previously annotated mCDS (161 vs. 603). The Harringtonine
9 data argue that many of the novel mCDS are correct predictions, thus we suggest that
10 some other property of the mCDS may explain their low detection rate. In particular, it is
11 possible that the novel proteins may have higher turnover rates than annotated proteins.
12 For example it is possible that the proteins translated from novel mCDS may have
13 substantially lower half-life than annotated proteins, or may be secreted, and thus have
14 too low concentrations within the cell to be detectable by mass spectrometry assays.

15

16 **Translation of short alternate coding sequences in addition to the mCDS**

17 Protein-coding transcripts in eukaryotes are typically annotated to have only one CDS
18 (i.e., they are monocistronic). However, a number of studies have demonstrated that
19 ribosomes can initiate translation at alternative start codons (Xu et al. 2010; Ingolia et al.
20 2011; Lee et al. 2012) and many others have identified transcripts with alternative CDSs
21 encoding functional peptides (Vanderperre et al. 2013; Kochetov 2008; Barbosa et al.
22 2013). Furthermore, anecdotal evidence has suggested that translation of the alternate
23 CDS serves as a mechanism to suppress translation of the main CDS (Lee et al. 2002;
24 Hernández-Sánchez et al. 2003; Lammich et al. 2004). However, assessing such a
25 mechanism genome-wide has been challenging, mainly due to a lack of appropriate
26 annotations (Calvo et al. 2009).

1 To this end, we adapted our approach to identify additional coding sequences within
 2 transcripts that are translated in LCLs. Assuming that the sub-codon structure of
 3 footprint abundance is similar between the main and alternate CDS, we identified 2,442
 4 novel CDSs upstream of the mCDS inferred by our method (FDR = 5%); we call them
 5 upstream alternate coding sequences or uaCDS (see **Materials and Methods** for more
 6 details; see also Figure S9). Figure 6A illustrates the ribosome footprint density within
 7 the uaCDS of the transmembrane gene TM7SF2, and its conservation across mammals.
 8 We find strong enrichment of harringtonine-treated ribosome footprints at the initiation
 9 sites of uaCDS similar to the initiation sites of mCDS in the same transcripts (Figure 6B).
 10 Using mass-spectrometry data, we identified 46 uaCDS that have at least one peptide
 11 hit, substantially lower than the expectation of 891 hits predicted by our model. Finally,
 12 comparing the substitution rates at inferred synonymous and nonsynonymous sites, we
 13 identified over 317 uaCDS with highly constrained coding function (Figure 6C).
 14

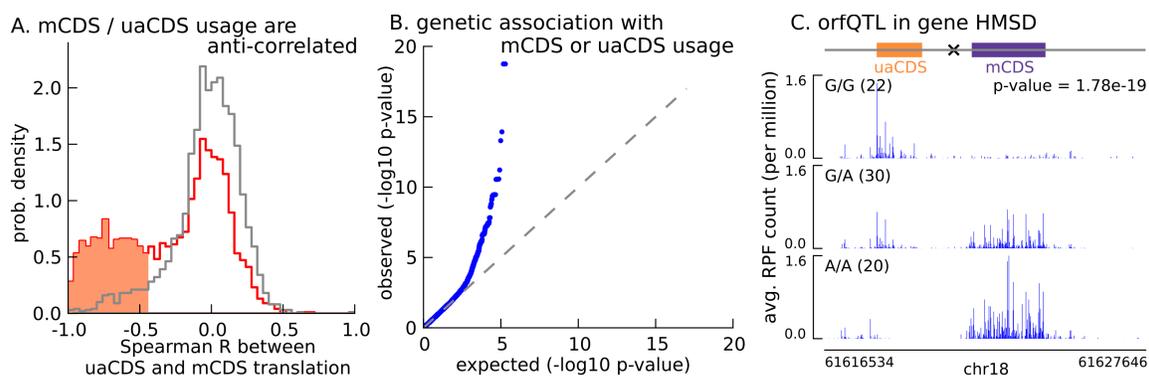


15
 16 **Figure 6: Short translated sequences identified upstream of thousands of translated main coding**
 17 **sequences. (A)** An alternate, novel CDS was identified upstream of the inferred main CDS in gene

1 TM7SF2. As shown in its protein sequence alignment across mammals, the uaCDS (in particular, the start
2 and stop codons) is highly conserved with dN/dS = 0.14. **(B)** Harringtonine-treated ribosome footprints show
3 strong enrichment at the inferred initiation sites of uaCDS, comparable to the enrichment at the initiation
4 sites of the corresponding mCDS, suggesting that ribosomes do initiate translation of these uaCDS. **(C)**
5 Using multiple sequence alignment across 100 vertebrate species, 317 uaCDS were identified to have
6 strong, significant long-term conservation.

8 Translation of uaCDS negatively correlates with translation of mCDS

9 With 2,442 uaCDS identified as translated in LCLs, we next tested the hypothesis that
10 uaCDS expression negatively correlates with mCDS for each pair. We observed that, at
11 10% FDR, 917 pairs of uaCDS and mCDS had significant negative correlations across
12 individuals between the proportion of footprints assigned to them (Figure 7A). Our
13 observation that nearly 40% of pairs of uaCDS and mCDS are significantly anti-
14 correlated, despite incomplete power due to low sample size, suggests that a key role of
15 alternate CDSs in a transcript is to regulate the translation of the main CDS.



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

27

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

13

14 **Discussion**

15 We developed riboHMM, a mixture of hidden Markov models to accurately resolve the
16 precise set of mRNA sequences that are being translated in a given cell type, using
17 sequenced RPFs from a ribosome profiling assay, sequenced reads from an RNA-seq
18 assay and the RNA sequence. When applied to human LCLs, this method was able to
19 accurately identify the translated frame in 96% of annotated coding genes that had a
20 high posterior mCDS. In addition, a key advantage of our framework is the ability to infer
21 novel translated sequences that may be missed by annotation pipelines that focus on
22 long CDSs (>100 amino acids), conservation based approaches that require long-term
23 purifying selection, or direct proteomics measurements that are biased toward highly
24 expressed, stable proteins. We used riboHMM to identify 7,273 novel CDSs, including
25 448 of novel translated sequences in pseudogenes and 2,442 bicistronic transcripts that

1 contain an upstream CDS in addition to a main CDS. We observed enrichment in
2 harringtonine-arrested ribosome occupancy at the inferred translation initiation sites,
3 suggesting that many of the novel mCDS are real. These novel sequences showed
4 significant differences in the amount of purifying selection acting on inferred
5 nonsynonymous versus synonymous sites, suggesting that many of these sequences
6 are conserved as functional peptides, including those mCDS identified in lncRNAs,
7 pseudogenes and novel transcripts.

8

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

22

23 While the precise function of these novel CDSs remains unclear, we found evidence
24 supporting a regulatory role for novel alternate CDSs identified upstream of the mCDS
25 (uaCDS). Although it is unclear whether the down regulation of mCDS by uaCDS is
26 dependent on the peptide sequences of uaCDS, our finding is consistent with previous

1 assertions under which translation of upstream ORFs regulates translation of the main
2 CDS in cap-dependent translation initiation (Morris and Geballe 2000).

3

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

1 **Materials and Methods**

2 **Assembling expressed transcripts in LCLs**

3 We mapped paired-end 75bp RNA-seq reads pooled across 85 Yoruba lymphoblastoid
4 cell lines (Lappalainen et al. 2013) to the Genome Reference Consortium Human
5 Reference 37 (GRCh37) assembly using STAR (Dobin et al. 2013), with the additional
6 flag --outSAMstrandField intronMotif to aid transcript assembly downstream, resulting in
7 2.8 billion uniquely mapped fragments. Using the mapped reads, we assembled models
8 of transcripts expressed in LCLs using StringTie v1.0.4 (Pertea et al. 2015), and used
9 GENCODE v.19 transcript models to guide the assembly. In addition, we required that
10 the lowest expressed isoform of a gene have no less than 1% the expression of the
11 highest expressed isoform (-f 0.01), and that each exon-exon junction be supported by
12 at least 2 spliced reads (-j 2). Since the RNA-seq protocol did not produce strand-
13 specific reads, we treated the forward strand and reverse strand of a transcript model
14 assembled by StringTie as distinct transcripts. Our final set of 430,754 expressed
15 transcripts included 122,168 GENCODE annotated transcript isoforms and 308,586
16 novel isoforms. (We did not consider novel isoforms of annotated genes identified by
17 StringTie.)

18

19 **Ribosome footprint profiling**

20 Ribosome footprint profiling experiments and sequencing data processing were
21 performed as previously described (Battle et al, 2015), with the exception of a
22 harringtonine treatment step to arrest ribosomes at the sites of translation initiation.
23 Briefly, lymphoblastoid cell lines, GM19204 and GM19238, were cultured at 37°C with
24 5% CO₂ in RPMI media with 15% FBS. The media were further supplemented with 2 mM
25 L-glutamate, 100 IU/ml penicillin, and 100 µg/ml streptomycin. Right before cell lysate
26 preparation, each culture was treated with 2 µg/ml harringtonine (final concentration in
27 media) for 2 minutes followed by 100 µg/ml cycloheximide (final concentration in media).
28 For ribosome profiling experiments, ARTseqTM Ribosome Profiling kit for mammalian
29 cells (RPHMR12126) was used following vendor's instructions. Sephacryl S400 spin
30 columns (GE; 27-5140-01) were used for monosome isolation. Libraries were
31 sequenced on an Illumina HiSeq 2500. For sequencing data processing and mapping,
32 adaptor sequences were removed from the 3' end of each read using the Clipper tool
33 from the FASTX-Toolkit. In addition, the 5' most nucleotide (commonly resulted from
34 non-templated additions) was removed using the Trimmer tool from the FASTX-Toolkit.
35 To increase mapping efficiency, we filtered out sequence reads that mapped to rRNA,
36 tRNA or snoRNA (FASTA files downloaded from Ensembl on 05/02/13) using Bowtie 2,
37 version 2.0.2 (Langmead and Salzberg 2012). Processed reads were aligned to genome
38 build hg19 (human) using TopHat v2.0.6 (Trapnell et al. 2009). The mapping step was
39 guided by transcriptome annotations (downloaded from Ensembl on 01/31/13).

