

What can Ribo-seq and proteomics tell us about the non-canonical proteome?

Author list

John R. Prensner,¹ Jennifer G. Abelin,² Leron W. Kok,³ Karl R. Clauser,² Jonathan M. Mudge,⁴ Jorge Ruiz-Orera,⁵ Michal Bassani-Sternberg,⁶⁻⁸ Eric W. Deutsch,⁹ Sebastiaan van Heesch³

Affiliations

¹Department of Pediatrics, Division of Pediatric Hematology/Oncology, University of Michigan Medical School, Ann Arbor, MI 48109, USA

²Broad Institute of MIT and Harvard, Cambridge, MA, 02142, USA

³Princess Máxima Center for Pediatric Oncology, Heidelberglaan 25, 3584 CS, Utrecht, the Netherlands

⁴European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

⁵Cardiovascular and Metabolic Sciences, Max Delbrück Center for Molecular Medicine in the Helmholtz Association (MDC), 13125 Berlin, Germany

⁶Ludwig Institute for Cancer Research, University of Lausanne, Agora Center Bugnon 25A, 1005 Lausanne, Switzerland

⁷Department of Oncology, Centre hospitalier universitaire vaudois (CHUV), Rue du Bugnon 46, 1005 Lausanne, Switzerland

⁸Agora Cancer Research Centre, 1011 Lausanne, Switzerland

⁹Institute for Systems Biology (ISB), Seattle, Washington 98109, USA

Word count: 7,827

Keywords: Ribo-seq, mass spectrometry, immunopectidomics, non-canonical open reading frame, microprotein

Running title: Ribo-seq and the non-canonical proteome [35 characters]

Abbreviations: HLA, human leukocyte antigen; MHC, major histocompatibility complex; CDS, coding sequence; LC-MS/MS, liquid chromatography with tandem mass spectrometry; MS, mass spectrometry

Address correspondence to:

John R. Prensner, MD, PhD

Department of Pediatrics

Medical Science Research Building II, Room 2560B

1150 Medical Center Drive

Ann Arbor, MI 48109

Email: prensner@umich.edu

Phone: 734-763-5939

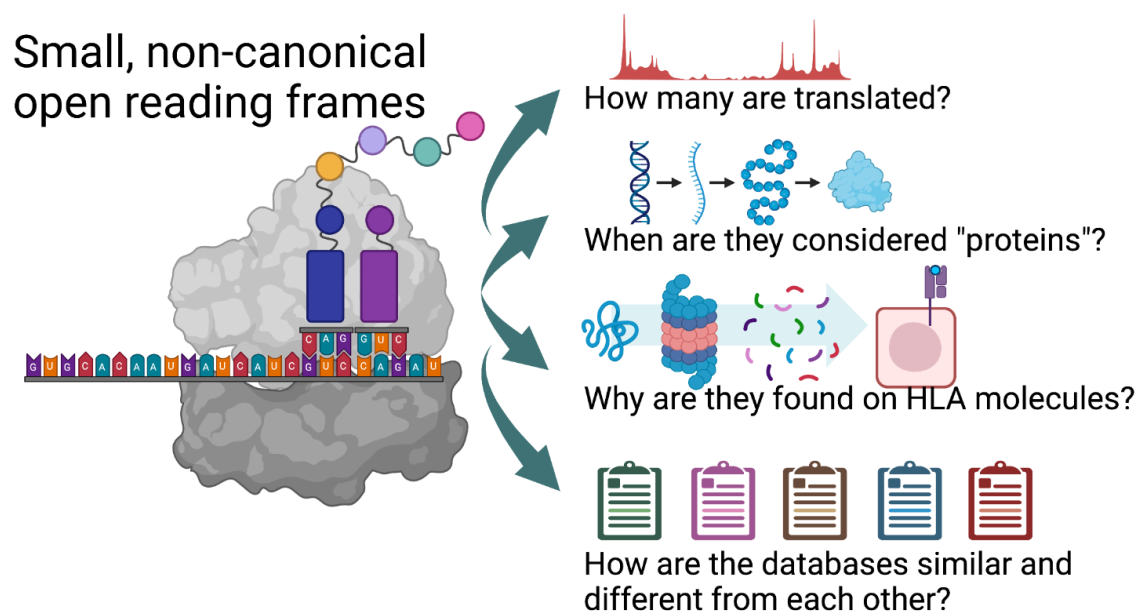
In brief

The human genome encodes thousands of non-canonical open reading frames (ORFs) in addition to protein-coding genes. As a nascent field, many questions remain regarding non-canonical ORFs. How many exist? Do they encode proteins? What level of evidence is needed for their verification? Central to these debates has been the advent of ribosome profiling (Ribo-seq) as a method to discern genome-wide ribosome occupancy, and immunopectidomics as a method to detect peptides that are processed and presented by MHC molecules and not observed in traditional proteomics experiments. This article provides a synthesis of the current state of non-canonical ORF research and proposes standards for their future investigation and reporting.

Highlights

- Combined use of Ribo-seq and proteomics-based methods enables optimal confidence in detecting non-canonical ORFs and their protein products.
- Ribo-seq can provide more sensitive detection of non-canonical ORFs, but data quality and analytical pipelines will impact results.
- Non-canonical ORF catalogs are diverse and span both high-stringency and low-stringency ORF nominations.
- A framework for standardized non-canonical ORF evidence will advance the research field.

Graphical Abstract



Abstract

Ribosome profiling (Ribo-seq) has proven transformative for our understanding of the human genome and proteome by illuminating thousands of non-canonical sites of ribosome translation outside of the currently annotated coding sequences (CDSs). A conservative estimate suggests that at least 7,000 non-canonical open reading frames (ORFs) are translated, which, at first glance, has the potential to expand the number of human protein-coding sequences by 30%, from ~19,500 annotated CDSs to over 26,000. Yet, additional scrutiny of these ORFs has raised numerous questions about what fraction of them truly produce a protein product and what fraction of those can be understood as proteins according to conventional understanding of the term. Adding further complication is the fact that published estimates of non-canonical ORFs vary widely by around 30-fold, from several thousand to several hundred thousand. The summation of this research has left the genomics and proteomics communities both excited by the prospect of new coding regions in the human genome, but searching for guidance on how to proceed. Here, we discuss the current state of non-canonical ORF research, databases, and interpretation, focusing on how to assess whether a given ORF can be said to be “protein-coding”.

Introduction

Defining the extent of RNA translation in the human genome – and the resulting proteins – has long been a major focus for biomedical research. Approximately 19,500 protein-coding genes, which produce ~80,000 annotated protein coding isoforms, constitute the canonical proteome (1–6). Yet, whether this catalog is comprehensive has recently undergone substantial debate spurred by sequencing-based advances in the analysis of ribosome translation, termed ribosome profiling (Ribo-seq). Based on classical techniques used to isolate ribosome-RNA complexes, Ribo-seq is a RNA sequencing-based approach that profiles ribosome-protected RNA fragments, precisely defining open reading frames (ORFs) actively engaged by translating ribosomes (7, 8). As a tool to detect the translation of RNA, the precision of this methodology is unprecedented: from individual ribosome footprints, the exact codon being translated in a purified ribosome-RNA complex can be determined. Through the sequencing of hundreds of millions of ribosome footprints, a single Ribo-seq experiment can therefore produce a detailed and accurate representation of a given sample’s translated RNAs, typically identifying ~11,000 - 12,000 translated genes per sample (9–11), which is more similar to the ~12,000 - 13,000 expressed protein-coding mRNAs detected in a given cell type (12) compared to the ~9,000 - 11,000 proteins per sample typically detected in mass spectrometry methods (13, 14).

In addition to confirming known protein-coding sequences (CDSs), the high predictive power of Ribo-seq has unveiled thousands of other genomic sites of ribosome translation. These are most commonly found within known mRNAs (i.e., different reading frames than canonical CDS regions), but also within transcripts annotated as long noncoding RNAs, pseudogenes, or retroviral elements in the genome (7, 9, 11, 15–23). Ribo-seq can also provide clues on previously missed N-terminal in-frame extensions to known CDSs, initiated at sites alternative to the classically annotated initiation codon (24–27). The nomenclature and estimated abundance of non-canonical ORFs are listed in **Figure 1A**. For clarity, these ORFs are termed “non-canonical” to distinguish them from CDSs included in reference gene annotation - i.e. Ensembl-GENCODE - even though their translation, to our knowledge, occurs through mechanisms of ribosome activity similar to that of CDSs. Throughout this text the term “non-canonical open reading frame” is therefore defined as any open reading frame that is not an annotated CDS, an in-frame extension or truncation (either N-terminal or C-terminal), or an in-frame intron retention of an annotated CDS. For our purposes, we will be focusing on upstream ORFs (uORFs), upstream overlapping ORFs (uoORFs), internal ORFs that overlap the CDS but are translated in a different frame (intORFs), downstream overlapping ORFs (doORFs), downstream ORFs (dORFs), and lncRNA-ORFs (as in **Figure 1A**). We will not discuss in depth ORFs that may be translated from pseudogenes (19), genomic retroviruses (28), or other repetitive sequences (29) (See **Limitations** section).

Given these observations, the genomics community has been faced with the fundamental question: does the genome actually encode far more than the ~19,500 protein-coding genes currently accepted as canonical? In response, there have been increasing efforts to corroborate the observations from Ribo-seq using mass spectrometry, with the overall conclusion that only a low percentage of non-canonical ORFs are detectable by conventional tryptic whole proteome methods employing liquid chromatography with tandem mass spectrometry (LC-MS/MS) techniques (9, 15, 30–34). Yet, far more non-canonical ORFs appear to be detectable with immunopeptidomic approaches that profile peptides presented by the class

I human leukocyte antigen (HLA-I) system (**Figure 1B**) (34–39). Moreover, independent of their protein-coding capacity, non-canonical ORFs may serve important roles in the regulation of mRNA translation (40–42). With these observations at hand, one of the central tasks for the proteomics and genomics communities alike is to develop a consensus understanding on what constitutes sufficient evidence of detection for a non-canonical ORF from each technology and how to standardize these assessments given the limitations of each methodology.

Types of evidence for non-canonical ORFs

Translated non-canonical ORFs can be detected by either Ribo-seq or LC-MS/MS approaches, with examples of transition to canonical annotated protein-coding genes emerging from both. For example, translation of the signaling proteins APELA (43), POLGARF (44, 45), TINCR (46) and the cardiac proteins MYMX (47) and MRLN (48) was first identified using Ribo-seq, while shotgun LC-MS/MS data provided the initial evidence for the translation products of uORFs in ASNSD1, MKKS, MIEF1, and SLC35A4 (30, 49).

Together, the combination of Ribo-seq and LC-MS/MS is a powerful way to identify translated CDSs and ORFs (21, 50–52). Ribo-seq does not directly detect proteins, but rather provides evidence of ongoing nucleotide translation. By contrast, LC-MS/MS evidence for non-canonical ORFs takes the form of direct detection of peptides. In the case of conventional LC-MS/MS of cellular lysates, these peptides are typically tryptic, meaning they were generated by protein cleavage at the C-terminal side of a lysine or arginine, or semi-tryptic, meaning they were generated by protein cleavage at the C-terminal side of a lysine or arginine at one end of the peptide but not the other. However, many ORFs have now been observed in mass spectrometry-based HLA-I immunopeptidomics data (18, 34, 36, 38, 53). Here, no tryptic digestion is employed. Instead, peptides containing the HLA-I peptide binding motifs of the HLA-I allele expressed by a specific cell line or tissue are observed. A variety of lower-throughput approaches have also been used to assess translation of non-canonical ORFs, including generation of custom antibodies, expression of epitope-tagged ORF cDNAs, selective reaction monitoring (SRM), and radiolabeled *in vitro* translation (9, 17, 54–56).

While high-quality Ribo-seq and LC-MS/MS whole proteome data on the same sample should be able to identify highly consistent sets of endogenous CDSs, Ribo-seq is not able to pinpoint the responsible translation event for exogenous proteins, which originate from sources other than the sample's own genetic material. Similarly, Ribo-seq cannot detect or predict protein stability, folding, or post-translational modification. If there is a substantial discrepancy with MS detecting many additional proteins, then the quality of the Ribo-seq library should be inspected (see below). It should also be noted that Ribo-seq, like all sequencing-based methods, may not be able to resolve translation events in repetitive genomic regions, such as retrotransposons, pseudogenes, or genes with very high homology.

By contrast, Ribo-seq will almost always detect many non-canonical ORFs that are not found by proteomics. This is due to several factors: both the nature of the data itself as well as technological differences in the methods that may impact the ability to detect lowly-expressed molecules with high confidence. For example, all mass spectrometry-based proteomics methods lack a PCR amplification step that is present in most sequencing-based methods, which enables higher sensitivity at lower sample

inputs. Regarding the nature of the data, Ribo-seq has the ability to identify translating ribosome signatures in an unbiased way, which may confidently find ORFs less than 8 amino acids long that are fundamentally challenging to identify by MS (15, 57). In fact, Ribo-seq can confidently identify an ORF that is simply a start codon followed by a stop codon (i.e. Met*), because the Ribo-seq reads remain sufficiently long for unique genomic mapping (58).

Second, since some non-canonical ORFs are located in GC-rich promoters (such as uORFs), these may encode amino acid sequences that are enriched in arginine (CGU/CGC/CGA/CGG codons) and thus would be excessively cleaved by trypsin to small peptides that cannot be uniquely mapped to a single ORF. Whether use of alternative proteases (59) could improve non-canonical ORF detection in whole lysate proteomics is unclear.

Considerations and quality control steps for the data-driven discovery of non-canonical human ORFs

Differences in the nature of Ribo-seq and LC-MS/MS-based whole proteome and immunopeptidome data collection also represent a source of substantial variability in the detection of non-canonical ORFs. Notably, while targeted whole proteome and immunopeptidome LC-MS/MS approaches may offer improved sensitivity, these require candidate non-canonical proteins of interest to be known prior to analysis. While each method uses high-throughput data generation to profile cellular translation comprehensively, the data have intrinsically different strengths and weaknesses that may result in discordance between them (**Table 1**).

Ribo-seq

The quality of a Ribo-seq dataset is most commonly evaluated using three considerations: codon periodicity, library complexity, and number of canonical CDSs identified.

Codon periodicity reflects the percentage of Ribo-seq reads that correctly identify the known reading frame of CDSs (**Figure 2A-C**). In a high-quality Ribo-seq dataset, $\geq 70\%$ of reads that are between 28 and 30 nucleotides in length map to the correct reading frame of known CDSs. The precise read length that displays the most preferable (the “cleanest”) signal can vary and depends on the sample type and the method of nuclease digestion used to eliminate cellular RNAs not bound within the translating ribosome. Because of limitations of the experimental technique as well as biological variation in ribosome occupancy, a codon periodicity above 90% is typically not attainable (60). A Ribo-seq dataset with a codon periodicity $< 60\%$ should ideally not be used for ORF discovery due to challenges with accurate identification of the reading frame (19, 60, 61). A periodicity between 60 – 70% is a gray zone where the data may be used in some cases with increased caution and stringency.

Library complexity refers to the number of unique RNA molecules sequenced and what fraction of these are ribosome footprints that map to CDSs. The challenge with a low complexity library is that the majority of the reads will be PCR duplicates. When the number of initially isolated footprints is limited (e.g. due to low quality of the input material or sub-optimal sample processing), ultimately many duplicate copies of this limited number of footprints will be sequenced. This means that deeper sequencing of this library will yield no or only minimally more biologically distinct footprints. Typically, the majority of reads in such low

complexity libraries will come from non-footprint sources, particularly intergenic and intronic contaminants (e.g. microsatellite repeat elements, ribosomal RNAs, or small RNAs that overlap gene regions), which are unintentionally isolated during the Ribo-seq procedure because these RNA species are of a similar size to the ribosomal footprint and may have certain RNA structures (62, 63). In general, a Ribo-seq library with sufficient complexity will have the majority of reads mapping to annotated and novel CDSs. In some cases, such as with degraded samples, there may be substantial intergenic noise or a higher fraction of RNA species that are normally restricted to the cell nucleus, but yet still sufficient codon periodicity and library complexity in terms of unique RNA molecules that map to CDSs. Here, the challenge is to achieve sufficient sequencing depth to ensure adequate sampling of unique RNA molecules. While 150 million reads typically suffices for the analysis of a high-quality Ribo-seq library, a “noisy” – yet usable – library may require very deep coverage (>400 million reads), which is mostly a consideration for the financial cost of the sequencing (60, 64, 65). For human Ribo-seq libraries, typically 15-30% of the sequenced reads can be classified as ribosome footprints and the rest is often discarded. For a library sequenced to a depth of 150 million reads, that would total to approximately 22.5 to 45 million ribosome footprints – a number comparable to a routinely sequenced RNA-seq library. Of these, >80% should map to annotated CDSs (60), leaving ~5 million ribosome footprints for ORF discovery.

The number of known CDSs identified is particularly important when one aims to provide a comprehensive view of all translated ORFs in a sample of interest. This metric relates both to the amount of noise in the library, the periodicity of the footprints, as well as the depth of the sequencing. A sufficiently sequenced Ribo-seq library for a human sample with high periodicity should detect at least >9,000 annotated CDSs, and often >10,000 (9–11, 18). Human sample Ribo-seq libraries that do not reach this threshold – despite sufficiently deep sequencing and periodicity – should be used with caution, as the false negative rate for detecting ORFs will be high (many ORFs will be missed). While Ribo-seq-based ORF detection tools theoretically have a low false-negative rate, the confidence (FDR) with which an ORF or CDS is detected, the number of independent samples in which it can be found, and the translation rate of the ORF should always inform research decision-making. For instance, direct comparison of non-canonical ORF FDRs and translation rates, compared to those of canonical CDSs, can inform both the relative abundance of the ORF’s translation product and the degree of certainty with which the algorithm could nominate it.

Because *de novo* and *ab initio* RNA assemblies are technically challenging with the short nucleotide sequences (28 – 30 nt) obtained during a Ribo-seq experiment, analysis of Ribo-seq data requires alignment of the reads to a reference transcriptome, most commonly Ensembl or RefSeq though custom transcriptomes are also used in some cases. Statistical assessment of a non-canonical ORF nomination is inconsistent across computational methods, with some approaches calculating a p-value for significance (e.g., Ribotaper (61), ORFquant (10), Ribo-TISH (66), PRICE (67), RiboCode (68)) and other approaches computing confidence scores (e.g., RibORF (19), Ribotricer (69), ORF-RATER (70)). In addition, these methods are often based on fundamentally different modeling approaches, including Hidden Markov (RiboHMM (20)), multitaper (Ribotaper (61)), transformer (DeebRibo, TIS Transformer (71, 72)), support vector machine (RibORF (19)), expectation-maximization (EM) (PRICE (19, 67)) models, among others. As such, different methods may be more appropriate for certain research questions, datasets or desired ORF types.

As a consequence, two different algorithms can have differing ORF outputs for the same gene. This can be due to the level of stringency or the strengths and weaknesses of a particular ORF caller for a certain type of ORF or certain quality of data. For example, some ORF callers cannot detect ORFs with near cognate start codons, whereas others are better suited for the detection of overlapping reading frames where periodic footprint signals are mixed and hard to dissect. Other tools handle alternative splicing better. Depending on the research question, input data quality, species of interest, or annotation goals, combinations of ORF callers followed by curation of called ORFs may be necessary (see below in **How many non-canonical ORFs are there?**).

