



DATABASE

SmProt: A Reliable Repository with Comprehensive Annotation of Small Proteins Identified from Ribosome Profiling



Yanyan Li^{1,2,#}, Honghong Zhou^{2,#}, Xiaomin Chen^{2,3,#}, Yu Zheng^{1,2}, Quan Kang², Di Hao², Lili Zhang^{2,3}, Tingrui Song², Huaxia Luo², Yajing Hao⁴, Runsheng Chen^{2,3,5,*}, Peng Zhang^{2,*}, Shunmin He^{1,2,*}

¹College of Life Sciences, University of Chinese Academy of Sciences, Beijing 100049, China

²Key Laboratory of RNA Biology, Center for Big Data Research in Health, Institute of Biophysics, Chinese Academy of Sciences, Beijing 100101, China

³University of Chinese Academy of Sciences, Beijing 100049, China

⁴Department of Cellular and Molecular Medicine, University of California, San Diego, La Jolla, CA 92093, USA

⁵Guangdong Geneway Decoding Bio-Tech Co. Ltd, Foshan 528316, China

Received 9 December 2020; revised 7 September 2021; accepted 8 September 2021

Available online 15 September 2021

Handled by Zhang Zhang

Abstract Small proteins specifically refer to proteins consisting of less than 100 amino acids translated from small open reading frames (sORFs), which were usually missed in previous genome annotation. The significance of small proteins has been revealed in current years, along with the discovery of their diverse functions. However, systematic annotation of small proteins is still insufficient. SmProt was specially developed to provide valuable information on small proteins for scientific community. Here we present the update of SmProt, which emphasizes reliability of translated sORFs, genetic variants in translated sORFs, disease-specific sORF translation events or sequences, and remarkably increased data volume. More components such as non-ATG translation initiation, function, and new sources are also included. SmProt incorporated 638,958 unique small proteins curated from 3,165,229 primary records, which were computationally predicted from 419 ribosome profiling (Ribo-seq) datasets or collected from literature and other sources from 370 cell lines or tissues in 8 species (*Homo sapiens*, *Mus musculus*, *Rattus norvegicus*, *Drosophila melanogaster*, *Danio rerio*, *Saccharomyces cerevisiae*, *Caenorhabditis elegans*, and *Escherichia coli*). In addition, small protein families identified from human microbiomes were also collected. All datasets in SmProt are free to access, and available for browse, search, and bulk downloads at <http://bigdata.ibp.ac.cn/SmProt/>.

KEYWORDS Ribosome profiling; Small open reading frame; Upstream open reading frame; Variants; Disease

Introduction

Genome annotation is fundamental to life science. In recent years, it has been found that small open reading frames (sORFs) widely exist in genomes of many organisms including humans [1] and human microbiomes [2], and

*Corresponding authors.

E-mail: heshunmin@ibp.ac.cn (He S), zhangp@ibp.ac.cn (Zhang P), crs@ibp.ac.cn (Chen R).

#Equal contribution.

Peer review under responsibility of Beijing Institute of Genomics, Chinese Academy of Sciences / China National Center for Bioinformation and Genetics Society of China.
<https://doi.org/10.1016/j.gpb.2021.09.002>

some are able to be translated into small proteins [3–5]. Small proteins are proteins with less than 100 amino acids, which may be derived from untranslated regions (UTRs) of mRNAs [6] or non-coding RNAs (ncRNAs) [7,8] including primary microRNAs (pri-miRNAs) [9,10], long ncRNAs (lncRNAs) [11], and circular RNAs (circRNAs) [12]. Small proteins were usually missed in previous coding sequence annotation, while their significance has been revealed in current years for diverse functions [13], such as embryonic development [14,15], cell apoptosis [16], muscle contraction [17], and antimicrobial activity [18]. Some small proteins play roles in multiple diseases [19,20] including tumors [9,11,12]. Despite the abundance of sORFs in genome, the number of well-studied small proteins is very limited. Annotation of numerous small proteins will contribute to studies on various physiological and pathological processes.

Identification of small proteins at proteomic level is challenging. Mass spectrometry (MS) can provide direct evidence of small proteins, but it relies much on the coverage of existing libraries, which mainly focus on large proteins rather than small proteins. Protease cleavage sites are lacking in small proteins due to the limited length. Besides, small proteins are usually of low abundance, and tend to be filtered out during enrichment process [21]. Ribosome profiling (also named as ribosomal footprinting or Ribo-seq) provides a more sensitive way for global detection of translation events based on the deep sequencing of ribosome-protected mRNA fragments (RPFs) [22,23], which allows for identifying the location of translated ORFs and translation initiation sites (TISs), the distribution of ribosomes on mRNA, and the speed of translating ribosomes [24]. Reference libraries for MS can also be constructed with Ribo-seq data. The regular Ribo-seq (rRibo-seq) utilizes cycloheximide (CHX) [25], a drug bound at the ribosome E-site [26], as a translation elongation inhibitor to freeze the translating ribosomes. Translation is principally regulated at the initiation stage. Translation initiation sequencing (TI-seq) is a variation of rRibo-seq technique that uses different translation inhibitors, usually lactomidomycin (LTM) [25] or harringtonine (HARR) [27], which can induce ribosome stasis at TISs. TI-seq enables the global mapping of TISs, and is more accurate in prediction of non-ATG start codons. Many sORFs are proved to use non-classical ATG start codons [28], which is also an important mechanism for generating protein isoforms [29,30]. rRibo-seq data usually show clear triplet periodicity [26]. Different computational analysis strategies [31–38] have been developed to identify translated sequences using Ribo-seq data.

Emerging evidence shows that many upstream ORFs (uORFs) act in *cis* to regulate the translation of downstream ORFs by leaky scanning [39], reinitiation [40], and ribosome stalling [41]. Recently, variants creating new upstream start codons or disrupting stop sites of existing uORFs (uORF-perturbing) are found to be under strong negative selection [42]. uORF-perturbing variants have

been demonstrated as an under-recognized functional class that contribute to human disease.

