

# Unraveling the small proteome of the plant symbiont *Sinorhizobium meliloti* by ribosome profiling and proteogenomics

Lydia Hadjeras<sup>1</sup>, Benjamin Heiniger<sup>2,†</sup>, Sandra Maaß<sup>3,†</sup>, Robina Scheuer<sup>4,†</sup>, Rick Gelhausen<sup>5</sup>, Saina Azarderakhsh<sup>4</sup>, Susanne Barth-Weber<sup>4</sup>, Rolf Backofen<sup>5</sup>, Dörte Becher<sup>3</sup>, Christian H. Ahrens<sup>2,\*</sup>, Cynthia M. Sharma<sup>1,\*</sup>, Elena Evguenieva-Hackenberg<sup>4,\*</sup>

<sup>1</sup>Institute of Molecular Infection Biology, University of Würzburg, 97080 Würzburg, Germany

<sup>2</sup>Molecular Ecology, Agroscope and SIB Swiss Institute of Bioinformatics, 8046 Zurich, Switzerland

<sup>3</sup>Institute of Microbiology, University of Greifswald, 17489 Greifswald, Germany

<sup>4</sup>Institute of Microbiology and Molecular Biology, University of Giessen, 35392 Giessen, Germany

<sup>5</sup>Bioinformatics Group, Department of Computer Science, University of Freiburg, 79110 Freiburg, Germany

\*Corresponding author. Molecular Ecology, Agroscope and SIB Swiss Institute of Bioinformatics, 8046 Zurich, Switzerland. E-mail:

[christian.ahrens@agroscope.admin.ch](mailto:christian.ahrens@agroscope.admin.ch); Institute of Molecular Infection Biology, University of Würzburg, 97080 Würzburg, Germany.

E-mail: [cynthia.sharma@uni-wuerzburg.de](mailto:cynthia.sharma@uni-wuerzburg.de); Institute of Microbiology and Molecular Biology, University of Giessen, 35392 Giessen, Germany.

E-mail: [Elena.Evguenieva-Hackenberg@mikro.bio.uni-giessen.de](mailto:Elena.Evguenieva-Hackenberg@mikro.bio.uni-giessen.de)

<sup>†</sup>These authors contributed equally to the study

## Abstract

The soil-dwelling plant symbiont *Sinorhizobium meliloti* is a major model organism of Alphaproteobacteria. Despite numerous detailed OMICS studies, information about small open reading frame (sORF)-encoded proteins (SEPs) is largely missing, because sORFs are poorly annotated and SEPs are hard to detect experimentally. However, given that SEPs can fulfill important functions, identification of translated sORFs is critical for analyzing their roles in bacterial physiology. Ribosome profiling (Ribo-seq) can detect translated sORFs with high sensitivity, but is not yet routinely applied to bacteria because it must be adapted for each species. Here, we established a Ribo-seq procedure for *S. meliloti* 2011 based on RNase I digestion and detected translation for 60% of the annotated coding sequences during growth in minimal medium. Using ORF prediction tools based on Ribo-seq data, subsequent filtering, and manual curation, the translation of 37 non-annotated sORFs with  $\leq 70$  amino acids was predicted with confidence. The Ribo-seq data were supplemented by mass spectrometry (MS) analyses from three sample preparation approaches and two integrated proteogenomic search database (iPtgxDB) types. Searches against standard and 20-fold smaller Ribo-seq data-informed custom iPtgxDBs confirmed 47 annotated SEPs and identified 11 additional novel SEPs. Epitope tagging and Western blot analysis confirmed the translation of 15 out of 20 SEPs selected from the translome map. Overall, by combining MS and Ribo-seq approaches, the small proteome of *S. meliloti* was substantially expanded by 48 novel SEPs. Several of them are part of predicted operons and/or are conserved from Rhizobiaceae to Bacteria, suggesting important physiological functions.

**Keywords:** Ribosome profiling, proteomics, proteogenomics, small proteins, small open reading frame, *Sinorhizobium meliloti*, Alphaproteobacteria

## Introduction

Over the last two decades, using next-generation sequencing and high throughput OMICS profiling technologies, the genomes of thousands of bacteria have been assembled. Moreover, the transcriptomes and proteomes of many of them have been analyzed under different conditions, with the aim of gaining insights into the genetic and molecular basis of their biology. Despite this wealth of data, information about small open reading frame (sORF)-encoded proteins (SEPs), which are proteins with less than 50 or 100 amino acids (aa), is scarce (Storz et al. 2014, Duval and Cossart 2017, Hemm et al. 2020, Orr et al. 2020, Gray et al. 2022). Recently, the small proteomes of eukaryotes, bacteria, and viruses have gained expanding interest, as a growing number of small proteins have been demonstrated to fulfill important physiological functions, such as in cell division, metabolism, transport, sig-

nal transduction, spore formation, cell communication, cellular stress responses, and virulence (Storz et al. 2014, Duval and Cossart 2017, Khitun and Slavoff 2019, Hemm et al. 2020, Melior et al. 2020, Patraquim et al. 2020, Aoyama et al. 2022, Song et al. 2022). Therefore, cataloging the full complement of small proteins is critical in achieving a more comprehensive and accurate description of the proteomes of bacterial model organisms and their potential functions.

Small protein identification is difficult due to several technical challenges. For instance, SEPs are difficult to detect using SDS-PAGE or mass spectrometry (MS) for various technical reasons (Storz et al. 2014, Ahrens et al. 2022, Fijalkowski et al. 2022). Limitations of standard shotgun proteomics workflows at the sample preparation, protease digestion, liquid chromatography, MS data acquisition, and bioinformatic data analysis steps affect compre-

Received: November 15, 2022. Revised: February 8, 2023. Accepted: March 7, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of FEMS. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

hensive MS-based SEP identification (Cassidy et al. 2021, Ahrens et al. 2022). Furthermore, variable length thresholds were typically used in the genome annotation step to minimize the number of spurious ORF predictions. As a result, sORFs encoding truly expressed small proteins are often missing from genome annotations (Storz et al. 2014, Hahn et al. 2016). Meanwhile, various strategies to achieve extensive proteome coverage of the notoriously under-represented classes of small and membrane proteins (novel small proteins are often membrane associated) have been applied for prokaryotes (Omasits et al. 2013, Zhang et al. 2013, Wiśniewski 2016). Methods to enrich bacterial SEPs in samples are further improved, for example, with the use of small pore-sized solid-phase materials (Cassidy et al. 2019, Bartel et al. 2020, Petruschke et al. 2020), and digestion with alternative/multiple proteases has been performed to increase the number of identified SEPs (Bartel et al. 2020, Kaulich et al. 2021, Petruschke et al. 2021). The obtained mass spectra are usually assigned to peptide or protein sequences by matching the determined fragment ion masses to the predictions derived from a sequence database (DB). Therefore, only peptides with sequences available in the protein search DB can be identified. Consequently, custom protein search DBs that try to capture the entire coding potential of prokaryotic genomes have been proposed, such as integrated proteogenomic search DBs (iPtgxDBs). They integrate and consolidate the differences among existing reference genome annotations, *ab initio* gene predictions, and a modified six-frame translation that considers alternative start sites, thereby enabling the detection of novel proteins, including SEPs (Omasits et al. 2017). Thus, proteogenomic studies that combine results from SEP-optimized MS data searched with iPtgxDBs or other custom search DBs and ribosome profiling (Ribo-seq) have great potential to detect more comprehensive and accurate compendia of novel small proteins.

Ribo-seq is a powerful method to study and annotate translomes globally, including sORFs (Ingolia 2016, Vazquez-Laslop et al. 2022). Compared with MS-based proteomics, Ribo-seq has the advantage of higher sensitivity for detecting translated ORFs (Storz et al. 2014, Duval and Cossart 2017, Hemm et al. 2020, Orr et al. 2020, Venturini et al. 2020, Ahrens et al. 2022, Gray et al. 2022). Ribo-seq relies on deep sequencing of approximately 30-nt-long 'footprint' regions of the mRNA bound by the ribosome during translation and protected against nuclease digestion. In addition to providing a global picture of translated mRNAs in the cell, Ribo-seq also reveals the specific location on the mRNA where the ribosome was bound, allowing the mapping of ORFs. For this, cells are lysed under certain conditions, allowing for the 'freezing' of ribosomes on mRNAs. mRNA parts that are not protected by the ribosomes are then digested to generate ribosome footprints that are sequenced and mapped to the genome (Ingolia et al. 2009, Ingolia 2016). While Ribo-seq-based detection of translated mRNA works well for eukaryotic cells at single codon resolution, this method is difficult to utilize for prokaryotes (Mohammad et al. 2019, Glaub et al. 2020, Vazquez-Laslop et al. 2022). Nevertheless, adapting and refining the Ribo-seq method has enabled the detection of many new, translated sORFs and corresponding SEPs not only in *Escherichia coli* but also in several other bacterial species and in halophilic archaea (Meydan et al. 2019, Mohammad et al. 2019, Weaver et al. 2019, Gelsinger et al. 2020, Vazquez-Laslop et al. 2022, Hadjeras et al., 2023). However, for many bacterial model organisms, Ribo-seq data are still missing, as the protocols typically have to be adapted and optimized for each bacterial organism (Storz et al. 2014, Duval and Cossart 2017, Hemm et al. 2020, Orr et al. 2020, Venturini et al. 2020, Gray et al. 2022).

*Sinorhizobium meliloti* is an agriculturally important bacterial species that lives in soil and can fix molecular nitrogen in symbiosis with legume plants (Jones et al. 2007). Due to its versatile lifestyle and ecological relevance, it is a major model organism for studying gene regulation in Alphaproteobacteria. In addition, its relatively close relationship to pathogens of the genus *Brucella* makes *S. meliloti* an attractive model for host-pathogen research (Marlow et al. 2009). Several OMICS datasets are available for *S. meliloti* 2011 and its sibling, *S. meliloti* 1021, the first strain of this species with a sequenced genome (Galibert et al. 2001). These comprise proteomics (Djordjevic 2004, Barra-Bily et al. 2010, Sobrero et al. 2012, Marx et al. 2016) and transcriptomic datasets, including differential RNA-seq that enables the annotation of transcription start sites, 5'- and 3'-UTRs, and novel transcripts (Becker et al. 2004, Sallet et al. 2013, Schlüter et al. 2013). The *S. meliloti* 2011 6.7 Mb genome harbors a 3.66 Mb chromosome and two megaplasmids, the 1.35 Mb pSymA and the 1.68 Mb pSymB (Sallet et al. 2013). As a proof of principle, an iPtgxDB created for *S. meliloti* 2011 has allowed the detection of the 14-aa-long leader peptide peTrpL in the proteomic data, for which a function in resistance to multiple antimicrobial compounds could subsequently be established (Melior et al. 2020). However, the identification of additional functional SEPs in *S. meliloti* and related Alphaproteobacteria has been limited by the lack of studies specifically targeting the small proteome and translome.

Here, we developed and then applied a Ribo-seq protocol on *S. meliloti* 2011 to map its translome globally, with a focus on the small proteome (data available at our interactive web-based genome-browser: <http://www.bioinf.uni-freiburg.de/ribobase>). The use of RNase I in our Ribo-seq showed successful trimming of mRNA regions that were not protected by ribosomes, allowing differentiation between translated and untranslated regions. Besides detecting the translation of annotated sORFs (some of which are available in recent updates of the genome annotation), we also uncovered 37 translated novel, non-annotated sORFs located on different replicons. The translation of several annotated and novel sORFs was further validated by MS-based proteomics using iPtgxDBs and/or epitope tagging and Western blot analysis, thereby confirming predictions based on Ribo-seq coverage. Eleven novel SEPs were uniquely identified by MS, showing that using both methods when mapping the small proteome is advantageous. Overall, our combined approach provided a set of 48 novel *S. meliloti* sORFs, many of which are conserved, as a resource to further elucidate their roles in bacterial physiology and symbiosis.

## Methods

### Growth and harvest of *S. meliloti* for Ribo-seq

*S. meliloti* 2011 (Casse et al. 1979) was first cultivated on TY (5 g of Bacto-Tryptone, 3 g of Bacto-yeast extract, and 0.3 g of CaC<sub>2</sub> per liter) agar plates (Beringer 1974). The plate cultures were used to inoculate liquid cultures, which were grown semi-aerobically (routinely, 30 ml of medium in a 50 ml Erlenmeyer flask under constant agitation at 140 rpm) at 30°C in GMS minimal medium (10 g of D-mannitol, 5 g of sodium glutamate, 5 g of K<sub>2</sub>HPO<sub>4</sub>, 0.2 g of MgSO<sub>4</sub> × 7H<sub>2</sub>O, and 0.04 g of CaCl<sub>2</sub> per liter; trace elements: 0.05 mg of FeCl<sub>3</sub> × 6H<sub>2</sub>O, 0.01 mg of H<sub>3</sub>BO<sub>3</sub>, 0.01 mg of ZnSO<sub>4</sub> × 7H<sub>2</sub>O, 0.01 mg of CoCl<sub>2</sub> × 6H<sub>2</sub>O, 0.01 mg of CuSO<sub>4</sub> × 5H<sub>2</sub>O, 1.35 mg of MnCl<sub>2</sub>, and 0.01 mg of Na<sub>2</sub>MoO<sub>4</sub> × 2H<sub>2</sub>O per liter; 10 µg of biotin and 10 mg of thiamine per liter) (Zevenhuizen and van Neerven 1983). As the strain exhibits chromosomally encoded strepto-

mycin resistance, 250 µg/ml streptomycin was added to the media. For Ribo-seq sample preparation, cells corresponding to 40 OD<sub>600</sub> equivalent units were harvested after rapid chilling in an ice bath to halt cell growth and translation. In brief, cultures in the exponential phase (OD<sub>600nm</sub> 0.5) were rapidly placed in a pre-chilled flask in an ice-water bath and incubated with gentle shaking for 3 min. Cells were then immediately pelleted by centrifugation (10 min at 6000 ×g) before snap-freezing in liquid N<sub>2</sub>. Before centrifugation, a culture aliquot was withdrawn for total RNA analysis, mixed with 0.2 vol stop mix (5% buffer-saturated phenol [Roth] in 95% ethanol), and snap-frozen in liquid N<sub>2</sub>. Even though translation elongation inhibitors have been extensively used in both eukaryotic and bacterial Ribo-seq workflows, using such chemicals can introduce bias into Ribo-seq coverage (Gerashchenko and Gladyshev 2014, Mohammad et al. 2019). Therefore, we chose to perform Ribo-seq without these inhibitors because we were able to recover sufficient polysomes using the fast-chilling method (see Fig. 1).

### Preparation of ribosome footprints

Ribo-seq was performed as previously described (Oh et al. 2011, Hadjeras et al. 2023), with some modifications. In brief, cell pellets were resuspended with cold lysis buffer (1 M NH<sub>4</sub>Cl, 150 mM MgCl<sub>2</sub>, 20 mM Tris-HCl, 5 mM CaCl<sub>2</sub>, 0.4% Triton X-100, 150 U DNase I [Fermentas], and 1000 U RNase Inhibitor [MoloX, Berlin] at pH 8.0) and lysed by sonication (constant power 50%, duty cycle 50%, and 3 × 30 s cycles with 30 s cooling on a water-ice bath between each sonication cycle to avoid heating of the sample). The lysate was clarified by centrifugation at 10,000 ×g for 12 min at 4°C. To approximately 15 A<sub>260</sub> of lysate, 200 U of RNase I (Thermo Fisher Scientific) was added. Polysome digestion was performed at 25°C with shaking at 650 rpm for 90 min. A mock-digested control (no enzyme added) was performed in parallel to confirm the presence of polysomes in the lysate. To analyze polysome profiles and recover digested monosomes, we layered 15 A<sub>260</sub> units onto a linear 10%–55% sucrose gradient prepared in 4× gradient buffer (10× gradient buffer: 100 mM MgCl<sub>2</sub>, 200 mM Tris-HCl, 1 M NH<sub>4</sub>Cl, and 20 mM dithiothreitol [DTT] at pH 8.0) in an ultracentrifuge tube (13.2 mL Beckman Coulter SW-41). Gradients were centrifuged in a SW40-Ti rotor at 35 000rpm for 2 h and 30 min at 4°C in a Beckman Coulter Optima XPN-80 ultracentrifuge. Gradients were processed using a gradient station (IP, Biocomp Instruments) fractionation system with continuous absorbance monitoring at 254 nm to resolve ribosomal subunit peaks. The 70S monosome fractions were collected and subjected to RNA extraction to purify the RNA footprints.

RNA was extracted from fractions or cell pellets for total RNA using hot phenol-chloroform-isoamyl alcohol (25:24:1, Roth) or hot phenol (Roth), respectively, as previously described (Sharma et al. 2007, Venturini et al. 2020). Ribosomal RNA (rRNA) was depleted from 5 µg of DNase I-digested total RNA by subtractive hybridization with the Pan-Bacteria riboPOOLS (siTOOLS, Germany) according to the manufacturer's protocol with Dynabeads MyOne Streptavidin T1 beads (Invitrogen). Total RNA was fragmented with an RNA fragmentation reagent (Ambion). Monosome RNA and fragmented total RNA were size selected (26–34 nt) on 15% polyacrylamide/7 M urea gels, as previously described (Ingolia et al. 2012) using RNA oligonucleotides NI-19 and NI-20 as guides. RNA was cleaned and concentrated by isopropanol precipitation with 15 µg of GlycoBlue (Ambion) and dissolved in H<sub>2</sub>O. cDNA libraries were prepared by Vertis Biotechnologie AG (Freising, Germany) using the adapter ligation protocol without

fragmentation. First, an oligonucleotide adapter was ligated to the 3' end of the RNA molecules. First-strand cDNA synthesis was performed using M-MLV reverse transcriptase and the 3' adapter as the primer. The first strand of cDNA was purified, and the 5' Illumina TruSeq sequencing adapter was ligated to the 3' end of the antisense cDNA. The resulting cDNA was PCR-amplified to approximately 10–20 ng/µl using a high-fidelity DNA polymerase. The DNA was purified using the Agencourt AMPure XP kit (Beckman Coulter Genomics) and analyzed by capillary electrophoresis. The primers used for PCR amplification were designed for TruSeq sequencing according to the instructions of Illumina. The following adapter sequences flank the cDNA inserts: TruSeq\_Sense\_primer: (NNNNNNNN = i5 Barcode for multiplexing) 5'-AATGATACGGCGACCACCGAGATCTACAC-NNNNNNNN-ACACTCTTCCCTACA CGACGCTCTCCGATCT-3'; TruSeq\_Antisense\_primer: (NNNNNNNN = i7 Barcode for multiplexing) 5'-CAAGCAGAAGACGGCATACGAGAT-NNNNNNNN-GTGACTGGAGTTCAGACGTGT GCTCTTCCGATCT-3'. cDNA libraries were pooled on an Illumina NextSeq 500 high-output flow cell and sequenced in single-end mode (75 cycles, with 20 million reads per library) at the Core Unit SysMed at the University of Würzburg.

