

¹ **XPRESSyourself: Enhancing, Standardizing, and
2 Automating Ribosome Profiling Computational
3 Analyses Yields Improved Insight into Data**

⁴ **Jordan A. Berg,^{1*} Jonathan R. Belyeu,² Jeffrey T. Morgan,¹ Yeyun Ouyang,¹ Alex J. Bott,¹
⁵ Aaron R. Quinlan,^{2,4,5} Jason Gertz,³ Jared Rutter^{1,6*}**

⁶

⁷ ¹Department of Biochemistry, University of Utah, Salt Lake City, UT, USA, 84112.

⁸ ²Department of Human Genetics, University of Utah, Salt Lake City, UT, USA, 84112.

⁹ ³Department of Oncological Sciences, University of Utah, Salt Lake City, UT, USA, 84112.

¹⁰ ⁴USTAR Center for Genetic Discovery, University of Utah, Salt Lake City, UT, USA, 84112.

¹¹ ⁵Department of Biomedical Informatics, University of Utah, Salt Lake City, UT, USA, 84112.

¹² ⁶Howard Hughes Medical Institute, University of Utah, Salt Lake City, UT, USA, 84112.

¹³ *Address correspondence to: jordan.berg@biochem.utah.edu, rutter@biochem.utah.edu

14 **Abstract**

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

28

29 **Introduction**

30 High-throughput sequencing data has revolutionized biomedical and biological research. One such application of this
31 consequential technology is ribosome profiling, which, coupled with bulk RNA-Seq, measures translation efficiency,
32 translation pausing, novel protein translation products, and more [1–3]. Though the experimental procedures for
33 ribosome profiling have matured, an abundance of biases and peculiarities associated with each analytical method
34 or tool are still present and may often be obscured to a new user of this methodology [4–8]. Additionally, standardized
35 methods for handling this unique data type remain elusive. This has been problematic and evidenced by various
36 studies using vague or opaque methods for data analysis (for examples, see [9–13]), or methods rely on outdated
37 tools [5]. Very few labs have the tools necessary to separate the biological signals in ribosome profiling data from the
38 inherent biases of the experimental measurements, and these tools are not readily accessible by the community. This
39 is a critical time in the rapidly expanding influence of ribosome profiling. For too long, the bioinformatic know-how of
40 this incredibly powerful technique has been limited to a small handful of labs. As more and more ribosome profiling
41 studies are performed, more and more labs will lack the ability to analyze their data with ease and fidelity. Few, if

42 any, extant pipelines or toolkits offer a thorough set of integrated tools for assessing standard quality control metrics
43 or performing proper reference curation to reduce systematic biases across any organism, particularly with ribosome
44 profiling data [14–18].

45 For example, one issue in ribosome profiling is the pile-up of ribosomes at the 5'- and 3'- ends of coding regions
46 within a transcript, a systematic biological signal arising from the slower kinetics of ribosome initiation and termination
47 compared to translation elongation and is generally regarded to not accurately reflect measurements of translation
48 efficiency. These signals are further exacerbated by pre-treatment with cycloheximide during ribosome footprint
49 harvesting [4, 19, 20]. These pile-ups can dramatically skew ribosome footprint quantification and measurements of
50 translational efficiency. Current practices in the field recommend excluding pile-up-prone regions when quantifying
51 ribosome profiling alignments as they lead to noisier estimations of translation efficiency [3, 21]; however, no publicly
52 available computational tools currently exist to facilitate these automated adjustments to reference transcripts. Curating
53 references properly and robustly requires advanced implementations. In addition, downstream data visualization
54 methods presently available are often not optimized to analyze and compare translation regulatory regions of a gene.

55 To address deficiencies in the public ribosome profiling computational toolkit, we developed XPRESSyourself, a
56 computational toolkit and adaptable, end-to-end pipeline that bridges these and other gaps in ribosome profiling data
57 analysis. XPRESSyourself implements the complete suite of tools necessary for comprehensive ribosome profiling
58 and bulk RNA-Seq data processing and analysis in a robust and easy-to-use fashion, often packaging tasks that
59 would typically require hundreds to thousands of lines of code into a single command. For instance, XPRESSyourself
60 creates the mRNA annotation files necessary to remove confounding systematic factors during quantification and
61 analysis of ribosome profiling data, allowing for accurate measurements of translation efficiency. It provides the
62 built-in capacity to quantify and visualize differential upstream open-reading frame (uORF) usage by generating
63 IGV-like, intron-less plots for easier visualization [22]. The ability to visualize (and in another XPRESSyourself
64 module, quantify) the usage of micro-uORFs is important in exploring regulatory events or mechanisms in a wide
65 array of biological responses and diseases. XPRESSyourself also introduces a tool for efficient identification of
66 the most problematic rRNA fragments for targeted depletion, which provides immense financial and experimental
67 benefits to the user by amplifying ribosome footprint signal over rRNA noise. Tools like this will become vital as
68 ribosome profiling moves into development in new organisms.

69 XPRESSyourself aims to address the lack of consensus in analytical approaches used to process ribosome
70 profiling data by acting as a reference implementation of current best practices for ribosome profiling analysis. While

71 a basic bioinformatic understanding is becoming more commonplace amongst the scientific community, the intricacies
72 of processing RNA-Seq data remain challenging for many. Moreover, many users are often not aware of the most
73 up-to-date tools or the appropriate settings for their application [23, 24]. Even for the experienced user, developing
74 robust automated pipelines that accurately process and assess the quality of these datasets can be laborious. The
75 variability that inevitably arises with each lab or core facility designing and using distinct pipelines is also a challenge
76 to reproducibility in the field. XPRESSyourself curates the state-of-the-art methods for use and where a required
77 functionality is unavailable, introduces a thoroughly tested module to fill that gap. While some tasks in these pipelines
78 may be considered mundane, we eliminate the need of each user to rewrite even simple functionality and promote
79 reproducibility between implementations. To aid users of any skill-level in using this toolkit, we provide thorough
80 documentation, walkthrough videos, and interactive command builders to make usage as easy as possible, while
81 allowing for broad use of this toolkit from personal computers to high-performance clusters.

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

89

90 **Design and Implementation**

91 **Architecture and Organization**

92 XPRESSyourself is currently partitioned into two software packages, XPRESSpipe and XPRESSplot. XPRESSpipe
93 contains automated, end-to-end pipelines tailored for ribosome profiling, single-end RNA-Seq, and paired-end RNA-Seq
94 datasets. Figure 1 outlines the tasks completed by these pipelines. Individual sub-modules can be run automatically
95 through a pipeline or manually step-by-step. Modules optimize available computational resources where appropriate
96 to deliver results as quickly as possible. XPRESSplot is available as a Python library and provides an array of
97 analytical methods specifically for sequence data, but tractable to other data types. For a comparison of how
98 XPRESSyourself compares to other available software packages available at the time of writing, we refer the reader

99 to Figure S1 [14, 15, 18, 25–49].

100 To make analysis as easy and accessible as possible, an integrated command builder for reference curation
101 and sample analysis can be run by executing `xpresspipe build`. This command builder will walk the user through
102 potential considerations based on their library preparation method and build the appropriate command for execution
103 on their personal computer or a supercomputing cluster. The builder will then output the requested command for use
104 on a computational cluster, or the command can be executed immediately on a personal computer.

105 The software is designed such that updating and testing of a new module, or updating dependency usage,
106 are facile tasks for a trained bioinformatician. More details on current and future capabilities can be found in each
107 package's documentation [50, 51] or their respective `versions` page on each toolkit's repository page [52].

108

109 **Automated Reference Curation**

110 The first step of RNA-Seq alignment is curating an organism reference to which the alignment software will map
111 sequence reads. XPRESSpipe uses STAR [53] for mapping reads as it has been shown consistently to be the best
112 performing RNA-Seq read aligner for the majority of cases [54, 55]. The appropriate reference files are automatically
113 curated by providing the appropriate GTF file saved as `transcripts.gtf` and the directory path to the genomic
114 FASTA file(s). Additional modifications to the GTF file required for ribosome profiling or desired for RNA-Seq are
115 discussed in the next section. We recommend organizing these files in their own directory per organism.

116

117 **GTF Modification**

118 For ribosome profiling, frequent read pile-ups are observed at the 5'- and 3'- ends of an open reading frame which are
119 largely uninformative to a gene's translational efficiency [4]. While these pile-ups can be indicative of true translation
120 dynamics [56], current best-practices have more recently settled on ignoring these regions during read quantification
121 and calculations of translation efficiency [3, 21]. By providing the `--truncate` argument during reference curation,
122 the 5'- and 3'- ends of each coding region will be recursively trimmed until the specified amounts are removed from
123 coding space. A recursive strategy is required here as GTF file-formats split the CDS record into regions separated
124 by introns. By default, 45 nt will be trimmed from the 5'-ends and 15 nt from the 3'-ends recursively until the full
125 length is removed from coding space, as is the current convention within the ribosome profiling field [3]. The resulting
126 output file will then be used to process ribosome footprint libraries and their corresponding bulk RNA-Seq libraries. If

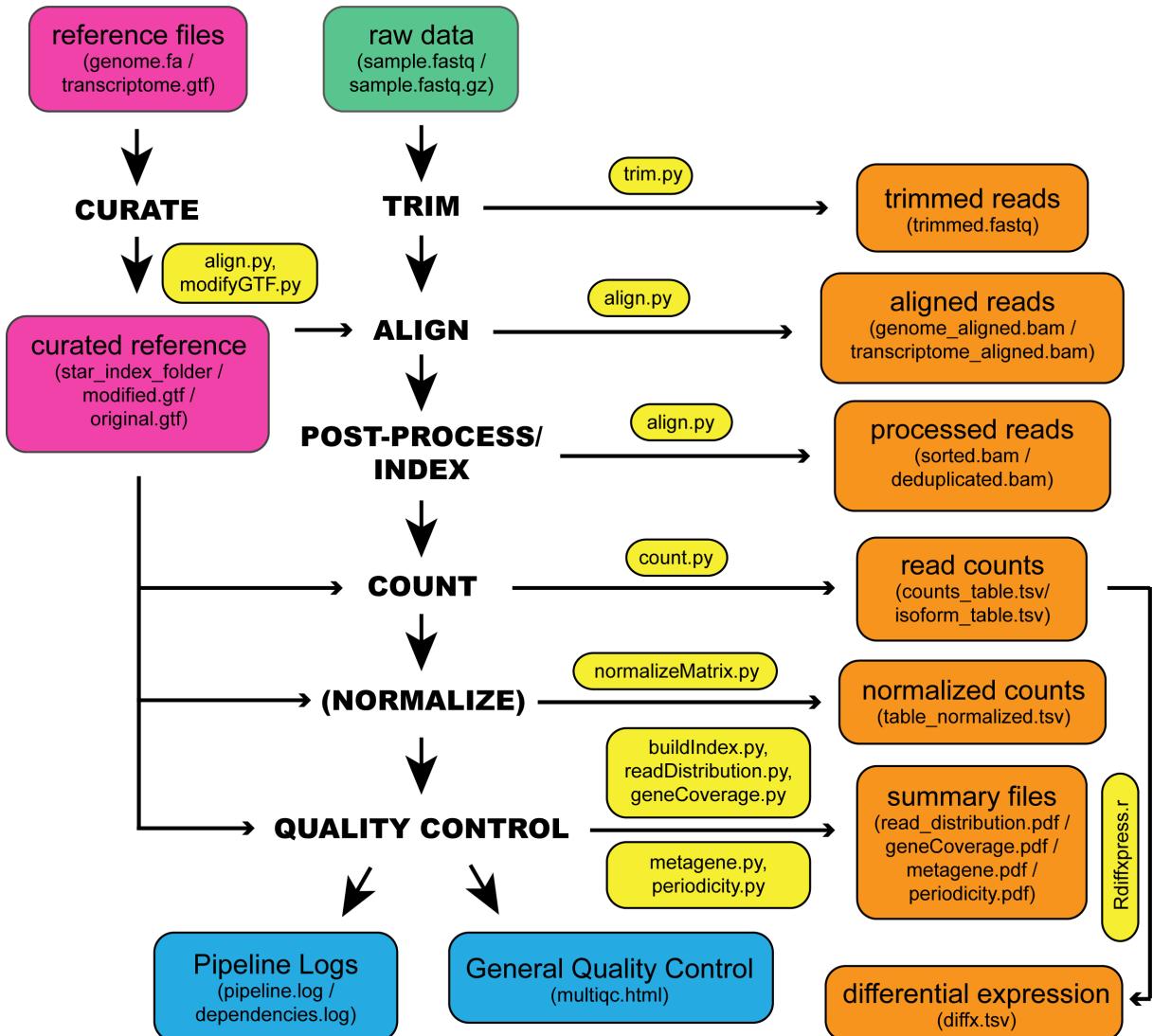


Figure 1: **Workflow schematic of the inputs, outputs, and organization of XPRESSpipe.** Representation of the general steps performed by XPRESSpipe with data and log outputs. Steps in parentheses are optional to the user. Input and output file types are in parentheses for each input or output block. The 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. Differential expression analysis is run independently from the pipeline as the user will need to ensure count table and metadata table formatting are correct before use.

