

Microproteins - discovery, structure and function

Jessica Mohsen¹, Alina Martel¹, and Sarah Slavoff¹

¹Yale University

July 5, 2023

Abstract

Advances in proteogenomic technologies have revealed hundreds to thousands of translated small open reading frames (sORFs) that encode microproteins in genomes across evolutionary space. While many microproteins have now been shown to play critical roles in biology and human disease, a majority of recently identified microproteins have little or no experimental evidence regarding their functionality. Computational tools have some limitations for analysis of short, poorly conserved microprotein sequences, so additional tools are needed to determine the role of each member of this recently discovered polypeptide class. A currently underexplored avenue in the study of microproteins is structure prediction and determination, which delivers a depth of functional information. In this review, we provide a brief overview of microprotein discovery methods, then examine examples of microprotein structures (and, conversely, intrinsic disorder) that have been experimentally determined using crystallography, cryo-electron microscopy, and NMR, which provide insight into their molecular functions and mechanisms. Additionally, we discuss examples of predicted microprotein structures that have provided insight or context regarding their function. Analysis of microprotein structure at the angstrom level, and confirmation of predicted structures, therefore, has potential to identify translated microproteins that are of biological importance and to provide molecular mechanism for their in vivo roles.

Title: Microproteins—Discovery, Structure, and Function

Authors: Jessica J. Mohsen^{1,2}, Alina A. Martel², Sarah A. Slavoff^{1,2,3}

1. Department of Chemistry, Yale University, New Haven, CT, USA
2. Institute of Biomolecular Design and Discovery, Yale University, West Haven, CT, USA
3. Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT, USA

Corresponding Author : Associate Professor of Chemistry Sarah A. Slavoff^{1,2,3}

sarah.slavoff@yale.edu

Yale University West Campus

Institute of Biomolecular Design and Discovery

Molecular Innovations Center 232

600 West Campus Drive

PO Box 27392

West Haven, CT 06516

Phone: 203-737-8670

Abbreviations: small open reading frame (sORF), sORF-encoded polypeptides (SEPs), long non-coding RNA (lncRNA), untranslated region (UTR), upstream ORF (uORF), overlapping uORF (o.uORF), downstream ORF (dORF), alternative ORF (alt-ORF), ribosome profiling (Ribo-seq), ribosome-protected fragment (RPF), translation initiation sequencing (TI-seq), cryo-electron microscopy (cryo-EM),

sarco/endoplasmic reticulum (SR/ER) calcium ATPase (SERCA), ribosomal RNA (rRNA), short linear interaction motifs (SLIMs), non-homologous end joining (NHEJ), prefoldin-like (PFDL)

Keywords: sORF, Microprotein, Structure, Mass Spectrometry – LC-MS/MS, Genome

Abstract:

Advances in proteogenomic technologies have revealed hundreds to thousands of translated small open reading frames (sORFs) that encode microproteins in genomes across evolutionary space. While many microproteins have now been shown to play critical roles in biology and human disease, a majority of recently identified microproteins have little or no experimental evidence regarding their functionality. Computational tools have some limitations for analysis of short, poorly conserved microprotein sequences, so additional tools are needed to determine the role of each member of this recently discovered polypeptide class. A currently underexplored avenue in the study of microproteins is structure prediction and determination, which delivers a depth of functional information. In this review, we provide a brief overview of microprotein discovery methods, then examine examples of microprotein structures (and, conversely, intrinsic disorder) that have been experimentally determined using crystallography, cryo-electron microscopy, and NMR, which provide insight into their molecular functions and mechanisms. Additionally, we discuss examples of predicted microprotein structures that have provided insight or context regarding their function. Analysis of microprotein structure at the angstrom level, and confirmation of predicted structures, therefore, has potential to identify translated microproteins that are of biological importance and to provide molecular mechanism for their in vivo roles.

1 | Introduction:

Small open reading frames (sORFs, also termed smORFs) below 100 codons were excluded by the FANTOM genome annotation consortium to filter out the high rate of false positive sORFs that were detected under this size in eukaryotic long noncoding RNAs. A similar small size cutoff of 50 codons was applied for prokaryotic gene annotation. These cutoffs were expedient, since the number of known genes pales in comparison to the number of background ORF-like sequences within a genome, most of which are not expressed, but resulted in systematic under-detection of functional sORFs. Many expressed sORFs have now been discovered: recent studies have converged on hundreds of previously unannotated sORFs in bacteria and thousands in human, and multiple CRISPR screens have suggested that hundreds of human sORFs are required for cell survival and proliferation. The emerging relevance of sORFs to infectious disease, the microbiome, and human disease opens the possibility of new therapeutic strategies, and, as such, consortium efforts to enter translated sORFs into the genome annotation are underway.

Early discoveries of functional sORF-encoded polypeptides, such as humanin in human, tarsal-less/polished rice in *Drosophila* and SgrT in bacteria, occurred individually. As a result, the global nature of sORF translation was not recognized until the seminal demonstration of ubiquitous translating ribosome occupancy outside canonical reading frames by Ingolia *et al.* and subsequent confirmation of the presence of a large number of unannotated sORF translation products with mass spectrometry. The products of sORF translation have been termed small proteins, microproteins, micropeptides, sORF-encoded polypeptides (SEPs) and, evocatively, ghost proteins; we will utilize the term microprotein throughout this review. In addition, longer, non-annotated proteins, in some cases referred to as alternative proteins, particularly when they overlap canonical proteins, have also been identified, but they will not be specifically discussed herein. For the purpose of this review, our definition of a eukaryotic microprotein will extend to previously unannotated proteins below 130 amino acids, as many previously undetected ORFs of this length have been reported in human cells. Prokaryotic microproteins are typically categorized as less than or equal to 50 amino acids in length; however, our definition in this work will extend to 70 amino acids since many unannotated microproteins of this size have been detected in multiple bacterial species.