Given the great importance of small proteins, in-depth investigations of small proteins across various species are in need. SmProt is dedicated to integrating knowledge of small proteins translated from various sources, especially for those from UTRs and ncRNAs. The annotation information and functional sections in the current release are much richer than those in the initial release [43], and the data volume and reliability are also greatly improved.

Data collection and processing

Data sources

rRibo-seq and TI-seq datasets derived from diverse tissues/cell lines were collected from GEO [44] and European Nucleotide Archive (ENA) [45] databases. The latest reference genomes and gene annotations were download from Ensembl [46], GENCODE [47], and NCBI-Genome database. Whole-genome sequencing (WGS) variants were collected from various websites. The construction pipeline of SmProt is summarized in **Figure 1**.

Ribo-seq data processing

The fastq files of 547 Ribo-seq datasets were downloaded from GEO and ENA databases. Each dataset was checked manually to confirm the sequencing adapters. The adapters were removed using cutadapt 1.18 [48] and only reads with 25–35 bp in length were retained. Then the sequences were mapped to the latest genome using STAR 2.5.2a [49] using EndToEnd mode with allowance of up to 2 mismatches.

Ribo-seq quality and P-site offsets were assessed by Ribo-TISH [34] quality module. For TI-seq data, more attention was put on TIS quality (-t). Manual checks were then carried out to verify offset values and eliminate datasets without obvious triplet periodicity. After the quality control, 419 Ribo-seq datasets (**Table S1**) were retained.

Translated ORFs were predicted by Ribo-TISH predict module. Biological and technical duplication data under the same treatment in one dataset were merged. Minimum amino acid length of candidate ORFs was set to 5. Considering both ATG and near-cognate start codons (with one base different from ATG), rRibo-seq datasets using only CHX without matched TI-seq data were analyzed twice. One is prediction of ORFs with canonical ATG start codon, the other is prediction of ORFs with near-cognate start codons (--alt). Preferring data evidence instead of prior assumption in our database, only the best frame test results from multiple candidate start codons in the same ORF were reported (--framebest). For datasets containing TI-seq data, alternative start codons were included (--alt), and different parameters were set for LTM-based TI-seq and

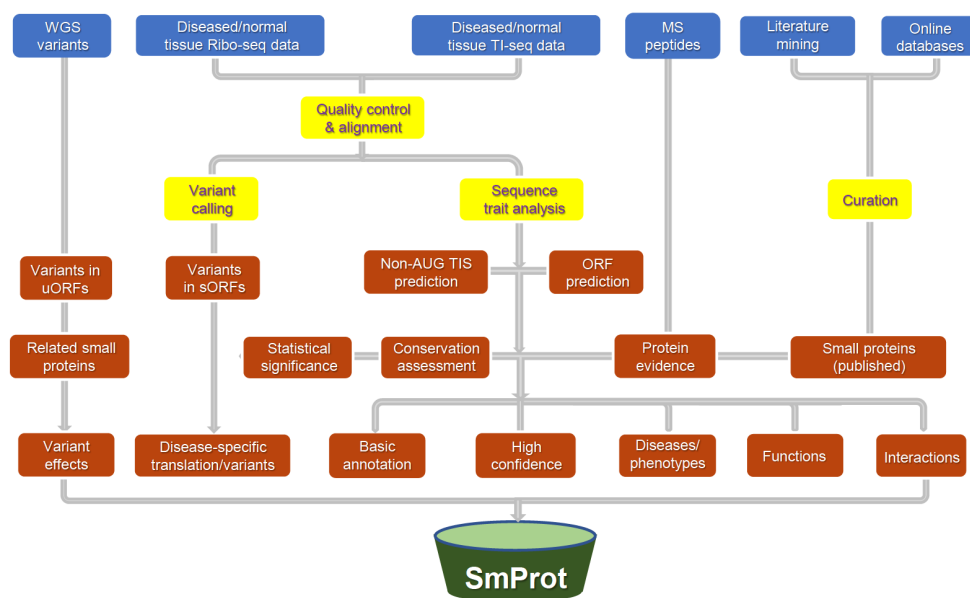


Figure 1 Construction of SmProt

Items in blue background represent data sources. Items in yellow background represent management processes. Items in red background represent results. WGS, whole genome sequencing; MS, mass spectrometry; TIS, translation initiation site; ORF, open reading frame; sORF, small ORF; uORF, upstream ORF.

HARR-based TI-seq (--harr).

sORFs with less than 100 amino acids were filtered from the prediction results above. Furthermore, we removed some prediction results that may be supported by RPFs from other classic proteins with more than 100 amino acids. These include ORFs marked as *known* (i.e., TIS annotated in another transcript), *CDSFrameOverlap* (i.e., ORF overlapping with annotated CDS in another transcript in the same reading frame), and *Truncated* (i.e., ORF as a part of annotated CDS in the same transcript) without translation initiation evidence (i.e., no significant results identified from paired TI-seq datasets).

In-frame reads of sORFs were counted and normalized by library sequencing depth (in-frame total reads count) and sORF length, a similar method with reads per kilobase per million mapped reads (RPKM) in RNA-seq but using ribosome profiling data that represents the translation levels.

Finally, 3,060,793 records (i.e., unmerged primary results from all datasets and tissues) were retained. Results with the identical genome loci in one species were merged as the same small protein, generating 577,206 unique IDs, while information derived from multiple datasets were retained, a similar integration method as for piRBase [50].

