

# XPRESSyourself: Automating and Democratizing High-Throughput Sequencing

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

\*To whom correspondence should be addressed; E-mail: [rutter@biochem.utah.edu](mailto:rutter@biochem.utah.edu).

**With the advent of high-throughput sequencing platforms, expression profiling is becoming common-place in medical research. However, for the general user, often someone who ends up outsourcing their bioinformatics, a computational overhead exists. The XPRESSyourself suite aims to remove these barriers and create a tools to help standardize and increase throughput of data processing and analysis. The XPRESSyourself suite is currently broken down into two software packages. The first, XPRESSpipe, automates the pre-processing, alignment, quantification, normalization, and quality control of single-end and paired-end RNAseq, as well as ribosome profiling sequence data. The second, XPRESStools, is a Python toolkit for expression data analysis, compatible with private or public microarray and RNAseq datasets. This software suite is designed where features can easily be modified, and additional packages can be included for processing of other data types in the future, such as CHIPseq or genome alignment. Currently, this package offers several new tools for ribosome profiling and general RNA-seq.**

XPRESSyourself is freely available on GitHub: <https://github.com/XPRESSyourself>

## 1 Introduction

High-throughput profiling of gene expression data has revolutionized biomedical, industrial and basic science research. Within the last two decades, RNA-seq as found itself the forerunner technology for highest quality expression profiling, as it can measure relative transcript abundance, differential splice variants, sequence polymorphisms, and more. This technology has also been adopted to create technologies such as single-cell RNA-seq, capable of assaying the transcriptional profile cell by cell; and ribosome profiling, which measures ribosome occupancy and translation efficiency.

While vast strides have been made to these technologies, various bottlenecks still exist. For example, while more and more researchers are becoming accustomed to these technologies, learning the bioinformatics portion of sequencing possesses its own learning curve and often efficiency is lacking in how sequencing reads are processed and analyzed. Also for these users, they may not be aware of which tools are accepted as the

standard in the field or which analyses they should be perform. They may also lack the experience to process their sequencing libraries rapidly or may not know all the in-between steps that are not always explicitly stated in protocols.

While several pipelines have emerged over the last several years that have been built to tackle various aspects of these bottlenecks, most are not widely used or usable by the average wet-bench researcher. Some are difficult to install or use, often they break easily or do not perform well. Rarely do these tools offer anything new to help overcome emerging challenges in the field.

In response to these issues surrounding the automation and democratization of sequencing technology, we created the XPRESSyourself bioinformatics suite for processing and analyzing high-throughput expression data. In creating this tool, we focused on five aspects in order to create an easy, reliable tool where large barriers-to-entry would be eliminated. These were create a tool that was useful, usable, reliable, efficient, and flexible.

1. We wanted the software we created to be useful for a broad audience, where the bulk of processing and analysis desired by a general user would be covered. We wanted to use pre-existing tools that were fast and accurate. We also wanted to provide additional, new tools that would be of use to the general RNA-seq community, which will be discussed in more detail later.
2. We wanted to create a software package that was easy to use. To do so, we made the tools installable by a single command in the command line interface (CLI) using the Conda and PyPi package managers. We also included thorough external documentation hosted on readthedocs that outlines use and considerations for each tool, as well as provides several examples of how to use each tool. Internally in the CLI-packages, summary documentation has been included by way of the help interface. Jupyter notebooks are also created and installed with the software that provide example analyses that can be easily modified and run.
3. To create a reliable pipeline and analysis package, we use the most current state-of-the-art software tools that have undergone robust benchmarking. We utilize a two-pass RNA-seq alignment process to provide the best coverage around splice sites. We also built the RNA-seq pipeline according to The Cancer Genome Atlas (TCGA) standards. While this technology will no doubt improve over the years, the software is structures in a way for easy modification for addition of tools or substitution of software.
4. In order to make the most efficient package possible, by default XPRESSyourself optimizes use of computing cores to ensure all available are utilized when possible. Additionally, for analysis tools processing large files, we utilize a data matrix chunking method, where a dataset is portioned off into a number equal to the number of cores available, and processes each parallelly before rejoining the data chunks.
5. Flexibility is paramount in creating a tool that can be widely used and built upon. The general structure of the software was designed to make it easy to add or remove features. We envision as this suite of tools is more widely adopted by the RNAseq community, modules will be added to handle other sequencing platforms, such as genome sequencing, CHIPseq, and so on.