Multiple classes and regions of RNA, both coding and noncoding, have been shown to harbor sORFs in prokaryotes and eukaryotes (Figure 1). Functional sORFs have been discovered in small and long noncoding RNAs (ncRNAs and lncRNAs), antisense lncRNAs, microRNA precursors, and circular RNAs[7,43] in bacteria, plants and other eukaryotes. Interestingly, an increasing number of genes have been shown to exert

functions both at the level of the RNA and of the encoded microprotein, such as *sgrST*, *azuCR*, Spot42/SpfP, and some micro RNAs (miRNAs). Additional classes of sORFs have been identified in multicistronic mRNAs alongside canonical protein coding sequences (CDS) in both prokaryotes and, surprisingly, eukaryotes. sORFs in 5' untranslated regions (UTRs) relative to an annotated CDS are referred to as upstream ORFs (uORFs). Importantly, while eukaryotic uORFs have long been regarded as cis-translational regulators that generally decrease translation efficiency of the downstream CDS, in some instances, uORFs encode microproteins with independent cellular functions in trans, such as MIEF1-MP, which regulates mitochondrial protein translation, and ASDURF, which is a previously unidentified component of the prefoldin-like module of the PAQosome. Some sORFs that initiate in the 5' leader extend into and overlap the CDS in an alternative reading frame, and can be termed overlapping uORFs (o.uORFs), such as human alt-RPL36, which overlaps ribosomal protein L36 and regulates the phospholipid transporter TMEM24. It is important to note that, because they are translated in a different reading frame, o.uORF polypeptide amino acid sequence is completely different from that of the downstream, overlapping annotated protein. At the other end of the mRNA, the 3' UTR has also been found to encode microproteins from downstream ORFs (dORF), which may also regulate CDS translation. An emerging class of frameshifted sORFs occur internally within a protein CDS. These nested sORFs lie completely within the main ORF with translation initiating downstream of the main ORF start codon, and translation terminating upstream of the main ORF stop codon. Nested sORFs can occur in the +2 or +3 (same-strand, frameshifted) reading frames (Figure 2), such as *E. coli* GndA and human alt-FUS. Surprisingly, these findings point to the fact that mammalian mRNAs may be multicistronic or dual coding. While prokaryotic organisms are known to express polycistronic mRNA transcripts termed operons, and compact viral genomes have long been known to contain overlapping open reading frames, eukaryotic transcripts have long been thought to be monocistronic as a result of the scanning model of translation initiation. Importantly, microproteins and longer alternative proteins encoded in each of these classes of sORFs have been shown to be functional. In summary, coding and noncoding regions of both prokaryotic and eukaryotic genomes encode functional sORFs in loci that are denser and more complex than previously presumed.

2 | Microprotein Discovery

2.1 | Computation

Accurate annotation of sORFs using computational tools is challenging not only due to their short lengths that impede statistical analyses, but also because they exhibit intermediate conservation relative to longer genes, which has been interpreted as evidence for the *de novo* evolution of some microproteins. Notwithstanding these challenges, algorithms and machine learning strategies are currently being developed to better find sORFs within genomes. Some computational efforts rely on phylogeny, nucleotide and amino acid homology, and secondary structure to identify unannotated sORFs with sequence or structural similarities to canonical proteins; examples include PhyloCSF and miPFinder. Additional dimensions of predictive information, including the presence of a ribosome binding site upstream of bacterial sORF start codons or a Kozak consensus sequence surrounding a eukaryotic sORF start codon, have been applied to sORF prediction. Ambitiously, OpenProt predicts all AUG-initiated sORFs and alternative ORFs (alt-ORFs) within all known mRNAs for several organisms, and curates experimental evidence (or lack thereof) for their expression. Finally, deep forest and deep learning models have been applied to sORF prediction, with application to individual microbial genomes, as well as the microbiome and metagenomes. These methods have highlighted new sORFs in intergenic regions, noncoding RNAs and in multicistronic/dual coding mRNAs.

2.2 | Ribosome Profiling

Deep sequencing of the protected mRNA footprints of actively translating ribosomes (ribosome profiling or Ribo-seq) has been extensively reviewed and provides a powerful technology for detection of translated sORFs. Ribo-seq is carried out by isolating translating ribosomes associated with mRNA transcripts, using either elongation inhibitors like cycloheximide or rapid freezing. Because translating 80S ribosomes protect bound mRNA fragments from digestion by RNase, sequencing the ribosome-protected fragments (RPFs) reports on translated regions of mRNA. Furthermore, the codon-by-codon elongation of 80S ribosome gives

RPFs a characteristic 3-nucleotide periodic distribution, which can be used to infer the reading frame and confidently differentiate translated ORFs from noise. Furthermore, translation efficiency can be assessed by comparing the frequency of ribosome footprint reads to mRNA transcript levels. Rigorous data analysis, high-resolution datasets, and analysis of replicates are essential for calling sORF translation using Ribo-seq, because their short lengths and translation by monosomes lead to lower signal-to-noise in sORF-mapped reads relative to longer canonical protein coding sequences.

While Ribo-seq is powerful in profiling the footprints of elongating ribosomes and identifying novel coding regions, elongation inhibitors like cycloheximide are not well-suited to deconvolute some translation initiation sites, especially for ORFs with multiple start sites or overlapping reading frames. As a result, a specialized method called translation initiation sequencing (TI-seq) has been developed for inhibition and profiling of the footprints of initiating ribosomes that leverage molecules like puromycin, harringtonine and lactimidomycin in eukaryotes, and retapamulin, tetracycline and Onc112 in prokaryotes. The enrichment of ribosome footprints at canonical and non-canonical start codons in TI-seq datasets generates peaks at the beginning of putative sORFs as well as canonical protein coding sequences. This allows deconvolution of sORF translation initiation from larger main ORFs in multicistronic mRNAs, and is especially important for detection of nested and out of frame sORFs. TI-Seq can also be combined with Ribo-Seq to call translated ORFs with higher confidence.

2.3 | Mass Spectrometry

Mass spectrometry proteomics is able to detect translational products of sORFs directly in biological samples using either bottom-up (from peptide fragments) or top-down (intact precursor) modalities. However, specialized sample preparation and computational methods must be applied for high-sensitivity detection of small, unannotated microproteins. For example, a standard bottom-up proteomics experiment begins with isolation of the proteome, during which small molecules and proteolytic fragments are typically removed by SDS-PAGE or filter-aided sample preparation. Furthermore, most peptide and protein identification from proteomics data is accomplished via spectral matching against the annotated proteome database. For these reasons, sORF-encoded polypeptides are both de-enriched from proteomic samples, and absent from databases, and therefore cannot be detected with standard proteomic workflows and searches.

Multiple recent reviews and protocols describing microprotein identification via proteomics are available, so we provide a brief overview highlighting only the key concerns here. Microprotein discovery methods are built on the same technologies used for standard shotgun proteomics, with several modifications (Figure 3). First, because sORF-encoded microproteins are small, most are identified by only a single proteotypic or fingerprint tryptic fragment in a typical proteomics experiment. A major factor complicating detection of microproteins is coelution and/or cofragmentation of the one or two detectable tryptic peptides derived from a given microprotein with abundant tryptic and/or proteolytic fragments of larger proteins. Resulting ion suppression and/or complex spectra preclude detection and/or identification of the microprotein fragment, regardless of its abundance; this consideration is less severe for larger, canonical proteins, which generate many tryptic peptides and thus detection of any individual fragment is not required. Therefore, the first critical step of any sORF proteomic experiment is to achieve proteome extraction in the absence of proteolysis of canonical proteins (e.g., via boiling in acidic solution or application of protease inhibitors) to minimize sample complexity, followed by or concomitant with enrichment of the small proteome and exclusion of large proteins. Small protein enrichment can be achieved via multiple chemical and biophysical methods, such as solid phase extraction, peptide gels, GELFrEE resolution, and organic solvent or surfactant extraction. When they have been compared head-to-head, these methods have typically been shown to offer comparable numbers, but non-overlapping sets, of microproteins detected. Depending on the experimental goals, the size selection approach for microprotein proteomics can therefore be optimally chosen: for the deepest coverage, multiple methods should be employed on replicate samples and the results combined; for a rapid, robust and economical approach, organic solvent extraction may prove attractive.