### Ribo-seq data analysis

*S. meliloti* Ribo-seq data were processed and analyzed using the published HRIBO workflow (version 1.6.0) (Gelhausen et al. 2021), which has previously been used for the analysis of bacterial Ribo-seq data (Venturini et al. 2020). In brief, sequencing read files were processed with a snakemake (Köster and Rahmann 2012) workflow, which downloads all required tools from bioconda (Grüning et al. 2018) and automatically determines the necessary processing steps. Adapters were trimmed from the reads with cutadapt (version 2.1) (Martin 2011) and then mapped against the *S. meliloti* 2011 genome with *segemehl* (version 0.3.4) (Otto et al. 2014). Reads corresponding to rRNA and other multiply mapping reads were removed with SAMtools (version 1.9) (Li et al. 2009). Quality control was performed by creating read count statistics for each processing step and RNA class with Subread featureCounts (1.6.3) (Liao et al. 2014). All processing steps were analyzed with FastQC (version 0.11.8) (Wingett and Andrews 2018), and the results were aggregated with MultiQC (version 1.7) (Ewels et al. 2016). Summary statistics are shown in Table S1.

Read coverage files were generated with HRIBO using different full-read mapping approaches (global or centered) and single-nucleotide mapping strategies (5' or 3' end). Read coverage files using two different normalization methods were created (mil and min). For the mil normalization, read counts were normalized by the total number of mapped reads within the sample and scaled by a per-million factor. For the min normalization, the read counts were normalized by the total number of mapped reads within the sample and scaled by the minimum number of mapped reads among all analyzed samples. The coverage files generated using the min normalization and the global mapping (full read) approach were used for genome browser visualization. Metagene analysis of ribosome density at start codons was performed as previously described (Becker et al. 2013).

### Ribo-seq-based ORF prediction, filtering, and manual curation

ORFs were called with an adapted variant of REPARATION (Ndah et al. 2017) using blast instead of usearch (see [https://github.com/RickGelhausen/REPARATION\\_blast](https://github.com/RickGelhausen/REPARATION_blast)) and DeepRibo (Clauwaert et

al. 2019). Generic feature format (GFF) track files with this information, plus potential start and stop codons and ribosome binding site information were created for in-depth manual genome browser inspection. Summary statistics for GenBank annotated and merged novel ORFs detected by REPARATION and DeepRibo were computed in a tabularized form, including, among other values, translation efficiency (TE), RPKM (reads per kilobase of transcript per million mapped reads) normalized read counts, codon counts, and nucleotide and aa sequences (see Table S2). Annotated sORFs were classified as translated if they fulfilled an arbitrary mean TE cut-off of  $\geq 0.5$  and RNA-seq and Ribo-seq RPKM of  $\geq 10$  (cut-offs chosen based on the lowest TE and RPKM values associated with housekeeping genes [i.e. ribosomal protein genes] and the genes detected by proteomics). To identify robust novel sORF candidates, we inspected HRIBO ORF predictions from DeepRibo and REPARATION. As DeepRibo is prone to a high rate of false positives (Gelhausen et al. 2022), we first generated a reasonable set of potential novel sORFs by applying the following expression cut-off filters: mean TE of  $\geq 0.5$  and RNA-seq and Ribo-seq RPKM of  $\geq 10$  (in both replicates) based on the 85 positively labeled translated sORFs (see Fig. 3). In addition, novel translated sORF candidates were required to have a DeepRibo prediction score of  $> -0.5$  that allows for ORF candidate ranking (Clauwaert et al. 2019). The filtered sORFs were then subjected to manual curation as described (Gelhausen et al. 2022). This manual inspection of paired Ribo-seq and RNA-seq read coverage files in a genome browser allowed for asserting the translation status for the filtered novel predicted sORFs. Briefly, the Ribo-seq and RNA-seq read coverage files were loaded in the Integrated Genome Browser (IGB) along the sequence of the reference genome and the GenBank 2014 annotation, which contains annotated 5'- and 3'-UTRs. The Ribo-seq and RNA-seq read coverage files (normalized to the lowest number of reads between the two) were visually inspected with similar scales. To assess translation of the predicted novel sORFs, we used the following criteria: (i) Ribo-seq read coverage within ORF boundaries with the detection of ribosome footprints in the UTRs (15–16 nucleotides) near the start and stop codons resulting from initiating and terminating ribosomes; (ii) exclusion of Ribo-seq read coverage from the residual 5'- and 3'-UTRs; (iii) the shape of the Ribo-seq read coverage; here the evenness of the read coverage was considered and predicted sORFs with uneven read coverage (exhibiting peaks with plateau, which resulted from either RNA structures or cDNA library preparation artifacts) were not taken into account; (iv) the Ribo-seq read signal was generally required to be comparable to or higher than the transcriptome signal from the RNA-seq library. We created an interactive web-based genome browser using JBrowse (<http://www.bioinf.uni-freiburg.de/ribobase>) (Buels et al. 2016), where the coverage files for the Ribo-seq replicates, the annotation, and the predicted sORF can be visualized.

### Sample preparation for MS

For MS analysis, cells of 1.5 l of an *S. meliloti* culture ( $OD_{600nm}$  0.5) were harvested by centrifugation at 6000 rpm and 4°C. The cell pellet was resuspended in 30 ml of buffer containing 20 mM Tris, 150 mM KCl, 1 mM  $MgCl_2$ , and 1 mM DTT at pH 7.5. After lysis by sonication and centrifugation at 13,000rpm for 30 min at 4°C, the cleared lysates were frozen in liquid nitrogen and stored at  $-80^\circ C$ . To generate a highly comprehensive small protein dataset, we used three complementary approaches for sample preparation: (i) tryptic in-solution digest of all proteins in the sample, (ii) solid-phase enrichment (SPE) of small proteins without any sub-

sequent digestion, and (iii) SPE of small proteins with subsequent digestion using Lys-C. Sample preparation was performed as previously described (Bartel et al. 2020) with some modifications. In brief, samples for tryptic in-solution digests were reduced and alkylated before trypsin was added in an enzyme-to-protein ratio of 1:100, and samples were incubated at 37°C for 14 h. The digest was stopped by acidifying the mixture with HCl. For SPE, samples were loaded on an equilibrated column packed with an 8.5 nm pore size, modified styrene-divinylbenzene resin (8B-S100-AAK, Phenomenex), which was then washed to remove large proteins. The enriched small protein fraction was eluted with 70% (v/v) acetonitrile and evaporated to dryness in a vacuum centrifuge. The SPE samples were either directly used for MS or in-solution digested as described above but with Lys-C instead of trypsin.

### Generation of standard and custom iPtgxDBs to identify novel SEPs

iPtgxDBs were generated based on the *S. meliloti* 2011 ASM34606v1 reference genome sequence as described (Omasits et al. 2017). Annotations from several reference genome centers and/or releases (GenBank 2014, RefSeq2017, Genoscope), two *ab initio* gene predictions (Prodigal, Hyatt et al. 2010; ChemGenome, Mishra et al. 2019), and *in silico* ORF predictions were hierarchically integrated for a trypsin-specific iPtgxDB as previously detailed (Melior et al. 2020), (<https://iptgxdb.expasy.org/database/annotations/s-meliloti-tryptic>; see Table S3.1). To capture data from all three experimental approaches, two more iPtgxDBs were created in a similar fashion using command-line utilities. For the LysC-specific iPtgxDB, the regular expression '(K)(.)' was used, allowing cleavage after every lysine. The iPtgxDB for the experiments without protease digestion was generated with a regular expression that did not allow any cleavages. In addition, three 20-fold smaller custom iPtgxDBs were created to improve search statistics/predictive potential. For these, instead of adding the ChemGenome and *in silico* predictions, 266 selected Ribo-seq translation products identified from the sORF prediction tools DeepRibo (Clauwaert et al. 2019) and Reparation (Ndah et al. 2017), as well as manual analysis, were converted to GFF format using a custom Python script and integrated along with the RefSeq, GenBank, Genoscope (Vallenet et al. 2013), and Prodigal predictions to create the respective iPtgxDBs (Tables S3.3 and S3.4). All six iPtgxDBs (downloadable from <https://iptgxdb.expasy.org>) also contained sequences of common laboratory contaminants (116 from CrapOme and 256 from the Functional Genomics Center Zurich). All peptides implying potentially novel proteins were subjected to a PeptideClassifier analysis (Qeli and Ahrens 2010) extended for proteogenomics in prokaryotes (Omasits et al. 2017). This procedure ensures that i) only unambiguous peptides were considered (class 1a or 3a; see below) or ii) annotation cluster-specific cases can be distinguished: Class 2a peptides imply a subset of all possible proteoforms (e.g. like an extension, reduction), class 2b peptides imply all isoforms, which means that the gene encoding the proteoforms, but not a specific proteoform, was identified. Class 3a peptides unambiguously imply a protein sequence, that however can be encoded by several identical gene copies. For more information about peptide evidence classes and annotation clusters, please see the iPtgxDB web server documentation ([https://iptgxdb.expasy.org/creating\\_iptgx\\_dbs/](https://iptgxdb.expasy.org/creating_iptgx_dbs/)).

### MS analysis

Samples were loaded on an EASY-nLC 1200 (Thermo-Fisher Scientific) equipped with an in-house-built 20 cm reversed-phase

column packed with 3  $\mu\text{m}$  Repronil-Pur 120 C18-AQ (Dr. Maisch) and an integrated emitter tip. Peptides were eluted by a 156 min non-linear gradient of solvent B (0.1% v/v acetic acid in acetonitrile) and injected online in an Orbitrap Velos (Thermo-Fisher Scientific). The survey scans were acquired in the Orbitrap (300–1700 Th; 60 000 resolution at 400 m/z;  $1 \times 10^6$  predictive automatic gain control target; activated lock mass correction). After collision-induced dissociation with a normalized collision energy of 35, fragment spectra were recorded in the LTQ (mass range dependent on precursor m/z;  $3 \times 10^4$  predictive automatic gain control) for the 20 most abundant ions. Fragmented ions were dynamically excluded from fragmentation for 30 s.

DB searches were performed with Sorcerer-SEQUEST 4 (Sage-N Research, Milpitas, USA), allowing two missed cleavages for samples derived from tryptic in solution digest or LysC digested SPE samples and with non-specified enzymes for SPE samples without proteolytic digest. No fixed modifications were considered, and oxidation of methionine was considered a variable modification. The mass tolerance for precursor ions was set to 10 ppm, and the mass tolerance for fragment ions was set to 1.0 Da. Validation of MS/MS-based peptide and protein identification was performed with Scaffold V4.8.7 (Proteome Software, Portland, USA), and peptide identifications were accepted if they exhibited at least deltaCn scores of  $> 0.1$  and XCorr scores of  $> 2.2$ , 3.3, and 3.75 for doubly, triply, and all high-charged peptides, respectively. Identifications for proteins of  $> 15$  kDa were only accepted if at least two unique peptides were identified. Proteins that contained ambiguous, non-unique peptides and could not be differentiated based on MS/MS analysis alone were grouped to satisfy the principles of parsimony (Sorcerer-SEQUEST). Identifications for annotated proteins of  $< 15$  kDa were accepted if at least one unique peptide was identified with at least two peptide spectrum matches (PSMs). To identify novel proteins, we required additional PSM evidence from predictions as described before (Varadarajan et al. 2020a,b), that is, 3 PSMs for *ab initio* predictions and 4 PSMs from *in silico* predictions. Here, we also allowed *in silico* candidates with 3 PSMs if they were observed in each of the three replicates. Similar to the RefSeq annotated proteins, novel proteins greater than 15 kDa (~150 aa) required two unique peptides (however, these were not the focus of this study). The application of these filter criteria kept the protein false discovery rate (FDR) below 1%. To facilitate overview and comparison, we integrated MS-identified proteins, Ribo-Seq, and Western blot analysis data in a 'master table' (Table S4).

## Cloning procedures

The oligonucleotides (Microsynth) used for cloning are listed in Table S5. Routinely, FastDigest Restriction Endonucleases and Phusion polymerase (Thermo Fisher Scientific) were used. PCR products were first ligated into pJet1.2/blunt (CloneJet PCR Cloning Kit, Thermo Fisher Scientific) and transformed into *E. coli* DH5- $\alpha$ . Subsequently, inserts were subcloned in conjugative plasmids originating from pSRKGm (Khan et al. 2008). Insert sequences were analyzed by Sanger sequencing with plasmid-specific primers (Microsynth Seqlab). *E. coli* S17-1 was used to transfer the plasmids to *S. meliloti* 2011 by diparental conjugation (Simon et al. 1983).

Plasmid pSW2 was used to clone the candidate sORFs. It was constructed using pRS1, a derivative of pSRKGm, in which the *E. coli* lac module was exchanged for a multiple cleavage site-containing cloning site for the restriction endonucleases NheI, HindIII, XbaI, SpeI, BamHI, PstI, and EcoRI. First, a transcription terminator  $T_{rm}$  from *Bradyrhizobium japonicum* USDA 110 was

cloned into the EcoRI restriction site of pRS1. For this, the terminator containing sequence was amplified with the forward primer Bj-Trrn-Fw-2019 and the reverse primer Bj-Trrn-Rv-2019 using plasmid pJH-O1 as a template (Čuklina et al. 2016). In the PCR product, an EcoRI restriction site was present downstream of the forward primer sequence. This restriction site and that in the reverse primer were used for the transcription terminator cloning. A clone with the desired orientation was selected, and the plasmid was named pRS1-Trrn (Fig. S1). Double-stranded DNA encoding a sequential peptide affinity (SPA) tag, which is composed of the calmodulin-binding peptide and three modified FLAG sequences separated by a TEV protease cleavage site (Zeghouf et al. 2004), was then cloned between the BamHI and EcoRI cleavage sites of pRS1-Trrn. The SPA-tag encoding sequence was designed without an ATG codon, without rare codons, and with Gly-Gly-Gly-Ser linker codons at the 5' end and adapted to the high GC content of *S. meliloti*. It was generated synthetically by Eurofins and provided on plasmid pEX-A128, which was used as a template for PCR amplification with primers SmSPA-Ct-BamFW and SmSPA-Ct-EcoRV. The resulting plasmid pSW1 can be used to clone an sORF in frame with the SPA-encoding sequence and under the control of its own promoter. Here, pSW1 was used to clone the promoter  $P_{sinI}$  between the NheI and XbaI restriction sites. The promoter sequence (McIntosh et al. 2008) was amplified using primers NheI-PsinI-FW and XbaI-PsinI-RV and *S. meliloti* 2011 genomic DNA as a template. The resulting pSW2 plasmid was used to clone candidate sORFs, each with a 15-nt upstream region potentially harboring a Shine-Dalgarno sequence between the XbaI and BamHI restriction sites (Fig. S1). In total, 20 sORF::SPA fusions were cloned and tested by Western blot analysis. The corresponding plasmids were designated from pSW2-SEP1 to pSW2-SEP20.

## Western blot analysis

Exponentially grown *S. meliloti* cells ( $OD_{600\text{nm}}$  0.5; minimal medium) were harvested ( $3500 \times g$  for 10 min at  $4^\circ\text{C}$ ) and resuspended in an SDS-loading buffer. After incubation for 5 min at  $95^\circ\text{C}$ , the crude lysate proteins were separated by Tricine-SDS PAGE (16%) and blotted onto a PVDF membrane (Amersham<sup>TM</sup>Hybond<sup>TM</sup>, 0.2  $\mu\text{m}$  PVDF; GE Healthcare Life Science, Chalfont St Giles, Great Britain) as described (Schägger 2006). For detection, monoclonal ANTI-FLAG M2-Peroxidase (HRP) antibodies (Merck, Darmstadt, Germany) and Lumi-Light Western-Blot-Substrate (Roche, Basel, Schweiz) were used. Signal visualization was performed with a chemiluminescence imager (Fusion SL4, Vilber, Eberhardzell, Germany). For fractionation, the cell pellets were resuspended in TKMDP buffer containing 20 mM Tris-HCl, 150 mM KCl, 1 mM  $\text{MgCl}_2$ , 1 mM DTT, and one protease inhibitor cocktail tablet at pH 7.5 (Sigma Aldrich, St. Louis, USA). Lysates prepared by three passages in a French press at 1000psi were cleared by centrifugation at  $14,000 \times g$  for 30 min at  $4^\circ\text{C}$ . The supernatant was subjected to ultracentrifugation at  $100,000 \times g$  for 1 h at  $4^\circ\text{C}$ . The supernatant (S100 fraction) was then removed, and the P100 pellet was resuspended in the same volume of TKMDP buffer.

## Conservation analysis, protein domain, and operon prediction

The identification of novel small protein homologues was performed using Blastp and tBlastn searches in bacteria on the National Center for Biotechnology Information DB (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>). The protein sequences for novel protein candidates identified by Ribo-seq and/or MS were used as

query sequences. For tBlastn, the following parameters were used: the filter for low complexity regions off, a seed length that initiates an alignment (word size) of 6, 60% coverage of the query sequence with at least 40% identity, an E-value (Expect value) of  $\leq 100$  to capture all potential orthologs, and an E-value below 0.1 for high-confidence hits (Allen et al. 2014). Moreover, novel small protein candidates were further analyzed for secondary structure, for predicted protein domains and lipoprotein signatures, as well as for potential subcellular localization using predictions from the Phyre2 v2.0 (<http://www.sbg.bio.ic.ac.uk/~phyre2/>), LipoP-1.0 (<https://services.healthtech.dtu.dk/service.php?LipoP-1.0>) TMHMM v2.0 (<https://services.healthtech.dtu.dk/service.php?TMHMM-2.0>), and PSORTb v3.0.2 servers (<https://www.psорт.org/psортb/>). For operon prediction, we used the publicly available OperonMapper software (Taboada et al. 2018) by creating a GFF file containing both the 48 novel sORFs as well as all CDS from the RefSeq 2022 annotation.