1 Mixture of HMMs to model translated coding sequences

2 Consider N transcripts where the n^{th} transcript has length of L_n assumed to be a multiple
 3 of three ($L_n = 3M_n$; see **Transcripts with length not a multiple of three** for details on
 4 how our model handles the remaining one or two base positions when L_n is not a multiple
 5 of three). Our data consist of RPF counts $T = (T^n)_{n=1}^N$, RNA sequence $S = (S^n)_{n=1}^N$, and
 6 transcript expression $E = (E^n)_{n=1}^N$ (in units of RNA-seq reads per base position per million
 7 sequenced reads) on N transcripts, where T^n and S^n are vector quantities and E^n is a
 8 scalar aggregated over the entire length of the transcript. Let $T^n = (T_1^n, \dots, T_{L_n}^n)$ and $S^n =$
 9 $(S_1^n, \dots, S_{L_n}^n)$, where T_b^n and S_b^n denote the RPF counts and the base at the b^{th} position
 10 in the n^{th} transcript, respectively. We model the footprint data T using a mixture of HMMs
 11 that incorporates S and E . Assuming independence across transcripts, the probability of
 12 T given S and E is written as $P(T|\Theta, S, E) = \prod_{n=1}^N P(T^n|\Theta, S^n, E^n)$, where Θ denotes the
 13 set of model parameters.

14 **Mixture of three reading frames for a transcript:** To capture the three-base struc-
 15 ture in RPF data within the CDS, we represent each transcript as a sequence of non-
 16 overlapping base triplets, some of which potentially represent codons. Since the CDS
 17 of the transcript could belong to one of three reading frames (as illustrated in Figure
 18 1B), we introduced a latent frame variable, $F^n \in \{1, 2, 3\}$, that specifies the reading
 19 frame for the n^{th} transcript. Then, given $F^n = f$, T^n can be represented as a se-
 20 quence of $M_n - 1$ triplets and three remaining base positions (see Figure 1B). Specifically,
 21 $T^n | F^n = f := (X_{f,1}^n, \dots, X_{f,(M_n-1)}^n, R_f^n)$, where $X_{f,m}^n = (T_{3m-3+f}^n, T_{3m-2+f}^n, T_{3m-1+f}^n)$
 22 and

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

1 The probability of T^n is then given by

$$P(T^n | \Theta, S^n, E^n) = \sum_{f=1}^3 P(T^n | F^n = f, \Theta, S^n, E^n) P(F^n = f | \Theta, S^n, E^n), \quad (2)$$

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

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

$$P(X_{f,1}^n, \dots, X_{f,(M_n-1)}^n, R_f^n | F^n = f) = P(X_{f,1}^n, \dots, X_{f,(M_n-1)}^n | F^n = f) P(R_f^n | F^n = f). \quad (4)$$

6 (For notation convenience, we have dropped highlighting the dependence of X^n and R^n
 7 on Θ, S^n , and E^n .) We modeled the probability of the data from the sequence of triplets,
 8 $P(X_{f,1}^n, \dots, X_{f,(M_n-1)}^n | F^n = f)$, using an HMM, and the probability of the data from the
 9 remaining positions, $P(R_f^n | F^n = f)$, using a Poisson-gamma model as described below.

10 **HMM for each frame of a transcript:** The pattern of RPF count data in triplets depends
 11 on whether the triplet is being translated or not. To model these patterns, we assumed that
 12 each triplet belongs to one of nine states (see Figure 1B): 5' Untranslated State (5'UTS),
 13 last untranslated triplet 5' to the CDS (5'UTS+), Translation Initiation State (TIS), state after
 14 TIS (TIS+), Translation Elongation State (TES), state before TTS (TTS-), Translation Ter-
 15 mination State (TTS), first untranslated triplet 3' to the CDS (3'UTS-), and 3' Untranslated
 16 State (3'UTS). The five states (TIS, TIS+, TTS, TTS-, TES) denote translated triplets and
 17 the remaining four states (5'UTS, 5'UTS+, 3'UTS, 3'UTS-) denote untranslated triplets. In
 18 particular, the start codon corresponds to the base triplet assigned to the TIS state and the
 19 stop codon corresponds to the base triplet assigned to the 3'UTS- state. The groups of
 20 states (5'UTS+, TIS, TIS+) and (TTS-, TTS, 3'UTS-) help model the distinct features of the
 21 footprint data around the translation initiation and termination sites, respectively. We intro-
 22 duced 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
 23 the n^{th} transcript, where $Z_{f,m}^n$ denotes the state for the m^{th} triplet in the f^{th} frame.

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

1 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

$$X_{f,m}^n | Y_{f,m}^n, Z_{f,m}^n = z \sim \text{multinomial}(Y_{f,m}^n, \pi_z), \quad (5)$$

$$Y_{f,m}^n | Z_{f,m}^n = z \sim \text{Poisson}(\mu_{zfm}^n E^n), \quad (6)$$

$$\mu_{zfm}^n \sim \text{gamma}(\alpha_z, \beta_z), \quad (7)$$

2 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
3 and variance equal to α and $\frac{\alpha}{\beta}$, respectively.

4 The periodicity of RPF counts within the CDS is captured by the multinomial distribu-
5 tion 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 $z \in$
6 $\{5'UTS, 5'UTS+, 3'UTS, 3'UTS-\}$ to capture the lack of periodicity in the RPF data in un-
7 translated regions. Furthermore, we allow the pattern of periodicity to differ across five
8 types of codons (TIS, TIS+, TTS, TTS-, TES).

9 The Poisson distribution for $Y_{f,m}^n$ captures the difference in RPF abundance between trans-
10 lated and untranslated regions (precisely, difference in abundance between triplets in dif-
11 ferent states). We corrected for differences in RPF abundance across transcripts due to
12 differences in transcript expression levels by using E^n as a transcript-specific normaliza-
13 tion factor (see Figure S10). To account for additional variation in the RPF counts across
14 triplets in the same state (e.g., due to varying translation rates across transcripts, and
15 translational pausing), we allowed for triplet-specific parameters μ_{zfm}^n in the Poisson inten-
16 sity and assumed that those parameters follow a gamma distribution. Under this model,
17 $E[Y_{f,m}^n | Z_{f,m}^n = z] = \alpha_z E^n$ and $\text{Var}[Y_{f,m}^n | Z_{f,m}^n = z] = \frac{\alpha_z}{\beta_z} (E^n)^2 + \alpha_z E^n$.

18 We assumed that the sequence of hidden variables Z_f^n follow a Markov chain. The as-
19 sumption of up to one CDS in each transcript leads to a transition probability shown in
20 Figure 1B, where $\rho_{f,m}^n = P(Z_{f,m+1}^n = 5'UTS+ | Z_{f,m}^n = 5'UTS)$ and $\zeta_{f,m}^n = P(Z_{f,m+1}^n =$
21 $TTS- | Z_{f,m}^n = TES)$ depend on the underlying RNA sequence and are given by

$$\rho_{f,m}^n = \begin{cases} \text{logistic}(\psi_\kappa K_{f,m+2}^n + \sum_{c \in \Omega_{\text{start}}} \psi_c \mathbb{1}[M_{f,m+2}^n = c]) & \text{if } M_{f,m+2}^n \in \Omega_{\text{start}}, \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

$$\zeta_{f,m}^n = \begin{cases} 1 & \text{if } M_{f,m+3}^n \in \Omega_{\text{stop}}, \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

1 where $\mathbb{1}[\cdot]$ 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
 2 base sequence of the m^{th} triplet, and $K_{f,m}^n$ denotes the log of ratio of likelihood under
 3 the Kozak model to likelihood under a background model of the base sequence flanking
 4 the m^{th} triplet (see **Kozak model** for details). In our analysis, Ω_{start} contained the canoni-
 5 cal start codon and all near-cognates, $\Omega_{\text{start}} = \{\text{AUG, CUG, GUG, UUG, AAG, ACG, AGG,}$
 6 $\text{AUA, AUC, AUU}\}$, and Ω_{stop} contained the canonical stop codons, $\Omega_{\text{stop}} = \{\text{UAA, UAG,}$
 7 $\text{UGA}\}$. The parameters, ψ_c and ψ_κ , indicate the importance of the triplet base sequence
 8 and the flanking base sequence in determining transition from untranslated triplets to trans-
 9 lated triplets. The current specification of $\zeta_{f,m}^n$ and Ω_{stop} forces the coding sequence to
 10 terminate at the first in-frame occurrence of a stop codon. This model can be extended to
 11 account for stop codon read-through by using a logistic function for $\zeta_{f,m}^n$ for the same set
 12 Ω_{stop} .

13 **Model for R_f^n :** We model R_f^n , the RPF counts at bases before or after the sequence
 14 of triplets (see Equation (1)), using the emission probabilities of the 5'UTS or 3'UTS
 15 states. Assuming that the three elements of R_f^n are independent, we have $P(R_f^n | F^n =$
 16 $f) = \prod_{i=1}^3 P(R_{f,i}^n | F^n = f)$. Each element can be modeled as

$$R_{f,i}^n \sim \text{Poisson}\left(\frac{1}{3} \lambda_{f,i}^n E^n\right), \quad (10)$$

$$\lambda_{f,i}^n \sim \text{gamma}(\alpha_z, \beta_z), \quad (11)$$

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

18 **Parameter estimation and inference:** We used an EM algorithm to compute the maxi-
 19 mum likelihood estimate for the model parameters $\Theta = \{\pi_z, \alpha_z, \beta_z, \psi_\kappa, \psi_c\}$, that is, $\hat{\Theta} :=$
 20 $\text{argmax}_\Theta P(T | \Theta, S, E)$.

