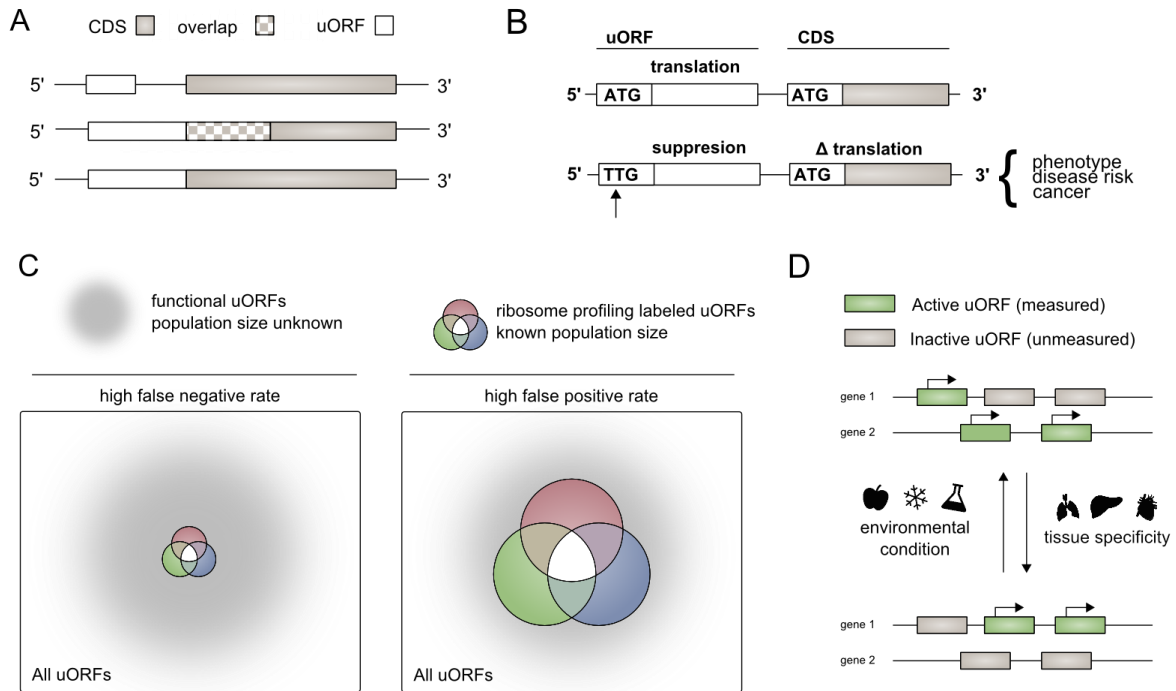


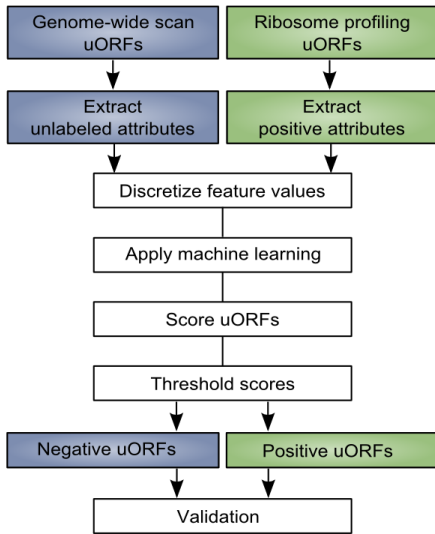
**Figures:**



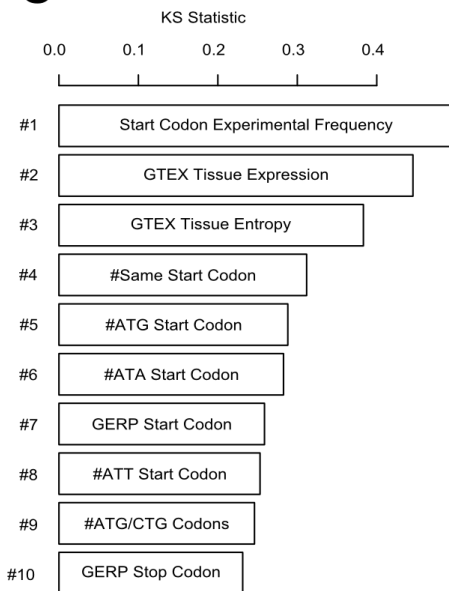
**Figure 1:**

**A. Structure of upstream open reading frames.** The stop codon of an uORF may be located before the CDS start codon [top], or downstream of the CDS start codon, if the uORF is frame-shifted relative to the CDS [middle]. If uORF and CDS share the same stop codon, the uORF acts as a 5' extension of the CDS [bottom]. **B. Effect of mutation or variation on upstream open reading frames.** Creation or destruction of an upstream open reading may have downstream effect on translation of the coding sequence. Change in translation of the coding sequence, may result in change in phenotype and disease risk. **C. Sensitivity and specificity of ribosome profiling for identifying upstream open reading frames.** It is possible that ribosome profiling studies have a high false positive rate (top), or a high false negative rate (bottom). We make the assumption that ribosome profiling studies have a high false negative rate for identifying translated upstream open reading frames (bottom). **D. Activity of uORFs varies according to cell type and environmental stimuli.** uORFs may not be detected in a ribosome profiling experiment, due to variation in uORF activity with cell type and cell environment.