HLA-I and HLA-II Immunopeptidomics

In the past decade, interest in HLA-I and HLA-II presented peptides has become widespread across many areas of biomedical research, as a subset of HLA-presented peptides demonstrate antigenic properties and represent a class of potential therapeutic targets (73–76). The application of HLA immunopeptidomics differs from tryptic whole proteome protocols, as these methods leverage native lysis buffer and antibody or affinity-tag enrichment steps to isolate HLA-peptide complexes from cell lysates (**Figure 1B**) (77, 78). The peptides are naturally produced following degradation of endogenously expressed source proteins by cellular proteases and peptidases and the proteasome. As such, no tryptic digestion is used in immunopeptidome analyses. Therefore, regarding detection of non-canonical proteins, HLA immunopeptidome analysis has three advantages over tryptic whole proteome analysis: 1) each HLA allele has a distinct peptide-binding motif that presents specific subsets of peptides, which can then be detected with mass spectrometry in the absence of digestion with a protease; 2) the HLA presentation pathway may have privileged access to proteins that are rapidly degraded as the half-life of HLA-peptide complexes (hours) are in general longer than the half-life of rapidly degraded proteins (minutes); and 3) HLA immunopeptidomics broadly samples endogenous proteins from all abundance levels including those from lower-abundance non-canonical ORFs (79–81). These advantages align with recent studies that have shown higher observation rates of non-canonical proteins in the HLA-I immunopeptidome compared to the tryptic whole proteome (39, 82).

Similar to tryptic whole proteome datasets, immunopeptidome datasets require strict quality control steps to ensure the data and analysis are of high-quality. Peptide length, the presence of peptide-binding motifs, and predicted binding to HLA molecules coded by specific alleles are common quality control steps in immunopeptidomics workflows. Because HLA-I and HLA-II molecules have unique peptide binding grooves that accommodate peptides of different lengths, peptide size is an important quality control metric of immunopeptidomics data. Specifically, HLA-I peptides are ~8-12 amino acids long (mostly 9mers), whereas HLA-II peptides are generally 12-25mers (77). HLA-II peptides are also typically found in nested sets, while this is not a global feature of HLA-I peptides, and can also be used to quality control HLA-II immunopeptidome datasets. Furthermore, each individual person expresses different HLA alleles with distinct HLA-binding motifs, which influence which peptides are presented. Therefore, it is common to confirm that HLA allele-specific binding motifs of the expressed HLA molecules are present in the immunopeptidome data, and that peptides derived from canonical and non-canonical ORFs in a given dataset are predicted to bind to the expressed HLA molecules to a similar extent. A number of

computational approaches (e.g. MHCflurry, NetMHCpan, MixMHCpred, ForestMHC, HLAthena) can be used to both predict HLA peptides and the strength of their binding to various HLA molecules (76, 83–88). It is important to note that HLA-I binding prediction is currently more accurate compared to HLA-II binding prediction, as HLA-II motifs are more complex and large subsets of diverse HLA-II heterodimers are in the process of being characterized and the associated prediction algorithms are being further improved (89–92).

Interestingly, peptides derived from non-canonical ORFs are much more abundant in HLA-I datasets compared to HLA-II datasets (18, 34, 36, 38, 39, 53, 93). HLA-I molecules usually present peptides derived from proteasome-mediated degradation of newly synthesized and other cellular proteins, and HLA-I presentation is tightly linked with protein synthesis and degradation rates. In contrast, HLA-II molecules, that are often expressed on professional antigen presenting cells (APCs), present peptides derived from degradation of extracellular proteins that were taken up by the APCs or from endogenous proteins that are destined to be degraded in specialized vacuolar compartments of the endosome-lysosome system. Both HLA-I and HLA-II systems require trafficking to ensure peptide loading in the right compartment. For HLA-I, the peptides themselves are transported into the ER by a transporter associated with antigen processing (TAP), while in case of HLA-II, the source proteins must first reach the acidic compartments for degradation, for example via receptor mediated internalization or recycling of transmembrane proteins. Hence, the sources of HLA-II presented peptides are often stable and abundant proteins.

Because of HLA-I binding constraints, and the short length of some non-canonical proteins, a non-canonical ORF is often represented by a single peptide in HLA-I immunopeptidome data, and therefore additional quality control measures should be taken to support these identifications. To this end, a non-canonical protein subset-specific FDR threshold should be applied to each individual ORF type, rather than a global FDR (82, 94) because non-canonical ORF peptides represent a small fraction (typically <5%) of the overall immunopeptidome and individual ORF types vary considerably in their frequency. Thus, a global FDR can be excessively permissive for a small subpopulation and lead to higher false-positive identifications.

Beyond leveraging known HLA specific peptide lengths, binding motifs and subset-specific FDR, there are further quality metrics that can be applied to immunopeptidomics datasets when the focus is the identification of rare, non-canonical proteins (95). The gold standard for supporting the identification of non-canonical peptides presented by HLA molecules is by comparing the retention time and MS/MS spectrum of an identified peptide with a synthetic peptide of the same amino acid sequence. However, it is often the case that hundreds of non-canonical peptides are identified in a single HLA-I immunopeptidome experiment, making the synthetic peptide confirmation for all potential non-canonical derived HLA-I peptides not feasible. To overcome this challenge, it is now possible to compare the observed MS/MS spectra with predicted MS/MS spectra with tools such as Prosit (96). The comparison of the predicted and observed MS/MS spectra provides additional support for non-canonical peptide identification (97, 98). In addition, there are also multiple algorithms that can predict peptide retention times. The predicted retention time, using tools such as DeepLC or DeepRescore, can be compared to measured retention time for all peptides in a sample (canonical and non-canonical), as the correlation

between predicted and observed retention time supports the LC-MS/MS identifications of non-canonical derived peptides in immunopeptidomes (99, 100). Overall, deep learning based prediction of peptide MS/MS spectra and retention time are powerful tools that help reduce the number of false positive non-canonical peptide identifications in immunopeptidome datasets.

Tryptic whole proteome LC-MS/MS

Rigorous standards for the analysis of LC-MS/MS tryptic whole proteome data have been established by the Human Proteome Organization/Human Proteome Project (HUPO/HPP) international consortium, as reviewed elsewhere (101–103), and these standards remain the expectation for researchers claiming identification of non-canonical ORF peptides (30). For claims of detection of proteins not previously detected, these guidelines require two non-nested, uniquely mapping peptides each of at least nine residues in length with a total extent of at least 18 amino acids and with high-quality peptide-spectrum matches (PSMs) upon manual inspection (30, 101, 103). These peptide-spectrum matches should be provided in the form of Universal Spectrum Identifiers (USIs) so that the spectra can be easily examined by others (104).

Yet, consistent application of high-quality tryptic whole proteome data collection and analysis guidelines remains non-uniform across the research community. Proteogenomic studies looking for non-canonical ORFs without Ribo-seq data – i.e., by predicting and including all ORFs in RNA transcripts – have been plagued by high false-positive rates (30, 49, 105–108), and initial efforts to inspect early claims of non-canonical ORF peptides concluded that “many of the spectral matches appear suspect” (30).

Moreover, while use of decoys is standard in tryptic whole proteome experiments to define global false discovery rates, decoys may be less useful for distinguishing true peptides for non-canonical ORFs. Indeed, Wacholder *et al.* have concluded that decoy bias among non-canonical ORF products leads to inaccurate FDR estimates for short ORFs when decoys are created by reversing the complete protein sequence, but not when excluding the initial Met from the reversal (109). Lastly, efforts to identify non-canonical ORFs in tryptic whole proteome data must account for peptides instead being derived from canonical variants including single amino acid variants (SAAVs) and splice-site peptides for alternative isoforms of known CDSs. The use of personalized proteogenomic database searches is not straightforward or used by all in the proteomics community.

Considering these factors, the general experience of the research community is that few non-canonical ORFs are found by conventional tryptic whole proteome LC-MS/MS analyses and some of those are ultimately false-positive peptides (110, 111). In some cases, such ORFs are “undiscoverable” by tryptic whole proteome approaches, either due to the short length of non-canonical ORFs or intrinsic sequence features that do not produce LC-MS/MS observable tryptic peptides. For example, translation of repetitive amino acid sequences (e.g. glycine-leucine) has recently been described (29). Nevertheless, even approaches aimed at enriching for small proteins from cell lysates result in only modest increases in non-canonical ORF detection, rather than exponential increases (33). On the other hand, other enrichment techniques focused on post-translational modifications (PTMs; i.e. the acetylome, phosphoproteome, and ubiquitylome) have also reported non-canonical proteins and may provide both

an alternative method to enrich for non-canonical proteins and also hint towards potential functional relevance of this subset of non-canonical proteins given the cellular roles of those PTMs (82).

Furthermore, data-independent acquisition mass spectrometry (DIA-MS) provides a potential opportunity to detect non-canonical ORF-derived peptides that have been reliably detected previously with high quality spectra obtained with narrow isolation windows from a data-dependent acquisition (DDA) approach. In DIA-MS, previously identified peptides are more reproducibly sampled by sequentially isolating and fragmenting peptides across the m/z range, which decreases stochastic sampling bias toward higher abundant species and may increase the chances of finding rare non-canonical ORFs (112). This approach has been used in conjunction with Ribo-seq to claim detection of microproteins from non-canonical ORFs (50).

Beyond technical limitations of mass spectrometry, there are also biological factors that may make non-canonical ORFs less frequently observed in tryptic whole proteome LC-MS/MS datasets. To this end, there is increasing evidence that points toward intrinsic instability of proteins translated from non-canonical ORFs, resulting in their immediate degradation. Kesner *et al.* used functional genomics approaches to demonstrate that the ribosome-associated BAG6 membrane protein may directly triage hydrophobic non-canonical ORF translations to the proteasome for degradation (113). Thus, it is possible that many non-canonical ORFs do not generate a stable protein product and might only be observable by immunopeptidomics or in whole proteome experiments with inhibition of the protein degradation mechanisms of a cell.

How many non-canonical human ORFs are there?

The number of non-canonical ORFs encoded in the human genome remains highly speculative. To date, a limited number of human tissues and cell lines have been analyzed by Ribo-seq, and proteogenomics studies that have aimed to incorporate ORFs derived from these datasets have been difficult to interpret due to numerous false positives. As such, while it is well-established that the human genome contains thousands of translated non-canonical ORFs, whether the precise number is closer to 10,000 or 100,000 remains a matter of debate. A further complication is that different research communities may not use a consistent definition of what types of ORFs we define as “non-canonical”. Yet, while analyses of more cell lines and tissues will certainly uncover additional non-canonical ORFs, there can be variable non-canonical ORF identifications even within analyses of the same cell line. Such variability reflects the equal – perhaps foremost – contribution of different analytical methods for non-canonical ORFs in the estimation of their prevalence.

The number of non-canonical ORFs

Most Ribo-seq studies focusing on non-canonical ORFs report detection of several thousand ORFs, typically between 2,000 and 8,000 (9, 11, 15, 16, 18–21, 51, 61, 114). Interestingly, this range seems relatively stable when comparing studies that employ only a few cell lines and broader analyses looking across many different human tissue types. To consolidate these findings, we have recently participated in an international consortium to aggregate 7,264 high-confidence non-canonical ORFs and provided formalized annotations for them within the GENCODE gene annotation database (16). This GENCODE set

demonstrates substantial overlap in the identification of certain types of ORFs, such as upstream ORFs (uORFs), across diverse datasets such as pancreatic progenitors, heart and stem cells, suggesting that perhaps the diversity of several ORF types may not be dramatically larger with the inclusion of more tissue types. In support of this, Ribo-seq profiling of five human tissue types and 6 primary human cell types similarly reported 7,767 ORFs in total (15). When subsetting this dataset for consistency with the inclusion criteria for the GENCODE catalog (i.e., removing ORFs below 16 amino acids in size, as well as ORFs without an AUG start codon), 2,475 out of 7,767 ORFs remained, of which 1,702 ($\pm 70\%$) were represented in the GENCODE catalog as well (**Supplementary Tables 1-4**).

While these studies have measured and determined non-canonical ORF translation directly from Ribo-seq data, there are many other databases that have aggregated larger numbers of ORFs from a variety of sources, including both Ribo-seq and *in silico* predictions. Among these, smProt (n = 327,995 human ORFs (115)), sORFs.org (n = 4,377,422 ORFs across humans, mouse and fruit flies (116)), RPFdb (117, 118) and smORFunction (n = 617,462 human ORFs (119)) have compiled reported or putative non-canonical ORFs. Notably, OpenProt (120, 121) has two aspects to their database workflow: one that collates all predicted ORFs (n = 488,956) and a second that proposes 33,836 translated ORFs identified by a re-analysis of over a hundred Ribo-seq datasets with the PRICE pipeline (67). When considering studies that have generated Ribo-seq datasets to measure non-canonical ORF translation, there are also several efforts that have proffered exceptionally large numbers of directly detected ORFs – specifically the nuORFdb (34) by Ouspenskaia *et al.* and the Human Brain Translatome Database (122) by Duffy *et al.*, which propose numbers of >230,000 and >75,000 ORFs, respectively.

Why is there such discordance in the number of non-canonical ORFs across databases?

The interpretation of such dramatically different accounts of non-canonical ORF abundance remains a challenge. Indeed, given that there are currently only ~60,000 Ensembl genes (including 19,827 protein-coding genes, 18,886 lncRNAs, 4,864 small ncRNAs, 15,241 pseudogenes, and 2,221 other RNAs in Ensembl v109.38), colossal datasets with >200,000 ORFs may be interpreted to suggest that every gene has upwards of 4 distinct ORFs. In practice, these large datasets may include isoform variants (e.g. N-terminal extensions, C-terminal extensions and intron retentions) that are not part of the reference proteome, and thus the number of non-canonical ORFs may be larger in some databases due to differences in how these isoforms are categorized.

While sample and data quality likely contribute to the variability in the numbers of non-canonical ORFs in some catalogs, differences in Ribo-seq data analysis also account for much variation in prospective non-canonical ORFs. For example: biologically, there is some amount of stochastic or pervasive translation across all RNAs, which may relate to leaky ribosomal scanning (123–125) or transient interactions between ribosomes and RNAs as the ribosomes locate coding sequences or RNAs accomplish proper folding (126, 127). Yet, the manner in which computational pipelines process Ribo-seq data results in ORF calls that may be more or less stringent (**Figure 2B**), resulting in different proportions of false-positive (stochastic) and false-negative (e.g., sample-specific) ORF calls (60, 128, 129). For example, RibORF (19), which uses a support vector machine and recommends a fixed cut-off score of 0.7, has been shown to produce the highest numbers of ORF calls of any tested algorithm in a recent benchmarking study (130). To confirm

these differences directly, we have re-analyzed published high-quality Ribo-seq data for six biological replicates of pancreatic progenitor cells differentiated from human embryonic stem cells (11) using four common ORF detection pipelines (ORFquant (10), PRICE (67), Ribo-TISH (66), and Ribotricer (69)), observing substantial variability in the number of ORFs called (~10-fold difference from ~50,000 to ~500,000), the types of ORFs called, the length of the called ORFs, and the reproducibility with which ORFs could be detected across all six replicates (**Figure 3** and **Methods**).

There may be specific reasons for the different performance characteristics of each algorithm. For example, the lower stringency of RibORF may be due to the fact that this pipeline considers uniformity of read coverage across the ORF, while Ribo-seq is known to have a 5' bias to read coverage. Therefore, RibORF may excessively promote intORFs and doORFs since the 5' ends of these ORFs overlap annotated CDSs, which typically have higher read coverage independent of a periodic footprint signal that matches the correct reading frame. This is evident in nuORFdb (34) and the Human Brain Translatome Database (122): when analyzing the fraction of ORFs with an AUG-start resulting in an ORF ≥ 16 amino acids, doORFs and intORFs are 173-fold and 18-fold (respectively) higher in abundance compared to other major datasets (**Figure 4**, **Supplementary Tables 5-8**). By contrast, uORFs are only 3 times more abundant (**Figure 4**).

It is also true that different computational pipelines may have different capacity to identify certain classes of non-canonical ORFs. For example, the deterministic multitaper-based statistical inference of significant periodic signal within predicted ORFs as performed by Ribotaper (61) and ORFquant (10) provides high-confidence detection of ORFs with an AUG start codon, but have not, to date, been optimized for non-AUG ORFs. In contrast, the probabilistic algorithm employed by PRICE (67) has enhanced ability to identify very short ORFs and non-AUG ORFs absent from other ORF callers (**Figure 3B+E**). Yet, when there are neighboring putative initiation codons (e.g. CUG and AUG), PRICE will generate larger numbers of putative ORFs that might require manual curation or further filtering. In addition, since annotated CDSs have generally more abundant Ribo-seq read coverage, low-abundance out-of-frame reads may be more readily interpreted as an intORF with a non-AUG start codon by PRICE, whereas other ORF callers are less likely to consider these reads as sufficient evidence for a translated ORF. Thus, when applied to biological replicates of the same sample, PRICE produces the least consistent ORF calls compared to other pipelines, independent of initiation codon variability (**Figure 3A-C**) (130). nuORFdb (34) and OpenProt (121) both employ PRICE in their analysis pipelines. It is important to note, however, that the specific research question being pursued should inform the types of ORF callers used: indeed, deterministic algorithms such as RiboTaper or ORFquant may miss intORFs or overlapping ORFs identified by PRICE because of the difficulty in resolving mixed periodicity signals of overlapping reading frames (**Figure 3A**).

In summary, depending on the type of ORF one aims to find and the desired inclusiveness of ORFs one aims to output, one ORF caller might be better suited than another. Certain ORF callers outperform others in detecting specific ORF categories such as intORFs (**Figure 3A**), very small ORFs (**Figure 3B+D**), or near cognate start codons (**Figure 3E**), whereas others handle exon-exon junctions and longer ORFs better and/or provide better replicate behavior. These differences then lend to substantially different results when producing non-canonical ORF catalogs (**Figure 4**).

Detection of translational start sites

Determining the translational start site of an open reading frame remains a nuanced problem. While conventionally proteins have been annotated with AUG start sites, exceptions to this rule have long been known (131, 132), and non-canonical ORFs are more likely to employ non-AUG start sites (124, 133). In a typical Ribo-seq experiment, identification of translational start sites from Ribo-seq data is inferred based on two factors: sequencing coverage and the intrinsic restrictions of the computational pipeline (e.g. some algorithms only consider AUG start codons, as discussed above). Yet, independent of the computational pipeline, there may be gaps in the sequencing coverage that lead to mis-identification of the main translational initiation site (**Figure 2C**). For experiments with cultured cells, use of small molecules that block ribosome elongation, such as homoharringtonine (134) or lactimidomycin (135), enables ribosome accumulation on translational initiation sites, which enables more precise determination of the start codon. Due to the difficulty in identifying non-canonical ORF start sites and the variability in computational approaches to start codon recognition (e.g., **Figure 3E**), use of homoharringtonine or lactimidomycin with cultured cells is highly recommended. In frozen tissue samples, these compounds are no longer effective.