With XPRESSyourself, the user is provided with a complete suite of software to handle pre-processing, aligning, and quantifying reads, performing quality control via various meta-analyses of pre- and post-processed reads, and tools to perform the bulk of sequence analysis with enough flexibility to generate professional, figure-worthy images.

2 Materials and Methods

2.1 XPRESSpipe

XPRESSpipe pipelines for single-end RNA-seq, paired-end RNA-seq, and ribosome profiling offer a handful of tunable parameters to the user, while keeping most hidden to maintain TCGA alignment standards. In the future it is feasible that additional tunable parameters will be added. Table 1 outlines these parameters.

Table 1: Summary of XPRESSpipe pipeline arguments.

Arguments	Description
<b>Required</b>	
-i, -input	Path to input directory
-o, -output	Path to output directory
-r, -reference	Path to parent organism reference directory
-t, -reference.type	GTF, refFlat, etc. type (i.e. "DEFAULT", "CODING", "CODING_TRUNCATED")
-e, -experiment	Experiment name
<b>Optional</b>	
-a, -adaptors	Specify adaptor as string (only one allowed) – if "None" is provided, software will attempt to auto-detect adaptors – if "POLYX" is provided as a single string in the list, polyX adaptors will be trimmed
-q, -quality	PHRED read quality threshold (default: 28)
-min.length	Minimum read length threshold to keep for reads (default: 18)
-output.bed	Include option to output BED files for each aligned file
-output.bigwig	Include flag to output bigwig files for each aligned file
-method	Normalization method to perform (options: "RPM", "RPKM", "FPKM", "LOG")
-batch	Include path and filename of dataframe with batch normalization parameters
-sjdbOverhang	Sequencing platform read-length for constructing splice-aware reference previously (see documentation for more information)
-downstream	Number of nucleotides to track after the landmark (default: 200)
-m, -max.processors	Number of max processors to use for tasks (default: No limit)

2.1.1 Installation

General installation of XPRESSpipe is handled by the package manager, Anaconda, which is easy to download and install (<https://www.anaconda.com/distribution/>). XPRESSpipe can be installed using the following command:

```
1 $ conda install -c bioconda XPRESSpipe
```

Listing 1: curateReference example

Table 2: Summary of dependency software, accession location, and purpose in the XPRESSpipe package.

Package	URL	Purpose
fastp	<a href="https://github.com/OpenGene/fastp">https://github.com/OpenGene/fastp</a>	Read pre-processing
ucsc-gtftogenepred	<a href="http://hgdownload.cse.ucsc.edu/admin/exe/">http://hgdownload.cse.ucsc.edu/admin/exe/</a>	Create refFlat reference file used in meta-analysis
STAR	<a href="https://github.com/alexdobin/STAR">https://github.com/alexdobin/STAR</a>	Reference curation and read alignment
samtools	<a href="http://www.htslib.org/">http://www.htslib.org/</a>	Alignment file manipulation
bedtools	<a href="https://github.com/arq5x/bedtools2">https://github.com/arq5x/bedtools2</a>	Alignment file manipulation
deeptools	<a href="https://github.com/deeptools/deepTools">https://github.com/deeptools/deepTools</a>	Alignment file manipulation
htseq	<a href="https://github.com/simon-anders/htseq">https://github.com/simon-anders/htseq</a>	Read quantification
fastqc	<a href="https://www.bioinformatics.babraham.ac.uk/projects/fastqc/">https://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a>	Quality Control
multiqc	<a href="https://multiqc.info/">https://multiqc.info/</a>	Quality Control
picard	<a href="https://broadinstitute.github.io/picard/">https://broadinstitute.github.io/picard/</a>	Transcript meta-analysis
plastid	<a href="https://github.com/joshuagryphon/plastid">https://github.com/joshuagryphon/plastid</a>	Periodicity meta-analysis

XPRESSpipe is built upon several pre-established software packages, listed in Table 2. These packages will be automatically installing using the Anaconda package manager; however, if installing manually, these packages will need to be installed by the user.