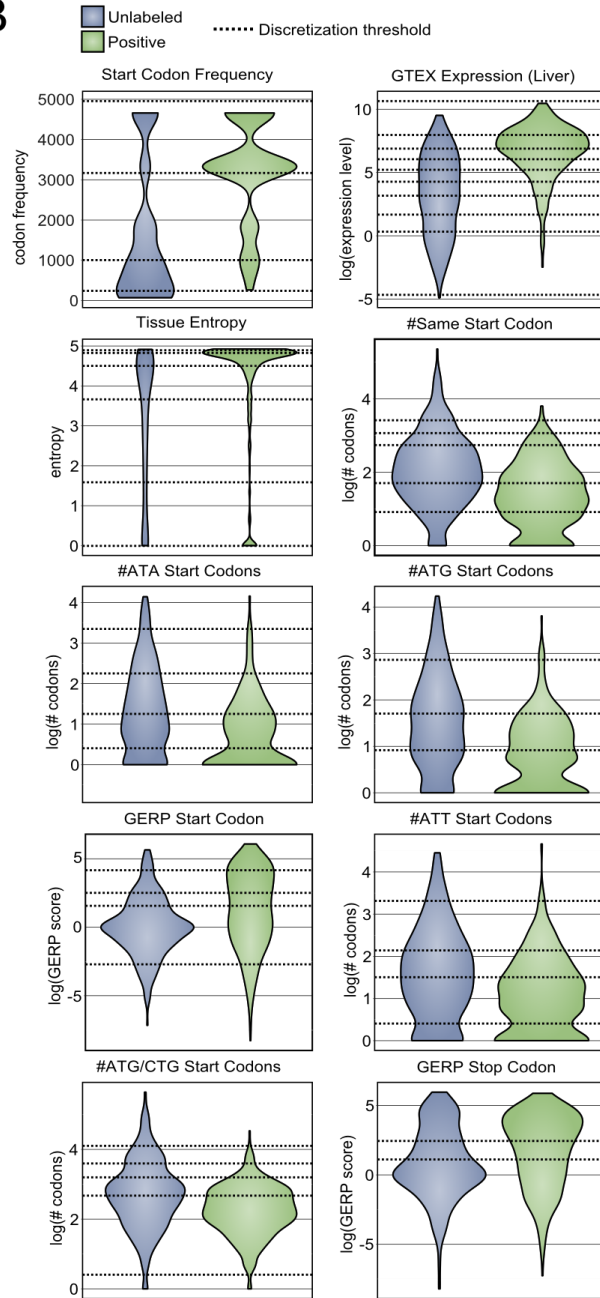
**A**



**C**



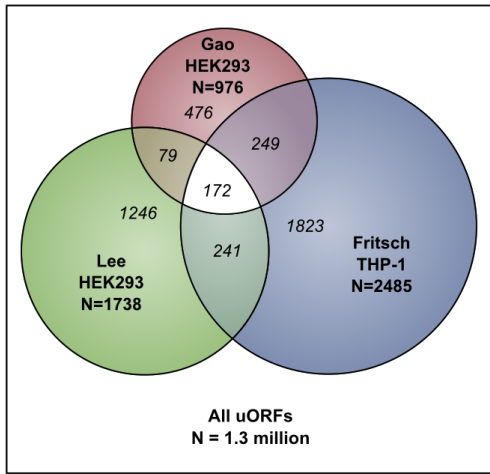
**B**



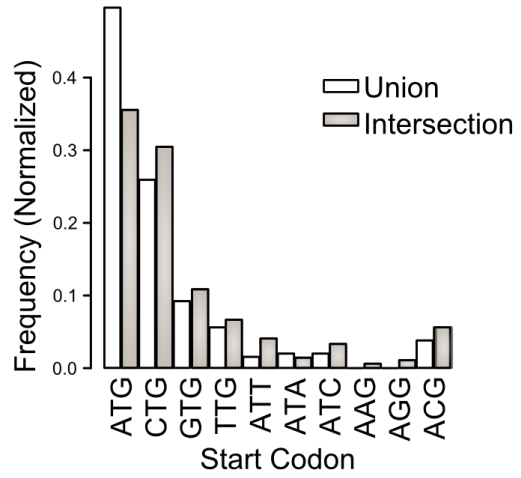
**Figure 2:**

**A. Methodology for distinguishing positive from unlabeled uORFs.** uORFs identified through genome-wide scan, and uORFs labeled in ribosome profiling experiments, were used to train a machine learning algorithm to identify uORFs that are likely active (positive predictions). **B. Distributions of attributes for positive and unlabeled uORFs.** uORF attributes are used to distinguish positive from unlabeled uORFs. Continuous distributions were discretized and optimized for machine learning using the minimum description length principle (MDLP) binning algorithm. Horizontal lines on the plot correspond to these binning intervals. The 10 attributes with the greatest difference in distribution (largest Kolmogorov Smirnov (KS) statistic) between positive and unlabeled uORFs are shown. **C. Upstream open reading frame attributes as classifiers.** Attributes are ranked, according to the difference in distribution between positive and unlabeled uORFs, using the KS statistic.