How to select an ORF sequence database for MS data analysis

Given the wide differences between the different databases for Ribo-seq ORFs, one central question is how to use these databases, or which to use for any specific analysis? Because the size of the ORF output in a given database can vary enormously, users should base their decision on what scientific question they intend to pursue and evaluate carefully the suitability of the input Ribo-seq data quality as well as the stringency with which ORF calling was performed. In general, high stringency databases provide high confidence Ribo-seq ORF detections, and thus peptides found mapping to these ORFs are more likely to reflect a true positive result. While these databases reduce false positives, it is at the expense of comprehensiveness, as the existing high stringency databases will yield more false negatives in the MS analysis. Low stringency databases provide a much larger set of Ribo-seq ORFs, but will yield more false positives – due to the lack of support from another orthogonal technique. If the ORFs are accompanied by Ribo-seq quality metrics, it may be tractable to estimate the proportion of false positives, and re-filter the ORFs to suit one's own purposes. These databases will provide a larger candidate search space for peptide alignment and may enable detection of true positive ORFs not present in the high stringency databases. Yet as described earlier, due to the concern for false positive nominations, ORFs detected by mass spectrometry searches should be closely inspected to verify integrity of both ORF call and peptide identification, as there will likely be cases of false-positive ORFs being supported by false-positive peptides. Ultimately, certain scientific questions may lend themselves to certain databases: for example, analyses of alternative N-terminal CDS extensions often emphasize non-AUG start sites (24), which may benefit from a Ribo-seq analysis that employs the PRICE algorithm. Research efforts aimed to identify a maximal space of potential translation events may also favor a lower stringency database, with the caveat that any individual result should receive additional scrutiny. Alternatively, if the goal is to characterize a high-confidence unannotated microprotein, a high stringency database may be more desirable. Likewise, for reference annotation purposes and functional studies we prefer more stringent workflows that yield reproducible ORF calls across samples (no false positives).

Are non-canonical ORFs proteins?

The term “protein” is conventionally used to refer to an amino acid sequence that produces a molecular structure that plays an intrinsic cellular role in maintaining normal cell biology. While some proteins may be unstable and rapidly degraded under certain conditions (e.g., beta-catenin), most proteins participate in cell biology when present in a stable form. Also, almost all annotated proteins show evidence of evolutionary conservation, structural folding and domain architecture, and frequently also protein-protein interactions and/or interactions with nucleic acids.

According to this understanding of the term “protein”, it could be inferred that the vast majority of non-canonical ORFs do not encode proteins on the basis that they lack these characteristics. However, we see two additional considerations. First, it may be incorrect to assume that a protein that exists in the cell – even one that is detectable by mass spectrometry – is therefore a functional molecule. It could be that the proteome contains a certain amount of non-functional translational “noise”. While it is difficult to prove the extent to which such translation occurs in normal cells, evidence from cancer cells shows abundant dysregulation of translation, exemplified by “aberrant” non-canonical proteins that lack evidence for function under normal physiological conditions (34, 35) as well as out-of-frame peptide byproducts of oncogene activity (136).

Second, the classical definition of protein “function” invokes the protein’s role in cellular processes that have been derived over time through evolution, which has been summarized as the maxim that “conservation = function”. This maxim has been central – but not universally required – for gene annotation projects, and the only canonical proteins currently within GENCODE that can be inferred to have evolved *de novo* in human or higher primates were initially detected in cancer cells (e.g. MYEOV (137) and HMHB1 (138)). Even so, evidence for the existence and function of *de novo* proteins under normal physiological conditions is accumulating (57, 139–141). Nonetheless, it remains true that most non-canonical ORFs display much higher rates of intrinsic disorder, fewer structural features, and lack amino acid constraint across evolution (17, 18, 139, 140, 142–148). While these features may be observed in diverse annotated proteins (e.g. intrinsically disordered regions of a given protein), their presence is predominant in non-canonical ORFs.

The interpretation of peptide-level evidence of Ribo-seq ORFs

How, then, should one interpret the peptide-level evidence for some non-canonical ORFs? High-quality whole proteome LC-MS/MS PSMs that survive rigorous manual inspection are strong evidence of true translation of a non-canonical ORF. With adequate evidence, therefore, whole proteome PSMs supporting non-canonical ORFs do indicate the possible existence of a translated protein, and these cases may reasonably be considered to be part of the cell proteome, similar to any other proteins.

When considering the larger number of non-canonical ORFs with peptide-level evidence in HLA immunopeptidomics but not shotgun LC-MS/MS (18, 34, 36, 38, 149), firm conclusions are more difficult to draw. These non-canonical ORFs cannot be said to generate a true protein based on immunopeptidomics alone, considering that the HLA system is expected to present peptides resulting

from translation products that are unstable and rapidly degraded, alongside those derived from canonical proteins. Yet, detection of an HLA-presented peptide does verify RNA translation in these cases, which distinguishes them from the majority of Ribo-seq-detected non-canonical ORFs that are detected in neither shotgun LC-MS/MS nor immunopeptidomics experiments. Therefore, these non-canonical ORFs can at least be said to be confirmed as both translated and presented by the HLA, as opposed to an artifact of the Ribo-seq protocol.

A related question is how to interpret PSMs matching non-canonical ORFs that are not detected by Ribo-seq, when the same sample is interrogated using both technologies. Because the sensitivity of Ribo-seq is generally higher than MS-based methods, and because Ribo-seq provides nucleotide-level precision for genome mapping, there are three possibilities here: first, these peptides may be false-positive identifications, second, the Ribo-seq data exhibit a false-negative, or third, they may be derived from another source not included in the search space (ex: aberrant splicing). None of these hypotheses has been rigorously evaluated at this time. One challenge is that many proteomics and immunopeptidomics experiments do not currently generate matched Ribo-seq data for their samples, and thus it cannot be directly known if Ribo-seq supports translation of that ORF. When considering unmatched analyses, it is also noted that, at present, proteomics and immunopeptidomics datasets cover a broader range of tissue and cell types than Ribo-seq datasets.

A proposed framework to classify the translation of non-canonical ORFs

Given the expanding volume of research on non-canonical ORFs, a shared vocabulary for the interpretation of their detection is a critical need in the genomics, translomics, proteomics, and immunopeptidomics communities. Notably, there has been no formalized initiative to annotate non-canonical ORFs as protein-coding genes by major genome databases, although recent collaborative work has raised this point as a topic of interest (16). Historically, protein-coding genes have been annotated one-by-one in a manual process of careful data inspection, which may or may not have included protein-level evidence. At this time, non-canonical ORFs detected by tryptic whole proteome data would potentially be eligible for manual annotation as protein-coding genes. Yet, given the paucity of non-canonical ORFs in tryptic whole proteome data and their much greater abundance in HLA immunopeptidomic datasets, there is uncertainty about whether most non-canonical ORFs produce proteins in the classical sense, and whether immunopeptidomic evidence is equivalent to tryptic whole proteome data for the purposes of protein annotation.

We advocate both a cautious but open-minded approach to non-canonical ORF classification, summarized in **Table 2**. Notably, although most annotated proteins show evidence of amino acid constraint across species and most non-canonical ORFs do not, it is also unquestionably true that at least some proteins are lineage- or species-specific. Thus, we propose that *de novo* translations should be considered for annotation as protein-coding. However, recognizing that evolutionary analysis is a core part of gene annotation workflows in projects like GENCODE, we have not included conservation or constraint metrics as part of this proposed framework. The framework itself is oriented toward harmonizing subsequent dataset generation and analysis. In practice it might be applied to classifying published datasets, and it is intended as a helpful tool for candidate prioritization rather than a guarantee that certain ORFs will be

annotated by a genome database. We stress that researchers looking to move forward with potential annotation of a protein encoded by a non-canonical ORF should be able to provide the raw LC-MS/MS spectra for review.

Our framework proposes these definitions for specific terminology:

- *“Protein candidate”*: a **Tier 1A** non-canonical ORF can be regarded as translated into a protein candidate if it satisfies current HUPO HPP guidelines for the detection of ≥ 2 uniquely mapping peptides, as well as having evidence of translation by Ribo-seq. Such candidates would be prioritized for further manual review by annotation groups.
- *“Presented”*: A presented non-canonical ORF (**Tier 1B**) is one with multiple lines of evidence for its translation and presentation on HLA molecules. These ORFs are detected with multiple high confidence peptides from multiple distinct samples for HLA immunopeptidomics data, as well as having evidence of translation by Ribo-seq.
- *“Detected”*: A detected non-canonical ORF is one with evidence of translation by Ribo-seq as well as evidence of protein production by either (**Tier 2A**) shotgun LC-MS/MS (1 peptide or >1 peptide not satisfying HUPO HPP guidelines for their spacing), or evidence of protein production by HLA immunopeptidomics with a single PSMs (**Tier 2B**).
- *“Putative”*: A putative non-canonical ORF (**Tier 3**) is one with evidence of translation with LC-MS/MS or HLA immunopeptidomics data, but no evidence of translation in Ribo-seq data. This discrepancy may alert to the possibility of false-positive MS identifications, or false-negative absence in Ribo-seq, and therefore requires more investigation.
- *“Ribo-seq ORF”*: A non-canonical ORF that is only detected in Ribo-seq data but not elsewhere is considered a “Ribo-seq ORF” (**Tier 4**). These are likely to be the majority of cases. The number of these ORF nominations may be variable based on the stringency of the Ribo-seq analysis and/or the quality of the input data.
- *“Predicted”*: A predicted non-canonical ORF (**Tier 5**) is one that is computationally predicted *in silico* on an expressed RNA transcript, but without current evidence in Ribo-seq or MS datasets.

Limitations

With this work, we have endeavored to clarify how Ribo-seq can be used for non-canonical ORF research. Yet, our focus has several important limitations. First, the vast majority of – but not all – translated peptides can be traced back to an RNA sequence. There may be peptides that derive from amino acid splicing within the proteasome during protein degradation (150), which would not be detectable in Ribo-seq data. Second, there are also well-established protein-coding sequences that are difficult to resolve with Ribo-seq and do not have optimized computational methods for their quantification. For example, translated pseudogenes, retroviruses, retrotransposons, and paralogous protein-coding genes may have high sequence homology that precludes unique mapping of the short ~30 bp reads from a Ribo-seq experiment, although multi-mapping reads will provide evidence of translation. These cases are not discussed here. Lastly, each individual’s genome (and particularly each cancer’s genome) has a unique range of germline or somatic single nucleotide variants that will impact the proteome: in this article, we have not addressed the importance of generating personalized reference genomes and proteomes for the analysis of microproteins and non-canonical ORFs.

Conclusions

The widespread description of non-canonical ORFs has sparked a paradigm shift in the perception of both the human genome and the proteome. Yet, as a field still in its infancy, this area of investigation is plagued by a lack of standardization, which may lead to imprecise analyses, ultimately leading to self-injurious confusion. While the proportion of non-canonical ORFs that encode a functional protein remains to be seen, a large fraction of them can be verified as translated by both MS-based and Ribo-seq-based approaches. A central effort for the research community is now to build reputable databases and analysis pipelines to ensure rigor in this quickly-expanding – and highly exciting – field. Here, we have considered the technologies used to detect non-canonical ORFs and attempted to provide a framework for categorizing differing levels of evidence for them. Our work aims to coalesce the research community around a common terminology and shared set of database resources for non-canonical ORFs. Ultimately, we believe that the study of non-canonical ORFs, if pursued with proper precision, will prove invaluable to the global community of biomedical researchers.

Acknowledgements

J.R.P. acknowledges funding from the National Institutes of Health/National Cancer Institute (K08-CA263552-01A1), the Alex's Lemonade Stand Foundation Young Investigator Award (#21-23983), the St. Baldrick's Foundation Scholar Award (#931638), the Musella Foundation for Brain Tumor Research, the DIPG/DMG Research Funding Alliance, and a Collaborative Pediatric Cancer Research Awards Program/Kids Join the Fight award (#22FN23). E.W.D. acknowledges funding from National Institutes of Health grant R01 GM087221 and National Science Foundation grant DBI-1933311. S.v.H. acknowledges funding from Fonds Cancers, Stichting Reggeborgh, Stichting Bergh in het Zadel, and Stichting Villa Joep. J.G.A. and K.R.C. were supported in part by grants P01CA206978 from the NIH, U24CA270823, U01CA271402 and U24CA271075 from National Cancer Institute (NCI) Clinical Proteomic Tumor Analysis Consortium program and from the Dr. Miriam and Sheldon G. Adelson Medical Research Foundation. J.M.M. is supported by the Wellcome Trust (grant number 108749/Z/15/Z), the National Human Genome Research Institute (NHGRI) of the US National Institutes of Health (NIH) under award number 2U41HG007234, and the European Molecular Biology Laboratory (EMBL). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. Ensembl is a registered trademark of EMBL.

Author Contributions

Conceptualization: E.W.D., J.R.P., S.v.H., J.R.-O., J.M.M., M.B.-S., J.G.A., K.R.C.

Methodology, E.W.D., J.R.P., S.v.H., J.R.-O., J.M.M., M.B.-S., J.G.A., K.R.C., L.W.K.

Formal analysis, J.R.P., L.W.K.

Data curation, J.R.P., J.R.-O., L.W.K.

Writing - E.W.D., J.R.P., S.v.H., J.R.-O., J.M.M., M.B.-S., J.G.A., K.R.C.

Writing - review & editing, E.W.D., J.R.P., S.v.H., J.R.-O., J.M.M., M.B.-S., J.G.A., K.R.C., L.W.K.

Visualization, J.R.P., L.W.K.

Supervision, E.W.D., J.R.P., S.v.H., J.R.-O., J.M.M., M.B.-S., J.G.A., K.R.C.

Funding acquisition, E.W.D., J.R.P., S.v.H.

Declaration of interests

The authors declare no competing interests.

Inclusion and diversity

We support inclusive, diverse and equitable research.

Methods

Benchmarking and comparing ORF caller performance on replicate Ribo-seq datasets

Ribo-seq data processing and mapping: Ribosome profiling data of late pancreatic progenitor cells obtained from six independent differentiations of H1 human embryonic stem cells (11) were collected from the GEO database (GSE144682). For all analyses, the Ensembl primary DNA assembly (GRCh38) and the Ensembl human reference transcriptome (Ensembl v102) were used as reference. Quality control and trimming of the Ribo-seq reads was done using Trim Galore 0.6.6 with the options ‘--length 25’, and ‘--trim-n’ (151). Next, contaminant RNA and DNA were removed using Bowtie2 2.4.2 by aligning reads to a contaminant file using the default options of Bowtie2 (152). The contaminant-depleted reads were aligned using STAR with the options ‘--twopassMode Basic’, ‘--outFilterMismatchNmax 2’, ‘--outFilterMultimapNmax 20’, ‘--limitOutSJcollapsed 10000000’, ‘--alignSJoverhangMin 1000’, and ‘--outSAMattributes All’ (153). For PRICE, the option ‘--alignEndsType EndToEnd’ was set as well. Also, the individual bamfiles were filtered using SAMtools 1.12 to exclude reads with a mapping quality lower than 5 (154).

ORF calling with ORFquant: The function RiboseQC_analysis from RiboseQC 1.1 was run in R 4.1.2 with the options ‘read_subset’ and ‘fast_mode’ set to false (155). The output was used by the function run_ORFquant from ORFquant 1.02 in R with the default options (10). ORF calling with PRICE: Before using PRICE, a reference genome was created with the IndexGenome function of the Gedi framework 1.0.2. After the creation of the reference genome, PRICE 1.0.3b was run (67). A filtered list of ORFs detected by PRICE and a list of P-sites (called activity values by PRICE) were extracted from the outputted ‘orfs.cit’ files using the Gedi Nashorn and ViewCIT functions respectively. Because the start codon prediction is a separate step in the PRICE program, ORF coordinates from both before and after start codon prediction were available. We used the coordinates after start codon prediction. PRICE can also be run in a multi sample mode by providing a text file with the bam file locations as input. This mode favours ORFs that occur in all samples during the ORF calling process and would likely enhance the reproducibility of ORF calls between replicates. To keep all ORF callers comparable we did not use this mode. ORF calling with Ribo-TISH: From Ribo-TISH 0.2.7, the predict function was used to infer ORFs with the option ‘--longest’ set (66). The output file contained only the genomic start and end coordinates and the transcript id of each ORF. The reference GTF was used to determine the exons within each ORF. ORF calling with Ribotricer: The Ribotricer 1.3.3 function prepare_orfs was first used with the options ‘--longest’, ‘--min_orf_length 9’ (69). The option ‘--start_codons’ was set to include all near cognate start codons with one base difference compared to ATG. Afterwards, the function “detect_orfs” was used with the option ‘--phase_score_cutoff 0.440’.

Comparing ORF callers: ORF calls were compared between algorithms for the types of ORF categories that were found, in how many replicates they were independently discovered, how ORF differed in length, and how reproducible and similar their detection was based on e.g., the percentage of ORF sequence overlap between replicate ORF calls. Before the analyses, data were converted to GRangesList objects in R with stop codons included in the coordinates. ORF categories were determined by comparing the start and end coordinates, and the transcript id of each ORF with the coding sequences in the ‘gtf.rannot’ object created by the ORFquant function ‘prepare_annotation_files’. ORFs were compared by their overlap, with

different thresholds set for the required percentage of overlap. Two ORFs were considered to be similar if the exons of one ORF were fully contained within the exons of a second ORF, both codons had the same stop codon, and the first ORF covered at least the required percentage of overlap of the length of the second ORF. These overlap relations were recursive, such that a parent ORF could be the child of another ORF, and all three would be counted as one unique ORF.

Code availability: All code used for these analyses as well as data visualization are available at https://bitbucket.org/vanHeeschLab/orfcaller_comparison.

Comparison of published Ribo-seq datasets

We used publicly-available datasets from GENCODE (16), Chothani et al. (15), Ouspenskaia et al. (34), and Duffy et al. (122) for comparisons of published reports of non-canonical ORFs that might encode microproteins. The GENCODE dataset itself is a metaanalysis of data from Ji et al. (19), Calviello et al. (61), Raj et al. (20), van Heesch et al. (9), Martinez et al. (21), Chen et al. (18), and Gaertner et al. (11); datasets employed are listed in **Supplementary Table 1**. Source data for these datasets is listed in **Supplementary Table 2**. To facilitate comparisons between studies, we extracted only non-canonical ORFs with a length of ≥ 16 amino acids and had an AUG start codon. For ORFs using a non-AUG start site, the first internal AUG start codon was identified and the amino acid sequence starting with that internal AUG was included for analysis if the resulting ORF was ≥ 16 amino acids long. ORFs were then analyzed for their replication across primary datasets. Since the GENCODE list represents a metaanalysis of other individual datasets, presence of an ORF in the GENCODE list was not used as part of the analysis for ORF replication across primary datasets. Next, ORF calls were associated with one of the following six categories: lncRNA-ORF, upstream ORF, upstream overlapping ORF, internal ORF, downstream overlapping ORF, or downstream ORF, according to schema by Mudge et al. (16). Duffy et al. used the nomenclature “external” for doORF, and these ORFs were reclassified as doORF for this analysis; they used “internal” for intORFs, which were reclassified as intORFs for this analysis. For lncRNAs, Duffy et al. used the term “non-coding” which included the biotypes “noncoding”, “lncRNA”, “antisense_RNA”, “misc_RNA”, “TEC” and “processed_transcript”, which were included as part of the lncRNA-ORF designation for this study. For Ouspenskaia et al., we analyzed ORFs according to the authors’ designation of ORF “plotType”, reflecting their final classification. Ouspenskaia et al. used the term “3’ dORF” for dORF, “3’ overlap dORF” for doORF, “5’ overlap uORF” of uORF, “5’ uORF” for uORF, “lncRNA” for lncRNA-ORF, and “out-of-frame” for intORF. Chothani et al. reported final ORF types of “dORF”, “doORF”, “ncORF”, “overlap_uORF”, “intORF” and “uORF”. For Chothani et al., Duffy et al. and Ouspenskaia et al., ORFs that had a final classification of pseudogene were excluded from this analysis; however, these datasets variably reclassified some ORFs on pseudogene transcript biotypes as non-coding or lncRNA, and we did not re-filter these ORFs beyond the original reclassifications provided by the authors. ORFs switch a classification corresponding to a small RNA, tRNA or rRNA species, such as “rRNA”, “snoRNA”, “tRNA”, “snRNA” or “miRNA”, were excluded from this analysis. The number of cell types and/or tissue types for analyses of each ORF dataset was extracted from the source publication.

