

# **XPRESSyourself: Enhancing, Standardizing, and Automating Ribosome Profiling Computational Analyses Yields Improved Insight into Data**

**Jordan A. Berg,<sup>1\*</sup> Jonathan R. Belyeu,<sup>2</sup> Jeffrey T. Morgan,<sup>1</sup> Yeyun Ouyang,<sup>1</sup> Alex J. Bott,<sup>1</sup> Aaron R. Quinlan,<sup>2,4,5</sup> Jason Gertz,<sup>3</sup> Jared Rutter<sup>1,6\*</sup>**

<sup>1</sup>Department of Biochemistry, University of Utah, Salt Lake City, UT, USA, 84112.

<sup>2</sup>Department of Human Genetics, University of Utah, Salt Lake City, UT, USA, 84112.

<sup>3</sup>Department of Oncological Sciences, University of Utah, Salt Lake City, UT, USA, 84112.

<sup>4</sup>USTAR Center for Genetic Discovery, University of Utah, Salt Lake City, UT, USA, 84112.

<sup>5</sup>Department of Biomedical Informatics, University of Utah, Salt Lake City, UT, USA, 84112.

<sup>6</sup>Howard Hughes Medical Institute, University of Utah, Salt Lake City, UT, USA, 84112.

\*Address correspondence to: [jordan.berg@biochem.utah.edu](mailto:jordan.berg@biochem.utah.edu), [rutter@biochem.utah.edu](mailto:rutter@biochem.utah.edu)

## Abstract

Ribosome profiling, an application of nucleic acid sequencing for monitoring ribosome activity, has revolutionized our understanding of protein translation dynamics. This technique has been available for a decade, yet the current state of publicly available computational tools for these data is bleak. We introduce XPRESSyourself, an analytical toolkit that eliminates barriers and bottlenecks associated with this specialized data type by filling gaps in the computational toolset for both experts and non-experts of ribosome profiling. XPRESSyourself automates and standardizes analysis procedures, decreasing time-to-discovery and increasing reproducibility. This toolkit acts as a reference implementation of current best practices in ribosome profiling analysis. We demonstrate this toolkit's performance on publicly available ribosome profiling data and associated bulk RNA-Seq data by rapidly identifying hypothetical mechanisms related to neurodegenerative phenotypes and neuroprotective mechanisms of the small-molecule ISRIB during acute cellular stress. XPRESSyourself brings robust, rapid analysis of ribosome-profiling data to a broad and ever-expanding audience and will lead to more reproducible and accessible measurements of translation regulation. The XPRESSyourself software is perpetually open-source under the GPL-3.0 license and is hosted at <https://github.com/XPRESSyourself>, where the user can find additional documentation and report issues.

## Introduction

High-throughput sequencing data has revolutionized biomedical and biological research. One such application of this consequential technology is ribosome profiling, which, coupled with bulk RNA-Seq, measures translation efficiency, translation pausing, novel protein translation products, and more. [1–3]. Though ribosome profiling has matured, there remains an abundance of biases and peculiarities associated with each analytical method or tool, which are often obscured to a user [4–8]. Additionally, standardized methods for handling this unique data type remain elusive. This has been problematic and evidenced by various studies using vague or opaque methods for data analysis (for several examples, see [9–13]), or methods rely on outdated tools, such as Bowtie [14], which continues to be used in various ribosome profiling studies for alignment. In this example, this tool is inappropriate as it is unable to map reads across splice junctions [5]. Very few labs have the tools to separate the biological signals in ribosome profiling data from inherent biases of the experimental measurements, and these tools are not readily accessible by the community. This is a critical time in the rapidly expanding influence of ribosome-footprint profiling. For too long, the bioinformatic know-how of this incredibly powerful technique has been limited to a small handful of labs. As more

and more ribosome profiling studies are performed, more and more labs will lack the ability to analyze their data with ease and fidelity. Few if any extant pipelines or toolkits offer a thorough set of integrated tools for assessing standard quality control metrics or performing proper reference curation and reducing systematic biases across any organism, particularly with ribosome profiling data [15–19].

For example, one issue in ribosome profiling is the pile-up of ribosomes at the 5'- and 3'- ends of coding regions within a transcript is a systematic biological signal arising from the slower kinetics of ribosome initiation and termination compared to translation elongation and is often not relevant to measurements of translation regulation [4, 20, 21]. These pile-ups can dramatically skew ribosome footprint quantification and measurements of translational efficiency. Leaders in the field currently recommend excluding these pile-up-prone regions when quantifying ribosome profiling alignments [3, 22]; however, no publicly available computational tools currently exist to facilitate these automated adjustments to reference transcripts. Curating references for this ability requires advanced implementations to achieve correctly and robustly. In addition, downstream data visualization methods presently available are not optimized to analyze and compare translationally regulation regions of a gene.

To address these deficiencies in the public ribosome profiling computational toolkit, we developed XPRESSyourself, a computational toolkit and server-less, adaptable pipeline that bridges these and other gaps in ribosome profiling data analysis. XPRESSyourself implements the complete suite of tools necessary for ribosome profiling and bulk RNA-Seq analysis in a robust and easy-to-use software package, often packaging tasks that would typically require hundreds to thousands of lines of code into a single command. This software will lead to more reproducible and more accessible measurements of translation regulation. For instance, XPRESSyourself creates the mRNA annotation files necessary to remove confounding systematic factors during quantification and analysis of ribosome profiling data to measure translation accurately. It also provides a built-in capacity to quantify and visualize differential upstream open-reading frame (uORF) usage by generating IGV-like, intron-less plots for easier visualization [23]. The ability to visualize (and in another XPRESSyourself module, quantify) the usage of micro-uORFs is important in exploring regulatory events or mechanisms in a wide array of biological responses and diseases. XPRESSyourself also introduces a tool for efficient identification of the most problematic rRNA fragments for targeted depletion, which provides immense financial and experimental benefits to the user as ribosome footprint signal can be better amplified over rRNA noise. Tools like this will become vital as ribosome profiling moves into development in new organisms.

XPRESSyourself aims to address the lack of consensus in analytical approaches used to process ribosome

profiling data by acting as a reference implementation of current best practices for ribosome profiling analysis. While a basic bioinformatic understanding is becoming more commonplace amongst the scientific community, the intricacies of processing RNA-Seq data remain challenging for many. Moreover, many users are often not aware of the most up-to-date tools or the appropriate settings for their application [24, 25]. Even for the experienced user, developing robust automated pipelines that accurately process and assess the quality of these datasets can be laborious. The variability that inevitably arises with each lab or core facility designing and using distinct pipelines is also a challenge to reproducibility in the field. XPRESSyourself curates the state-of-the-art methods for use and where a required functionality is unavailable, introduces a thoroughly tested module to fill that gap. In addition to the new tools described above, the toolkit provides the user with a complete suite of software to handle pre-processing, aligning, and quantifying of sequence reads, performing quality control via various meta-analyses of pre- and post-processed reads, in many cases tailored for ribosome profiling but flexible enough to handle general RNA-Seq datasets. While some of these steps may be considered more mundane, we eliminate the need of each user to rewrite even simple functionality and promote reproducibility between implementations. To aid users of any skill-level in using this toolkit, we provide thorough documentation, walkthrough videos, and interactive command builders to make usage as easy as possible, while allowing for broad use of this toolkit on high-performance clusters.

Finally, the most broadly relevant aspect of our update and streamlining of ribosome-profiling analysis is the novel biological insights we are able to obtain from published datasets. We highlight this in the ISRIB ribosome-profiling study discussed in this manuscript, where we are able to observe significant translation regulation that was missed previously when the data were initially analyzed using now outdated techniques. This analysis generates novel hypotheses for genes potentially involved in neurodegeneration in humans, but more broadly emphasizes the benefit of analysis and re-analysis of data using the complete and up-to-date benchmarked methodology provided within XPRESSyourself.

## **Design and Implementation**

### **Architecture and Organization**

XPRESSyourself is currently partitioned into two software packages, XPRESSpipe and XPRESSplot. XPRESSpipe contains automated pipelines tailored for ribosome profiling, single-end RNA-Seq, and paired-end RNA-Seq datasets. Figure 1 outlines tasks completed within the pipelines. Individual sub-modules can be run automatically through a

pipeline or manually step-by-step. Modules optimize available computational resources to deliver results as quickly as possible. XPRESSplot is available as a Python library and provides an array of analytical methods specifically for sequence data, but tractable to other data types. For a comparison of how XPRESSyourself compares to other available software packages available at time of writing, we refer the reader to Figure S1 [15, 16, 19, 26–50].

XPRESSyourself aims to make ribosome profiling and sequence analysis as easy and accessible as possible to all users. As such, an integrated command builder for reference curation and sample analysis can be run by executing `xpresspipe build`. This command builder will walk the user through potential considerations based on their library preparation method and build the appropriate command for execution on their personal computer or a supercomputing cluster.

The software is designed such that updating and testing of a new module, or updating dependency usage, are facile tasks for a trained bioinformatician. More details on current and future capabilities can be found in each package's documentation [51, 52] or their respective `versions` page on their respective repository pages [53].

## **Automated Reference Curation**

The first step of RNA-Seq alignment is curating an organism reference to which the alignment software will map sequence reads. XPRESSpipe uses STAR [54] for mapping reads as it has been shown consistently to be the best performing read aligner for RNA-Seq data [55, 56]. The appropriate reference files are automatically curated by providing the appropriate GTF file saved as `transcripts.gtf` and the directory path to the genomic FASTA file(s). Additional modifications to the GTF file required for ribosome profiling or desired for RNA-Seq are discussed in the next section.

## **GTF Modification**

For ribosome profiling, frequent read pile-ups are observed at the 5'- and 3'- ends of an open reading frame which are largely uninformative as to the translational efficiency of the gene [4]. While these pile-ups can be indicative of true translation dynamics [57], more recently leaders in the field have determined that these regions should be ignored when quantifying reads and calculating translation efficiency [3, 22]. By providing the `--truncate` argument during reference curation, the 5'- and 3'- ends of each coding region will be recursively trimmed until the specified amounts are removed from coding space. A recursive strategy is required here as GTF file formats section out