21 To infer the translated CDS for the n^{th} transcript, we identified the frame and state se-
 22 quence that maximizes the joint posterior probability

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

23 We first computed the maximum *a posteriori* (MAP) state sequence for each reading frame
 24 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 $f = 1, 2, 3$.
 25 Then, the MAP state sequence and frame is given as

$$(z^{n*}, f^{n*}) = \text{argmax}_f P(Z_f^n = z_f^{n*} | F^n = f, T^n, S^n, E^n, \hat{\Theta}) P(F^n = f | T^n, S^n, E^n, \hat{\Theta}), \quad (13)$$

1 where z_f^{n*} is a function of f , $P(F^n = f|T^n, S^n, E^n, \hat{\Theta}) \propto P(T^n|F^n = f, S^n, E^n, \hat{\Theta})P(F^n =$
 2 $f)$ and $P(T^n|F^n = f, S^n, E^n, \hat{\Theta})$ is the probability of the data marginalized over the latent
 3 states.

4 In our analyses, we estimated the model parameters using the top five thousand highly
 5 expressed genes. Then, we inferred the translated CDS for those transcripts in which
 6 each exon has at least five distinct ribosome footprints mapping to it. We restricted our
 7 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)
 8 the MAP state sequence z^{n*} contains a TIS state and a TTS state (i.e., a pair of initiation
 9 and termination sites), (3) more than 50% of base positions within the inferred CDS are
 10 mappable, and (4) the coding sequence encodes a peptide more than 6 amino acids long
 11 – we call these translated sequences as main coding sequences or mCDS.

12 **Modeling ribosome footprints of different lengths:** We observed that ribosome foot-
 13 prints with different lengths, arising due to incomplete nuclease digestion, show slightly
 14 different patterns of abundance when aggregated across transcripts (see Figure S11). To
 15 model these differences, we partition the footprints into multiple groups based on length,
 16 and model the data in each group with a separate set of parameters in the emission prob-
 17 ability (all groups share the same state sequence along a transcript). Specifically, for G
 18 groups of footprints, the data at the m^{th} triplet in f^{th} reading frame $X_{f,m}^n$ can be partitioned
 19 into G components, $X_{f,m}^n = (X_{g,f,m}^n)_{g=1}^G$, where $X_{g,f,m}^n$ denotes the triplet of RPF counts
 20 from g^{th} group. Assuming that the RPF counts from different groups at a given triplet are
 21 independent, conditional on the state of the triplet, the emission probability can be written
 22 as $P(X_{f,m}^n|Z_{f,m}^n = z) = \prod_{g=1}^G P(X_{g,f,m}^n|Z_{f,m}^n = z)$ and

$$X_{g,f,m}^n|Y_{g,f,m}^n, Z_{f,m}^n = z \sim \text{multinomial}(Y_{g,f,m}^n, \pi_{g,z}), \quad (14)$$

$$Y_{g,f,m}^n|Z_{f,m}^n = z \sim \text{Poisson}(\mu_{g,zf,m}^n E^n), \quad (15)$$

$$\mu_{g,zf,m}^n \sim \text{gamma}(\alpha_{g,z}, \beta_{g,z}), \quad (16)$$

23 where group-specific parameters, $(\pi_{g,z}, \alpha_{g,z}, \beta_{g,z})$, capture the distinct patterns in each
 24 group. The RPF data used in our analyses had four groups of footprints of lengths 28, 29,
 25 30, and 31 bases.

26 **Base positions with missing data:** Approximately 15% of the transcriptome have un-
 27 mappable base positions, in part due to the short lengths of ribosome footprints. Consider

- 1 the m^{th} base triplet in frame f in the n^{th} transcript. If $\mathcal{J}_{g,fm}^n$ is the set of positions in this
- 2 triplet that are unmappable for footprints corresponding to group g , the emission probabili-
- 3 ties become

$$X_{g,fm}^n | Y_{g,fm}^n, Z_{f,m}^n = z \sim \text{multinomial}(Y_{g,fm}^n, \tilde{\pi}_{g,z}), \quad (17)$$

$$Y_{g,fm}^n | Z_{f,m}^n = z \sim \text{Poisson}(\psi_{g,zfm}^n \mu_{g,zfm}^n E^n), \quad (18)$$

$$\mu_{g,zfm}^n \sim \text{gamma}(\alpha_{g,z}, \beta_{g,z}), \quad (19)$$

- 4 where

$$\psi_{g,zfm}^n = \sum_{j \notin \mathcal{J}_{g,fm}^n} \pi_{g,zj}, \quad (20)$$

$$\tilde{\pi}_{g,zj} = \begin{cases} 0 & \text{if } j \in \mathcal{J}_{g,fm}^n, \\ \frac{\pi_{g,zj}}{\psi_{g,zfm}^n} & \text{otherwise.} \end{cases} \quad (21)$$

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

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

16 **Transcripts with length not a multiple of three:** The length of such a transcript can be
 17 written as $L_n = 3M_n + B$, where $B \in \{1, 2\}$. We assumed that the RPF data on the first
 18 $3M_n$ bases ($T_{1:3M_n}^n$) and the data on the remaining B bases ($T_{3M_n+1:L_n}^n$) are independent.
 19 We modeled $T_{1:3M_n}^n$ using a mixture of HMMs as described above, and modeled $T_{3M_n+1:L_n}^n$

1 using the emission probability of the 3'UTS state as follows.

$$P(T_{3M_n+1:L_n}^n | E^n, \alpha_z, \beta_z) = \prod_{m=3M_n+1}^{L_n} P(T_m^n | E^n, \alpha_z, \beta_z), \quad (22)$$

$$T_m^n \sim \text{Poisson}\left(\frac{1}{3}\tau_m^n E^n\right), \quad (23)$$

$$\tau_m^n \sim \text{gamma}(\alpha_z, \beta_z), \quad (24)$$

2 where $z = 3'$ UTS.

3 **Quantifying false discoveries of riboHMM**

4 We characterize the performance of riboHMM by addressing two scenarios: (1) How often
5 does riboHMM identify an mCDS in transcripts with no signal of translation? (2) How often
6 does riboHMM identify an incorrect initiation site in transcripts with signal for translation?
7 To address the first question, we started with the transcripts for which riboHMM was able to
8 identify an mCDS and generated a set of “null transcripts” by permuting the footprint counts
9 among base positions within each transcript. Applying a posterior cutoff of 0.8, riboHMM
10 incorrectly identified an mCDS in 4.5% of these null transcripts. We used this estimate of
11 the Type 1 error rate to compute the false discovery rate for novel mCDS and uaCDS iden-
12 tified by riboHMM. To address the second question, we started with the set of annotated
13 coding transcripts for which riboHMM was able to recover the precise CDS (i.e., the mCDS
14 matched the annotated CDS exactly). We generated a set of “simulated transcripts” using
15 the following strategy: (1) randomly select a new TIS downstream and in-frame to the an-
16 notated TIS, ensuring that the codon underlying the new TIS belonged to the set Ω_{start} , (2)
17 permute the footprint counts among bases upstream of the new TIS. Among the simulated
18 transcripts in which riboHMM could identify an mCDS, the inferred TIS matched the new
19 TIS exactly in 62% of transcripts; this corresponds to a false discovery proportion of 38%.

1 **Translated mCDS in pseudogenes**

2 Starting with 14,065 pseudogenes that have been identified and categorized in humans
3 (Pei et al. 2012), 9,375 pseudogenes were identified by StringTie to be expressed in
4 LCLs. Using a very stringent posterior cutoff of 99.99%, we inferred mCDS in 448 of
5 these expressed pseudogenes. Using pairwise alignment of the pseudogene and parent
6 gene transcript, we observed that although the pseudogene mCDS typically code for
7 shorter protein sequences compared with the parent protein, a large fraction of the
8 pseudogene mCDS share coding-frame with their parent gene (see Figure S6A).

9

10 **Validation with Harringtonine-treated data**

11 Harringtonine-treated ribosome footprints were measured in LCLs with a total
12 sequencing depth of 21 million reads. In Figure 4, we illustrate the aggregate proportion
13 of treated ribosome footprints centered at the inferred start codon for all novel mCDS,
14 and compare it with the aggregate proportion of treated footprints around the start codon
15 of an equal number of annotated CDSs that have a posterior probability greater than 0.8
16 under our model. In Figure S6B, we illustrate the aggregate proportion of treated
17 footprints for mCDS inferred in pseudogenes alone, and in Figure 6B, we compare the
18 aggregate treated footprint proportions at the start codons of inferred uaCDS and their
19 corresponding mCDS.

20

21 **Identifying translated alternate ORFs**

22 For each transcript that had a mCDS with posterior greater than 0.8 and more than 50
23 base pairs of RNA sequence in the 5'UTS state, we defined an "upstream-restricted
24 transcript" consisting of the exons within the 5'UTS state. Using a random set of 5000
25 non-overlapping upstream-restricted transcripts in which more than 80% of base
26 positions were mappable, we computed the maximum likelihood estimates of the
27 transition parameters and occupancy parameters to identify additional translated
28 sequences within these upstream-restricted transcripts. Assuming that the fine-scale
29 structure of footprint counts within these translated sequences would be similar to that
30 within the mCDS, we kept the periodicity parameters fixed to their previously estimated
31 values. With these parameter estimates, we inferred the MAP frame and state
32 sequences with posterior greater than 0.8 and filtered out inferences where less than
33 50% of the inferred CDS was mappable. These additional translated sequences within

1 the upstream-restricted transcripts were called upstream alternate coding sequences or
2 uaCDS.