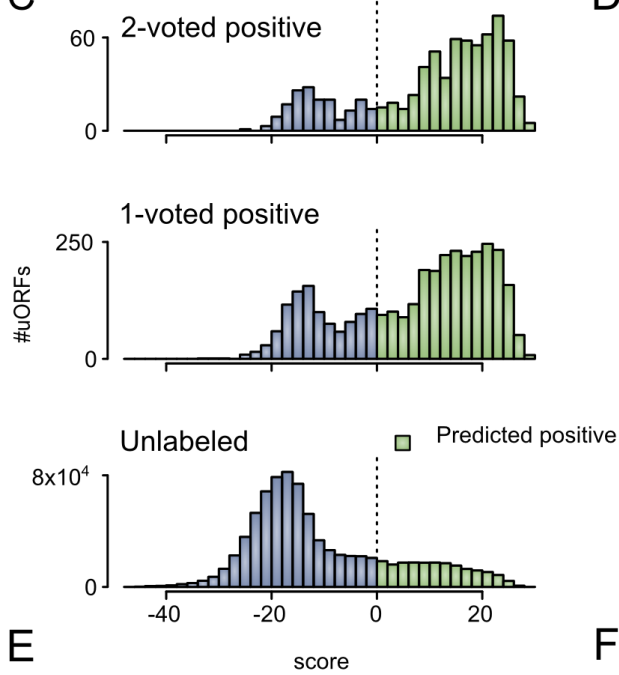
**A**



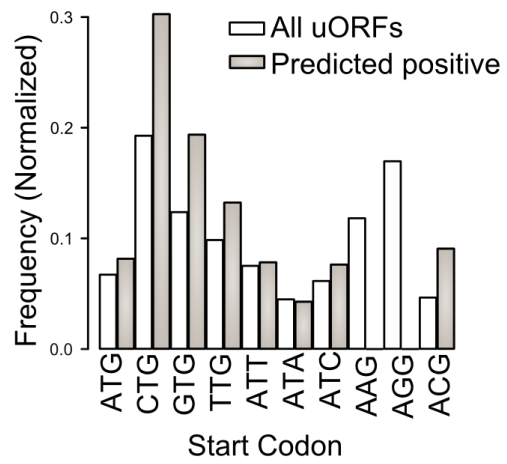
**B**



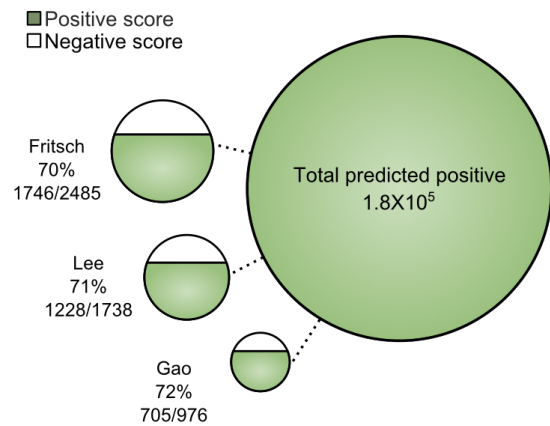
**C**



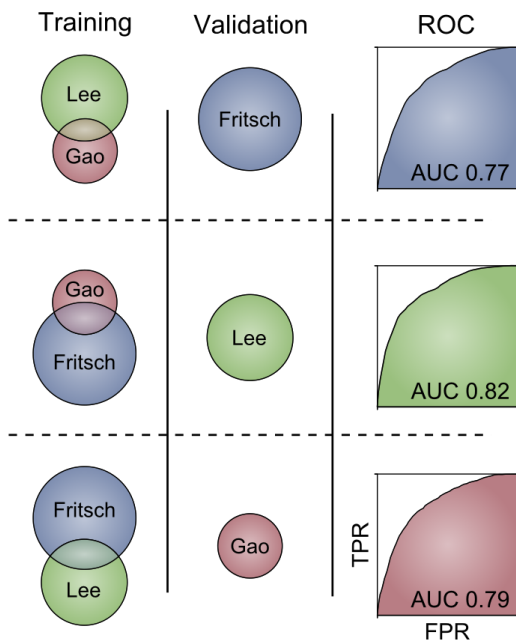
**D**



**E**

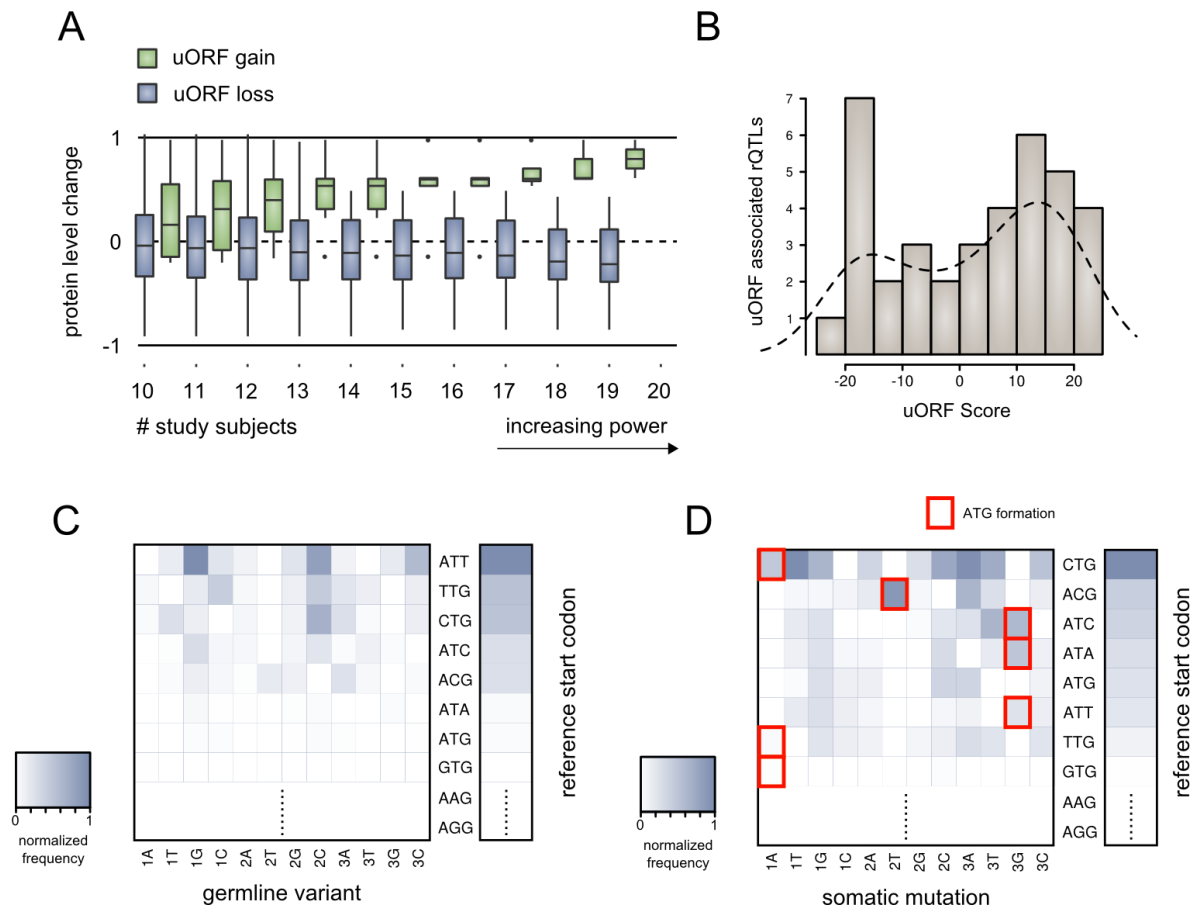


**F**



**Figure 3:**

**A. Frequency of translated uORF ATG start codons, and near-cognate start codons, from ribosome profiling experiments.** Frequency for uORFs translated in any experiment (union), or in more than one experiment (intersection). **B. Ribosome profiling identified uORFs as a subset of all uORFs.** The universe of all uORFs is identified through comprehensive search of the GENCODE human genome annotation [outer border]. Ribosome profiling studies of Fritsch et al., Lee et al., and Gao et al. are shown as overlapping subsets of this universe. Pair-wise and three-way intersections between these experiments are highlighted. **C. Score distributions for upstream open reading frames.** Score distributions for ~~positive-2-voted~~ positive\_uORFs that are translated in two or more ribosome profiling experiments (top), ~~neutral-1-voted~~ positive\_uORFs that are translated in one ribosome profiling experiment (middle), and unlabeled uORFs uncovered through genome-wide search (bottom). **D. The frequency of uORF ATG start codons, and near-cognate start codons, for predicted positive upstream open reading frames.** Frequency is given for all uORFs genome-wide, and for the subset of uORFs that are predicted to be active (predicted positive). **E. uORFs predicted as positive from genome-wide scan and ribosome profiling experiments.** Approximately 180 000 uORFs in the genome are predicted as active upstream open reading frames. This large set includes substantial proportions uORFs identified in the ribosome profiling experiments (~70% each). **F. Performance of the machine learning algorithm.** The machine learning algorithm was trained on two of three ribosome profiling data sets, and used to extract the third data set from among unlabeled examples. The ROC curve is shown for each of the three combinations: 1) Train Lee et al. and Fritsch et al. – extract Gao et al. (AUC = 0.79), 2) Train Lee et al. and Gao et al. – extract Fritsch et al. (AUC = 0.77). 3) Train Fritsch et al. and Gao et al. - extract Lee et al. (AUC = 0.82).



**Figure 4:**

[[PDM to MG: Emphasis in figure 4.D. on uORF creation (ATG formation), Figure 4.A now compares gain of uORF vs. loss of uORF (more appropriate comparison)]]

**A: Gene level protein expression change between individuals with variants interrupting predicted positive uORFs and wild type individuals. uORF gain is associated with increased protein expression, while uORF loss is associated with decreased protein expression. This difference in protein level is shown for different ratios of variant possessing individuals (+/, /) to wild type individuals (+/+). Larger numbers of individuals with the variant allele allow for larger statistical power in calculating the effect of the variant on protein level.** **B: rQTLs interrupting uORFs, according to score of the corresponding uORF.** rQTLs show bias towards interrupting positively predicted uORFs. **C: Density matrix showing the distribution of 1000 Genomes variants that interrupt predicted positive uORF start codons.** The vertical axis displays the reference start codon, the horizontal axis shows the interrupting variant (position – 1,2,3 – and codon – A,T,G,C). **D: Density matrix showing the distribution of somatic mutations found in exomic tumor samples that interrupt predicted positive uORF start codons.** The vertical axis displays the reference start codon, the horizontal axis shows the interrupting variant (position – 1,2,3 – and codon – A,T,G,C). ATG forming mutations are highlighted.

**Title:**

## A comprehensive catalog of predicted functional upstream open reading frames.

Patrick McGillivray, Russell Ault, Mayur Pawashe, Rob Kitchen, Suganthi Balasubramanian, Mark Gerstein

### Abstract

The activity of an upstream open reading frame (uORF) latent in an mRNA transcript is thought to modify translation of coding sequences in that same transcript by modifying local ribosome activity. Not all uORFs are thought to be active in such a process. It represents a challenge to estimate the impact and scope of the role uORFs play in regulation of translation.

[[PDM to MG and SB: More explicit language concerning procedure followed, and content of paper]] We first circumscribed the universe of all uORFs based on coding gene sequence. This universe includes over one million unique uORFs. In order to determine which of these uORFs are likely to be biologically relevant, we built a classifier using 89 attributes of uORFs labeled as active in experiment. Many of these attributes contribute toward accurate identification of active uORFs. This classifier allowed us to extrapolate to a catalog of uORFs that are likely active from the universe of all uORFs.

~~We first circumscribed the universe of all uORFs. This universe includes over one million unique uORFs. We compared patterns of structure in this complete set of uORFs, to the attributes of uORFs labeled as active in experiment. A classifier built using these attributes, was used to extrapolate a catalog of uORFs that are likely active. This is a substantially larger catalog of uORFs than has previously been associated with active function. Our ranked list of likely active uORFs, allows researchers to test their hypotheses regarding the role of upstream open reading frames in health and disease. We demonstrate several interesting examples of biological relevance through application of our catalog.~~

### **Intro**

Upstream open reading frames (uORFs) consist of a start codon in the 5' untranslated region of a gene (UTR) and an associated stop codon appearing before the stop codon of the main coding sequence (CDS). An uORF may begin and end before the main gene coding sequence.