Subsequent to small proteome isolation, most microprotein studies to date have employed bottom-up proteomic analysis, in which microproteins are enzymatically digested into peptide fragments (typically with

trypsin, though multienzyme digests have been shown to improve small proteome coverage), followed by liquid chromatography-tandem mass spectrometry, often with multi-dimensional separation. This experiment provides thousands of raw peptide fragmentation spectra corresponding both to known canonical small proteins and microproteins, which must then be identified and distinguished. This is typically accomplished via peptide-spectral matching against expanded databases comprising the canonical proteome as well as candidate sORF sequences. For eukaryotes, databases can be derived from three-frame transcriptome translations, ribosome profiling-derived translomes, or publicly available noncanonical ORF databases such as OpenProt and sORFs.org; six-frame genomic translation can be employed for prokaryotes. Peptide-spectral matching against any of these databases affords identifications of both canonical small proteins and unannotated microproteins. It is important to note that discrimination of false-positive identifications that arise from searching expanded databases is critical. One important consideration is use of a contaminants database to prevent aberrant matching of artefactual peptides (e.g., fragments of trypsin or keratin in dust) to sORF sequences. Another method commonly applied for this purpose is application of a stringent false-discovery rate of less than or equal to 1%, estimated by querying hits to a decoy database constructed from reversed amino acid sequences of the search database entries. However, the expansion of the decoy database also decreases sensitivity for true positive matches, as documented in work from Fournier and colleagues . An alternative approach is to employ permissive false discovery rates, followed by either manual inspection of fragmentation spectra or a secondary algorithm like PepQuery to exclude false positive spectra better explained by peptides arising from canonical, mutant or post-translationally modified proteins. After exclusion of peptides matching (or near-matching) annotated proteins, the resulting list of identifications represent candidate unannotated microproteins, which can be computationally mapped to the sORFs that encode them and experimentally validated.

Mass spectrometry typically detects one to two orders of magnitude fewer microproteins in a given experiment than ribosome profiling. This may be due to the abovementioned challenge in detecting single microprotein-derived fingerprint peptides; the relative insensitivity of mass spectrometry to some classes of microproteins, including membrane-localized, positively charged, and low-abundance species; the instability of some sORF translation products; reduced sensitivity for true-positive detections as a result of expanded decoy databases applied for stringent false discovery rate estimation; or all of these factors. Nonetheless, mass spectrometry offers several advantages. First, enrichment strategies, such as membrane fractionation and chemical labeling, can enable identification of microproteins that are refractory to shotgun analysis of whole-cell tryptic digests, thus beginning to address one of the major limitations of microprotein proteomics while at the same time affording functional information about microproteins (e.g., chemical reactivity, subcellular localization) that is inaccessible to sequencing methods. Second, without specialized analysis pipelines, ribosome profiling with elongation inhibitors is refractory to confident detection of sORFs that overlap canonical protein coding sequences in alternative reading frames, due to the requirement for three-nucleotide periodicity for ORF calling. In contrast, mass spectrometry can readily detect and identify microproteins derived from overlapping ORFs, which can represent as much as 30% of microproteins identified in a proteomic experiment. Given the complementary nature of genomics, ribosome profiling and mass spectrometry, it is likely that the combination of these methods offers the greatest power for large-scale, high-confidence microprotein identification.

3 | Microprotein Structure and Function

Dozens of human microproteins, and many more in model organisms, have now been assigned function at the molecular, cellular, and/or organismal level. CRISPR screens have implicated hundreds of sORFs in cell survival and proliferation. Experimental approaches are yielding insights into the roles of microproteins in biological processes and disease, which have been extensively reviewed. Recently emerging trends in microprotein function include roles in immunity and inflammation, mitochondrial functions and energetics, adiposity, microbial carbon metabolism, and cancer initiation and progression, among others. Nonetheless, the vast majority of recently discovered microproteins remain entirely uncharacterized in mechanistic detail. This is in large part because bioinformatic predictions of sORF function are challenging—even when they exhibit signatures of conservation in multiple species, microproteins tend to lack primary sequence homology

to proteins of known function. While three-dimensional structure prediction and elucidation are likely to provide important insights into microprotein functions, structures of microproteins have not yet been examined on scale. However, the number of experimentally determined structures of microproteins, in isolation or in complex with their effectors, is growing, and general trends have begun to emerge, which we will describe in this section. First, we discuss a subclass of single-pass alpha-helical transmembrane microproteins, many of which are evolutionarily novel, and some of which bind to and regulate important transporters. Next, we consider examples of microproteins with solved or predicted structures and the potential relevance to their functions. Last, we will examine several intrinsically disordered microproteins that undergo regulatory protein-protein interactions.

3.1 | *Alpha-helical transmembrane microproteins*

Intergenic regions of eukaryotic genomes are rich in A/T residues relative to genes, which are G/C rich. When microproteins are expressed from “noncoding” regions, they therefore tend to contain predicted transmembrane helices arising from the preponderance of T/U residues within codons that correspond to hydrophobic and aromatic amino acids. This intergenic sequence bias therefore affects the amino acid composition of evolutionarily young, species-specific microproteins, that arise *de novo* from previously noncoding regions of the genome. A recent study demonstrated that C-terminal hydrophobic patches tend to target evolutionarily young microproteins to the BAG6 membrane protein triage complex, resulting either in membrane insertion or, if mislocalized or improperly folded, proteasomal degradation. Interestingly, species-specific transmembrane microproteins that exhibit low expression can nonetheless contribute fitness advantage to cells, and examples have been shown to function in processes such as yeast mating. Not all membrane-associated microproteins are evolutionarily novel; a large and growing number of well-characterized, conserved transmembrane microproteins are predicted to contain transmembrane helices, such as the lysosomal membrane-localized polypeptide regulator of mTORC1, SPAR, and the plasma membrane localized micropeptide Myomixer, which is required for myoblast fusion during skeletal muscle development. The class of alpha-helical transmembrane microproteins is therefore large, and of outsize biological importance. We turn our attention in this section to those membrane-associated microproteins that have been subjected to experimental structure determination.

