

Pipelinier: A Nextflow-based framework for the definition of sequencing data processing pipelines

Anthony Federico*^{1,2}, Tanya Karagiannis¹, Kritika Karri¹, Dileep Kishore¹, Yusuke Koga²,
Joshua D. Campbell^{1,2} and Stefano Monti*^{1,2}

¹ Bioinformatics Program, Boston University, 24 Cummington Mall, Boston, MA 02215

² Division of Computational Biomedicine, Boston University School of Medicine, 75 East Newton Street, Boston, MA, 02118

* Correspondence: anfed@bu.edu (A.F.) smonti@bu.edu (S.M.)

Abstract

The advent of high-throughput sequencing technologies has led to the need for flexible and user-friendly data pre-processing platforms. The Pipelinier framework provides an out-of-the-box solution for processing various types of sequencing data. It combines the Nextflow scripting language and Anaconda package manager to generate modular computational workflows. We have used Pipelinier to create several pipelines for sequencing data processing including bulk RNA-seq, single-cell RNA-seq (scRNA-seq), as well as Digital Gene Expression (DGE) data. This report highlights the design methodology behind Pipelinier which enables the development of highly flexible and reproducible pipelines that are easy to extend and maintain on multiple computing environments. We also provide a quick start user guide demonstrating how to setup and execute available pipelines with toy datasets.

Introduction

High-throughput sequencing (HTS) technologies are vital to the study of genomics and related fields. Breakthroughs in cost efficiency have made it common for studies to obtain millions of raw sequencing reads. However, processing this data requires a series of computationally intensive tools which can be unintuitive to use, difficult to combine into stable workflows that can handle large number of samples, and challenging to maintain over long periods of time in different environments.

Pipelinier is a framework for the definition of sequencing data processing pipelines that aims to solve these issues. Pipelines developed within the framework are platform independent, fully reproducible, and inherit automated job parallelization and failure recovery. Their flexibility and modular architecture allows users to easily customize and modify processes based on their needs. Pipelinier also provides additional resources that allow developers to rapidly build and test their own pipelines in an efficient and scalable manner. Pipelinier is a complete and user-friendly solution to meet the demands of processing large amounts and various types of sequencing data.

Materials and Methods

Design and Features

Pipelinier is a suite of tools and methods for defining sequencing pipelines. It uses Nextflow, a portable, scalable, and parallelizable domain-specific language, to define data workflows (Di Tommaso *et al.*, 2017). Using Nextflow, each pipeline is modularized, consisting of a

configuration file as well as a series of processes. These processes define the major steps in each pipeline, and can be written in Linux-executable scripting languages such as Bash, Python, Ruby, etc. Nextflow processes are connected through channels – asynchronous FIFO queues – which allow data to be passed between the different steps in each pipeline using a Dataflow programming model. Using this architecture, pipelines developed within the Pipeliner framework inherit multiple features that contribute to their flexibility, reproducibility, and extensibility.