Alternatively, if the upstream reading frame is out of frame with the CDS, it may overlap with the CDS [Figure 1.A]. uORFs are latent in mRNA transcripts and may undergo partial or complete translation.

An initial survey of the human genome identified uORFs contained in approximately 10% of mRNA transcripts (1). More recent analyses identify uORFs in association with nearly half of all mRNA transcripts (2). The discovery that many translated uORFs utilize near-cognate start codons to the canonical ATG start codon, has broadened estimates of uORF prevalence further (3–6).

Presence of functional uORFs is generally thought to suppress translation of downstream genes (7–12). Proposed molecular mechanisms for modification of CDS translation by uORFs are

numerous. These include *translation reinitiation* -- the uORF and CDS are translated by the same ribosome in series -- *leaky-scanning* -- ribosome recognition of an uORF and subsequent CDS translation, without uORF translation -- and *ribosome-stalling* -- decreased translation of the CDS, due to ribosome retention at the upstream uORF (3,13,14). Differential translation of multiple protein products may occur in consequence to an uORF (15). It is also possible for an uORF to function as short open reading frame, encoding a short functional peptide (16–19). uORF function is not necessarily constant -- uORFs may display differential function in stressed cells, compared with non-stressed controls (20–25).

Study of uORF translation and function, was historically limited to the experimental evaluation of individual uORFs (7,26). Genome-scale ribosome profiling studies have allowed for the identification of large populations of uORFs known to undergo translation (4,27,28). This mapping of translation initiation is sufficient for association between ribosomes and particular start codons and reading frames (29–31).

We proceed on the assumption that the total universe of active uORFs is much larger than that identified through ribosome profiling experiments. In other words, we assume that ribosome profiling experiments have high specificity in identifying functional uORFs with a high false-negative rate [Figure 1.C]. **[[PDM to MG: attempting more direct/clear language]] Ribosome profiling experiments follow a challenging technical procedure, and it is uncertain that all potentially active uORFs are measurable in a given sample** ~~Consistent with this perspective, is the hypothesis that uORFs display differential activity according to environmental condition or organ tissue. Ribosome profiling experiments may suffer from a form of sampling bias, incapable of detecting functional uORFs of transiently or locally decreased activity~~ [Figure 1.D]. ~~This is consistent with a high false-negative rate.~~ Other researchers have implicitly endorsed this hidden assumption, when predicting translated uORFs in *Saccharomyces cerevisiae* and *Arabidopsis thaliana*, on the basis of DNA sequence and ribosome profiling data (32,33). A similar assumption is the basis for using patterns of ribosome profiling occupancy to maximize the number of inferred translation products in humans (34,35).

For our investigation of the prevalence of active uORFs in humans, we began with a genome wide scan, searching for uORFs associated with protein coding genes listed in the GENCODE genome annotation (36). All possible uORFs beginning with ATG, or a single nucleotide variant of ATG, were identified. This scan yields a universe of all possible uORFs, numbering nearly 1.3 million.

uORFs in this large set were classified as active according to similarity to uORFs occupied in ribosome profiling experiments. This classification was accomplished using a Naïve-Bayes classifier, trained on 89 uORF attributes. We validated our predicted uORFs using a cross-validation method where two ribosome profiling experiments are used to predict the uORFs translated in a third experiment. We also validated our predictions by examining how **[[PDM to MG: simplifying + clarifying language]] gene expression and ribosome activity varies in response to genetic variants that alter uORFs.** ~~individual genotype altering uORF sequence affects parameters related to gene level control of translation by uORFs: protein level from the downstream gene, and ribosome occupancy.~~

MURZ



The 1000 Genomes Project's database of human variation (37) and the NHGRI-EBI GWAS catalog (38) were used to provide a baseline for the functional consequence of our predicted active uORFs. The predictions we generated were also used to measure the functional impact of somatic mutations affecting uORFs, in tissue-matched tumor samples (39).

We provide a resource of predicted active uORFs for other scientists to use in their effort to understand uORF function in health and disease.

## **Methods:**

### *Extracting uORFs from GENCODE:*

uORFs were identified through genome-wide search, performed on v19 of GENCODE's human genome annotation (36). uORFs were defined as a start codon within the 5'UTR and a downstream stop codon before the end of the CDS. All three possible reading frames were examined. ATG and near cognate start codons were included in this search [ATG, TTG, GTG, CTG, AAG, AGG, ACG, ATA, ATT, ATC].

### *Ribosome profiling experiments as a reference set:*

The ribosome profiling experiments of Lee et al. (2012), Fritsch et al. (2012) and Gao et al. (2014), were used to obtain an experimentally validated set of translated upstream open reading frames. These studies identify translation initiation sites (TIS) through treatment of human cell lines with antibiotic translation inhibitors. These treatments reliably halt ribosomes in predictable proximity to the start codon (12-13 nucleotides downstream). As such, these experiments provide high resolution information about translation initiation sites in the human genome.

We employed the read alignments and identification of the translation initiation sites as provided by these three groups of researchers. The cell lines, treatment protocols, and TIS identification mechanism employed by each of these three research groups is summarized in *Supplement - Methods*.

### *Literature review of translated human uORFs:*

In addition to ribosome profiling studies, confirmed translated uORFs were obtained from the biomedical literature (7,40,41). uORFs studied in humans that displayed functionality -- demonstrated regulation of the CDS product -- were added to the set of positive uORFs. In total, 33 uORFs, associated with 33 separate genes, were included from this literature review.

### *Cleansing the data set, by removal of N-terminal extensions and aTISs, and isolation of unique transcript IDs:*

N-terminal extensions of the CDS sequence, may retain some functional activity of the primary gene protein product, and were removed from the data set. Any uORF start codon annotated as an alternative translation initiation site (aTIS) for the CDS, was also removed from the data set.

Multiple transcripts may share the same uORF. In order to avoid over-counting, only one transcript ID is attributed to a given uORF. This selection was made randomly, from among transcripts with identical chromosomal coordinates.

Positive1-voted, neutral2-voted, and unlabeled data sets:

uORFs were divided into three separate sets, according to their experimental translation status:

[[PDM to MG: change in terminology, to decrease ambiguity]]Positive2-voted: uORFs identified as translated in two or more ribosome profiling experiments, or through literature review.

Neutral1-voted: uORFs identified as translated in not more than one ribosome profiling experiment.

Unlabeled: uORFs that were not identified as translated in any ribosome profiling experiment, or through literature review.

Estimating the total population of active uORFs:

Based on observed overlap among ribosome profiling experiments, an estimate for the total number of active uORFs was made using methods borrowed from population biology.

Ribosome profiling experiments are treated as independent population samplings, and the Schnabel equation (Eq. 1) or Schumacher and Eschmeyer equation (Eq. 2) provide a population size estimate:

$$\hat{N} = \frac{\sum_{t=1}^S (C_t M_t)}{\sum_{t=1}^S R_t} \quad (1)$$

$$\hat{N} = \frac{\sum_{t=1}^S (C_t M_t^2)}{\sum_{t=1}^S R_t M_t} \quad (2)$$

Where  $\hat{N}$  is an estimate of the number of individuals in a population, given a series of  $S$  samplings taken at times  $t \in \{1 \dots S\}$ , with  $C_t$  the number of individuals 'captured' in a sample,  $M_t$  the total number of marked individuals prior to sampling at time  $t$ , and  $R_t$  the number of marked individuals 'recaptured' at sampling  $t$ .

*Extraction of attributes associated with uORFs:*

Feature data was extracted for each uORF. Features were chosen to cover a broad range of categories of data, including features associated with uORF structure, uORF evolutionary conservation, and genomic context. 89 features were used. A complete listing of these features, including details relating to the extraction and calculation of each feature, is included in *Methods Supplement*.