## Data availability

The MS-based proteomics data were deposited to the ProteomeXchange Consortium at the PRIDE partner repository, with dataset identifier PXD034931. The iPTxDBs can be downloaded from <https://iptgxdb.expasy.org/>. Ribo-seq and RNA-seq data were deposited in GEO, with accession number GSE206492. The Ribo-seq and RNA-seq data of *S. meliloti* 2011 can be viewed with an interactive online JBrowse instance (<http://www.bioinf.uni-freiburg.de/ribobase>).

## Results

### Establishing Ribo-seq in *S. meliloti* to map its translome

To provide a genome-wide map of translated annotated sORFs and to reveal new sORFs in the plant symbiont *S. meliloti*, we first adapted the Ribo-seq protocol (Oh et al. 2011, Hadjeras et al. 2023) to this organism (Fig. 1A). For this purpose, several steps, including cell harvest, lysis, and footprint generation, were optimized (see Methods). *S. meliloti* 2011 cells were grown to the mid-log phase in minimal medium, and samples were rapidly cooled and harvested to avoid polysome run-off. Polysome profile analysis after lysate fractionation on a sucrose gradient showed successfully captured translating ribosomes (Fig. 1B, black profile). The mRNA should be ribonucleolytically digested outside ribosomes to produce ribosome footprints. Since the broad-range ribonuclease RNase I, which is often used for eukaryotic Ribo-seq analysis, is inactive on polysomes from enteric bacteria (Datta and Burma 1972, Bartholomäus et al. 2016), most prokaryotic Ribo-seq protocols mainly use micrococcal nuclease (MNase) instead. Since MNase preferentially cleaves at pyrimidines, it typically introduces periodicity artifacts, and generates footprints that are more heterogeneous in length than those from RNase I (Ingolia 2016, Vazquez-Laslop et al. 2022). Therefore, we used RNase I to convert *S. meliloti* polysomes into monosomes (Fig. 1B) and to generate ribosome footprints (Fig. 1C and E). By comparing Ribo-seq read coverage data and expression signals from a paired RNA-seq library generated from fragmented total RNA, features, such as coding potential, ORF boundaries, and 5'- and 3'-UTRs, can be defined (Fig. 1C and E).

Inspection of Ribo-seq coverage for translated ORFs and known non-coding transcripts further demonstrated the successful setup of Ribo-seq in *S. meliloti*. For example, the protein-coding genes *rpsO* and *icd* showed higher cDNA read coverage in the Ribo-

seq library compared with the paired RNA-seq library (Fig. 1C), whereas the RNase P RNA gene *mpB* showed high cDNA read coverage only in the RNA-seq library (Fig. 1D). Furthermore, the cDNA read coverages of the 5'- and 3'-UTRs of *rpsO* and *icd* were higher in the RNA-seq library than in the Ribo-seq library (Fig. 1C), showing successful digestion of non-translated or unprotected mRNA regions by RNase I. Similarly, the protein-coding polycistronic *fixN1OQP* mRNA showed high read coverage in the Ribo-seq library along its four ORFs. In contrast, the 5'-leader and 3'-trailer mainly showed coverage in the RNA-seq library, suggesting that they were digested by RNase I (Fig. 1E).

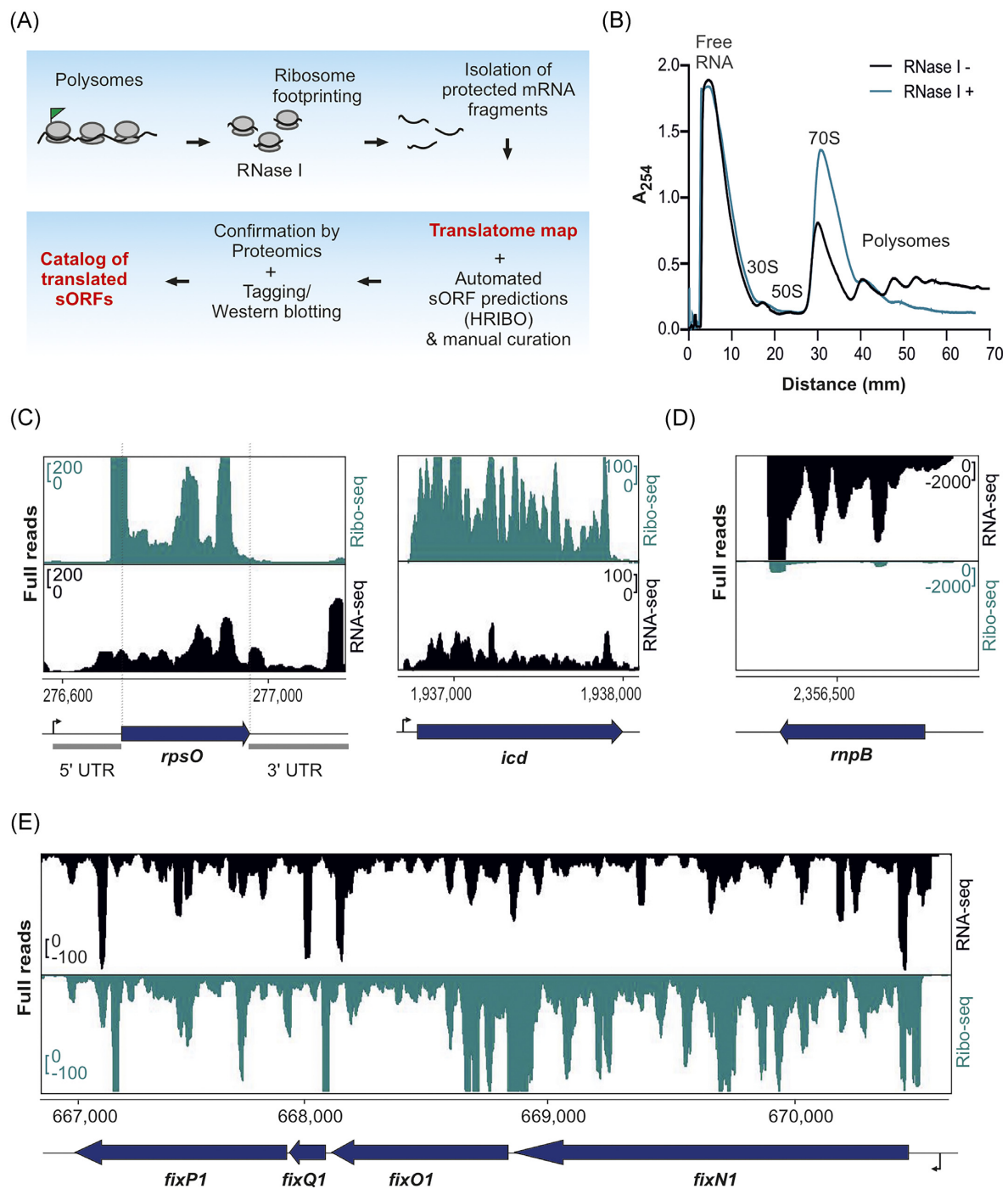
The high ribosome density in the Ribo-seq library, which covers the 14 and 12-nt-long intergenic regions between *fixN1*–*fixO1* and *fixO1*–*fixQ1*, probably represents the footprints of ribosomes that terminate the translation of the upstream ORF and initiate the translation of the downstream ORF. Such events are slower than elongation at most codons in an ORF (Oh et al. 2011). The latter example indicates the translation of the sORF *fixQ1*, which encodes a 50 aa protein (Fig. 1E).

Metagenome analysis of ribosome occupancy near all annotated start codons (i.e. ATG, GTG, and TTG) showed an enriched ribosome density at the  $-16$  nt upstream (mapping of the 5' ends of the footprints) and at  $+16$  nt downstream (mapping of the 3' ends of the footprints) (Fig. S2A and S2B; note:  $+1$  is the first nucleotide of the start codon), in line with the expected position of initiating ribosomes waiting to engage in elongation. This feature is a characteristic of translated bacterial ORFs identified by Ribo-seq (Oh et al. 2011, Mohammad et al. 2019). In contrast to MNase-generated Ribo-seq libraries in *E. coli* (Mohammad et al. 2019), no differences in the assignment of ribosome position using the 5' end or 3' end mapping approaches were observed (Fig. S2A and S2B). In the Ribo-seq libraries, we consistently recovered footprints between 27 and 33 nt (mean at 30 nt), with enrichment of ribosome density strongest at the start codon for the 32 nt footprints (Fig. S2C and D).

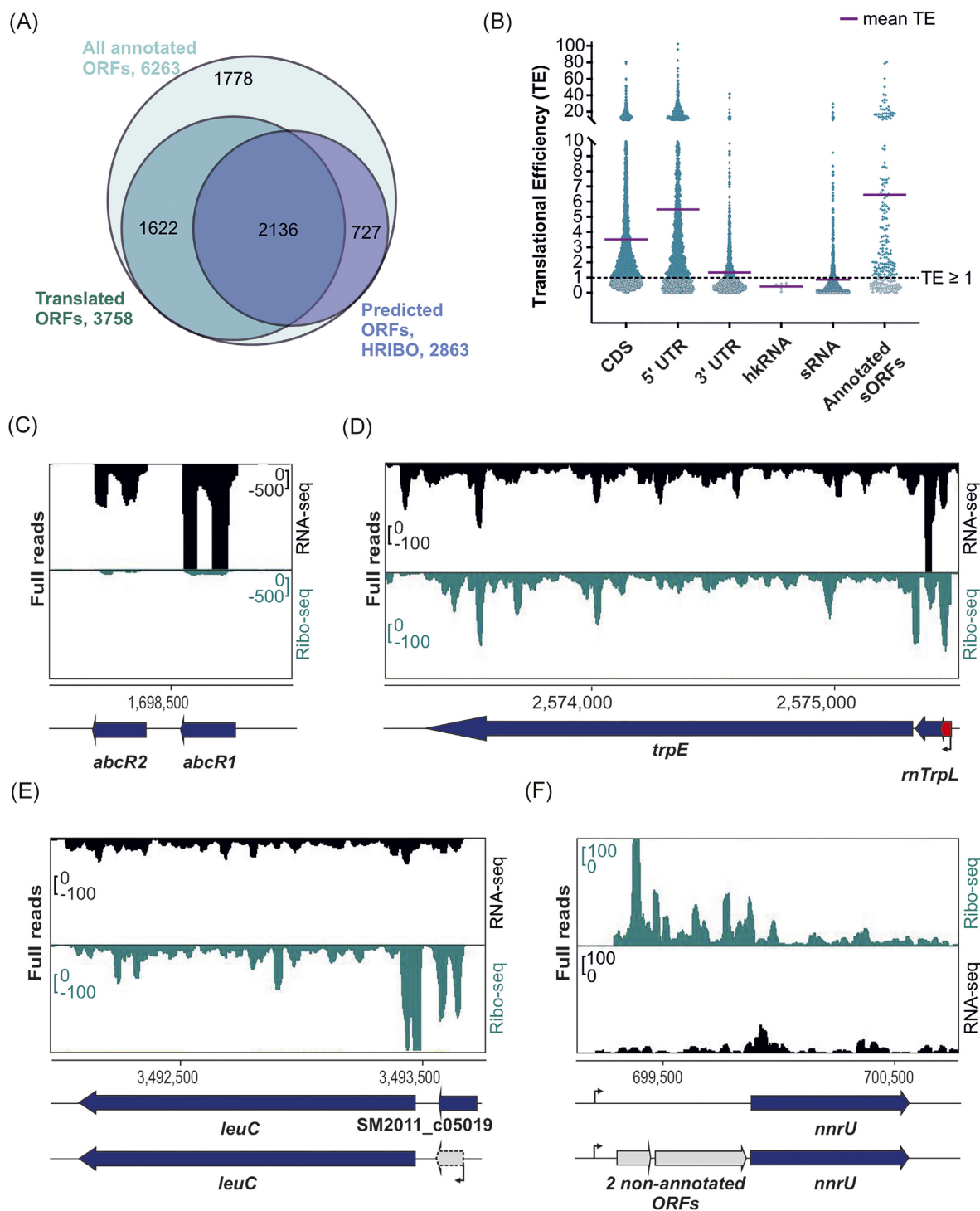
### Ribo-seq captures the translome of *S. meliloti* and reveals features at the single gene level

By comparing the signals of the Ribo-seq and RNA-seq libraries, the TE (ratio Ribo-seq/total RNA coverage) can be estimated at a given locus. This method allowed us to derive a genome-wide estimate of the translome in minimal medium, where 3758 of the 6263 annotated coding sequences (CDS) (60%; GenBank 2014 annotation) had a Ribo-seq signal above the arbitrarily chosen TE cut-off of  $\geq 0.5$  and RNA-seq and Ribo-seq RPKM of  $\geq 10$  (see Methods, Fig. 2A, Table S6). In contrast, the ORF prediction tools implemented in HRIBO (Gelhausen et al. 2021, 2022) detected translation for 2136 of the 3758 ORFs (57%), suggesting an average performance in predicting long translated ORFs in *S. meliloti* (Fig. 2A, Table S6).

Inspection of the TE for different annotated gene classes and untranslated mRNA regions (all CDS, 5'- and 3'-UTRs, non-coding RNAs, and sORFs) revealed that annotated ORFs exhibited a higher mean TE (TE  $\geq 1$ ) compared with non-coding genes, such as housekeeping RNA genes (hkrRNA, e.g. tmRNA, 6S, ffs, *mpB* and *incA1/2* RNA, mean TE  $< 1$ ) (Fig. 2B), again corroborating the ability of our Ribo-seq data to differentiate between coding and non-coding genes. The 5'-UTR regions of translated mRNAs generally had a mean TE of  $\geq 1$ , which possibly resulted from protection from RNase I trimming of the  $-16$  nt region upstream of the start codon by the initiating ribosomes (Fig. S2). This feature was particularly prominent in the leader regions of mRNAs with short

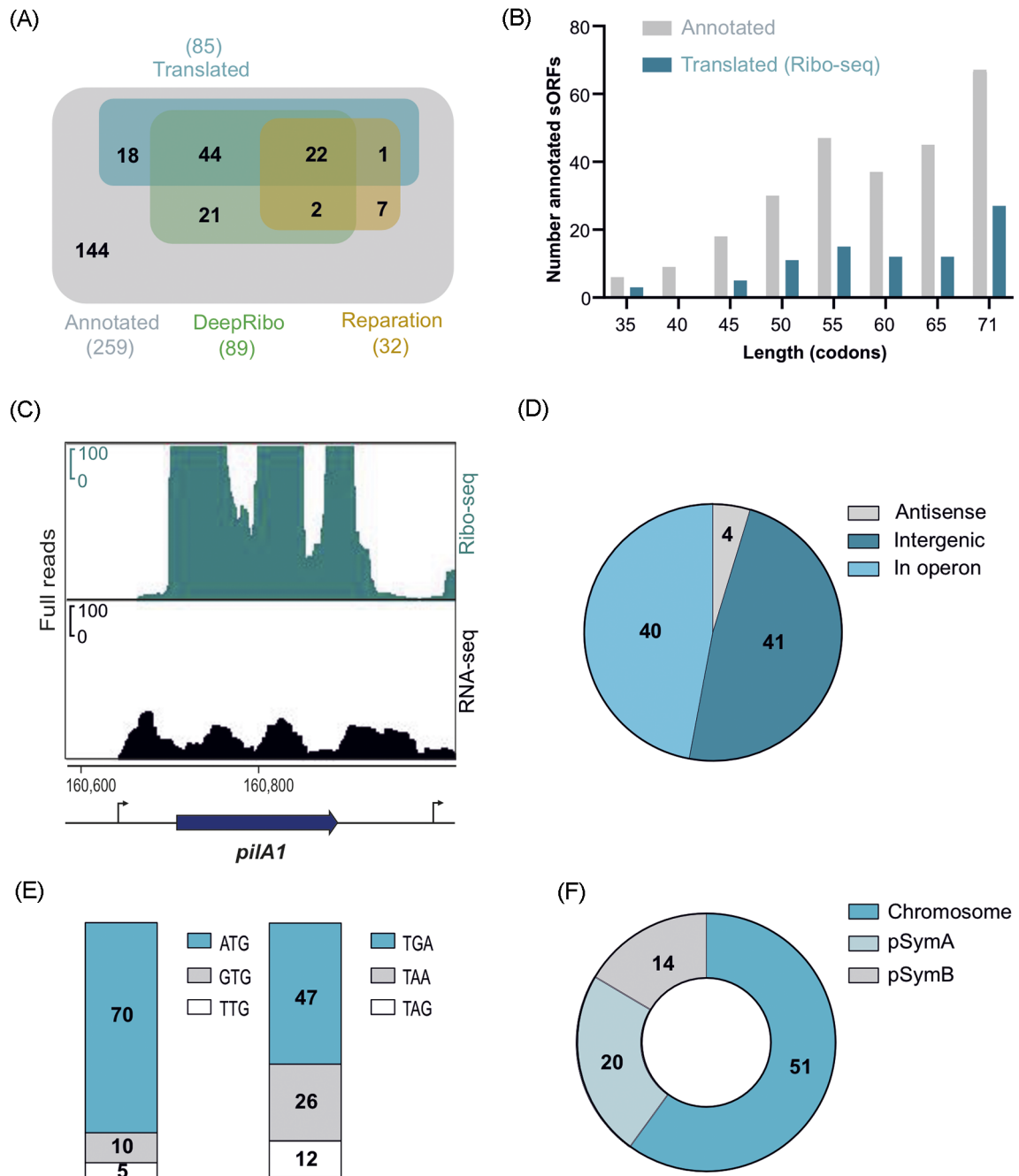


**Figure 1.** Establishment of ribosome profiling (Ribo-seq) for *Sinorhizobium meliloti*. **(A)** Schematic Ribo-seq workflow to map the *S. meliloti* 2011 translome. Translating ribosomes (indicated by the polysome fraction) were first captured on the mRNAs. Unprotected mRNA regions were digested by RNase I, converting polysomes to monosomes. Approximately 30-nt-long footprints protected by and co-purified with 70S ribosomes were then subjected to cDNA library preparation and deep sequencing to identify the translome under the used conditions. The small proteome was identified using HRIBO automated predictions and manual curation. Mass spectrometry and Western blot analysis of recombinant, tagged small open reading frame (sORF)-encoded proteins were used to validate the translated sORFs. **(B)** Sucrose gradient fractionation of the lysates. Cells were harvested at the exponential growth phase by a fast-chilling method to avoid polysome run-off. RNase I digestion led to enrichment of monosomes (70S peak in the green profile) in contrast to the untreated sample (Mock, black profile). Absorbance at 254 nm was measured. **(C)** Integrated genome browser screenshots depicting reads from Ribo-seq and RNA-seq libraries for two annotated ORFs: *rpsO* encoding ribosomal protein S15 and *icd* encoding isocitrate dehydrogenase. They show read coverage enrichment in the Ribo-seq library along their coding parts in contrast to the RNA-seq library but not in the ribosome-non-protected regions (UTRs). The UTRs of *rpsO* are marked. **(D)** Read coverage for *mpb* corresponding to the housekeeping RNase P RNA. Reads are mostly restricted to the RNA-seq library, suggesting that this RNA is not translated. **(E)** The *fixN1OQP* operon shows read coverage in both the RNA-seq library and Ribo-seq library, the latter indicating that this operon contains translated genes. Genomic locations and coding regions are indicated below the image. Bent arrow indicates the transcription start site based on (Sallet et al. 2013).