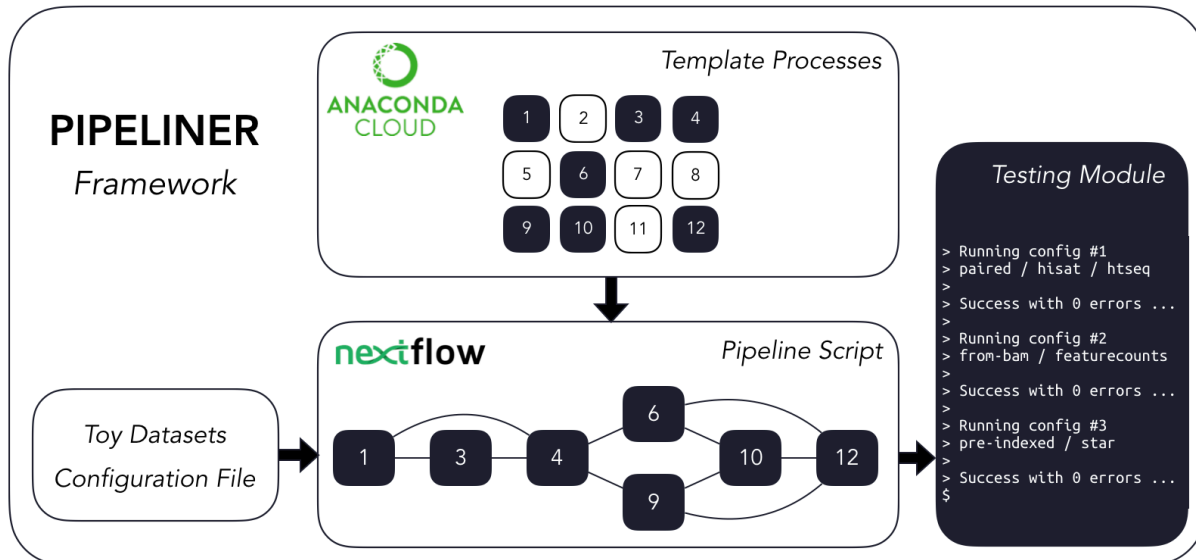


Fig 1. The Pipeliner framework employs reusable template processes strung together via Nextflow’s scripting language to create workflows in addition to developer tools such as toy datasets and testing modules.

Pipeline Flexibility

Pipeliner enables flexible customization of pipeline options and parameters. Pipeliner currently offers three pipelines to demonstrate its applicability in processing different types of data, including bulk RNA-seq, single-cell RNA-seq (scRNA-seq), as well as Digital Gene Expression (DGE) data (Soumillon *et al.*, 2014). For the RNA-seq pipeline, sequencing reads are checked for quality with FastQC (Andrews, 2010), trimmed with TrimGalore (Andrews, 2012), mapped to a reference genome with either STAR (Dobin *et al.*, 2012) or HISAT2 (Kim *et al.*, 2015), and quantified with either StringTie (Pertea *et al.*, 2015), HTSeq (Anders *et al.*, 2015), or featureCounts (Liao *et al.*, 2014). After alignment, mapping quality is checked with RSeQC (Wang *et al.*, 2012) and a comprehensive summary report of all processes is generated with MultiQC (Ewels *et al.*, 2016). The scRNA-seq and DGE pipelines adopt a similar methodology, and the development of additional pipelines for microRNA-seq (miRNA-seq) and RNA-seq Variant Calling is currently underway.

Parameter Configuration

All pipeline options and process parameters are set from a single configuration file (Fig. 2). Users have the option to select and skip various steps as well as customize parameters and allocate computing resources for specific processes. This flexibility gives rise to many different use cases. For example, a user may opt to provide a pre-indexed reference genome, or start the pipeline after

the mapping step with saved alignment files, or output an ExpressionSet data structure with count and phenotypic data. Thus, each pipeline is multi-purpose and allows users to frequently tweak settings without adding complexity or sacrificing reproducibility.

```
1 process {
2   executor = 'sge'
3
4   mapping.clusterOptions = "--P project -pe omp 16 -l mem_total=94G"
5   counting.clusterOptions = "--P project -pe omp 8 -l mem_total=16G"
6
7   indir = "/Users/pipeliner/pipelines/toy_data/rna-seq"
8   outdir = "/Users/pipeliner/pipelines/rna-seq-results"
9   fasta = "${params.indir}/genome_reference.fa"
10  gtf = "${params.indir}/genome_annotation.gtf"
11
12  // General pipeline parameters
13  paired = true
14  aligner = "hisat"
15  quantifier = "featurecounts"
16
17  index.use_existing = true
18  index.path = "${params.indir}/alignment_indices/hisat_index/index/part"
19
20  // Process-specific parameters
21  feature_counts.cpus = 8
22  feature_counts.type = "exon"
23  feature_counts.id = "gene_id"
24  feature_counts.xargs = ""
25  feature_counts.ainj = ""
26 }
```

Fig 2. A shortened example of a configuration file, highlighting the key components. This configuration includes resource allocations for cluster executions, input and output paths to data, general pipeline parameters, as well as process-specific parameters.