Variants from ribosome profiling data

We performed germline variants detection on 96 human ribosome profiling datasets, referring to the workflow for processing RNA data for germline short variant discovery with GATK v4.1.8 [51–54]. Duplicate reads were identified using MarkDuplicates tool after alignment, then reads with

unidentified nucleotide (N) in Cigar were split using SplitNCigarReads tool. Base quality score recalibration was carried out based on true sites in training sets using BaseRecalibrator tool and applied using ApplyBQSR tool. Variants were called individually in each sample using the HaplotypeCaller tool. Variants with QualByDepth (QD) < 2 were removed using VariantFiltration tool. Germline single nucleotide variants (SNVs) were linked to small proteins in SmProt according to genomic positions.

Variants from WGS data

Variants from 1KGP3 [55], GAS [56], TOPMed [57], gnomAD3 [42,58], and NyuWa [59] were collected. VCF files were lifted over from old genome version to GRCh38 using LiftoverVcf tool of GATK with allowance to recover swapped reference and alternative alleles.

Variants in 5' UTRs were evaluated for their effects on uORFs using VEP [60] with plugin UTRannotator [42,61], and classified by their functional consequences. Translation evidence of uORFs was based on small proteins recorded in SmProt.

Disease-specific small proteins

Small proteins identified only from diseased cell lines/tissues but not from corresponding normal cell lines/tissues were predicted as disease-specific translation events. If there were matched data of normal and diseased groups in the same dataset, small proteins derived uniquely from diseased group were screened as disease-specific ones. If

there was no matched control group in the same dataset, the same type of healthy tissue/cell line in other datasets were used as control. If there was no matched same tissue/cell line, all data from diverse normal tissues/cell lines were merged for comparisons (Table S2), and small proteins identified only from the diseased cell lines/tissues were predicted as tissue-specific. Disease-specific or tissue-specific translation events require Ribo *P* value in disease groups lower than 0.01 while similar proteins with different TISs at the same loci in control group not detected (Ribo *P* value higher than 0.05).

SNVs in diseased cell lines/tissues derived from ribosome profiling data and located within the genomic region of small proteins were regarded as diseased variant sets. SNVs detected only in diseased variant sets but not in normal sets were predicted as disease-specific SNVs. SNVs in corresponding normal cell lines/tissues (Table S2) derived from ribosome profiling data were combined with all variants derived from multiple WGS projects, as control variant sets for comparison.

Function domain prediction

Besides function of small proteins collected from literature mining, we used InterProScan [62] to predict function domain of small proteins, which focuses on combination of protein family membership and the functional domains/sites, and has been extensively used by genome sequencing projects and the UniProt Knowledgebase [63]. Default

thresholds and additional parameters *-goterms -pa* were adopted for gene oncology and pathway annotations.

PhyloCSF calculation

Pre-calculated BigWig data of PhyloCSF [64] scores at each base across the whole genome were downloaded from Broad Institute (<https://data.broadinstitute.org/compbio1/PhyloCSFtracks/>), and the score for genomic region of each small protein was extracted with our script using pyBigWig (<https://github.com/deeptools/pyBigWig>).

Database implementation

Database website was organized with HTML (<https://html.spec.whatwg.org/>), JavaScript (<https://www.javascript.com/>), PHP (<https://www.php.net/>), and MYSQL (<https://www.mysql.com/>). UCSC Genome Browser (<http://genome.ucsc.edu/>) was used to visualize the small proteins and variants. NCBI BLAST (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) was used for sequence similarity searches.

Database content and usage

Overview

SmProt was constructed by pipeline described in Figure 1. Multiple ways were provided to search, browse, visualize, and study small proteins (Figure 2). Small proteins were

Search small proteins

The screenshot shows the search interface with three main sections:

- ID Search:** A form to search by species (Human) and ID type (small protein ID) with a text input for 'SPRCHSA197474' and a 'Submit' button.
- Location Search:** A form to search for proteins overlapping with a location, including species (Human), chromosome (chr10), start location (61034338), and stop location (61959438).
- Sequence Search:** A section with a 'go to BLAST module' button.

Browse small proteins by related diseases

The screenshot shows the 'Browse small proteins by related diseases' section. It includes a dropdown for 'Disease' (set to 'All'), a 'Detected' dropdown (set to 'Predicted'), and a 'Start Codon' dropdown (set to 'ATG'). A list of diseases is shown on the left, and a table of results is displayed below.

Disease	Detected	Start Codon
SPRCHS	Predicted by Ribo-seq	ATG
SPRCHS	Predicted by Ribo-seq	ATG
SPRCHS141390	Predicted by Ribo-seq	ATG

Browse variants related to small/upstream ORFs

The screenshot shows the 'Browse variants related to small/upstream ORFs' section. It includes a 'Variant Type' dropdown (set to 'All'), an 'Effect' dropdown (set to 'All'), and a 'Data Source' dropdown (set to 'WGS'). A table of variants is displayed below.

Variant	Gene	Distance to CDS (bp)	Variant Type
9-76384190-A-C	RFK	21	uAUG_gained
19-43934847-T-G	ZNF45	.	uSTOP_lost
9-122264840-A-C	MRRF	.	uSTOP_lost
9-34336268-C-T	NUD2	226	uAUG_gained

Visualize small proteins and variants in genome browser

The screenshot shows the 'Visualize small proteins and variants in genome browser' section. It displays a genomic track for chromosome 10, showing the location of small proteins and variants. The track includes a scale bar, a gene track, and a variant track. A 'Genes and Gene Predictions' section at the bottom allows users to manually change tracks to show or hide.

Figure 2 Usage of SmProt

SmProt provides multiple ways to search, browse, and visualize small proteins, as well as related diseases and variants.

found mainly from rRibo-seq and TI-seq data. All information for small proteins from different data sources and datasets were integrated. General information for small proteins was provided such as sequence, mass, location, blocks, tissue or cell line, predicted functions, conservation, and multiple IDs including small protein ID, Ensembl ID, and NONCODE [65] ID. Translation level (in frame counts and Ribo RPKM) of small proteins identified from each dataset and record was provided. Details for their related variants and diseases were also provided (Figure 3). SmProt now has 638,958 unique small proteins and 3,165,229 small protein records in total (Table 1; Table S3).