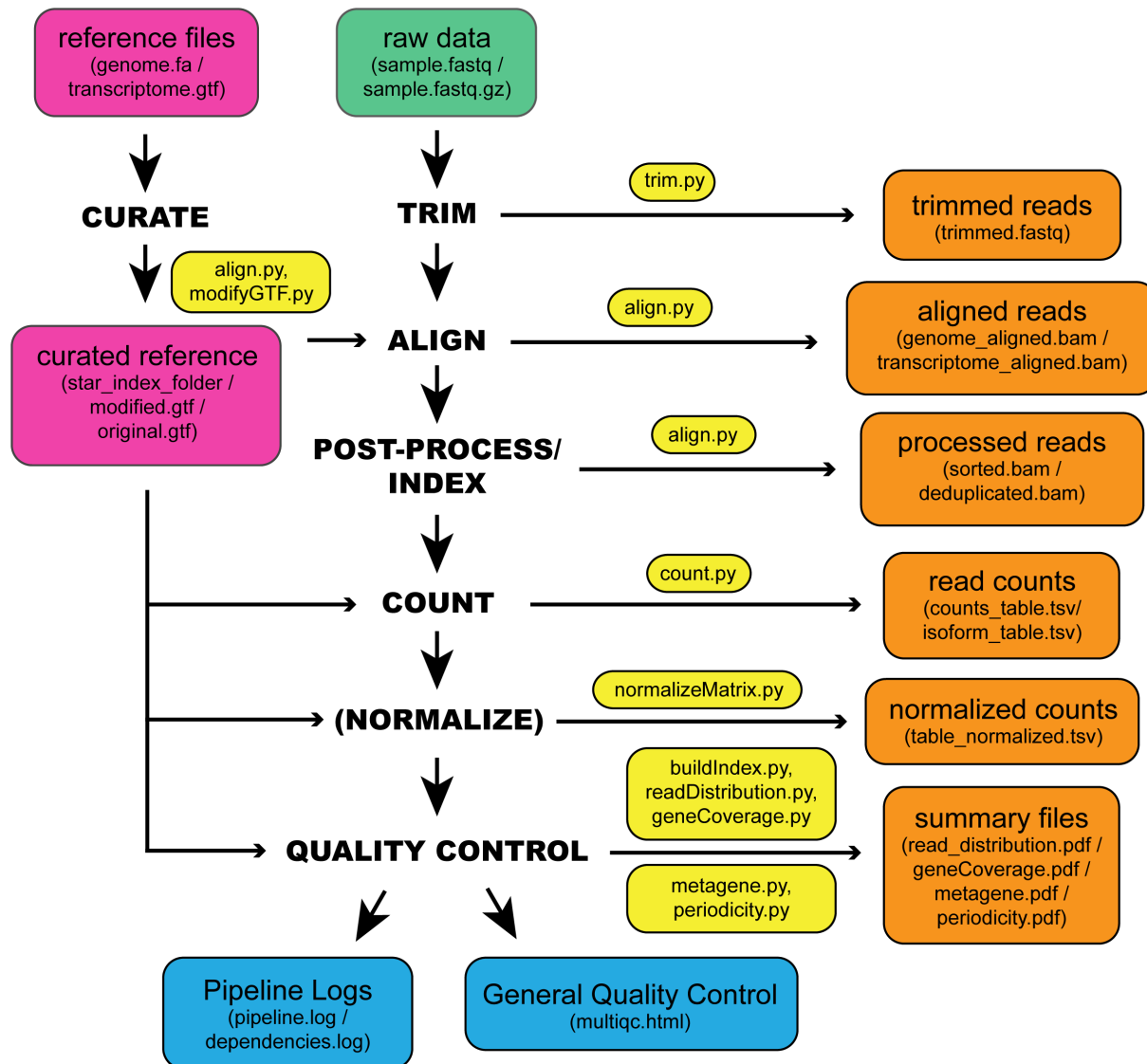


Figure 1: **Workflow schematic of the inputs, outputs, and organization of XPRESSpipe.** Representation of the general steps performed by XPRESSpipe with data and log outputs. Steps in parentheses are optional to the user. Input and output file types are in parentheses for each input or output block. Main script(s) used for a given step are in yellow blocks. The green block indicates input sequence file(s). Pink blocks indicate reference input files and curated reference. Orange blocks indicate output files. Blue blocks indicate general quality control and log file outputs.

each CDS region on each exon as separate records. By default, 45 nt will be trimmed from the 5'-ends and 15 nt from the 3'-ends, as is the convention within the ribosome profiling field [3], but these amounts can be modified. The resulting output file should be used to process ribosome footprint libraries, along with their corresponding bulk RNA-Seq libraries. If generating a GTF file for use solely with general bulk RNA-Seq datasets, this file should not be truncated.

Additionally, parameters that can be provided during this step to retain only protein-coding genes records. As ribosomal RNAs and other non-coding RNAs can be highly abundant in RNA-Seq experiments, it is often recommended to not include these sequences for quantification. This acts as a read masking step to exclude non-protein coding transcripts from downstream analyses. Parameters can also be provided to retain only the Ensembl canonical transcript record. This can be useful for some tools that penalize reads that overlap multiple isoforms of the same gene. If using HTSeq with default XPRESSpipe parameters or Cufflinks to quantify reads, this is not necessary as they do not penalize a read mapping to multiple isoforms of the same gene or are capable of handling quantification of different isoforms of a gene [58,59].

## **Read Processing**

**Pre-Processing.** In order for sequence reads to be mapped to the genome, reads generally need to be cleaned of artifacts from library creation. These include adaptors, unique molecular identifier (UMI) sequences, and technical errors in the form of low-quality base calls. Parameters, like minimum acceptable quality or length, can be modified, or features such as UMIs can be specified to identify and group PCR artifacts for later removal [60,61].

**Mapping.** Reads are aligned to the reference genome with STAR, which, despite being more memory-intensive, is one of the fastest and most accurate sequence alignment options currently available [54–56]. XPRESSpipe is capable of performing a single-pass, splice-aware, GTF-guided alignment or a two-pass alignment of reads wherein novel splice junctions are determined and built into the reference, followed by alignment of reads to the new reference. A coordinate-sorted and indexed BAM file is output by STAR. We abstain from rRNA negative alignment at this step as downstream analysis of these mapped reads could be of interest to some users. When rRNA alignment is preferred, a protein-coding-only GTF file should be provided during quantification.

**Quantification.** XPRESSpipe further processes alignment files by optionally parsing for unique alignments for downstream analyses. PCR duplicates are also detected and marked or removed for downstream analyses; however,

these files are only used for relevant downstream steps or if the user specifies to use these de-duplicated files in all downstream steps. Use of de-duplicated alignment files may be advisable in situations where the library complexity profiles (discussed below) exhibit high duplication frequencies. However, generally the abundance of PCR-duplicates is low in properly-prepared sequencing libraries; thus, doing so may be overly stringent and unnecessary [60]. Optionally, BED coverage files can also be output.

**Post-Processing.** XPRESSpipe quantifies read alignments for each input file using HTSeq with the `intersection-nonempty` method by default [58,62]. We use this quantification method as it conforms to the current TCGA processing standards and is favorable on the majority of applications [63]. If masking of non-coding RNAs is desired, a `protein_coding` modified GTF file should be provided for the `--gtf` argument. HTSeq is recommended for processing ribosome profiling data as it allows selection of feature type across which to quantify, thus allowing for quantification across the CDSs of a transcript instead of entire exons. If a user is interested in quantifying ribosome occupancy of transcript uORFs in ribosome footprint samples, they can provide `five_prime_utr` or `three_prime_utr` for the `--feature_type` parameter if such annotations exist in the organism of interest's GTF file. If the user is interested in isoform abundance estimation, Cufflinks is alternatively available for quantification [59,62]. If Cufflinks is provided for a ribosome profiling dataset, XPRESSpipe will use Ribomap to provide isoform abundance of ribosome profiles [45,46].

**Normalization.** Methods for count normalization are available within XPRESSpipe by way of the XPRESSplot package. For normalizations correcting for transcript length, the appropriate GTF must be provided. Sample normalization methods available include reads-per-million (RPM), Reads-per-kilobase-million (RPKM) or Fragments-per-kilobase-million (FPKM), and transcripts per million (TPM) normalization [64]. For samples sequenced on different flow cells, prepared by different individuals, or on different days, the `--batch` argument should be provided along with the appropriate metadata matrix [65].

## Quality Control

**Read Length Distribution.** The lengths of all reads are analyzed after trimming. By assessing the read distribution of each sample, the user can ensure the expected read size was sequenced. This is particularly helpful in ribosome profiling experiments for verifying the requisite 17-33 nt ribosome footprints were selectively captured during library preparation [3,66]. Metrics here, as in all other quality control sub-modules, are compiled into summary figures for



easy experiment-wide assessment by the user.

**Library Complexity.** Measuring library complexity is an effective method for analyzing the robustness of a sequencing experiment in capturing various, unique RNA species. As the majority of RNA-Seq preparation methods involve a PCR step, sometimes particular fragments will be favored and over-amplified in contrast to others. By plotting the number of PCR replicates versus expression level for each gene, one can monitor any effects of limited transcript capture diversity and/or high estimated PCR duplication rate on the robustness of their libraries. This analysis is performed using dupRadar [67] where inputs are PCR duplicate-tagged BAM files output during post-processing. Metrics are then compiled and plotted by XPRESSpipe.

**Metagene Estimation Profile.** To identify any general biases for the preferential capture of the 5'- or 3'- ends of transcripts, metagene profiles are generated for each sample. This is performed by determining the meta-genomic coordinate for each aligned read in exon space. Coverage is calculated for each transcript, normalized, and combined to eliminate greediness of super-expressors in profile coverage. Required inputs are an indexed BAM file and an un-modified GTF reference file. Outputs include metagene metrics, individual plots, and summary plots.

**Gene Coverage Profile.** Extending the metagene estimation analysis, the user can focus on the coverage profile across a single gene. Although traditional tools like IGV [23] offer the ability to perform such tasks, XPRESSpipe offers the ability to collapse the introns to observe coverage over exon space only. This is helpful in situations where massive introns spread out exons and make it difficult to visualize exon coverage for the entire transcript in a concise manner. CDS region is identified to aid ribosome profiling data users in identifying CDS coverage and uORF translation events. When running an XPRESSpipe pipeline, a housekeeping gene will be processed and output for the user's reference. Figure S2 provides a comparison with the output of IGV [23] and XPRESSpipe's `geneCoverage` module over a similar region for two genes to demonstrate the compatibility between the methods. We note that while the tool `superTranscripts` offers similar functionality, it lacks integration and automation and must be paired with IGV for visualization [32]. XPRESSpipe's `geneCoverage` module offers easy and automated functionality for this task.

**Codon Phasing/Periodicity Estimation Profile.** In ribosome profiling, a useful measure of a successful experiment is obtained by investigating the codon phasing of ribosome footprints [3]. To do so, the P-site positions relative to the start codon of each mapped ribosome footprint are calculated using `riboWaltz` [68]. The same inputs are required as in the `metagene` sub-module.

**Identify Problematic rRNA Fragments from Ribosome Footprinting for Depletion.** rRNA depletion is intrinsically complicated during the preparation of ribosome-footprint profiling libraries: poly(A) selection is irrelevant, and kit-based

rRNA depletion is grossly insufficient. Especially in the case of ribosome profiling experiments, where RNA is digested by an RNase to create ribosome footprints, many commercial depletion kits will not target the most abundant rRNA fragment species produces during the footprinting step of ribosome profiling. The sequencing of these RNAs becomes highly repetitive, wasteful, and typically biologically uninteresting in the context of gene expression and translation efficiency. The depletion of these sequences is therefore desired to increase the depth of coverage of ribosome footprints. Depending on the species and condition being profiled, custom rRNA-depletion probes for a small subset of rRNA fragments (generally 2-5) can easily account for more than 90% of sequenced reads [1, 3]. `rrnaProbe` analyzes the over-represented sequences within a collection of footprint sequence files that have already undergone adaptor and quality trimming, compiles conserved sequences across the overall experiment, and outputs a rank-ordered list of these sequences for probe design.

## **Analysis**

XPRESSpipe provides a DESeq2 command line wrapper for performing differential expression analysis of count data. We refer users to the original publication for more information about uses and methodology [69].

More analytical tools are available in XPRESSplot, which requires as input a gene count table as output by XPRESSpipe and a meta-sample table (explained in the documentation [52]). Analyses not currently available in other Python libraries include principle components plotting with confidence intervals and automated volcano plot creation for RNA-Seq or other data. Other instances of analyses can be found in the documentation [52].

## **Results**

### **Benchmarking Against Published Ribosome Profiling Data and New Insights**

The integrated stress response (ISR) is a signaling mechanism used by cells and organisms in response to a variety of cellular stresses [70]. Although acute ISR activation is essential for cells to properly respond to stresses, long periods of sustained ISR activity can be damaging. These prolonged episodes contribute to a variety of diseases, including many that result in neurological decline [71]. A recently discovered small-molecule inhibitor of the ISR, ISRIB, has been demonstrated to potentially be a safe and effective therapeutic for traumatic brain injury and other neurological diseases. Interestingly, ISRIB can suppress the damaging chronic low activation of the ISR, while it does not interfere with a cytoprotective acute, high-grade ISR. It has also been shown to be neuroprotective in mouse

models of traumatic brain injury, adding to its wide pharmacological interest [9, 72–77].

A recent study (data available under Gene Expression Omnibus accession number GSE65778) utilized ribosome profiling to better define the mechanisms of ISRIB action on the ISR, modeled by 1-hour tunicamycin (Tm) treatment in HEK293T cells [9]. A key finding of this study is that a specific subset of stress-related transcription factor mRNAs exhibit increased translational efficiency (TE) compared to untreated cells during the tunicamycin-induced ISR. However, when cells were co-treated with tunicamycin and ISRIB, the TE of these stress-related mRNAs showed no significant increase compared to untreated cells, which indicates that ISRIB can counteract the translational responses associated with the ISR.

To showcase the utility of XPRESSpipe in analyzing ribosome profiling and sequencing datasets, we re-processed and analyzed this dataset using the more current *in silico* techniques included in the XPRESSpipe package to further query the translational mechanisms of the ISR and ISRIB. All XPRESSpipe-processed biological replicate samples exhibited a strong correlation between read counts per gene when thresholded similarly to count data available with the original publication (Spearman  $\rho$  values 0.991-0.997) (Figure 2A shows representative plots; Figure S3A shows all replicate comparisons; Figure S4B shows the corresponding plots using the count data provided with the original publication for reference).