127 generating a GTF file for use solely with general bulk RNA-Seq datasets, this file should not be truncated.

128 Optionally, the GTF can be parsed to retain only protein-coding genes records. This acts as a read masking
129 step to exclude non-protein coding transcripts. In particular, overabundant ribosomal RNAs resulting during library
130 preparation are excluded from downstream analyses using this modified reference file. Parameters can also be
131 provided to retain only the Ensembl canonical transcript record. This can be useful for some tools that penalize reads
132 that overlap multiple isoforms of the same gene. If using HTSeq with default XPRESSpipe parameters or Cufflinks to
133 quantify reads, this is not necessary as they do not penalize a read mapping to multiple isoforms of the same gene
134 or are capable of handling quantification of different isoforms of a gene [57, 58].

135

136 **Read Processing**

137 **Pre-Processing.** In order for sequence reads to be mapped to the genome, reads generally need to be cleaned of
138 artifacts from library creation. These include adaptors, unique molecular identifier (UMI) sequences, and technical
139 errors in the form of low-quality base calls. Parameters, like minimum acceptable quality or length, can be modified,
140 or features such as unique molecular identifiers (UMIs) can be specified to identify and group PCR artifacts for later
141 removal [59, 60].

142 **Alignment.** Reads are aligned to the reference genome with STAR, which, despite being more memory-intensive, is
143 one of the fastest and most accurate sequence alignment options currently available [53–55]. XPRESSpipe is capable
144 of performing a single-pass, splice-aware, GTF-guided alignment or a two-pass alignment of reads wherein novel
145 splice junctions are determined and built into the genome index, followed by alignment of reads using the updated
146 index. Both coordinate- and transcriptome-aligned BAM files are output by STAR. We abstain from rRNA negative
147 alignment at this step as downstream analysis of these mapped reads could be of interest to some users. When rRNA
148 alignment is preferred, a protein-coding-only GTF file should be provided during quantification. A STAR-compatible
149 VCF file can also be passed to this step to allow for genomic variant consideration during alignment.

150 **Post-Processing.** XPRESSpipe further processes alignment files by sorting, indexing, and optionally parsing unique
151 alignments based on UMIs for downstream analyses. PCR duplicates are also detected and marked or removed for
152 downstream analyses; however, these files are only used for relevant downstream steps or if the user specifies to use
153 these de-duplicated files in all downstream steps. Use of de-duplicated alignment files may be advisable in situations
154 where the library complexity profiles (discussed below) exhibit high duplication frequencies. However, generally the

155 abundance of PCR-duplicates is low in properly-prepared sequencing libraries; thus, doing so may be overly stringent
156 and unnecessary [59]. Optionally, BED coverage files can also be output.

157 **Quantification.** XPRESSpipe quantifies read alignments for each input file using HTSeq with the
158 intersection-nonempty method by default [57, 61]. We use this quantification method as it conforms to the current
159 TCGA processing standards and is favorable in the majority of applications [62]. If masking of non-coding RNAs
160 or quantification to truncated CDS records is desired, a protein_coding modified GTF file should be provided to
161 the --gtf argument. HTSeq importantly allows selection of feature type across which to quantify, thus allowing for
162 quantification across the CDSs of a transcript instead of entire exons. If a user is interested in quantifying ribosome
163 occupancy of transcript uORFs in ribosome footprint samples, they can provide five_prime_utr or three_prime_utr
164 for the --feature_type parameter if such annotations exist in the organism of interest's GTF file. If the user is
165 interested in isoform abundance estimation of reads, Cufflinks is alternatively available for quantification [58, 61].

166 **Normalization.** Methods for count normalization are available within XPRESSpipe by way of the XPRESSplot
167 package. For normalizations correcting for transcript length, the appropriate GTF must be provided. Sample
168 normalization methods available include reads-per-million (RPM), Reads-per-kilobase-million (RPKM) or
169 Fragments-per-kilobase-million (FPKM), and transcripts per million (TPM) normalization [63]. For samples sequenced
170 on different chips, prepared by different individuals, or on different days, the --batch argument should be provided
171 along with the appropriate metadata matrix [64].

172

173 **Quality Control**

174 **Read Length Distribution.** The lengths of all reads are analyzed after trimming. By assessing the read distribution
175 of each sample, the user can ensure the expected read size was sequenced. This is particularly helpful in ribosome
176 profiling experiments for verifying the requisite 17-33 nt ribosome footprints were selectively captured during library
177 preparation [3, 65]. Metrics here, as in all other quality control sub-modules, are compiled into summary figures for
178 easy pan-sample assessment by the user.

179 **Library Complexity.** Measuring library complexity is an effective method for analyzing the robustness of a sequencing
180 experiment in capturing various, unique RNA species. As the majority of RNA-Seq preparation methods involve a
181 PCR step, sometimes particular fragments will be favored and over-amplified in contrast to others. By plotting the
182 number of PCR replicates versus expression level for each gene, one can monitor any effects of limited transcript

183 capture diversity and high estimated PCR duplication rate on the robustness of their libraries. This analysis is
184 performed using dupRadar [66] using the duplicate-tagged alignment files output during post-processing. Metrics
185 are then compiled and plotted by XPRESSpipe.

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

192 **Gene Coverage Profile.** Extending the metagene estimation analysis, the user can focus on the coverage profile
193 across a single gene. Although traditional tools like IGV [22] offer the ability to perform such tasks, XPRESSpipe
194 offers the ability to collapse the introns to observe coverage over exon space only. This is helpful in situations where
195 massive introns spread out exons and make it difficult to visualize exon coverage for the entire transcript in a concise
196 manner. CDS feature annotations are displayed to aid ribosome profiling data users in identifying CDS coverage
197 and uORF translation events. When running a XPRESSpipe pipeline, a housekeeping gene will be automatically
198 processed and output for the user's reference. Figure S2 provides a comparison with the output of IGV [22] and
199 XPRESSpipe's geneCoverage module over a similar region for two genes to demonstrate the compatibility between
200 the methods. We note that while the published superTranscripts tool offers similar functionality, it lacks integration and
201 automation and must be manually paired with IGV for multi-sample comparisons and visualization [31]. Other tools,
202 such as Rfeet and riboStreamR [25, 47], suffer from similar integration and automation shortcomings. XPRESSpipe's
203 geneCoverage module offers easy and automated functionality for this task.

204 **Codon Phasing/Periodicity Estimation Profile.** In ribosome profiling, a useful measure of a successful experiment
205 is obtained by investigating the codon phasing of ribosome footprints [3]. To do so, the P-site positions relative to the
206 start codon of each mapped ribosome footprint are calculated using riboWaltz [67]. The same inputs are required as
207 in the metagene sub-module.

208 **Identify Problematic rRNA Fragments from Ribosome Footprinting for Depletion.** rRNA depletion is intrinsically
209 complicated during the preparation of ribosome-footprint profiling libraries: poly(A) selection is irrelevant, and kit-based
210 rRNA depletion is grossly insufficient. Especially in the case of ribosome profiling experiments, where RNA is
211 digested by an RNase to create ribosome footprints, many commercial depletion kits will not target the most abundant

212 rRNA fragment species produces during the footprinting step of ribosome profiling. The sequencing of these RNAs
213 becomes highly repetitive, wasteful, and typically biologically uninteresting in the context of gene expression and
214 translation efficiency. The depletion of these sequences is therefore desired to increase the depth of coverage of
215 ribosome footprints. Depending on the species and condition being profiled, custom rRNA-depletion probes for a
216 small subset of rRNA fragments (generally 2-5) can easily account for more than 90% of sequenced reads [1, 3].
217 `rrnaProbe` analyzes the over-represented sequences within a collection of footprint sequence files that have already
218 undergone adaptor and quality trimming, compiles conserved sequences across the overall experiment, and outputs
219 a rank-ordered list of these sequences for probe design.

220

221 **Analysis**

222 XPRESSpipe provides a DESeq2 command-line wrapper for performing differential expression analysis of count data.
223 We refer users to the original publication for more information about uses and methodology [68].

224 More analytical features are available in XPRESSplot, which requires as input a gene count table as output by
225 XPRESSpipe and a meta-sample table (explained in the documentation [51]). Analyses with limited to no options in
226 Python libraries include principle components plotting with confidence intervals and automated volcano plot creation
227 for RNA-Seq or other data. Other instances of analyses can be found in the documentation [51].

228

229 **Results**

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

231 The integrated stress response (ISR) is a signaling mechanism used by cells and organisms in response to a variety
232 of cellular stresses [69]. Although acute ISR activation is essential for cells to properly respond to stresses, long
233 periods of sustained ISR activity can be damaging. These prolonged episodes contribute to a variety of diseases,
234 including many resulting in neurological decline [70]. A recently discovered small-molecule inhibitor of the ISR, ISRib,
235 has been demonstrated to be a potentially safe and effective neuroprotective therapeutic for traumatic brain injury and
236 other neurological diseases. Interestingly, ISRib can suppress the damaging chronic low activation of the ISR, while it
237 does not interfere with a cytoprotective acute, high-grade ISR, adding to its wide pharmacological interest [9, 71–76].

238 A recent study (data available under Gene Expression Omnibus accession number GSE65778) utilized ribosome
239 profiling to better define the mechanisms of ISRib action on the ISR, modeled by 1-hour tunicamycin (Tm) treatment

240 in HEK293T cells [9]. A key finding of this study is that a specific subset of stress-related transcription factor mRNAs
241 exhibits increased translational efficiency (TE) compared to untreated cells during the tunicamycin-induced ISR.
242 However, when cells were co-treated with tunicamycin and ISRIB, the TE of these stress-related mRNAs showed no
243 significant increase compared to untreated cells, which indicates that ISRIB can counteract the translational regulation
244 associated with the ISR.