*Feature discretization:*

The minimum description length principle (MDLP) algorithm was used to discretize each of our chosen attributes (42). The MDLP algorithm minimizes information lost through discretization. MDLP discretization was implemented using the 'discretization' package available for R (<http://cran.r-project.org/web/packages/discretization/index.html>).

*Prioritization of feature data:*

The distribution for each feature was compared between positive and unlabeled uORFs using the Kolmogorov-Smirnov (KS) statistic. A greater KS statistic, suggests greater ability of that attribute, to distinguish between positive and unlabeled features.

*Classifying uORFs, according to attributes:*

We determined that attributes of an uORF were consistent with an active uORF, according to a Naive-Bayes machine learning algorithm applied to positive and unlabeled examples (43):

$$P_{pos} \sum_{i=1}^N p(A_i|pos) = p_{pos} \tag{3}$$

$$P_{neg} \sum_{i=1}^N p(A_i|unl) = p_{neg} \tag{4}$$

Where:

$$P_{neg} + P_{pos} = 1 \tag{5}$$

$P_{pos}$  is the prior probability associated with positive uORFs.  $P_{pos}$  was chosen as the F1 score maximizing value (0.61).  $p(A_i|pos)$ , and  $p(A_i|unl)$  represent the frequency of that attribute value among the positive and unlabeled sets respectively.  $p_{pos}$  represents the probability the uORF is positive.  $p_{neg}$  represents the probability the uORF is negative. We label an uORF as positive or negative according to the greater value between  $p_{pos}$  and  $p_{neg}$ . We note likely violation of the feature independence requirement of Naive-Bayes. However, empirical and theoretical study has demonstrated optimal classification performance, even where feature independence does not hold (44,45).

*Model validation:*

Our model was serially trained on two of three ribosome profiling data sets, using the trained model to extract a third withheld ribosome profiling data set from among the unlabeled

examples. The success of differentially trained models in this cross-validation, was evaluated using ROC curves, with area under the curve (AUC) calculated for each curve.

As biologic validation of our predicted uORFs we examined the effect of alteration of a predicted active uORF start codon on gene protein levels and local ribosome occupancy. Protein levels and local ribosome quantitative trait loci (cis-rQTL) for 47 individuals were obtained from the ribosome profiling and proteomic experiments of Battle et al. 2015 (46). Individual genotype information for 47 individuals in the Battle et al. study, is provided by the 1000 Genomes Project. Gene expression change was evaluated in association with both gain of predicted positive uORFs (ATG and CTG), and loss of predicted positive uORFs.

#### *Natural variation affecting predicted positive uORFs:*

Natural variant SNPs affecting the start codons of predicted positive uORFs, were obtained from the 1000 Genomes project. The subset of these SNPs that are associated with differential disease susceptibility was identified through search of the NHGRI-EBI GWAS database. Measurement of comparative frequency of mutation among uORF start codons, was taken as a measure of evolutionary conservation and functional significance of predicted positive uORFs.

#### *Cancer mutation affecting predicted positive uORFs:*

The study of Alexandrov et al. 2012 (39) provides a set of exomic somatic mutations according to patient sample, and cancer type. We used these mutations, as a comparison standard for the healthy 1000 Genomes Project population. We identified start codons of our predicted positive uORFs altered by somatic mutation in cancer.

### **Results:**

Genome-wide search yielded 1 270 265 unique uORFs. Within this large set, we isolated the subset of uORFs identified as translated in the studies of Lee et al. 2012, Fritsch et al. 2012, and Gao et al. 2014. We further stratified this set of translated uORFs according to shared representation of uORFs among the three studies. uORFs identified in the intersection between two or more of these studies were used as the reference standard for functional uORFs. Literature review yielded 33 additional examples of active uORFs that were also included in the set of positive, functional uORFs.

We followed the procedure outlined in Figure 2.A to isolate uORFs that are likely to be active. Distributions of attributes for positive, translated uORFs were compared with distributions of those same attributes observed in the set of unlabeled uORFs [Figure 2.B]. The KS statistic and corresponding p-value for each of the 89 attributes assessed in this study are provided in *Supplement Table 2*. The top 10 attributes listed according to magnitude of KS statistic are given in Figure 2.C. From this prioritization of attributes, we can draw insights into the relationship between uORF structure and function. The presence of large numbers of start codons within a single uORF is a high priority attribute for positive classification, as is a shorter distance

between the uORF and the CDS. ATG is the start codon associated with greatest functional significance. Start and stop codons of functional uORFs are generally located in evolutionarily conserved sites suggesting a meaningful physiologic role.

Overlap between the three ribosome profiling experiments was found to be low, with pairwise intersections of 12.2% (Gao  $\cap$  Fritsch), 9.2% (Gao  $\cap$  Lee), and 9.8% (Lee  $\cap$  Fritsch). The number of uORFs shared between all three sets represents only 3.3% of uORFs identified in these studies [Figure 3.A]. If independent ribosome profiling experiments represent resampling of the same population, repeat identification of uORFs among experiments yields an estimate of the total number of functional uORFs. 10 000 functional uORFs are estimated in this way to be present in the human genome using the Schnabel equation (Eq. 1) or Schumacher and Eschmeyer equation (Eq. 2) (47,48).

CTG (28.2%) and ATG (46.1%) are the most prevalent start codons identified in ribosome profiling experiments. CTG (30.5%) and ATG (34.6%) continue to represent the majority of start codons in intersection between ribosome profiling experiments [Figure 3.B.]. Representation of every near-cognate start codon was found in intersections between studies, with the exception of AAG and AGG. This indicates that uORFs do not generally employ AAG and AGG as start codons. Identification of uORFs beginning with AAG or AGG in ribosome profiling experiments, may represent false-positives.

Discretized attributes of positive and unlabeled sets of uORFs were used to build a statistical classifier within a Naive-Bayes framework. The result of application of the classifier is shown in figure 3.C. 76.8% of **positive-2-voted positive\_uORFs** [590/768], 67.1% of **neutral-1-voted positive\_uORFs** [2379/3543], and 14.7% of unlabeled uORFs [185833/1265954] are ultimately classified as likely active. A total of 14.9% of all uORFs are identified as likely active [188802/1270265]. A complete list of upstream open reading frames predicted to be active, is provided in *Supplement -- Results*. The 10% highest probability examples are also specified.

A large proportion of uORFs in the human genome begin with CTG start codons (19.3%). The greatest number of predicted positive uORFs are also initiated with a CTG start codon (11.8%). ATG has a lower comparative prevalence in the human genome and in the predicted positive set (6.7% and 8.2% respectively) [Figure 3.D]. 8 genes are associated with greater than 200 positively scored uORFs (FAM156B, FAM156A, EEF1D, UBA1, C6orf62, HMGB1, HP1BP3, TBC1D5), suggesting that these genes are under strong and redundant translational regulation mediated by uORFs. The proportion of uORFs ultimately identified as positive from each ribosome profiling study, is shown in Figure 3.E. The results were similar for each of the ribosome profiling experiments, approximately 70% in each case (72% of Gao, 71% of Lee, 70% of Fritsch).

**[[PDM to MG: next 3 paragraphs transition to validation info.]]**

As a validation of our technique, we serially excluded one of three ribosome profiling experiments from the positive training set, instead including the excluded set among unlabeled

examples for subsequent retrieval [Figure 3.F]. The AUC for each of the ROC curves corresponding to these trials is similar: 0.82, 0.79, and 0.77. This suggests a high false-negative rate for ribosome profiling studies; predicted active uORFs, reflect those uORFs that additional experiments would discover are translated.