AcrZ, previously named YbhT, was reported in a seminal study identifying unannotated small protein genes in *E. coli* utilizing computational tools that incorporate ribosome binding site prediction. AcrZ is a 49-amino acid microprotein that is conserved in many Gram-negative bacteria and localizes to the *E. coli* inner and outer membranes by virtue of an N-terminal transmembrane helix. AcrZ binds to the AcrB subunit of the AcrAB/*tolC* multidrug efflux pump, increasing the efficiency of transport of (and, thus, resistance of *E. coli* to) a subset of its substrates. Multiple structures of AcrZ in complex with the AcrB homotrimer have been solved, including crystal structures of detergent-solubilized complexes, as well as a cryo-electron microscopy (cryo-EM) structure of the complex reconstituted in lipid discs (Figure 4A). AcrZ binds to a transmembrane groove within each molecule of AcrB. The cryo-EM structure revealed that AcrZ exhibits a profound bend between positions 10-15, conferred by a helix-breaking proline residue. Mutagenesis studies revealed that the proline is required for interaction of AcrZ with AcrB. At the same time, proline, or an equally helix-breaking glycine residue, can be moved to any position within the AcrZ interaction motif while retaining its association with AcrB. Several of these mutations that retain AcrB binding also recapitulate the selective drug transport-promoting phenotype of wild-type AcrZ. While the precise effects of AcrZ binding on cargo occupancy and transport are not fully clear, allosteric modulation of binding sites in AcrB is evident by comparing the AcrB vs. AcrBZ structures. Furthermore, AcrZ promotes cardiolipin association with AcrB, likely contributing to allosteric modulation of cargo binding pockets in the transporter. Taken together, these results indicate that the bend in the transmembrane helical shape of AcrZ, and not its amino acid sequence, is essential for interaction and modulation of AcrB.

E. coli CydX was originally identified as YbgT, a predicted 37-amino acid microprotein encoded downstream of the cytochrome bd oxidase operon genes *cydA* and *cydB*. Cytochrome bd oxidases operate as terminal electron acceptors in the electron transport chain under hypoxic conditions due to their high oxygen affinity.

The two canonical subunits, CydA and CydB, form a pseudosymmetric heterodimer, of which the CydA subunit contains all three heme residues responsible for reduction of molecular oxygen to water, as well as the Q loop that is responsible for binding an electron donor quinol. CydX is a single-pass alpha helical transmembrane protein that copurifies with the CydAB complex and is required for the assembly, stability, and/or activity of cytochrome bd oxidase in multiple species. Several atomic structures of cytochrome bd oxidases have revealed the role of CydX homologs in the complex (Figure 4B). First, the presence of an unannotated, CydX homolog, CydS, was serendipitously discovered in a crystal structure of cytochrome bd oxidase purified from the gram-positive bacterium *Geobacillus thermodenitrificans*. CydS forms an alpha helix that binds between helices 5 and 6 of CydA, leading the authors to speculate that it may stabilize the heme cofactor when the Q loop undergoes dynamic movement during catalysis. A subsequent cryo-EM structure of the *E. coli* cytochrome BD oxidase revealed CydX bound to CydA between helices 1 and 6, again suggesting a structural role. Interestingly, the *E. coli* CydAB unexpectedly revealed the presence of another single-pass transmembrane microprotein, CydH, which is encoded in the *ynhF* gene that is not located within the cytochrome bd oxidase operon. CydH binds between transmembrane helices 1 and 8 of CydA, on the opposite face of CydA relative to CydX. CydH is proposed to occlude the proposed oxygen-conducting channel from the *Geobacillus* complex structure, which has been replaced with a hydrophobic channel that traverses CydB directly to the heme d site. The CydH oxygen channel rearrangement was proposed to be required due to the swapped positions of two heme cofactors in the *E. coli* enzyme relative to the *Geobacillus* structure, and, accordingly, CydH homologs are found in Proteobacteria. Overall, cytochrome bd oxidase is a unique system in which microproteins are required for activity, structure and stability of a critical complex of proteins.

In another well-characterized example, a class of microproteins (also called micropeptides) termed “regulins” regulate the activity of the sarco/endoplasmic reticulum (SR/ER) calcium ATPase (SERCA). During muscle contraction, including the contraction of the heart and calcium-dependent signaling processes, calcium is released from the SR/ER into the cytosol; then, to terminate signaling or contraction, calcium is pumped back into the SR/ER against its concentration gradient using the energy of ATP hydrolysis by SERCA. Regulins colocalize with SERCA in the SR/ER membrane, and each micropeptide is expressed in the same, specific tissue as the SERCA isoform that it regulates. The first known regulins, phospholamban and sarcolipin, were identified as inhibitors of SERCA in cardiac and skeletal muscle, respectively. Structural analysis of these canonical regulins, both of which are <100 amino acids, reveals that they are small, single-pass membrane proteins bearing a single transmembrane alpha-helix. The crystal structure of the SERCA-sarcolipin complex reveals that the micropeptide binds in a transmembrane groove in the SERCA channel between helices 2, 6 and 9, where it allosterically alters the conformation of SERCA to decrease its apparent calcium affinity. Phospholamban binds to the same regulatory groove (Figure 4C). One seminal discovery of novel SERCA regulating micropeptides came from a study in *Drosophila*. In this work, Couso and colleagues analyzed putative lncRNAs associated with polysomes, suggesting that they are translated. Of these lncRNAs, one contained an sORF encoding a peptide predicted to be homologous to phospholamban and sarcolipin, which was accordingly given the name sarcolamban. Sarcolamban may have arisen via duplication of an ancestral phospholamban/sarcolipin gene in insects, which subsequently diverged to the sarcolamban sequence. Sarcolamban was demonstrated to bind SERCA in flies and its deletion caused heart arrhythmias, consistent with a role in regulating SERCA. Docking the predicted structure of sarcolamban onto SERCA was consistent with a similar binding mode as that observed for phospholamban and sarcolipin. Just as importantly, additional novel regulins have also been discovered in mammals. In analyses of mammalian lncRNAs to identify potential micropeptides expressed in skeletal muscle and other tissues lacking known regulin expression, translated sORFs were identified that encode the novel SERCA binding micropeptides myoregulin, endoregulin, and another-regulin, all of which bind to the same transmembrane groove of SERCA, exhibit similar inhibition of SERCA to phospholamban, and are predicted to have similar single-pass transmembrane alpha-helical structures. Interestingly, an unannotated, SERCA-activating micropeptide, DWORF, was identified in yet another long noncoding RNA in mouse. DWORF is expressed in skeletal muscle, and ectopic over-expression of DWORF in heart tissue enhances contractility and reverses heart failure in a model of heart failure. However, the mechanism by which DWORF activates