2.1.2 Inputs

While inputs will vary sub-module to sub-module, and further information can be found in the documentation (<https://xpresspipe.readthedocs.io/en/latest/>) or by entering “xpresspipe <sub-module name> –help”, a few points of guidance are important to consider.

- Single-end reads should end in “.fa”, “.fasta”, “.txt”
- Paired-end reads should end in “.read1/2.suffix” or “.r1/2.suffix”
- The transcriptome reference file should be a valid GTF file and should be named “transcripts.gtf”
- If specifying a group of fasta files to use for alignment or reference curation, the directory containing this files cannot contain any other files ending in “.txt” or “.fa”

2.1.3 Reference Curation

One of the first preparatory steps of RNAseq alignment is curating a reference transcriptome for the aligner to map reads to. For the purposes of XPRESSpipe, a STAR reference should be created. However, for the purposes of several meta-analyses and more specific mapping, a STAR reference alone is not enough. For example, if one is aligning ribosome profiling reads and wishes to avoid mapping reads to the first 45 nucleotides in order to not quantify the inherent 5’ biases in library preparation, a truncated transcriptome reference needs to be used when counting reads to transcripts. Additionally, if one wishes to only quantify reads to the coding portions of the transcriptome, a custom reference file needs to be created. While each of these steps are available as stand-alone sub-modules, all potential reference files, including STAR files, modified transcriptome files, and other meta-analysis files can be created by running the “curateReference” command. An example is shown below for creating a ribosome profiling-ready reference XPRESSpipe directory. The following assumes one is avoiding mapping to the first 45 nucleotides, will only quantify to protein coding regions, and is tailored for mapping 50-bp single-end RNA-seq reads. As this can be a time-consuming process, we will leave the “–max\_processors” argument as default in order to utilize all cores available to the computing unit.

```
1 $ xpresspipe curateReference -o /path/to/output/location /
2                               -f /path/to/fasta/genome/files /
3                               -g /path/transcripts.gtf
4                               --truncate_amount 45
5                               --sjdbOverhang 49
```

Listing 2: curateReference example

While current available arguments are limited, the design is simple enough where arguments could be added easily to give more control over STAR reference creation; however, current setting align with TCGA standards.

2.1.4 Read Processing

While all intermediate steps of the pipelines can be run singly, we will describe the outline of the software in the context of the pipelines.

1. Trim: The first step in read processing is removing reads of artifacts for the library preparation process. These are most often adaptor sequences that were ligated to each read during the preparatory steps, which will not have homology to the organism’s reference sequence. As sequencers are only so accurate,

another common artifact are base calls that are statistically not confident. Therefore, in the first step of read processing, it is important to remove these sequences from each read to allow for proper alignment of reads to the reference. XPRESSpipe begins by trimming reads of adaptors, low quality bases, and reads that are too short using fastp, a more recent, faster trimming package that has improved alignable read output. Adaptor sequence, base quality, and read length are all adjustable parameters available to the user. Descriptions of the options can be found in Table 1.