As experimental validation of our technique, we examined how natural variation affecting our predicted active uORFs, alters protein level and ribosome localization in humans. We hypothesized that an active uORF altered by naturally occurring variants, should create observable effect on ribosome occupancy and protein levels from that gene. The results of Battle et al. 2015, supplemented by genotype information from the 1000 Genomes Project, provide the basis for validation of our predictions in 47 human subjects. In this natural study, variants causing gain of predicted positive ATG or CTG uORFs are associated with increase in downstream gene expression. Variants that cause loss of predicted positive uORFs, are associated with decrease in downstream gene expression. In this natural study, alteration of a predicted active uORF start codon results in a decrease in protein levels from downstream genes [Figure 4.A]. There is statistically significant difference in gene expression between variants causing uORF gain compared with uORF loss, among variants with approximate balance between individuals with and without the variant (increased power).

For these same 47 human subjects, cis-rQTLs provide an inventory of variants with statistically significant effect on local ribosome occupancy. There is significant enrichment for rQTLs interrupting positively scored start codons [Figure 4.B]. [[PDM to MG: clarifying meaning of the 14.9% expectation]] If mutations hit uORFs randomly, 14.9% of the time they would hit a positively scored uORF. While the effect we would expect due to random mutation is 14.9%, However, we observe that 48% of these rQTLs (21/44) interrupt positively scored start codons - a 3x higher rate. This indicates that many rQTLs may measure the direct effect of disruption of functional uORFs.

The ATG start codon is relatively conserved among predicted positive start codons -- it is rarely interrupted by 1000 Genomes Project variants (relative rate (RR) 0.03), suggesting its functional importance. The CTG start codon, although more prevalent among predicted positive uORFs, is altered relatively frequently by natural human variants (RR 0.52) [Figure 4.C]. In exomic tumor samples from cancer patients, CTG is the most commonly modified predicted positive uORF start codon. ATG is interrupted at a RR of 0.25 in comparison to CTG [Figure 4.D]. The higher RR of interruption of both ATG and CTG in cancer as compared to germline variants – 8 fold higher, and 2 fold higher respectively – further suggests functional consequences attributable to these uORFs.

Exomic cancer mutations breaking the highest scored uORFs, are listed in *Supplemental Table 3*. These mutations interrupt uORFs associated with well-studied oncogenes and tumor suppressors. MYC and BCL2 are the two genes associated with the greatest recurrence of uORF interruptions, and we identify recurrent mutation of positively scored uORFs associated with PTEN, TP53, ERCC1, and MSH5. GWAS SNPs listed in the NHGRI-EBI GWAS database that impact our predicted uORFs are listed in *Supplemental Table 4*. GWAS diseases associated with

SNVs interrupting positively scored uORFs include prevalent chronic conditions like obesity (rs11603334), osteoporosis (rs3755955), asthma (rs3771180), and type 2 diabetes (rs1552224). Additional variants associated with susceptibility and prognosis in cancer are found to interrupt positively scored uORFs, like rs779805 upstream of the VHL gene, and rs34330 upstream of CDKN1B. **[[PDM to SB: added text to convey uncertainty about GWAS results]]** Although linkage disequilibrium and overlap among regulatory elements complicates interpretation of these GWAS studies, these disease associated SNVs, may owe their functional consequence to alteration of a translated uORF.

## Discussion:

In this study, we identify 188 802 likely active upstream open reading frames, from a genome-wide set of 1 270 265 unique uORFs. We further highlight the 10% of our predictions that are most likely to be functional, as a high reliability subset.

We began by assuming that ribosome profiling experiments have a high false negative rate for identification of functional uORFs. Our method applied the intersection of three ribosome profiling studies, to form a reference set of known active uORFs. The low overlap between ribosome profiling experiments suggests a high false-negative rate in individual experiments. The finding that pairs of ribosome profiling experiments may be used to correctly identify the uORFs translated in a third experiment also suggests a high false negative rate. The large number of uORFs we identified as likely functional is consistent with this premise, but **remarkable/significant** in comparison to other studies on the topic.

There is precedent for our findings, in comparisons of large-scale parallel experiments of interaction between biomolecules. The protein-protein interaction experiments of Uetz et al. employed a comprehensive, genome-wide scope (49). Subsequent experiments by Ito et al., with similar technique and scope, showed low overlap with results of the prior project (50). **[[PDM to MG: simplify, and cite Current Opinion in Microbiology paper]]** It became clear that both experiments had relatively high false-negative rates. The the universe of possible protein-protein interactions, is much larger than identified in either experiment individually. Benefit in identifying these interactions, is achieved by combining datasets (51).

Our use of an intersection between ribosome profiling experiments, provides some control against differences experimental conditions and tissue specific results (both HEK293 and THP-1 cells were examined). However, just as protein levels vary widely across cell-types (52) it may prove that the activity of uORFs varies considerably across cell types and cellular conditions. Analysis of cell-type specific and condition specific activity of uORFs may further expand estimates of the population of uORFs.

Our study helps clarify how attributes of structure and context of a given uORF -- including start codon, base composition, and relative position to the CDS -- likely contribute to varying functionality among uORFs. Although ATG is the most common uORF start codon identified in ribosome profiling experiments, lower affinity near cognate-start codons may have great functional impact on the landscape of translation, due to their overall abundance.



An important validation of our predictions, is the finding that alteration of predicted functional uORFs as a consequence of germline genetic variation, impacts ribosome binding and protein level in humans. ~~This is contrary to common view that uORFs act as translational repressors. It is of interest, that~~ ~~Generally we assume that uORFs act as translational repressors. However,~~ ~~the overall effect of uORF loss, appears to be a decrease in downstream protein level.~~ ~~This is contrary to common view that uORFs act as translational repressors.~~ Mechanisms have been studied, where uORFs act to up-regulate expression of a downstream coding sequence (e.g. leaky-scanning, and translation reinitiation). Our analysis suggests that this effect is a more common consequence for upstream open reading frames than is previously credited.

Applications of our results, suggest avenues for future research. Identification of human germline variants altering predicted positive uORFs, reveals locations where the creation or destruction of an uORF, is likely to alter protein levels. Employing this method, we identified disease associated SNVs -- including a number of GWAS SNVs -- that likely owe their significance to alteration of a functional uORF. Among diseases, our work could be used to help broaden knowledge of the role of uORFs in cancer beyond recently identified individual examples (53).

We provide a catalog that can serve as a point of reference for other researchers engaged in the investigation of uORF function.

1. Kozak M. An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic Acids Res* [Internet]. 1987 [cited 2016 Aug 16];15(20):8125–48. Available from: <http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/15.20.8125>
2. Kochetov A V., Sarai A, Rogozin IB, Shumny VK, Kolchanov NA. The role of alternative translation start sites in the generation of human protein diversity. *Mol Genet Genomics* [Internet]. 2005 Jul 15 [cited 2016 Aug 16];273(6):491–6. Available from: <http://link.springer.com/10.1007/s00438-005-1152-7>
3. Ingolia NT, Lareau LF, Weissman JS. Ribosome Profiling of Mouse Embryonic Stem Cells Reveals the Complexity and Dynamics of Mammalian Proteomes. *Cell*. 2011;147(4):789–802.
4. Ingolia NT, Ghaemmaghami S, Newman JRS, Weissman JS. Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science* [Internet]. 2009 Apr 10 [cited 2016 Aug 16];324(5924):218–23. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/19213877>
5. Ivanov IP, Loughran G, Atkins JF. uORFs with unusual translational start codons autoregulate expression of eukaryotic ornithine decarboxylase homologs. *Proc Natl Acad Sci* [Internet]. 2008 Jul 22 [cited 2016 Aug 23];105(29):10079–84. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.0801590105>
6. Ivanov IP, Firth AE, Michel AM, Atkins JF, Baranov P V. Identification of evolutionarily conserved non-AUG-initiated N-terminal extensions in human coding sequences. *Nucleic Acids Res* [Internet]. 2011 May 1 [cited 2016 Aug 17];39(10):4220–34. Available from:



<http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/gkr007>

7. Calvo SE, Pagliarini DJ, Mootha VK. Upstream open reading frames cause widespread reduction of protein expression and are polymorphic among humans. *Proc Natl Acad Sci* [Internet]. 2009 May 5 [cited 2016 Aug 16];106(18):7507–12. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.0810916106>
8. Johnstone TG, Bazzini AA, Giraldez AJ, Abramoff M, Magalhães P, Ram S, et al. Upstream ORFs are prevalent translational repressors in vertebrates. *EMBO J* [Internet]. 2016 Apr 1 [cited 2016 Aug 17];35(7):706–23. Available from: <http://emboj.embopress.org/lookup/doi/10.15252/embj.201592759>
9. Somers J, Pöyry T, Willis AE. A perspective on mammalian upstream open reading frame function. *Int J Biochem Cell Biol*. 2013;45(8):1690–700.
10. Meijer HA, Thomas AAM. Control of eukaryotic protein synthesis by upstream open reading frames in the 5'-untranslated region of an mRNA. *Biochem J* [Internet]. 2002 Oct 1 [cited 2016 Aug 17];367(Pt 1):1–11. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/12117416>
11. Barbosa C, Peixeiro I, Romão L, Morris D, Geballe A, Calvo S, et al. Gene Expression Regulation by Upstream Open Reading Frames and Human Disease. Fisher EMC, editor. *PLoS Genet* [Internet]. 2013 Aug 8 [cited 2016 Aug 17];9(8):e1003529. Available from: <http://dx.plos.org/10.1371/journal.pgen.1003529>
12. Morris DR, Geballe AP. Upstream Open Reading Frames as Regulators of mRNA Translation. *Mol Cell Biol* [Internet]. 2000 Dec 1 [cited 2016 Aug 17];20(23):8635–42. Available from: <http://mcb.asm.org/cgi/doi/10.1128/MCB.20.23.8635-8642.2000>
13. Sonenberg N, Hinnebusch AG. Regulation of translation initiation in eukaryotes: mechanisms and biological targets. *Cell* [Internet]. 2009 Feb 20 [cited 2016 Aug 30];136(4):731–45. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/19239892>
14. Hinnebusch AG, Ivanov IP, Sonenberg N, Hinnebusch AG, Kozak M, Starck SR, et al. Translational control by 5'-untranslated regions of eukaryotic mRNAs. *Science* [Internet]. 2016 Jun 17 [cited 2016 Aug 17];352(6292):1413–6. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/27313038>
15. Chua JJE, Schob C, Rehbein M, Gkogkas CG, Richter D, Kindler S, et al. Synthesis of two SAPAP3 isoforms from a single mRNA is mediated via alternative translational initiation. *Sci Rep* [Internet]. 2012 Jul 2 [cited 2016 Aug 16];2:277–98. Available from: <http://www.nature.com/articles/srep00484>
16. Mackowiak SD, Zauber H, Bielow C, Thiel D, Kutz K, Calviello L, et al. Extensive identification and analysis of conserved small ORFs in animals. *Genome Biol* [Internet]. 2015 Dec 14 [cited 2016 Aug 23];16(1):179. Available from: <http://genomebiology.com/2015/16/1/179>
17. Andrews SJ, Rothnagel JA. Emerging evidence for functional peptides encoded by short

- open reading frames. *Nat Rev Genet* [Internet]. 2014 Mar [cited 2016 Aug 17];15(3):193–204. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/24514441>
18. Oyama M, Itagaki C, Hata H, Suzuki Y, Izumi T, Natsume T, et al. Analysis of Small Human Proteins Reveals the Translation of Upstream Open Reading Frames of mRNAs. *Genome Res* [Internet]. 2004 Oct 15 [cited 2016 Aug 17];14(10b):2048–52. Available from: <http://www.genome.org/cgi/doi/10.1101/gr.2384604>
  19. Bergeron D, Lapointe C, Bissonnette C, Tremblay G, Motard J, Roucou X. An Out-of-frame Overlapping Reading Frame in the Ataxin-1 Coding Sequence Encodes a Novel Ataxin-1 Interacting Protein. *J Biol Chem* [Internet]. 2013 Jul 26 [cited 2016 Aug 17];288(30):21824–35. Available from: <http://www.jbc.org/cgi/doi/10.1074/jbc.M113.472654>
  20. Starck SR, Tsai JC, Chen K, Shodiya M, Wang L, Yahiro K, et al. Translation from the 5' untranslated region shapes the integrated stress response. *Science* [Internet]. 2016 Jan 29 [cited 2016 Aug 17];351(6272):aad3867. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/26823435>
  21. Andreev DE, O'Connor PBF, Zhdanov A V, Dmitriev RI, Shatsky IN, Papkovsky DB, et al. Oxygen and glucose deprivation induces widespread alterations in mRNA translation within 20 minutes. *Genome Biol* [Internet]. 2015 Dec 6 [cited 2016 Aug 17];16(1):90. Available from: <http://genomebiology.com/2015/16/1/90>
  22. Shalgi R, Hurt JA, Krykbaeva I, Taipale M, Lindquist S, Burge CB. Widespread Regulation of Translation by Elongation Pausing in Heat Shock. *Mol Cell*. 2013;49(3):439–52.
  23. Wiita AP, Ziv E, Wiita PJ, Urisman A, Julien O, Burlingame AL, et al. Global cellular response to chemotherapy-induced apoptosis. *Elife* [Internet]. 2013 Jul [cited 2016 Aug 17];2(1):e01236. Available from: <http://linkinghub.elsevier.com/retrieve/pii/S1097276507004005>
  24. Gerashchenko M V., Lobanov A V., Gladyshev VN. Genome-wide ribosome profiling reveals complex translational regulation in response to oxidative stress. *Proc Natl Acad Sci* [Internet]. 2012 Oct 23 [cited 2016 Aug 17];109(43):17394–9. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.1120799109>
  25. Liu B, Han Y, Qian S-B. Cotranslational Response to Proteotoxic Stress by Elongation Pausing of Ribosomes. *Mol Cell*. 2013;49(3):453–63.
  26. Kozak M. Point mutations define a sequence flanking the AUG initiator codon that modulates translation by eukaryotic ribosomes. *Cell* [Internet]. 1986 Jan [cited 2016 Aug 17];44(2):283–92. Available from: <http://linkinghub.elsevier.com/retrieve/pii/0092867486907622>
  27. Brar G a, Weissman JS. Ribosome profiling reveals the what, when, where and how of protein synthesis. *Nat Rev Mol Cell Biol* [Internet]. 2015;16(11):651–64. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/26465719>