**Figure 2.** Ribosome profiling (Ribo-seq) captures the translome of *Sinorhizobium meliloti* 2011 and reveals some features at the single-gene level. **(A)** Comparison of all annotated open reading frames (ORFs), annotated translated ORFs detected by Ribo-seq, and ORFs predicted to be translated by tools included in the HRIBO pipeline. To detect translation, we used the following parameters on the Ribo-seq data: TE of  $\geq 0.5$  and RNA-seq and Ribo-seq RPKM of  $\geq 10$ . The numbers of ORFs per category are shown and represented by area size. Diagrams were prepared with BioVenn ([www.biovenn.nl](http://www.biovenn.nl)). **(B)** Scatter plot showing global TEs (TE = Ribo-seq/RNA-seq) computed from *S. meliloti* Ribo-seq replicates for all annotated coding sequences (CDS), annotated 5'- and 3'-UTRs, annotated housekeeping RNAs (hkRNA), annotated small RNAs (sRNAs) with (putative) regulatory functions, and annotated sORFs encoding proteins of  $\leq 70$  amino acids (aa). The purple lines indicate the mean TE for each transcript class. **(C)** Analysis of the two well-characterized sRNAs AbcR1 and AbcR2 by Ribo-seq. These two sRNAs show read coverage mostly in the RNA-seq library. **(D)** Ribo-seq reveals the active translation of the *trpE* leader peptide peTrpL (14 aa, encoded by the leaderless sORF *trpL* in the 5'-UTR (red arrow) and/or by the attenuator sRNA *mTrpL*). In addition, the coverage of the Ribo-seq library shows that the biosynthetic gene *trpE* is translated in minimal medium, as expected. **(E)** Re-annotation of sORF SM2011\_c05019 (50 aa). The GenBank 2014 annotation does not fit the RNA-seq and Ribo-seq read coverages. HRIBO predicts a shorter leaderless sORF (38 aa) that corresponds to the read coverage in both libraries. **(F)** Two ORFs missing from the GenBank 2014 annotation are revealed by the *nnrU* gene related to denitrification. Genomic locations and coding regions are indicated below the image. Bent arrows indicate transcription start sites based on (Sallet et al. 2013).





**Figure 3.** Ribo-seq reveals translated annotated small open reading frames (sORFs) in *Sinorhizobium meliloti* 2011. **(A)** Venn diagrams showing the overlap between all annotated sORFs (259 sORFs, GenBank 2014), the sORFs detected as translated by Ribo-seq (benchmark set, TE of  $\geq 0.5$ , RNA-seq and Ribo-seq RPKM of  $\geq 10$ , and extensive manual curation), and sORFs predicted by the automated ORF prediction tools Repairation or DeepRibo. **(B)** Histogram showing the length distribution of the 85 annotated sORFs identified as translated by Ribo-seq in comparison with the 259 annotated sORFs. **(C)** Integrated genome browser screenshot depicting reads from the Ribo-seq and RNA-seq libraries for the annotated sORF *pilA1* (60 amino acids, encoding a pilin subunit). The genomic position and the coding region are indicated below the image. Bent arrows indicate transcription start sites based on (Sallet et al. 2013). **(D)** Genomic context for the translated annotated sORFs relative to the annotated neighboring genes. **(E)** Start (left) and stop (right) codon usage of the translated annotated sORFs. **(F)** Replicon distribution of the translated annotated sORFs.

5'-UTRs, indicating that they are partially protected from digestion by initiating ribosomes (Fig. S3A). In addition, some 5'-UTRs might contain translated upstream sORFs, such as *trpL* upstream of *trpE* (marked in red in Fig. 2D) (Melior et al. 2020). Although less pronounced than at the start codon, the translation-terminating ribosome also protects a certain 3'-UTR region from RNase digestion (Oh et al. 2011), explaining the slightly higher mean TE of 3'-UTRs (Fig. 2B). Furthermore, a few of the 3'-UTRs might also con-

tain translated downstream sORFs (Fig. S3B; Dodbele and Wilusz 2020, Wu et al. 2020), which may explain the slightly higher mean TE of 3'-UTRs.

Most of the annotated sRNAs had a mean TE of  $< 1$ , indicating that they are in fact non-coding, such as the sRNAs AbcR1 (TE = 0.2) and AbcR2 (TE = 0.09) (Fig. 2C) (Torres-Quesada et al. 2013). However, some annotated sRNAs had a mean TE of  $\geq 1$ , suggesting that they may be small mRNAs or dual-function sRNAs (Fig. S3C).

For example, Fig. 2D shows the recently described dual-function sRNA mTrpL (TE = 1.16), which corresponds to the tryptophan attenuator and contains the *trpL* sORF encoding the functional 14 aa leader peptide peTrpL (Melior et al. 2019, Melior et al. 2021). Since mTrpL is a small, leaderless mRNA starting with the AUG of *trpL*, Fig. 2D also exemplifies how our Ribo-seq analysis can capture leaderless translated ORFs. Furthermore, as expected, we detected translation of the biosynthetic genes *trpE* and *leuC* under growth in minimal medium lacking tryptophan and leucine (Fig. 2D and E).

Finally, we used our Ribo-seq data to curate the annotation of *S. meliloti*. For example, Ribo-seq, RNA-seq data, and our computational ORF predictions based on Ribo-seq all indicated that the start of the sORF SM2011\_c05019 (50 aa) is likely located downstream of the one in the GenBank 2014 annotation, implying a shorter sORF of 38 aa (Fig. 2E; this gene is missing in the latest RefSeq 2022 annotation). Additional sORFs whose annotation should be adjusted are reported in Table S4. Moreover, our data revealed additional ORFs that should be added to the genome annotation. For example, the RNA-seq and Ribo-seq read coverages indicate expression (transcription and translation) upstream of the *nnrU* gene. However, no gene was predicted in this region of the GenBank 2014 annotation. HRIBO's prediction tools indicated the potential for two non-annotated ORFs encoding 51 and 132 aa proteins upstream of the *nnrU* gene (Fig. 2F). The 51 aa ORF is annotated in the related *Sinorhizobium medicae* and *Ensifer adhaerens*, and in the latter, a homologous 142 aa ORF is annotated between the 51 aa sORF and *nnrU*. Notably, while both ORFs were contained in the *S. meliloti* RefSeq 2017 annotation, the 132 aa ORF was removed again from the latest version (June 2022). This observation underlines the need for and value of integrative approaches that can capture and consolidate reference genome annotations from different annotation centers and even from different releases, which can differ substantially. The iPTgxDB approach (Omasits et al. 2017) represents one strategy to readily capture and visualize such differences, as we show here and for a number of additional cases below.

### Ribo-seq reveals translated annotated small proteins in *S. meliloti*

Among the 6263 annotated CDS in the *S. meliloti* 2011 genome (the annotation from GenBank 2014 has been used in the laboratory as a reference point for several years), 259 (roughly 4%) correspond to SEPs, with sizes ranging between 30 (the smallest annotated SEP) and 70 aa (Table S6). To benchmark our Ribo-seq data for its capacity for global identification of translated sORFs, we analyzed the Ribo-seq read coverage of these 259 annotated sORFs. By applying the TE of  $\geq 0.5$  and RNA-seq and Ribo-seq RPKM of  $\geq 10$  cut-off criteria, 131 of them were suggested to be translated (Table S6). However, we further included an extensive manual inspection (see Methods) of the Ribo-seq read coverage on top of these cut-offs to derive a high-confidence dataset of 85 (33%) translated sORFs (Fig. 3A, Table S6).

We then used this set of manually curated, translated sORFs as a benchmark sORF data set to evaluate the performance of two machine learning-based, automated, Ribo-seq-based ORF prediction tools included in our HRIBO pipeline (Gelhausen et al. 2021, 2022), REPARATION (Ndah et al. 2017), and DeepRibo (Clauwaert et al. 2019). REPARATION predicted the translation for 23 of the 85 benchmark sORFs (26%; Fig. 3A), even missing some highly translated sORFs, such as those encoding ribosomal proteins (SM2011\_c04434 encoding 50S ribosomal protein

L34, mean TE = 5.47) and proteins with housekeeping functions (SM2011\_c04884 encoding an anti-sigma factor, mean TE = 2.02, and SM2011\_c03850 encoding the heme exporter D, a cytochrome C-type biogenesis protein, mean TE = 0.88). In contrast, DeepRibo predicted translation for 66 of the 85 benchmark sORFs (78%; Fig. 3A).

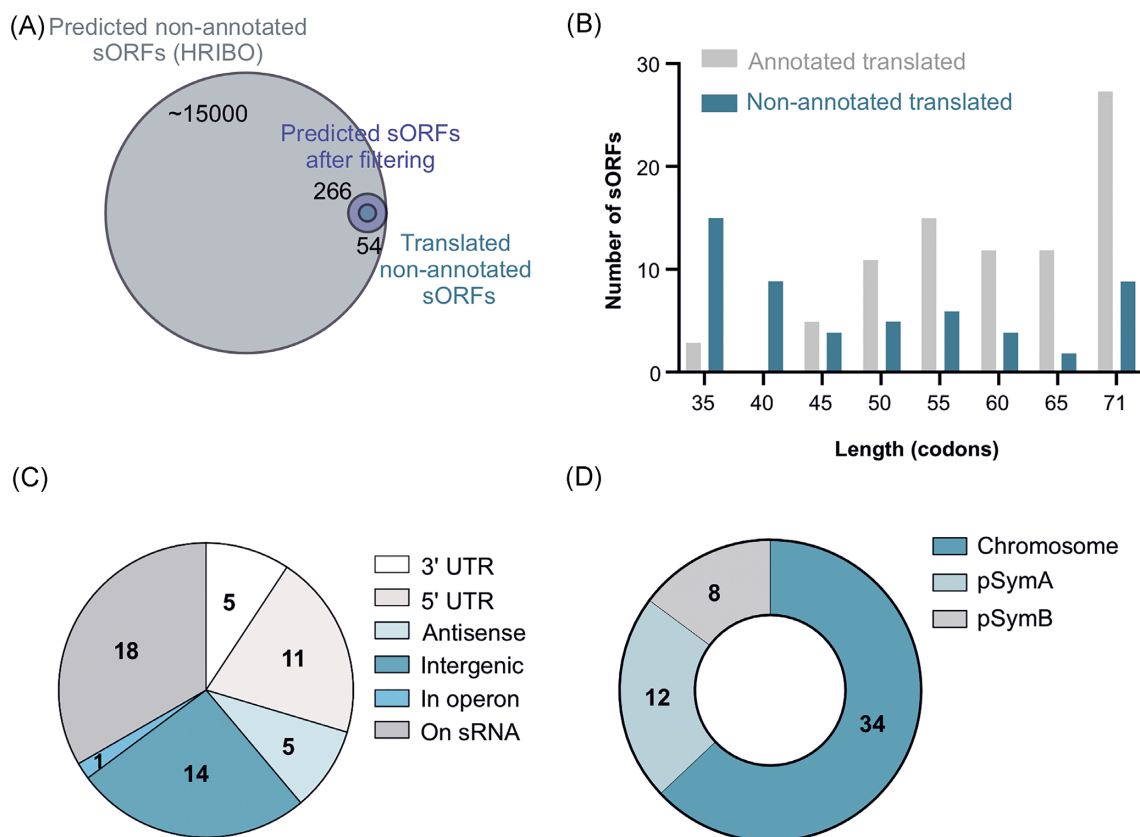
The majority of the 259 annotated (76%) and the subset of 85 translated sORFs (78%) encode SEPs of  $\geq 50$  aa (Fig. 3B), in line with the expected poor annotation of very short ORFs.

Figure 3C shows read coverage from the Ribo-seq and RNA-seq libraries for the sORF encoding a 60 aa pilin subunit (TE = 23.3), which illustrates the successful RNase I digestion of parts of the 5'- and 3'-UTR regions not covered by ribosomes, thus allowing us to define sORF borders. In terms of type of genomic location, most of the translated annotated sORFs are located in intergenic regions and operons, and only a few were found in antisense transcripts (Fig. 3D). The vast majority of the translated annotated sORFs were found to start with ATG, followed by GTG and TTG. The stop codon preference, although less pronounced, was TGA > TAA > TAG (Fig. 3E). Finally, 60% of the 85 translated annotated sORFs were located on the chromosome, 23.5% on the megaplasmid pSymA, and 16.5% on the megaplasmid pSymB (Fig. 3F).

### Ribo-seq further expands the small proteome of *S. meliloti*

We then aimed to exploit the sensitivity of Ribo-seq to identify potential novel *S. meliloti* 2011 sORFs missing from the GenBank 2014 annotation and thereby provide a more complete catalog of its small proteome. The two machine learning-based, automated, Ribo-seq-based ORF prediction tools integrated into the HRIBO pipeline produced a large number of predictions (approximately 15,000) for potential non-annotated sORFs (Fig. 4A), as previously shown in other bacterial species (Gelhausen et al. 2022). Given that these ORF prediction tools neither consider RNA-seq data nor TE but only utilize ribosome occupancy, we decided to filter the predictions for those with RNA-seq and Ribo-seq RPKM values of  $\geq 10$  and mean TE of  $\geq 0.5$ . In addition, we applied a stringent cut-off for the DeepRibo score (see Methods) that allowed an ORF candidate ranking, which led to 266 candidates of translated non-annotated sORFs. Manual curation of all candidates based on their Ribo-seq coverage left us with a list of 54 non-annotated sORFs, which we proposed with high confidence to be translated during growth of *S. meliloti* in minimal medium (Fig. 4A; Table S7). Overall, the 54 non-annotated sORFs were shorter than the annotated ones: 33 of them (61%) correspond to SEPs with lengths between 10 and 49 aa, and nine of them (17%) represent SEPs shorter than 30 aa (the shortest annotated ORF in the *S. meliloti* annotation). A comparison to the length distribution of the 85 annotated and 54 non-annotated translated sORFs (Fig. 4B) illustrates the potential of Ribo-seq to detect very short translated sORFs.

The 54 non-annotated sORFs are encoded in diverse genomic contexts (Fig. 4C): 33% were located on annotated sRNAs, suggesting that these sRNAs are small mRNAs or dual-function sRNAs, 26% were in the intergenic regions, thus defining small mRNAs, and 20% were in the 5'-UTRs and may correspond to regulatory upstream ORFs (Evguenieva-Hackenberg 2022). Only a few were located in 3'-UTRs, on antisense transcripts and inside an operon (Fig. 4C). Moreover, the majority of the 54 sORFs (63%) were located on the chromosome, 22% on pSymA, and 15% on pSymB (Fig. 4D), a distribution comparable to that of the annotated sORFs (Fig. 3F). Similar to the annotated sORFs, ATG was also the preferred start codon among the 54 non-annotated translated



**Figure 4.** Ribo-seq uncovers a repertoire of small open reading frames (sORFs) missing from the *Sinorhizobium meliloti* 2011 genome annotation. **(A)** sORF predictions from HRIBO included a high number of potential non-annotated sORFs (approximately 15,000). These sORFs were first filtered (TE of  $\geq 0.5$ , RNA-seq and Ribo-seq RPKM of  $\geq 10$ , DeepRibo score of  $> -0.5$ ) to generate a set of 266 translated sORF candidates that were additionally manually curated by inspection of the Ribo-seq read coverage in a genome browser. Overall, 54 high-confidence non-annotated sORFs displayed translation during growth in minimal medium. A Venn diagram shows the respective number of proteins from each category (scaled with area size). Diagrams were prepared with BioVenn ([www.biovenn.nl](http://www.biovenn.nl)). **(B)** Histogram showing the length distribution of the 54 non-annotated versus the 85 annotated sORFs identified as translated by Ribo-seq. **(C)** Genomic context of the translated non-annotated sORFs. **(D)** Replicon distribution of the translated non-annotated sORFs.