2. **Align:** After sequencing reads have been trimmed of artifacts, the next step is to determine what transcripts these reads originated from to quantitate expression levels. This is done through alignment of each read to the reference transcriptome. XPRESSpipe uses STAR for this process as it has consistently proven itself as a fast and accurate alignment tool, such that it has been chosen as the TCGA standard alignment software. Alignment is made easy by generating a reference directory using XPRESSpipe’s “curateReference” sub-module. XPRESSpipe then performs a two-pass alignment to maximize alignment to splice junctions. In the first pass, STAR identified potential splice junctions for each sample, generates a reference taking these sites into account, then uses this reference to remap reads to the reference. Afterwards, some file processing is performed to generate the appropriate output files and formats for downstream steps.
3. **Count:** Once reads have been aligned and genome coordinates have been determined, these reads need to be quantified. XPRESSpipe uses htseq for this purpose. In this step, reads coordinates along the genome are mapped to the reference transcriptome. If a coding-only or truncated reference was created during the reference curation, this will be used for the quantification of transcripts. By default, htseq behavior conforms to TCGA standards by being strand agnostic, mapping assuming reads were sorted by name, and by using the intersection-nonempty method for handling reads overlapping multiple genes.
4. **Normalization:** As RNA-seq measurements are only relative, they are subject to sample and batch effects. These can arise from different people preparing libraries, one person preparing libraries on different days, or are just inherent in differences in total reads per sample a given chip sequences. In order to remove sample effects arising per sequence chip, it is important to perform normalization. Reads-per-million (RPM) normalization controls for these sample effects. Reads-per-kilobase-million (RPKM) or Fragments-per-kilobase-million (FPKM) performs the same RPM normalization, but additionally normalizes for transcript length as longer transcripts will appear to have more reads aligning to these transcripts normally. Batch effect normalization controls for library to library variances. These normalization options are available in the XPRESSpipe arguments for each pipeline (see Table 1).
5. **Quality Control:** It is important to perform quality control of sequencing sample to ensure they are reliable and their biological insight can be trusted. XPRESSpipe performs a variety of quality control measures. fastqc runs a gambit of quality control measures and multiqc summarizes these and metrics from trimming, aligning, and so on. The “readDistribution” sub-module will take the read size distributions for each library, as this can be informative, especially in the case of ribosome profiling, to ensure the proper read length was enriched in the sequencing library. The “metagene” sub-module will summarize read distribution along a representative transcript to highlight any sequencing biases, such as an enrichment of reads on the 3’ end indicative of poor library preparation. The “periodicity” sub-module will take the most enriched read