References

1. Aebersold, R., Agar, J. N., Amster, I. J., Baker, M. S., Bertozzi, C. R., Boja, E. S., Costello, C. E., Cravatt, B. F., Fenselau, C., Garcia, B. A., Ge, Y., Gunawardena, J., Hendrickson, R. C., Hergenrother, P. J., Huber, C. G., Ivanov, A. R., Jensen, O. N., Jewett, M. C., Kelleher, N. L., Kiessling, L. L., Krogan, N. J., Larsen, M. R., Loo, J. A., Ogorzalek Loo, R. R., Lundberg, E., MacCoss, M. J., Mallick, P., Mootha, V. K., Mrksich, M., Muir, T. W., Patrie, S. M., Pesavento, J. J., Pitteri, S. J., Rodriguez, H., Saghatelian, A., Sandoval, W., Schlüter, H., Sechi, S., Slavoff, S. A., Smith, L. M., Snyder, M. P., Thomas, P. M., Uhlén, M., Van Eyk, J. E., Vidal, M., Walt, D. R., White, F. M., Williams, E. R., Wohlschläger, T., Wysocki, V. H., Yates, N. A., Young, N. L., and Zhang, B. (2018) How many human proteoforms are there? *Nat. Chem. Biol.* **14**, 206–214
2. Tress, M. L., Abascal, F., and Valencia, A. (2017) Alternative Splicing May Not Be the Key to Proteome Complexity. *Trends Biochem. Sci.* **42**, 98–110
3. Blencowe, B. J. (2017) The Relationship between Alternative Splicing and Proteomic Complexity. *Trends in Biochemical Sciences.* **42**, 407–408
4. Sinitcyn, P., Richards, A. L., Weatheritt, R. J., Brademan, D. R., Marx, H., Shishkova, E., Meyer, J. G., Hebert, A. S., Westphall, M. S., Blencowe, B. J., Cox, J., and Coon, J. J. (2023) Global detection of human variants and isoforms by deep proteome sequencing. *Nat. Biotechnol.* 10.1038/s41587-023-01714-x
5. Frankish, A., Carbonell-Sala, S., Diekhans, M., Jungreis, I., Loveland, J. E., Mudge, J. M., Sisu, C., Wright, J. C., Arnan, C., Barnes, I., Banerjee, A., Bennett, R., Berry, A., Bignell, A., Boix, C., Calvet, F., Cerdán-Vélez, D., Cunningham, F., Davidson, C., Donaldson, S., Dursun, C., Fatima, R., Giorgetti, S., Giron, C. G., Gonzalez, J. M., Hardy, M., Harrison, P. W., Hourlier, T., Hollis, Z., Hunt, T., James, B., Jiang, Y., Johnson, R., Kay, M., Lagarde, J., Martin, F. J., Gómez, L. M., Nair, S., Ni, P., Pozo, F., Ramalingam, V., Ruffier, M., Schmitt, B. M., Schreiber, J. M., Steed, E., Suner, M.-M., Sumathipala, D., Sycheva, I., Uszczynska-Ratajczak, B., Wass, E., Yang, Y. T., Yates, A., Zafrulla, Z., Choudhary, J. S., Gerstein, M., Guigo, R., Hubbard, T. J. P., Kellis, M., Kundaje, A., Paten, B., Tress, M. L., and Flicek, P. (2023) GENCODE: reference annotation for the human and mouse genomes in 2023. *Nucleic Acids Res.* **51**, D942–D949
6. UniProt Consortium (2023) UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Res.* **51**, D523–D531
7. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S., and Weissman, J. S. (2009) Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science.* **324**, 218–223
8. McGlincy, N. J., and Ingolia, N. T. (2017) Transcriptome-wide measurement of translation by ribosome profiling. *Methods.* **126**, 112–129
9. van Heesch, S., Witte, F., Schneider-Lunitz, V., Schulz, J. F., Adami, E., Faber, A. B., Kirchner, M., Maatz, H., Blachut, S., Sandmann, C.-L., Kanda, M., Worth, C. L., Schafer, S., Calviello, L., Merriott, R., Patone, G., Hummel, O., Wyler, E., Obermayer, B., Mücke, M. B., Lindberg, E. L., Trnka, F., Memczak, S., Schilling, M., Felkin, L. E., Barton, P. J. R., Quaife, N. M., Vanezis, K., Diecke, S., Mukai, M., Mah, N., Oh, S.-J., Kurtz, A., Schramm, C., Schwinge, D., Sebode, M., Harakalova, M., Asselbergs, F. W., Vink, A., de Weger, R. A., Viswanathan, S., Widjaja, A. A., Gärtner-Rommel, A., Milting, H., Dos Remedios, C., Knosalla, C., Mertins, P., Landthaler, M., Vingron, M., Linke, W. A., Seidman, J. G., Seidman, C. E., Rajewsky, N., Ohler, U., Cook, S. A., and Hubner, N. (2019) The Translational Landscape of the Human Heart. *Cell.* **178**, 242–260.e29
10. Calviello, L., Hirsekorn, A., and Ohler, U. (2020) Quantification of translation uncovers the functions of the alternative transcriptome. *Nature Structural & Molecular Biology.* **27**, 717–725
11. Gaertner, B., van Heesch, S., Schneider-Lunitz, V., Schulz, J. F., Witte, F., Blachut, S., Nguyen, S.,

- Wong, R., Matta, I., Hübner, N., and Sander, M. (2020) A human ESC-based screen identifies a role for the translated lncRNA in pancreatic endocrine differentiation. *Elife*. 10.7554/eLife.58659
12. Fagerberg, L., Hallström, B. M., Oksvold, P., Kampf, C., Djureinovic, D., Odeberg, J., Habuka, M., Tahmasebpour, S., Danielsson, A., Edlund, K., Asplund, A., Sjöstedt, E., Lundberg, E., Szgyarto, C. A.-K., Skogs, M., Takanen, J. O., Berling, H., Tegel, H., Mulder, J., Nilsson, P., Schwenk, J. M., Lindskog, C., Danielsson, F., Mardinoglu, A., Sivertsson, A., von Feilitzen, K., Forsberg, M., Zwahlen, M., Olsson, I., Navani, S., Huss, M., Nielsen, J., Ponten, F., and Uhlen, M. (2014) Analysis of the human tissue-specific expression by genome-wide integration of transcriptomics and antibody-based proteomics. *Mol. Cell. Proteomics*. **13**, 397–406
13. Krug, K., Jaehnig, E. J., Satpathy, S., Blumenberg, L., Karpova, A., Anurag, M., Miles, G., Mertins, P., Geffen, Y., Tang, L. C., Heiman, D. I., Cao, S., Maruvka, Y. E., Lei, J. T., Huang, C., Kothadia, R. B., Colaprico, A., Birger, C., Wang, J., Dou, Y., Wen, B., Shi, Z., Liao, Y., Wiznerowicz, M., Wyczalkowski, M. A., Chen, X. S., Kennedy, J. J., Paulovich, A. G., Thiagarajan, M., Kinsinger, C. R., Hiltke, T., Boja, E. S., Mesri, M., Robles, A. I., Rodriguez, H., Westbrook, T. F., Ding, L., Getz, G., Clauser, K. R., Fenyo, D., Ruggles, K. V., Zhang, B., Mani, D. R., Carr, S. A., Ellis, M. J., Gillette, M. A., and Clinical Proteomic Tumor Analysis Consortium (2020) Proteogenomic Landscape of Breast Cancer Tumorigenesis and Targeted Therapy. *Cell*. **183**, 1436–1456.e31
14. Cao, L., Huang, C., Cui Zhou, D., Hu, Y., Lih, T. M., Savage, S. R., Krug, K., Clark, D. J., Schnaubelt, M., Chen, L., da Veiga Leprevost, F., Egeuz, R. V., Yang, W., Pan, J., Wen, B., Dou, Y., Jiang, W., Liao, Y., Shi, Z., Terekhanova, N. V., Cao, S., Lu, R. J.-H., Li, Y., Liu, R., Zhu, H., Ronning, P., Wu, Y., Wyczalkowski, M. A., Easwaran, H., Danilova, L., Mer, A. S., Yoo, S., Wang, J. M., Liu, W., Haibe-Kains, B., Thiagarajan, M., Jewell, S. D., Hostetter, G., Newton, C. J., Li, Q. K., Roehrl, M. H., Fenyo, D., Wang, P., Nesvizhskii, A. I., Mani, D. R., Omenn, G. S., Boja, E. S., Mesri, M., Robles, A. I., Rodriguez, H., Bathe, O. F., Chan, D. W., Hruban, R. H., Ding, L., Zhang, B., Zhang, H., and Clinical Proteomic Tumor Analysis Consortium (2021) Proteogenomic characterization of pancreatic ductal adenocarcinoma. *Cell*. **184**, 5031–5052.e26
15. Chothani, S. P., Adami, E., Widjaja, A. A., Langley, S. R., Viswanathan, S., Pua, C. J., Zhihao, N. T., Harmston, N., D’Agostino, G., Whiffin, N., Mao, W., Ouyang, J. F., Lim, W. W., Lim, S., Lee, C. Q. E., Grubman, A., Chen, J., Kovalik, J. P., Tryggvason, K., Polo, J. M., Ho, L., Cook, S. A., Rackham, O. J. L., and Schafer, S. (2022) A high-resolution map of human RNA translation. *Mol. Cell*. **82**, 2885–2899.e8
16. Mudge, J. M., Ruiz-Orera, J., Prensner, J. R., Brunet, M. A., Calvet, F., Jungreis, I., Gonzalez, J. M., Magrane, M., Martinez, T. F., Schulz, J. F., Yang, Y. T., Albà, M. M., Aspdén, J. L., Baranov, P. V., Bazzini, A. A., Bruford, E., Martin, M. J., Calviello, L., Carvunis, A.-R., Chen, J., Couso, J. P., Deutsch, E. W., Flicek, P., Frankish, A., Gerstein, M., Hubner, N., Ingolia, N. T., Kellis, M., Menschaert, G., Moritz, R. L., Ohler, U., Roucou, X., Saghatelian, A., Weissman, J. S., and van Heesch, S. (2022) Standardized annotation of translated open reading frames. *Nat. Biotechnol.* **40**, 994–999
17. Prensner, J. R., Enache, O. M., Luria, V., Krug, K., Clauser, K. R., Dempster, J. M., Karger, A., Wang, L., Stumbraite, K., Wang, V. M., Botta, G., Lyons, N. J., Goodale, A., Kalani, Z., Fritchman, B., Brown, A., Alan, D., Green, T., Yang, X., Jaffe, J. D., Roth, J. A., Piccioni, F., Kirschner, M. W., Ji, Z., Root, D. E., and Golub, T. R. (2021) Noncanonical open reading frames encode functional proteins essential for cancer cell survival. *Nat. Biotechnol.* **39**, 697–704
18. Chen, J., Brunner, A.-D., Cogan, J. Z., Nuñez, J. K., Fields, A. P., Adamson, B., Itzhak, D. N., Li, J. Y., Mann, M., Leonetti, M. D., and Weissman, J. S. (2020) Pervasive functional translation of noncanonical human open reading frames. *Science*. **367**, 1140–1146
19. Ji, Z., Song, R., Regev, A., and Struhl, K. (2015) Many lncRNAs, 5’UTRs, and pseudogenes are translated and some are likely to express functional proteins. *Elife*. **4**, e08890
20. Raj, A., Wang, S. H., Shim, H., Harpak, A., Li, Y. I., Engelmann, B., Stephens, M., Gilad, Y., and

- Pritchard, J. K. (2016) Thousands of novel translated open reading frames in humans inferred by ribosome footprint profiling. *Elife*. 10.7554/eLife.13328
21. Martinez, T. F., Chu, Q., Donaldson, C., Tan, D., Shokhirev, M. N., and Saghatelian, A. (2020) Accurate annotation of human protein-coding small open reading frames. *Nat. Chem. Biol.* **16**, 458–468
22. Aspden, J. L., Eyre-Walker, Y. C., Phillips, R. J., Amin, U., Mumtaz, M. A. S., Brocard, M., and Couso, J.-P. (2014) Extensive translation of small Open Reading Frames revealed by Poly-Ribo-Seq. *Elife*. **3**, e03528
23. Douka, K., Birds, I., Wang, D., Kosteletos, A., Clayton, S., Byford, A., Vasconcelos, E. J. R., O’Connell, M. J., Deuchars, J., Whitehouse, A., and Aspden, J. L. (2021) Cytoplasmic long noncoding RNAs are differentially regulated and translated during human neuronal differentiation. *RNA*. **27**, 1082–1101
24. Fedorova, A. D., Kiniry, S. J., Andreev, D. E., Mudge, J. M., and Baranov, P. V. (2022) Thousands of human non-AUG extended proteoforms lack evidence of evolutionary selection among mammals. *Nat. Commun.* **13**, 7910
25. Van Damme, P., Gawron, D., Van Crielinge, W., and Menschaert, G. (2014) N-terminal proteomics and ribosome profiling provide a comprehensive view of the alternative translation initiation landscape in mice and men. *Mol. Cell. Proteomics*. **13**, 1245–1261
26. Koch, A., Gawron, D., Steyaert, S., Ndah, E., Crappé, J., De Keulenaer, S., De Meester, E., Ma, M., Shen, B., Gevaert, K., Van Crielinge, W., Van Damme, P., and Menschaert, G. (2014) A proteogenomics approach integrating proteomics and ribosome profiling increases the efficiency of protein identification and enables the discovery of alternative translation start sites. *Proteomics*. **14**, 2688–2698
27. Menschaert, G., Van Crielinge, W., Notelaers, T., Koch, A., Crappé, J., Gevaert, K., and Van Damme, P. (2013) Deep proteome coverage based on ribosome profiling aids mass spectrometry-based protein and peptide discovery and provides evidence of alternative translation products and near-cognate translation initiation events. *Mol. Cell. Proteomics*. **12**, 1780–1790
28. Griffin, G. K., Wu, J., Iracheta-Vellve, A., Patti, J. C., Hsu, J., Davis, T., Dele-Oni, D., Du, P. P., Halawi, A. G., Ishizuka, J. J., Kim, S. Y., Klaeger, S., Knudsen, N. H., Miller, B. C., Nguyen, T. H., Olander, K. E., Papanastasiou, M., Rachimi, S., Robitschek, E. J., Schneider, E. M., Yearly, M. D., Zimmer, M. D., Jaffe, J. D., Carr, S. A., Doench, J. G., Haining, W. N., Yates, K. B., Manguso, R. T., and Bernstein, B. E. (2021) Epigenetic silencing by SETDB1 suppresses tumour intrinsic immunogenicity. *Nature*. **595**, 309–314
29. Al-Turki, T. M., and Griffith, J. D. (2023) Mammalian telomeric RNA (TERRA) can be translated to produce valine–arginine and glycine–leucine dipeptide repeat proteins. *Proceedings of the National Academy of Sciences*. **120**, e2221529120
30. Omenn, G. S., Lane, L., Lundberg, E. K., Overall, C. M., and Deutsch, E. W. (2017) Progress on the HUPO Draft Human Proteome: 2017 Metrics of the Human Proteome Project. *J. Proteome Res.* **16**, 4281–4287
31. Schwaid, A. G., Shannon, D. A., Ma, J., Slavoff, S. A., Levin, J. Z., Weerapana, E., and Saghatelian, A. (2013) Chemoproteomic discovery of cysteine-containing human short open reading frames. *J. Am. Chem. Soc.* **135**, 16750–16753
32. Cao, X., Khitun, A., Na, Z., Dumitrescu, D. G., Kubica, M., Olatunji, E., and Slavoff, S. A. (2020) Comparative Proteomic Profiling of Unannotated Microproteins and Alternative Proteins in Human Cell Lines. *J. Proteome Res.* **19**, 3418–3426
33. Ma, J., Diedrich, J. K., Jungreis, I., Donaldson, C., Vaughan, J., Kellis, M., Yates, J. R., 3rd, and Saghatelian, A. (2016) Improved Identification and Analysis of Small Open Reading Frame Encoded Polypeptides. *Anal. Chem.* **88**, 3967–3975
34. Ouspenskaia, T., Law, T., Clauser, K. R., Klaeger, S., Sarkizova, S., Aguet, F., Li, B., Christian, E.,