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

252 Compared to the count data made available with the original manuscript, when XPRESSpipe-processed samples
253 were thresholded as in the original published count data, samples showed generally comparable read counts per
254 gene between the two analytical regimes (Spearman ρ values 0.937-0.951) (Figure 2B shows representative plots;
255 Figure S3B shows all comparisons). This is in spite of the fact that the methods section of the original publication
256 employed software that was current at the time but is now outdated, such as TopHat2 [77], which has a documented
257 higher false-positive alignment rate, generally lower recall, and lower precision at correctly aligning multi-mapping
258 reads compared to STAR [53–55]. Many of the genes over-represented in the original count data as compared to
259 data processed by XPRESSpipe appear to be due to the over-estimation of pseudogenes or other gene paralogs.
260 Figure S4A highlights a sampling of some extreme cases where particular genes with paralogs are consistently
261 over-represented between samples in the original processed data. This suggests a programmatic difference in how
262 these transcripts are being treated. As these genes share high sequence similarity with each other, reads mapping
263 to these regions are difficult to attribute to a specific genomic locus and are often excluded from further analyses due
264 to their multi-mapping nature. The benchmarking study [54] that examined these and other aligners described how
265 TopHat2 had a disproportionately high rate of incorrectly aligned bases or bases that were aligned uniquely when they
266 should have been aligned ambiguously, at least partially explaining the observed overcounting effect with TopHat2.
267 Had TopHat2 marked problematic reads as ambiguous, they would have been excluded from later quantification.

268 Additionally, when TopHat2 and STAR were tested using multi-mapper simulated test data of varying complexity,

269 TopHat2 consistently suffered in precision and recall. These calls are increasingly more difficult to make with smaller
270 reads as well, and this is evident from Figure 2B, where ribosome footprint samples consistently showed more
271 over-counted genes than the corresponding RNA-Seq samples. When dealing with a ribosome footprint library of
272 about 50-100 million reads, and with TopHat2's simulated likelihood of not marking an ambiguous read as such
273 being about 0.5% higher than STAR, this would lead to around 250,000 to 500,000 spuriously aligned reads, which
274 is in line with our observations (statistics were derived from [54]; analyses are available in the manuscript scripts
275 repository [78]).

276 Another potential contributor to this divergence is that the alignment and quantification within XPRESSpipe use
277 a current human transcriptome reference, which no doubt contains updates and modifications to annotated canonical
278 transcripts and so forth when compared to the version used in the original study. However, in practice, these effects
279 are modest for this dataset (Figure S5). Additionally, the usage of the now outdated DESeq1 [79] appears to contribute
280 significantly to the outcome in differential expression analysis (Figure S6). While differences in processing between
281 the outdated and current methods may not always create systematic differences in output, key biological insights
282 may be missed. The analysis that follows is exploratory and only meant to suggest putative targets identifiable by
283 re-analyzing pre-existing, publicly available data.

284 We first looked at the canonical targets of translation regulation during ISR, as identified in the original study
285 within the XPRESSpipe-processed data. These targets include ATF4, ATF5, PPP1R15A, and DDIT3 (Figure 3A-C,
286 highlighted in purple) [9]. Of note, the fold-change in ribosome occupancy of ATF4 (6.83) from XPRESSpipe-processed
287 samples closely mirrored the estimate from the original publication (6.44). Other targets highlighted in the original
288 study [9], such as ATF5, PPP1R15A, and DDIT3 also demonstrated comparable increases in their ribosome occupancy
289 fold-changes to the original publication count data (XPRESSpipe: 5.90, 2.47, and 3.94; respectively. Original: 7.50,
290 2.70, and 3.89; respectively) (Figure 3A). Similar to the originally processed data, all of these notable changes in
291 ribosome occupancy return to untreated levels during Tm + ISRB co-treatment (Figure 3B). Additional ISR targets
292 containing micro-ORFs described in the study (highlighted in green in Figure 3A-C) were also similar in translational
293 and transcriptional regulation across conditions between the two analytical regimes.

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

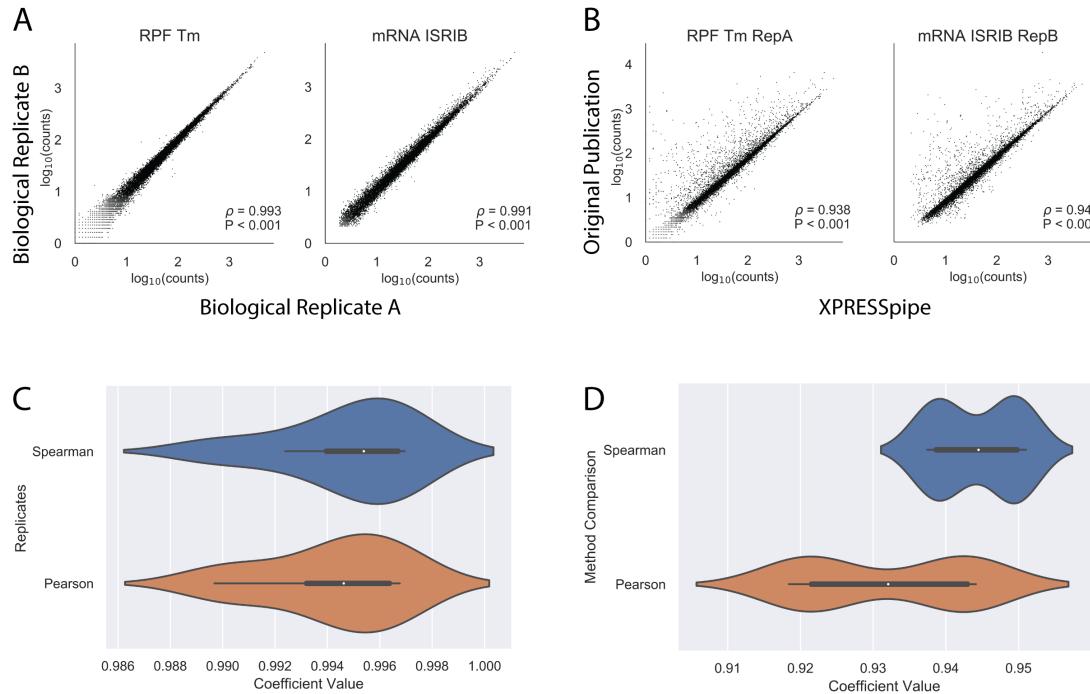


Figure 2: Representative comparisons between processed data produced by XPRESSpipe and original study. Genes were eliminated from analysis if any RNA-Seq sample for that gene had fewer than 10 counts. A) Representative comparisons of biological replicate read counts processed by XPRESSpipe. B) Representative 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 ρ and Pearson r values for biological replicate comparisons. D) Boxplot summaries of Spearman ρ and Pearson r values for between method processing. RPF, ribosome-protected fragments. Tm, tunicamycin. All ρ 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.

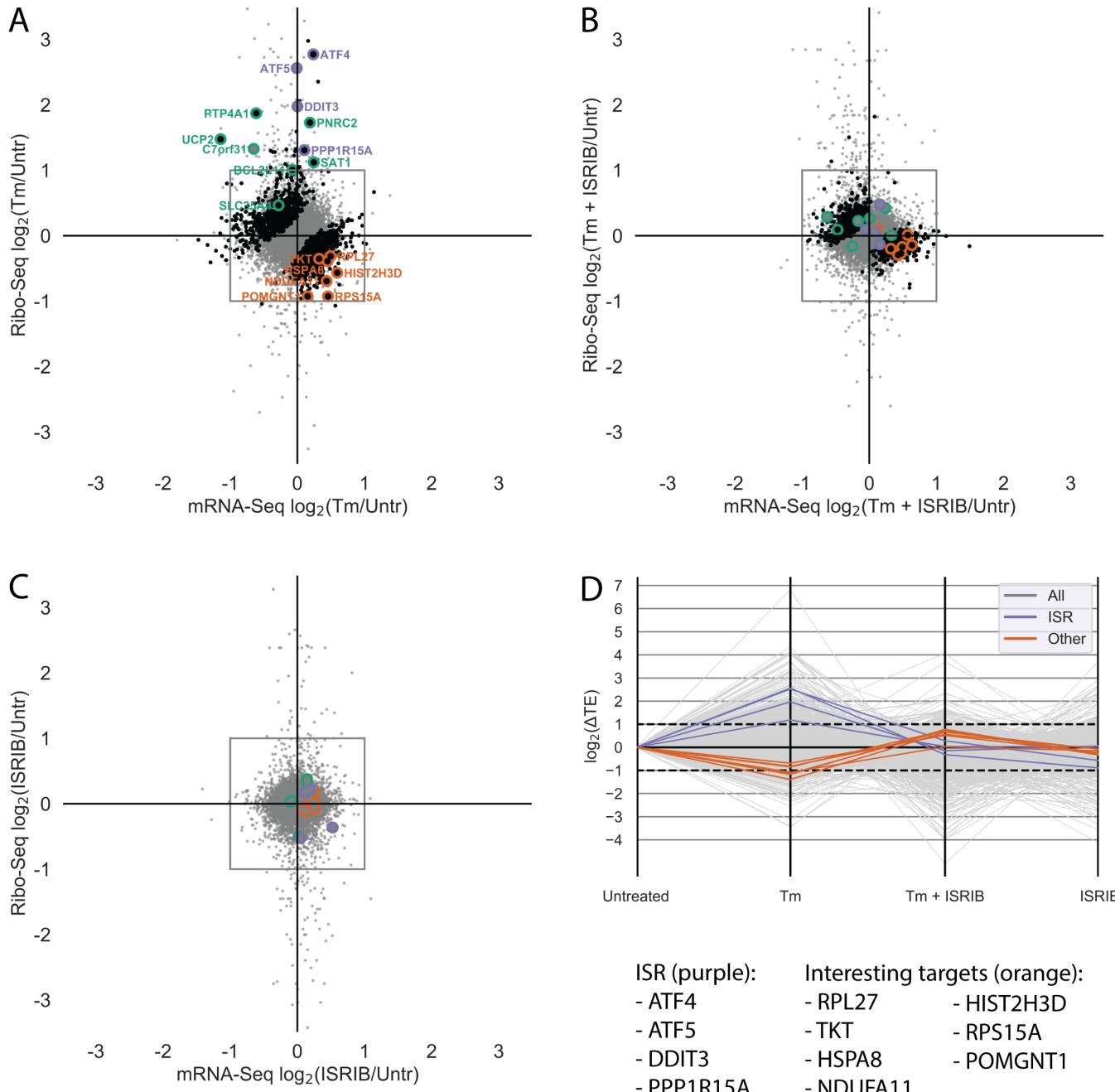


Figure 3: Analysis of previously published ISR TE data using XPRESSpipe. A-C) $\log_2(\text{Fold Change})$ for each drug condition compared to untreated for the ribo-seq and RNA-Seq data. Purple, ISR canonical targets highlighted in the original study. Green, genes with uORFs affected by ISR as highlighted in the original study. Orange, genes fitting a strict TE 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 $\log_2(\Delta TE)$ for each drug condition compared to untreated control. Grey, all genes. Purple, ISR targets identified in the original study. Orange, genes fitting a strict TE 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.

298 focused on in the study. However, re-analysis of these data with the updated XPRESSpipe methodology identifies
299 genes with apparent translational down-regulation that may play a role in the neurodegenerative effects of ISR and
300 the neuroprotective properties of ISRB [73–76]. Importantly, several of these genes were not identified as having
301 significantly down-regulated TEs in the original analysis, which suggests a rationale for not focusing on translational
302 downregulation. In all, we identified seven genes with the regulatory paradigm of interest: significant decreases in
303 TE during tunicamycin-induced ISR that are restored in the ISR + ISRB condition (Table 1, descriptions sourced
304 from [80–82] (Figure 3D)). RNA-Seq and ribosome-footprint coverage across these genes show that the significant
305 changes in their TE are due to neither spurious, high-abundance fragments differentially present across libraries nor
306 variance from an especially small number of mapped reads (Figure S7). This is an important consideration as the
307 commonly suggested use of the CircLigase enzyme in published ribosome profiling library preparation protocols,
308 which circularizes template cDNA before sequencing, can bias certain molecules' incorporation into sequencing
309 libraries based on read-end base content alone [83].