The default configuration file defines variables for common parameters of third-party software tools used in each pipeline. These tools are wrapped into templates – one for each process – which are executed sequentially within the pipeline script. Because some software tools have hundreds of arguments, users have the option to insert code injections from the configuration file. These code injections can be used to pass uncommon keyword arguments or to append ad hoc processing steps (Fig. 3). These features provide unrestricted control over each step in the execution of a pipeline. Furthermore, since all modifications are made within the configuration file – which is copied with each run – the pipeline script is left intact, preserving the reproducibility of each run regardless of any execution-specific changes the user may make.

```
1 featureCounts \\  
2  
3 # Common flags directly defined by the user  
4 -T ${params.feature_counts.cpus} \\  
5 -t ${params.feature_counts.type} \\  
6 -g ${params.feature_counts.id} \\  
7  
8 # Flags handled by the pipeline  
9 -a ${gtf} \\  
10 -o "counts.raw.txt" \\  
11  
12 # Arguments indirectly defined by the user  
13 ${feature_counts_sargs} \\  
14  
15 # Extra arguments  
16 ${params.feature_counts.xargs} \\  
17  
18 # Input data  
19 ${bamfiles};  
20  
21 # After injection  
22 ${params.feature_counts.ainj}
```

Fig 3. A code example of a template defined for the software tool featureCounts. The template wraps user-defined parameters, paths to data, as well as code injections into an executable bash script used in one of the pipeline steps.

Workflow Reproducibility

Pipelinier is designed to create reproducible workflows. An abstraction layer between Nextflow and Pipelinier logic enables platform independence and seamless compatibility with high performance cloud computing executors such as Amazon Web Services (AWS). Pipelinier also uses Anaconda – a multi-platform package and environment manager – to manage all third-party software dependencies and handle pre-compilation of all required tools before a pipeline is executed (Continuum Analytics, 2016).

Pipelinier is bundled with a pre-packaged environment hosted on Anaconda Cloud which contains all software packages necessary to run any of the three pipelines available. This virtual environment ensures consistent versioning of all software tools used during each pipeline execution. Additionally, all file paths, pipeline options, and process parameters are recorded, time-stamped, and copied into a new configuration file with each run, ensuring pipelines are fully reproducible regardless of where and when they are executed.

Extensibility

Pipelinier makes the development of Bioinformatics pipelines more efficient. The configuration file and processes that makeup each pipeline are inherited from shared blocks of code called template processes. For example, if a major update to an alignment tool requires modification to its template process, these changes propagate to all pipelines inheriting it. This property also minimizes the amount of code introduced as new pipelines are created, making them quicker to develop and easier to maintain. If a pipeline can inherit all of its processes with pre-defined templates, the user is only required to link these processes via Nextflow's scripting language and create a basic configuration file.

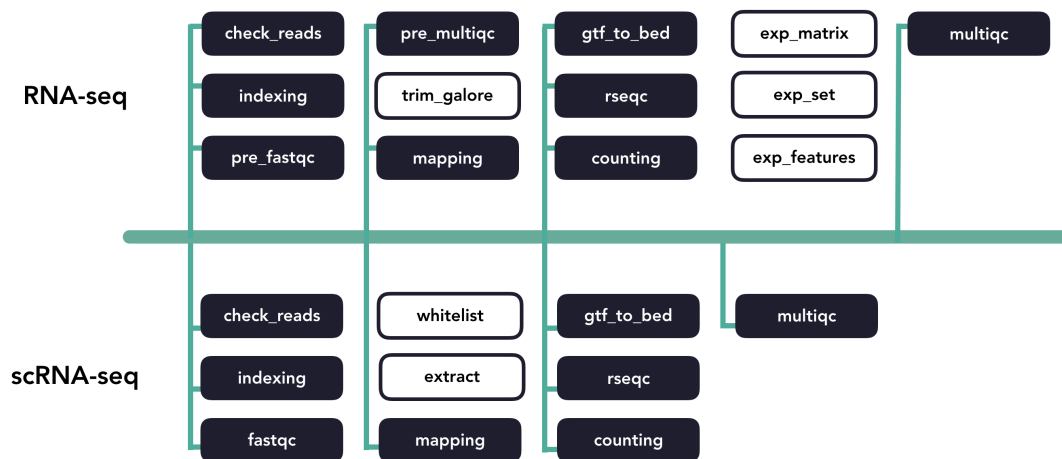


Fig 4. A diagram of template code sharing between the RNA-seq and scRNA-seq pipelines. Each block represents an individual workflow step. Shaded blocks share template code while unshaded blocks are unique.