- Knisbacher, B. A., Le, P. M., Hartigan, C. R., Keshishian, H., Apffel, A., Oliveira, G., Zhang, W., Chen, S., Chow, Y. T., Ji, Z., Jungreis, I., Shukla, S. A., Justesen, S., Bachireddy, P., Kellis, M., Getz, G., Hacohen, N., Keskin, D. B., Carr, S. A., Wu, C. J., and Regev, A. (2022) Unannotated proteins expand the MHC-I-restricted immunopeptidome in cancer. *Nat. Biotechnol.* **40**, 209–217
35. Chong, C., Coukos, G., and Bassani-Sternberg, M. (2022) Identification of tumor antigens with immunopeptidomics. *Nat. Biotechnol.* **40**, 175–188
36. Laumont, C. M., Vincent, K., Hesnard, L., Audemard, É., Bonneil, É., Laverdure, J.-P., Gendron, P., Courcelles, M., Hardy, M.-P., Côté, C., Durette, C., St-Pierre, C., Benhammadi, M., Lanoix, J., Vobecky, S., Haddad, E., Lemieux, S., Thibault, P., and Perreault, C. (2018) Noncoding regions are the main source of targetable tumor-specific antigens. *Sci. Transl. Med.* 10.1126/scitranslmed.aau5516
37. Laumont, C. M., and Perreault, C. (2018) Exploiting non-canonical translation to identify new targets for T cell-based cancer immunotherapy. *Cell. Mol. Life Sci.* **75**, 607–621
38. Laumont, C. M., Daouda, T., Laverdure, J.-P., Bonneil, É., Caron-Lizotte, O., Hardy, M.-P., Granados, D. P., Durette, C., Lemieux, S., Thibault, P., and Perreault, C. (2016) Global proteogenomic analysis of human MHC class I-associated peptides derived from non-canonical reading frames. *Nat. Commun.* **7**, 10238
39. Ruiz Cuevas, M. V., Hardy, M.-P., Hollý, J., Bonneil, É., Durette, C., Courcelles, M., Lanoix, J., Côté, C., Staudt, L. M., Lemieux, S., Thibault, P., Perreault, C., and Yewdell, J. W. (2021) Most non-canonical proteins uniquely populate the proteome or immunopeptidome. *Cell Rep.* **34**, 108815
40. Calvo, S. E., Pagliarini, D. J., and Mootha, V. K. (2009) Upstream open reading frames cause widespread reduction of protein expression and are polymorphic among humans. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 7507–7512
41. Johnstone, T. G., Bazzini, A. A., and Giraldez, A. J. (2016) Upstream ORFs are prevalent translational repressors in vertebrates. *EMBO J.* **35**, 706–723
42. Wu, Q., Wright, M., Gogol, M. M., Bradford, W. D., Zhang, N., and Bazzini, A. A. (2020) Translation of small downstream ORFs enhances translation of canonical main open reading frames. *EMBO J.* **39**, e104763
43. Pauli, A., Norris, M. L., Valen, E., Chew, G.-L., Gagnon, J. A., Zimmerman, S., Mitchell, A., Ma, J., Dubrulle, J., Reyon, D., Tsai, S. Q., Joung, J. K., Saghatelian, A., and Schier, A. F. (2014) Toddler: an embryonic signal that promotes cell movement via Apelin receptors. *Science.* **343**, 1248636
44. Khan, Y. A., Jungreis, I., Wright, J. C., Mudge, J. M., Choudhary, J. S., Firth, A. E., and Kellis, M. (2020) Evidence for a novel overlapping coding sequence in POLG initiated at a CUG start codon. *BMC Genet.* **21**, 25
45. Loughran, G., Zhdanov, A. V., Mikhaylova, M. S., Rozov, F. N., Datskevich, P. N., Kovalchuk, S. I., Serebryakova, M. V., Kiniry, S. J., Michel, A. M., O'Connor, P. B. F., Papkovsky, D. B., Atkins, J. F., Baranov, P. V., Shatsky, I. N., and Andreev, D. E. (2020) Unusually efficient CUG initiation of an overlapping reading frame in mRNA yields novel protein POLGARF. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 24936–24946
46. Boix, O., Martinez, M., Vidal, S., Giménez-Alejandro, M., Palenzuela, L., Lorenzo-Sanz, L., Quevedo, L., Moscoso, O., Ruiz-Orera, J., Ximénez-Embún, P., Ciriaco, N., Nuciforo, P., Stephan-Otto Attolini, C., Albà, M. M., Muñoz, J., Tian, T. V., Varela, I., Vivancos, A., Ramón Y Cajal, S., Muñoz, P., Rivas, C., and Abad, M. (2022) pTINCR microprotein promotes epithelial differentiation and suppresses tumor growth through CDC42 SUMOylation and activation. *Nat. Commun.* **13**, 6840
47. Bi, P., Ramirez-Martinez, A., Li, H., Cannavino, J., McAnally, J. R., Shelton, J. M., Sánchez-Ortiz, E., Bassel-Duby, R., and Olson, E. N. (2017) Control of muscle formation by the fusogenic micropeptide myomixer. *Science.* **356**, 323–327
48. Anderson, D. M., Anderson, K. M., Chang, C.-L., Makarewich, C. A., Nelson, B. R., McAnally, J. R.,

- Kasaragod, P., Shelton, J. M., Liou, J., Bassel-Duby, R., and Olson, E. N. (2015) A micropeptide encoded by a putative long noncoding RNA regulates muscle performance. *Cell*. **160**, 595–606
49. Slavoff, S. A., Mitchell, A. J., Schwaid, A. G., Cabili, M. N., Ma, J., Levin, J. Z., Karger, A. D., Budnik, B. A., Rinn, J. L., and Saghatelian, A. (2013) Peptidomic discovery of short open reading frame-encoded peptides in human cells. *Nat. Chem. Biol.* **9**, 59–64
50. Martinez, T. F., Lyons-Abbott, S., Bookout, A. L., De Souza, E. V., Donaldson, C., Vaughan, J. M., Lau, C., Abramov, A., Baquero, A. F., Baquero, K., Friedrich, D., Huard, J., Davis, R., Kim, B., Koch, T., Mercer, A. J., Misquith, A., Murray, S. A., Perry, S., Pino, L. K., Sanford, C., Simon, A., Zhang, Y., Zipp, G., Bizarro, C. V., Shokhirev, M. N., Whittle, A. J., Searle, B. C., MacCoss, M. J., Saghatelian, A., and Barnes, C. A. (2023) Profiling mouse brown and white adipocytes to identify metabolically relevant small ORFs and functional microproteins. *Cell Metab.* **35**, 166–183.e11
51. Mackowiak, S. D., Zauber, H., Bielow, C., Thiel, D., Kutz, K., Calviello, L., Mastrobuoni, G., Rajewsky, N., Kempa, S., Selbach, M., and Obermayer, B. (2015) Extensive identification and analysis of conserved small ORFs in animals. *Genome Biol.* **16**, 179
52. Bazzini, A. A., Johnstone, T. G., Christiano, R., Mackowiak, S. D., Obermayer, B., Fleming, E. S., Vejnar, C. E., Lee, M. T., Rajewsky, N., Walther, T. C., and Giraldez, A. J. (2014) Identification of small ORFs in vertebrates using ribosome footprinting and evolutionary conservation. *EMBO J.* **33**, 981–993
53. Chong, C., Müller, M., Pak, H., Harnett, D., Huber, F., Grun, D., Leleu, M., Auger, A., Arnaud, M., Stevenson, B. J., Michaux, J., Bilic, I., Hirsekorn, A., Calviello, L., Simó-Riudalbas, L., Planet, E., Lubiński, J., Bryśkiewicz, M., Wiznerowicz, M., Xenarios, I., Zhang, L., Trono, D., Harari, A., Ohler, U., Coukos, G., and Bassani-Sternberg, M. (2020) Integrated proteogenomic deep sequencing and analytics accurately identify non-canonical peptides in tumor immunopeptidomes. *Nat. Commun.* **11**, 1293
54. Huang, N., Li, F., Zhang, M., Zhou, H., Chen, Z., Ma, X., Yang, L., Wu, X., Zhong, J., Xiao, F., Yang, X., Zhao, K., Li, X., Xia, X., Liu, Z., Gao, S., and Zhang, N. (2021) An Upstream Open Reading Frame in Phosphatase and Tensin Homolog Encodes a Circuit Breaker of Lactate Metabolism. *Cell Metab.* **33**, 454
55. Na, Z., Dai, X., Zheng, S.-J., Bryant, C. J., Loh, K. H., Su, H., Luo, Y., Buhagiar, A. F., Cao, X., Baserga, S. J., Chen, S., and Slavoff, S. A. (2022) Mapping subcellular localizations of unannotated microproteins and alternative proteins with MicroID. *Mol. Cell.* **82**, 2900–2911.e7
56. Jayaram, D. R., Frost, S., Argov, C., Liju, V. B., Anto, N. P., Muraleedharan, A., Ben-Ari, A., Sinay, R., Smoly, I., Novoplansky, O., Isakov, N., Toiber, D., Keasar, C., Elkabets, M., Yeger-Lotem, E., and Livneh, E. (2021) Unraveling the hidden role of a uORF-encoded peptide as a kinase inhibitor of PKCs. *Proc. Natl. Acad. Sci. U. S. A.* 10.1073/pnas.2018899118
57. Sandmann, C.-L., Schulz, J. F., Ruiz-Orera, J., Kirchner, M., Ziehm, M., Adami, E., Marczenke, M., Christ, A., Liebe, N., Greiner, J., Schoenenberger, A., Muecke, M. B., Liang, N., Moritz, R. L., Sun, Z., Deutsch, E. W., Gotthardt, M., Mudge, J. M., Prensner, J. R., Willnow, T. E., Mertins, P., van Heesch, S., and Hubner, N. (2023) Evolutionary origins and interactomes of human, young microproteins and small peptides translated from short open reading frames. *Mol. Cell.* 10.1016/j.molcel.2023.01.023
58. Tanaka, M., Sotta, N., Yamazumi, Y., Yamashita, Y., Miwa, K., Murota, K., Chiba, Y., Hirai, M. Y., Akiyama, T., Onouchi, H., Naito, S., and Fujiwara, T. (2016) The Minimum Open Reading Frame, AUG-Stop, Induces Boron-Dependent Ribosome Stalling and mRNA Degradation. *Plant Cell.* **28**, 2830–2849
59. Dau, T., Bartolomucci, G., and Rappsilber, J. (2020) Proteomics Using Protease Alternatives to Trypsin Benefits from Sequential Digestion with Trypsin. *Anal. Chem.* **92**, 9523–9527
60. Calviello, L., and Ohler, U. (2017) Beyond Read-Counts: Ribo-seq Data Analysis to Understand the

- Functions of the Transcriptome. *Trends Genet.* **33**, 728–744
61. Calviello, L., Mukherjee, N., Wyler, E., Zauber, H., Hirsekorn, A., Selbach, M., Landthaler, M., Obermayer, B., and Ohler, U. (2016) Detecting actively translated open reading frames in ribosome profiling data. *Nat. Methods.* **13**, 165–170
 62. Fremin, B. J., and Bhatt, A. S. (2020) Structured RNA Contaminants in Bacterial Ribo-Seq. *mSphere*. 10.1128/mSphere.00855-20
 63. Chung, B. Y., Hardcastle, T. J., Jones, J. D., Irigoyen, N., Firth, A. E., Baulcombe, D. C., and Brierley, I. (2015) The use of duplex-specific nuclease in ribosome profiling and a user-friendly software package for Ribo-seq data analysis. *RNA.* **21**, 1731–1745
 64. Hsu, P. Y., Calviello, L., Wu, H.-Y. L., Li, F.-W., Rothfels, C. J., Ohler, U., and Benfey, P. N. (2016) Super-resolution ribosome profiling reveals unannotated translation events in. *Proc. Natl. Acad. Sci. U. S. A.* **113**, E7126–E7135
 65. Diamant, A., and Tuller, T. (2016) Estimation of ribosome profiling performance and reproducibility at various levels of resolution. *Biol. Direct.* **11**, 24
 66. Zhang, P., He, D., Xu, Y., Hou, J., Pan, B.-F., Wang, Y., Liu, T., Davis, C. M., Ehli, E. A., Tan, L., Zhou, F., Hu, J., Yu, Y., Chen, X., Nguyen, T. M., Rosen, J. M., Hawke, D. H., Ji, Z., and Chen, Y. (2017) Genome-wide identification and differential analysis of translational initiation. *Nat. Commun.* **8**, 1749
 67. Erhard, F., Halenius, A., Zimmermann, C., L'Hernault, A., Kowalewski, D. J., Weekes, M. P., Stevanovic, S., Zimmer, R., and Dölken, L. (2018) Improved Ribo-seq enables identification of cryptic translation events. *Nat. Methods.* **15**, 363–366
 68. Xiao, Z., Huang, R., Xing, X., Chen, Y., Deng, H., and Yang, X. (2018) De novo annotation and characterization of the translatome with ribosome profiling data. *Nucleic Acids Res.* **46**, e61
 69. Choudhary, S., Li, W., and D Smith, A. (2020) Accurate detection of short and long active ORFs using Ribo-seq data. *Bioinformatics.* **36**, 2053–2059
 70. Fields, A. P., Rodriguez, E. H., Jovanovic, M., Stern-Ginossar, N., Haas, B. J., Mertins, P., Raychowdhury, R., Hacohen, N., Carr, S. A., Ingolia, N. T., Regev, A., and Weissman, J. S. (2015) A Regression-Based Analysis of Ribosome-Profiling Data Reveals a Conserved Complexity to Mammalian Translation. *Mol. Cell.* **60**, 816–827
 71. Clauwaert, J., Menschaert, G., and Waegeman, W. (2019) DeepRibo: a neural network for precise gene annotation of prokaryotes by combining ribosome profiling signal and binding site patterns. *Nucleic Acids Res.* **47**, e36
 72. Clauwaert, J., McVey, Z., Gupta, R., and Menschaert, G. (2023) TIS Transformer: remapping the human proteome using deep learning. *NAR Genom Bioinform.* **5**, lqad021
 73. Freudenmann, L. K., Marcu, A., and Stevanović, S. (2018) Mapping the tumour human leukocyte antigen (HLA) ligandome by mass spectrometry. *Immunology.* **154**, 331–345
 74. Bassani-Sternberg, M., Bräunlein, E., Klar, R., Engleitner, T., Sinitsyn, P., Audehm, S., Straub, M., Weber, J., Slotta-Huspenina, J., Specht, K., Martignoni, M. E., Werner, A., Hein, R., H Busch, D., Peschel, C., Rad, R., Cox, J., Mann, M., and Krackhardt, A. M. (2016) Direct identification of clinically relevant neoepitopes presented on native human melanoma tissue by mass spectrometry. *Nat. Commun.* **7**, 13404
 75. Shapiro, I. E., and Bassani-Sternberg, M. (2023) The impact of immunopectidomics: From basic research to clinical implementation. *Semin. Immunol.* **66**, 101727
 76. Abelin, J. G., Keskin, D. B., Sarkizova, S., Hartigan, C. R., Zhang, W., Sidney, J., Stevens, J., Lane, W., Zhang, G. L., Eisenhaure, T. M., Clauser, K. R., Hacohen, N., Rooney, M. S., Carr, S. A., and Wu, C. J. (2017) Mass Spectrometry Profiling of HLA-Associated Peptidomes in Mono-allelic Cells Enables More Accurate Epitope Prediction. *Immunity.* **46**, 315–326
 77. Purcell, A. W., Ramarathnam, S. H., and Ternette, N. (2019) Mass spectrometry-based identification of MHC-bound peptides for immunopectidomics. *Nat. Protoc.* **14**, 1687–1707

78. Bassani-Sternberg, M., Pletscher-Frankild, S., Jensen, L. J., and Mann, M. (2015) Mass spectrometry of human leukocyte antigen class I peptidomes reveals strong effects of protein abundance and turnover on antigen presentation. *Mol. Cell. Proteomics*. **14**, 658–673
79. Yewdell, J. W. (2003) Immunology. Hide and seek in the peptidome. *Science*. **301**, 1334–1335
80. Blaha, D. T., Anderson, S. D., Yoakum, D. M., Hager, M. V., Zha, Y., Gajewski, T. F., and Kranz, D. M. (2019) High-Throughput Stability Screening of Neoantigen/HLA Complexes Improves Immunogenicity Predictions. *Cancer Immunol Res*. **7**, 50–61
81. Prevosto, C., Usmani, M. F., McDonald, S., Gumienny, A. M., Key, T., Goodman, R. S., Gaston, J. S. H., Deery, M. J., and Busch, R. (2016) Allele-Independent Turnover of Human Leukocyte Antigen (HLA) Class Ia Molecules. *PLoS One*. **11**, e0161011
82. Abelin, J. G., Bergstrom, E. J., Taylor, H. B., Rivera, K. D., Klaeger, S., Xu, C., Jackson White, C., Olive, M. E., Maynard, M., Harry Kane, M., Rachimi, S., Mani, D. R., Gillette, M. A., Clauser, K. R., Udeshi, N. D., and Carr, S. A. (2022) MONTE enables serial immunopeptidome, ubiquitylome, proteome, phosphoproteome, acetylome analyses of sample-limited tissues. *bioRxiv*. 10.1101/2021.06.22.449417
83. Boehm, K. M., Bhinder, B., Raja, V. J., Dephore, N., and Elemento, O. (2019) Predicting peptide presentation by major histocompatibility complex class I: an improved machine learning approach to the immunopeptidome. *BMC Bioinformatics*. 10.1186/s12859-018-2561-z
84. Abelin, J. G., Harjanto, D., Malloy, M., Suri, P., Colson, T., Goulding, S. P., Creech, A. L., Serrano, L. R., Nasir, G., Nasrullah, Y., McGann, C. D., Velez, D., Ting, Y. S., Poran, A., Rothenberg, D. A., Chhangawala, S., Rubinsteyn, A., Hammerbacher, J., Gaynor, R. B., Fritsch, E. F., Greshock, J., Oslund, R. C., Barthelme, D., Addona, T. A., Arieta, C. M., and Rooney, M. S. (2019) Defining HLA-II Ligand Processing and Binding Rules with Mass Spectrometry Enhances Cancer Epitope Prediction. *Immunity*. **51**, 766–779.e17
85. Alvarez, B., Reynisson, B., Barra, C., Buus, S., Ternette, N., Connelley, T., Andreatta, M., and Nielsen, M. (2019) NNAlign_MA; MHC Peptidome Deconvolution for Accurate MHC Binding Motif Characterization and Improved T-cell Epitope Predictions. *Mol. Cell. Proteomics*. **18**, 2459–2477
86. O'Donnell, T. J., Rubinsteyn, A., and Laserson, U. (2020) MHCflurry 2.0: Improved Pan-Allele Prediction of MHC Class I-Presented Peptides by Incorporating Antigen Processing. *Cell Syst*. **11**, 418–419
87. Jurtz, V., Paul, S., Andreatta, M., Marcatili, P., Peters, B., and Nielsen, M. (2017) NetMHCpan-4.0: Improved peptide-MHC class I interaction predictions integrating eluted ligand and peptide binding affinity data. *J. Immunol*. **199**, 3360–3368
88. Sarkizova, S., Klaeger, S., Le, P. M., Li, L. W., Oliveira, G., Keshishian, H., Hartigan, C. R., Zhang, W., Braun, D. A., Ligon, K. L., Bachiredy, P., Zervantonakis, I. K., Rosenbluth, J. M., Ouspenskaia, T., Law, T., Justesen, S., Stevens, J., Lane, W. J., Eisenhaure, T., Zhang, G. L., Clauser, K. R., Hacohen, N., Carr, S. A., Wu, C. J., and Keskin, D. B. (2020) A large peptidome dataset improves HLA class I epitope prediction across most of the human population. *Nature Biotechnology*. **38**, 199–209
89. Taylor, H. B., Klaeger, S., Clauser, K. R., Sarkizova, S., Weingarten-Gabbay, S., Graham, D. B., Carr, S. A., and Abelin, J. G. (2021) MS-Based HLA-II Peptidomics Combined With Multiomics Will Aid the Development of Future Immunotherapies. *Mol. Cell. Proteomics*. **20**, 100116
90. Chen, B., Khodadoust, M. S., Olsson, N., Wagar, L. E., Fast, E., Liu, C. L., Muftuoglu, Y., Swarder, B. J., Diehn, M., Levy, R., Davis, M. M., Elias, J. E., Altman, R. B., and Alizadeh, A. A. (2019) Predicting HLA class II antigen presentation through integrated deep learning. *Nat. Biotechnol*. **37**, 1332–1343
91. Racle, J., Michaux, J., Rockinger, G. A., Arnaud, M., Bobisse, S., Chong, C., Guillaume, P., Coukos, G., Harari, A., Jandus, C., Bassani-Sternberg, M., and Gfeller, D. (2019) Robust prediction of HLA class II epitopes by deep motif deconvolution of immunopeptidomes. *Nat. Biotechnol*. **37**, 1283–1286
92. Shao, X. M., Bhattacharya, R., Huang, J., Sivakumar, I. K. A., Tokheim, C., Zheng, L., Hirsch, D.,