3

4 **Identifying stable peptides with mass spectrometry data**

5 To identify stable proteins translated from the novel CDSs (mCDS and uaCDS), we
6 analyzed quantitative, high-resolution mass spectrometry data derived from 60 LCLs,
7 with MaxQuant v1.5.0.30 (Cox and Mann 2008) and the Andromeda (Cox et al. 2011)
8 search engine. Sample labeling, processing and data collection details can be found
9 elsewhere (Battle et al. 2015; Khan et al. 2013). Peptides were identified using a
10 database that contained 63,904 GENCODE annotated protein sequences and 7,271
11 novel CDSs identified by our method. For all searches, up to two missed tryptic
12 cleavages were allowed, carbamidomethylation of cysteine was entered as a fixed
13 modification, and N-terminal acetylation and oxidation of methionine were included as
14 variable modifications for all searches. A 'first search' tolerance of 40 ppm with a score
15 threshold of 70 was used for time-dependent mass recalibration followed by a main
16 search MS1 tolerance of 6 ppm and an MS2 tolerance of 20 ppm. The 're-quantify'
17 option was used to aid the assignment of isotopic patterns to labeling pairs. The 'match
18 between runs' option was enabled to match identifications across samples using a
19 matching time window of 42 seconds and an alignment time window of 20 min. Peptide
20 and protein false discovery rates were set to 10% using a reverted version of the search
21 database. Protein group quantifications were taken as the median
22 $\log_2(\text{sample}/\text{standard})$ ratio for all groups containing at least two independent unique or
23 'razor' peptide quantifications (including multiple measurements of the same peptide in
24 different fractions) without a modified peptide counterpart.

25

26 **Bias correction to compute expected number of peptide hits**

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

- 29 1. The median footprint density of annotated coding genes with at least one peptide
30 match is about 125 fold higher than that of coding genes with no peptide match (see
31 Figure S12A).
- 32 2. The median length of coding genes with at least one peptide match is 20% higher
33 than that of coding genes without a peptide match (see Figure S12B).

- 1 3. The trypsin cleavage step of the protocol ensures that nearly all observable peptides
2 have a C-terminal lysine or arginine residue, and up to two additional lysine or
3 arginine residues within the peptide sequence (called “tryptic peptides”). This step
4 imposes a strict constraint on the set of unique peptide sequences that can be
5 observed from a protein sequence, and genes with fewer tryptic peptides are less
6 likely to have a mass-spectrum match.
- 7 4. All tryptic peptides in an expressed protein are not equally likely to be observed. The
8 probability of detecting a tryptic peptide depends on its electrostatic properties
9 relative to other tryptic peptides from all expressed proteins, which in turn depends
10 on the amino acid composition of the tryptic peptides (see Figure S12 C-F).

11

12 To account for these biases, we developed a predictive model to estimate the probability
13 that a protein has at least one peptide hit in a mass-spectrometry experiment. The
14 predicted label for a protein is whether the protein has at least one mass-spectrum
15 match ($H_n = 1$) or no mass-spectrum match ($H_n = 0$). The predictive features of a
16 protein used in the model are (1) the ribosome footprint density of the corresponding
17 transcript (D_n), (2) the protein length (S_n), and (3) the counts of amino acids within each
18 of the K tryptic peptides that can be generated from the protein ($L_n = \{L_{n1}, \dots, L_{nK}\}$).
19 Since the relevant feature of an amino acid is its charge, we partitioned the set of amino
20 acids into four groups – positively charged (R, H, K), negatively charged (D, E), polar
21 uncharged (S, T, N Q), and others. The amino acid count vector L_{nk} was then collapsed
22 into a vector of the counts of each of these four groups. Conditional on $H_n = 1$, we
23 introduced a latent variable for each tryptic peptide that indicates whether the peptide
24 was matched to a mass-spectrum or not ($Z_{nk} \in \{1,0\}$); this latent variable accounts for
25 differences between matched and unmatched peptides.

26

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

29

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

30

1 We learn the predictive model using annotated coding genes and partitioning them into
2 those that have at least one peptide hit (“hit genes”) and those that do not have a
3 peptide hit (“no-hit genes”). We computed $p(D_n|H_n)$ using an empirical distribution of
4 footprint density within coding genes, $p(S_n|H_n)$ using an empirical distribution of the
5 lengths of coding genes, $p(L_{nk}|Z_{nk} = 1, H_n = 1)$ using tryptic peptides within hit genes
6 matched to mass-spectra, $p(L_{nk}|Z_{nk} = 0, H_n = 1)$ using unmatched tryptic peptides
7 within hit genes, and $p(L_{nk}|Z_{nk} = 0, H_n = 0)$ using tryptic peptides within no-hit genes.
8 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
9 peptide hits in annotated proteins, we evaluated the accuracy of this model by holding
10 out some annotated proteins as test data, learning the predictive distributions using the
11 remaining training data and computing the expected number of test proteins that had a
12 mass-spectrum match. We estimated the expected number of held-out annotated
13 proteins with at least one mass-spectrum match to be 1,206 (s.d.=34), while the actual
14 number of held-out proteins with a match was 1,387 (s.d.=36).

15

16 **Test for long-term purifying selection**

17 In order to quantify whether the novel mCDS are evolutionary conserved in terms of their
18 amino acid sequence, we first extracted DNA sequences orthologous to the mCDS from
19 a 100-way vertebrate whole-genome alignments (UCSC), restricting to genomes aligned
20 with either Syntenic net or Reciprocal best net. We next performed a 3-frame translation
21 on each orthologous sequence and a multiple alignment to obtain the correct codon
22 alignments. More specifically, for each orthologous sequence, we kept the frame with
23 the highest amino acid identity compared to the human peptide, requiring at least 60%
24 identity for alignable positions and no more than 50% of the alignment as gaps. Finally,
25 we used codeML/PAML (Yang 2007) to estimate dN and dS rates across the trees
26 consisting of all remaining peptides, first using a model allowing variable omega and
27 then a model with omega fixed to one. To determine whether a specific peptide is under
28 purifying selection or not, we compared the two models using a likelihood ratio test and
29 reported peptides that satisfied a Bonferroni-corrected p -value threshold.

30

31 **Correlation between uaCDS and mCDS**

32 We computed the correlation across LCLs between the proportion of footprints mapped
33 to a transcript that fall within its uaCDS and the proportion that fall within its mCDS. We
34 evaluated the statistical significance of these correlations using an empirical null

1 distribution of Spearman correlations computed using random pairs of mCDS and
2 uaCDS. A random pair of mCDS and uaCDS was obtained by randomly shifting the
3 coordinates of an observed pair of mCDS and uaCDS, matching for their respective
4 lengths and the distance between them.

5

6 **Data Release**

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

9

10 **Acknowledgements**

11 We thank Adam Frankish from the GENCODE group for discussions on the sources of
12 information used for annotating coding genes, Zia Khan for discussions on the mass-
13 spectrometry data analysis, and Audrey Fu for discussions on evaluating the correlation
14 between mCDS and uaCDS. We also thank members of the Pritchard, Gilad and
15 Stephens labs for comments and suggestions. This work was funded by grants from the
16 NIH (HG007036 to J.K.P., MH084703 to Y.G. and J.K.P., and HG02585 to M.S.), and by
17 the Howard Hughes Medical Institute. The funders had no role in study design, data
18 collection and analysis, decision to publish, or preparation of the manuscript.

1 **References**

- 2
- 3 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang
4 HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, et al. 2015. A global
5 reference for human genetic variation. *Nature* **526**: 68–74.
- 6 Bairoch A, Apweiler R. 2000. The SWISS-PROT protein sequence database and its
7 supplement TrEMBL in 2000. *Nucleic Acids Res* **28**: 45–48.
- 8 Barbosa C, Peixeiro I, Romão L. 2013. Gene Expression Regulation by Upstream Open
9 Reading Frames and Human Disease. *PLoS Genet* **9**: e1003529.
- 10 Battle A, Khan Z, Wang SH, Mitrano A, Ford MJ, Pritchard JK, Gilad Y. 2015. Genomic
11 variation. Impact of regulatory variation from RNA to protein. *Science* **347**: 664–
12 667.
- 13 Calvo SE, Pagliarini DJ, Mootha VK. 2009. Upstream open reading frames cause
14 widespread reduction of protein expression and are polymorphic among humans.
15 *Proc Natl Acad Sci USA* **106**: 7507–7512.
- 16 Camby I, Le Mercier M, Lefranc F, Kiss R. 2006. Galectin-1: a small protein with major
17 functions. *Glycobiology* **16**: 137R–157R.
- 18 Clark MB, Amaral PP, Schlesinger FJ, Dinger ME, Taft RJ, Rinn JL, Ponting CP, Stadler
19 PF, Morris KV, Morillon A, et al. 2011. The reality of pervasive transcription.
20 *PLoS Biol* **9**: e1000625; discussion e1001102.
- 21 Cox J, Mann M. 2008. MaxQuant enables high peptide identification rates, individualized
22 p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat*
23 *Biotech* **26**: 1367–1372.
- 24 Cox J, Neuhauser N, Michalski A, Scheltema RA, Olsen JV, Mann M. 2011. Andromeda:
25 a peptide search engine integrated into the MaxQuant environment. *J Proteome*
26 *Res* **10**: 1794–1805.
- 27 Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J,
28 Lin W, Schlesinger F, et al. 2012. Landscape of transcription in human cells.
29 *Nature* **489**: 101–108.
- 30 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M,
31 Gingeras TR. 2013. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*
32 **29**: 15–21.
- 33 ENCODE Project Consortium, Birney E, Stamatoyannopoulos JA, Dutta A, Guigó R,
34 Gingeras TR, Margulies EH, Weng Z, Snyder M, Dermitzakis ET, et al. 2007.
35 Identification and analysis of functional elements in 1% of the human genome by
36 the ENCODE pilot project. *Nature* **447**: 799–816.
- 37 Evans SN, Shvets Y, Slatkin M. 2007. Non-equilibrium theory of the allele frequency
38 spectrum. *Theoretical Population Biology* **71**: 109–119.