Reliability of small proteins

SmProt emphasizes reliability of small proteins, which is ensured mainly by the significance of 3-nt periodicity in RPF P-site profile.

Firstly, we constructed new pipeline based on independently

published toolkit Ribo-TISH [34], which allows for accurate detection of ORFs and TISs using rRibo-seq and TI-seq. Ribo-TISH uses rank sum test to detect 3-nt periodicity, and uses negative binomial test to detect TISs, which outperforms other established methods in prediction accuracy.

Secondly, in addition to the quality control based on Ribo-TISH quality module, manual checks were also carried out to ensure clear triplet periodicity and unambiguous offset of Ribo-seq data, which further eliminates noises.

Thirdly, we provided several evaluations as supporting evidence. These include 1) *P* values of small proteins called from multiple ribosome profiling datasets, which indicate the confidence in different samples and conditions; 2) PhyloCSF conservation of genomic regions, which reflects coding potential; and 3) peptide evidence derived from mass spectrum data. All these lines of evidences are exhibited in the small protein page. Moreover, a set with evidence of both translation events and protein fragments is provided on download page.

General annotation of small proteins

General Information						
Small Protein ID	SPROHSA193481					
Organism	Human (<i>Homo sapiens</i>)					
Small Protein Sequence	MATRSGGTLVLVGLGSEMTTVPPLHAAIREVDIKGVFRYCNTPWPAWSMLASKSVNWKPLVHRFLEKALEAFETFKKGLKMLKCDPSDQNP*					
RNA Sequence	ATGGCCACTCGCTGGTGGGCCCTCGTGCCTGTGGGGCTGGGCTCTGAGATGACCACCGTACCCTACTGCATGCAGCATCCGGGAG					
Protein Length	95					
Start Codon	ATG					
Location	chr15:44827055-44832239					
Blocks	44827055-44827221, 44828449-44828571, 44830236-44830239					
Mean PhyloCSF	-2.35629554385					
Data Source	Ribosome profiling; Literature					
Related Genes	ENSG00000259479; SORD2P; ENSG00000259187; AC122108.1; NONHSAG016753; NONHSAG016754					
Mass (Da)	mono: 10478.6, avg: 10485.3					

Information from other sources

Mass Spectrometry Information						
MS ID	Seq	Length	Chr	Start	Stop	Strand
MSHSA415549	SGGTLVLVGLGSEMTTVPPLHAAIR	25	chr15	44828488	44828562	-

Function domain prediction

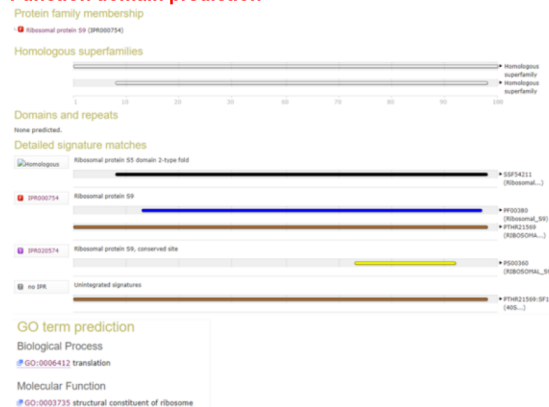


Figure 3 Contents of SmProt

Detailed information for small proteins are provided, including general annotation, information from ribosome profiling data, literature, other databases, MS, function domain prediction, related diseases, and related variants from WGS projects, as well as corresponding effects.

Detailed information of the small proteins in each dataset

Ribosome profiling										
Ribo-seq ID	Ensembl Transcript ID	Symbol	Gene Type	TIS Type	Ribo P value	TIS P value	Start On Trans	Stop On Trans	In-frame Count	Ribo-seq RPKM
SRR4045276	ENST00000564140.5	SORD2P	Transcribed_unprocessed_pseudogene	Novel	2.541E-12	None	198	489	76	61.8553950
SRR4045276_alt	ENST00000564140.5	SORD2P	Transcribed_unprocessed_pseudogene	Novel	2.541E-12	None	198	489	76	61.8553950
SRR320821	ENST00000564140.5	SORD2P	Transcribed_unprocessed_pseudogene	Novel	0.000205	None	198	489	15	11.4579517
Min Ribo P value	7.903E-24									
Min TIS P value	None									
Ribo-seq ID	Cell or Tissue	Phenotype	Ribo-seq Source Details			PMID				
GSE83493_alt	HeLa S3	Cervical cancer	GSM2204389; GSM2204390; GSM2204391; GSM2204392			28460002;				
GSE94454	Huh7	Hepatoma	GSM2476003; GSM2476004; GSM2476005			28323820;				
GSE94454_alt	Huh7	Hepatoma	GSM2476003; GSM2476004; GSM2476005			28323820;				
GSE94613_1	MOLM13	Acute myeloid leukemia	GSM2481052			29186125;				

Related Small Proteins with Different TISs

Related Small Proteins with Different TISs				
SmProt ID	Small Protein Length	Start Codon	Strand	Blocks
SPROHSA117403	12	ATC	+	45073491-45073530
SPROHSA387281	16	TTG	-	44827055-44827106
SPROHSA147008	79	ATG	-	44827055-44827221, 44828449-44828523
SPROHSA302475	83	CTG	-	44827055-44827221, 44828449-44828535

Variants on RNA sequence of the small protein

Related Variants			
Variant ID	Consequence to sORF	rsID	Ribo-seq ID
12-96334805-G-A	Non-Synonymous p.T111	-	SRR3208921