- Kaminow, B., Omdahl, A., Bonsack, M., Riemer, A. B., Velculescu, V. E., Anagnostou, V., Pagel, K. A., and Karchin, R. (2020) High-Throughput Prediction of MHC Class I and II Neoantigens with MHCnuggets. *Cancer Immunol Res.* **8**, 396–408
93. Lozano-Rabella, M., Garcia-Garijo, A., Palomero, J., Yuste-Estevanez, A., Erhard, F., Martín-Liberal, J., de Olza, M. O., Matos, I., Gartner, J. J., Ghosh, M., Canals, F., Vidal, A., Piulats, J. M., Matias-Guiu, X., Braña, I., Muñoz-Couselo, E., Garraza, E., Schlosser, A., and Gros, A. (2022) Immunogenicity of non-canonical HLA-I tumor ligands identified through proteogenomics. *bioRxiv*. 10.1101/2022.11.07.514886
94. Erhard, F., Dölken, L., Schilling, B., and Schlosser, A. (2020) Identification of the Cryptic HLA-I Immunopeptidome. *Cancer Immunol Res.* **8**, 1018–1026
95. Lichti, C. F., Vigneron, N., Clauser, K. R., Van den Eynde, B. J., and Bassani-Sternberg, M. (2022) Navigating Critical Challenges Associated with Immunopeptidomics-Based Detection of Proteasomal Spliced Peptide Candidates. *Cancer Immunol Res.* **10**, 275–284
96. Gessulat, S., Schmidt, T., Zolg, D. P., Samaras, P., Schnatbaum, K., Zerweck, J., Knaute, T., Rechenberger, J., Delanghe, B., Huhmer, A., Reimer, U., Ehrlich, H.-C., Aiche, S., Kuster, B., and Wilhelm, M. (2019) Prosit: proteome-wide prediction of peptide tandem mass spectra by deep learning. *Nat. Methods.* **16**, 509–518
97. Wilhelm, M., Zolg, D. P., Graber, M., Gessulat, S., Schmidt, T., Schnatbaum, K., Schwencke-Westphal, C., Seifert, P., de Andrade Krätzig, N., Zerweck, J., Knaute, T., Bräunlein, E., Samaras, P., Lautenbacher, L., Klaeger, S., Wenschuh, H., Rad, R., Delanghe, B., Huhmer, A., Carr, S. A., Clauser, K. R., Krackhardt, A. M., Reimer, U., and Kuster, B. (2021) Deep learning boosts sensitivity of mass spectrometry-based immunopeptidomics. *Nat. Commun.* **12**, 3346
98. Declercq, A., Bouwmeester, R., Chiva, C., Sabidó, E., Hirschler, A., Carapito, C., Martens, L., Degroove, S., and Gabriels, R. (2023) Updated MS²PIP web server supports cutting-edge proteomics applications. *Nucleic Acids Res.* 10.1093/nar/gkad335
99. Bouwmeester, R., Gabriels, R., Hulstaert, N., Martens, L., and Degroove, S. (2021) DeepLC can predict retention times for peptides that carry as-yet unseen modifications. *Nat. Methods.* **18**, 1363–1369
100. Li, K., Jain, A., Malovannaya, A., Wen, B., and Zhang, B. (2020) DeepRescore: Leveraging Deep Learning to Improve Peptide Identification in Immunopeptidomics. *Proteomics.* **20**, e1900334
101. Deutsch, E. W., Lane, L., Overall, C. M., Bandeira, N., Baker, M. S., Pineau, C., Moritz, R. L., Corrales, F., Orchard, S., Van Eyk, J. E., Paik, Y.-K., Weintraub, S. T., Vandenbrouck, Y., and Omenn, G. S. (2019) Human Proteome Project Mass Spectrometry Data Interpretation Guidelines 3.0. *J. Proteome Res.* **18**, 4108–4116
102. Adhikari, S., Nice, E. C., Deutsch, E. W., Lane, L., Omenn, G. S., Pennington, S. R., Paik, Y.-K., Overall, C. M., Corrales, F. J., Cristea, I. M., Van Eyk, J. E., Uhlén, M., Lindskog, C., Chan, D. W., Bairoch, A., Waddington, J. C., Justice, J. L., LaBaer, J., Rodriguez, H., He, F., Kostrzewa, M., Ping, P., Gundry, R. L., Stewart, P., Srivastava, S., Srivastava, S., Nogueira, F. C. S., Domont, G. B., Vandenbrouck, Y., Lam, M. P. Y., Wennersten, S., Vizcaino, J. A., Wilkins, M., Schwenk, J. M., Lundberg, E., Bandeira, N., Marko-Varga, G., Weintraub, S. T., Pineau, C., Kusebauch, U., Moritz, R. L., Ahn, S. B., Palmblad, M., Snyder, M. P., Aebersold, R., and Baker, M. S. (2020) A high-stringency blueprint of the human proteome. *Nat. Commun.* **11**, 5301
103. Deutsch, E. W., Overall, C. M., Van Eyk, J. E., Baker, M. S., Paik, Y.-K., Weintraub, S. T., Lane, L., Martens, L., Vandenbrouck, Y., Kusebauch, U., Hancock, W. S., Hermjakob, H., Aebersold, R., Moritz, R. L., and Omenn, G. S. (2016) Human Proteome Project Mass Spectrometry Data Interpretation Guidelines 2.1. *J. Proteome Res.* **15**, 3961–3970
104. Deutsch, E. W., Perez-Riverol, Y., Carver, J., Kawano, S., Mendoza, L., Van Den Bossche, T., Gabriels, R., Binz, P.-A., Pullman, B., Sun, Z., Shofstahl, J., Bittremieux, W., Mak, T. D., Klein, J., Zhu, Y., Lam,

- H., Vizcaíno, J. A., and Bandeira, N. (2021) Universal Spectrum Identifier for mass spectra. *Nat. Methods*. **18**, 768–770
105. Kim, M.-S., Pinto, S. M., Getnet, D., Nirujogi, R. S., Manda, S. S., Chaerkady, R., Madugundu, A. K., Kelkar, D. S., Isserlin, R., Jain, S., Thomas, J. K., Muthusamy, B., Leal-Rojas, P., Kumar, P., Sahasrabudhe, N. A., Balakrishnan, L., Advani, J., George, B., Renuse, S., Selvan, L. D. N., Patil, A. H., Nanjappa, V., Radhakrishnan, A., Prasad, S., Subbannayya, T., Raju, R., Kumar, M., Sreenivasamurthy, S. K., Marimuthu, A., Sathe, G. J., Chavan, S., Datta, K. K., Subbannayya, Y., Sahu, A., Yelamanchi, S. D., Jayaram, S., Rajagopalan, P., Sharma, J., Murthy, K. R., Syed, N., Goel, R., Khan, A. A., Ahmad, S., Dey, G., Mudgal, K., Chatterjee, A., Huang, T.-C., Zhong, J., Wu, X., Shaw, P. G., Freed, D., Zahari, M. S., Mukherjee, K. K., Shankar, S., Mahadevan, A., Lam, H., Mitchell, C. J., Shankar, S. K., Satishchandra, P., Schroeder, J. T., Sirdeshmukh, R., Maitra, A., Leach, S. D., Drake, C. G., Halushka, M. K., Prasad, T. S. K., Hruban, R. H., Kerr, C. L., Bader, G. D., Iacobuzio-Donahue, C. A., Gowda, H., and Pandey, A. (2014) A draft map of the human proteome. *Nature*. **509**, 575–581
106. Oyama, M., Kozuka-Hata, H., Suzuki, Y., Semba, K., Yamamoto, T., and Sugano, S. (2007) Diversity of Translation Start Sites May Define Increased Complexity of the Human Short ORFeome* S. *Mol. Cell. Proteomics*. **6**, 1000–1006
107. Volders, P. J., Verheggen, K., Menschaert, G., Vandepoele, K., Martens, L., Vandesompele, J., and Mestdag, P. (2015) An update on LNCipedia: a database for annotated human lncRNA sequences. *Nucleic Acids Res.* **43**, 4363–4364
108. Iyer, M. K., Niknafs, Y. S., Malik, R., Singhal, U., Sahu, A., Hosono, Y., Barrette, T. R., Prensner, J. R., Evans, J. R., Zhao, S., Poliakov, A., Cao, X., Dhanasekaran, S. M., Wu, Y.-M., Robinson, D. R., Beer, D. G., Feng, F. Y., Iyer, H. K., and Chinnaiyan, A. M. (2015) The landscape of long noncoding RNAs in the human transcriptome. *Nat. Genet.* **47**, 199–208
109. Wacholder, A., and Carvunis, A.-R. Rare detection of noncanonical proteins in yeast mass spectrometry studies. [10.1101/2023.03.09.531963](https://doi.org/10.1101/2023.03.09.531963)
110. Verheggen, K., Volders, P.-J., Mestdag, P., Menschaert, G., Van Damme, P., Gevaert, K., Martens, L., and Vandesompele, J. (2017) Noncoding after All: Biases in Proteomics Data Do Not Explain Observed Absence of lncRNA Translation Products. *J. Proteome Res.* **16**, 2508–2515
111. Bogaert, A., Fijalkowska, D., Staes, A., Van de Steene, T., Demol, H., and Gevaert, K. (2022) Limited Evidence for Protein Products of Noncoding Transcripts in the HEK293T Cellular Cytosol. *Mol. Cell. Proteomics*. **21**, 100264
112. Cassidy, L., Kaulich, P. T., and Tholey, A. (2023) Proteoforms expand the world of microproteins and short open reading frame-encoded peptides. *iScience*. **26**, 106069
113. Kesner, J. S., Chen, Z., Shi, P., Aparicio, A. O., Murphy, M. R., Guo, Y., Trehan, A., Lipponen, J. E., Recinos, Y., Myeku, N., and Wu, X. (2023) Noncoding translation mitigation. *Nature*. [10.1038/s41586-023-05946-4](https://doi.org/10.1038/s41586-023-05946-4)
114. Fritsch, C., Herrmann, A., Nothnagel, M., Szafranski, K., Huse, K., Schumann, F., Schreiber, S., Platzer, M., Krawczak, M., Hampe, J., and Brosch, M. (2012) Genome-wide search for novel human uORFs and N-terminal protein extensions using ribosomal footprinting. *Genome Res.* **22**, 2208–2218
115. Li, Y., Zhou, H., Chen, X., Zheng, Y., Kang, Q., Hao, D., Zhang, L., Song, T., Luo, H., Hao, Y., Chen, R., Zhang, P., and He, S. (2021) SmProt: A Reliable Repository with Comprehensive Annotation of Small Proteins Identified from Ribosome Profiling. *Genomics Proteomics Bioinformatics*. **19**, 602–610
116. Olexiouk, V., Crappé, J., Verbruggen, S., Verheggen, K., Martens, L., and Menschaert, G. (2016) sORFs.org: a repository of small ORFs identified by ribosome profiling. *Nucleic Acids Res.* **44**, D324–9
117. Wang, H., Yang, L., Wang, Y., Chen, L., Li, H., and Xie, Z. (2019) RPFdb v2. 0: an updated database for genome-wide information of translated mRNA generated from ribosome profiling. *Nucleic Acids*

Res. **47**, D230–D234

118. Xie, S.-Q., Nie, P., Wang, Y., Wang, H., Li, H., Yang, Z., Liu, Y., Ren, J., and Xie, Z. (2016) RPFdb: a database for genome wide information of translated mRNA generated from ribosome profiling. *Nucleic Acids Res.* **44**, D254–8
119. Ji, X., Cui, C., and Cui, Q. (2020) smORFunction: a tool for predicting functions of small open reading frames and microproteins. *BMC Bioinformatics.* **21**, 455
120. Brunet, M. A., Brunelle, M., Lucier, J.-F., Delcourt, V., Levesque, M., Grenier, F., Samandi, S., Leblanc, S., Aguilar, J.-D., Dufour, P., Jacques, J.-F., Fournier, I., Ouangraoua, A., Scott, M. S., Boisvert, F.-M., and Roucou, X. (2019) OpenProt: a more comprehensive guide to explore eukaryotic coding potential and proteomes. *Nucleic Acids Res.* **47**, D403–D410
121. Brunet, M. A., Lucier, J.-F., Levesque, M., Leblanc, S., Jacques, J.-F., Al-Saedi, H. R. H., Guilloy, N., Grenier, F., Avino, M., Fournier, I., Salzet, M., Ouangraoua, A., Scott, M. S., Boisvert, F.-M., and Roucou, X. (2021) OpenProt 2021: deeper functional annotation of the coding potential of eukaryotic genomes. *Nucleic Acids Res.* **49**, D380–D388
122. Duffy, E. E., Finander, B., Choi, G., Carter, A. C., Pritisanac, I., Alam, A., Luria, V., Karger, A., Phu, W., Sherman, M. A., Assad, E. G., Khitun, A., Crouch, E. E., Ganesh, S., Berger, B., Sestan, N., O'Donnell-Luria, A., Huang, E., Griffith, E. C., Forman-Kay, J. D., Moses, A. M., Kalish, B. T., and Greenberg, M. E. (2022) Developmental Dynamics of RNA Translation in the Human Brain. *Nat Neurosci.* **25**, 1353–1365
123. Smirnova, V. V., Shestakova, E. D., Nogina, D. S., Mishchenko, P. A., Prikazchikova, T. A., Zatsepin, T. S., Kulakovskiy, I. V., Shatsky, I. N., and Terenin, I. M. (2022) Ribosomal leaky scanning through a translated uORF requires eIF4G2. *Nucleic Acids Res.* **50**, 1111–1127
124. Andreev, D. E., Loughran, G., Fedorova, A. D., Mikhaylova, M. S., Shatsky, I. N., and Baranov, P. V. (2022) Non-AUG translation initiation in mammals. *Genome Biol.* **23**, 111
125. Stacey, S. N., Jordan, D., Williamson, A. J., Brown, M., Coote, J. H., and Arrand, J. R. (2000) Leaky scanning is the predominant mechanism for translation of human papillomavirus type 16 E7 oncoprotein from E6/E7 bicistronic mRNA. *J. Virol.* **74**, 7284–7297
126. Duss, O., Stepanyuk, G. A., Puglisi, J. D., and Williamson, J. R. (2019) Transient Protein-RNA Interactions Guide Nascent Ribosomal RNA Folding. *Cell.* **179**, 1357–1369.e16
127. Karamyshev, A. L., and Karamysheva, Z. N. (2018) Lost in Translation: Ribosome-Associated mRNA and Protein Quality Controls. *Front. Genet.* **9**, 431
128. Gelhausen, R., Müller, T., Svensson, S. L., Alkhnbashi, O. S., Sharma, C. M., Eggenhofer, F., and Backofen, R. (2022) RiboReport - benchmarking tools for ribosome profiling-based identification of open reading frames in bacteria. *Brief. Bioinform.* 10.1093/bib/bbab549
129. Kiniry, S. J., Michel, A. M., and Baranov, P. V. (2020) Computational methods for ribosome profiling data analysis. *WIREs RNA.* 10.1002/wrna.1577
130. Lei, T., Chang, Y., Yao, C., and Zhang, H. (2022) A systematic evaluation revealed that detecting translated non-canonical ORFs from ribosome profiling data remains challenging. *bioRxiv.* 10.1101/2022.12.11.520003
131. Blackwood, E. M., Lugo, T. G., Kretzner, L., King, M. W., Street, A. J., Witte, O. N., and Eisenman, R. N. (1994) Functional analysis of the AUG- and CUG-initiated forms of the c-Myc protein. *Mol. Biol. Cell.* **5**, 597–609
132. Prats, H., Kaghad, M., Prats, A. C., Klagsbrun, M., Lélías, J. M., Liauzun, P., Chalon, P., Tauber, J. P., Amalric, F., and Smith, J. A. (1989) High molecular mass forms of basic fibroblast growth factor are initiated by alternative CUG codons. *Proc. Natl. Acad. Sci. U. S. A.* **86**, 1836–1840
133. Cao, X., and Slavoff, S. A. (2020) Non-AUG start codons: Expanding and regulating the small and alternative ORFeome. *Exp. Cell Res.* **391**, 111973
134. Ingolia, N. T., Lareau, L. F., and Weissman, J. S. (2011) Ribosome profiling of mouse embryonic stem

- cells reveals the complexity and dynamics of mammalian proteomes. *Cell*. **147**, 789–802
135. Lee, S., Liu, B., Lee, S., Huang, S.-X., Shen, B., and Qian, S.-B. (2012) Global mapping of translation initiation sites in mammalian cells at single-nucleotide resolution. *Proc. Natl. Acad. Sci. U. S. A.* **109**, E2424–32
136. Champagne, J., Pataskar, A., Blommaert, N., Nagel, R., Wernaart, D., Ramalho, S., Kenski, J., Bleijerveld, O. B., Zaal, E. A., Berkers, C. R., Altelaar, M., Peeper, D. S., Faller, W. J., and Agami, R. (2021) Oncogene-dependent sloppiness in mRNA translation. *Mol. Cell*. **81**, 4709–4721.e9
137. Janssen, J. W., Vaandrager, J. W., Heuser, T., Jauch, A., Kluin, P. M., Geelen, E., Bergsagel, P. L., Kuehl, W. M., Drexler, H. G., Otsuki, T., Bartram, C. R., and Schuurin, E. (2000) Concurrent activation of a novel putative transforming gene, myeov, and cyclin D1 in a subset of multiple myeloma cell lines with t(11;14)(q13;q32). *Blood*. **95**, 2691–2698
138. Dolstra, H., Fredrix, H., Preijers, F., Goulmy, E., Figdor, C. G., de Witte, T. M., and van de Wiel-van Kemenade, E. (1997) Recognition of a B cell leukemia-associated minor histocompatibility antigen by CTL. *J. Immunol.* **158**, 560–565
139. Ruiz-Orera, J., Villanueva-Cañás, J. L., and Albà, M. M. (2020) Evolution of new proteins from translated sORFs in long non-coding RNAs. *Exp. Cell Res.* **391**, 111940
140. Vakirlis, N., Vance, Z., Duggan, K. M., and McLysaght, A. (2022) De novo birth of functional microproteins in the human lineage. *Cell Rep.* **41**, 111808
141. Broeils, L. A., Ruiz-Orera, J., Snel, B., Hubner, N., and van Heesch, S. (2023) Evolution and implications of de novo genes in humans. *Nat Ecol Evol*. 10.1038/s41559-023-02014-y
142. Erady, C., Boxall, A., Puntambekar, S., Suhas Jagannathan, N., Chauhan, R., Chong, D., Meena, N., Kulkarni, A., Kasabe, B., Bhayankaram, K. P., Umrana, Y., Andreani, A., Nel, J., Wayland, M. T., Pina, C., Lilley, K. S., and Prabakaran, S. (2021) Pan-cancer analysis of transcripts encoding novel open-reading frames (nORFs) and their potential biological functions. *npj Genomic Medicine*. 10.1038/s41525-020-00167-4
143. Na, Z., Luo, Y., Cui, D. S., Khitun, A., Smelyansky, S., Loria, J. P., and Slavoff, S. A. (2021) Phosphorylation of a Human Microprotein Promotes Dissociation of Biomolecular Condensates. *J. Am. Chem. Soc.* **143**, 12675–12687
144. D’Lima, N. G., Ma, J., Winkler, L., Chu, Q., Loh, K. H., Corpuz, E. O., Budnik, B. A., Lykke-Andersen, J., Saghatelian, A., and Slavoff, S. A. (2017) A human microprotein that interacts with the mRNA decapping complex. *Nat. Chem. Biol.* **13**, 174–180
145. Ruiz-Orera, J., Messeguer, X., Subirana, J. A., and Alba, M. M. (2014) Long non-coding RNAs as a source of new peptides. *Elife*. **3**, e03523
146. Schlesinger, D., and Elsässer, S. J. (2022) Revisiting sORFs: overcoming challenges to identify and characterize functional microproteins. *FEBS J.* **289**, 53–74
147. Vakirlis, N., Acar, O., Hsu, B., Castilho Coelho, N., Van Oss, S. B., Wacholder, A., Medetgul-Ernar, K., Bowman, R. W., 2nd, Hines, C. P., Iannotta, J., Parikh, S. B., McLysaght, A., Camacho, C. J., O’Donnell, A. F., Ideker, T., and Carvunis, A.-R. (2020) De novo emergence of adaptive membrane proteins from thymine-rich genomic sequences. *Nat. Commun.* **11**, 781
148. Heames, B., Buchel, F., Aubel, M., Tretyachenko, V., Loginov, D., Novák, P., Lange, A., Bornberg-Bauer, E., and Hlouchová, K. (2023) Experimental characterization of de novo proteins and their unevolved random-sequence counterparts. *Nat Ecol Evol*. **7**, 570–580
149. Kesner, J. S., Chen, Z., Aparicio, A. A., and Wu, X. (2022) A unified model for the surveillance of translation in diverse noncoding sequences. *bioRxiv*. 10.1101/2022.07.20.500724
150. Paes, W., Leonov, G., Partridge, T., Chikata, T., Murakoshi, H., Frangou, A., Brackenridge, S., Nicastrì, A., Smith, A. G., Learn, G. H., Li, Y., Parker, R., Oka, S., Pellegrino, P., Williams, I., Haynes, B. F., McMichael, A. J., Shaw, G. M., Hahn, B. H., Takiguchi, M., Ternette, N., and Borrow, P. (2019) Contribution of proteasome-catalyzed peptide -splicing to viral targeting by CD8 T cells in HIV-1