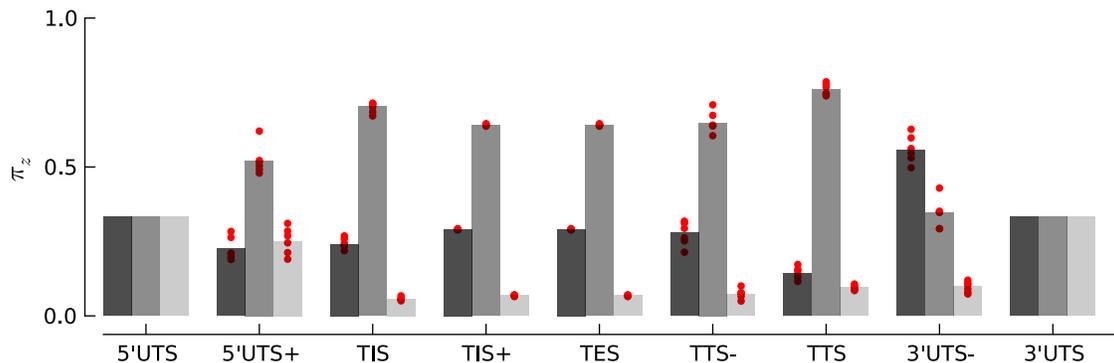
- 1 Furuno M, Kasukawa T, Saito R, Adachi J, Suzuki H, Baldarelli R, Hayashizaki Y,
2 Okazaki Y. 2003. CDS Annotation in Full-Length cDNA Sequence. *Genome Res*
3 **13**: 1478–1487.
- 4 Galindo MI, Pueyo JI, Fouix S, Bishop SA, Couso JP. 2007. Peptides encoded by short
5 ORFs control development and define a new eukaryotic gene family. *PLoS Biol*
6 **5**: e106.
- 7 Guttman M, Russell P, Ingolia NT, Weissman JS, Lander ES. 2013. Ribosome profiling
8 provides evidence that large noncoding RNAs do not encode proteins. *Cell* **154**:
9 240–251.
- 10 Hangauer MJ, Vaughn IW, McManus MT. 2013. Pervasive transcription of the human
11 genome produces thousands of previously unidentified long intergenic noncoding
12 RNAs. *PLoS Genet* **9**: e1003569.
- 13 Hernández-Sánchez C, Mansilla A, de la Rosa EJ, Pollerberg GE, Martínez-Salas E, de
14 Pablo F. 2003. Upstream AUGs in embryonic proinsulin mRNA control its low
15 translation level. *EMBO J* **22**: 5582–5592.
- 16 Ingolia NT, Brar GA, Stern-Ginossar N, Harris MS, Talhouarne GJS, Jackson SE, Wills
17 MR, Weissman JS. 2014. Ribosome profiling reveals pervasive translation
18 outside of annotated protein-coding genes. *Cell Rep* **8**: 1365–1379.
- 19 Ingolia NT, Ghaemmaghami S, Newman JRS, Weissman JS. 2009. Genome-wide
20 analysis in vivo of translation with nucleotide resolution using ribosome profiling.
21 *Science* **324**: 218–223.
- 22 Ingolia NT, Lareau LF, Weissman JS. 2011. Ribosome profiling of mouse embryonic
23 stem cells reveals the complexity and dynamics of mammalian proteomes. *Cell*
24 **147**: 789–802.
- 25 Jung HW, Tschaplinski TJ, Wang L, Glazebrook J, Greenberg JT. 2009. Priming in
26 systemic plant immunity. *Science* **324**: 89–91.
- 27 Kapranov P, Cheng J, Dike S, Nix DA, Dutttagupta R, Willingham AT, Stadler PF, Hertel
28 J, Hackermüller J, Hofacker IL, et al. 2007. RNA maps reveal new RNA classes
29 and a possible function for pervasive transcription. *Science* **316**: 1484–1488.
- 30 Kawase T, Akatsuka Y, Torikai H, Morishima S, Oka A, Tsujimura A, Miyazaki M,
31 Tsujimura K, Miyamura K, Ogawa S, et al. 2007. Alternative splicing due to an
32 intronic SNP in HMSD generates a novel minor histocompatibility antigen. *Blood*
33 **110**: 1055–1063.
- 34 Khan Z, Ford MJ, Cusanovich DA, Mitrano A, Pritchard JK, Gilad Y. 2013. Primate
35 transcript and protein expression levels evolve under compensatory selection
36 pressures. *Science* **342**: 1100–1104.
- 37 Kochetov AV. 2008. Alternative translation start sites and hidden coding potential of
38 eukaryotic mRNAs. *Bioessays* **30**: 683–691.

- 1 Kondo T, Hashimoto Y, Kato K, Inagaki S, Hayashi S, Kageyama Y. 2007. Small peptide
2 regulators of actin-based cell morphogenesis encoded by a polycistronic mRNA.
3 *Nat Cell Biol* **9**: 660–665.
- 4 Kondo T, Plaza S, Zanet J, Benrabah E, Valenti P, Hashimoto Y, Kobayashi S, Payre F,
5 Kageyama Y. 2010. Small peptides switch the transcriptional activity of
6 Shavenbaby during *Drosophila* embryogenesis. *Science* **329**: 336–339.
- 7 Kozak M. 1987. An analysis of 5'-noncoding sequences from 699 vertebrate messenger
8 RNAs. *Nucleic Acids Res* **15**: 8125–8148.
- 9 Lammich S, Schöbel S, Zimmer A-K, Lichtenthaler SF, Haass C. 2004. Expression of
10 the Alzheimer protease BACE1 is suppressed via its 5'-untranslated region.
11 *EMBO Rep* **5**: 620–625.
- 12 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Meth*
13 **9**: 357–359.
- 14 Lappalainen T, Sammeth M, Friedländer MR, Hoen PAC 't, Monlong J, Rivas MA,
15 González-Porta M, Kurbatova N, Griebel T, Ferreira PG, et al. 2013.
16 Transcriptome and genome sequencing uncovers functional variation in humans.
17 *Nature* **501**: 506–511.
- 18 Laressergues D, Couzigou J-M, Clemente HS, Martinez Y, Dunand C, Bécard G,
19 Combiér J-P. 2015. Primary transcripts of microRNAs encode regulatory
20 peptides. *Nature* **520**: 90–93.
- 21 Lee J, Park EH, Couture G, Harvey I, Garneau P, Pelletier J. 2002. An upstream open
22 reading frame impedes translation of the huntingtin gene. *Nucl Acids Res* **30**:
23 5110–5119.
- 24 Lee S, Liu B, Lee S, Huang S-X, Shen B, Qian S-B. 2012. Global mapping of translation
25 initiation sites in mammalian cells at single-nucleotide resolution. *Proc Natl Acad*
26 *Sci USA* **109**: E2424–2432.
- 27 Lin MF, Jungreis I, Kellis M. 2011. PhyloCSF: a comparative genomics method to
28 distinguish protein coding and non-coding regions. *Bioinformatics* **27**: i275–282.
- 29 Ma B. 2015. Novor: Real-Time Peptide de Novo Sequencing Software. *J Am Soc Mass*
30 *Spectrom* 1–10.
- 31 Michel AM, Choudhury KR, Firth AE, Ingolia NT, Atkins JF, Baranov PV. 2012.
32 Observation of dually decoded regions of the human genome using ribosome
33 profiling data. *Genome Res* **22**: 2219–2229.
- 34 Morris DR, Geballe AP. 2000. Upstream open reading frames as regulators of mRNA
35 translation. *Mol Cell Biol* **20**: 8635–8642.
- 36 Nesvizhskii AI. 2014. Proteogenomics: concepts, applications and computational
37 strategies. *Nat Methods* **11**: 1114–1125.

- 1 Nielsen R. 2005. Molecular signatures of natural selection. *Annu Rev Genet* **39**: 197–
2 218.
- 3 Oelkers K, Goffard N, Weiller GF, Gresshoff PM, Mathesius U, Frickey T. 2008.
4 Bioinformatic analysis of the CLE signaling peptide family. *BMC Plant Biol* **8**: 1.
- 5 Pei B, Sisu C, Frankish A, Howald C, Habegger L, Mu XJ, Harte R, Balasubramanian S,
6 Tanzer A, Diekhans M, et al. 2012. The GENCODE pseudogene resource.
7 *Genome Biol* **13**: R51.
- 8 Pertea M, Pertea GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL. 2015.
9 StringTie enables improved reconstruction of a transcriptome from RNA-seq
10 reads. *Nat Biotechnol* **33**: 290–295.
- 11 Pevtsov S, Fedulova I, Mirzaei H, Buck C, Zhang X. 2006. Performance Evaluation of
12 Existing De Novo Sequencing Algorithms. *J Proteome Res* **5**: 3018–3028.
- 13 Trapnell C, Pachter L, Salzberg SL. 2009. TopHat: discovering splice junctions with
14 RNA-Seq. *Bioinformatics* **25**: 1105–1111.
- 15 van Bakel H, Nislow C, Blencowe BJ, Hughes TR. 2010. Most “dark matter” transcripts
16 are associated with known genes. *PLoS Biol* **8**: e1000371.
- 17 Vanderperre B, Lucier J-F, Bissonnette C, Motard J, Tremblay G, Vanderperre S,
18 Wisztorski M, Salzet M, Boisvert F-M, Roucou X. 2013. Direct detection of
19 alternative open reading frames translation products in human significantly
20 expands the proteome. *PLoS ONE* **8**: e70698.
- 21 Weinberg DE, Shah P, Eichhorn SW, Hussmann JA, Plotkin JB, Bartel DP. 2015.
22 Improved ribosome-footprint and mRNA measurements provide insights into
23 dynamics and regulation of yeast translation. *bioRxiv* 021501.
- 24 Xu H, Wang P, Fu Y, Zheng Y, Tang Q, Si L, You J, Zhang Z, Zhu Y, Zhou L, et al.
25 2010. Length of the ORF, position of the first AUG and the Kozak motif are
26 important factors in potential dual-coding transcripts. *Cell Res* **20**: 445–457.
- 27 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**:
28 1586–1591.
- 29

1 **Supplementary Figures**

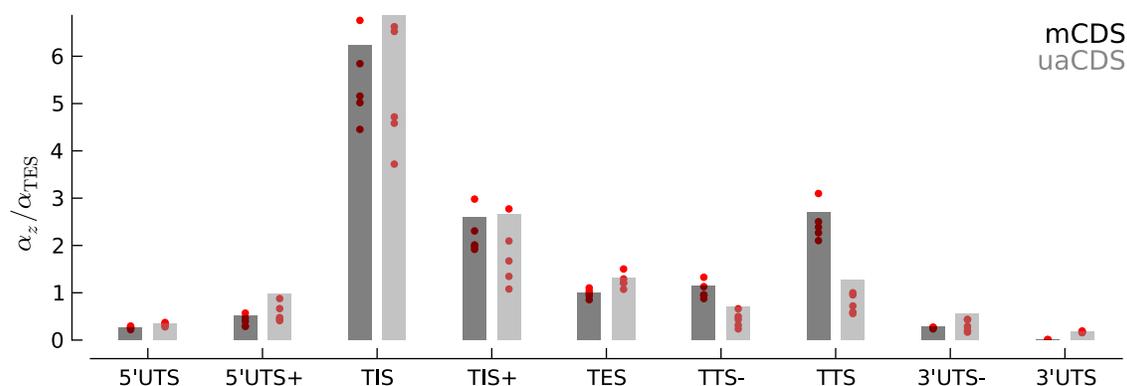
2



3

4 **Figure S1: Robustness of periodicity parameter estimates.** Bar plots of the periodicity parameters (π_z),
 5 for 29 bp long footprints, estimated using the top 5000 expressed genes. Each red circle indicates the
 6 parameter values estimated by using a random set of 1000 genes.