310 Five (POMGNT1, RPL27, TKT, HSPA8, NDUFA11) out of the seven identified genes have annotated neurological
311 functions or mutations that cause severe neurological disorders. Mutations in one other gene (RPS15A) generally
312 result in metabolic disorders. While none of these genes were identified as being of interest in the original study using
313 the original methods, by re-processing the original manuscript count data with DESeq2 [68] and the same expression
314 pattern thresholding, four of these genes are now present in the analysis (RPL27, TKT, HSPA8, RPS15A) (see Figure
315 S6 for a systematic comparison). NDUFA11 and TKT are protein-coding genes whose functions are integrally tied
316 to successful central carbon metabolism and mitochondrial electron transport chain function, respectively. NDUFA11
317 encodes a subunit of mitochondrial respiratory complex I [84], and TKT encodes a thiamine-dependent enzyme
318 that channels excess sugars phosphates into glycolysis as part of the pentose phosphate pathway [85]. Mutations
319 in NDUFA11 cause severe neurodegenerative phenotypes such as brain atrophy and encephalopathy [84], and
320 mutations in TKT cause diseases associated with neurological phenotypes [86]. These regulatory and phenotypic
321 observations raise the possibility that their role may be functionally relevant to the neurodegenerative effects of ISR
322 and the neuroprotective properties of ISRB.

323 While at this stage speculative, it is interesting that the processing of these data with updated methods provides
324 a very conservative list of differentially expressed genes, and that the majority of those genes are associated with
325 severe neurological phenotypes. It is therefore easy to speculate that TE regulation of these targets' abundance
326 might be important in the neurodegeneration observed in prolonged ISR conditions. ISRB's neuroprotective effects

Table 1: **Translatonally down-regulated genes during acute Tm treatment with restored regulation during Tm + ISRB treatment.** Gene names succeeded by an asterisk indicate these genes were identified in the original data when re-analyzed with DESeq2 [68]. Gene names succeeded by an ampersand indicate genes with strong neurological phenotype annotations. None of these genes were present in the original analysis tables.

Gene Name	Relevant Description
POMGNT1 ^{&}	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 ^{*&}	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 ^{*&}	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 ^{*&}	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 ^{&}	Encodes subunit of mitochondrial complex I, a vital component of the electron transport chain. Mutations are 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 nucleosome structure. No neurological phenotypes currently annotated.
RPS15A*	Subunit of ribosome catalyzing protein synthesis. Diseases associated include Diamond-Blackfan Anemia, an inborn error of metabolism disease.

327 may stem from a restoration of one or more of these entities' protein expression. Though speculative without further
 328 validation, these ISRB-responsive neuronal targets act as interesting cases for further validation and study in a
 329 model more representative of neurotoxic injury and disease than the HEK-293T model used in the original study.
 330 In all, this comparison demonstrates the utility of XPRESSpipe for rapid, user-friendly analysis and re-analysis of
 331 ribosome-profiling experiments in the pursuit of biological insights and hypothesis generation.

332

333 Cost Analysis and Performance

334 XPRESSpipe functions can be computationally intensive. Super-computing resources are recommended, especially
 335 when handling large datasets or when aligning to larger, more complex genomes. Many universities provide super-computing
 336 resources to their affiliates; however, in cases where these resources are not available, servers such as Amazon Web
 337 Services (AWS) [87] can be used to process sequencing data using XPRESSpipe. Table 2 summarizes the runtime
 338 statistics for the ISRB dataset used in this study. The ISRB ribosome profiling dataset contained a total of 32
 339 raw sequence files that were aligned to *Homo sapiens*; thus it acts as a high-end estimate of the time required to
 340 process data with XPRESSpipe. For a comparable dataset, cost to use an AWS computational node with similar

341 specifications for the above elapsed time would be approximately 44.39 USD using an Amazon EC2 On-Demand
 342 m5.12xlarge node (however, significantly reduced rates are available if using Spot instances or by using the free
 343 tier; also, while relatively slower, one could use a m5.4xlarge node at about a third the rate as the m5.12xlarge) and
 344 storage cost would amount to around 17.41 USD/month for all input and output data on Amazon S3 storage (storage
 345 costs could be reduced as much of the intermediate data may not be relevant for certain users; however, raw input
 346 data should always be archived by the user).

347

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. Input raw FASTQ files were uncompressed. The metagene and geneCoverage sub-modules used a conservative BAM file multiprocessing limit to avoid out-of-memory errors.

Process	Command	Wallclock Time	Max RAM
Curate STAR reference	curateReference	00h 34m 01s	34.03 GB
Truncate GTF	modifyGTF -t	00h 02m 45s	03.25 GB
Read Pre-processing	trim	00h 18m 21s	00.48 GB
Alignment and Post-processing	align	06h 16m 18s	38.00 GB
Read Quantification	count -c htseq	04h 13m 04s	00.16 GB
Isoform Abundance	count -c cufflinks	01h 09m 32s	02.36 GB
Differential Expression (n=9)	diffxpress	00h 07m 50s	00.65 GB
Read Distributions	readDistribution	00h 06m 19s	00.28 GB
Metagene Analysis	metagene	01h 34m 09s	35.52 GB
Gene Coverage (n=1)	geneCoverage	03h 49m 03s	19.22 GB
Periodicity	periodicity	01h 00m 57s	61.89 GB
Library Complexity	complexity	00h 48m 46s	01.98 GB
rRNA probe	rrnaProbe	00h 00m 55s	00.15 GB
Pipeline	riboseq	19h 16m 08s	61.89 GB

Attribute	Value
Total Raw Input	257 GB
Total Output	500 GB
Allocated CPUs	20
Allocated Memory	64 GB

348 Availability and Future Directions

349 We have described the software suite, XPRESSyourself, an automated reference implementation of best-practices
 350 in ribosome profiling data analysis built upon a synthesis of new tools, old tools, and pipelines. XPRESSyourself is
 351 perpetually open source and protected under the GPL-3.0 license. Updates to the software are version controlled
 352 and maintained on GitHub [52]. Jupyter notebooks and video walkthroughs are included within the README files
 353 at [52]. Documentation is hosted on readthedocs [88] at [50] and [51]. Source code for associated analyses and

354 figures for this manuscript can be accessed at [78]. The data used in this manuscript are available under the Gene
355 Expression Omnibus persistent identifier GSE65778 [89] for ribosome profiling data and under the dbGaP Study
356 Accession persistent identifier phs000178 [90] for the TCGA data.

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

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

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

378 With the adoption of this flexible pipeline, the field of high-throughput sequencing, particularly ribosome profiling,
379 can continue to standardize the processing protocol for associated sequence data and eliminate the variability that
380 comes from the availability of a variety of software packages for various steps during sequence read processing.
381 Additionally, XPRESSpipe consolidates various tools used by the ribosome profiling and RNA-Seq communities
382 into a single, end-to-end pipeline. With these tools, genome reference formatting and curation are automated and

Table 3: **Software Description**

Project Name	XPRESSyourself
Project Home Page	https://github.com/XPRESSyourself
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

383 accessible to the public. Adoption of this tool will allow scientists to process and access their data independently
 384 and quickly, guide them in understanding key considerations in processing their data, and standardize protocols for
 385 ribosome profiling and other RNA-Seq applications, thus increasing reproducibility in sequencing analyses.

386

387 **List of Abbreviations**

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

398

399 **Ethics Approval and Consent to Participate**

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

402

403 **Consent for Publication**

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

406

407 **Competing Interests**

408 The authors declare that they have no competing interests.

409

410 **Funding**

411 J.A.B. received support from the National Institute of Diabetes and Digestive and Kidney Diseases (NIDDK) Inter-disciplinary
412 Training Grant T32 Program in Computational Approaches to Diabetes and Metabolism Research, 1T32DK11096601
413 to Wendy W. Chapman and Simon J. Fisher. J.T.M. received support as an HHMI Fellow of the Jane Coffin
414 Childs Memorial Fund for Medical Research. A.J.B received support from the National Cancer Institute (NCI)
415 Predoctoral to Postdoctoral Fellow Transition Award, K00CA212445. This work was supported by NIDDK fellowship
416 1T32DK11096601 (to J.A.B.) and NIH grant R35GM13185 (to J.R.). The computational resources used were partially
417 funded by the NIH Shared Instrumentation Grant 1S10OD021644-01A1.

418

419 **Contributions**

420 J.A.B. conceptualized and administered the project; performed all investigation, analysis, visualization, and data
421 curation; provisioned resources; and wrote the original draft of this manuscript. J.A.B., J.R.B., J.T.M., and J.G. and
422 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.
423 performed software validation. J.A.B. and J.R. acquired funding. J.R., A.R.Q., and J.G. supervised the study. All
424 authors reviewed and edited this manuscript.

425

426 **Acknowledgments**

427 The authors wish to thank Michael T. Howard for helpful discussions concerning ribosome profiling and sequencing
428 analysis. The authors also wish to thank Mark E. Wadsworth, Ryan Miller, and Michael J. Cormier for helpful

Table 4: Author ORCIDs

Author	ORCID
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

429 discussions on pipeline design. They also wish to thank T. Cameron Waller for helpful discussions related to pipeline
 430 design and biological analysis. The support and resources from the Center for High-Performance Computing at the
 431 University of Utah are gratefully acknowledged. The results published here are in whole or part based upon data
 432 generated by the TCGA Research Network [62].

433

434 References

- 435 [1] N. Ingolia, S. Ghaemmaghami, J. Newman, J. Weissman. Genome-wide analysis in vivo of
 436 translation with nucleotide resolution using ribosome profiling. *Science* **324**, 218 (2009). Available from:
 437 <https://doi.org/10.1126/science.1168978>.
- 438 [2] G. Brar, J. Weissman. Ribosome profiling reveals the what, when, where and how of protein synthesis. *Nat Rev
 439 Mol Cell Biol* **16**, 651 (2015). Available from: <https://doi.org/10.1038/nrm4069>.
- 440 [3] N. McGlincy, N. Ingolia. Transcriptome-wide measurement of translation by ribosome profiling. *Methods* **126**,
 441 112 (2017). Available from: [https://doi.org/10.1016/j.ymeth.2017.05.028](https://doi.org/10.1016/jymeth.2017.05.028).
- 442 [4] M. Gerashchenko, V. Gladyshev. Translation inhibitors cause abnormalities in ribosome profiling experiments.
 443 *Nucleic Acids Res* **42** (2014). Available from: <https://doi.org/10.1093/nar/gku671>.
- 444 [5] A. Bartholomäus, C. D. Campo, I. Z. Mapping the non-standardized biases of ribosome profiling. *Biol Chem* **397**
 445 (2016). Available from: <https://doi.org/https://doi.org/10.1515/hsz-2015-0197>.
- 446 [6] J. Hussmann, S. Patchett, A. Johnson, S. Sawyer, W. Press. Understanding Biases in Ribosome Profiling
 447 Experiments Reveals Signatures of Translation Dynamics in Yeast. *PLoS Genet* **11** (2015). Available from:
 448 <https://doi.org/https://doi.org/10.1371/journal.pgen.1005732>.