SERCA was unclear, since it is predicted to bear a similar alpha-helical transmembrane domain and binds to the same SERCA groove as previously characterized regulins, which are all inhibitory. Some evidence from fluorescence resonance energy transfer suggests that DWORF binding can directly activate SERCA. A recent NMR structural study demonstrated that the alpha helix of DWORF is kinked at a unique proline residue, creating a significant bend in the transmembrane region without disrupting its binding to SERCA (Figure 4C). Mutating this proline residue diminished the bend angle between the two alpha helical regions of DWORF, and not only prevented its activation of SERCA, but converted it into a SERCA inhibitor. Therefore, activation of SERCA by DWORF appears to require its proline-induced kink, and, by extension, inhibition of SERCA by phospholamban, sarcolipin, myoregulin, endoregulin and another-regulin may be hypothesized to require binding of their uninterrupted transmembrane helices to the regulatory groove of SERCA. It is also fascinating to note the parallels between DWORF and AcrZ (see above), both of which utilize kinked transmembrane alpha-helices to allosterically regulate the membrane transporters SERCA and AcrB, respectively.

3.2 | *Humanin and its disorder-to-order transition*

Humanin is a secreted 24-amino acid polypeptide found in human serum that protects neurons from cell death in the presence of familial early onset-Alzheimer's disease-associated mutants of amyloid precursor protein. Interestingly, the humanin coding sequence was mapped to a polyadenylated cDNA that was expressed in the surviving brain tissue of an Alzheimer's disease patient, and is derived from the mitochondrial 16S ribosomal RNA (rRNA). Given that another mitochondrial peptide, MOTS-C, is encoded in a region overlapping the mitochondrial 12S rRNA, this raises the intriguing possibility that the mitochondrial rRNA genes may be polycistronic, though the molecular mechanisms by which microprotein-encoding transcripts are generated or processed are not yet defined. Humanin's neuroprotective effects have been proposed to occur through multiple intracellular and cell-surface interaction partners, including BAX, IGFBP3, FPRL1, and CNTF Receptor α /WSX-1/gp130, though the relative contributions of these pathways to its in vivo activity remain to be determined. A circular dichroism and NMR study of humanin revealed that it does not adopt a stable secondary structure in aqueous solution, although through-space interactions consistent with turns at the N- and C-termini of the peptide were observed. In contrast, in 30% organic solvent, humanin forms an alpha helix spanning residues G5 to L18 (Figure 4D). This suggests that humanin may fold in hydrophobic environments such as cell membranes or in complex with interaction partners. Testing this hypothesis could provide deeper insight into its localization and associations with functional interaction partner(s).

3.3 | *Ubiquitin-like microproteins*

Several groups recently reported the discovery of ubiquitin-like microproteins. In one example, the ubiquitin pseudogene *UBBP4* was reported to be translated. Interestingly, *UBBP4* encodes three ubiquitin variants within two independent open reading frames, and mass spectrometric evidence uniquely identifying all three have been previously obtained. The *UBBP4* ubiquitin-like proteins exhibit high sequence similarity to canonical ubiquitin, with 1 (variant Ubbp4^{A1}), 4 (Ubbp4^{B1} or Ub^{KEKS}), or 8 amino acid substitutions (Ubbp4^{A2}). Ubbp4^{A2} and Ub^{KEKS} retain a functional C-terminal diglycine motif and can be covalently conjugated to high molecular weight cellular proteins, while Ubbp4^{A1} was predominantly observed as a monomer. Despite being ~700-fold less abundant than canonical ubiquitin, Ub^{KEKS} modifies a specific subset of cellular proteins including lamins, and, rather than promoting proteasomal degradation, may be important for regulating target protein localization and/or function.

In 2020 the *TINCR* RNA, which was previously classified as noncoding, was shown to encode an 87-amino acid microprotein with 85% sequence homology to ubiquitin. The microprotein translated from *TINCR* RNA, termed pTINCR or TUBL, was predicted to adopt a ubiquitin-like fold (Figure 4D). This prediction was confirmed in a recent crystal structure of pTINCR, which revealed an overall ubiquitin-like fold with a positively charged N-terminal domain hypothesized to enable interaction with other biomolecules (Figure 4E). Due to the lack of a C-terminal diglycine motif, pTINCR is a type II ubiquitin-like protein that associates with ubiquitin-binding proteins rather than being covalently attached to proteins. pTINCR is expressed in skin, and mice lacking pTINCR exhibit a mild delay in wound healing. Importantly, two

reports have identified pTINCR as a tumor suppressor in cutaneous squamous cell carcinoma and other epithelial cancers. pTINCR is upregulated after DNA damage-induced p53 activation, and it is frequently lost or mutated in squamous cell carcinoma. It normally promotes differentiation of keratinocytes and other epithelial cell types via its interaction with SUMOylated Cdc42. Consequently, mouse embryonic stem cell-derived teratomas overexpressing pTINCR exhibit decreased growth and increased keratin deposition consistent with involvement in differentiation of skin cells. Along the same lines, pTINCR overexpression inhibits the proliferation of squamous cell carcinoma cells in culture and in xenografts. Additionally, mice heterozygous for *Xpc* that lack pTINCR are DNA damage repair-deficient and exhibit increased formation of invasive skin papillomas and squamous cell carcinomas relative to *Xpc*heterozygous/pTINCR wild-type mice upon UV exposure. Overall, pTINCR is a type II ubiquitin-like microprotein that is required for keratinocyte differentiation and acts as a tumor suppressor in squamous cell carcinoma.

3.4 | Microproteins with predicted structures

With the advent of three-dimensional macromolecular structure prediction tools such as Rosetta, iTasser, Phyre, and, most recently, AlphaFold, many recently discovered, now-annotated microproteins have been subjected to computational structure prediction, and these structural models are publicly available. For microproteins that remain unannotated, computational tools can be used to generate testable structural predictions. For example, analysis of the recently identified *E. coli* cold-shock microprotein YmcF using iTasser led to the hypothesis that YmcF may adopt a folded structure consisting of an alpha helix and 2-3 beta strands separated by a turn, homologous to a zinc-binding domain of aspartate transcarbamoylase (Figure 4G). While no functional data for YmcF yet exists, this predicted structural model, if correct, may have implications in the cold shock response, which requires RNA binding proteins—some of which coordinate zinc—to chaperone RNA secondary structures that become hyper-stable at low temperature. In another example, plant microProteins are specifically defined as proteins predicted to fold into single domains that bind to and generally antagonize the functions of their effectors, such as transcription factors.