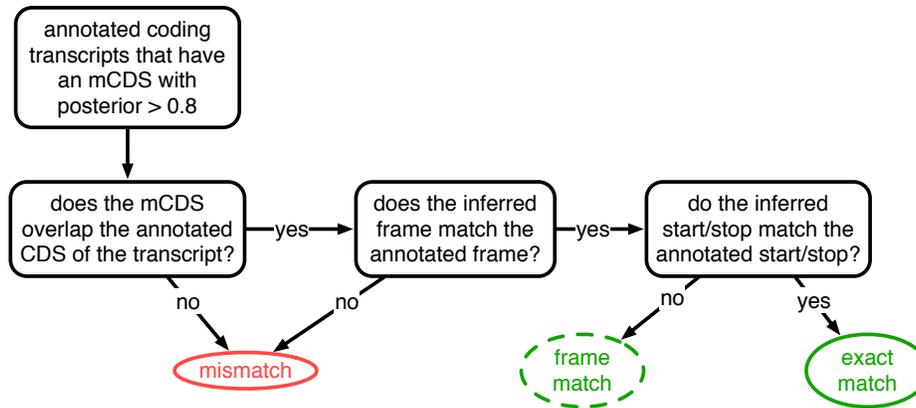
7



8

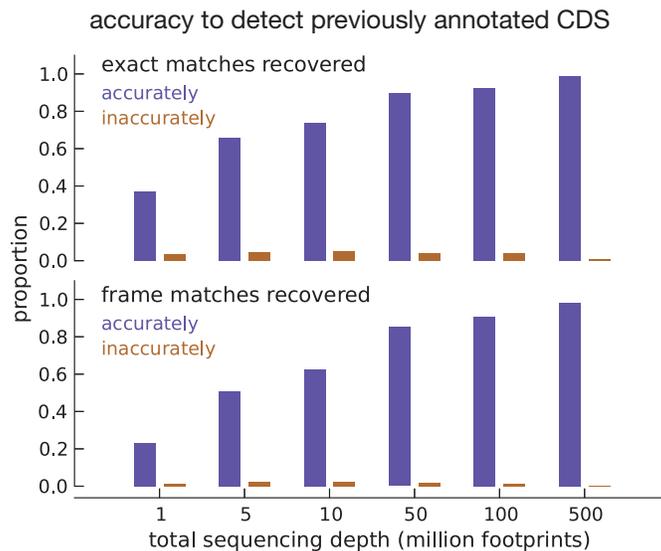
9 **Figure S2: Robustness of occupancy parameter estimates.** Bar plots of the maximum likelihood
 10 estimates of the occupancy parameters in each state, α_z (scaled by α_{TES} for mCDS). For mCDS (dark gray),
 11 the parameters were estimated using the top 5000 expressed genes, and the red circles indicate parameter
 12 values estimated using a random set of 1000 genes. For uaCDS (light gray), the parameters were estimated
 13 using a random set of 5000 upstream-restricted transcripts, and each red circle indicates values estimated
 14 using random sets of 1000 upstream-restricted transcripts.

15



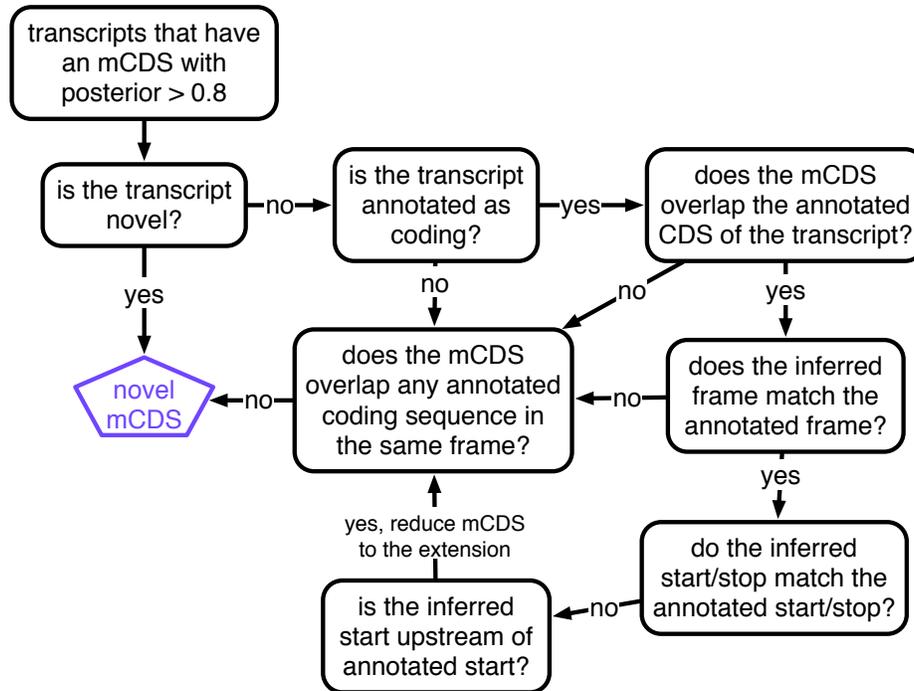
1
2
3
4
5
6
7
8

Figure S3: Decision rules to identify matches and mismatches of mCDS to annotation. Illustrating the decisions by which mCDS inferred on each transcript are identified as an exact match, a frame match) or a mismatch to annotated coding sequences. Matches and mismatches are only defined when the transcript is annotated by GENCODE as coding and the classification depends on agreement with the annotated CDS. A gene is considered to have an exact (or frame) match if at least one of its isoforms is labeled an exact (or frame) match. In all other cases, the inference for the coding gene is considered a mismatch.



9
10
11
12
13
14
15
16
17

Figure S4: Model accuracy. Quantifying model accuracy in terms of the fraction of previously annotated CDSs accurately recovered as a function of total footprint sequencing depth. We performed the entire analysis on the same set of assembled transcripts (parameter estimation and inference) after subsampling the data. Starting with inferences using the complete footprint data set that exactly match annotated CDSs (top subpanel), we show the fraction of these CDSs that were accurately (blue) and inaccurately (brown) recovered with a high posterior for varying sequencing depths. Starting with inferences that only match the frame of annotated CDSs (bottom subpanel), we show the fraction of accurately and inaccurately recovered CDSs for varying sequencing depths.