- 449 [7] A. Diamant, T. Tuller. Estimation of ribosome profiling performance and reproducibility at various levels of
450 resolution. *Biol Direct* **11** (2016). Available from: <https://doi.org/https://doi.org/10.1186/s13062-016-0127-4>.
- 451 [8] M. Gerashchenko, V. Gladyshev. Ribonuclease selection for ribosome profiling. *Nucleic Acids Res* **45** (2017).
452 Available from: <https://doi.org/https://doi.org/10.1093/nar/gkw822>.
- 453 [9] C. Sidrauski, A. McGeachy, N. Ingolia, P. Walter. The small molecule ISRIB reverses the effects
454 of eIF2α phosphorylation on translation and stress granule assembly. *eLife* (2015). Available from:
455 <https://doi.org/10.7554/eLife.05033>.
- 456 [10] F. Mohammad, C. Woolstenhulme, R. Green, A. Buskirk. Clarifying the Translational
457 Pausing Landscape in Bacteria by Ribosome Profiling. *Cell Rep* **14** (2016). Available from:
458 <https://doi.org/https://doi.org/10.1016/j.celrep.2015.12.073>.
- 459 [11] G. Li, E. Oh, J. Weissman. The anti-Shine-Dalgarno sequence drives translational pausing and codon choice in
460 bacteria. *Nature* **484** (2012). Available from: <https://doi.org/https://doi.org/10.1038/nature10965>.
- 461 [12] A. Lecanda, *et al.*. Dual randomization of oligonucleotides to reduce the bias in ribosome-profiling libraries.
462 *Methods* **107** (2016). Available from: <https://doi.org/https://doi.org/10.1016/j.ymeth.2016.07.011>.
- 463 [13] X. Gao, *et al.*. Quantitative profiling of initiating ribosomes in vivo. *Nat Methods* **12** (2015). Available from:
464 <https://doi.org/https://doi.org/10.1038/nmeth.3208>.
- 465 [14] E. Afgan, *et al.*. The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2018
466 update. *Nucleic Acids Res.* **46**, W537 (2018). Available from: <https://doi.org/10.1093/nar/gky379>.
- 467 [15] A. Michel, *et al.*. RiboGalaxy: A browser based platform for the alignment, analysis and visualization of ribosome
468 profiling data. *RNA Biol* **13**, 316 (2016). Available from: <https://doi.org/10.1080/15476286.2016.1141862>.
- 469 [16] Nextflow. <https://www.nextflow.io/example4.html>.
- 470 [17] DNAnexus. https://github.com/dnanexus/tophat_cufflinks_rnaseq.
- 471 [18] O. Carja, T. Xing, E. Wallace, J. Plotkin, P. Shah. riboviz: analysis and visualization of ribosome profiling datasets.
472 *BMC Bioinformatics* **18** (2017). Available from: <https://doi.org/10.1186/s12859-017-1873-8>.

- 473 [19] C. Artieri, H. Fraser. Accounting for biases in riboprofiling data indicates a major role for proline in stalling
474 translation. *Genome Res* **24**, 2011 (2014). Available from: <https://doi.org/10.1101/gr.175893.114>.
- 475 [20] J. Hussmann, S. Patchett, A. Johnson, S. Sawyer, W. Press. Understanding Biases in Ribosome Profiling
476 Experiments Reveals Signatures of Translation Dynamics in Yeast. *PLoS Genet* **11** (2015). Available from:
477 <https://doi.org/10.1371/journal.pgen.1005732>.
- 478 [21] D. Weinberg, *et al.*. Improved Ribosome-Footprint and mRNA Measurements Provide Insights
479 into Dynamics and Regulation of Yeast Translation. *Cell Rep* **14**, 1787 (2016). Available from:
480 <https://doi.org/10.1016/j.celrep.2016.01.043>.
- 481 [22] J. Robinson, *et al.*. Integrative Genomics Viewer. *Nat Biotechnol* **29**, 24 (2011). Available from:
482 <https://doi.org/10.1038/nbt.1754>.
- 483 [23] Z. Costello, H. Martin. A machine learning approach to predict metabolic pathway dynamics from time-series
484 multiomics data. *NPJ Syst Biol Appl* **4** (2018). Available from: <https://doi.org/10.1038/s41540-018-0054-3>.
- 485 [24] V. Funari, S. Canosa. The Importance of Bioinformatics in NGS: Breaking the Bottleneck in Data Interpretation.
486 *Science* **344**, 653 (2014). Available from: <https://doi.org/10.1126/science.344.6184.653-c>.
- 487 [25] R. Kumari, A. Michel, P. Baranov. PausePred and Rfeet: webtools for inferring ribosome pauses
488 and visualizing footprint density from ribosome profiling data. *RNA* **24** (2018). Available from:
489 <https://doi.org/10.1261/rna.065235.117>.
- 490 [26] C. Oertlin, *et al.*. Generally applicable transcriptome-wide analysis of translation using anota2seq. *Nucleic Acids
491 Res* **47** (2019). Available from: <https://doi.org/10.1093/nar/gkz223>.
- 492 [27] A. Popa, *et al.*. RiboProfiling: a Bioconductor package for standard Ribo-seq pipeline processing. *F1000Res* **5**
493 (2016). Available from: <https://doi.org/10.12688/f1000research.8964.1>.
- 494 [28] W. Li, W. Wang, P. Uren, L. Penalva, A. Smith. Riborex: fast and flexible identification of differential translation
495 from Ribo-seq data. *Bioinformatics* **33** (2017). Available from: <https://doi.org/10.1093/bioinformatics/btx047>.
- 496 [29] S. Verbruggen, G. Menschaert. mQC: A post-mapping data exploration tool for ribosome profiling. *Comput
497 Methods Programs Biomed* (2018). Available from: <https://doi.org/10.1016/j.cmpb.2018.10.018>.

- 498 [30] Å. Birkeland, K. ChyŻyńska, E. Valen. Shoelaces: an interactive tool for ribosome profiling processing and
499 visualization. *BMC Genomics* **19** (2018). Available from: <https://doi.org/10.1186/s12864-018-4912-6>.
- 500 [31] N. Davidson, A. Hawkins, A. Oshlack. SuperTranscripts: a data driven reference for analysis and visualisation of
501 transcriptomes. *Genome Biol* **18** (2017). Available from: <https://doi.org/10.1186/s13059-017-1284-1>.
- 502 [32] T. Backman, T. Girke. systemPipeR: NGS workflow and report generation environment. *BMC Bioinformatics* **17**
503 (2016). Available from: <https://doi.org/10.1186/s12859-016-1241-0>.
- 504 [33] H. Tjeldnes, K. Labun. ORFik: Open Reading Frames in Genomics. <https://github.com/JokingHero/ORFik>
505 (2017). Available from: <https://doi.org/10.18129/B9.bioc.ORFik>.
- 506 [34] T. Martin, I. Erte, P. Tsai, J. Bell. coMET: an R plotting package to visualize regional plots of epigenome-wide
507 association scan results. *QG14* (2014). Available from: <http://quantgen.soc.srccf.net/qg14/>.
- 508 [35] T. Martin, I. Yet, P. Tsai, J. Bell. coMET: visualisation of regional epigenome-wide association
509 scan results and DNA co-methylation patterns. *BMC Bioinformatics* **16** (2015). Available from:
510 <https://doi.org/10.1186/s12859-015-0568-2>.
- 511 [36] T. Hardcastle. riboSeqR. Available from: <https://doi.org/10.18129/B9.bioc.roboSeqR>.
- 512 [37] F. Ramírez, F. Dündar, S. Diehl, B. Grüning, T. Manke. deepTools: a flexible platform for exploring
513 deep-sequencing data. *Nucleic Acids Res* **42** (2014). Available from: <https://doi.org/10.1093/nar/gku365>.
- 514 [38] Picard. <https://broadinstitute.github.io/picard/>.
- 515 [39] S. Zhang, *et al.*. Analysis of Ribosome Stalling and Translation Elongation Dynamics by Deep Learning. *Cell
516 Syst* **5** (2017). Available from: <https://doi.org/10.1016/j.cels.2017.08.004>.
- 517 [40] P. O'Connor, D. Andreev, P. Baranov. Comparative survey of the relative impact of mRNA features on local
518 ribosome profiling read density. *Nat Commun* **7** (2016). Available from: <https://doi.org/10.1038/ncomms12915>.
- 519 [41] Z. Xiao, Q. Zou, Y. Liu, X. Yang. Genome-wide assessment of differential translations with ribosome profiling
520 data. *Nat Commun* **7** (2016). Available from: <https://doi.org/10.1038/ncomms11194>.
- 521 [42] Y. Zhong, *et al.*. RiboDiff: detecting changes of mRNA translation efficiency from ribosome footprints.
522 *Bioinformatics* **33** (2017). Available from: <https://doi.org/10.1093/bioinformatics/btw585>.

- 523 [43] L. Calviello, *et al.*. Detecting actively translated open reading frames in ribosome profiling data. *Nat Methods* **13**
524 (2016). Available from: <https://doi.org/10.1038/nmeth.3688>.
- 525 [44] H. Wang, J. McManus, C. Kingsford. Isoform-level ribosome occupancy estimation
526 guided by transcript abundance with Ribomap. *Bioinformatics* **32** (2016). Available from:
527 <https://doi.org/10.1093/bioinformatics/btw085>.
- 528 [45] P. Spealman, H. Wang, G. May, C. Kingsford, C. McManus. Exploring Ribosome Positioning
529 on Translating Transcripts with Ribosome Profiling. *Methods Mol Biol* **1358** (2016). Available from:
530 https://doi.org/10.1007/978-1-4939-3067-8_5.
- 531 [46] J. Dunn, J. Weissman. Plastid: nucleotide-resolution analysis of next-generation sequencing and genomics data.
532 *BMC Genomics* **17** (2016). Available from: <https://doi.org/https://doi.org/10.1186/s12864-016-3278-x>.
- 533 [47] P. Perkins, S. Mazzoni-Putman, A. Stepanova, J. Alonso, S. Heber. RiboStreamR: a web application for
534 quality control, analysis, and visualization of Ribo-seq data. *BMC Genomics* **20** (2019). Available from:
535 <https://doi.org/https://doi.org/10.1186/s12864-019-5700-7>.
- 536 [48] H. Fang, *et al.*. Scikit-ribo Enables Accurate Estimation and Robust Modeling of Translation Dynamics at Codon
537 Resolution. *Cell Syst* **6** (2018). Available from: <https://doi.org/https://doi.org/10.1016/j.cels.2017.12.007>.
- 538 [49] S. Chun, C. Rodriguez, P. Todd, R. Mills. SPECtre: a spectral coherence-based classifier of actively
539 translated transcripts from ribosome profiling sequence data. *BMC Bioinformatics* **17** (2016). Available from:
540 <https://doi.org/https://doi.org/10.1186/s12859-016-1355-4>.
- 541 [50] XPRESSpipe documentation. <https://xpresspipe.readthedocs.io/en/latest/>.
- 542 [51] XPRESSplot documentation. <https://xpressplot.readthedocs.io/en/latest/>.
- 543 [52] XPRESSyourself. <https://github.com/XPRESSyourself/>.
- 544 [53] A. Dobin, *et al.*. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15 (2013). Available from:
545 <https://doi.org/10.1093/bioinformatics/bts635>.
- 546 [54] G. Baruzzo, *et al.*. Simulation-based comprehensive benchmarking of RNA-seq aligners. *Nat Methods* **14**, 135
547 (2017). Available from: <https://doi.org/10.1038/nmeth.4106>.