Pipeliner requires configuration of paths to input data such as fastq reads, bam alignments, references files, etc. When cloning Pipeliner to a new machine, all paths must be reconfigured. This process can be automated by running a script which will reconfigure any paths to the same directory of your clone.

```
$ python pipeliner/scripts/paths.py
```

The final step is to download a Nextflow executable package in the same directory as the available pipelines.

```
$ cd pipeliner/pipelines
$ curl -s https://get.nextflow.io | bash
```

With the setup complete, any of the available pipelines can be executed with their respective toy datasets with the following commands.

```
$ ./nextflow rnaseq.nf -c rnaseq.config
$ ./nextflow scrnaseq.nf -c scrnaseq.config
$ ./nextflow dge.nf -c dge.config
```

Proof of Concept

To display the applicability of Pipeliner to real-world datasets, we reprocessed 48 RNA-seq paired read files for the Lymphoid Neoplasm Diffuse Large B-Cell Lymphoma (DLBC) cohort from The Cancer Genome Atlas (TCGA). For each cohort, the TCGA uses a standardized pipeline where reads are mapped to a reference genome with STAR and quantified by HTSeq. While the TCGA provides open access to the count matrix, some researchers have opted to use alignment and quantification algorithms specific to their research interests (Mumtahena *et al.*, 2015). For this reason, the TCGA also provides raw sequencing data, however its large size requires parallelization on a high-performance computing platform. We argue Pipeliner is a suitable choice for users looking for alternative reprocessing of TCGA datasets with minimal pipeline development.

Pipeliner makes alternative processing of TCGA and other publicly available data straightforward. In processing raw RNA-seq data for DLBC, paired fastq reads were downloaded from the Genomic Data Commons (GDC) Data Portal. For each sample, Pipeliner requires an absolute file path to reads. After specifying this information, Pipeliner was able to successfully process all data with HISAT2, featureCounts, and the remaining settings left to default (Fig. 6). The flexibility provided by Pipeliner is ideal for users experimenting with different tools and parameters. For example, because Pipeline is capable of taking aligned bam files as input and skipping preceding steps, we were able to rapidly try all three quantification options without re-running unrelated processes. This level of control is critical for downstream analysis of the processed data. To help researchers extend this example to other datasets, we provide the scripts used to obtain and organize TCGA

data from the GDG as well as the configuration file used by Pipeliner to process the data in the supplementary information.