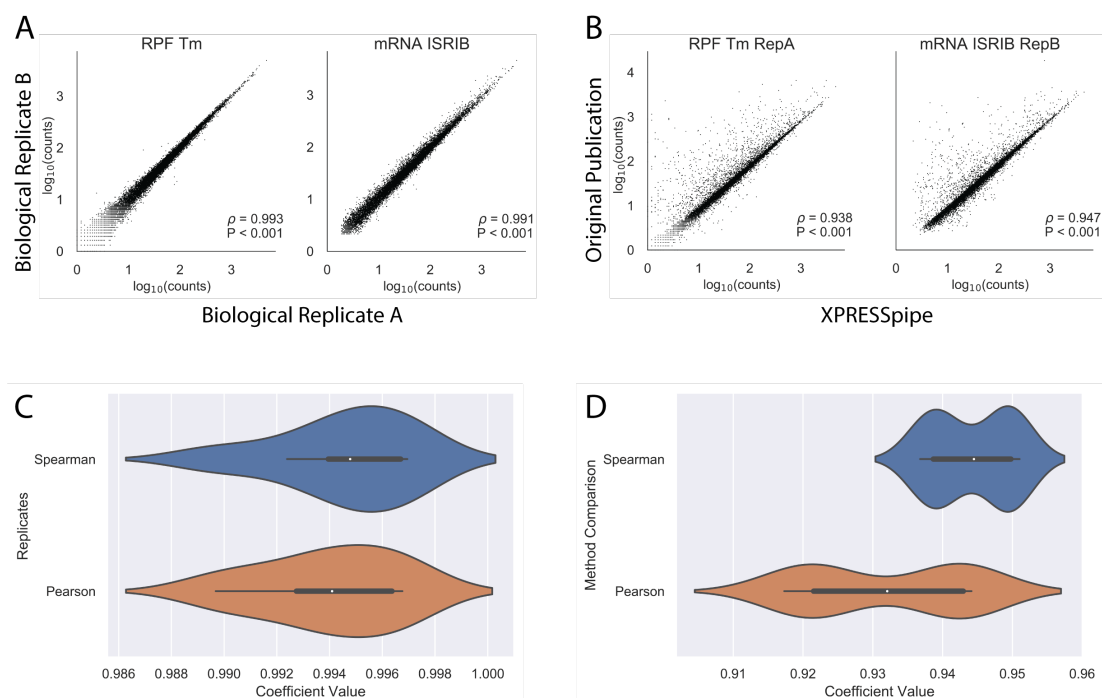
Compared to the count data made available with the original manuscript, when XPRESSpipe-processed samples were thresholded as in the original published count data, samples showed generally comparable read counts per gene between the two analytical regimes (Spearman  $\rho$  values 0.937-0.950) (Figure 2B shows representative plots; Figure S3B shows all comparisons). This is in spite of the fact that the methods section of the original publication employed software that was current at the time but is now outdated, such as TopHat2 [78], which has a documented higher false positive alignment rate, generally lower recall, and lower precision at correctly aligning multi-mapping reads compared to STAR [54–56]. Many of the genes over-represented in the original count data as compared to data processed by XPRESSpipe are genes that have pseudogenes or other paralogs (Figure S4A highlights a sampling of some extreme cases). As these genes share high sequence similarity with each other, reads mapping to these regions are difficult to attribute to a specific genomic locus and are often excluded from further analyses due to their multi-mapping nature. The benchmarking study [55] that examined these and other aligners described how TopHat2 had a disproportionately high rate of incorrectly aligned bases, or bases that were aligned uniquely when they should have been aligned ambiguously, at least partially explaining the observed overcounted effect with TopHat2. Had TopHat2 marked problematic reads as ambiguous, they would have been excluded from later quantification.

Additionally, when TopHat2 and STAR were tested using multi-mapper simulated test data of varying complexity, TopHat2 consistently suffered in precision and recall. These calls are increasingly more difficult to make with smaller reads as well, and this is evident from Figure 2B, where ribosome footprint samples consistently showed more over-counted genes than the corresponding RNA-Seq samples. When dealing with a ribosome footprint library of about 50-100 million reads, and with TopHat2's simulated likelihood of not marking an ambiguous read as such being about 0.5% higher than STAR, this would lead to around 250,000 to 500,000 spuriously aligned reads, which is in line with our observations (statistics were derived from [55]).

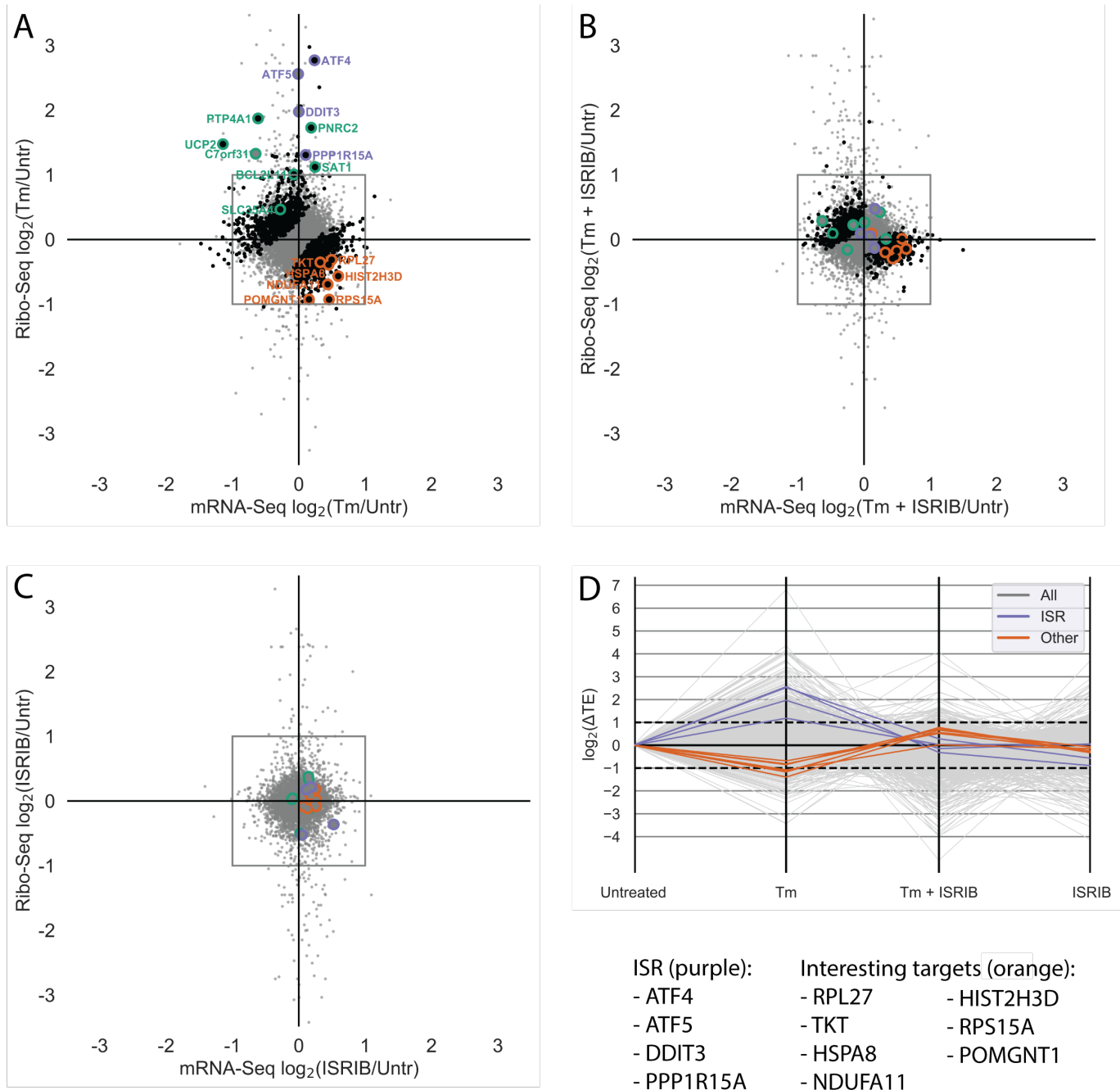
Another potential contributor to this divergence is that the alignment and quantification within XPRESSpipe use a current human transcriptome reference, which no doubt contains updates and modifications to annotated canonical transcripts and so forth when compared to the version used in the original study. However, in practice, these effects are modest for this dataset (Figure S5). Additionally, the usage of the now outdated DESeq1 [79] appears to contribute significantly to the outcome in differential expression analysis (Figure S7). While differences in processing between the outdated and current methods may not always create broad differences in output, key biological insights may be missed. The analysis that follows is exploratory and only meant to suggest putative targets identifiable by re-analyzing pre-existing, publicly available data.

Similar canonical targets of translation regulation during ISR were identified in the XPRESSpipe-processed data as were identified in the original study. These targets include ATF4, ATF5, PPP1R15A, and DDIT3 (Figure 3A-C, highlighted in purple) [9]. Of note, the fold-change in ribosome occupancy of ATF4 (6.83) from XPRESSpipe-processed samples closely mirrored the estimate from the original publication (6.44). Other targets highlighted in the original study [9], such as ATF5, PPP1R15A, and DDIT3 also demonstrated comparable increases in their ribosome occupancy fold-changes to the original publication count data (XPRESSpipe: 5.90, 2.47, and 3.94; respectively. Original: 7.50, 2.70, and 3.89; respectively) (Figure 3A). Similar to the originally processed data, all of these notable changes in ribosome occupancy return to untreated levels during Tm + ISRIB co-treatment (Figure 3B). Additional ISR targets containing micro-ORFs described in the study (highlighted in green in Figure 3A-C) were also similar in translational and transcriptional regulation across conditions between the two analyses.

Both the original study and our XPRESSpipe-based re-analysis show that ISRIB can counteract the significant increase in TE for a set of genes during ISR. To further explore TE regulation during ISR, we asked if ISRIB has a similar muting effect on genes with significant decreases in TE induced by the ISR. In the original study, genes with significant decreases in TE were reported in a source-data table but not a focus in the study. However, re-analysis



**Figure 2: Representative comparisons between processed data produced by XPRESSpipe and original study.** Genes were eliminated from analysis if any RNA-Seq sample for that gene had fewer than 10 counts. A) Representative comparisons of biological replicate read counts processed by XPRESSpipe. B) Representative comparisons of read counts per gene between count data from the original study and the same raw data processed and quantified by XPRESSpipe. C) Boxplot summaries of Spearman  $\rho$  and Pearson  $r$  values for biological replicate comparisons. D) Boxplot summaries of Spearman  $\rho$  and Pearson  $r$  values for between method processing. RPF, ribosome-protected fragments. Tm, tunicamycin. All  $\rho$  values reported in A and B are Spearman correlation coefficients using RPM-normalized count data. Pearson correlation coefficients were calculated using  $\log_{10}(\text{rpm}(\text{counts}) + 1)$  transformed data. XPRESSpipe-processed read alignments were quantified to *Homo sapiens* build CRCh38v98 using a protein-coding only, truncated GTF.



**Figure 3: Analysis of previously published ISR TE data using XPRESSpipe.** A-C)  $\log_2(\text{Fold Change})$  for each drug condition compared to untreated for the ribosome profiling and RNA-Seq data. Purple, ISR canonical targets highlighted in the original study. Green, genes with uORFs affected by ISR as highlighted in the original study. Orange, genes fitting a strict thresholding paradigm to identify genes that display a 2-fold or greater increase in TE in Tm + ISRIB treatment compared to Tm treatment. Black, genes with statistically significant changes in TE. Grey, all genes. Changes in ribo-seq and mRNA-Seq were calculated using DESeq2. TE was calculated using DESeq2. Points falling outside of the plotted range are not included. D) Changes in  $\log_2(\text{TE})$  for each drug condition compared to untreated control. Grey, all genes. Purple, ISR targets identified in the original study. Orange, genes fitting a strict thresholding paradigm to identify genes that display a 2-fold or greater increase in TE in Tm + ISRIB treatment compared to Tm treatment. XPRESSpipe-processed read alignments were quantified to *Homo sapiens* build CRCh38v98 using a protein-coding only, truncated GTF.

of these data with the updated XPRESSpipe methodology identifies genes whose translational down-regulation may play a role in the neurodegenerative effects of ISR and the neuroprotective properties of ISRIB [74–77]. Importantly, several of these genes were not identified as having significantly down-regulated TEs in the original analysis, which suggests why a focus on translational downregulation may have been foregone. In all, we identified seven genes with the regulatory paradigm of interest: significant decreases in TE during tunicamycin-induced ISR that are rescued or overexpressed in the ISR + ISRIB condition (Table 1, descriptions sourced from [80–82] (Figure 3D). RNA-Seq and ribosome-footprint coverage across these genes show that the significant changes in their TE are due to neither spurious, high-abundance fragments differentially present across libraries nor variance from an especially small number of mapped reads (Figure S6). This is an important consideration as the commonly suggested use of the CircLigase enzyme in published ribosome profiling library preparation protocols, which circularizes template cDNA before sequencing, can bias certain molecules' incorporation into sequencing libraries based on read-end base content alone [83].