- 548 [55] I. Raplee, A. Evsikov, C. M. de Evsikova. Aligning the Aligners: Comparison of RNA Sequencing Data Alignment
549 and Gene Expression Quantification Tools for Clinical Breast Cancer Research. *J Pers Med* **9** (2019). Available
550 from: <https://doi.org/10.3390/jpm9020018>.
- 551 [56] T. Tuller, H. Zur. Multiple roles of the coding sequence 5' end in gene expression regulation. *Nucleic Acids Res*
552 **42** (2015). Available from: <https://doi.org/10.1093/nar/gku1313>.
- 553 [57] S. Anders, P. Pyl, W. Huber. HTSeq—a Python framework to work with high-throughput sequencing data.
554 *Bioinformatics* **31**, 166 (2015). Available from: <https://doi.org/10.1093/bioinformatics/btu638>.
- 555 [58] C. Trapnell, *et al.*. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and
556 Cufflinks. *Nat Protoc* **7** (2012). Available from: <https://doi.org/10.1038/nprot.2012.016>.
- 557 [59] Y. Fu, P. Wu, T. Beane, P. Zamore, Z. Weng. Elimination of PCR duplicates in RNA-seq
558 and small RNA-seq using unique molecular identifiers. *BMC Genomics* **19** (2018). Available from:
559 <https://doi.org/10.1186/s12864-018-4933-1>.
- 560 [60] T. Smith, A. Heger, I. Sudbery. UMI-tools: Modelling sequencing errors in Unique Molecular Identifiers to improve
561 quantification accuracy. *Genome Res* **27** (2017). Available from: <https://doi.org/10.1101/gr.209601.116>.
- 562 [61] C. Robert, M. Watson. Errors in RNA-Seq quantification affect genes of relevance to human disease. *Genome
563 Biol* **16** (2015). Available from: <https://doi.org/10.1186/s13059-015-0734-x>.
- 564 [62] The Cancer Genome Atlas. <https://portal.gdc.cancer.gov>.
- 565 [63] C. Evans, J. Hardin, D. Stoebel. Selecting between-sample RNA-Seq normalization methods
566 from the perspective of their assumptions. *Brief Bioinform* **19**, 776–792 (2018). Available from:
567 <https://doi.org/10.1093/bib/bbx008>.
- 568 [64] J. Leek, W. Johnson, H. Parker, A. Jaffe, J. Storey. The sva package for removing batch effects
569 and other unwanted variation in high-throughput experiments. *Bioinformatics* **28** (2012). Available from:
570 <https://doi.org/10.1093/bioinformatics/bts034>.
- 571 [65] C. Wu, B. Zinshteyn, K. Wehner, R. Green. High-Resolution Ribosome Profiling Defines Discrete Ribosome
572 Elongation States and Translational Regulation during Cellular Stress. *Mol Cell* **73** (2019). Available from:
573 <https://doi.org/10.1016/j.molcel.2018.12.009>.

- 574 [66] S. Sayols, D. Scherzinger, H. Klein. dupRadar: a Bioconductor package for the assessment of PCR artifacts in
575 RNA-Seq data. *BMC Bioinformatics* **17** (2016). Available from: <https://doi.org/10.1186/s12859-016-1276-2>.
- 576 [67] F. Lauria, *et al.*. riboWaltz: Optimization of ribosome P-site positioning in ribosome profiling data. *PLoS Comput
577 Biol* **14** (2018). Available from: <https://doi.org/10.1371/journal.pcbi.1006169>.
- 578 [68] M. Love, W. Huber, S. Anders. Moderated estimation of fold change and dispersion for RNA-seq data with
579 DESeq2. *Genome Biol* **15** (2014). Available from: <https://doi.org/10.1186/s13059-014-0550-8>.
- 580 [69] H. Harding, *et al.*. An integrated stress response regulates amino acid metabolism and resistance to oxidative
581 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).
- 582 [70] D. Santos-Ribeiro, L. Godinas, C. Pilette, F. Perros. The integrated stress response system in cardiovascular
583 disease. *Drug Discov Today* **23** (2018). Available from: <https://doi.org/10.1016/j.drudis.2018.02.008>.
- 584 [71] H. Rabouw, *et al.*. Small molecule ISRIB suppresses the integrated stress response within a defined window of
585 activation. *Proc Natl Acad Sci U S A* **116** (2019). Available from: <https://doi.org/10.1073/pnas.1815767116>.
- 586 [72] J. Tsai, *et al.*. Structure of the nucleotide exchange factor eIF2B reveals mechanism of memory-enhancing
587 molecule. *Science* **359** (2018). Available from: <https://doi.org/10.1126/science.aaq0939>.
- 588 [73] A. Choua, *et al.*. Inhibition of the integrated stress response reverses cognitive deficits after traumatic brain
589 injury. *Proc Natl Acad Sci U S A* **114** (2017). Available from: <https://doi.org/10.1073/pnas.1707661114>.
- 590 [74] M. Halliday, *et al.*. Partial restoration of protein synthesis rates by the small molecule ISRIB
591 prevents neurodegeneration without pancreatic toxicity. *Cell Death Dis* **6** (2015). Available from:
592 <https://doi.org/10.1038/cddis.2015.49>.
- 593 [75] C. Sidrauski, *et al.*. Pharmacological brake-release of mRNA translation enhances cognitive memory. *eLife* **2**
594 (2013). Available from: <https://doi.org/10.7554/eLife.00498>.
- 595 [76] Y. Sekine, *et al.*. Stress responses. Mutations in a translation initiation factor identify the target of a
596 memory-enhancing compound. *Science* **348** (2015). Available from: <https://doi.org/10.1126/science.aaa6986>.
- 597 [77] D. Kim, *et al.*. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene
598 fusions. *Genome Biol* **14** (2013). Available from: <https://doi.org/10.1186/gb-2013-14-4-r36>.

- 599 [78] Manuscript code. https://github.com/XPRESSyourself/xpressyourself_manuscript/tree/master/
600 supplemental_files. Available from: <https://doi.org/DOI: 10.5281/zenodo.3337599>.
- 601 [79] S. Anders, W. Huber. Differential expression analysis for sequence count data. *Genome Biol* **11** (2010). Available
602 from: <https://doi.org/10.1186/gb-2010-11-10-r106>.
- 603 [80] GeneCards. <https://www.genecards.org/>. Accessed 20 October 2019.
- 604 [81] National Center for Biotechnology Information. <https://www.ncbi.nlm.nih.gov/gene/>. Accessed 20 October
605 2019.
- 606 [82] UniProt. <https://www.uniprot.org/uniprot/>. Accessed 20 October 2019.
- 607 [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>.
- 608 [84] I. Berger, *et al.*. Mitochondrial complex I deficiency caused by a deleterious NDUFA11 mutation. *Ann Neurol* **63**
609 (2008). Available from: <https://doi.org/https://doi.org/10.1002/ana.21332>.
- 610 [85] L. Mitschke, *et al.*. The crystal structure of human transketolase and new insights into its mode of action. *J Biol Chem* **285** (2010). Available from: <https://doi.org/https://doi.org/10.1074/jbc.M110.149955>.
- 611 [86] L. Boyle, *et al.*. The crystal structure of human transketolase and new insights into its mode of action. *Am J Hum Genet* **98** (2016). Available from: <https://doi.org/https://doi.org/10.1016/j.ajhg.2016.03.030>.
- 612 [87] Amazon Web Services. <https://aws.amazon.com>.
- 613 [88] Read the Docs. <https://readthedocs.org/>.
- 614 [89] Ribosome Profiling GEO Accession. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE65778>.
- 615 [90] TCGA dbGaP Accession. (https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000178.v10.p8).

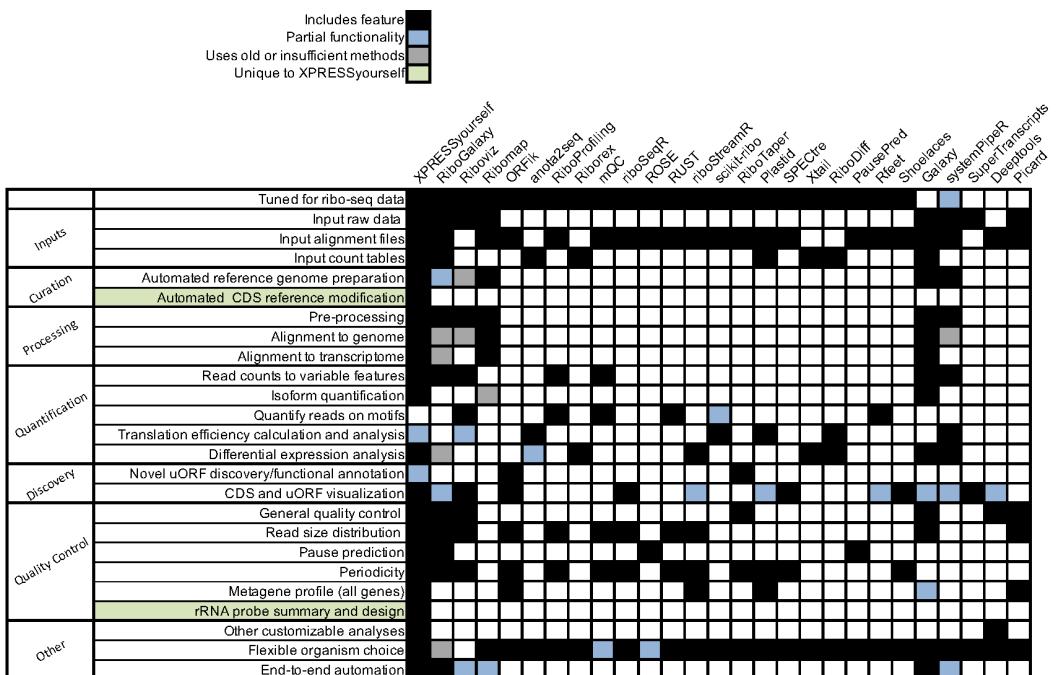


Figure S1: Comparison between XPRESSyourself and other available software packages for ribosome profiling data analysis. Black boxes indicate full functionality, blue boxes indicate partial functionality, grey boxes indicate incomplete or outdated functionality, and blank boxes indicate no functionality for the specified task. Rankings were compiled using the tools' documentation, manuscript, and codebase. If a function was not clearly described in any of these venues, a blank box was given.