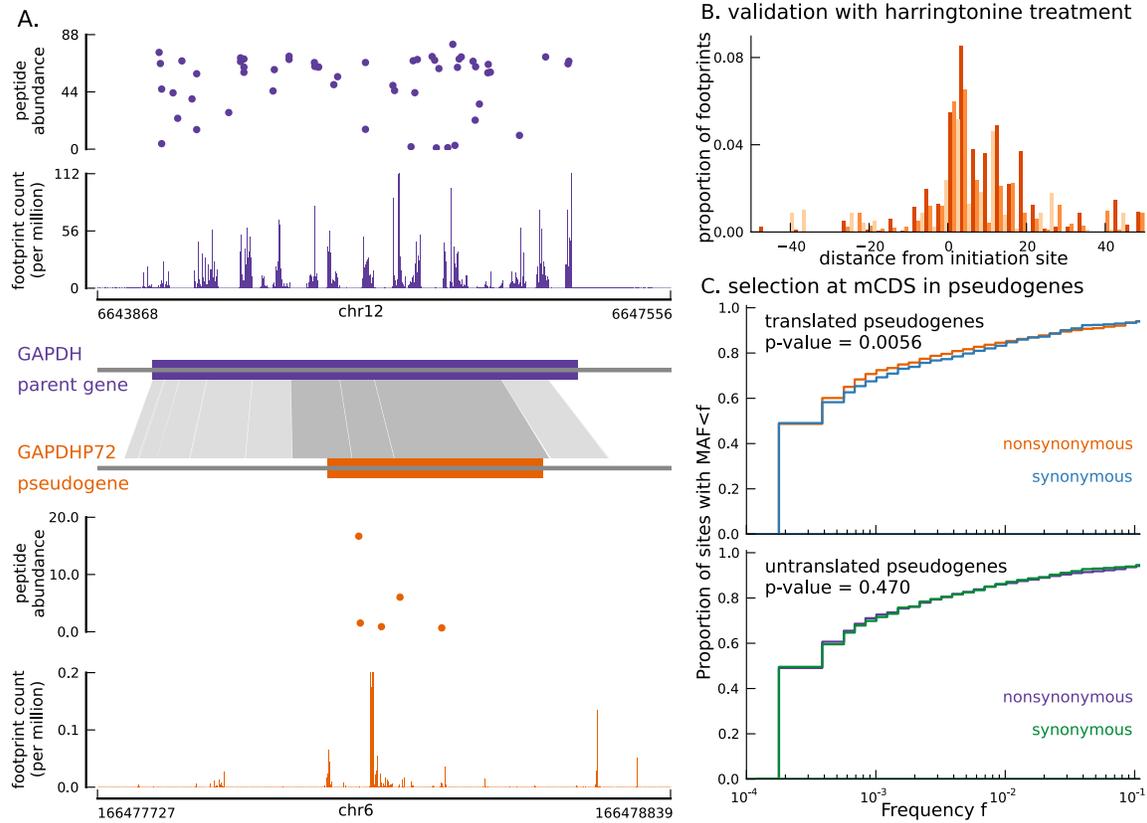


1

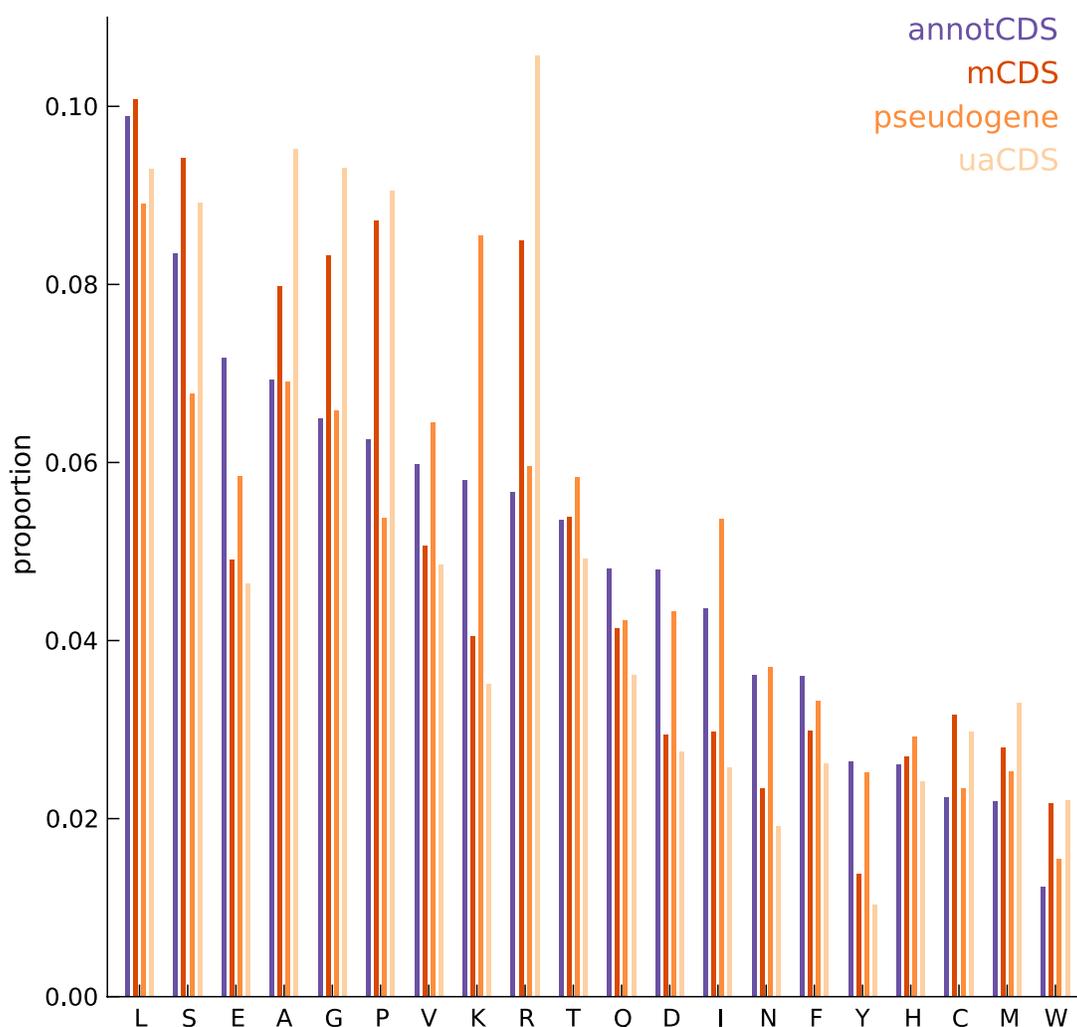
2 **Figure S5: Decision rules to identify novel mCDS.** Illustrating the decisions that identify novel mCDS.

3 The mCDS that do not overlap any known CDS (pooling GENCODE, UCSC and CCDS annotations) are

4 labeled as novel mCDS -- these include mCDS from both novel and annotated transcripts.



1
 2 **Figure S6: Translated coding sequences identified in hundreds of pseudogenes.** (A) Comparing the
 3 inferred CDS in the pseudogene GAPDHP72 (orange) with the annotated CDS (purple) of its parent gene
 4 GAPDH (introns have been removed for better visualization). Gray shaded boxes show the alignment
 5 between the pseudogene and parent gene transcripts, and dark shaded boxes indicate that the inferred
 6 coding frame of the pseudogene matches that of the parent gene. Although the pseudogene and parent
 7 gene share coding frames, the underlying sequences are sufficiently different. (B) Enrichment of
 8 harringtonine-arrested ribosome occupancy at the inferred translation initiation sites validates our inferred
 9 mCDS in pseudogenes. (C) Nonsynonymous variants (orange) segregate at significantly lower frequencies
 10 than synonymous variants (blue) in those pseudogenes predicted to have a translated CDS. Pseudogenes
 11 with ribosome occupancy, but predicted to have no translated CDS, do not show any significant difference
 12 between the site frequency spectra of synonymous and nonsynonymous variants.



1

2

Figure S7: Comparing the amino acid content between annotated and novel CDS. The overall amino

3

acid content is largely similar between annotated CDSs and novel CDSs (mCDS, pseudogenes and

4

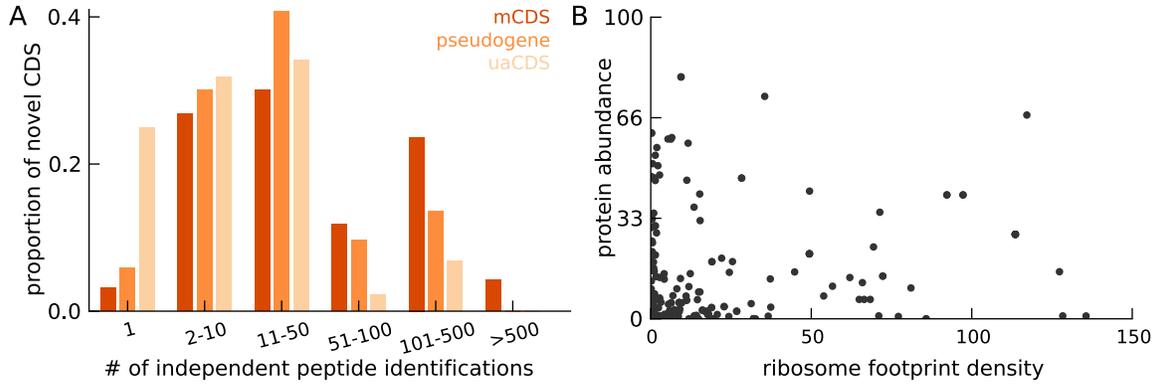
uaCDS), with a substantially higher proportion of lysine and isoleucine residues in CDSs within

5

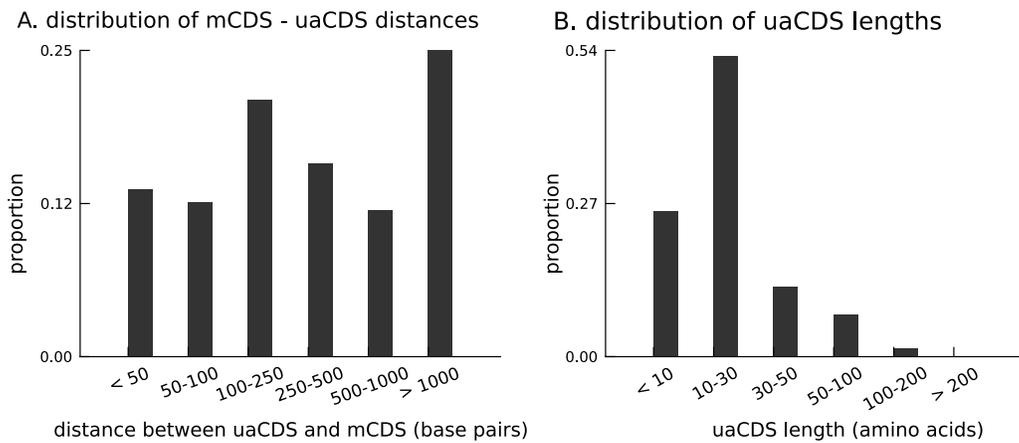
pseudogenes, and a higher proportion of alanine, arginine, glycine, and proline residues within mCDS and

6

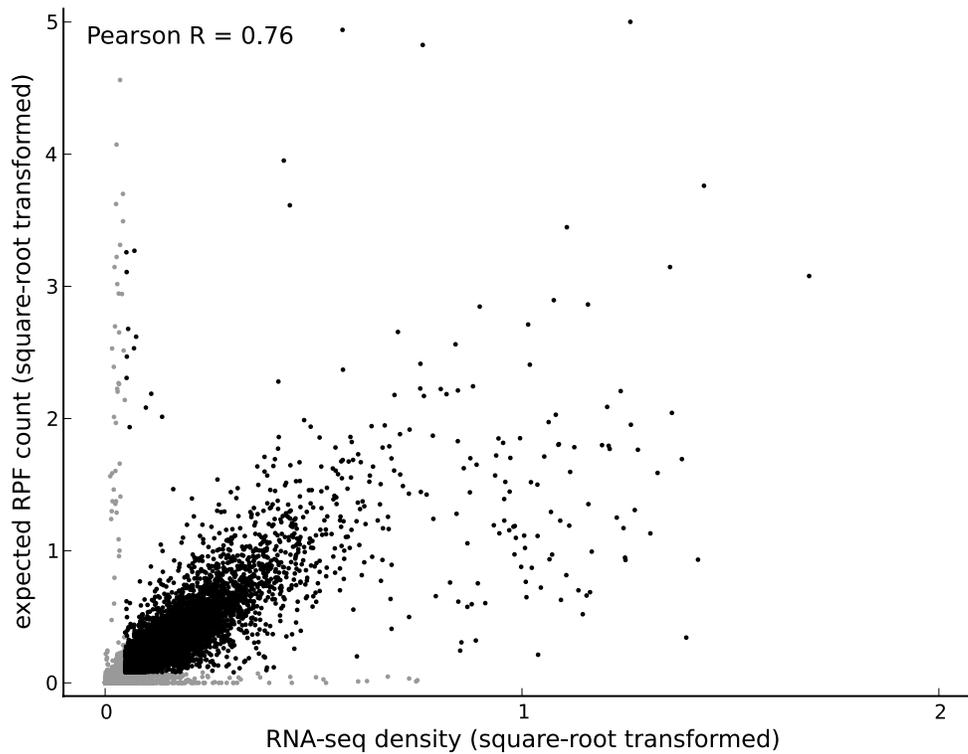
uaCDS.