Five (POMGNT1, RPL27, TKT, HSPA8, NDUFA11) out of the seven identified genes have annotated neurological functions or mutations in which lead to severe neurological disorders. Several of these disorders are directly tied to central carbon metabolism. Mutations in one other gene (RPS15A) generally present with metabolic disorders. None of these genes were identified as of interest in the original study using original methods; however, when re-processing the original count data provided with the original manuscript with DESeq2 [69] and the same expression pattern thresholding, four of these genes are present in the analysis (RPL27, TKT, HSPA8, RPS15A) S7. These observations suggest their regulation may be functionally important for the neurodegenerative effects of ISR and the neuroprotective properties of ISRIB. For example, NDUFA11 and TKT are two protein-coding genes whose functions are integrally tied to central carbon metabolism. NDUFA11 forms a subunit of mitochondrial complex I and TKT encodes a thiamine-dependent enzyme which channels excess sugars phosphates into glycolysis. Mutations in TKT present diseases associated with neurological phenotypes, and mutations in NDUFA11 present severe neurodegenerative phenotypes such as brain atrophy and encephalopathy. While at this stage speculative, it is interesting that the processing of these data for this manuscript provide a very conservative list of differentially expressed genes, and that the majority of which are associated with severe neurological phenotypes. It is therefore easy to speculate that TE regulation of these targets' abundance might be important in the neurodegeneration observed in prolonged ISR conditions. ISRIB's neuroprotective effects may, therefore, stem from a recovery of one or several of these entities' protein expression to wild-type or better levels. For example, NDUFA11 knockout models are

Table 1: **Translationally down-regulated genes during acute Tm treatment with recovered regulation during Tm + ISRIB treatment.** Gene names succeeded by an asterisk indicate these genes were identified in the original data when re-analyzed with DESeq2 [69]. Gene names succeeded by an ampersand indicate genes with strong neurological phenotypes.

Gene Name	Relevant Description
POMGNT1 <sup>&amp;</sup>	Participates in O-mannosyl glycosylation. Mutations have been associated with muscle-eye-brain diseases and congenital muscular dystrophies. Expressed especially in astrocytes, as well as in immature and mature neurons. Expressed across brain.
RPL27 <sup>*&amp;</sup>	Subunit of ribosome catalyzing protein synthesis. Expressed in cerebral cortex in embryonic tissue and/or stem cells. Mutations associated with Diamond-Blackfan Anemia 16, a metabolic disease, which may present with microcephaly.
TKT <sup>*&amp;</sup>	Encodes thiamine-dependent enzyme that channels excess sugars phosphates to glycolysis. Mutations associated with developmental delays and Wernicke-Korsakoff Syndrome, a metabolic and neuronal disease and associated with encephalopathy and dementia-like characteristics.
HSPA8 <sup>*&amp;</sup>	Encodes heat shock protein 70 member. Facilitates protein folding and localization. Diseases associated with mutations include Auditory System Disease and Brain Ischemia, both neurological disorders. Expressed in cerebral cortex in embryonic tissue and/or stem cells.
NDUFA11 <sup>&amp;</sup>	Encodes subunit of mitochondrial complex I, a vital part of the electron transport chain. Mutations associated with severe mitochondrial complex I deficiency. Related pathways include the GABAergic synapse. Associated diseases include brain atrophy, encephalopathy, and leber hereditary optic neuropathy. Overexpressed in frontal cortex.
HIST2H3D	Responsible for the nucleosome structure. No neurological descriptions currently known.
RPS15A <sup>*</sup>	Subunit of ribosome catalyzing protein synthesis. Diseases associated include Diamond-Blackfan Anemia, an inborn error of metabolism disease.

routinely used to model mitochondrial stress, which itself activates the ISR. Therefore, some mechanism of ISRIB may be allowing for recovery of NDUFA11 translation levels, explaining the neuroprotective capabilities of ISRIB. Though speculative without further validation, these ISRIB-responsive neuronal targets act as interesting cases for further validation and study in a model more representative of neurotoxic injury and disease than the HEK-293T model used in the original study. In all, this comparison demonstrates the utility of XPRESSpipe for rapid, user-friendly analysis and re-analysis of ribosome-profiling experiments in the pursuit of biological insights and hypothesis generation.

## Cost Analysis and Performance

XPRESSpipe functions can be computationally intensive, and thus, super-computing resources are recommended, especially when handling large datasets or when aligning to larger, more complex genomes. Many universities provide super-computing resources to their affiliates; however, in cases where these resources are not available, servers such as Amazon Web Services (AWS) [84] can be used to process sequencing data using XPRESSpipe. Table 2 outlines runtime statistics for the ISRIB dataset used in this study. The ISRIB ribosome profiling dataset contained a total of 32 raw sequence files that were aligned to *Homo sapiens*, thus it acts as a high-end estimate of the time required to process data with XPRESSpipe. For a comparable dataset, cost to use an AWS computational



node with similar specifications for the above elapsed time would be approximately 13.33 USD using an Amazon EC2 On-Demand m5.4xlarge node (however, significantly reduced rates are available if using Spot instances or by using the free tier) and storage cost would amount to around 11.5 USD/month on Amazon S3 storage (although much of the intermediate data does not need to be stored long term; however, input raw data should always be archived).

Table 2: **XPRESSpipe sub-module statistics for dataset GSE65778.** `geneCoverage` module performed on high-coverage gene. While some memory footprints are large in this test case, steps will scale based on available user resources.

Process	Command	Wallclock Time	Max RAM
Curate STAR reference	<code>curateReference</code>	00h 34m 01s	34.03gb
Truncate GTF	<code>modifyGTF -t</code>	00h 02m 45s	03.25gb
Read Pre-processing	<code>trim</code>	00h 18m 21s	00.48gb
Alignment and Post-processing	<code>align</code>	06h 16m 18s	38.00gb
Read Quantification	<code>count -c htseq</code>	04h 13m 04s	0.16gb
Isoform Abundance	<code>count -c cufflinks</code>	01h 09m 32s	2.36gb
Differential Expression (n=9)	<code>diffxpress</code>	00h 07m 50s	0.65gb
Read Distributions	<code>readDistribution</code>		
Metagene Analysis	<code>metagene</code>		
Gene Coverage (n=1)	<code>geneCoverage</code>	03h 49m 03s	19.22gb
Periodicity	<code>periodicity</code>	01h 00m 57s	61.89gb
Library Complexity	<code>complexity</code>	00h 48m 46s	1.98gb
rRNA probe	<code>rrnaProbe</code>		
Pipeline	<code>riboseq</code>		
Attribute		Value	
Total Raw Input		257 GB	
Total Output		500 GB	
Allocated CPUs		20	
Allocated Memory		64GB	

## Availability and Future Directions

We have described a new software suite, XPRESSyourself, a reference implementation of correct ribosome profiling data analysis built upon of a synthesis of new tools, old tools, and automated pipelines to aid in the processing and analysis of these data. XPRESSyourself is perpetually open source and protected under the GPL-3.0 license. Updates to the software are version controlled and maintained on GitHub [53]. Jupyter notebooks and video walkthroughs are included within the README files at [53]. Documentation is hosted on readthedocs [85] at [51] and [52]. Source code for associated analyses and figures for this manuscript can be accessed at [86]. The data used in this manuscript are available under the Gene Expression Omnibus persistent identifier GSE65778 [87] for ribosome profiling data and under the dbGaP Study Accession persistent identifier phs000178 [88] for the TCGA data.

Although RNA-Seq technologies are quite advanced, standardized computational protocols are much less established for ribosome profiling. As we discussed in this manuscript, this becomes problematic when individuals or groups are not using the most up-to-date or complete methods or may not be aware of particular biases or measures of quality control required to produce a reliable, high-quality sequencing study. XPRESSpipe handles these issues through on-going curation of benchmarked software tools and by simplifying the required user input. It also outputs all necessary quality control metrics so that the user can quickly assess the quality of their data and identify any systematic problems or technical biases that may compromise their analysis. Video walkthroughs, example scripts, and interactive command builders are available within this software suite to make these analyses accessible to experienced and inexperienced users alike. XPRESSyourself will enable individuals and labs to process and analyze their own data, which will result in quicker turnaround times of experiments and financial savings.

One particular benefit of XPRESSyourself is that it consolidates, streamlines, and/or introduces many tools specific to ribosome profiling processing and analysis. This includes curating GTF files with 5'- and 3'- truncated CDS annotations, rRNA probe design for subtractive hybridization of abundant rRNA contaminants, automated quality-control analyses to report on ribosome footprint periodicity and metagene coverage, and intron-less gene coverage profiles. These tools will help to democratize aspects of ribosome profiling for which software have not been previously publicly available.

We demonstrated the utility of the XPRESSyourself toolkit by re-analyzing a publicly available ribosome profiling dataset. From this analysis, we identified putative translational regulatory targets of the integrated stress response (ISR) that may contribute to its neurodegenerative effects and their rescue by the small-molecule ISR inhibitor, ISRIB. This highlights the importance of re-analyzing published datasets with more current methods, as improved analysis methodologies and updated organism genome references may result in new interpretations and hypotheses.

With the adoption of this flexible pipeline, the field of high-throughput sequencing, particularly ribosome profiling, can continue to standardize the processing protocol for associated sequence data and eliminate the variability that comes from the availability of a variety of software packages for various steps during sequence read processing. Additionally, XPRESSpipe consolidates various tools used by the ribosome profiling and RNA-Seq communities. With these tools, genome reference formatting and curation is automated and accessible to the public. Adoption of this tool will allow scientists to quickly process and access their data independently, guide them in understanding key considerations in processing their data, and standardize protocols for ribosome profiling and other RNA-Seq applications, thus increasing the reproducibility of sequencing analyses.

Table 3: **Software Description**

Project Name	XPRESSyourself
Project Home Page	<a href="https://github.com/XPRESSyourself">https://github.com/XPRESSyourself</a>
Archived Versions DOIs	10.5281/zenodo.3338669, 10.5281/zenodo.3337897
Operating Systems	macOS, Linux, CentOS
Programming Languages	Python, R
Other Requirements	Anaconda
License	GNU General Public License v3.0

## List of Abbreviations

AWS - Amazon Web Services, BAM - Binary Sequence Alignment Map, BED - Browser Extensible Data, cDNA - complementary DNA, CDS - coding sequence of gene, ChIP-seq - chromatin immunoprecipitation sequencing, CPU - central processing unit, dbGaP - Database of Genotypes and Phenotypes, DNA - deoxyribonucleic acid, FDR - false discovery rate, FPKM - fragments per kilobase of transcript per million, GEO - Gene Expression Omnibus, GTF - General Transfer Format, IGV - Integrative Genomics Viewer, ISR - integrated stress response, ISRIB - ISR inhibitor, mRNA - messenger RNA, nt - nucleotide, PCA - principal component analysis, PCR - polymerase chain reaction, RAM - random access memory, RNA - ribonucleic acid, RNA-Seq - RNA sequencing RPKM - reads per kilobase of transcript per million, RPM - reads per million, rRNA - ribosomal RNA, TCGA - The Cancer Genome Atlas, TE - translation efficiency, TPM - transcripts per million, UMI - unique molecular identifier, UTR - untranslated region

## Ethics Approval and Consent to Participate

Protected TCGA data were obtained through dbGaP project number 21674 and utilized according to the associated policies and guidelines.

## Consent for Publication

Protected TCGA data were obtained through dbGaP project number 21674 and utilized according to the associated policies and guidelines.