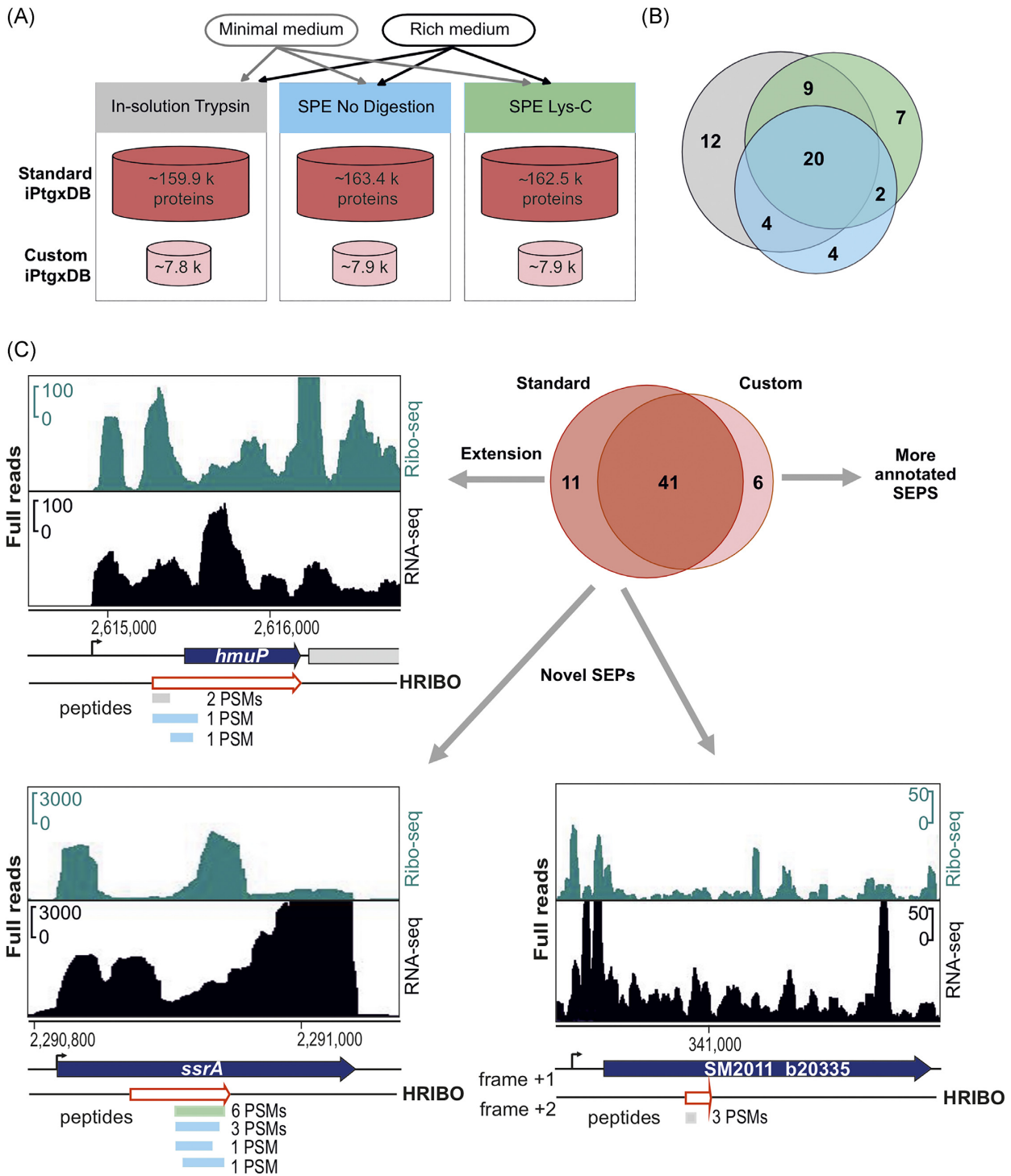
sORFs, and only five and four sORFs started with GTG or TTG, respectively; their stop codon preference was also similar to that of the annotated sORFs (Table S6). Importantly, as the iPtgxDB integrates and consolidates different reference genome annotations and various predictions, we could readily deduce that 11 of the 54 translated sORFs were contained in the RefSeq 2017 annotation, precisely matching their predicted start and stop codons (Table S7). Five candidates matched a RefSeq annotation, but they were shorter. One candidate matched the stop but was only 1 aa longer than the RefSeq annotation. Finally, three candidates matched a GenBank stop codon, but they were shorter than annotated (one of which was in fact again removed in the RefSeq annotation). In summary, Ribo-seq uncovered 37 translated sORF candidates that were novel compared to both GenBank 2014 and RefSeq 2017 annotations (Table S7).

### Both standard and small custom iPtgxDBs informed by Ribo-seq data facilitate novel SEP identification by MS

To validate sORF translation and identify novel SEPs of *S. meliloti* 2011, we then conducted MS-based proteomics using experimental strategies to increase the coverage of the MS-detectable small proteome and two types of search DBs. Cells were cultured either in minimal GMS medium (same as for Ribo-seq) or in rich TY medium, and three complementary sample preparation ap-

proaches were used: 1) tryptic in-solution digest of all proteins (a standard proteomics approach), 2) solid phase enrichment (SPE) of small proteins with subsequent Lys-C digestion, and 3) SPE of small proteins without subsequent digestion (Fig. 5A). Approaches 2 and 3 can identify SEPs whose peptides are not within the detectable range (approximately 7 aa to 40 aa) upon a tryptic digest (Tyanova et al. 2016).

For the DB searches, we first relied on a standard (full) iPtgxDB (Omasits et al. 2017) that hierarchically integrates reference genome annotations, *ab initio* gene predictions, and *in silico* ORF predictions (see Methods). The overlap and differences of all annotation sources were captured and consolidated in a composite gene identifier. Moreover, a large but minimally redundant protein search DB (for more details, see [https://iptgxdb.expasy.org/creating\\_iptgxdb/](https://iptgxdb.expasy.org/creating_iptgxdb/)) is created, as well as a GFF that allows researchers to overlay experimental evidence, such as RNA-seq, Ribo-seq, or proteomics data. Individual iPtgxDBs must be prepared for different proteases (see Methods). For trypsin, the standard iPtgxDB contained close to 160k protein entries of approximately 103k annotation clusters (Table S3.1), that is, genomic loci that share the stop codon but have different predicted protein start sites. Approximately 92% of the peptides unambiguously identify one protein entry, which are called class 1a peptides that facilitate downstream data analysis and allow to swiftly identify novel proteoforms or SEPs. Although standard iPtgxDBs are very large, when



**Figure 5.** Mass spectrometry-based identification of known and novel small open reading frame-encoded proteins (SEPs). **(A)** Experimental set-up for the proteomics analyses. Bacteria were grown in minimal and rich media, and protein extracts were further processed with tryptic in-solution digest (gray), solid-phase enrichment (SPE) of small proteins with subsequent Lys-C digestion (green), or without further digestion (blue). **(B)** Overlap of the identified SEPs by experimental approach; trypsin identified 45 SEPs; compared with the trypsin approach, Lys-C identified 38 SEPs (nine novel, 24%), and the approach without digestion found 30 SEPs (six novel, 20%). **(C)** Novel/unique identifications uncovered by the standard integrated proteogenomic search databases (iPtgxDB) and the small custom iPtgxDB. Standard iPtgxDB: Three peptides imply a 14 aa longer proteoform (60 aa) for HmuP than annotated; four peptides of the tmRNA-encoded proteolysis tag were identified; one peptide (3 peptide spectrum matches [PSMs]) implied a novel SEP (34 aa) internal to the genomic region that also encodes SM2011\_b20335 but in a different frame. Spectra identifying these peptides are shown in Fig. S5. These identifications were also predicted by HRIBO based on Ribo-seq. Finally, six annotated proteins (GenBank 2014 and/or RefSeq 2017) were identified only in the search against the small custom iPtgxDB, as they did not accumulate enough spectral evidence in the search against the standard iPtgxDB (Table S4).

combined with stringent FDR filtering, they have provided convincing results in the past for the identification of novel SEPs that withstood independent validation efforts (Omasits et al. 2017, Bartel et al. 2020, Melior et al. 2020). However, as large DBs inflate the search space, they complicate protein inference and FDR estimation, resulting in a large likelihood of a random hit, especially for SEPs (Burger 2018, Nesvizhskii 2010, Fancello and Burger 2022). Importantly, the 266 top Ribo-seq-implied novel candidates (Fig. 4A) allowed us to explore whether a much smaller custom iPtgxDB may provide additional value for the identification of annotated or novel SEPs. Adding these 266 candidates to the three reference genome annotations (RefSeq, GenBank, Genoscope) and the Prodigal *ab initio* gene predictions resulted in a 20-fold smaller custom iPtgxDB (Fig. 5A) (approximately 8000 protein entries in 7300 annotation clusters), with a higher percentage of class 1a peptides (nearly 98%; Table S3.3).

The acquired MS-spectra were searched against the standard and small iPtgxDBs, and the results were compiled and stringently filtered, requiring more PSM evidence (see Methods) for *ab initio* and *in silico* predictions (Varadarajan et al. 2020a,b). Overall, more than 1200 annotated proteins were detected at an estimated protein FDR of approximately 1%. The SPE-based small protein enrichment steps uniquely identified 160 of these proteins (Fig. S4A). Notably, the search against the small custom DB accounted for 112 unique identifications (Fig. S4B) due to improved search statistics. The MS-identified proteins included 58 SEPs, with  $\leq 70$  aa, 47 of which were annotated (GenBank 2014 and/or Refseq 2017) (Table S4). Similar to the overall results, the two SPE approaches also added unique SEPs: while 45 of the 58 MS-detected SEPs were identified with standard trypsin-based digestion, 13 SEPs were uniquely identified after processing the samples with SPE and either a Lys-C digest (9 of 38 not covered by trypsin) or no proteolytic digest (6 of 30 not covered by trypsin) (Fig. 5B). Most MS-identified SEPs were between 60 and 70 aa long (67%), and the smallest detected SEP was 20 aa long. They include abundantly expressed proteins (the cold shock proteins SM2011\_RS25125 and SM2011\_RS00515, and SM2011\_RS31025, a 50S ribosomal protein L32) (Table S4) down to candidates identified by only 2 PSMs, such as a 59 aa hypothetical protein, which we refer to as SEP7 (see next section). Among the 85 GenBank-annotated SEPs identified with high confidence as translated by Ribo-seq (Fig. 3A), 31 were identified by MS. Among the 54 SEPs missing from the GenBank 2014 annotation and identified as translated by Ribo-seq (Fig. 4A), five were identified by MS, and those are present in the Refseq 2017 annotation (Table S4).

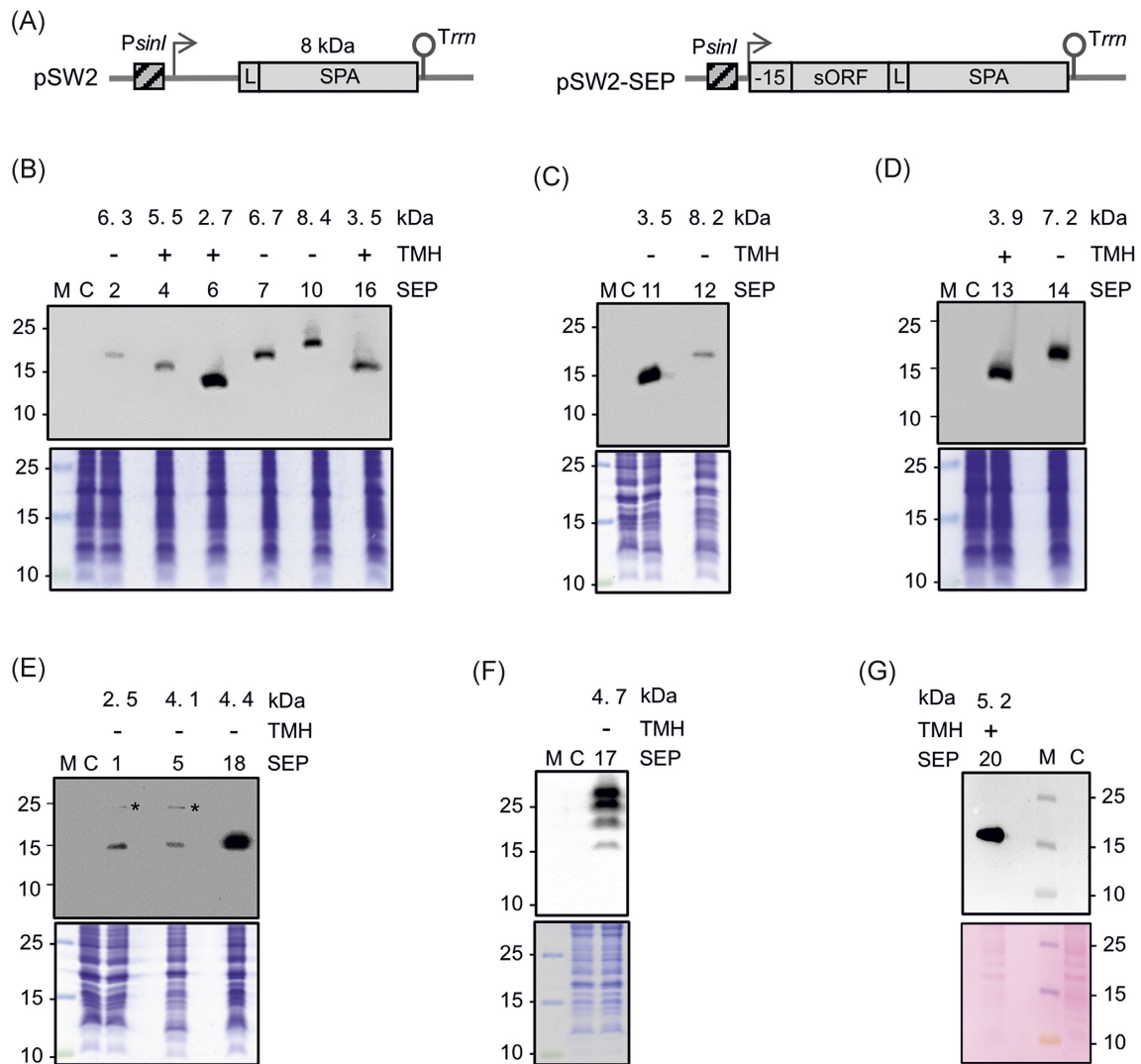
Importantly, both searches added unique identifications. The search versus the full iPtgxDB added 11 potential novel SEPs or longer proteoforms than annotated, which were *in silico* predictions that were excluded from the small custom iPtgxDB. A 14 aa longer proteoform of HmuP was identified by three peptides with 4 PSMs (Fig. S5A). Here, when manually inspecting the Ribo-seq data, it perfectly agreed with the extension of the 46 aa GenBank annotation (Fig. 5C). This finding exemplifies how proteomics and Ribo-seq jointly identify a novel proteoform. Furthermore, the tmRNA-encoded 12 aa proteolysis tag peptide was uniquely identified, which marks incompletely translated proteins for degradation (Karzai et al. 2000) (Fig. 5C). The tag peptide was identified as a C-terminal part of an *in silico* predicted 23 aa SEP included in the standard iPtgxDB. It was only detected in the minimal medium by four peptides: one in the Lys-C digest and three from the search without protease (Fig. 5C and S5B). Mutation of the start codon of the 23 aa sORF had no effect on the translation of the proteolysis tag peptide, in line with the mechanism proposed for this

split tmRNA (Keiler et al. 2000, Ulvé et al. 2007) (Fig. S6). An example of a completely novel 34 aa SEP is shown in the third panel of Fig. 5C (see also Fig. S5C); it is located in a genomic region that harbors an annotated CDS and is translated in a different frame. The novel sORF has Ribo-seq support (TE 0.4) but did not pass our stringent Ribo-seq cut-offs. Notably, the search against the small custom iPtgxDB added six unique SEP identifications (again, due to better search statistics) (Fig. 5C). Four of them were also among the 85 GenBank-annotated SEPs identified by Ribo-seq data (SM2011\_RS33030, SM2011\_RS33620, SM2011\_RS33980, and SM2011\_a6027), lending independent support for their expression (Table S4). SM2011\_RS33620 belongs to the arginine-rich DUF1127 family of proteins, the members of which are involved in phosphate and carbon metabolism in *Agrobacterium tumefaciens* (Kraus et al. 2020), and in RNA maturation and turnover in *Rhodobacter* (Grützner et al. 2021). In addition, the abovementioned RefSeq-annotated SEP7 was identified (Fig. S5D). Two other SEPs (one of them novel) were identified with only 1 PSM (Fig. S5E and S5F), which was below our threshold, but had strong Ribo-seq support (SEP1, SEP20; see next section).

### Validation of a subset of Ribo-seq-implied small proteins by Western blot analysis

Since out of the 54 high-confidence Ribo-seq-implied sORFs that were not contained in the GenBank 2014 annotation only five were detected with at least 2 PSMs in the MS analysis (Table S4), we attempted additional validation by epitope tagging and Western blot analysis (Fig. 6). Nineteen sORFs were selected that (i) cover a broad range of TE values, (ii) start with one of the three main start codons (ATG: 16 sORFs, GTG: two sORFs, or TTG: one sORF), and (iii) were either added in the RefSeq 2017 annotation (five sORFs) or were novel with respect to these two annotations (14 sORFs). The corresponding proteins were designated SEP1 to SEP19 (Table S4). They included three of the candidates that were also detected by MS (SEP7: 2 PSMs; SEP 10: 59 PSMs; SEP17: 29 PSMs). SEP1 was only identified by 1 PSM, that is, below the threshold (Table S4; Fig. S5E), but with strong Ribo-seq support (highest TE among the 54 high-confidence Ribo-seq candidates; Table S7). Moreover, three SEP candidates below 30 aa (SEP1, SEP3, and SEP6) and four candidates with a predicted transmembrane helix (TMH) (SEP4, SEP6, SEP13, and SEP16; Table S4) were analyzed. As a 20th candidate (SEP20), we included a conserved annotated sORF located in the cytochrome C oxidase cluster *ctaCDBGE* between *ctaB* and *ctaG* (GenBank 2014 annotation), which also contains a predicted TMH. SEP20 was identified by 1 PSM in the MS analysis (Fig. S5F) and did not pass the stringent HRIBO criteria for translated candidate sORFs (Table S3, TE = 6.99, RPKM of < 10 in replicate 1) but showed strong read coverage in the Ribo-seq library (Table S6).

Each sORF was cloned together with its  $-15$  nt 5'-UTR region into plasmid pSW2, thus containing its putative ribosome binding site in frame to the SPA-tag encoding sequence (Fig. 6A and Fig. S1). Transcription of the sORF::spa fusion was under the control of a *S. meliloti sinI* promoter ( $P_{sinI}$ ) of moderate strength, which is constitutively active (Charoenpanich et al. 2013). Thus, the detection of a SEP-SPA fusion protein by Western blot analysis would indicate sORF translation. The Western blot analysis of crude lysates of cultures grown in minimal medium using FLAG-directed antibodies revealed signals for 15 of the 20 candidates, including SEP20 (Fig. 6B and G). For 12 candidates, one band consistent with their predicted SEP length was detected. For SEP1 and SEP5, on top of the expected SEP-SPA bands, slow migrating bands at approximately 25 kDa (see asterisks in Fig. 6E) were detected, which



**Figure 6.** Detection of 15 sequential peptide affinity (SPA)-tagged small open reading frame-encoded proteins (SEPs) in *Sinorhizobium meliloti* crude lysates. **(A)** Schematic representation of the empty plasmid pSW2 (contains no promoter and no ribosome-binding site upstream of the linker [L] and SPA-encoding sequence) and a pSW2-SEP plasmid for the analysis of sORF translation. The constitutive  $P_{sinI}$  promoter (hatched box), the corresponding TSS (flexed arrow), the sORF coding sequence with its -15-nt-long region, the SPA-tag (with its molecular size indicated) preceded by a linker (L) (gray boxes), and the *Trm* terminator (hairpin) are depicted. **(B)** to **(F)** Western blot analysis of crude lysates (upper panels) and the corresponding Coomassie-stained gels, and **(G)** corresponding Ponceau-stained membrane for selected SEPs. Monoclonal FLAG-directed antibodies were used. Migration of marker proteins (in kDa) is shown on the left side. \*Unspecific signal. Above the panels, the numbers of the analyzed SEP protein (Table S7), the presence (+) or absence (-) of a predicted TMH, and the molecular size (in kDa) of the SEP without the SPA tag are given. M: protein marker. C: empty vector control, lysate from a strain containing pSW2.

probably corresponded to a non-specific signal, as they were also detected in some EVC samples after lysate fractionation (Fig. S7).