1
2 **Figure S8: Characteristics of peptides matched to novel CDS.** (A) Histogram of number of independent
3 identification of unique peptides across novel mCDS that have at least one peptide match. (B) Scatter plot of
4 the untreated ribosome footprint density and the protein abundance (median abundance across all unique
5 peptides) across all novel mCDS that have at least one peptide match (Spearman $R = 0.24$; p -value =
6 2×10^{-2}).
7



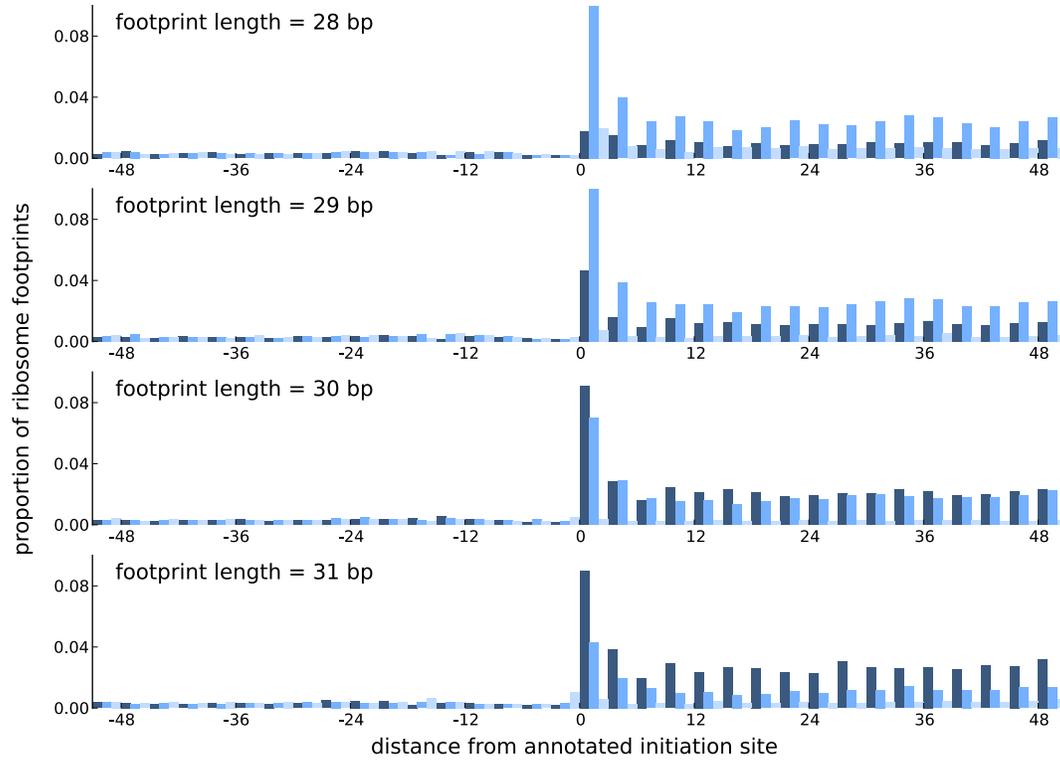
8
9 **Figure S9: Characteristics of novel uaCDS.** (A) Histogram of the distance between uaCDS and their
10 corresponding mCDS. (B) Histogram of the lengths of uaCDS (median length of 16 amino acids).



1

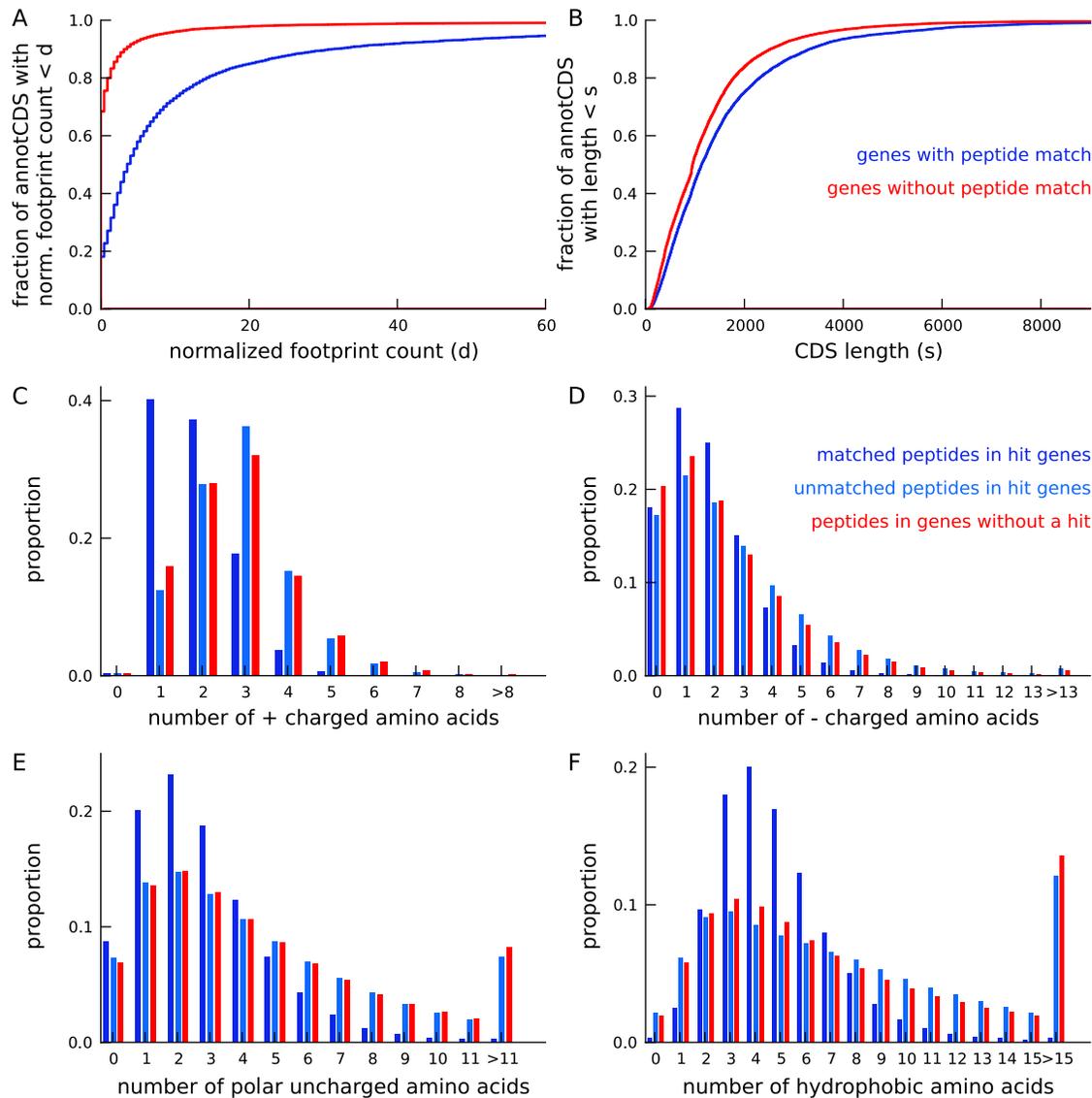
2 **Figure S10: Comparing footprint abundance and gene expression.** Scatter plot of the RNA-seq density
3 (reads per base per million sequenced reads) within annotated coding transcripts vs. the expected ribosome
4 footprint count (normalized by sequencing depth) in triplets within the CDS of the transcripts. The top 50th
5 percentile of expressed genes is shown in black and the bottom 50th percentile is shown in gray. For highly
6 expressed genes, it is reasonable to assume that the expected footprint count in a triplet scales linearly with
7 RNA-seq density. However, for lowly expressed genes and outlier genes (i.e., genes with high expected
8 footprint count and low RNA-seq density, and genes with high RNA-seq density and low expected footprint
9 count), this assumption may not be valid.

10



1
2
3
4
5
6
7

Figure S11: Comparing the periodicity in ribosome footprint counts for footprints of different lengths. Proportion of footprint counts aggregated across 1000 highly expressed annotated coding transcripts, centered at their translation initiation sites. Each subplot illustrates the proportions for footprints of a given length, with the periodicity in the proportions showing distinct features between footprints of different lengths.



1
2 **Figure S12: Annotated genes with peptide hits tend to be longer, have higher expression and a**
3 **distinct amino acid composition.** (A) Cumulative distribution of footprint density for genes with at least one
4 unique peptide hit (blue) and genes with no unique peptide hit (red). The median footprint density of genes
5 with a peptide hit is about 125 fold higher than the median footprint density of genes without a peptide hit.
6 (B) Cumulative distribution of protein length for genes with at least one unique peptide hit (blue) and genes
7 with no unique peptide hit (red). Genes with a peptide hit tend to code for proteins that are 20% longer than
8 proteins encoded by genes without a peptide hit. (C-F) Comparing amino acid composition within tryptic
9 peptides with a mass-spectrum match and tryptic peptides without a mass-spectrum match. Amino acids,
10 grouped by their electrostatic properties, have distinct compositions between matched and unmatched
11 peptides. Matched peptides tend to be significantly shorter than unmatched peptides, and have a distinct
12 composition of charged amino acids.