Table 4: **Author ORCIDs**

<b>Author</b>	<b>ORCID</b>
J.A.B.	0000-0002-5096-0558
J.R.B.	0000-0001-5470-8299
J.T.M.	0000-0002-3017-8665
Y.O.	0000-0001-9523-1044
A.J.B.	0000-0003-2273-8922
A.R.Q.	0000-0003-1756-0859
J.G.	0000-0001-7568-6789
J.R.	0000-0002-2710-9765

## Competing Interests

The authors declare that they have no competing interests.

## Funding

J.A.B. received support from the National Institute of Diabetes and Digestive and Kidney Diseases (NIDDK) Inter-disciplinary Training Grant T32 Program in Computational Approaches to Diabetes and Metabolism Research, 1T32DK11096601 to Wendy W. Chapman and Simon J. Fisher. J.T.M. received support as an HHMI Fellow of the Jane Coffin Childs Memorial Fund for Medical Research. A.J.B received support from the National Cancer Institute (NCI) Predoctoral to Postdoctoral Fellow Transition Award, K00CA212445. This work was supported by NIDDK fellowship 1T32DK11096601 (to J.A.B.) and NIH grant R35GM13185 (to J.R.). The computational resources used were partially funded by the NIH Shared Instrumentation Grant 1S10OD021644-01A1.

## Contributions

J.A.B. conceptualized and administered the project; performed all investigation, analysis, visualization, and data curation; provisioned resources; and wrote the original draft of this manuscript. J.A.B., J.R.B., J.T.M., and J.G. and developed the methodology. J.A.B. and J.R.B. designed and wrote the software. J.A.B., J.T.M., A.J.B., and Y.O. performed software validation. J.A.B. and J.R. acquired funding. J.R., A.R.Q., and J.G. supervised the study. All authors reviewed and edited this manuscript.

## Acknowledgments

The authors wish to thank Michael T. Howard for helpful discussions concerning ribosome profiling and sequencing analysis. The authors also wish to thank Mark E. Wadsworth, Ryan Miller, and Michael J. Cormier for helpful discussions on pipeline design. They also wish to thank T. Cameron Waller for helpful discussions related to pipeline design and biological analysis. The support and resources from the Center for High-Performance Computing at the University of Utah are gratefully acknowledged. The results published here are in whole or part based upon data generated by the TCGA Research Network [63].

## References

- [1] N. Ingolia, S. Ghaemmaghami, J. Newman, J. Weissman. Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science* **324**, 218 (2009). Available from: <https://doi.org/10.1126/science.1168978>.
- [2] G. Brar, J. Weissman. Ribosome profiling reveals the what, when, where and how of protein synthesis. *Nat Rev Mol Cell Biol* **16**, 651 (2015). Available from: <https://doi.org/10.1038/nrm4069>.
- [3] N. McGlincy, N. Ingolia. Transcriptome-wide measurement of translation by ribosome profiling. *Methods* **126**, 112 (2017). Available from: <https://doi.org/10.1016/j.ymeth.2017.05.028>.
- [4] M. Gerashchenko, V. Gladyshev. Translation inhibitors cause abnormalities in ribosome profiling experiments. *Nucleic Acids Res* **42** (2014). Available from: <https://doi.org/10.1093/nar/gku671>.
- [5] A. Bartholomäus, C. D. Campo, I. Z. Mapping the non-standardized biases of ribosome profiling. *Biol Chem* **397** (2016). Available from: <https://doi.org/https://doi.org/10.1515/hsz-2015-0197>.
- [6] J. Hussmann, S. Patchett, A. Johnson, S. Sawyer, W. Press. Understanding Biases in Ribosome Profiling Experiments Reveals Signatures of Translation Dynamics in Yeast. *PLoS Genet* **11** (2015). Available from: <https://doi.org/https://doi.org/10.1371/journal.pgen.1005732>.
- [7] A. Diamant, T. Tuller. Estimation of ribosome profiling performance and reproducibility at various levels of resolution. *Biol Direct* **11** (2016). Available from: <https://doi.org/https://doi.org/10.1186/s13062-016-0127-4>.

- [8] M. Gerashchenko, V. Gladyshev. Ribonuclease selection for ribosome profiling. *Nucleic Acids Res* **45** (2017). Available from: <https://doi.org/https://doi.org/10.1093/nar/gkw822>.
- [9] C. Sidrauski<sup>1</sup>, A. McGeachy, N. Ingolia, P. Walter. The small molecule ISRIB reverses the effects of eIF2 $\alpha$  phosphorylation on translation and stress granule assembly. *eLife* (2015). Available from: <https://doi.org/10.7554/eLife.05033>.
- [10] F. Mohammad, C. Woolstenhulme, R. Green, A. Buskirk. Clarifying the Translational Pausing Landscape in Bacteria by Ribosome Profiling. *Cell Rep* **14** (2016). Available from: <https://doi.org/https://doi.org/10.1016/j.celrep.2015.12.073>.
- [11] G. Li, E. Oh, J. Weissman. The anti-Shine-Dalgarno sequence drives translational pausing and codon choice in bacteria. *Nature* **484** (2012). Available from: <https://doi.org/https://doi.org/10.1038/nature10965>.
- [12] A. Lecanda, *et al.*. Dual randomization of oligonucleotides to reduce the bias in ribosome-profiling libraries. *Methods* **107** (2016). Available from: <https://doi.org/https://doi.org/10.1016/j.ymeth.2016.07.011>.
- [13] X. Gao, *et al.*. Quantitative profiling of initiating ribosomes in vivo. *Nat Methods* **12** (2015). Available from: <https://doi.org/https://doi.org/10.1038/nmeth.3208>.
- [14] B. Langmead, C. Trapnell, M. Pop, S. Salzberg. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10** (2009). Available from: <https://doi.org/https://doi.org/10.1186/gb-2009-10-3-r25>.
- [15] E. Afgan, *et al.*. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic Acids Res.* **46**, W537 (2018). Available from: <https://doi.org/10.1093/nar/gky379>.
- [16] A. Michel, *et al.*. RiboGalaxy: A browser based platform for the alignment, analysis and visualization of ribosome profiling data. *RNA Biol* **13**, 316 (2016). Available from: <https://doi.org/10.1080/15476286.2016.1141862>.
- [17] Nextflow. <https://www.nextflow.io/example4.html>.
- [18] DNAnexus. [https://github.com/dnanexus/tophat\\_cufflinks\\_rnaseq](https://github.com/dnanexus/tophat_cufflinks_rnaseq).
- [19] O. Carja, T. Xing, E. Wallace, J. Plotkin, P. Shah. riboviz: analysis and visualization of ribosome profiling datasets. *BMC Bioinformatics* **18** (2017). Available from: <https://doi.org/10.1186/s12859-017-1873-8>.

- [20] C. Artieri, H. Fraser. Accounting for biases in riboprofiling data indicates a major role for proline in stalling translation. *Genome Res* **24**, 2011 (2014). Available from: <https://doi.org/10.1101/gr.175893.114>.
- [21] J. Hussmann, S. Patchett, A. Johnson, S. Sawyer, W. Press. Understanding Biases in Ribosome Profiling Experiments Reveals Signatures of Translation Dynamics in Yeast. *PLoS Genet* **11** (2015). Available from: <https://doi.org/10.1371/journal.pgen.1005732>.
- [22] D. Weinberg, *et al.*. Improved Ribosome-Footprint and mRNA Measurements Provide Insights into Dynamics and Regulation of Yeast Translation. *Cell Rep* **14**, 1787 (2016). Available from: <https://doi.org/10.1016/j.celrep.2016.01.043>.
- [23] J. Robinson, *et al.*. Integrative Genomics Viewer. *Nat Biotechnol* **29**, 24 (2011). Available from: <https://doi.org/10.1038/nbt.1754>.
- [24] Z. Costello, H. Martin. A machine learning approach to predict metabolic pathway dynamics from time-series multiomics data. *NPJ Syst Biol Appl* **4** (2018). Available from: <https://doi.org/10.1038/s41540-018-0054-3>.
- [25] V. Funari, S. Canosa. The Importance of Bioinformatics in NGS: Breaking the Bottleneck in Data Interpretation. *Science* **344**, 653 (2014). Available from: <https://doi.org/10.1126/science.344.6184.653-c>.
- [26] R. Kumari, A. Michel, P. Baranov. PausePred and Rfeet: webtools for inferring ribosome pauses and visualizing footprint density from ribosome profiling data. *RNA* **24** (2018). Available from: <https://doi.org/10.1261/rna.065235.117>.
- [27] C. Oertlin, *et al.*. Generally applicable transcriptome-wide analysis of translation using anota2seq. *Nucleic Acids Res* **47** (2019). Available from: <https://doi.org/10.1093/nar/gkz223>.
- [28] A. Popa, *et al.*. RiboProfiling: a Bioconductor package for standard Ribo-seq pipeline processing. *F1000Res* **5** (2016). Available from: <https://doi.org/10.12688/f1000research.8964.1>.
- [29] W. Li, W. Wang, P. Uren, L. Penalva, A. Smith. Riborex: fast and flexible identification of differential translation from Ribo-seq data. *Bioinformatics* **33** (2017). Available from: <https://doi.org/10.1093/bioinformatics/btx047>.
- [30] S. Verbruggen, G. Menschaert. mQC: A post-mapping data exploration tool for ribosome profiling. *Comput Methods Programs Biomed* (2018). Available from: <https://doi.org/10.1016/j.cmpb.2018.10.018>.

- [31] Å. Birkeland, K. Chyżyńska, E. Valen. Shoelaces: an interactive tool for ribosome profiling processing and visualization. *BMC Genomics* **19** (2018). Available from: <https://doi.org/10.1186/s12864-018-4912-6>.
- [32] N. Davidson, A. Hawkins, A. Oshlack. SuperTranscripts: a data driven reference for analysis and visualisation of transcriptomes. *Genome Biol* **18** (2017). Available from: <https://doi.org/10.1186/s13059-017-1284-1>.
- [33] T. Backman, T. Girke. systemPipeR: NGS workflow and report generation environment. *BMC Bioinformatics* **17** (2016). Available from: <https://doi.org/10.1186/s12859-016-1241-0>.
- [34] H. Tjeldnes, K. Labun. ORFik: Open Reading Frames in Genomics. <https://github.com/JokingHero/ORFik> (2017). Available from: <https://doi.org/10.18129/B9.bioc.ORFik>.
- [35] T. Martin, I. Erte, P. Tsai, J. Bell. coMET: an R plotting package to visualize regional plots of epigenome-wide association scan results. *QG14* (2014). Available from: <http://quantgen.soc.srcf.net/qg14/>.
- [36] T. Martin, I. Yet, P. Tsai, J. Bell. coMET: visualisation of regional epigenome-wide association scan results and DNA co-methylation patterns. *BMC Bioinformatics* **16** (2015). Available from: <https://doi.org/10.1186/s12859-015-0568-2>.
- [37] T. Hardcastle. riboSeqR. Available from: <https://doi.org/10.18129/B9.bioc.riboSeqR>.
- [38] F. Ramírez, F. Dündar, S. Diehl, B. Grüning, T. Manke. deepTools: a flexible platform for exploring deep-sequencing data. *Nucleic Acids Res* **42** (2014). Available from: <https://doi.org/10.1093/nar/gku365>.
- [39] Picard. <https://broadinstitute.github.io/picard/>.
- [40] S. Zhang, *et al.* Analysis of Ribosome Stalling and Translation Elongation Dynamics by Deep Learning. *Cell Syst* **5** (2017). Available from: <https://doi.org/10.1016/j.cels.2017.08.004>.
- [41] P. O'Connor, D. Andreev, P. Baranov. Comparative survey of the relative impact of mRNA features on local ribosome profiling read density. *Nat Commun* **7** (2016). Available from: <https://doi.org/10.1038/ncomms12915>.
- [42] Z. Xiao, Q. Zou, Y. Liu, X. Yang. Genome-wide assessment of differential translations with ribosome profiling data. *Nat Commun* **7** (2016). Available from: <https://doi.org/10.1038/ncomms11194>.
- [43] Y. Zhong, *et al.* RiboDiff: detecting changes of mRNA translation efficiency from ribosome footprints. *Bioinformatics* **33** (2017). Available from: <https://doi.org/10.1093/bioinformatics/btw585>.