Sources and effects of variants

Data sources			5'UTR Effect	
Source	Allele Count	Allele Frequency	Variant Type	uAUG_gained
gnomAD3	143170	0.99873	Gene	RPK
1KGP3	5005	0.9994	Context	GGGAATG
TOPMed	125408	0.99873	Kozak Sequence	CCCATGC
GAsP	3474	0.999	Kozak Strength	Weak
NyuWa	5998	1	Effect	CDS_elongated
Ribosome profiling	SRR2818787; SRR2818791; SRR3208885; SRR3208921		Distance to CDS	21

Table 1 Statistics of unique small proteins in SmProt

Source	Start codon	Human	Mouse	Fruit fly	Rat	<i>C. elegans</i>	Yeast	<i>E. coli</i>	Zebrafish	All species examined
Ribo-seq	ATG	70,931	48,909	5269	3560	4334	4535	1881	1924	141,343
	Near-cognate codons	229,653	133,037	29,679	9910	9894	12,339	10,004	1347	435,863
Literature	ATG and near-cognate codons	38,157	8875	22,228	163	4	355	296	3612	73,690
Databases	ATG and near-cognate codons	786	797	100	271	120	336	955	64	3429
MS	ATG and near-cognate codons	768	51	66	38	0	3	0	1	927
All IDs examined	ATG and near-cognate codons	327,995	189,433	56,574	13,829	14,255	17,312	12,881	6679	638,958

Note: Small protein families from human microbiomes are not included. Near-cognate codons refer to non-ATG start codons that differ from the canonical ATG start codon by a single base but are able to initiate translation, such as TTG, GTG, CTG, AAG, AGG, ACG, ATA, ATT, and ATC. ID refers to a unique entry with identical genomic loci in one species. Ribo-seq, ribosome profiling; MS, mass spectrometry.

In addition, information of small protein derived from multiple sources is also integrated in small protein information page.

Variants related to small proteins

In total, 25,475 variants located on translated sORFs were provided, which are on display in the related small protein page. Given that uORF-perturbing variants are likely to impact translation of downstream proteins [42], variants from multiple WGS projects and ribosome profiling data were evaluated for their effects on uORFs. These include creating a new start codon ATG, removing an existing start codon ATG, creating a new stop codon within an existing uORF, removing the stop codon of an existing uORF, and creating a frameshift mutation in an existing uORF, which can be found at variants page.

Disease-specific small proteins

Disease-specific small proteins are potential candidates of molecular markers or targets for diagnosis and treatment. Disease-specific translation events as well as disease-specific SNVs of small proteins in 16 types of diseases were identified (see "Data collection and processing" section) (Table S4). Besides, small proteins that have been verified experimentally in certain diseases were also documented through literature mining.

Human microbiome small proteins

Over 4000 conserved small protein families identified from human microbiomes were collected [2]. A new section *HumanMicroBio* was created to integrate and display selected information of these small protein families.

Other sources

We use a set of keywords (File S1) to search articles about small proteins in PubMed database. High-confidence small proteins in CCDS [66] and Swiss-Prot [67] were also

integrated. Literature mining is processed in stages, and the newly-published data from other sources will be released continuously after completion of manual review and curation.

Function domain prediction

For successfully predicted functions of small proteins derived from ribosome profiling and literature mining, SmProt provides graph for visualization and prediction details including Gene Ontology (GO) and pathway annotations. Users can choose *predicted functions* on *Browse* page to filter the results with function domain prediction.

Inner BLAST

The abundant small proteins across multiple species allow for sequence similarity searches at both nucleotide and protein levels. Users can search for sequences of interests using BLASTp and BLASTx (NCBI BLAST 2.2.24 release) online.

Visualization using UCSC Genome Browser

SmProt incorporates UCSC Genome Browser [68] to visualize all the information including genomic loci of small proteins, as well as variants from ribosome profiling data and multiple WGS projects related to small proteins, MS data, and gene annotations. The latest genome versions including hg38, mm10, rn6, dm6, ce11, sacCer3, and danRer11 were provided.

Comparison with other databases

SmProt currently includes 419 Ribo-seq datasets derived from 116 cell lines/tissues, compared to 60 datasets derived from 37 cell lines/tissues in the initial version. The number of small protein records identified from ribosome profiling in the current release is 60 times that of the initial release (3 million vs. 0.05 million). The current release of SmProt

combined a large amount of duplicate records in the initial release [43], and Ribo-seq analysis pipeline was optimized to ensure the reliability of our results. Variants in translated sORFs identified from Ribo-seq data as well as uORF-perturbing variants identified from WGS projects were provided. Disease-specific small proteins may provide new perspectives for clinical studies.

Currently, there are a few databases for small proteins such as ARA-PEPs [69], PsORF [70], and sORFs.org [71]. ARA-PEPs and PsORF only harbor small proteins in plants. sORFs.org developed simple inner TIS-calling algorithm not based on triplet periodicity, which should be the most important feature of Ribo-seq. SmProt emphasizes high confidence using our Ribo-TISH pipeline that is more accurate than previous methods. In total, 419 Ribo-seq datasets have been analyzed in SmProt, while there were only 78 Ribo-seq datasets in sORFs.org. Moreover, SmProt pays special attention to function, variants, and related diseases of small proteins. Furthermore, WGS data resources are also integrated in SmProt, which are not covered in other databases.

Other proteomic databases such as UniProt, neXtProt [72], and OpenProt [73] are not specifically designed for small proteins. neXtProt only harbors proteins of humans while SmProt harbors small proteins of 8 species. Similar to SmProt, OpenProt also used ribosome profiling and mass spectrum to predict proteins including some small proteins longer than 30 amino acids. Nonetheless, SmProt has analyzed many more ribosome profiling datasets (419), which are about 5 times that in OpenProt (87), and provides information for small proteins longer than 5 amino acids.