28. Ingolia NT. Ribosome profiling: new views of translation, from single codons to genome scale. *Nat Rev Genet* [Internet]. 2014 Jan 28 [cited 2016 Aug 17];15(3):205–13. Available from: <http://www.nature.com/doi/10.1038/nrg3645>
29. Gao X, Wan J, Liu B, Ma M, Shen B, Qian S-B. Quantitative profiling of initiating ribosomes in vivo. *Nat Methods* [Internet]. 2014 Dec 8 [cited 2016 Aug 17];12(2):147–53. Available from: <http://www.nature.com/doi/10.1038/nmeth.3208>
30. Fritsch C, Herrmann A, Nothnagel M, Szafranski K, Huse K, Schumann F, et al. Genome-wide search for novel human uORFs and N-terminal protein extensions using ribosomal footprinting. *Genome Res* [Internet]. 2012 Nov 1 [cited 2016 Aug 17];22(11):2208–18. Available from: <http://genome.cshlp.org/cgi/doi/10.1101/gr.139568.112>
31. Lee S, Liu B, Lee S, Huang S-X, Shen B, Qian S-B. Global mapping of translation initiation sites in mammalian cells at single-nucleotide resolution. *Proc Natl Acad Sci* [Internet]. 2012 Sep 11 [cited 2016 Aug 17];109(37):E2424–32. Available from: <http://www.pnas.org/cgi/doi/10.1073/pnas.1207846109>
32. Selpi S, Bryant CH, Kemp GJ, Sarv J, Kristiansson E, Sunnerhagen P, et al. Predicting functional upstream open reading frames in *Saccharomyces cerevisiae*. *BMC Bioinformatics* [Internet]. 2009 [cited 2016 Aug 17];10(1):451. Available from: <http://www.biomedcentral.com/1471-2105/10/451>
33. Hu Q, Merchante C, Stepanova AN, Alonso JM, Heber S. Genome-Wide Search for Translated Upstream Open Reading Frames in *Arabidopsis Thaliana*. *IEEE Trans Nanobioscience* [Internet]. 2016 Mar [cited 2016 Aug 31];15(2):148–57. Available from: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=7404026>
34. Fields AP, Rodriguez EH, Jovanovic M, Stern-Ginossar N, Haas BJ, Mertins P, et al. A Regression-Based Analysis of Ribosome-Profiling Data Reveals a Conserved Complexity to Mammalian Translation. *Mol Cell*. 2015;60(5):816–27.
35. Raj A, Wang SH, Shim H, Harpak A, Li YI, Engelmann B, et al. Thousands of novel translated open reading frames in humans inferred by ribosome footprint profiling. *Elife* [Internet]. 2016 [cited 2016 Aug 31];5. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/27232982>
36. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et al. GENCODE: The reference human genome annotation for The ENCODE Project. *Genome Res* [Internet]. 2012 Sep 1 [cited 2016 Aug 21];22(9):1760–74. Available from: <http://genome.cshlp.org/cgi/doi/10.1101/gr.135350.111>
37. Project Consortium G, Consortium Participants are arranged by project role G, by institution alphabetically then, alphabetically within institutions except for Principal Investigators finally, Leaders P, indicated as, et al. An integrated map of genetic variation from 1,092 human genomes. *Nature*. 2012;490.
38. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P, Manolio T, Hindorff L and PH. The NHGRI GWAS Catalog, a curated resource of SNP-trait

- associations. *Nucleic Acids Res.* 2014;Vol. 42((Database issue)):D1001–6.
39. Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SAJR, Behjati S, Biankin A V., et al. Signatures of mutational processes in human cancer. *Nature* [Internet]. 2013 Aug 14 [cited 2016 Aug 17];500(7463):415–21. Available from: <http://www.nature.com/doi/10.1038/nature12477>
  40. Wen Y, Liu Y, Xu Y, Zhao Y, Hua R, Wang K, et al. Loss-of-function mutations of an inhibitory upstream ORF in the human hairless transcript cause Marie Unna hereditary hypotrichosis. *Nat Genet* [Internet]. 2009 Feb 4 [cited 2016 Aug 31];41(2):228–33. Available from: <http://www.nature.com/doi/10.1038/ng.276>
  41. Raveh-Amit H, Maissel A, Poller J, Marom L, Elroy-Stein O, Shapira M, et al. Translational Control of Protein Kinase C by Two Upstream Open Reading Frames. *Mol Cell Biol* [Internet]. 2009 Nov 15 [cited 2016 Aug 31];29(22):6140–8. Available from: <http://mcb.asm.org/cgi/doi/10.1128/MCB.01044-09>
  42. Fayyad U, Irani K. Multi-Interval Discretization of Continuous-Valued Attributes for Classification Learning. *donga.ac.kr* [Internet]. [cited 2016 Aug 17]; Available from: [http://web.donga.ac.kr/kjunwoo/files/Multi interval discretization of continuous valued attributes for classification learning.pdf](http://web.donga.ac.kr/kjunwoo/files/Multi%20interval%20discretization%20of%20continuous%20valued%20attributes%20for%20classification%20learning.pdf)
  43. Liu B, Dai Y, Li X, Lee WS, Yu PS. Building text classifiers using positive and unlabeled examples. In: *Third IEEE International Conference on Data Mining* [Internet]. IEEE Comput. Soc; 2003 [cited 2016 Aug 17]. p. 179–86. Available from: <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1250918>
  44. Rish I. An empirical study of the naive Bayes classifier. *researchgate.net* [Internet]. [cited 2016 Sep 19]; Available from: [https://www.researchgate.net/profile/Irina\\_Rish/publication/228845263\\_An\\_Empirical\\_Study\\_of\\_the\\_naive\\_Bayes\\_Classifier/links/00b7d52dc3ccd8d692000000.pdf](https://www.researchgate.net/profile/Irina_Rish/publication/228845263_An_Empirical_Study_of_the_naive_Bayes_Classifier/links/00b7d52dc3ccd8d692000000.pdf)
  45. Zhang H. The Optimality of Naive Bayes. *AA* [Internet]. 2004 [cited 2016 Sep 19]; Available from: [http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/Optimality\\_of\\_Naive\\_Bayes.pdf](http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/Optimality_of_Naive_Bayes.pdf)
  46. Battle A, Khan Z, Wang SH, Mitrano A, Ford MJ, Pritchard JK, et al. Genomic variation. Impact of regulatory variation from RNA to protein. *Science* [Internet]. 2015 Feb 6 [cited 2016 Aug 31];347(6222):664–7. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/25657249>
  47. Schnabel ZE. The Estimation of Total Fish Population of a Lake. *Am Math Mon* [Internet]. 1938 Jun [cited 2016 Sep 19];45(6):348. Available from: <http://www.jstor.org/stable/2304025?origin=crossref>
  48. Schumacher, F. X. and Eschmeyer RW. The estimation of fish populations in lakes and ponds. *J Tennessee Acad Sci* . 1943;18(228–249).
  49. Fields S, Uetz P, Giot L, Cagney G, Mansfield TA, Judson RS, et al. A comprehensive

analysis of protein [ndash] protein interactions in *Saccharomyces cerevisiae*. *Nature* [Internet]. 2000 Feb 10 [cited 2016 Oct 17];403(6770):623–7. Available from: <http://www.nature.com/doi/10.1038/35001009>

50. Ito T, Chiba T, Ozawa R, Yoshida M, Hattori M, Sakaki Y. A comprehensive two-hybrid analysis to explore the yeast protein interactome. *Proc Natl Acad Sci U S A* [Internet]. 2001 Apr 10 [cited 2016 Oct 17];98(8):4569–74. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/11283351>
51. Jansen R, Gerstein M. Analyzing protein function on a genomic scale: the importance of gold-standard positives and negatives for network prediction. *Curr Opin Microbiol*. 2004;7(5):535–45.
52. Pontén F, Gry M, Fagerberg L, Lundberg E, Asplund A, Berglund L, et al. A global view of protein expression in human cells, tissues, and organs. *Mol Syst Biol* [Internet]. 2009 Dec 22 [cited 2016 Aug 17];5(1):799–816. Available from: <http://msb.embopress.org/cgi/doi/10.1038/msb.2009.93>
53. Wethmar K, Schulz J, Muro EM, Talyan S, Andrade-Navarro MA, Leutz A. Comprehensive translational control of tyrosine kinase expression by upstream open reading frames. *Oncogene* [Internet]. 2016 Mar 31 [cited 2016 Aug 17];35(13):1736–42. Available from: <http://www.nature.com/doi/10.1038/onc.2015.233>