- [44] L. Calviello, *et al.*. Detecting actively translated open reading frames in ribosome profiling data. *Nat Methods* **13** (2016). Available from: <https://doi.org/10.1038/nmeth.3688>.
- [45] H. Wang, J. McManus, C. Kingsford. Isoform-level ribosome occupancy estimation guided by transcript abundance with Ribomap. *Bioinformatics* **32** (2016). Available from: <https://doi.org/10.1093/bioinformatics/btw085>.
- [46] P. Spealman, H. Wang, G. May, C. Kingsford, C. McManus. Exploring Ribosome Positioning on Translating Transcripts with Ribosome Profiling. *Methods Mol Biol* **1358** (2016). Available from: [https://doi.org/10.1007/978-1-4939-3067-8\\_5](https://doi.org/10.1007/978-1-4939-3067-8_5).
- [47] J. Dunn, J. Weissman. Plastid: nucleotide-resolution analysis of next-generation sequencing and genomics data. *BMC Genomics* **17** (2016). Available from: <https://doi.org/https://doi.org/10.1186/s12864-016-3278-x>.
- [48] P. Perkins, S. Mazzoni-Putman, A. Stepanova, J. Alonso, S. Heber. RiboStreamR: a web application for quality control, analysis, and visualization of Ribo-seq data. *BMC Genomics* **20** (2019). Available from: <https://doi.org/https://doi.org/10.1186/s12864-019-5700-7>.
- [49] H. Fang, *et al.*. Scikit-ribo Enables Accurate Estimation and Robust Modeling of Translation Dynamics at Codon Resolution. *Cell Syst* **6** (2018). Available from: <https://doi.org/https://doi.org/10.1016/j.cels.2017.12.007>.
- [50] S. Chun, C. Rodriguez, P. Todd, R. Mills. SPECtre: a spectral coherence-based classifier of actively translated transcripts from ribosome profiling sequence data. *BMC Bioinformatics* **17** (2016). Available from: <https://doi.org/https://doi.org/10.1186/s12859-016-1355-4>.
- [51] XPRESSpipe documentation. <https://xpresspipe.readthedocs.io/en/latest/>.
- [52] XPRESSplot documentation. <https://xpressplot.readthedocs.io/en/latest/>.
- [53] XPRESSyourself. <https://github.com/XPRESSyourself/>.
- [54] A. Dobin, *et al.*. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15 (2013). Available from: <https://doi.org/10.1093/bioinformatics/bts635>.
- [55] G. Baruzzo, *et al.*. Simulation-based comprehensive benchmarking of RNA-seq aligners. *Nat Methods* **14**, 135 (2017). Available from: <https://doi.org/10.1038/nmeth.4106>.

- [56] I. Raplee, A. Evsikov, C. M. de Evsikova. Aligning the Aligners: Comparison of RNA Sequencing Data Alignment and Gene Expression Quantification Tools for Clinical Breast Cancer Research. *J Pers Med* **9** (2019). Available from: <https://doi.org/10.3390/jpm9020018>.
- [57] T. Tuller, H. Zur. Multiple roles of the coding sequence 5' end in gene expression regulation. *Nucleic Acids Res* **42** (2015). Available from: <https://doi.org/https://doi.org/10.1093/nar/gku1313>.
- [58] S. Anders, P. Pyl, W. Huber. HTSeq—a Python framework to work with high-throughput sequencing data. *Bioinformatics* **31**, 166 (2015). Available from: <https://doi.org/10.1093/bioinformatics/btu638>.
- [59] C. Trapnell, *et al.*. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* **7** (2012). Available from: <https://doi.org/10.1038/nprot.2012.016>.
- [60] Y. Fu, P. Wu, T. Beane, P. Zamore, Z. Weng. Elimination of PCR duplicates in RNA-seq and small RNA-seq using unique molecular identifiers. *BMC Genomics* **19** (2018). Available from: <https://doi.org/10.1186/s12864-018-4933-1>.
- [61] T. Smith, A. Heger, I. Sudbery. UMI-tools: Modelling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy. *Genome Res* **27** (2017). Available from: <https://doi.org/10.1101/gr.209601.116>.
- [62] C. Robert, M. Watson. Errors in RNA-Seq quantification affect genes of relevance to human disease. *Genome Biol* **16** (2015). Available from: <https://doi.org/10.1186/s13059-015-0734-x>.
- [63] The Cancer Genome Atlas. <https://portal.gdc.cancer.gov>.
- [64] C. Evans, J. Hardin, D. Stoebe. Selecting between-sample RNA-Seq normalization methods from the perspective of their assumptions. *Brief Bioinform* **19**, 776–792 (2018). Available from: <https://doi.org/10.1093/bib/bbx008>.
- [65] J. Leek, W. Johnson, H. Parker, A. Jaffe, J. Storey. The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics* **28** (2012). Available from: <https://doi.org/10.1093/bioinformatics/bts034>.
- [66] C. Wu, B. Zinshteyn, K. Wehner, R. Green. High-Resolution Ribosome Profiling Defines Discrete Ribosome Elongation States and Translational Regulation during Cellular Stress. *Mol Cell* **73** (2019). Available from: <https://doi.org/10.1016/j.molcel.2018.12.009>.

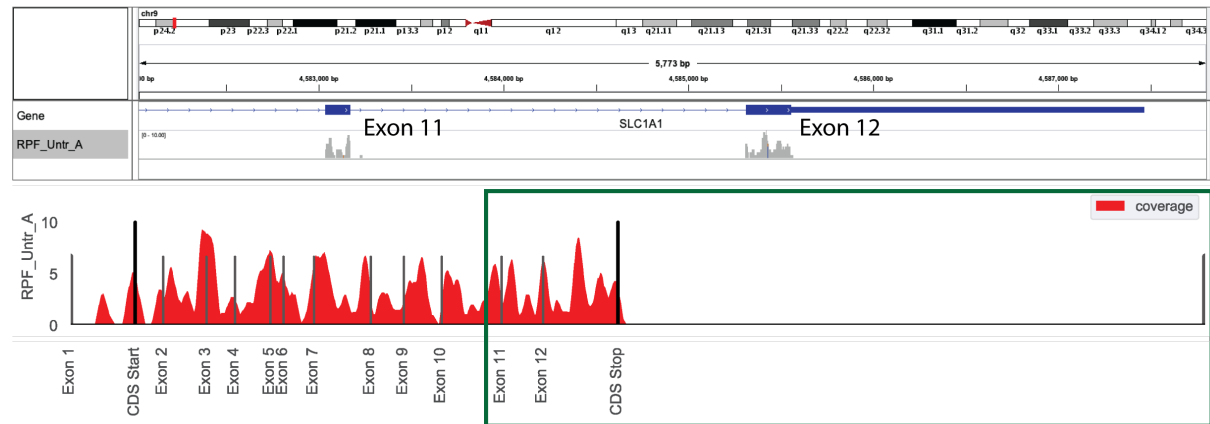
- [67] S. Sayols, D. Scherzinger, H. Klein. dupRadar: a Bioconductor package for the assessment of PCR artifacts in RNA-Seq data. *BMC Bioinformatics* **17**, 428 (2016). Available from: <https://doi.org/10.1186/s12859-016-1276-2>.
- [68] F. Lauria, *et al.*. riboWaltz: Optimization of ribosome P-site positioning in ribosome profiling data. *PLoS Comput Biol* **14** (2018). Available from: <https://doi.org/10.1371/journal.pcbi.1006169>.
- [69] M. Love, W. Huber, S. Anders. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* **15** (2014). Available from: <https://doi.org/10.1186/s13059-014-0550-8>.
- [70] H. Harding, *et al.*. An integrated stress response regulates amino acid metabolism and resistance to oxidative stress. *Mol Cell* **11** (2003). Available from: [https://doi.org/10.1016/S1097-2765\(03\)00105-9](https://doi.org/10.1016/S1097-2765(03)00105-9).
- [71] D. Santos-Ribeiro, L. Godinas, C. Pilette, F. Perros. The integrated stress response system in cardiovascular disease. *Drug Discov Today* **23** (2018). Available from: <https://doi.org/10.1016/j.drudis.2018.02.008>.
- [72] H. Rabouw, *et al.*. Small molecule ISRIB suppresses the integrated stress response within a defined window of activation. *Proc Natl Acad Sci U S A* **116** (2019). Available from: <https://doi.org/10.1073/pnas.1815767116>.
- [73] J. Tsai, *et al.*. Structure of the nucleotide exchange factor eIF2B reveals mechanism of memory-enhancing molecule. *Science* **359** (2018). Available from: <https://doi.org/10.1126/science.aag0939>.
- [74] A. Choua, *et al.*. Inhibition of the integrated stress response reverses cognitive deficits after traumatic brain injury. *Proc Natl Acad Sci U S A* **114** (2017). Available from: <https://doi.org/10.1073/pnas.1707661114>.
- [75] M. Halliday, *et al.*. Partial restoration of protein synthesis rates by the small molecule ISRIB prevents neurodegeneration without pancreatic toxicity. *Cell Death Dis* **6** (2015). Available from: <https://doi.org/10.1038/cddis.2015.49>.
- [76] C. Sidrauski, *et al.*. Pharmacological brake-release of mRNA translation enhances cognitive memory. *Elife* **2** (2013). Available from: <https://doi.org/10.7554/eLife.00498>.
- [77] Y. Sekine, *et al.*. Stress responses. Mutations in a translation initiation factor identify the target of a memory-enhancing compound. *Science* **348** (2015). Available from: <https://doi.org/10.1126/science.aaa6986>.
- [78] D. Kim, *et al.*. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol* **14** (2013). Available from: <https://doi.org/10.1186/gb-2013-14-4-r36>.

- [79] S. Anders, W. Huber. Differential expression analysis for sequence count data. *Genome Biol* **11** (2010). Available from: <https://doi.org/10.1186/gb-2010-11-10-r106>.
- [80] GeneCards. <https://www.genecards.org/>. Accessed 27 June 2019.
- [81] National Center for Biotechnology Information. <https://www.ncbi.nlm.nih.gov/gene/>. Accessed 27 June 2019.
- [82] UniProt. <https://www.uniprot.org/uniprot/>. Accessed 27 June 2019.
- [83] R. Tunney, *et al.*. Accurate design of translational output by a neural network model of ribosome distribution. *Nat Struct Mol Biol* **25**, 577 (2018). Available from: <https://doi.org/10.1038/s41594-018-0080-2>.
- [84] Amazon Web Services. <https://aws.amazon.com>.
- [85] Read the Docs. <https://readthedocs.org/>.
- [86] Manuscript code. [https://github.com/XPRESSyourself/xpressyourself\\_manuscript/tree/master/supplemental\\_files](https://github.com/XPRESSyourself/xpressyourself_manuscript/tree/master/supplemental_files). Available from: <https://doi.org/DOI: 10.5281/zenodo.3337599>.
- [87] Ribosome Profiling GEO Accession. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE65778>.
- [88] TCGA dbGaP Accession. ([https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\\_id=phs000178.v10.p8](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\_id=phs000178.v10.p8)).



A

SLC1A1



B

TSPAN33

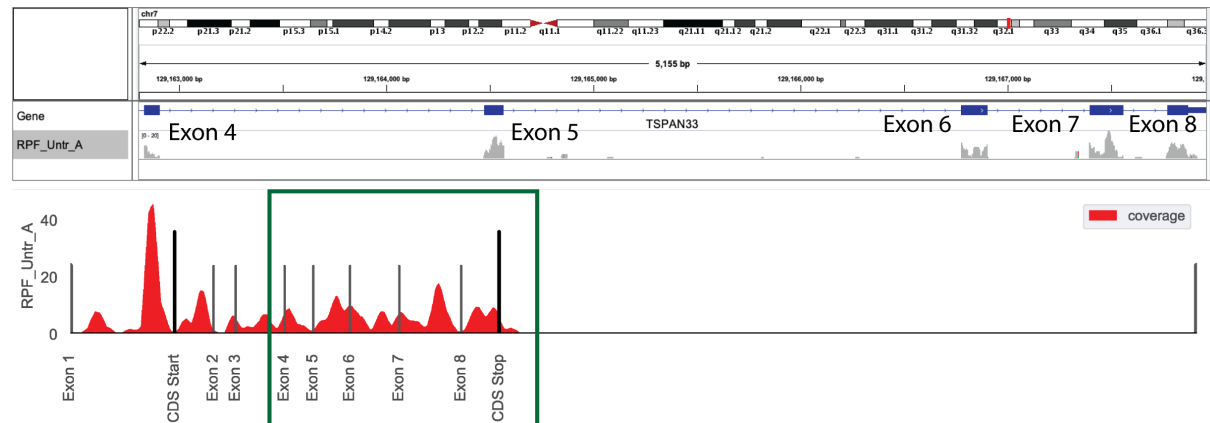


Figure S2: **Comparison between IGV browser and geneCoverage output.** A) Gene coverage from IGV (above) and XPRESSpipe (below) for SLC1A1. B) Gene coverage from IGV (above) and XPRESSpipe (below) for TSPAN33. Introns collapsed by XPRESSpipe. Green box, region shown in corresponding IGV window comparing outputs between the two programs.