- infection. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 24748–24759
151. Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal*. **17**, 10–12
152. Langmead, B., and Salzberg, S. L. (2012) Fast gapped-read alignment with Bowtie 2. *Nat. Methods*. **9**, 357–359
153. Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and Gingeras, T. R. (2013) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. **29**, 15–21
154. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*. **25**, 2078–2079
155. Calviello, L., Sydow, D., Harnett, D., and Ohler, U. (2019) Ribo-seQC: comprehensive analysis of cytoplasmic and organellar ribosome profiling data. *bioRxiv*. 10.1101/601468

Tables

Table 1: Features and characteristics of methods to detect non-canonical ORF translation

	Molecule detected	Digestion step?	Target size of analyte	Number of annotated CDSs detected	Number of ORFs detected	Strengths	Weaknesses
Ribo-seq	RNA bound within ribosomes	RNase and DNase	28 - 30 nt	10,000 - 13,000	2,000 - 200,000	<ul style="list-style-type: none"> -Genome-wide -No bias due to trypsin -Detects small and large CDSs -Nucleotide-level precision -Defines exact reading frame of ORF 	<ul style="list-style-type: none"> -Does not detect proteins directly -Cannot detect PTMs -Cannot inform post-translational protein regulation -Analysis pipelines may be discordant
LC-MS/MS	Tryptic peptides	Trypsin	8 - 25 amino acids	9,000 - 11,000	10s to 100s	<ul style="list-style-type: none"> -Direct protein detection -Informs protein abundance -May detect PTMs -Proteome-wide 	<ul style="list-style-type: none"> -High false-positive rate without Ribo-seq -Biased against small proteins -Trypsin may bias protein representation -Does not provide nucleotide-level precision
HLA Immuno-peptidomics	HLA-presented peptide antigens	None	8 - 12 amino acids	8,000 - 10,000	1000 - 5000	<ul style="list-style-type: none"> -Direct protein detection -Enrichment for low-abundance, strong binders -Proteome-wide -Can detect unstable translations -Does not require tryptic sites 	<ul style="list-style-type: none"> -Does not inform protein stability -Does not indicate intracellular abundance -HLA allele expression limits peptide representation -Does not provide nucleotide-level precision

Table 2: A proposed framework to standardize levels of evidence of non-canonical ORFs

Tier	Required Supporting Evidence	Standardized Outcome
<i>Tier 1A</i>	<ul style="list-style-type: none"> Whole proteome LC-MS/MS (≥ 2 peptides according to HUPO HPP criteria) Ribo-seq* 	<i>"Protein candidate"</i> Consider discussing research findings with genome annotation databases for possible annotation.
<i>Tier 1B</i>	<ul style="list-style-type: none"> HLA immunopeptidomics MS (≥ 2 observations; multiple high confidence peptides from multiple distinct samples) Ribo-seq* 	<i>"Presented"</i>
<i>Tier 2A</i>	<ul style="list-style-type: none"> Whole proteome LC-MS/MS (≥ 2 peptides not satisfying HUPO HPP spacing criteria) Whole proteome LC-MS/MS (1 peptide) Ribo-seq* 	<i>"Detected"</i>
<i>Tier 2B</i>	<ul style="list-style-type: none"> HLA immunopeptidomics MS (1 observation) Ribo-seq* 	<i>"Detected"</i>
<i>Tier 3</i>	<ul style="list-style-type: none"> Any HLA immunopeptidomics or whole proteome LC-MS/MS evidence without Ribo-seq* evidence 	<i>"Putative", consider alternative sources.</i>
<i>Tier 4</i>	<ul style="list-style-type: none"> Ribo-seq* evidence without any proteomic evidence 	<i>"Ribo-seq ORF"</i>
<i>Tier 5</i>	<ul style="list-style-type: none"> <i>In silico</i> prediction of an ORF on an expressed transcript without any Ribo-seq* or proteomic evidence 	<i>"Predicted"</i>

*from credible Ribo-seq data with quality metrics meeting the guidelines suggested in this manuscript. Ribo-seq need not be performed on aliquots of the same samples analyzed by proteomics.

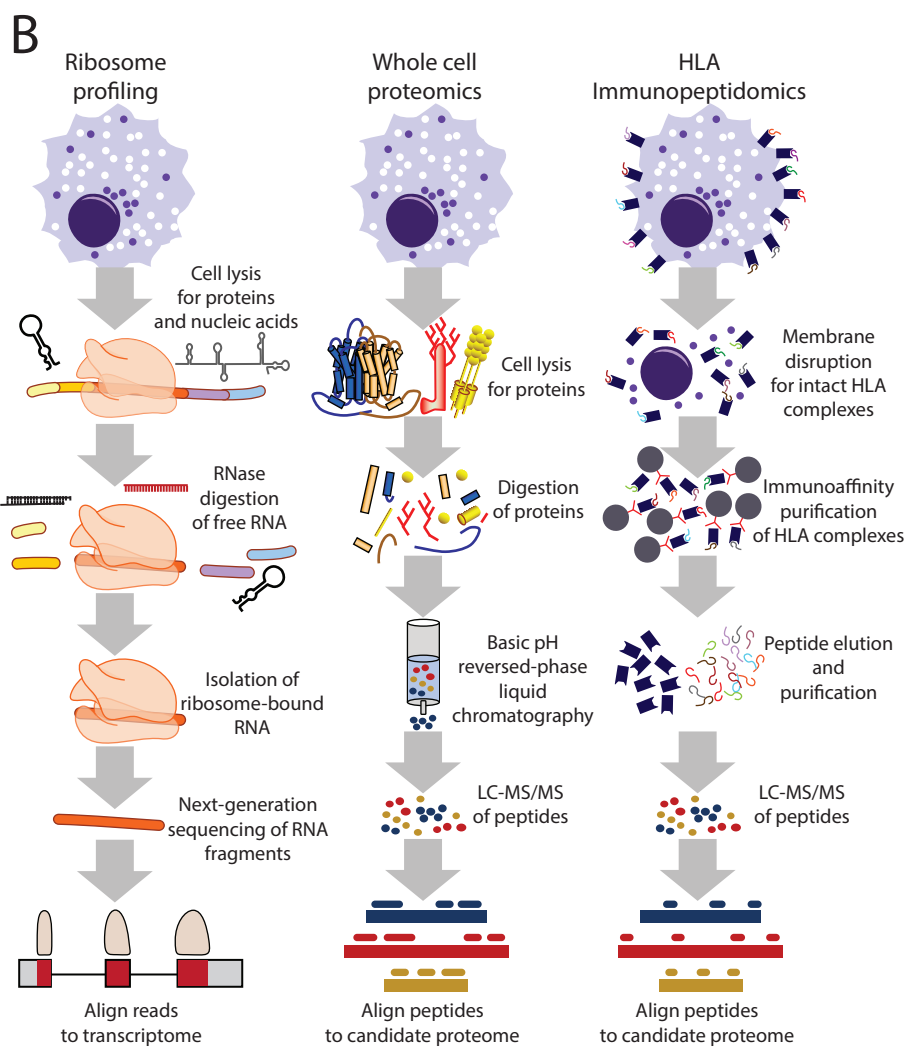
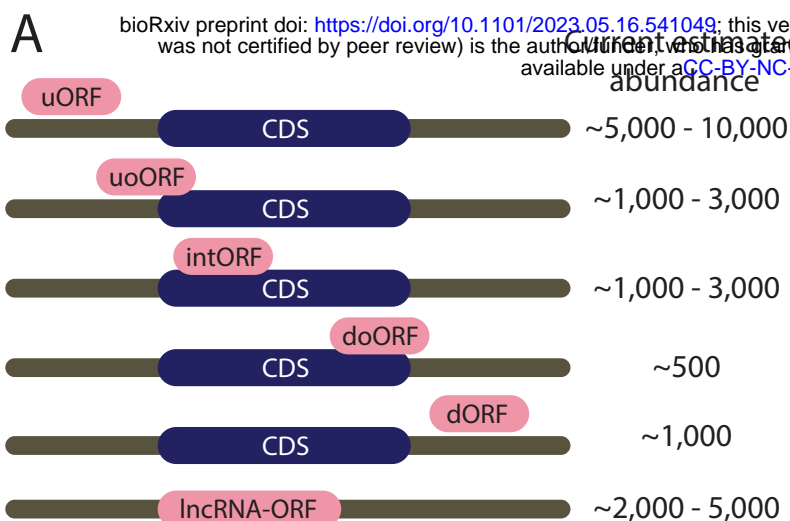
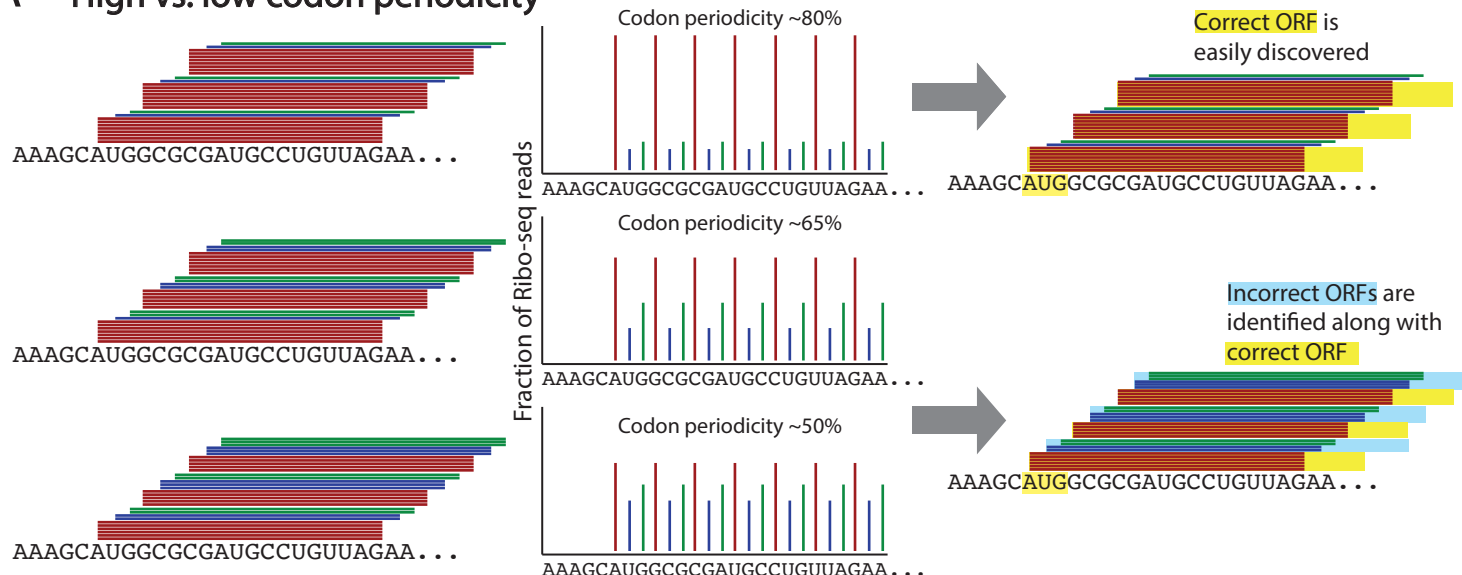


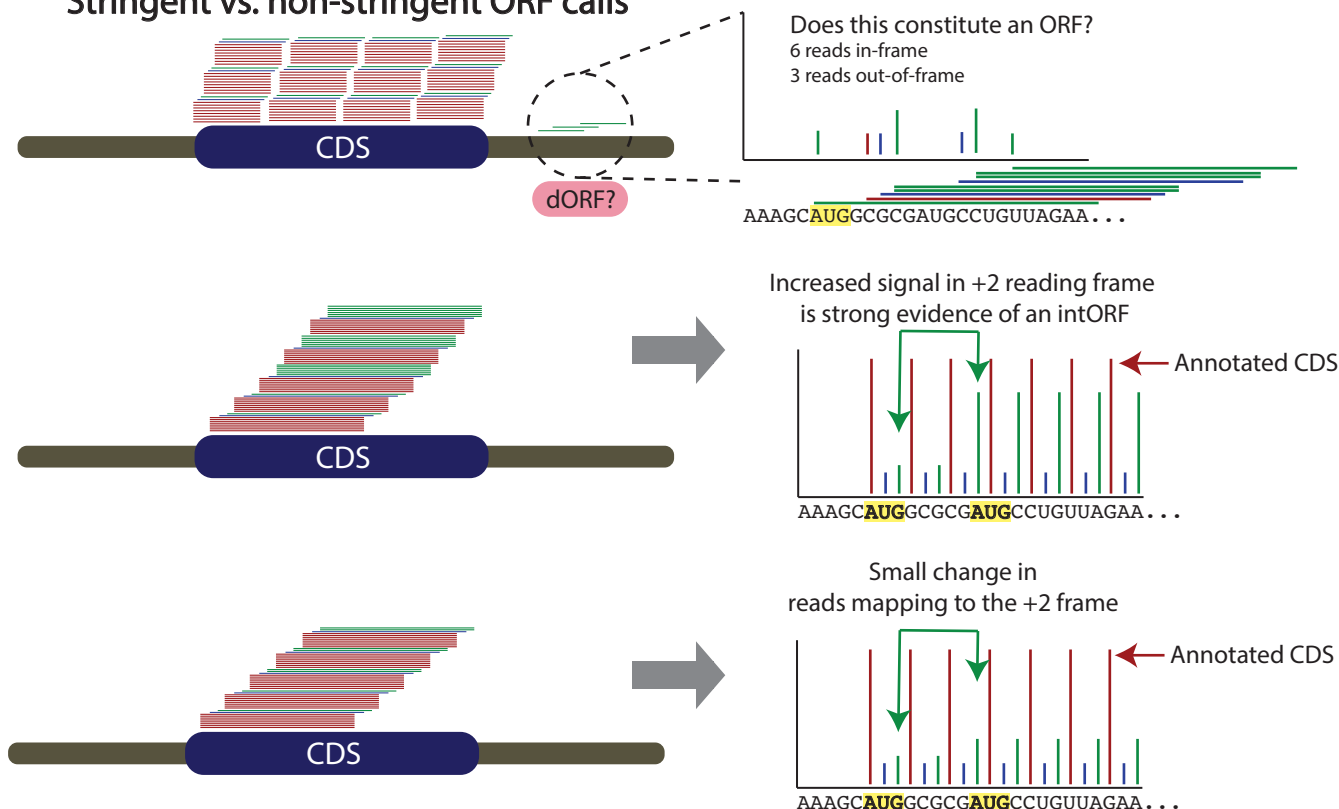
Figure 1: An overview of non-canonical ORF types and detection methods

- A) A schematic illustrating the standardized names of non-canonical ORF types, their relationship to known mRNAs, and current estimations of their abundance. CDS, protein-coding sequence; uORF, upstream open reading frame (ORF); uoORF, upstream overlapping ORF; intORF, internal ORF; doORF, downstream overlapping ORF; dORF, downstream ORF; lncRNA-ORF, ORF residing within an annotated lncRNA.
- B) Generalized workflows for ribosome profiling (Ribo-seq), tryptic whole cell mass spectrometry, and HLA immunopeptidomics. The schematic indicates general properties of sample preparation for these data types.