Conclusion

In brief, SmProt integrates small proteins from large amount of ribosome profiling data, and provides more abundant details. We strongly believe that SmProt will provide valuable and accurate information on small proteins for scientific community. Moreover, SmProt provides a new resource for users interested in functional and mechanistic studies, and a reference for construction of MS libraries of small proteins.

Data availability

SmProt is publicly available at <http://bigdata.ibp.ac.cn/SmProt/>.

CRedit author statement

Yanyan Li: Conceptualization, Methodology, Investigation, Formal analysis, Data curation, Writing - original draft, Software, Visualization. **Honghong Zhou:** Investigation,

Data curation, Funding acquisition. **Xiaomin Chen:** Investigation, Data curation. **Yu Zheng:** Data curation, Software, Visualization. **Quan Kang:** Software, Visualization. **Di Hao:** Data curation, Software. **Lili Zhang:** Visualization. **Tingrui Song:** Visualization. **Huaxia Luo:** Writing - review & editing. **Yajing Hao:** Writing - review & editing. **Runsheng Chen:** Resources, Supervision, Funding acquisition. **Peng Zhang:** Conceptualization, Methodology, Investigation, Software, Writing - review & editing, Visualization, Project administration, Funding acquisition. **Shunmin He:** Conceptualization, Methodology, Resources, Investigation, Writing - review & editing, Supervision, Funding acquisition. All authors have read and approved the final manuscript.

Competing interests

The authors have declared no competing interests.

Acknowledgments

This work was supported by the National Key R&D Program of China (Grant No. 2016YFC0901702); National Natural Science Foundation of China (Grant Nos. 81902519, 91940306, 31871294, 31701117, and 31970647); the National Key R&D Program of China (Grant Nos. 2017YFC0907503, 2016YFC0901002, and 2018YFA0106901); the Strategic Priority Research Program of Chinese Academy of Sciences (Grant No. XDB38040300); the 13th Five-year Informatization Plan of Chinese Academy of Sciences (Grant No. XXH13505-05); Special Investigation on Science and Technology Basic Resources, Ministry of Science and Technology, China (Grant No. 2019FY100102); and the National Genomics Data Center, China. We thank Center for Big Data Research in Health (<http://bigdata.ibp.ac.cn/>), Institute of Biophysics, Chinese Academy of Sciences, for supporting data analysis and computing resource. We thank Prof. Yiwen Chen from The University of Texas MD Anderson Cancer Center, USA for thoughtful discussions and valuable comments on the database.

Supplementary material

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.gpb.2021.09.002>.

ORCID

0000-0001-5256-6696 (Yanyan Li)
 0000-0001-7409-3092 (Honghong Zhou)
 0000-0002-0633-2984 (Xiaomin Chen)
 0000-0003-4936-8407 (Yu Zheng)

0000-0001-6790-5259 (Quan Kang)
 0000-0003-0082-0730 (Di Hao)
 0000-0002-3601-0150 (Lili Zhang)
 0000-0003-2967-7704 (Tingrui Song)
 0000-0001-9944-0345 (Huaxia Luo)
 0000-0003-1384-4176 (Yajing Hao)
 0000-0001-6049-8347 (Runsheng Chen)
 0000-0001-9303-1639 (Peng Zhang)
 0000-0002-7294-0865 (Shunmin He)