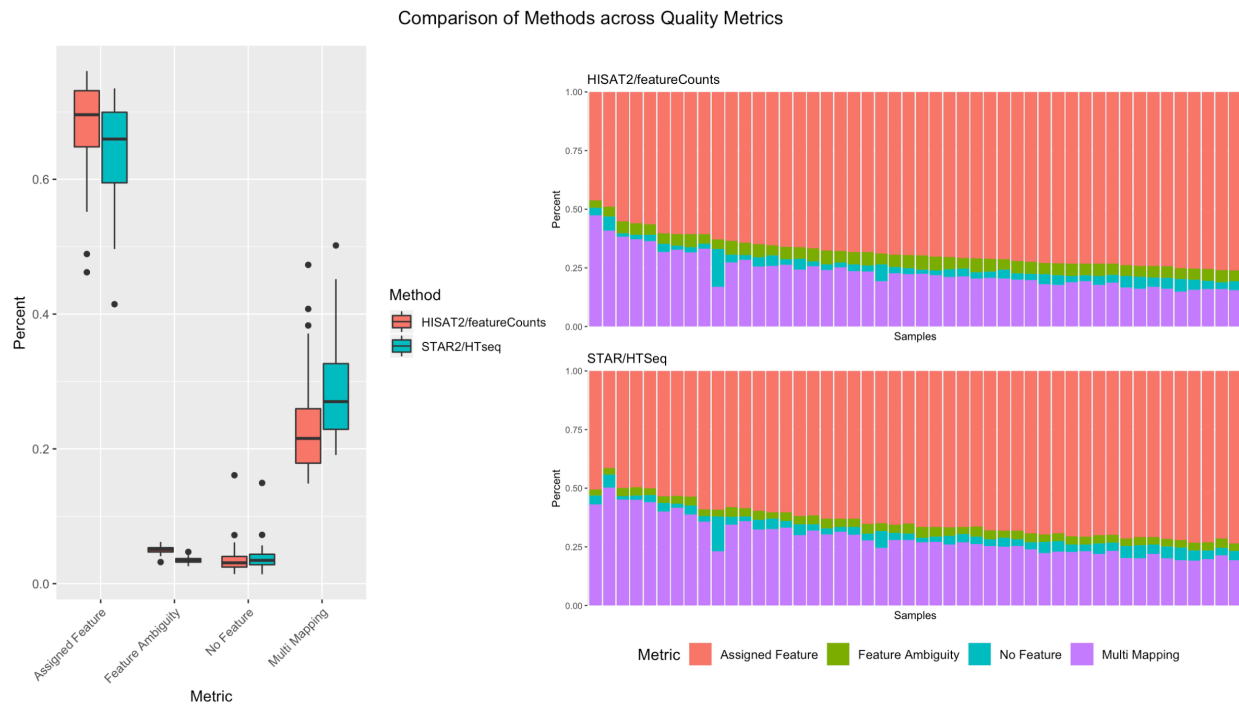


Fig 6. A summary of counting and mapping quality metrics when processing with Pipeliner using HISAT2 and featureCounts compared with processed counts from TCGA-DLBC using STAR and HTSeq.

Conclusions

Together with Nextflow and Anaconda, Pipeliner enables users to process large and complex sequencing datasets with pipelines that are customizable, reproducible, and extensible. The framework provides a set of user-friendly tools for rapidly developing and testing new pipelines for various types of sequencing data which will inherit valuable design features of existing pipelines. We apply the RNA-seq pipeline to real-world data by processing raw sequencing reads from the DLBC cohort provided by the TCGA and provide supplementary files which can be used to repeat the analysis or serve as a template for applying Pipeliner to other publicly available datasets.

Availability and Future Directions

Pipeliner is implemented in Nextflow, Python, R, and Bash and released under an MIT license. It is publicly available at <https://github.com/montilab/pipeliner> and supports Linux and OS X operating systems. Comprehensive documentation is generated with Sphinx and hosted by Read the Docs at <https://pipeliner.readthedocs.io/>. We will continue to develop the Pipeliner framework as the Nextflow programming language matures and we plan to provide additional pipelines for other types of sequencing data and analysis workflows in the future.

Conflict of Interest Statement

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Acknowledgements

The authors would like to thank P. Di Tommaso for his assistance with Nextflow-related inquiries and A. Gower for his advice for improving and testing Pipeliner.

Funding

This work was supported by a Superfund Research Program grant P42ES007381 (S.M.) and the LUNGeVity Career Development Award (J.D.C.)

References

Anders S. *et al.* (2015) HTSeq - a Python framework to work with high-throughput sequencing data. *Bioinformatics*, 31(2):166-9.

Andrews S. (2010) FastQC: a quality control tool for high throughput sequence data. Available online at: [https:// bioinformatics.babraham.ac.uk/projects/fastqc](https://bioinformatics.babraham.ac.uk/projects/fastqc).

Andrews S. (2012) Trim Galore: quality control and quality/adaptor trimming for Next-Gen sequencing applications. Available online at: https://bioinformatics.babraham.ac.uk/projects/trim_galore.

Continuum Analytics (2016) Anaconda Software Distribution: Version 2-2.4.0. Available online at: <https://continuum.io>.

Di Tommaso P. *et al.* (2017) Nextflow enables reproducible computational workflows. *Nature Biotechnology*, 35(4), 316-319.

Dobin A. *et al.* (2012) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, 1;29(1): 15-21.

Ewels P. *et al.* (2016) MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics*, 32(19): 3047-3048.

Kim D. *et al.* (2015) HISAT: a fast spliced aligner with low memory requirements. *Nature Methods*, 12(4):357-60.

Liao Y. *et al.* (2014) featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics*, 30(7): 923-30.

Mumtahena R. *et al.* (2015) Alternative preprocessing of RNA-Sequencing data in The Cancer Genome Atlas leads to improved analysis results. *Bioinformatics*. 31(22): 3666–3672.

Pertea M. *et al.* (2015) StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature Biotechnology*, 33(3): 290-295.

Soumillon M. *et al.* (2014) Characterization of directed differentiation by high-throughput single-cell RNA-seq. *BioRxiv*, doi.org/10.1101/003236.

Wang, L. *et al.* (2012) RSeQC: quality control of RNA-seq experiments. *Bioinformatics*, 28(16): 21845-5.