Predicted structures of microproteins have already begun to aid in determining their molecular and cellular functions. A translated upstream ORF (uORF) encoding a 96-amino acid microprotein within the 5' untranslated region (UTR) of the human *ASNSD1* gene was reported by Oyama et al. in 2007 and in subsequent proteomic analyses, leading to the annotation of the microprotein as ASDURF (ASNSD1 upstream open reading frame). As discussed above, evidence is accumulating that uORF microproteins can function *in trans*. Remarkably, Coulombe and colleagues recently implicated ASDURF as the “missing” subunit of a chaperone complex termed the PAQosome. Proximity biotinylation and pull-down experiments with PAQosome subunits revealed ASDURF as an interaction partner, and in vitro reconstitution assays suggested that it is an integral member of a PAQosome subcomplex. The PAQosome is a recently discovered chaperone that is essential for assembling complicated macromolecular complexes in the cell, including RNA polymerases, components of the spliceosome, and protein phosphatases. The PAQosome consists of two modules, one of which is termed the prefoldin-like (PFDL) module. The PFDL module shares some subunits and putative structural homology to prefoldin, another cellular chaperone required for folding cytoskeletal proteins and other clients. Prefoldin and the PFDL module are both hexameric, consisting of three alpha- and three beta-prefoldin subunits, which both contain an alpha-helical coiled-coil separated by either one (beta) or two (alpha) hairpins; however, only five of the six PFDL subunits (three alpha and two beta) had been identified. Tertiary structure modeling with Phyre suggested that ASDURF is a beta-prefoldin bearing a single beta hairpin and coiled-coil (Figure 4H), consistent with its potential identification as the undiscovered beta subunit of the PFDL module of the PAQosome – suggesting it had been missed because it was not part of the proteome annotation at the time of the PAQosome’s discovery. Many additional interesting questions are raised by the ASDURF microprotein: Why is it encoded in an upstream ORF within the *ASNSD1* gene? Does its 5' UTR location confer stress responsiveness via translational regulation, as suggested by Cloutier et al.? Is its function or regulation related to the downstream ASNSD1 protein, per the model of Chen, Weissman and colleagues that co-encoded microproteins and proteins tend to function in the same pathways? Regardless, while the structural model requires experimental validation, it appears that ASDURF is a particularly compelling example of a microprotein for which structure prediction informs its interactions

and likely function.

3.5 | *Intrinsically disordered microproteins*

Microproteins are much shorter than annotated proteins, and they tend to exhibit limited conservation to protein domains of known function. As a result, it is challenging to perform bioinformatic analyses, for example of predicted structure or intrinsic disorder, of microproteins with confidence, particularly because many of these predictive algorithms rely, at least in part, on homology to structures of known, larger proteins on which they are trained. Nonetheless, some studies have suggested that microproteins may be enriched in intrinsic disorder relative to canonical proteins (though an alternative analysis suggests that evolutionarily young microproteins are de-enriched in intrinsic disorder), which, if true, suggests that some microproteins could carry out cellular functions associated with intrinsically disordered proteins, such as regulating signaling and other processes by binding to protein partners via short linear interaction motifs (SLIMs). In this section we discuss two human microproteins that have been experimentally confirmed to be predominantly intrinsically disordered.

MRI (Modulator of retroviral infection) was first identified in a cDNA library screen for host proteins that could complement resistance to retroviral infection of human cells, but it remained annotated as a predicted or uncharacterized protein-coding gene (*C7ORF49*) in the early 2010s. While the long isoform of MRI (MRI-1 hereafter) is 157 amino acids long and therefore not a microprotein, a 2013 peptidomics study identified an unannotated, sORF encoded isoform (MRI-2) of 69 amino acids. Follow-up work demonstrated that the long MRI-1 and short MRI-2 proteins could interact with a complex of proteins essential for the non-homologous end joining pathway (NHEJ), which is essential for repairing DNA double strand breaks in G1 phase of the cell cycle, as well as for B and T cell receptor gene diversification via V(D)J recombination. Specifically, MRI-1 interacts with the double-strand break binding adaptor proteins Ku70/80 (Ku) and DNA-PKcs (DNA-dependent protein kinase catalytic subunit), while MRI-2 binds to Ku. Both of these MRI isoforms contain an N-terminal Ku-binding motif, explaining their association with Ku, while MRI-1 also contains a C-terminal XLF-like motif (XLM) that associates with additional, distinct NHEJ factors. The XLM of MRI-1 is absent in the frameshifted, truncated MRI-2 isoform. One study suggests that MRI inhibits aberrant NHEJ at telomeres during S phase, while two studies to date are consistent with a positive role for MRI in NHEJ during most phases of the cell cycle, suggesting that the activity of MRI may be context-dependent. Purified MRI-2 was shown to promote NHEJ in vitro. However, abrogating all isoforms via knockout of the MRI gene in vivo and in pre-B cells increases sensitivity to ionizing radiation and inhibits NHEJ when coupled with knockout of the NHEJ “sentinel” gene XLF. Purified MRI-1 was shown to be predominantly intrinsically disordered via hydrogen-deuterium exchange; while MRI-2 was not directly investigated in this study, it is likely to have a similar degree of intrinsic disorder because these proteins share substantial sequence identity until the frameshift that truncates MRI-2. Interestingly, the N-terminal and C-terminal motifs of MRI-1 alone can nucleate separate complexes of NHEJ factors, and MRI-1 can recruit NHEJ factors to chromatin in the presence of DNA double strand breaks. It is interesting to speculate MRI-2 may therefore be able to serve the same nucleating function in NHEJ via its Ku-binding motif even in the absence of the C-terminal XLM. Sleckman and colleagues proposed that MRI-1 serves as an adaptor protein for NHEJ, promoting stable association of active NHEJ complexes at sites of double strand breaks as a result of its (1) intrinsic disorder, (2) independent linear interaction motifs, and (3) its potential to multimerize. While better understanding of the contributions of individual MRI isoforms to their function in vivo is required, MRI-1 and MRI-2 appear to be paradigmatic examples of intrinsically disordered (micro)proteins that promote assembly of a functional protein interaction network.