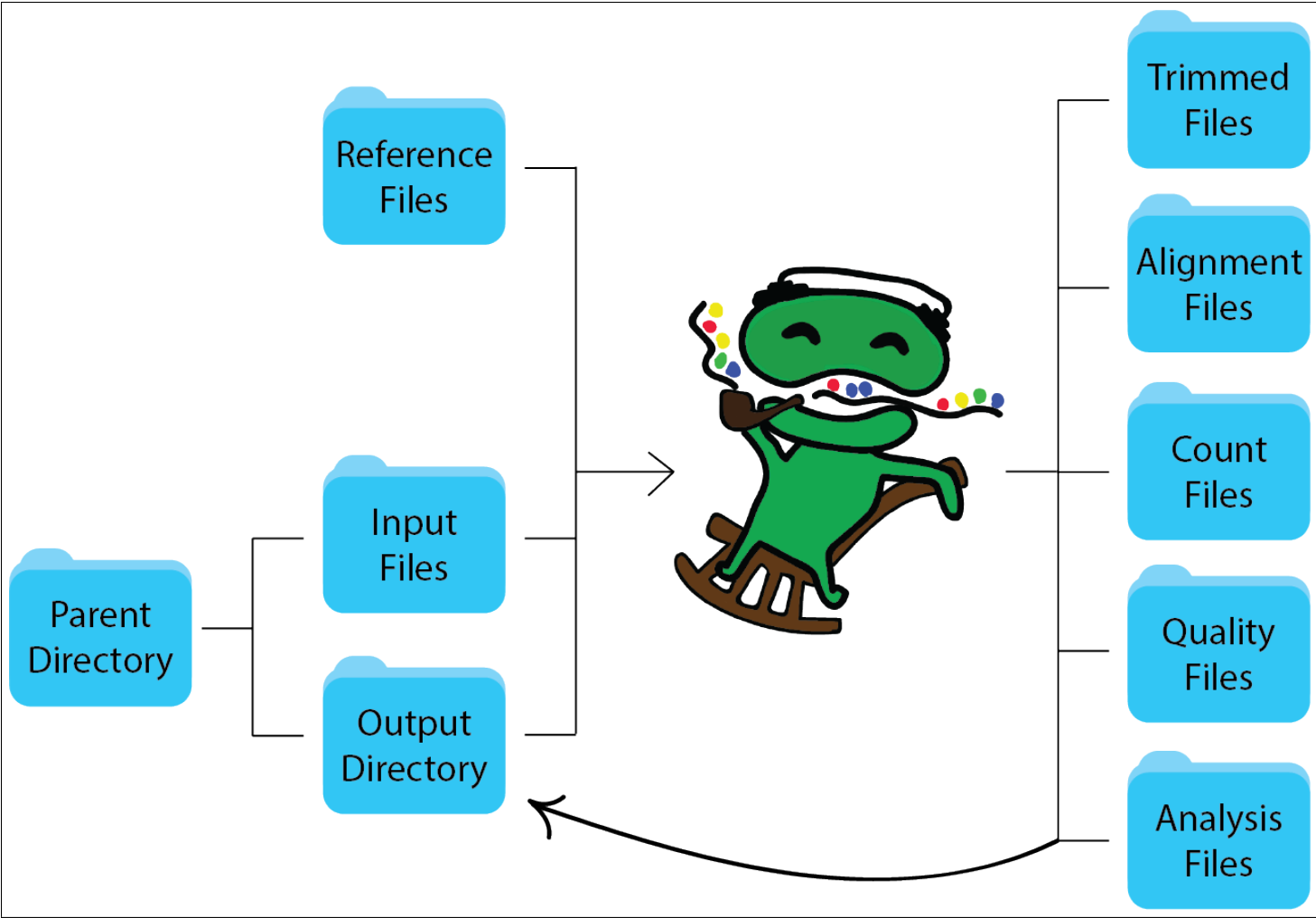


Figure 1: A schematic of XPRESSpipe outputs. Input, output, and often reference directories are specified and each stage of processing or analysis is given a separate directory.

length from ribosome profiling footprint libraries and map the P-site of the ribosome to investigate the 3 nucleotide phasing characteristic of the translating ribosome. Each of these sub-modular quality functions will summarize their results in a single PDF for quick reference by the user.

6. Analyses: prober

diffex

pipeline run singly or all together normalization / batch effect

2.1.5 Outputs

output minimal but enough to get clear picture optional outputs

Figure 1 provides an example of output file scheme for XPRESSpipe.

2.1.6 Quality Control

quality control read distribution meta-gene periodicity

2.1.7 Analyses

prober Deseq

2.2 XPRESStools

2.2.1 Getting Data

2.2.2 Normalizing and Formatting Data

2.2.3 Analyzing Data

2.3 Unit Testing and Code Coverage

New tools will require new tests to maintain code Coverage

2.4 Availability

Open source community GitHub Version Control Singularity

Results and Discussion

2.5 Benchmarking

2.6 Example Data Walkthrough

2.7 Cost Analysis

2.8 Summary

References

Acknowledgments

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.

Contributions

Conceptualization	J.A.B.
Supervision	M.T.H., J.G., A.R.Q., J.P.R.
Project Administration	J.A.B.
Investigation	J.A.B.
Formal Analysis	J.A.B.
Software	J.A.B.
Methodology	J.A.B.
Validation	J.A.B., A.J.B. Y.O.
Data Curation	J.A.B.
Resources	J.A.B., J.P.R.
Funding Acquisition	J.A.B., J.P.R.
Writing - Original Draft	J.A.B.
Writing - Review & Editing	J.A.B., M.T.H., J.G., A.R.Q., J.P.R.
Visualization	J.A.B.