References

- [1] Basrai MA, Hieter P, Boeke JD. Small open reading frames: beautiful needles in the haystack. *Genome Res* 1997;7:768–71.
- [2] Sberro H, Fremin BJ, Zlitni S, Edfors F, Greenfield N, Snyder MP, et al. Large-scale analyses of human microbiomes reveal thousands of small, novel genes. *Cell* 2019;178:1245–59.e14.
- [3] Bazzini AA, Johnstone TG, Christiano R, Mackowiak SD, Obermayer B, Fleming ES, et al. Identification of small ORFs in vertebrates using ribosome footprinting and evolutionary conservation. *EMBO J* 2014;33:981–93.
- [4] Smith JE, Alvarez-Dominguez JR, Kline N, Huynh NJ, Geisler S, Hu W, et al. Translation of small open reading frames within unannotated RNA transcripts in *Saccharomyces cerevisiae*. *Cell Rep* 2014;7:1858–66.
- [5] van Heesch S, Witte F, Schneider-Lunitz V, Schulz JF, Adami E, Faber AB, et al. The translational landscape of the human heart. *Cell* 2019;178:242–60.e29.
- [6] Calvo SE, Pagliarini DJ, Mootha VK. Upstream open reading frames cause widespread reduction of protein expression and are polymorphic among humans. *Proc Natl Acad Sci U S A* 2009;106:7507–12.
- [7] Zhu S, Wang J, He Y, Meng N, Yan GR. Peptides/proteins encoded by non-coding RNA: a novel resource bank for drug targets and biomarkers. *Front Pharmacol* 2018;9:1295.
- [8] Li LJ, Leng RX, Fan YG, Pan HF, Ye DQ. Translation of non-coding RNAs: focus on lncRNAs, pri-miRNAs, and circRNAs. *Exp Cell Res* 2017;361:1–8.
- [9] Fang J, Morsalin S, Rao VN, Reddy ESP. Decoding of non-coding DNA and non-coding RNA: pri-micro RNA-encoded novel peptides regulate migration of cancer cells. *J Pharm Sci Pharmacol* 2017;3:23–7.
- [10] Razoooky BS, Obermayer B, O'May JB, Tarakhovskiy A. Viral infection identifies micropeptides differentially regulated in smORF-containing lncRNAs. *Genes (Basel)* 2017;8:206.
- [11] Huang JZ, Chen M, Chen D, Gao XC, Zhu S, Huang H, et al. A peptide encoded by a putative lncRNA HOXB-AS3 suppresses colon cancer growth. *Mol Cell* 2017;68:171–84.e6.
- [12] Zhang M, Zhao K, Xu X, Yang Y, Yan S, Wei P, et al. A peptide encoded by circular form of LINC-PINT suppresses oncogenic transcriptional elongation in glioblastoma. *Nat Commun* 2018;9:4475.
- [13] Couso JP, Patraquim P. Classification and function of small open reading frames. *Nat Rev Mol Cell Biol* 2017;18:575–89.
- [14] Freyer L, Hsu CW, Nowotschin S, Pauli A, Ishida J, Kuba K, et al. Loss of Apela peptide in mice causes low penetrance embryonic lethality and defects in early mesodermal derivatives. *Cell Rep* 2017;20:2116–30.
- [15] Galindo MI, Pueyo JI, Fouix S, Bishop SA, Couso JP. Peptides encoded by short ORFs control development and define a new eukaryotic gene family. *PLoS Biol* 2007;5:e106.
- [16] Guo B, Zhai D, Cabezas E, Welsh K, Nouraini S, Satterthwait AC, et al. Humanin peptide suppresses apoptosis by interfering with Bax activation. *Nature* 2003;423:456–61.
- [17] Anderson DM, Anderson KM, Chang CL, Makarewich CA, Nelson BR, McAnally JR, et al. A micropeptide encoded by a putative long noncoding RNA regulates muscle performance. *Cell* 2015;160:595–606.
- [18] Knappe D, Goldbach T, Hatfield MPD, Palermo NY, Weinert S, Sträter N, et al. Proline-rich antimicrobial peptides optimized for binding to *Escherichia coli* chaperone DnaK. *Protein Pept Lett* 2016;23:1061–71.
- [19] Wen Y, Liu Y, Xu Y, Zhao Y, Hua R, Wang K, et al. Loss-of-function mutations of an inhibitory upstream ORF in the human hairless transcript cause Marie Unna hereditary hypotrichosis. *Nat Genet* 2009;41:228–33.
- [20] Cheng W, Wang S, Mestre AA, Fu C, Makarem A, Xian F, et al. C9ORF72 GGGGCC repeat-associated non-AUG translation is upregulated by stress through eIF2 α phosphorylation. *Nat Commun* 2018;9:51.
- [21] Hsu PY, Benfey PN. Small but mighty: functional peptides encoded by small ORFs in plants. *Proteomics* 2018;18:1700038.
- [22] Ingolia NT, Ghaemmaghami S, Newman JRS, Weissman JS. Genome-wide analysis *in vivo* of translation with nucleotide resolution using ribosome profiling. *Science* 2009;324:218–23.
- [23] Ingolia NT, Lareau LF, Weissman JS. Ribosome profiling of mouse embryonic stem cells reveals the complexity and dynamics of mammalian proteomes. *Cell* 2011;147:789–802.
- [24] Weiss RB, Atkins JF. Translation goes global. *Science* 2011;334:1509–10.
- [25] Schneider-Poetsch T, Ju J, Eyler DE, Dang Y, Bhat S, Merrick WC, et al. Inhibition of eukaryotic translation elongation by cycloheximide and lactimidomycin. *Nat Chem Biol* 2010;6:209–17.
- [26] Calviello L, Ohler U. Beyond read-counts: Ribo-seq data analysis to understand the functions of the transcriptome. *Trends Genet* 2017;33:728–44.
- [27] Ingolia NT, Brar GA, Rouskin S, McGeachy AM, Weissman JS. The ribosome profiling strategy for monitoring translation *in vivo* by deep sequencing of ribosome-protected mRNA fragments. *Nat Protoc* 2012;7:1534–50.
- [28] Lee S, Liu B, Lee S, Huang SX, Shen B, Qian SB. Global mapping of translation initiation sites in mammalian cells at single-nucleotide resolution. *Proc Natl Acad Sci U S A* 2012;109:E2424–32.
- [29] Kochetov AV, Sarai A, Rogozin IB, Shumny VK, Kolchanov NA. The role of alternative translation start sites in the generation of human protein diversity. *Mol Genet Genomics* 2005;273:491–6.
- [30] Oyama M, Kozuka-Hata H, Suzuki Y, Semba K, Yamamoto T, Sugano S. Diversity of translation start sites may define increased complexity of the human short ORFome. *Mol Cell Proteomics* 2007;6:1000–6.
- [31] Calviello L, Mukherjee N, Wyler E, Zauber H, Hirsekorn A, Selbach M, et al. Detecting actively translated open reading frames in ribosome profiling data. *Nat Methods* 2016;13:165–70.
- [32] Fields AP, Rodriguez EH, Jovanovic M, Stern-Ginossar N, Haas BJ, Mertins P, et al. A regression-based analysis of ribosome-profiling data reveals a conserved complexity to mammalian translation. *Mol Cell* 2015;60:816–27.
- [33] Ji Z, Song R, Regev A, Struhl K. Many lncRNAs, 5' UTRs, and pseudogenes are translated and some are likely to express functional proteins. *eLife* 2015;4:e08890.
- [34] Zhang P, He D, Xu Y, Hou J, Pan BF, Wang Y, et al. Genome-wide identification and differential analysis of translational initiation. *Nat Commun* 2017;8:1749.
- [35] Malone B, Atanassov I, Aeschmann F, Li X, Großhans H, Dietrich C. Bayesian prediction of RNA translation from ribosome profiling. *Nucleic Acids Res* 2017;45:2960–72.
- [36] Raj A, Wang SH, Shim H, Harpak A, Li YI, Engelmann B, et al. Thousands of novel translated open reading frames in humans inferred by ribosome footprint profiling. *eLife* 2016;5:e13328.
- [37] Chun SY, Rodriguez CM, Todd PK, Mills RE. SPECTre: a spectral