Another example of an experimentally validated, intrinsically disordered microprotein is NBDY. NBDY is a 68-amino acid microprotein expressed from a previously misannotated lncRNA (*LOC550643*). NBDY associates with members of the cytoplasmic mRNA decapping complex. The interaction partners of NBDY, EDC4 and DCP1A, are coactivators required for allosteric activation of DCP2, which catalyzes the first step in 5'-to-3' mRNA decay (removal of the 7-methylguanosine cap), thus regulating the stability of thousands of specific mRNA substrates. Genetic ablation or silencing of NBDY stabilizes a majority of DCP2

substrates, consistent with the requirement of NBDY for their effective decapping, including transcripts encoding proteins involved in immune responses – a pathway previously reported to be regulated by DCP2. However, at the same time, a number of DCP2 substrates are destabilized by NBDY ablation, suggesting that the microprotein may act as a specificity factor for recruitment of mRNA targets to the decapping complex. In particular, in the presence of NBDY, DCP2 substrate mRNAs with shorter 5'UTRs decay more rapidly, suggesting that there may be a requirement for NBDY for efficient recognition of transcripts with short leader sequences by DCP2. While the molecular mechanism by which NBDY regulates the mRNA decapping complex is not yet known, mRNA decapping proteins have previously been reported to associate via SLIMs within disordered regions, and it is likely that NBDY participates in this network. NMR experiments indicated that NBDY is largely intrinsically disordered in solution, consistent with its ability to phase-separate in the presence of RNA to form liquid droplets *in vitro*. Within the intrinsically disordered NBDY sequence, two independent SLIMs interact with the WD40 domain of EDC4 and the EVH1 domain of DCP1A. The interaction between EDC4 and NBDY appears to be more important for NBDY function in mutagenesis experiments, but, given the relatively low affinity of NBDY for EDC4 ($K_D \sim 1$ micromolar), the interaction with DCP1A could speculatively be important for increasing avidity of NBDY for the mRNA decapping complex, retaining it at interaction sites. Importantly, NBDY also partially localizes to and regulates phase-separated RNA granules termed P-bodies in cells, consistent with a role for intrinsically disordered microproteins in biological phase separation. NBDY is phosphorylated downstream of EGFR and cyclin-dependent kinase signaling, and this phosphorylation is required for dissociation of P-bodies – likely via electrostatic repulsion of negatively charged P-body components that promotes liquid-phase remixing and cell proliferation. Taken together, NBDY's intrinsic disorder enables its SLIM-mediated protein-protein interactions, phase separation and regulation of P-bodies, providing a well-defined example of the functional significance of intrinsic disorder in a microprotein.

4 | Conclusion

As microproteins are increasingly linked with roles in human health and disease, elucidating their numbers and biological roles will be ever more essential. Regarding the complete annotation of microproteins, while there are still inconsistencies in the specific sORF loci identified across ribosome profiling studies, most recent studies detect comparable numbers of translated sORFs in a given organism. This developing consensus suggests that meta-analysis of ribosome profiling data has the potential to resolve complete sORF translomes in the near future. As this effort advances, large-scale CRISPR screens and other methods can be (re-)employed to identify functional sORFs on scale. However, it is important to note that most sORF functional screens to date have focused on cell proliferation/survival, protein-protein interactions, and/or conservation, thus potentially screening out sORFs with roles beyond these readouts. For example, microproteins with clear involvement in yeast mating and cellular responses to stress have been reported, but can be species-specific, nonessential and may not undergo long-lived interactions with other proteins, and thus would not appear as hits in most functional screens to date. Thus, alternative avenues to identify microproteins with potential functions are needed. Given the exquisite link between protein three-dimensional structure and function, investigation of microprotein structure holds tremendous promise to address this need. The advent of AlphaFold, combined with the rapidly increasing number of solved microprotein structures and experimentally characterized intrinsically disordered microproteins, including those described above, are already contributing to the improved power of structural prediction to generate functional hypotheses about uncharacterized microproteins. Experimental structural investigations are also providing critical mechanistic insights into how microproteins exert their functions, for example in allosteric regulation of target proteins. Combined with insights into disease-associated microprotein mutations and dysregulation, structural and mechanistic information may also pave the way to determining whether microproteins and/or their binding partners are druggable in the future.

Acknowledgments

This work was supported by a Mark Foundation for Cancer Research Emerging Leader Award, a Paul G. Allen Frontiers Group Distinguished Investigator Award, and a Sloan Research Fellowship (FG-2022-18417)

(to S.A.S.), a Yale University fellowship associated with the NIGMS Chemistry-Biology Interface Training Program (5T32GM067543), and a Roberts Fellowship from the Yale University Department of Chemistry (to J.J.M.).

Conflict of Interest

The authors have no conflicts of interest to declare.

Data Availability Statement

Data sharing is not applicable to this article as no new data were created or analyzed in this study.

References

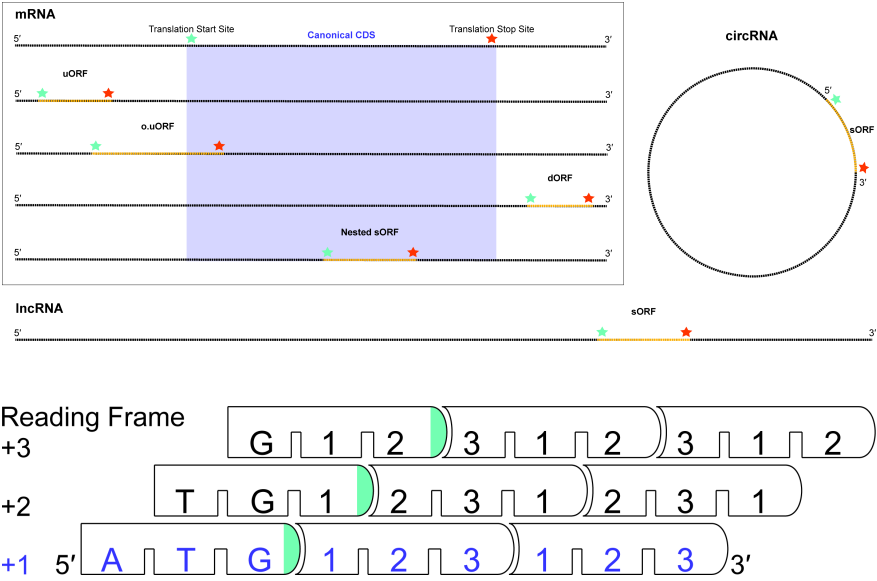
Figure Legends

FIGURE 1 Small Open Reading Frames (sORFs) and RNA. Box: Within mRNA that encodes canonical protein coding sequences (CDS), sORFs can appear in the 5' UTR (upstream ORF, uORF), initiating in the 5' UTR and extending into the CDS in an alternative reading frame (upstream overlapping ORF, u.oORF), in the 3' UTR (downstream ORF, dORF), or nested within the CDS in an alternative reading frame. sORFs can also be found in long noncoding RNA (lncRNA, bottom) and circular RNA (circRNA, right), as well as additional classes of RNA not pictured.

FIGURE 2 Alternative Reading Frames for Same-Strand Overlapping (Nested) sORFs. The +1 reading frame corresponds to the canonical coding sequence and is always the frame of reference. Frameshifted translation in the +2 or +3 reading frames generates protein products with completely different amino acid sequences because the codon identities are changed in alternative reading frames.