A High vs. low codon periodicity



B Stringent vs. non-stringent ORF calls



C Start site calling

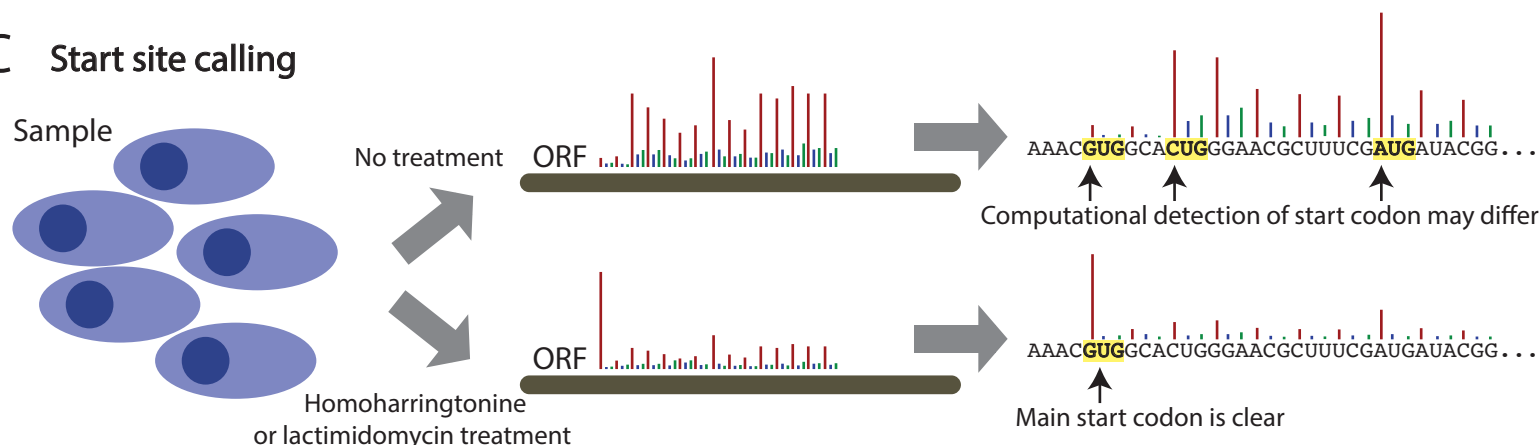
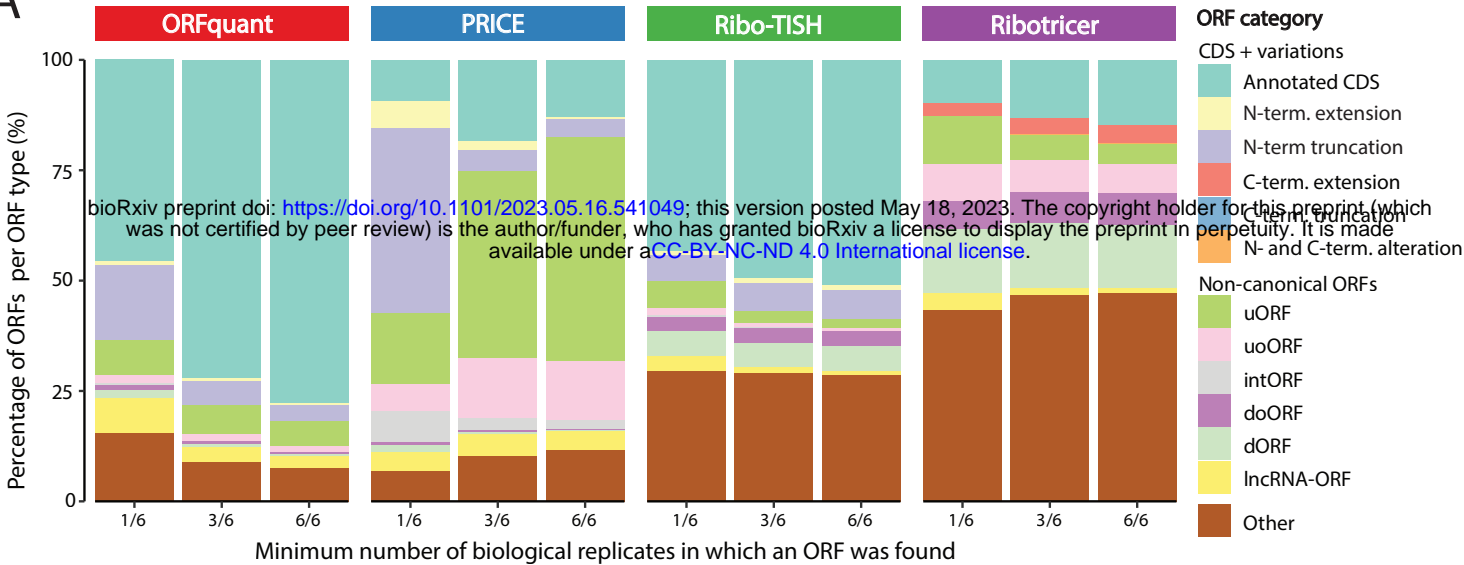


Figure 2: Quality metrics of Ribo-seq and stringency of ORF calling

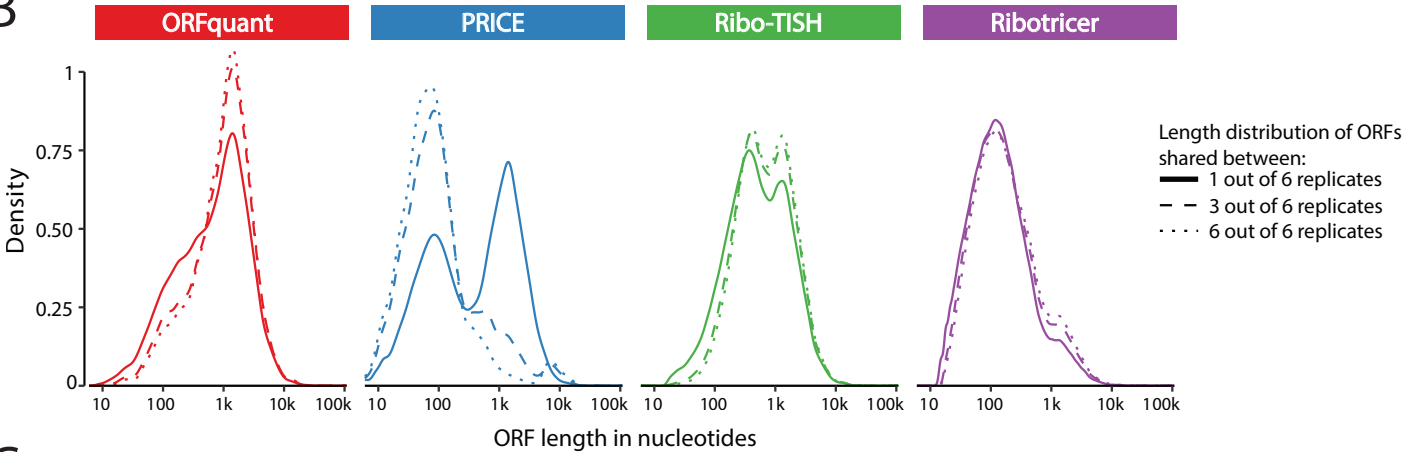
- A) An illustration showing codon periodicity as a central metric of Ribo-seq library generation. Three illustrations indicate high-quality, borderline, and poor-quality Ribo-seq libraries.
- B) An illustration representing high-stringency and low-stringency ORF calling. In the top case, a small number of reads map the the 3'UTR of an annotated mRNA, and only two-thirds of those 3'UTR reads support the same reading frame of a potential dORF nomination. In the middle and bottom cases, a potential intORF has varying read support evidence. The middle case shows clear evidence of an intORF by a large increase in reads mapping to the +2 reading frame midway through the CDS. In the bottom case, there is a smaller change in the reads mapping to the +2 reading frame.
- C) Use of ribosome-stalling drug treatments to clarify translational start sites. Cultured cells are treated with homoharringtonine or lactimidomycin to stall ribosomes at the main translational start site of a given ORF, leading to a clearer resolution of the specific start codon.

FIGURE 3

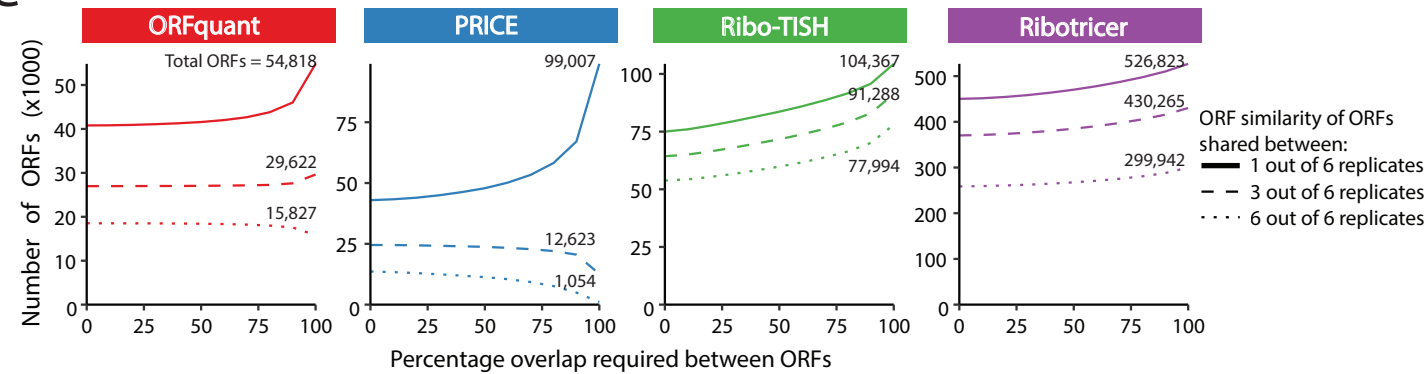
A



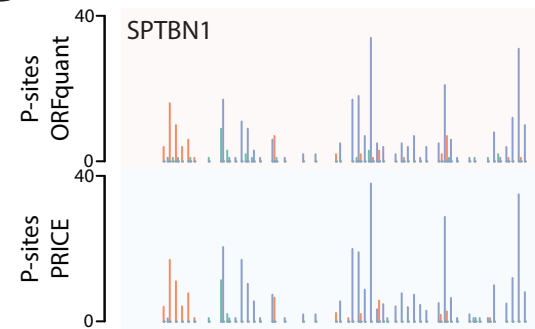
B



C



D



E

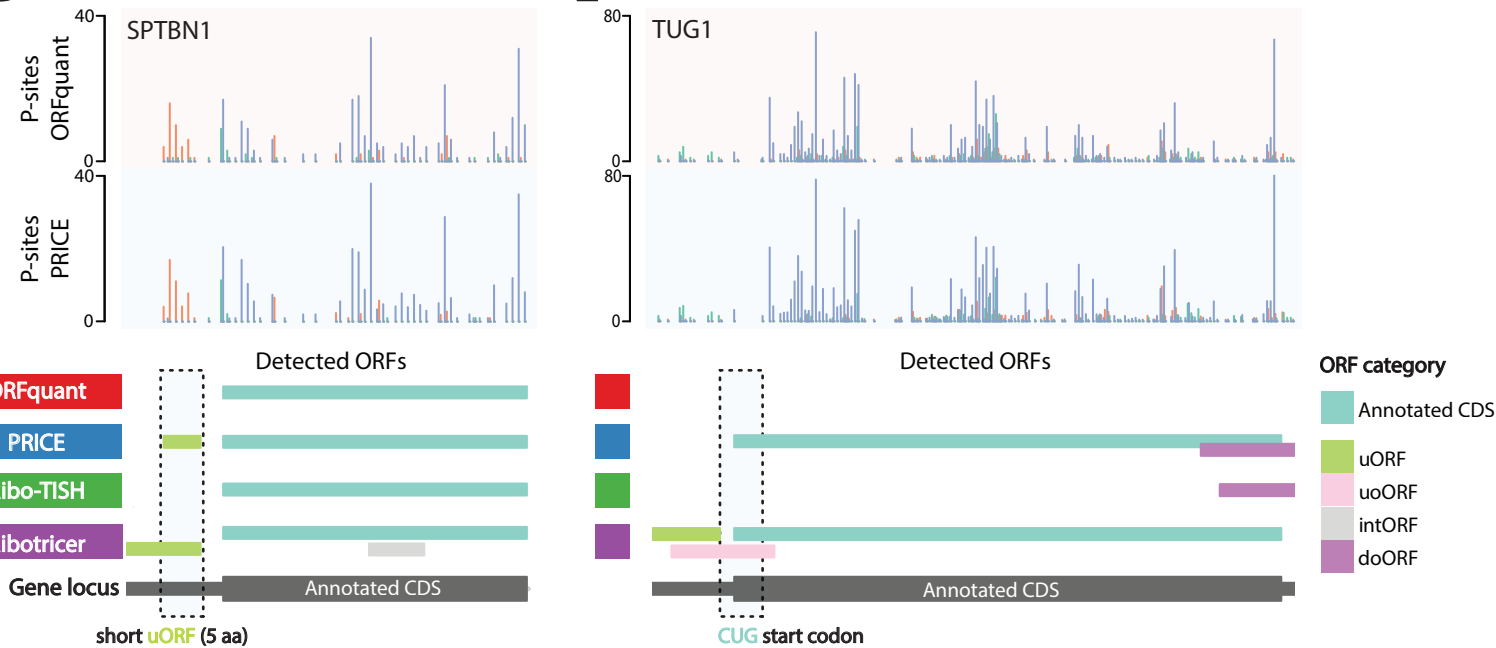


Figure 3: ORF callers have different specialties and variable performance

- A) Stacked bar plot displaying all detected ORF categories per ORF caller. For each, the percentage of unique ORFs shared between at least one, three, or six replicates is shown. Please note that these are relative contributions to the total number of ORFs. The absolute numbers of ORF identifications can be inferred from **Figure 3C**.
- B) Density plots displaying the distribution of ORF lengths in nucleotides (excluding the stop codon) for unique ORFs shared between at least one, three, or six replicates.
- C) Line graphs showing the numbers of unique ORFs detected by each tool shared between at least one, three or six replicates. The x-axis denotes the percentage of overlap used to consider two ORFs being similar or not, with 100% overlap meaning that the detected ORF was fully identical between [x] number of replicates. Please note that the total numbers of ORFs detected per algorithm (y-axis) can differ by an order of magnitude. These numbers are given for each line, with numbers reflecting the total ORFs with 100% similarity between replicates (i.e., the end of each curve).
- D) Genomic view of a short upstream ORF (uORF) in the STPBN1 gene indicating that ORF callers have variable affinity for certain types of ORFs. The top two tracks show the ribosomal P-site positions derived from the sequenced ribosome footprints, as processed independently from the sequencing data by the deterministic ORF caller ORFquant (top; red shading) and the probabilistic ORF caller PRICE (bottom; blue shading). The differently colored P-site bars indicate different reading frames (0, +1, +2) on the same transcript, with bars in the same color indicating a shared in-frame codon movement by the ribosome. For this visualization, newly found ORF variations of the annotated CDS that could be assigned to predicted non-coding RNA isoforms (e.g., transcript biotype: "processed_transcript"), but matched the CDS of SPTBN1 is not displayed.
- E) Genomic view of a near-cognate start codon ORF in TUG1. Image and track details as in (E) above.

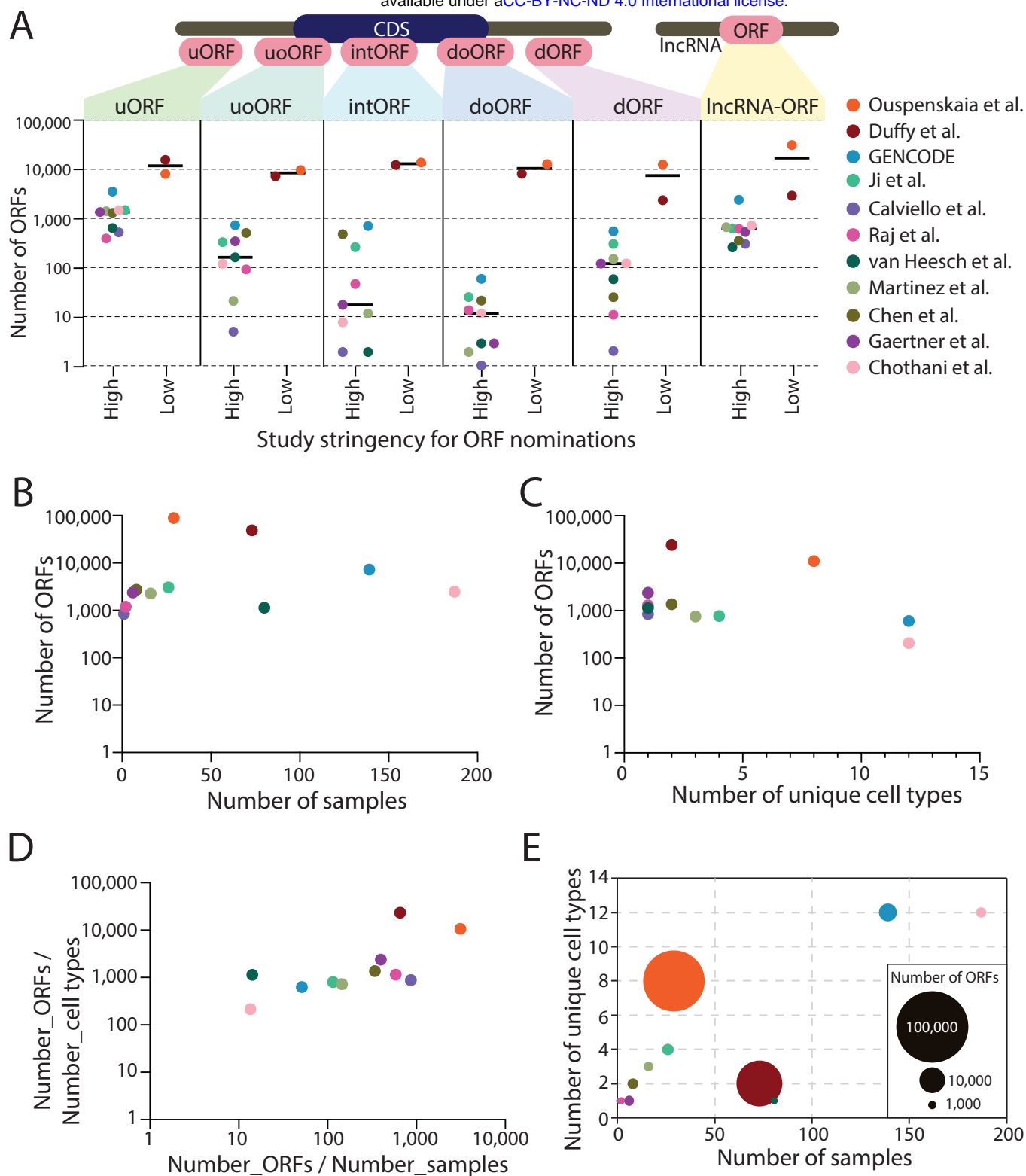


Figure 4: An analysis of major non-canonical ORF databases

- A) Here, each dot reflects a dataset, and the Y axis uses a log-10 scale to show the number of ORFs included that are ≥ 16 amino acids long and contain an AUG start codon. The GENCODE catalog reflects the summation of the Ji et al. (19), Calviello et al. (61), Raj et al. (20), van Heesch et al. (9), Martinez et al. (21), Chen et al. (18) and Gaertner et al. (11) datasets as described in (16).
- B) The number of ORFs per dataset compared to the number of samples profiled by Ribo-seq.
- C) The number of ORFs per dataset compared to the number of unique cell types profiled by Ribo-seq
- D) The ratio of the number of ORFs per cell type compared to the number of ORFs per number of samples for each dataset.
- E) A bubble plot integrating the number of samples, number of different cell or tissue types, and the number of non-canonical ORFs found in each dataset.

Supplementary Materials

What can Ribo-seq and proteomics tell us about the non-canonical proteome?

Author list

John R. Prensner,¹ Jennifer G. Abelin,² Leron W. Kok,³ Karl R. Clauser,² Jonathan M. Mudge,⁴ Jorge Ruiz-Orera,⁵ Michal Bassani-Sternberg,⁶⁻⁸ Eric W. Deutsch,⁹ Sebastiaan van Heesch³

Affiliations

¹Department of Pediatrics, Division of Pediatric Hematology/Oncology, University of Michigan Medical School, Ann Arbor, MI 48109, USA

²Broad Institute of MIT and Harvard, Cambridge, MA, 02142, USA

³Princess Máxima Center for Pediatric Oncology, Heidelberglaan 25, 3584 CS, Utrecht, the Netherlands

⁴European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

⁵Cardiovascular and Metabolic Sciences, Max Delbrück Center for Molecular Medicine in the Helmholtz Association (MDC), 13125 Berlin, Germany

⁶Ludwig Institute for Cancer Research, University of Lausanne, Agora Center Bugnon 25A, 1005 Lausanne, Switzerland

⁷Department of Oncology, Centre hospitalier universitaire vaudois (CHUV), Rue du Bugnon 46, 1005 Lausanne, Switzerland

⁸Agora Cancer Research Centre, 1011 Lausanne, Switzerland

⁹Institute for Systems Biology (ISB), Seattle, Washington 98109, USA

Address correspondence to:

John R. Prensner, MD, PhD

Department of Pediatrics

Medical Science Research Building II, Room 2560B

1150 Medical Center Drive

Ann Arbor, MI 48109

Email: prensner@umich.edu

Phone: 734-763-5939

List of Supplementary Tables

Supplementary Table 1: Description of aggregated data

Supplementary Table 2: A description of ORF primary datasets and data sources

Supplementary Table 3: Phase I Ribo-seq ORFs detected in only 1 of 8 studies

Supplementary Table 4: A metaanalysis of non-canonical ORFs present in 8 different high-stringency datasets

Supplementary Table 5: ORFs included for this analysis from the Duffy et al. dataset.

Supplementary Table 6: ORFs included for this analysis from the Ouspenskaia et al. dataset.

Supplementary Table 7: The number of ORFs for comparative analysis across studies

Supplementary Table 8: The number of ORFs per sample and per cell/tissue type according to each study