The bands of the tagged SEP1 and SEP5 ran similarly, although SEP1 is smaller than SEP5, as indicated above the panel (Fig. 6E). Probably, the aberrant migration of SEP1 is due to its acidic aa composition (pI of 4.18) (Guan et al. 2015). SEP17 showed multiple bands, with a weak and fast migrating band at approximately 15 kDa, which probably corresponds to the monomeric SEP17-SPA protein, and three strong and slow migrating bands, which could indicate protein oligomerization (Fig. 6F). Overall, the translation of SEPs with alternative start codons, that is, GTG (SEP10 and SEP14) and TTG (SEP7), and of the five candidates missed in the GenBank 2014 annotation but included by Refseq (2017) (SEPs Nr. 4, 7, 10, 17, and 18; SEP18 corresponds to the sORF upstream of *leuC*, Fig. 2E), was validated. Importantly, this analysis confirmed the translation of six novel SEPs (SEPs Nr. 1, 6, 11, 13, 14, and 16),

including two of the three SEP candidates shorter than 30 aa. Finally, our observation that 11 (out of 16) sORFs without MS support but with high-confidence Ribo-seq data were validated in the Western blot analysis shows the power of Ribo-seq to detect novel translated sORFs.

Since the analysis of exclusive or predominant subcellular localization is valuable for linking hypothetical proteins without any annotation to some potential function (Stekhoven et al. 2014), we decided to investigate the subcellular localization of the validated SPA-tagged SEPs by Western blot analysis of the supernatant (S100) and pellet (P100) fractions (see Methods) (Fig. S7). As expected, the predicted TMH-containing proteins SEP4, SEP6, SEP13, SEP16, and SEP20 were detected exclusively or predominantly in P100, which contains ribosomes and membranes, whereas the predicted cytoplasmic proteins SEP5 and SEP12 were detected exclusively in the S100 fraction (Fig. S7). The remain-

ing eight SEPs were detected exclusively or partially in the P100 fraction, suggesting that they could be associated with membrane complexes or ribosomes (SEP10 and SEP18 show similarities to the ribosomal proteins S21 and L7/12) or be prone to aggregation in their recombinant, tagged form.

### Conservation and potential functions of *S. meliloti* novel small proteins

As described above, we detected the translation of 48 sORFs missing in the GenBank 2014 and Refseq 2017 annotations (37 identified by Ribo-seq and additional 11 by MS), which we refer to as novel. Since conserved SEPs are likely to be functional, we used tBLASTn (Gertz et al. 2006) to examine the conservation of the proteins encoded by the 48 novel sORFs (Fig. 7; Table S8). The tBLASTn searches were conducted in bacteria with parameters previously established to identify conserved bacterial sORFs (Allen et al. 2014) (see Methods). We found a wide range of conservation, from an sORF detected in only four *S. meliloti* strains overall, to sORFs conserved at different higher taxonomic levels, to highly conserved sORFs present in different bacterial phyla (Fig. 7). Among the 14 sORFs encoding SEPs with < 30 aa (excluding the tmRNA sORF64), four are conserved beyond *S. meliloti*. One of the most widely conserved novel SEPs is a 64 aa small protein detected only by MS (sORF61 in Fig. 7). It was identified as a product of an *in silico* predicted sORF, with 3 PSMs in lysates from MM cultures (Fig. S5G; Table S4). However, no expression at the level of RNA was detected at its locus, possibly suggesting high protein stability. sORF61 has homologs in several bacterial phyla and multiple paralogs, with a maximal aa sequence identity of 64% on each replicon in *S. meliloti* 2011. Despite its wide distribution and strong conservation, its function is unknown. Overall, excluding the tmRNA, we detected seven sORFs conserved beyond the family *Rhizobiaceae*, suggesting that the corresponding SEPs may have important general functions.

Furthermore, we used TMHMM (Krogh et al. 2001) and PSORTb (Yu et al. 2010) to predict the presence of transmembrane helices and the subcellular localization of the 48 novel SEPs. Localization in the cytoplasmic membrane was predicted for seven SEPs using at least one of the tools (Fig. 7; Table S8). Among them are the Ribo-seq-identified and Western blot analysis-validated SEP6 (prediction by both TMHMM and PSORTb) and SEP16 (prediction by TMHMM only), which were detected with strong signals predominantly in the P100 fraction (see Fig. S7). The corresponding sORF6 and sORF16 are conserved in *Hyphomicrobiales* (Fig. 7). No proteins with predicted membrane localization were found among the 11 MS-detected SEPs (Fig. 7). Notably, two of the 48 novel SEPs harbor a predicted SpII cleavage site and are thus probably lipoproteins (Fig. 7). Lipoproteins play important roles in physiology, signaling, cell envelope structure, virulence, and antibiotic resistance (Kovacs-Simon et al. 2011); however, as previously reported, they are often missed in prokaryotic genome annotations (Omasits et al. 2017).

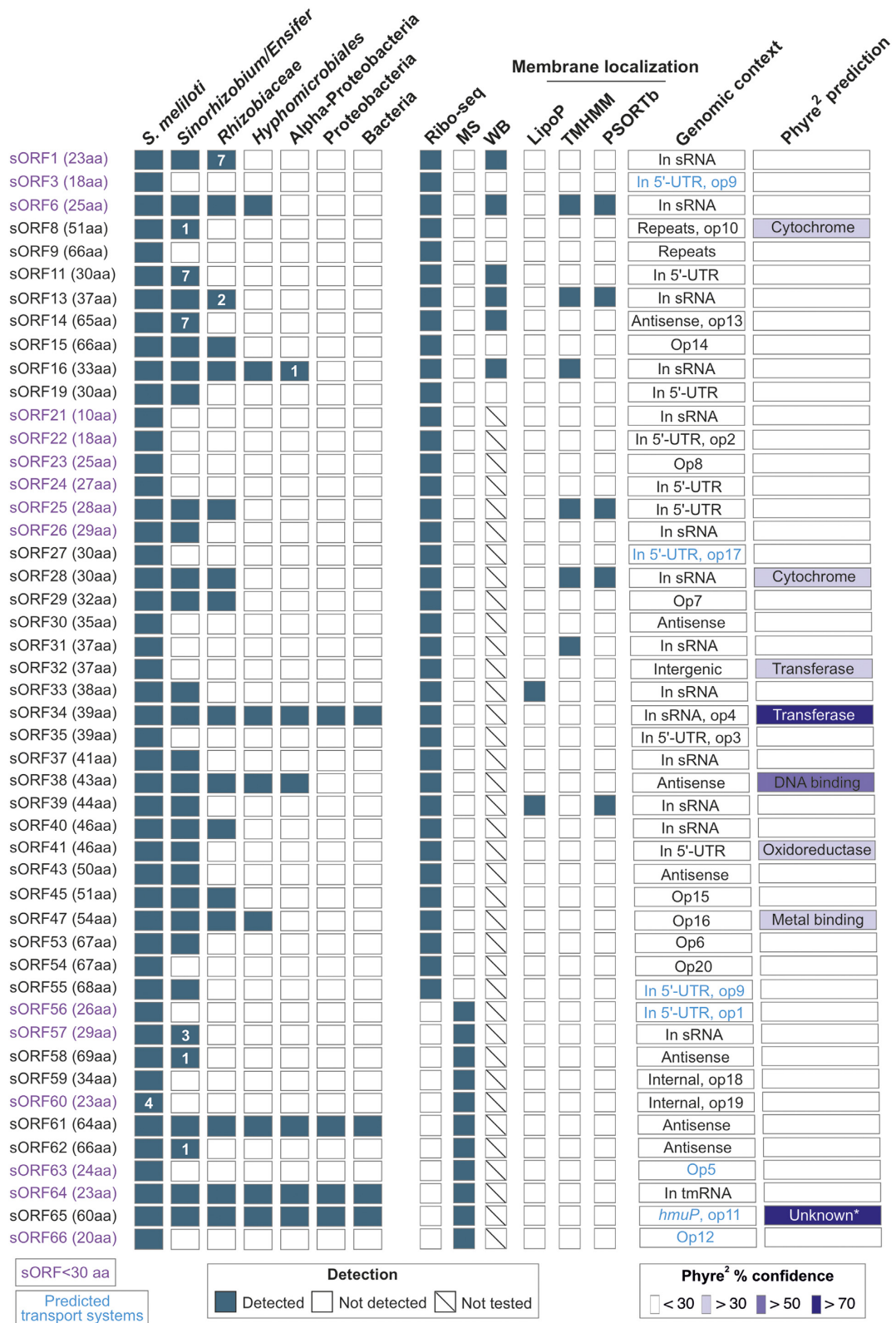
Moreover, we used Phyre2 (Kelley et al. 2015) to gain insights into the potential functions of novel SEPs with  $\geq 30$  aa by analyzing their similarity to proteins with known tertiary structures (Fig. 7; Table S8). Best hits with a confidence homology of  $\geq 30\%$  were obtained for eight novel SEPs (Fig. 7). The highest confidence homology suggesting a function was obtained for the SEPs encoded by sORF38 (DNA binding; 18 of the 43 aa residues were modeled with 66% confidence homology; conserved in Alphaproteobacteria) and sORF34 (bleomycin resistance; 37 of the 39 aa residues were modeled with 92% confidence homology; conserved

among Bacteria). The HmuP extension (sORF65 in Fig. 7; see also Fig. 5C) was modeled with 98% confidence along 59 of its 60 aa residues; however, according to Phyre2, the function is unknown. Overall, obtaining clear functional predictions was not possible even for conserved SEPs, most probably due to their small size.

Additionally, to assess potential functions of the 48 novel translated sORFs and/or SEPs, we used the *S. meliloti* 2011 RNA-Seq data by Sallet et al. (2013), who annotated coding and non-coding (e.g. ncRNAs and UTRs flanking CDSs) transcripts, and compared RNA levels under three different conditions (exponential growth, stationary phase and symbiosis). Out of the 48 novel sORFs, nine do not overlap in sense with annotated transcripts (seven anti-sense and two intergenic sORFs; Table S8). For them and for two sORFs overlapping with repeat elements (sORF8 and sORF9; Table S8), information about RNA levels could not be retrieved. The levels of most transcripts, which comprise the remaining 37 novel sORFs, showed specific abundance changes in the study by Sallet et al. (2013) (Table S8). Eleven of these sORFs are located in 5'-UTRs (Fig. 4C; Table S8), suggesting that at least some of them could act as upstream ORFs (uORFs) and could play a role in the regulation of the expression of the downstream genes. Noteworthy are sORF25, located in the 5'-UTR of *dnaE1* encoding a subunit of the replicative DNA polymerase, and sORF37 in the 5'-UTR of *rpoD* encoding the vegetative sigma factor. These sORFs may have functions in important pathways during bacterial growth. Finally, 13 novel sORFs are encoded in previously annotated ncRNAs. While only the sORF26-transcript had constant expression levels under the three conditions, the remaining 12 novel small mRNAs showed differential expression indicative of possible SEP-regulation and/or potential functions during exponential growth or the stationary phase (Table S8). A specific up-regulation under symbiosis was detected for the 5'-UTR transcripts that contain sORF11 and sORF22, and for the Smb20335 transcript, which overlaps with sORF59 (Table S8; see Fig. 5C above). We note that changes in the RNA levels do not necessarily directly correspond to similar changes in ribosome occupancy and SEP accumulation, which were not tested by Sallet et al. (2013).

Moreover, an operon prediction was carried out to possibly assign a function to some of the novel sORFs (SEPs) by guilt-by-association. We could retrieve a predicted operon assignment for 21 of the 48 novel small protein ORFs (as part of 20 putative operons; Table S8 and Table S9). Notably, seven sORFs are part of operons encoding predicted transport systems (mainly ABC transporters, an MFS transporter and an ion channel; Table S9). These likely represent higher priority targets for an experimental elucidation of their function.

Finally, we suggest the functions for three annotated sORFs/SEPs with validated translation. SEP5 (added in the RefSeq 2017 annotation) is conserved only in *Sinorhizobium*. Its translation was detected with Ribo-seq and Western blot analysis (Fig. 6E; Table S4). The SEP5 sORF contains a cluster of six threonine and three lysine codons near its 3'-end and is located in the 5'-UTR of the aspartate dehydrogenase-encoding gene. Since aspartate is a part of the threonine and lysine biosynthesis pathway (Vitreschak et al. 2004), our observation suggests that this sORF can be involved in the post-transcriptional regulation of the aspartate dehydrogenase gene in *Sinorhizobium*, and SEP5 is possibly a leader peptide. Furthermore, among the annotated SEPs with functional assignment, which were detected by MS, an entericidin A/B family lipoprotein was found (CP004140.1:2141558–2141716, Table S4). The 52 aa protein has a predicted TMH and is conserved in Alphaproteobacteria. Its *A. tumefaciens* homolog, the lipoprotein Atu8019, is involved in specific cell-cell interactions as a



**Figure 7.** Conservation analysis, functional prediction and operon assignment for 48 novel small open reading frames (sORFs) of *Sinorhizobium meliloti* 2011. The conservation analysis was conducted using tBLASTn. The respective hits (see methods for parameters and cutoffs) are broadly summarized at the level of different taxonomic groups. The number of species outside the lower taxonomic unit, which harbors a hit, is given, if at < 10. In addition, the method by which the respective sORF was detected or confirmed is shown (Ribo-seq: ribosome profiling, MS: proteomics, WB: Western blot), as well as the results of predictions for membrane localization (by TMHMM and PSORTb), signal peptide II cleavage sites of lipoproteins (by LipoP), and function (by Phyre2; only hits with confidence levels greater than 30% are shown). For details on Phyre2 prediction and genomic context including operon prediction, see Table S8 and Table S9. sORF1 to sORF55 are a subset of the Ribo-seq-detected, translated sORFs, which are listed in Table S7, and sORF56 to sORF66 represent the novel sORFs identified by proteomics. sORFs encoding small proteins below 30 amino acids are shown in red. The putative sORF64, present in tmRNA, contains the proteolytic tag sequence. The sORF65 corresponds to the N-terminal HmuP extension; outside of Proteobacteria, it is conserved in many genera of Planctomycetes. \*Structural genomics (92% confidence homology to protein of unknown function).



part of outer membrane vesicles (Knoke et al. 2020). Finally, an annotated small protein validated in this work by Ribo-seq and Western blot analysis (1 PSM in the MS) is the above mentioned SEP20 (Fig. 6G, Fig. S5E; Table S4). It contains a predicted TMH and is conserved in the Alphaproteobacteria, and its sORF is part of an uncharacterized cytochrome oxidase operon. This synteny suggests that SEP20 can participate in the assembly and/or function of the corresponding cytochrome oxidase complex, as previously shown for SEP CydX and cytochrome bd oxidase in *Brucella abortus* (Sun et al. 2012).

## Discussion

In this work, we have developed and applied a Ribo-seq workflow to comprehensively map the translome of *S. meliloti* 2011 under free-living conditions in a minimal medium. By combining Ribo-seq and MS-based proteomics in a proteogenomic approach, we added 48 novel SEPs below 70 aa to the *S. meliloti* annotation, that is, an increase in the number of annotated SEPs by approximately 15% compared to the RefSeq 2017 annotation.

Ribo-seq is a powerful technique for detecting translation on a global scale with high sensitivity (Ingolia et al. 2019). However, in contrast to eukaryotic model systems, codon resolution has not yet been achieved in Ribo-seq analyses of bacteria (Mohammad et al. 2019, Venturini et al. 2020, Cianciulli Sesso et al. 2021, Vazquez-Laslop et al. 2022). Trapping ribosomes on mRNA and generating ribosome footprints have remained challenging, requiring careful optimization for each bacterial species. Our Ribo-seq workflow for *S. meliloti* includes ribosome trapping by rapid cooling of the culture without using antibiotics and cell lysis in an adapted buffer, followed by digestion of unprotected RNA by RNase I, which is not inactivated by the ribosomes of *S. meliloti* (Fig. 1B). RNase I has the advantage of precisely cleaving at both 5' and 3' ends of ribosome-protected mRNA without sequence specificity, in contrast to the routinely used MNase (Bartholomäus et al. 2016). The digestion of 5' and 3' regions of translated mRNAs (Fig. 1C, Fig. 3C), higher TEs of annotated CDS in comparison to non-coding RNAs (Fig. 2B), and pronounced ribosome protection up to 16 nt upstream and downstream of start codons (Fig. S2; Fig. 1C and 1D) show the successful establishment of Ribo-seq for *S. meliloti*.

In addition to providing the first genome-wide ribosome-binding map of a *Hyphomicrobiales* member, our Ribo-seq analysis uncovered translation for 85 annotated sORFs and identified 37 novel sORFs missing in the GenBank 2014 and RefSeq 2017 annotations of the *S. meliloti* genome (17 of the overall 54 novel sORFs identified compared to Genbank 2014 were subsequently added in RefSeq2017; this underlines the high quality of our data; Fig. 8; Table S4). The translated sORFs were found on all three replicons and had similar preferences for start and stop codons independently of whether they were annotated or novel (Fig. 3 and Fig. 4). The novel sORFs were generally shorter than the annotated ones (Fig. 4B; Table S4), clearly showing the advantage of the Ribo-seq method for SEP discovery. Many of the novel sORFs were probably not annotated due to their location in short transcripts considered as non-coding RNAs or asRNAs or in 5'- and 3'-UTRs (Fig. 4C).