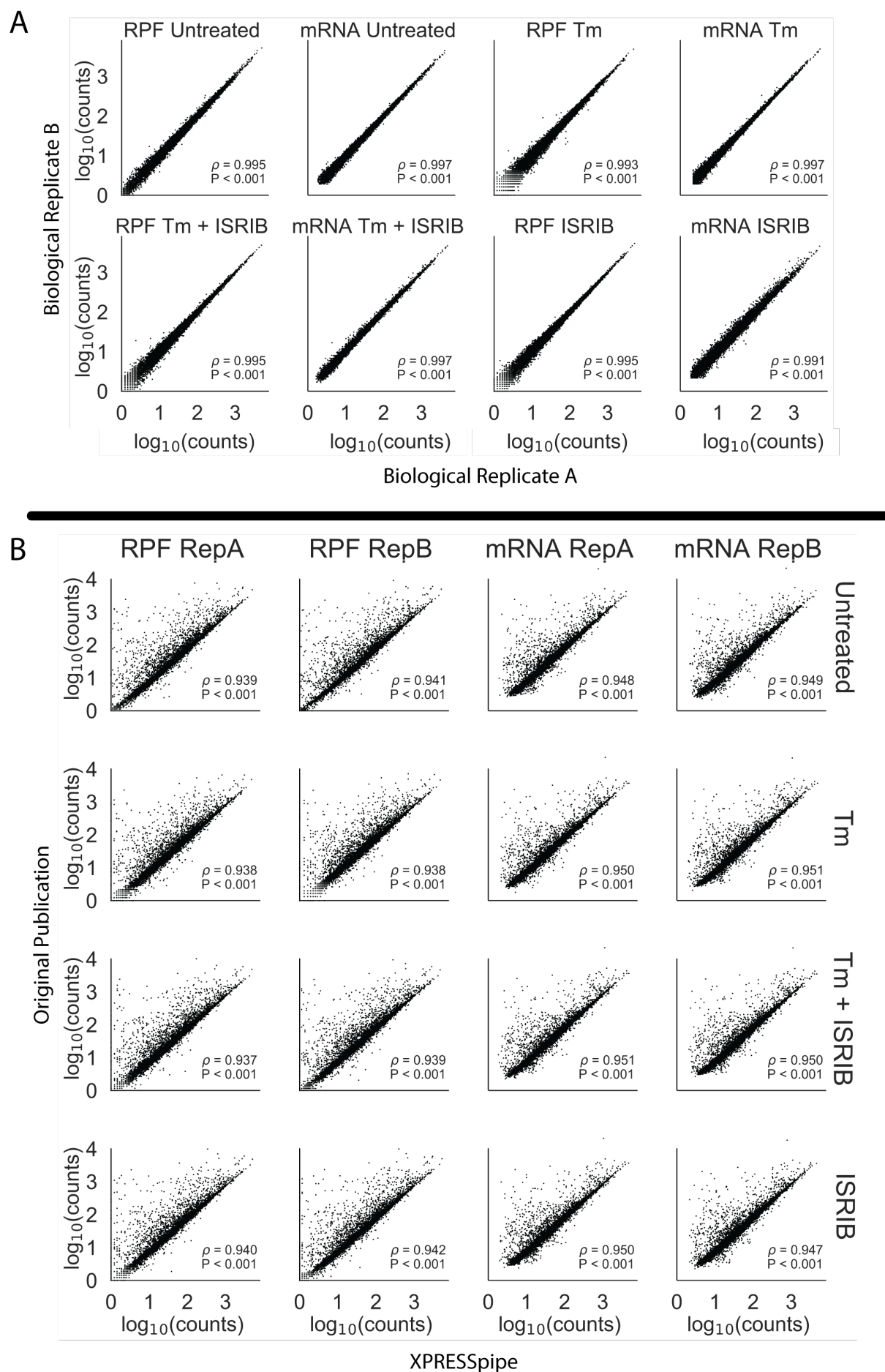


Figure S3: **Comparison between processed data produced by XPRESSpipe and original study.** Genes were eliminated from analysis if any RNA-Seq sample for that gene had fewer than 10 counts. A) Comparison of biological replicate read counts processed by XPRESSpipe. B) Comparison of read counts per gene between count data from the original study and the same raw

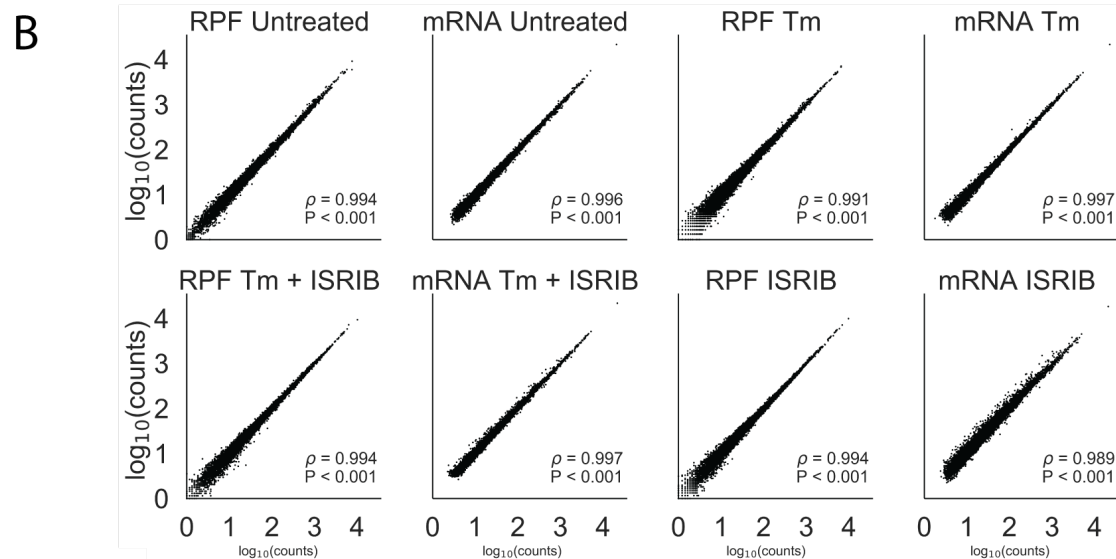
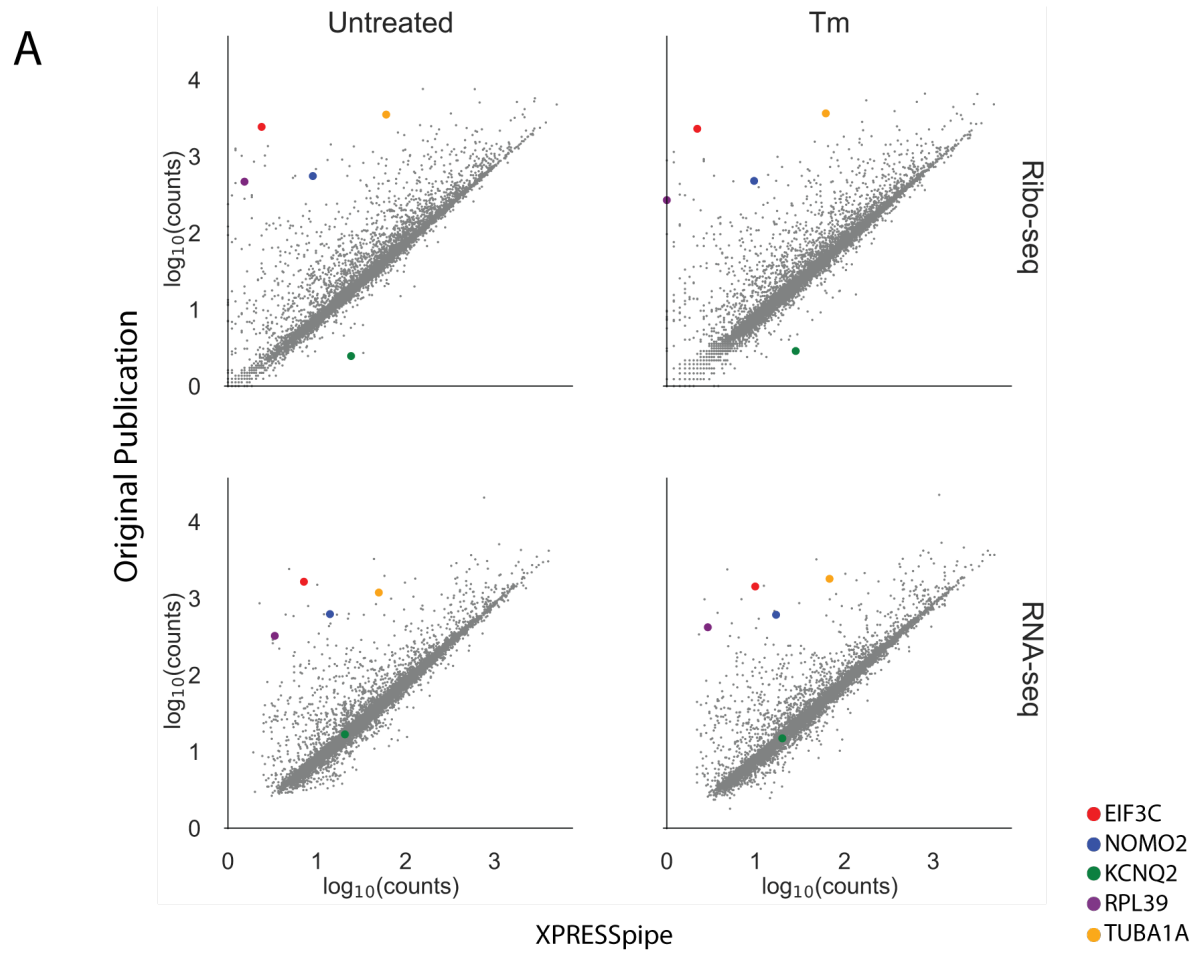


Figure S4: **Original ISRIB count data plotted against XPRESSpipe-processed data reveals systematic differences between the analytical regimes.** A) Selected highlighted genes show consistent differences between processing methods. B) Spearman correlation plots using the data table provided as supplementary data with the original ISRIB manuscript comparing biological replicates. RPF, ribosome-protected footprint. Tm, tunicamycin. All  $\rho$  values reported are Spearman correlation coefficients.