- coherence-based classifier of actively translated transcripts from ribosome profiling sequence data. *BMC Bioinformatics* 2016;17:482.
- [38] Crappe J, Ndah E, Koch A, Steyaert S, Gawron D, De Keulenaer S, et al. PROTEOFORMER: deep proteome coverage through ribosome profiling and MS integration. *Nucleic Acids Res* 2015;43:e29.
- [39] Wang XQ, Rothnagel JA. 5'-Untranslated regions with multiple upstream AUG codons can support low-level translation via leaky scanning and reinitiation. *Nucleic Acids Res* 2004;32:1382–91.
- [40] Gunišová S, Valášek LS. Fail-safe mechanism of GCN4 translational control—uORF2 promotes reinitiation by analogous mechanism to uORF1 and thus secures its key role in GCN4 expression. *Nucleic Acids Res* 2014;42:5880–93.
- [41] Ishimura R, Nagy G, Dotu I, Zhou H, Yang XL, Schimmel P, et al. Ribosome stalling induced by mutation of a CNS-specific tRNA causes neurodegeneration. *Science* 2014;345:455–9.
- [42] Whiffin N, Karczewski KJ, Zhang X, Chothoni S, Smith MJ, Evans DG, et al. Characterising the loss-of-function impact of 5' untranslated region variants in 15,708 individuals. *Nat Commun* 2020;11:2523.
- [43] Hao Y, Zhang L, Niu Y, Cai T, Luo J, He S, et al. SmProt: a database of small proteins encoded by annotated coding and non-coding RNA loci. *Brief Bioinform* 2018;19:636–43.
- [44] Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomshesky M, et al. NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res* 2013;41:D991–5.
- [45] Silvester N, Alako B, Amid C, Cerdeño-Tarraga A, Clarke L, Cleland I, et al. The European Nucleotide Archive in 2017. *Nucleic Acids Res* 2017;46:D36–40.
- [46] Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, et al. Ensembl 2018. *Nucleic Acids Res* 2018;46:D754–61.
- [47] Frankish A, Diekhans M, Ferreira AM, Johnson R, Jungreis I, Loveland J, et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res* 2019;47:D766–73.
- [48] Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet J* 2011;17:10.
- [49] Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;29:15–21.
- [50] Wang J, Zhang P, Lu Y, Li Y, Zheng Y, Kan Y, et al. piRBase: a comprehensive database of piRNA sequences. *Nucleic Acids Res* 2019;47:D175–80.
- [51] Poplin R, Ruano-Rubio V, DePristo MA, Fennell TJ, Carneiro MO, Van der Auwera GA, et al. Scaling accurate genetic variant discovery to tens of thousands of samples. *bioRxiv* 2018;201178.
- [52] Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, et al. From FastQ data to high-confidence variant calls: the genome analysis toolkit best practices pipeline. *Curr Protocols Bioinformatics* 2013;43:11.10.1–33.
- [53] DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet* 2011;43:491–8.
- [54] McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* 2010;20:1297–303.
- [55] Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An integrated map of structural variation in 2,504 human genomes. *Nature* 2015;526:75–81.
- [56] GenomeAsia100K Consortium. The GenomeAsia 100K Project enables genetic discoveries across Asia. *Nature* 2019;576:106–11.
- [57] Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* 2021;590:290–9.
- [58] Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* 2020;581:434–43.
- [59] Zhang P, Luo H, Li Y, Wang Y, Wang J, Zheng Y, et al. NyuWa Genome resource: a deep whole-genome sequencing-based variation profile and reference panel for the Chinese population. *Cell Rep* 2021;37:110017.
- [60] McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, et al. The Ensembl Variant Effect Predictor. *Genome Biol* 2016;17:122.
- [61] Zhang X, Wakeling M, Ware J, Whiffin N. Annotating high-impact 5'untranslated region variants with the UTRannotator. *Bioinformatics* 2021;37:1171–3.
- [62] Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, et al. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 2014;30:1236–40.
- [63] UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res* 2015;43:D204–12.
- [64] Lin MF, Jungreis I, Kellis M. PhyloCSF: a comparative genomics method to distinguish protein coding and non-coding regions. *Bioinformatics* 2011;27:i275–82.
- [65] He S, Liu C, Skogerboe G, Zhao H, Wang J, Liu T, et al. NONCODE v2.0: decoding the non-coding. *Nucleic Acids Res* 2008;36:D170–2.
- [66] Pujar S, O'Leary NA, Farrell CM, Loveland JE, Mudge JM, Wallin C, et al. Consensus coding sequence (CCDS) database: a standardized set of human and mouse protein-coding regions supported by expert curation. *Nucleic Acids Res* 2018;46:D221–8.
- [67] UniProt Consortium. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* 2019;47:D506–15.
- [68] Haussler M, Zweig AS, Tyner C, Speir ML, Rosenbloom KR, Raney BJ, et al. The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res* 2019;47:D853–8.
- [69] Hazarika RR, De Coninck B, Yamamoto LR, Martin LR, Cammue BPA, van Noort V. ARA-PEPs: a repository of putative sORF-encoded peptides in *Arabidopsis thaliana*. *BMC Bioinformatics* 2017;18:37.
- [70] Chen Y, Li D, Fan W, Zheng X, Zhou Y, Ye H, et al. PsORF: a database of small ORFs in plants. *Plant Biotechnol J* 2020;18:2158–60.
- [71] Olexiouk V, Van Criekinge W, Menschaert G. An update on sORFs.org: a repository of small ORFs identified by ribosome profiling. *Nucleic Acids Res* 2018;46:D497–502.
- [72] Gaudet P, Michel PA, Zahn-Zabal M, Britan A, Cusin I, Domagalski M, et al. The neXtProt knowledgebase on human proteins: 2017 update. *Nucleic Acids Res* 2017;45:D177–82.
- [73] Brunet MA, Brunelle M, Lucier JF, Delcourt V, Levesque M, Grenier F, et al. OpenProt: a more comprehensive guide to explore eukaryotic coding potential and proteomes. *Nucleic Acids Res* 2019;47:D403–10.