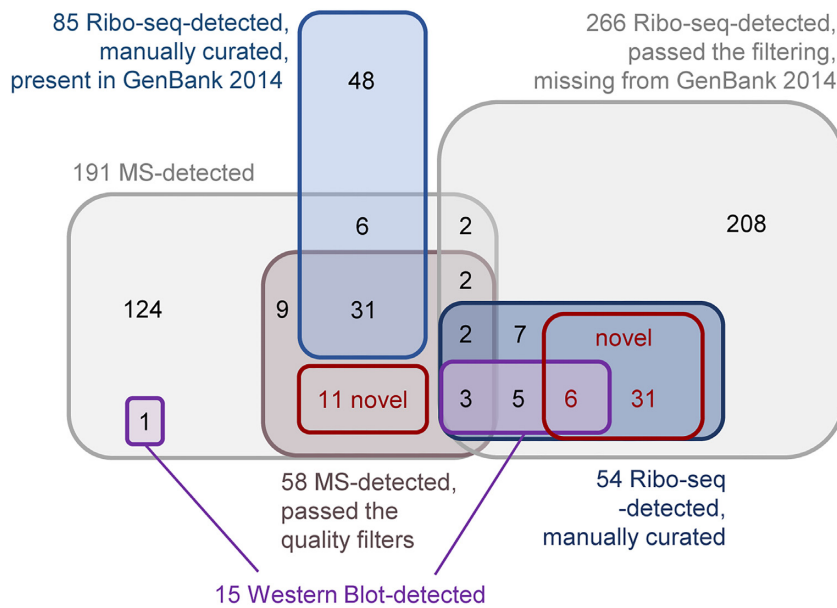
Several translated novel sORFs internal to annotated genes (nested ORFs; Gray et al. 2022) were also predicted by our Ribo-seq data. However, they were excluded from the analysis as additional evidence is needed to confirm their existence. Targeted detection of translation initiation sites is useful in uncovering such sORFs by Ribo-seq (Meydan et al. 2019, Weaver et al. 2019), a strategy beyond the scope of our study. However, the existence of an internal sORF with Ribo-seq coverage was supported by the MS de-

tection of a novel, 34-aa-long SEP translated in a different frame in the genomic region encoding SM2011\_b20335 (Fig. 5C and Fig. S5G; sORF59 in Fig. 7). The MS-detected SEPs encoded by sORF56 and sORF60 are also internal to annotated genes (Table S8).

A challenge in defining novel sORFs for any genome is that annotations from different reference genome annotation centers can differ substantially for an identical sequence and change over time (Omasits et al. 2017); that is, CDS are being added but are also removed in more recent annotations (see Fig. 2E and F, and the 'master' Table S4). Accordingly, two of the 48 novel SEPs are now bona fide-predicted CDS in the latest RefSeq 2022 annotation, with MS-evidence of a single PSM found with the custom iPTgxDB, whereas two other Ribo-seq-identified sORFs have variable pseudogene status in different annotation releases (see Table S4). iPTgxDBs, which integrate existing reference annotations and add *in silico* predicted ORFs in all six frames to virtually cover the entire protein coding potential of a prokaryote, can be used to overcome such problems and enable MS-based detection of novel SEPs (Omasits et al. 2017). Here, in addition to a standard large iPTgxDB of *S. meliloti* (Melior et al. 2020), we applied the concept of a small, custom iPTgxDB lacking *in silico* predictions and including the top predictions from our experimental Ribo-seq data. This custom iPTgxDB is approximately 20-fold smaller and benefits statistics and FDR estimation (Blakeley et al. 2012, Li et al. 2016). Notably, although the identification of 11 *in silico* predicted novel sORFs was possible only with the standard iPTgxDB, the small iPTgxDB contributed substantially to the validation of annotated sORFs, increasing the number of SEPs with experimental support by 10% (Fig. S4). The detection of more SEPs was also facilitated by applying three experimental approaches, two of which included enrichment of small proteins. The MS detection of enriched SEPs without a proteolytic digest, including, for example, the 12 aa proteolysis tag encoded by tmRNA (Fig. 5C and Fig. S5), shows that this method can be useful for the identification of SEPs.

The validation of translation by Western blot analysis for 15 out of 20 analyzed novel SEPs with Ribo-seq support (Fig. 6B–F; Table S7), only three of which were detected by MS with at least 2 PSMs (Table S4), underlines the power of the Ribo-seq technique for identification of translated sORFs. The example of SEP7 (Fig. 6B; 59 aa, restriction endonuclease-like, conserved in *Rhizobiaceae*), which was added to the RefSeq 2017 annotation and was detected by 2 PSMs using the small, custom iPTgxDB, again illustrates the added value of the latter. Detection of translation for the novel SEP1 (23 aa, conserved in *Rhizobiaceae*) by Western blot analysis (Fig. 6E) and Ribo-seq (highest TE among the non-annotated translated sORFs, Table S7), even though it was identified by only 1 PSM in the MS analysis (Fig. S5E), suggests that putative SEPs with 1 PSM can be truly expressed, real small proteins. Similarly, the annotated SEP20 (46 aa, conserved in Alphaproteobacteria) was confirmed by Western blot analysis (Fig. 6G), although it had only 1 PSM (Fig. S5F) and did not pass the filtering of the Ribo-seq data (Table S6). We suggest that the conservation analysis of putative SEPs, which have minimal MS evidence (e.g. 1 PSM) and/or correspond to sORFs that did not pass the very stringent manual curation of the Ribo-seq data, can help define SEP candidates with potentially important functions that can be validated and analyzed in the future.

Despite the lower sensitivity of MS compared with Ribo-seq, using MS we detected 16 additional SEPs that were not identified as translated by Ribo-seq. Eleven of them were novel, showing the importance of complementary methods for comprehensive analysis of bacterial small proteomes. The reported numbers of validated and novel sORFs and their encoded SEPs are affected by



**Figure 8.** Translated sORF (SEP) candidates and their detection by different methods. Overlap between the 191 MS-detected SEP candidates (annotated and non-annotated), the 85 Ribo-seq-detected, manually curated sORFs present in the Genbank 2014 annotation and the 266 Ribo-seq-detected sORF candidates, which are missing in the GenBank 2014 annotation (Table S4). SEPs and translated sORFs, which are missing from both the GenBank 2014 and Refseq 2017 annotations, were designated 'novel'. Two of the 11 Ribo-seq-detected novel sORFs are present in the RefSeq 2022 annotation (Table S4). Passing the stringent filtering criteria and (in the case of Ribo-seq) the manual curation, and detection by more than one method increases the confidence in sORF translation (for details see the master Table S4).

the somewhat arbitrary cut-off of 70 aa. In fact, our data provide evidence for the translation of three additional proteins below 100 aa, which are considered small in other studies (Baumgartner et al. 2016, VanOrsdel et al. 2018, Kaulich et al. 2021) (see Table S4). One of them (identifier CP004140.1:3367861–3368100) corresponds to a 79 aa ChemGenome predicted protein, the N-terminus of which is encoded by the pseudogene SM2011\_RS34080 (annotated as transcriptional regulator with a frame shift). Two additional exact copies of this ChemGenome-predicted sORF and the matching pseudogene (SM2011\_RS34090 and SM2011\_RS34095) are also present in this genomic region, which differs between the *S. meliloti* strains 1021 and 2011 (Sallet et al. 2013). Their promoters and the 5'-terminal CDS parts corresponding to the pseudogenes evolved by duplications of *fixK*, a gene controlled by the symbiotically relevant transcriptional regulator *FixJ* (Ferrières et al. 2004).

For the 48 novel sORFs listed in Fig. 7 and Table S8, we suggest that they are translated. Our high confidence in the translation of the 37 novel, Ribo-seq-detected sORFs relies on passing the stringent filtering criteria and careful manual curation, while the identification of the 11 novel, MS-detected SEPs is based on passing the MS quality filters and their detection with at least three PSMs. The 11 MS-detected SEPs do not correspond to sORFs that passed the Ribo-seq filtering (e.g. sORF59, out of frame in SMB20335; see also Fig. 5C) or to transcripts detected in our study (e.g. the highly conserved sORF61). The latter can be explained by the generally very short half-lives of bacterial mRNAs compared to protein half-lives (Bonnefoy et al. 1989, Bernstein et al. 2002, Chai et al. 2016). Therefore, it can be expected to detect some proteins without detecting their corresponding mRNAs (Omasits et al. 2013), an important argument for using both MS and Ribo-seq for sORF identification. Similarly, the failure to detect a recombinant, tagged SEP by our Western blot approach does not lower the confidence in the Ribo-seq-detected ribosome occupancy of its sORF. To understand why the Western blot result was negative, additional experiments

are needed. For example, the translation product could be below the detection limit, possibly because of its recombinant form or ribosome occupancy could have regulatory function and a non-functional SEP-product could be short-lived.

As mentioned above, our datasets include many sORF candidates that did not pass our stringent criteria, but also sORFs, whose translation was suggested by more than one method, likely increasing the confidence in the SEP existence (summarized in Fig. 8 and Table S4). However, we point out that additional efforts allow to assign a higher confidence to annotated and novel SEPs. For example, the Ribo-seq detection of translated sORFs can be additionally supported by targeted detection of translation initiation sites (Meydan et al. 2019, Weaver et al. 2019), which is still not established for *S. meliloti*. Further, increased confidence in the MS-detection can e.g. be achieved by validation using another, more sensitive mass spectrometry technique called parallel reaction monitoring (Omasits et al. 2017) or by matching experimentally observed spectra to those obtained from synthetic peptides (Petruschke et al. 2021). Both these approaches are quite expensive when many SEP candidates are analyzed, but would provide additional support on top of the Western blot analysis carried out here. Finally, the Western blot results can be further validated by using (i) additional reporter constructs, which cover potential transcriptional and post-transcriptional regulatory regions of the gene of interest (Scheuer et al. 2022), (ii) a marker-less tag insertion in the original genomic locus (Hemm et al. 2010), and (iii) detailed characterization of the subcellular localization of the tagged SEP (Fontaine et al. 2011).

The functions of small proteins are difficult to predict *in silico*, often because they are too small to harbor known protein domains or motifs (Ahrens et al. 2022). In addition, for SEPs smaller than 30 aa *in silico* analysis by Phyre2 is still impossible. Keeping these limitations in mind, we present a list of putative functions corresponding to Phyre2 best hits (Table S8). Since modeling of a partial SEP sequence by Phyre2 may provide a hint of potential

interactions with other proteins or protein complexes, we mention predictions based on greater than 30% confidence homology in Fig. 7, including the predicted DNA-binding function of the 43 aa SEP38 and a potential role in bleomycin resistance of the 39 aa SEP34. SEP function can also be predicted based on gene synteny (Ahrens et al. 2022), for example for SEPs encoded in 5'-UTRs (e.g. the RefSeq 2017-annotated SEP5, which is a potential leader peptide; see also Table S8) or in operons with predicted functions (e.g. encoding ABC transporters or ion channels, Table S9; SEP20 encoded in a cytochrome oxidase operon). Our findings show that, excluding the tmRNA sORF, 13 out of the 48 novel SEPs (sORFs) are conserved in *Rhizobiaceae*, seven in *Hyphomicrobiales*, and three in at least two bacterial phyla, which likely suggests physiological relevance. Most of the translated sORFs or SEPs were detected in logarithmic cultures grown in a minimal medium, where bacteria synthesize virtually all metabolites for cell reproduction. Thus, some of these SEPs can be of general importance for growth or are needed for survival and competitiveness under oligotrophic conditions in soil and rhizosphere.

In summary, our work shows that a combination of methods can increase the number of experimentally validated SEPs. Using Ribo-seq, MS, and Western blot analysis of C-terminally tagged proteins, we provide evidence for the translation of 48 SEPs with  $\leq 70$  aa to be added to the annotation of *S. meliloti*, thus substantially increasing the number of cataloged SEPs. With the MS data, the corresponding full and small custom iPTgxDBs, and importantly, the first Ribo-seq analysis of a *Hyphomicrobiales* member, which can be viewed with an interactive online JBrowse instance (<http://www.bioinf.uni-freiburg.de/ribobase>), our study provides valuable resources for future studies on and beyond the small proteome.

## Author contributions

CMS, RB, and EEH initiated the project. LH, RG, SM, RS, SA, CHA, BH, and EEH designed the experiments and analyzed the data. LH established Ribo-seq for *S. meliloti*, performed the Ribo-seq analysis, manually examined the implied novel sORFs, explored evolutionary conservation, and predicted the function of novel SEPs. RG performed the bioinformatic processing of the Ribo-seq data. SM conducted the mass-spectrometry experiments, database searches and manually evaluated spectra implying novel SEPs. BH and CHA contributed to the proteogenomic analysis to identify known and novel SEPs, explored the value of iPTgxDBs for Ribo-seq data, consolidated experimental results in a master table, and identified annotation differences. RS, SA, and SBW performed the cloning and Western blot analyses. LH, CHA, and EEH mainly wrote the manuscript with input and feedback from the other authors. DB, RB, CMS, CHA, and EEH supervised the research and provided resources and funding. All authors approved the submitted version.

## Acknowledgements

We thank Stephanie Färber for the technical assistance in growing *S. meliloti* to establish Ribo-seq, and Pierre-Alexander Mücke, as well as Claudia Hirschfeld, for MS measurements. The cell-free lysates for the MS analysis were provided by Hendrik Mellior. We thank Sarah L. Svensson for the critical feedback on the manuscript.

## Supplementary data

Supplementary data are available at [FEMSML](https://www.femsml.org) online.

**Conflict of interest statement.** The authors declare that they have no conflicts of interest.

## Funding

This work was supported by the German Research Foundation (DFG) priority program SPP2002 'Small Proteins in Prokaryotes, an Unexplored World' grants (grant SH580/7-1 and SH580/7-2 to CMS, grant BA 2168/21-2 to RB, grant BE 3869/5-1 and BE 3869/5-2 to DB, and grant Ev 42/7-1 to EEH). Additional funding was received from the Swiss National Science Foundation (grant 197391) to CHA and from DFG (GRK2355 project number 325443116) to EEH and Germany's Excellence Strategy (CIBSS—EXC-2189—Project ID 390939984) to RB. Computational resources were provided by the BMBF-funded de.NBI Cloud within the German Network for Bioinformatics Infrastructure (de.NBI) (031A532B, 031A533A, 031A533B, 031A534A, 031A535A, 031A537A, 031A537B, 031A537C, 031A537D, 031A538A).

## References

- Ahrens CH, Wade JT, Champion MM et al. A practical guide to small protein discovery and characterization using mass spectrometry. *J Bacteriol* 2022;**204**:e0035321.
- Allen RJ, Brenner EP, VanOrsdel CE et al. Conservation analysis of the CydX protein yields insights into small protein identification and evolution. *BMC Genomics* 2014;**15**:946.
- Aoyama JJ, Raina M, Zhong A et al. Dual-function Spot 42 RNA encodes a 15-amino acid protein that regulates the CRP transcription factor. *Proc Natl Acad Sci USA* 2022;**119**:e2119866119.
- Barra-Bily L, Fontenelle C, Jan G et al. Proteomic alterations explain phenotypic changes in *Sinorhizobium meliloti* lacking the RNA chaperone Hfq. *J Bacteriol* 2010;**192**:1719–29.
- Bartel J, Varadarajan AR, Sura T et al. Optimized proteomics workflow for the detection of small proteins. *J Proteome Res* 2020;**19**:4004–18.
- Bartholomäus A, Del Campo C, Ignatova Z. Mapping the non-standardized biases of ribosome profiling. *Biol Chem* 2016;**397**:23–35.
- Baumgartner D, Kopf M, Klähn S et al. Small proteins in cyanobacteria provide a paradigm for the functional analysis of the bacterial micro-proteome. *BMC Microbiol* 2016;**16**:285.
- Becker A, Bergès H, Krol E et al. Global changes in gene expression in *Sinorhizobium meliloti* 1021 under microoxic and symbiotic conditions. *Mol Plant Microbe Interact* 2004;**17**:292–303.
- Becker AH, Oh E, Weissman JS et al. Selective ribosome profiling as a tool for studying the interaction of chaperones and targeting factors with nascent polypeptide chains and ribosomes. *Nat Protoc* 2013;**8**:2212–39.
- Beringer JE. R factor transfer in *Rhizobium leguminosarum*. *J Gen Microbiol* 1974;**84**:188–98.
- Bernstein JA, Khodursky AB, Lin PH et al. Global analysis of mRNA decay and abundance in *Escherichia coli* at single-gene resolution using two-color fluorescent DNA microarrays. *Proc Natl Acad Sci U S A*. 2002;**99**:9697–702.
- Blakeley P, Overton IM, Hubbard SJ. Addressing statistical biases in nucleotide-derived protein databases for proteogenomic search strategies. *J Proteome Res* 2012;**11**:5221–34.
- Bonnefoy E, Almeida A, Rouviere-Yaniv J. Lon-dependent regulation of the DNA binding protein HU in *Escherichia coli*. *Proc Natl Acad Sci USA* 1989;**86**:7691–5.
- Buels R, Yao E, Diesh CM et al. JBrowse: a dynamic web platform for genome visualization and analysis. *Genome Biol* 2016;**17**:66.