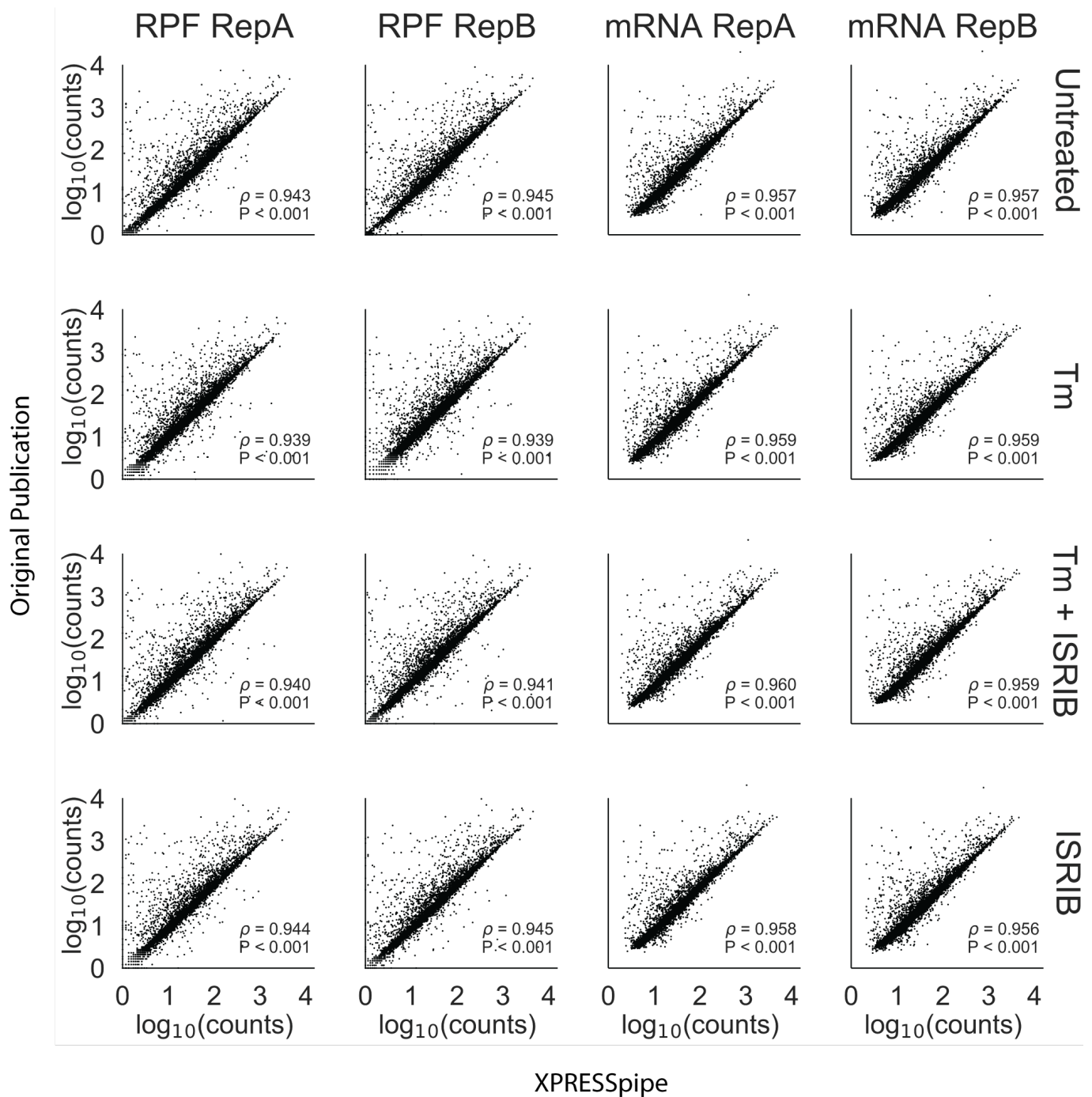


Figure S5: **Original ISRIB count data plotted against XPRESSpipe-processed data quantified using same reference version reveals mild improvement in comparability between the analytical regimes.** Original samples were processed using Ensembl human build GRCh38v72, as in the original manuscript, and compared with the original count data provided with the manuscript. XPRESSpipe-prepared counts were thresholded similarly as the original data (each gene needed to have at least 10 counts across all mRNA samples). RepA, biological replicate A. RepB, biological replicate B. RPF, ribosome-protected footprint. Tm, tunicamycin. All  $\rho$  values reported are Spearman correlation coefficients.

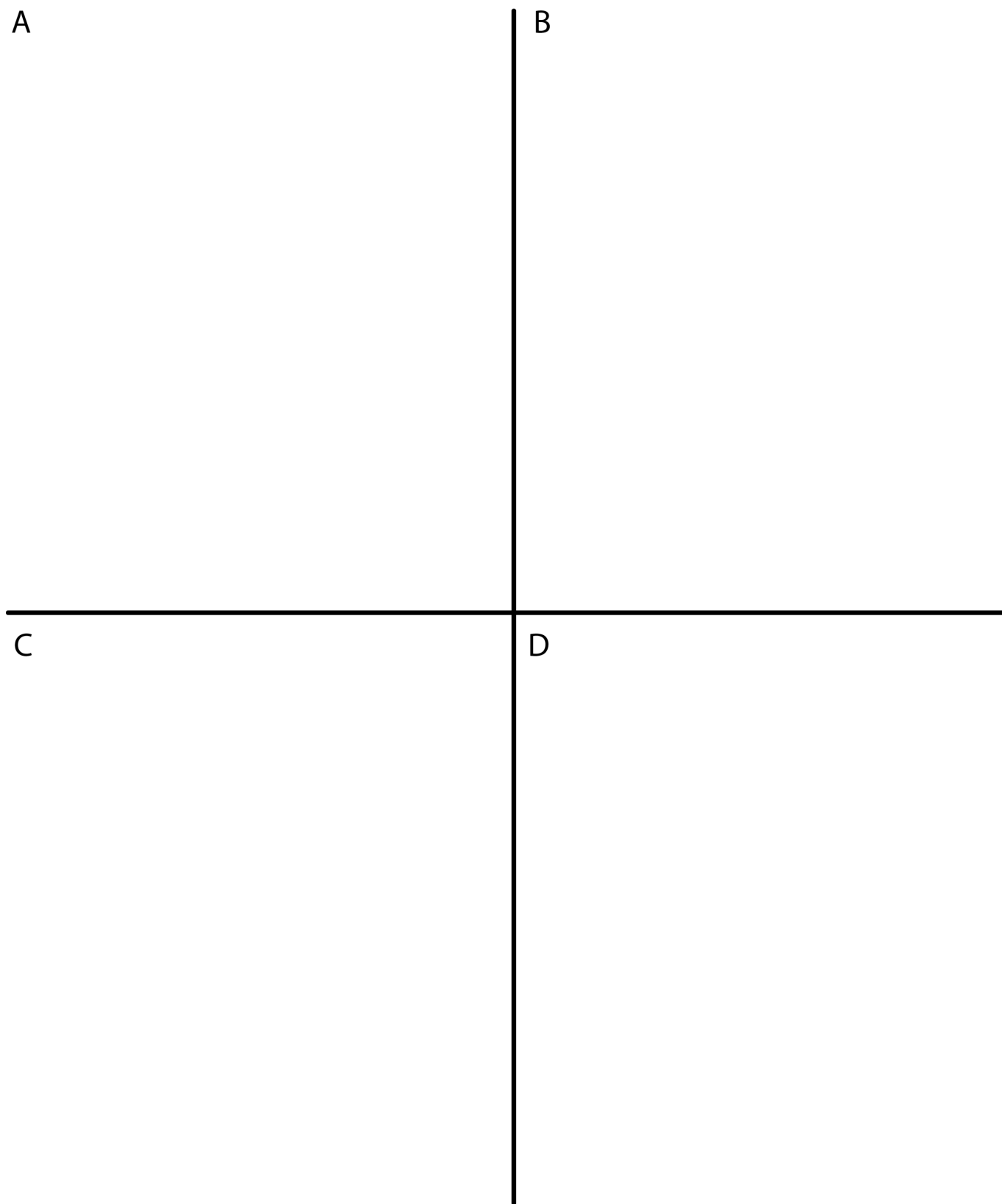


Figure S6: **Gene coverage plots for neurologically annotated genes passing strict thresholding.** Coverage plots were generated using XPRESSpipe's geneCoverage module, which collapses introns within the representation.

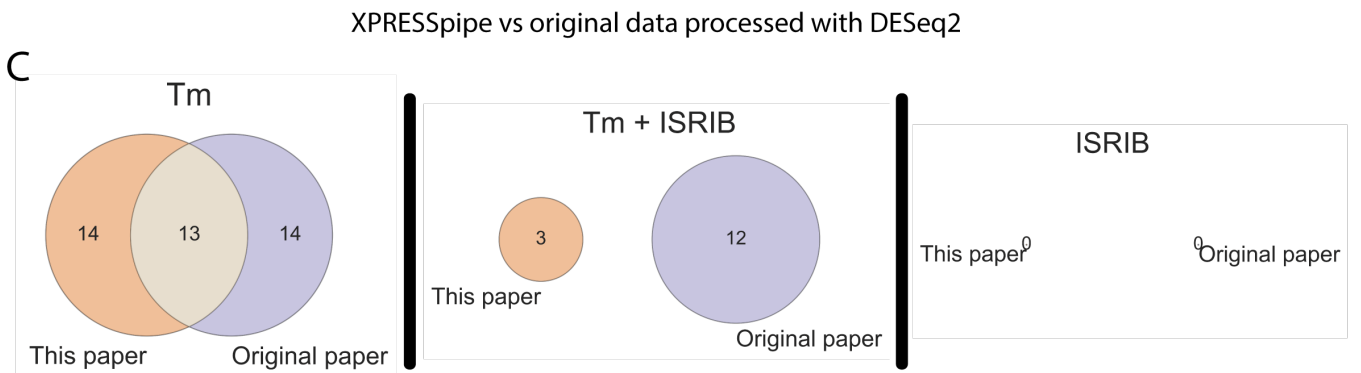
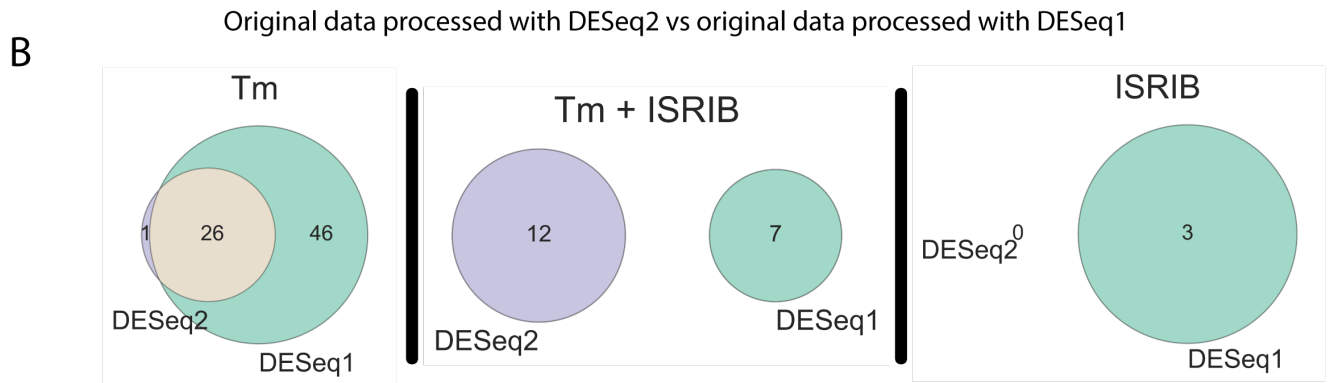
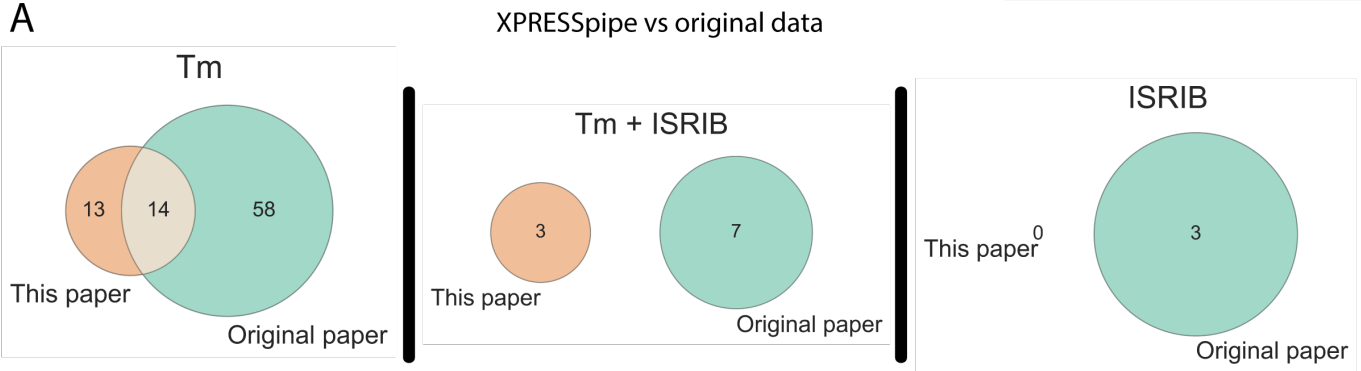


Figure S7: **Cross-method analysis comparisons.** A) XPRESSpipe-processed data (orange) versus data as originally presented within original manuscript using original methods (green). B) Comparison of analyses using provided count table in original publication using DESeq2 (purple) versus original analysis provided in manuscript using DESeq1 (green). C) XPRESSpipe-processed (orange) versus originally-processed data (purple), both using DESeq2 for differential expression analysis. Thresholds used were the same as those used in the original study:  $|\log_2(\text{Fold Change})| > 1$ ,  $\text{FDR} < 0.1$ .