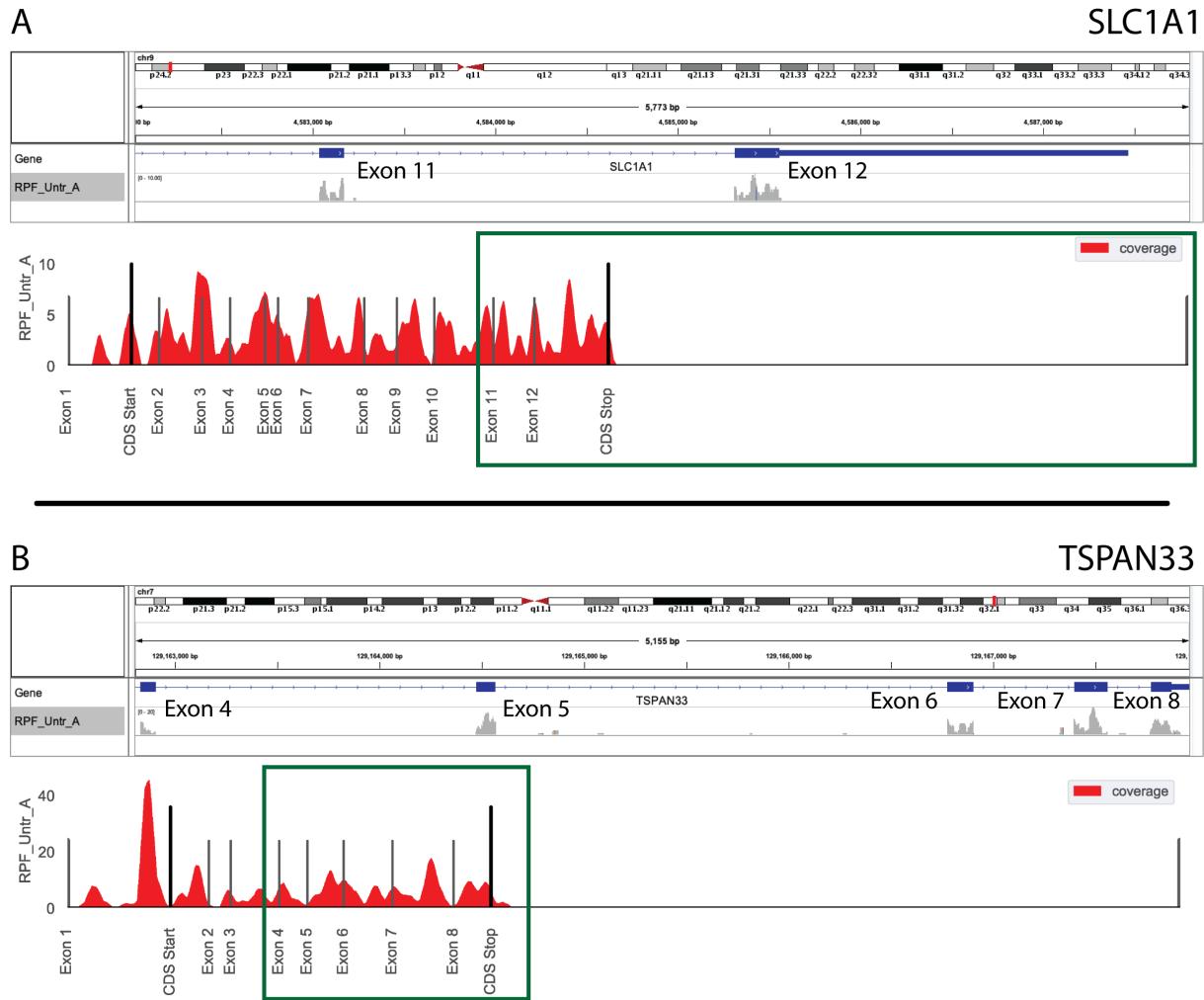


Figure S2: **Comparison between IGV browser and geneCoverage output.** A) Gene coverage from IGV (above) and XPRESSpipe (below) for SLC1A1. B) Gene coverage from IGV (above) and XPRESSpipe (below) for TSPAN33. Introns collapsed by XPRESSpipe. Green box, region displayed in corresponding IGV window.

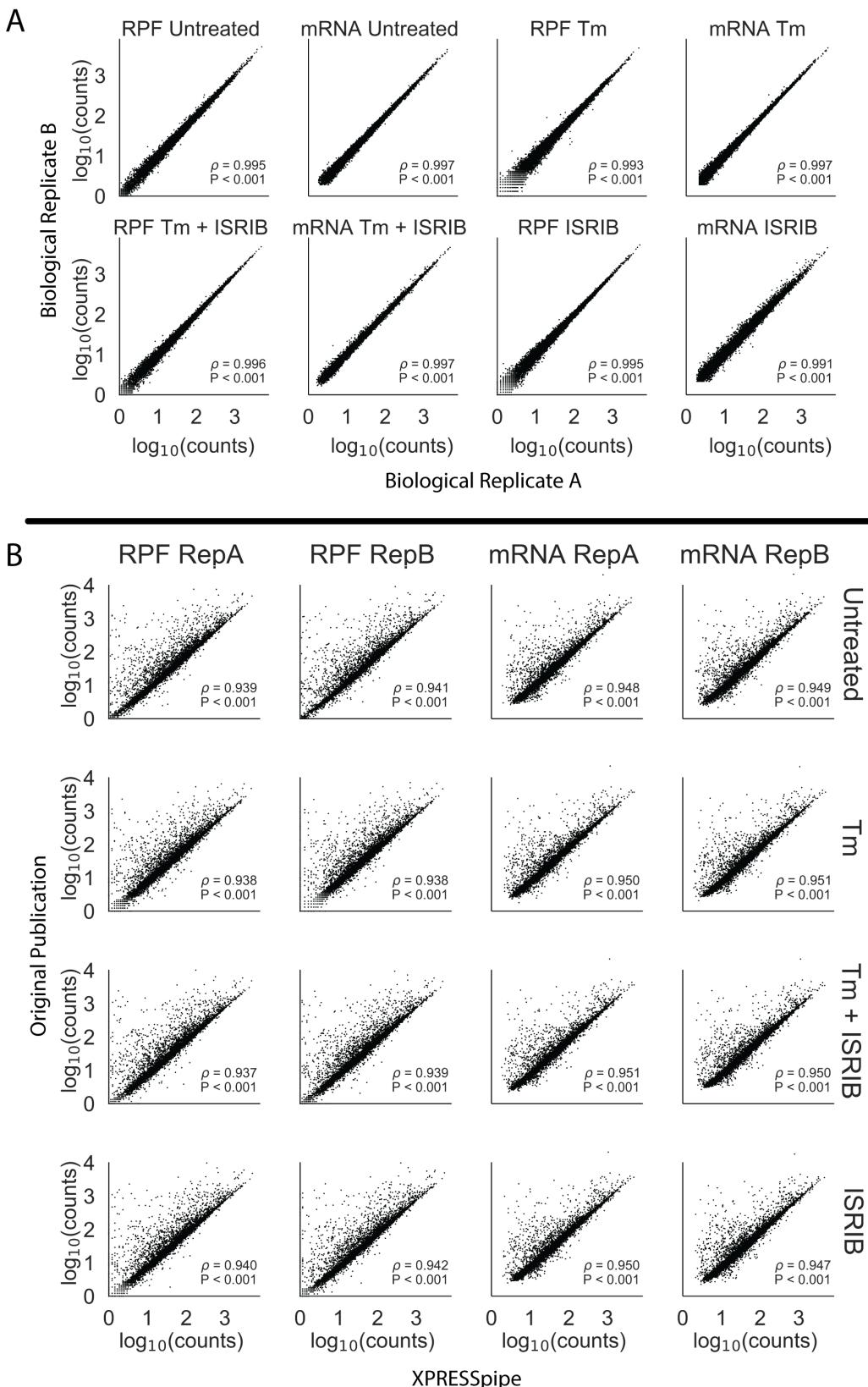


Figure S3: Comparison between processed data produced by XPRESSpipe and original study. Genes were eliminated from analysis if any RNA-Seq sample for that gene had fewer than 10 counts. A) Comparison of biological replicate read counts processed by XPRESSpipe. B) Comparison of read counts per gene between count data from the original study and the same raw data processed and quantified by XPRESSpipe. RPF, ribosome-protected fragments. Tm, tunicamycin. All ρ values reported are Spearman correlation coefficients. XPRESSpipe-processed read alignments were quantified to *Homo sapiens* build CRCh38v98 using a protein-coding only, truncated GTF.

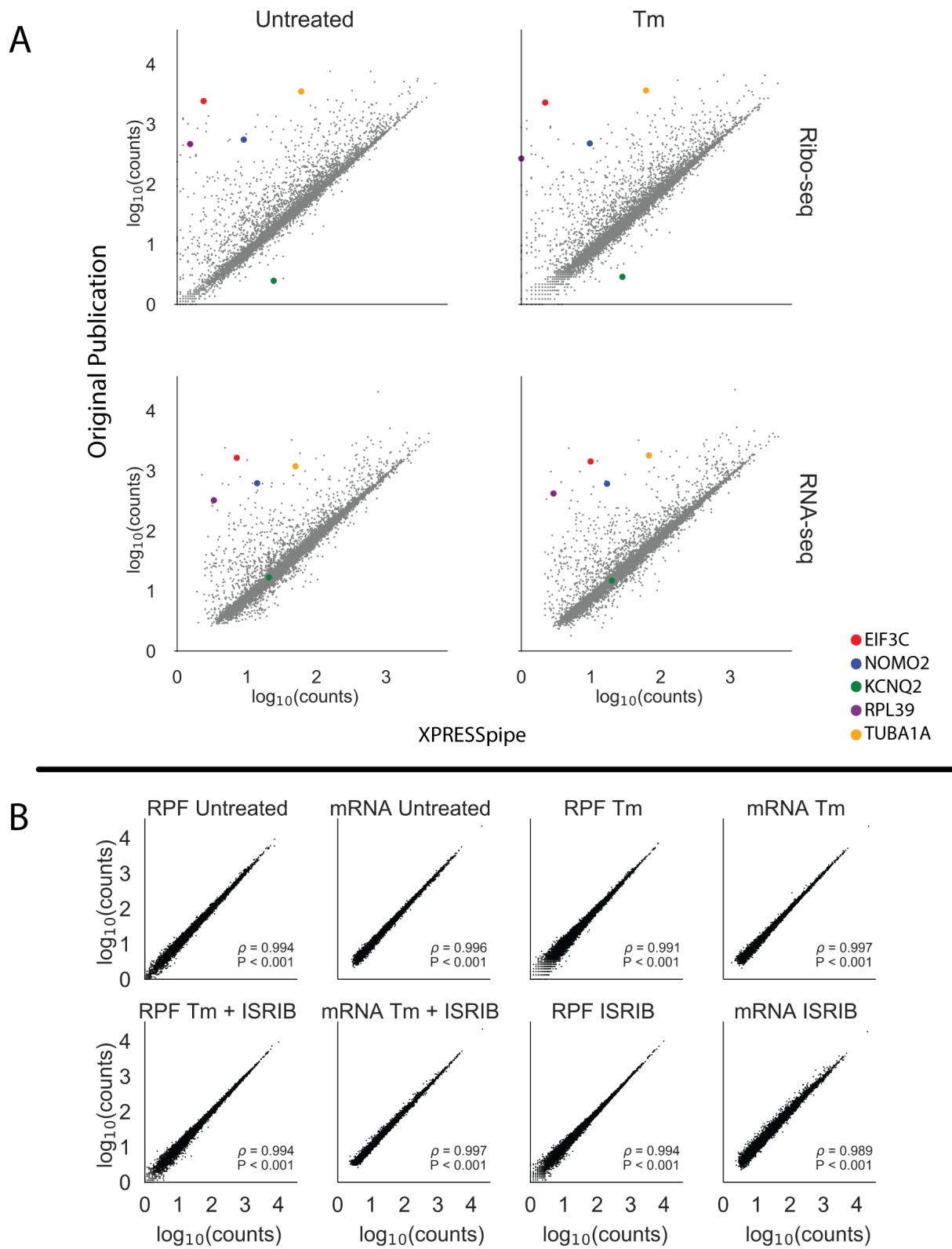


Figure S4: Original ISRIB count data plotted against XPRESSpipe-processed data reveals systematic differences between the analytical regimes. A) Selected highlighted genes show consistent differences between processing methods. B) Spearman correlation plots using the data table provided as supplementary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