- Burger T. Gentle introduction to the statistical foundations of false discovery rate in quantitative proteomics. *J Proteome Res* 2018;**17**:12–22.
- Casse F, Boucher C, Julliot JS *et al.* Identification and Characterization of Large Plasmids in *Rhizobium meliloti* using Agarose Gel Electrophoresis. *J Gen Microbiol* 1979;**113**:229–42.
- Cassidy L, Kaulich PT, Maaß S *et al.* Bottom-up and top-down proteomic approaches for the identification, characterization and quantification of the low molecular weight proteome with focus on short open reading frame-encoded peptides. *Proteomics* 2021;**21**:e2100008.
- Cassidy L, Kaulich PT, Tholey A. Depletion of High-Molecular-Mass Proteins for the Identification of Small Proteins and Short Open Reading Frame Encoded Peptides in Cellular Proteomes. *J Proteome Res* 2019;**18**:1725–34.
- Chai Q, Webb SR, Wang Z *et al.* Study of the degradation of a multidrug transporter using a non-radioactive pulse chase method. *Anal Bioanal Chem.* 2016;**408**:7745–51.
- Charoenpanich P, Meyer S, Becker A *et al.* Temporal expression program of quorum sensing-based transcription regulation in *Sinorhizobium meliloti*. *J Bacteriol* 2013;**195**:3224–36.
- Cianciulli Sesso A, Lilić B, Amman F *et al.* Gene Expression Profiling of *Pseudomonas aeruginosa* Upon Exposure to Colistin and Tobramycin. *Front Microbiol* 2021;**12**:626715.
- Clauwaert J, Menschaert G, Waegeman W. DeepRibo: a neural network for precise gene annotation of prokaryotes by combining ribosome profiling signal and binding site patterns. *Nucleic Acids Res* 2019;**47**:e36.
- Čuklina J, Hahn J, Imakaev M *et al.* Genome-wide transcription start site mapping of *Bradyrhizobium japonicum* grown free-living or in symbiosis - a rich resource to identify new transcripts, proteins and to study gene regulation. *BMC Genomics* 2016;**17**:302.
- Datta AK, Burma DP. Association of ribonuclease I with ribosomes and their subunits. *J Biol Chem.* 1972;**247**:6795–801.
- Djordjevic MA. *Sinorhizobium meliloti* metabolism in the root nodule: a proteomic perspective. *Proteomics* 2004;**4**:1859–72.
- Dodbele S, Wilusz JE. Ending on a high note: downstream ORFs enhance mRNA translational output. *EMBO J* 2020;**39**:e105959.
- Duval M, Cossart P. Small bacterial and phagic proteins: an updated view on a rapidly moving field. *Curr Opin Microbiol* 2017;**39**:81–88.
- Evguenieva-Hackenberg E. Riboregulation in bacteria: from general principles to novel mechanisms of the *trp* attenuator and its sRNA and peptide products. *Wiley Interdiscip Rev RNA* 2022;**13**:e1696.
- Ewels P, Magnusson M, Lundin S *et al.* MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* 2016;**32**:3047–8.
- Fancello L, Burger T. An analysis of proteogenomics and how and when transcriptome-informed reduction of protein databases can enhance eukaryotic proteomics. *Genome Biol* 2022;**23**:132.
- Ferrières L, Francez-Charlot A, Gouzy J *et al.* FixJ-regulated genes evolved through promoter duplication in *Sinorhizobium meliloti*. *Microbiology* 2004;**150**:2335–45.
- Fijalkowski I, Willems P, Jonckheere V *et al.* Hidden in plain sight: challenges in proteomics detection of small ORF-encoded polypeptides. *MicroLife* 2022;**3**:uqac005, <https://doi.org/10.1093/misml/uqac005>.
- Fontaine F, Fuchs RT, Storz G. Membrane localization of small proteins in *Escherichia coli*. *J Biol Chem.* 2011;**286**:32464–74.
- Galibert F, Finan TM, Long SR *et al.* The composite genome of the legume symbiont *Sinorhizobium meliloti*. *Science* 2001;**293**:668–72.
- Gelhausen R, Müller T, Svensson SL *et al.* RiboReport - benchmarking tools for ribosome profiling-based identification of open reading frames in bacteria. *Brief Bioinformatics* 2022;**23**:bbab549, doi:10.1093/bib/bbab549.
- Gelhausen R, Svensson SL, Froschauer K *et al.* HRIBO: high-throughput analysis of bacterial ribosome profiling data. *Bioinformatics* 2021;**37**:2061–3.
- Gelsinger DR, Dallon E, Reddy R *et al.* Ribosome profiling in archaea reveals leaderless translation, novel translational initiation sites, and ribosome pausing at single codon resolution. *Nucleic Acids Res* 2020;**48**:5201–16.
- Gerashchenko MV, Gladyshev VN. Translation inhibitors cause abnormalities in ribosome profiling experiments. *Nucleic Acids Res* 2014;**42**:e134.
- Gertz EM, Yu Y-K, Agarwala R *et al.* Composition-based statistics and translated nucleotide searches: improving the TBLASTN module of BLAST. *BMC Biol* 2006;**4**:41.
- Glaub A, Huptas C, Neuhaus K *et al.* Recommendations for bacterial ribosome profiling experiments based on bioinformatic evaluation of published data. *J Biol Chem* 2020;**295**:8999–9011.
- Gray T, Storz G, Papenfort K. Small proteins; big questions. *J Bacteriol* 2022;**204**:e0034121.
- Grüning B, Dale R, Sjödin A *et al.* Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nat Methods* 2018;**15**:475–6.
- Grützner J, Billenkamp F, Spanka D-T *et al.* The small DUF1127 protein CcaF1 from *Rhodobacter sphaeroides* is an RNA-binding protein involved in sRNA maturation and RNA turnover. *Nucleic Acids Res* 2021;**49**:3003–19.
- Guan Y, Zhu Q, Huang D *et al.* An equation to estimate the difference between theoretically predicted and SDS PAGE-displayed molecular weights for an acidic peptide. *Sci Rep* 2015;**5**:13370.
- Hadjeras L, Bartel J, Maier LK *et al.* Revealing the small proteome of *Haloferax volcanii* by combining ribosome profiling and small-protein optimized mass spectrometry. *MicroLife* 2023;**4**:uqad001.
- Hahn J, Tsouy OV, Thalmann S *et al.* Small Open Reading Frames, Non-Coding RNAs and Repetitive Elements in *Bradyrhizobium japonicum* USDA 110. *PLoS One* 2016;**11**:e0165429.
- Hemm MR, Paul BJ, Miranda-Ríos J *et al.* Small stress response proteins in *Escherichia coli*: proteins missed by classical proteomic studies. *J Bacteriol* 2010;**192**:46–58.
- Hemm MR, Weaver J, Storz G. *Escherichia coli* Small Proteome. *Ecosal Plus* 2020;**9**. doi: 10.1128/ecosalplus.ESP-0031-2019.
- Hyatt D, Chen GL, Locascio PF *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinf* 2010;**11**:119.
- Ingolia NT, Brar GA, Rouskin S *et al.* The ribosome profiling strategy for monitoring translation in vivo by deep sequencing of ribosome-protected mRNA fragments. *Nat Protoc* 2012;**7**:1534–50.
- Ingolia NT, Ghaemmaghani S, Newman JRS *et al.* Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science* 2009;**324**:218–23.
- Ingolia NT, Hussmann JA, Weissman JS. Ribosome profiling: global views of translation. *Cold Spring Harb Perspect Biol* 2019;**11**:a032698.
- Ingolia NT. Ribosome Footprint Profiling of Translation throughout the Genome. *Cell* 2016;**165**:22–33.
- Jones KM, Kobayashi H, Davies BW *et al.* How rhizobial symbionts invade plants: the *Sinorhizobium-Medicago* model. *Nat Rev Microbiol* 2007;**5**:619–33.
- Karzai AW, Roche ED, Sauer RT. The SsrA-SmpB system for protein tagging, directed degradation and ribosome rescue. *Nat Struct Biol* 2000;**7**:449–55.
- Kaulich PT, Cassidy L, Bartel J *et al.* Multi-protease Approach for the Improved Identification and Molecular Characterization of Small

- Proteins and Short Open Reading Frame-Encoded Peptides. *J Proteome Res* 2021;**20**:2895–903.
- Keller KC, Shapiro L, Williams KP. tmRNAs that encode proteolysis-inducing tags are found in all known bacterial genomes: a two-piece tmRNA functions in *Caulobacter*. *Proc Natl Acad Sci USA* 2000;**97**:7778–83.
- Kelley LA, Mezulis S, Yates CM et al. The Phyre2 web portal for protein modeling, prediction and analysis. *Nat Protoc* 2015;**10**:845–58.
- Khan SR, Gaines J, Roop RM et al. Broad-host-range expression vectors with tightly regulated promoters and their use to examine the influence of TraR and TraM expression on Ti plasmid quorum sensing. *Appl Environ Microbiol* 2008;**74**:5053–62.
- Khitun A, Slavoff SA. Proteomic detection and validation of translated small open reading frames. *Curr Protoc Chem Biol* 2019;**11**:e77.
- Knoke LR, Abad Herrera S, Götz K et al. *Agrobacterium tumefaciens* Small Lipoprotein Atu8019 Is Involved in Selective Outer Membrane Vesicle (OMV) Docking to Bacterial Cells. *Front Microbiol* 2020;**11**:1228.
- Köster J, Rahmann S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* 2012;**28**:2520–2.
- Kovacs-Simon A, Titball RW, Michell SL. Lipoproteins of bacterial pathogens. *Infect Immun* 2011;**79**:548–61.
- Kraus A, Weskamp M, Zierles J et al. Arginine-Rich Small Proteins with a Domain of Unknown Function, DUF1127, Play a Role in Phosphate and Carbon Metabolism of *Agrobacterium tumefaciens*. *J Bacteriol* 2020;**202**:e00309–20.
- Krogh A, Larsson B, von Heijne G et al. Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J Mol Biol* 2001;**305**:567–80.
- Liao Y, Smyth GK, Shi W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* 2014;**30**:923–30.
- Li H, Handsaker B, Wysoker A et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 2009;**25**:2078–9.
- Li H, Joh YS, Kim H et al. Evaluating the effect of database inflation in proteogenomic search on sensitive and reliable peptide identification. *BMC Genomics* 2016;**17**:1031.
- Marlow VL, Haag AF, Kobayashi H et al. Essential role for the BacA protein in the uptake of a truncated eukaryotic peptide in *Sinorhizobium meliloti*. *J Bacteriol* 2009;**191**:1519–27.
- Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet j* 2011;**17**:10.
- Marx H, Minogue CE, Jayaraman D et al. A proteomic atlas of the legume *Medicago truncatula* and its nitrogen-fixing endosymbiont *Sinorhizobium meliloti*. *Nat Biotechnol* 2016;**34**:1198–205.
- McIntosh M, Krol E, Becker A. Competitive and cooperative effects in quorum-sensing-regulated galactoglucan biosynthesis in *Sinorhizobium meliloti*. *J Bacteriol* 2008;**190**:5308–17.
- Melior H, Li S, Madhugiri R et al. Transcription attenuation-derived small RNA mTrpL regulates tryptophan biosynthesis gene expression in trans. *Nucleic Acids Res* 2019;**47**:6396–410.
- Melior H, Li S, Stötzel M et al. Reprograming of sRNA target specificity by the leader peptide peTrpL in response to antibiotic exposure. *Nucleic Acids Res* 2021;**49**:2894–915.
- Melior H, Maaß S, Li S et al. The Leader Peptide peTrpL Forms Antibiotic-Containing Ribonucleoprotein Complexes for Posttranscriptional Regulation of Multiresistance Genes. *MBio* 2020;**11**:e01027–20.
- Meydan S, Marks J, Klepacki D et al. Retapamulin-Assisted Ribosome Profiling Reveals the Alternative Bacterial Proteome. *Mol Cell* 2019;**74**:481–493.e6.
- Mishra A, Siwach P, Singhal P et al. ChemGenome2.1: an Ab Initio Gene Prediction Software. *Methods Mol Biol* 2019;**1962**:121–38.
- Mohammad F, Green R, Buskirk AR. A systematically-revised ribosome profiling method for bacteria reveals pauses at single-codon resolution. *Elife* 2019;**8**:e42591.
- Ndah E, Jonckheere V, Giess A et al. REPARATION: ribosome profiling assisted (re-)annotation of bacterial genomes. *Nucleic Acids Res* 2017;**45**:e168.
- Nesvizhskii AI. A survey of computational methods and error rate estimation procedures for peptide and protein identification in shotgun proteomics. *J Proteomics* 2010;**73**:2092–123.
- Oh E, Becker AH, Sandikci A et al. Selective ribosome profiling reveals the cotranslational chaperone action of trigger factor in vivo. *Cell* 2011;**147**:1295–308.
- Omasits U, Quebatte M, Stekhoven DJ et al. Directed shotgun proteomics guided by saturated RNA-seq identifies a complete expressed prokaryotic proteome. *Genome Res* 2013;**23**:1916–27.
- Omasits U, Varadarajan AR, Schmid M et al. An integrative strategy to identify the entire protein coding potential of prokaryotic genomes by proteogenomics. *Genome Res* 2017;**27**:2083–95.
- Orr MW, Mao Y, Storz G et al. Alternative ORFs and small ORFs: shedding light on the dark proteome. *Nucleic Acids Res* 2020;**48**:1029–42.
- Otto C, Stadler PF, Hoffmann S. Lacking alignments? The next-generation sequencing mapper segemehl revisited. *Bioinformatics* 2014;**30**:1837–43.
- Patraquim P, Mumtaz MAS, Pueyo JI et al. Developmental regulation of canonical and small ORF translation from mRNAs. *Genome Biol* 2020;**21**:128.
- Petruschke H, Anders J, Stadler PF et al. Enrichment and identification of small proteins in a simplified human gut microbiome. *J Proteomics* 2020;**213**:103604.
- Petruschke H, Schori C, Canzler S et al. Discovery of novel community-relevant small proteins in a simplified human intestinal microbiome. *Microbiome* 2021;**9**:55.
- Qeli E, Ahrens CH. PeptideClassifier for protein inference and targeted quantitative proteomics. *Nat Biotechnol* 2010;**28**:647–50.
- Sallet E, Roux B, Sauviac L et al. Next-generation annotation of prokaryotic genomes with EuGene-P: application to *Sinorhizobium meliloti* 2011. *DNA Res* 2013;**20**:339–54.
- Schägger H. Tricine-SDS-PAGE. *Nat Protoc* 2006;**1**:16–22.
- Scheuer R, Dietz T, Kretz J et al. Incoherent dual regulation by a SAM-II riboswitch controlling translation at a distance. *RNA Biol* 2022;**19**:980–95.
- Schlüter J-P, Reinkensmeier J, Barnett MJ et al. Global mapping of transcription start sites and promoter motifs in the symbiotic  $\alpha$ -proteobacterium *Sinorhizobium meliloti* 1021. *BMC Genomics* 2013;**14**:156.
- Sharma CM, Darfeuille F, Plantinga TH et al. A small RNA regulates multiple ABC transporter mRNAs by targeting C/A-rich elements inside and upstream of ribosome-binding sites. *Genes Dev* 2007;**21**:2804–17.
- Simon R, Priefer U, Pühler A. A broad host range mobilization system for *in vivo* genetic engineering: transposon mutagenesis in gram negative bacteria. *Nat Biotechnol* 1983;**1**:784–91.
- Sobrero P, Schlüter J-P, Lanner U et al. Quantitative proteomic analysis of the Hfq-regulon in *Sinorhizobium meliloti* 2011. *PLoS One* 2012;**7**:e48494.
- Song K, Baumgartner D, Hagemann M et al. Atp $\theta$  is an inhibitor of FOF1 ATP synthase to arrest ATP hydrolysis during low-energy conditions in cyanobacteria. *Curr Biol* 2022;**32**:136–48.e5.

- Stekhoven DJ, Omasits U, Quebatte M et al. Proteome-wide identification of predominant subcellular protein localizations in a bacterial model organism. *J Proteomics* 2014;**99**:123–37.
- Storz G, Wolf YI, Ramamurthi KS. Small proteins can no longer be ignored. *Annu Rev Biochem* 2014;**83**:753–77.
- Sun Y-H, de Jong MF, den Hartigh AB et al. The small protein CydX is required for function of cytochrome bd oxidase in *Brucella abortus*. *Front Cell Infect Microbiol* 2012;**2**:47.
- Taboada B, Estrada K, Ciria R et al. Operon-mapper: a web server for precise operon identification in bacterial and archaeal genomes. *Bioinformatics* 2018;**34**:4118–20.
- Torres-Quesada O, Millán V, Nisa-Martínez R et al. Independent activity of the homologous small regulatory RNAs AbcR1 and AbcR2 in the legume symbiont *Sinorhizobium meliloti*. *PLoS One* 2013;**8**:e68147.
- Tyanova S, Temu T, Cox J. The MaxQuant computational platform for mass spectrometry-based shotgun proteomics. *Nat Protoc* 2016;**11**:2301–19.
- Ulvé VM, Chéron A, Trautwetter A et al. Characterization and expression patterns of *Sinorhizobium meliloti* tmRNA (ssrA). *FEMS Microbiol Lett* 2007;**269**:117–23.
- Vallenet D, Belda E, Calteau A et al. MicroScope—an integrated microbial resource for the curation and comparative analysis of genomic and metabolic data. *Nucleic Acids Res* 2013;**41**:D636–647.
- VanOrsdel CE, Kelly JP, Burke BN et al. Identifying New Small Proteins in *Escherichia coli*. *Proteomics* 2018;**18**:e1700064.
- Varadarajan AR, Allan RN, Valentin JDP et al. An integrated model system to gain mechanistic insights into biofilm-associated antimicrobial resistance in *Pseudomonas aeruginosa* MPAO1. *Npj Biofilms and Microbiomes* 2020a;**6**:46.
- Varadarajan AR, Goetze S, Pavlou MP et al. A Proteogenomic Resource Enabling Integrated Analysis of *Listeria* Genotype-Proteotype-Phenotype Relationships. *J Proteome Res* 2020b;**19**:1647–62.
- Vazquez-Laslop N, Sharma CM, Mankin A et al. Identifying Small Open Reading Frames in Prokaryotes with Ribosome Profiling. *J Bacteriol* 2022;**204**:e0029421.
- Venturini E, Svensson SL, Maaß S et al. A global data-driven census of *Salmonella* small proteins and their potential functions in bacterial virulence. *MicroLife* 2020 doi: 10.1093/femsm/luqaa002.
- Vitreschak AG, Lyubetskaya EV, Shirshin MA et al. Attenuation regulation of amino acid biosynthetic operons in proteobacteria: comparative genomics analysis. *FEMS Microbiol Lett* 2004;**234**:357–70.
- Weaver J, Mohammad F, Buskirk AR et al. Identifying Small Proteins by Ribosome Profiling with Stalled Initiation Complexes. *MBio* 2019;**10**:e02819–18.
- Wingett SW, Andrews S. FastQ Screen: a tool for multi-genome mapping and quality control. [version 2; peer review: 4 approved]. *F1000Res* 2018;**7**:1338.
- Wiśniewski JR Quantitative evaluation of filter aided sample preparation (FASP) and multienzyme digestion FASP protocols. *Anal Chem* 2016;**88**:5438–43.
- Wu Q, Wright M, Gogol MM et al. Translation of small downstream ORFs enhances translation of canonical main open reading frames. *EMBO J* 2020;**39**:e104763.
- Yu NY, Wagner JR, Laird MR et al. PSORTb 3.0: improved protein subcellular localization prediction with refined localization subcategories and predictive capabilities for all prokaryotes. *Bioinformatics* 2010;**26**:1608–15.
- Zeghouf M, Li J, Butland G et al. Sequential Peptide Affinity (SPA) system for the identification of mammalian and bacterial protein complexes. *J Proteome Res* 2004;**3**:463–8.
- Zevenhuizen LPTM, van Neerven ARW. (1→2)- $\beta$ -d-glucan and acidic oligosaccharides produced by *Rhizobium meliloti*. *Carbohydr Res* 1983;**118**:127–34.
- Zhang Y, Fonslow BR, Shan B et al. Protein analysis by shotgun/bottom-up proteomics. *Chem Rev* 2013;**113**:2343–94.