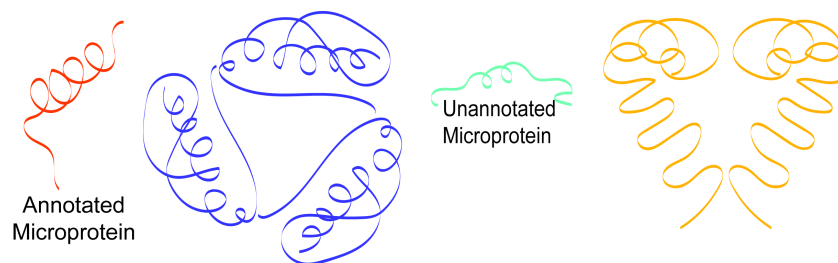
FIGURE 3 Mass Spectrometry Workflow for Detection of Unannotated Microproteins. To search for novel microproteins in a sample of interest, low molecular weight proteins are isolated from total protein after cell lysis. Size-exclusion techniques include, but are not limited to, solid-phase extraction and polyacrylamide gel electrophoresis techniques. Low molecular weight protein is digested with a protease, producing a sample of uniform peptide length appropriate for mass spectrometric (MS) analysis. Experimental spectra are generated and matched to theoretical spectra from a custom database using proteomics software. Detection of annotated microproteins known to be expressed in the system of interest can serve as a positive control for success of small protein enrichment and known small proteome coverage, but these spectra are otherwise computationally excluded. Peptides deriving from proteolysis of canonical proteins before size-exclusion are computationally identified and excluded from consideration. High scoring experimental spectra without any matches to known microproteins can be subjected to further molecular validation, leading to annotation of novel microproteins.

FIGURE 4 Experimentally Determined Microprotein Structures. (A) Crystal structure of AcrB (grayscale) of the *tolC* efflux pump in complex with microprotein AcrZ (cyan). PDB: 5NC5. (B) Cryo-EM structure of bacterial microprotein CydX (cyan) in complex with transmembrane cytochrome bd-I oxidase (grayscale). PDB: 6RKO. (C) Crystal structure of SERCA1a calcium pump (grayscale) with bound single-pass transmembrane microprotein phospholamban (cyan), which downregulates SERCA activity. PDB: 4Y3U. Solid-state NMR structure of helix-loop-helix microprotein DWORF (cyan) modeled into SERCA1a calcium pump (grayscale) based on Venkateswara et al. 2022. PDB: 4Y3U, 7MPA. (D) NMR structure of wild-type humanin in 30% 2,2,2-trifluoroethanol (organic) solution. PDB: 1Y32. (E) Crystal structure of Ubiquitin monomer. PDB: 1AAR. (F) Crystal Structure of ubiquitin-like TINCR microprotein with additional N-terminal alpha helix. PDB: 7MRJ. (G) Predicted structure of bacterial microprotein YmcF generated with AlphaFold, obtained from UniProt[166] (green). Five cysteines (orange) in the YmcF sequence are predicted to form a zinc-finger domain common to RNA binding proteins. (H) Predicted structure of PAQosome binding microprotein ASDURF generated with AlphaFold, obtained from UniProt[166].

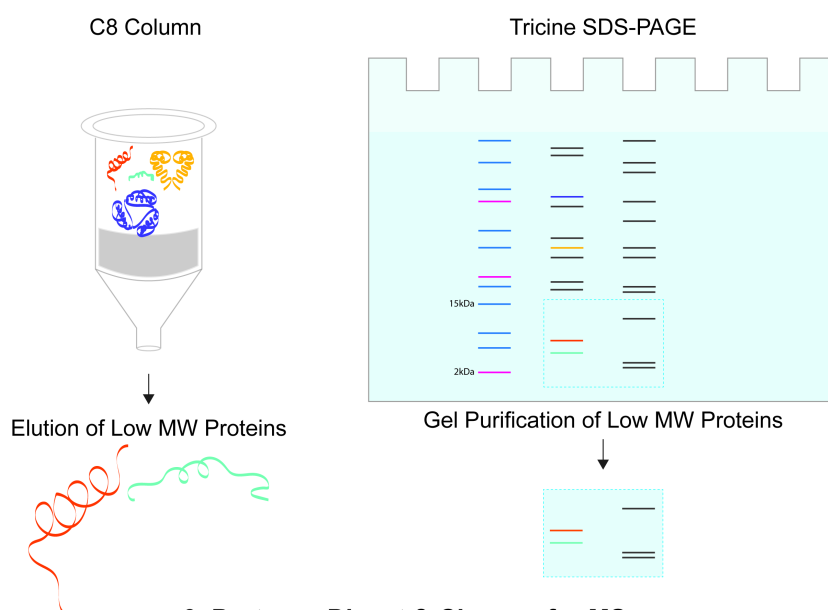


Microprotein Discovery Mass Spectrometry Workflow

1. Lysis and Protein Extraction

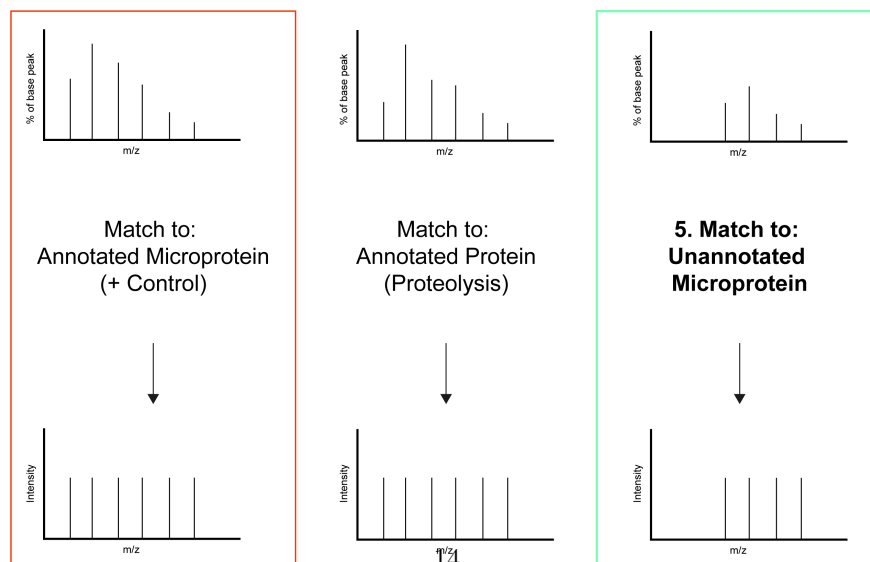


2. Size-Exclusion Techniques



3. Protease Digest & Cleanup for MS

4. MS & Peptide Spectral Matching



Theoretical Spectra from 6-Frame (Prokaryotic) or 3-Frame (Eukaryotic)
in silico Digest of Annotated Genome Sequence

Quality Control

6. Molecular Validation