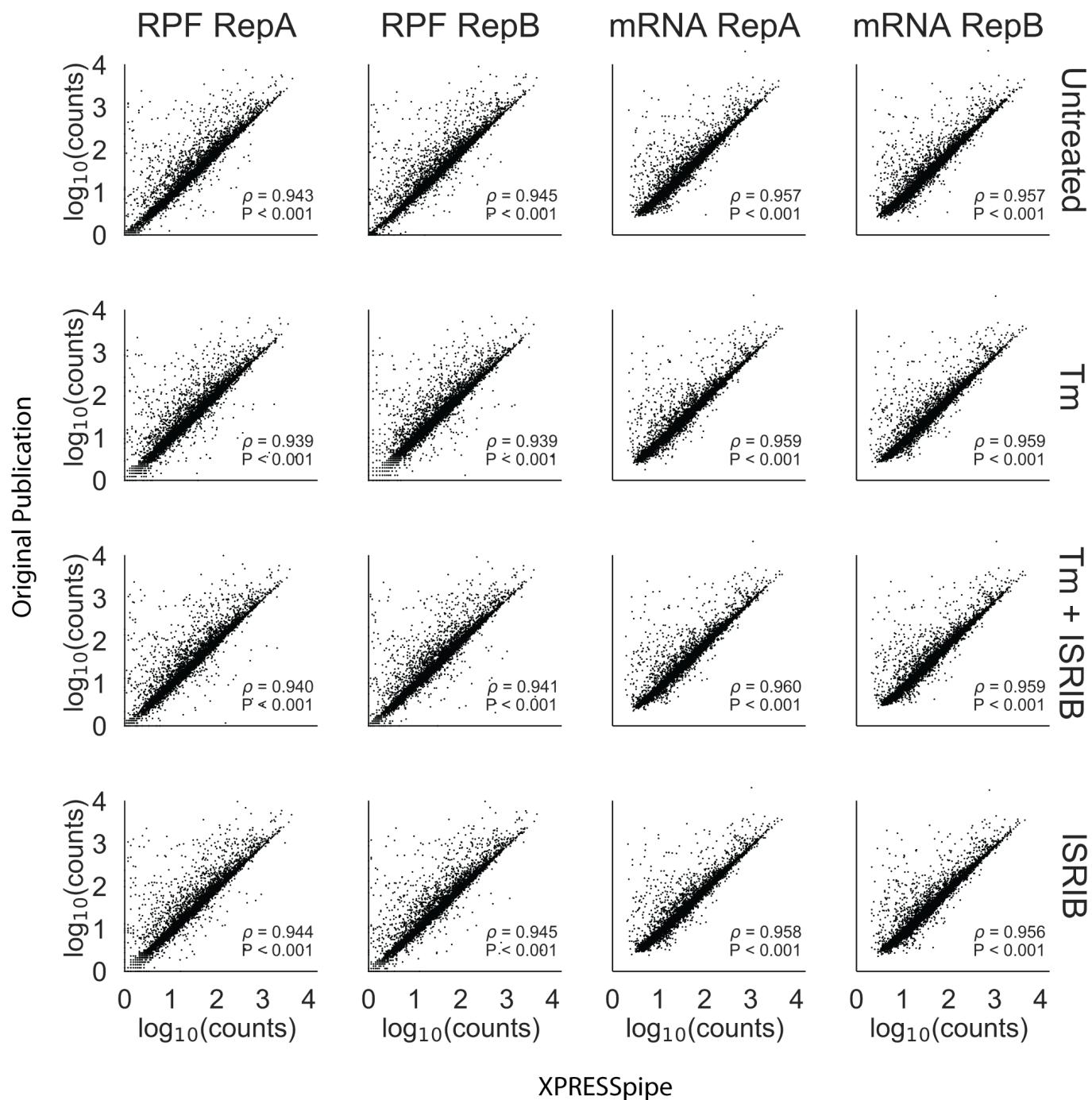


Figure S5: Original ISRIB count data plotted against XPRESSpipe-processed data quantified using same reference version reveals mild improvement in comparability between the analytical regimes. Original samples were processed using Ensembl human build GRCh38v72, as in the original manuscript, and compared 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 ρ values reported are Spearman correlation coefficients.

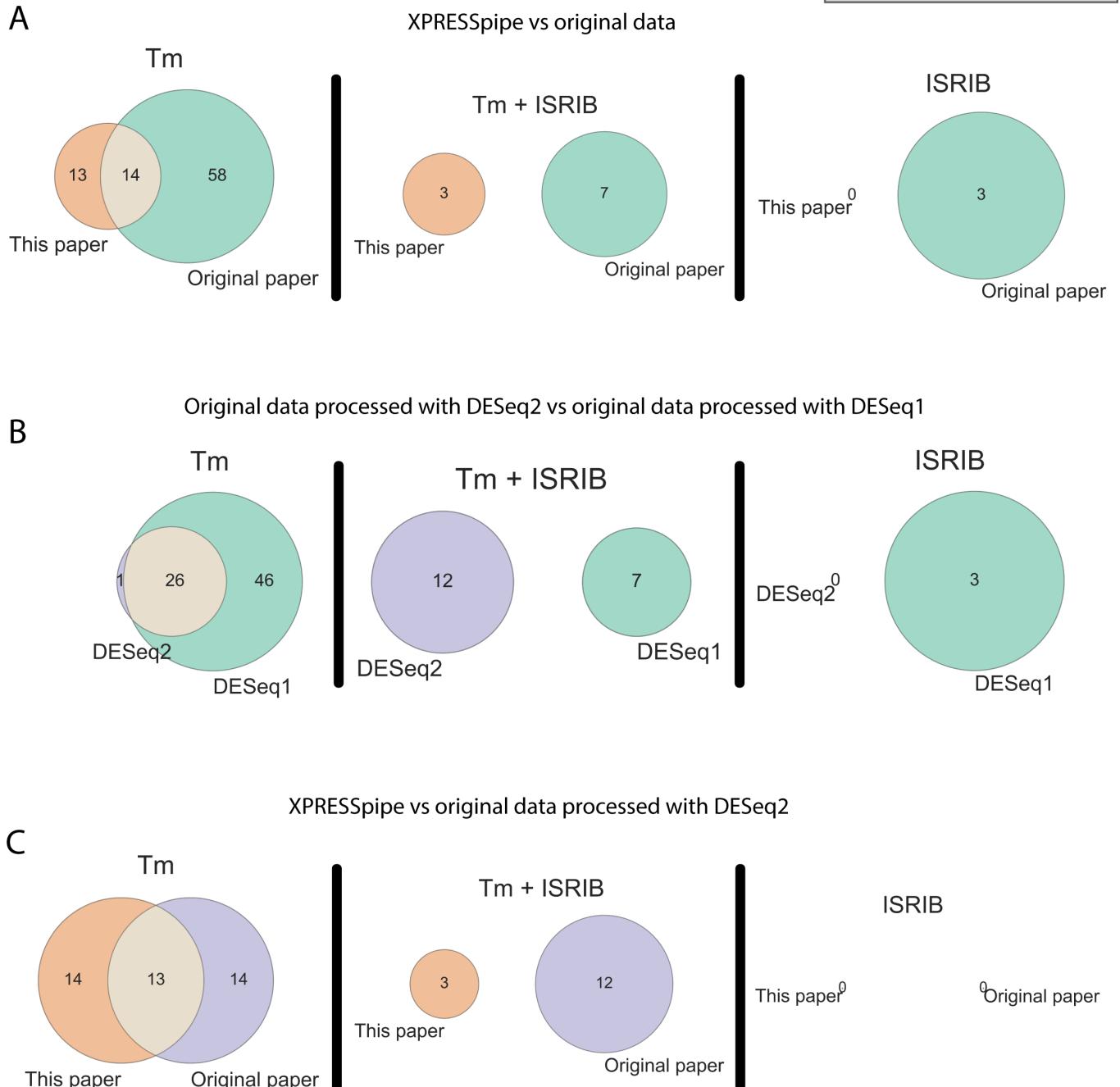


Figure S6: Cross-method analysis comparisons. A) XPRESSpipe-processed data (orange) versus data as originally presented within original manuscript using original methods (green). B) Comparison of analyses using provided count table in original publication using DESeq2 (purple) versus original analysis provided in manuscript using DESeq1 (green). C) XPRESSpipe-processed (orange) versus originally-processed data (purple), both using DESeq2 for differential expression analysis. Tan regions indicate overlap between gene lists. Thresholds used were the same as those used in the original study: $|\log_2(\text{Fold Change})| > 1$, FDR < 0.1.

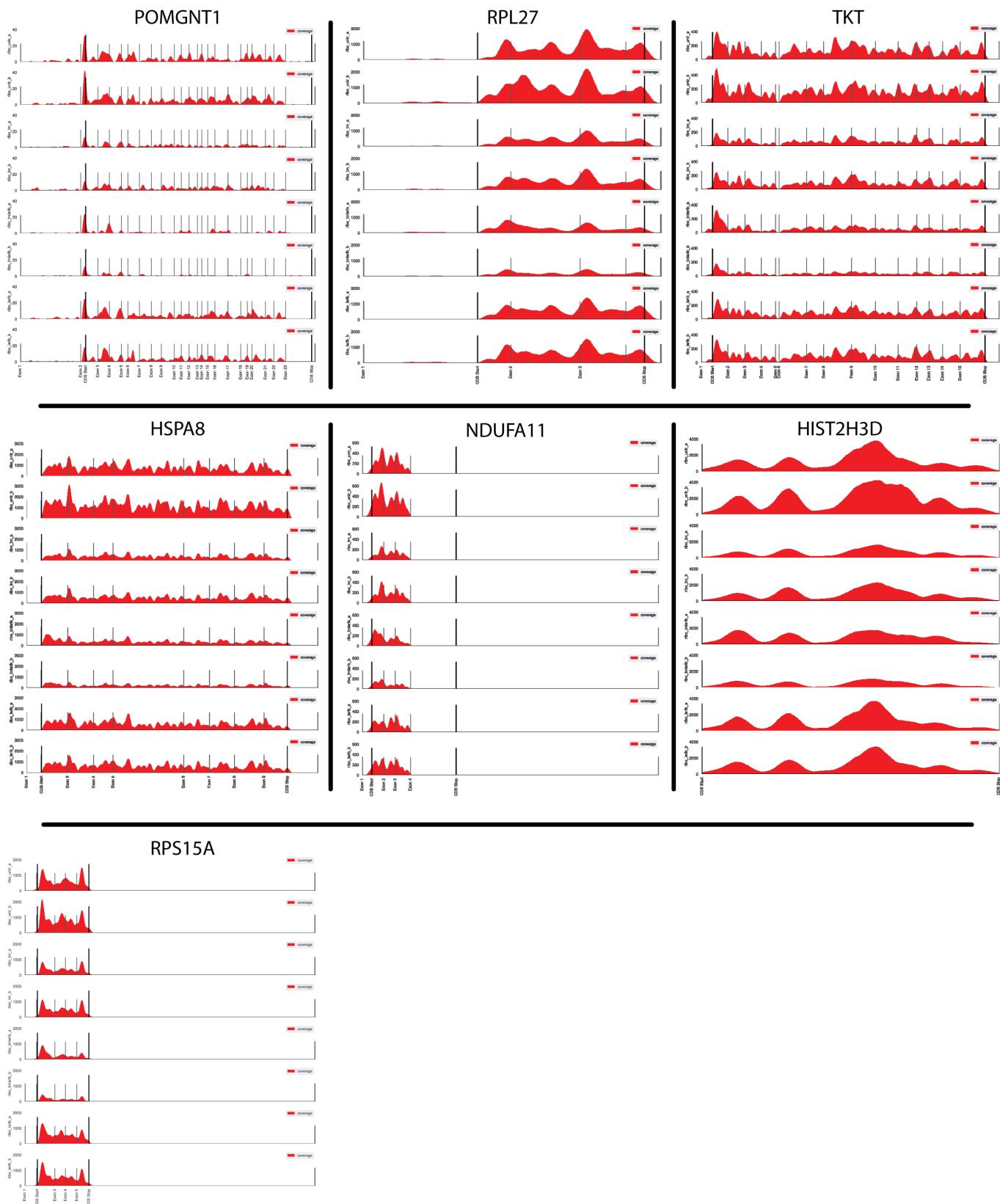


Figure S7: Gene coverage plots for neurologically annotated genes passing strict thresholding. Coverage plots were generated using XPRESSpipe's geneCoverage module, which collapses introns within the representation.