

RESEARCH ARTICLE

GeneTEFlow: A Nextflow-based pipeline for analysing gene and transposable elements expression from RNA-Seq data

Xiaochuan Liu¹, Jadwiga R. Bienkowska^{2*}, Wenyan Zhong^{1*}

1 Oncology Research & Development, Pfizer Worldwide Research and Development, Pearl River, NY, United States of America, **2** Oncology Research & Development, Pfizer Worldwide Research and Development, San Diego, CA, United States of America

* wenyan_zhong@yahoo.com (WZ); Jadwiga.R.Bienkowska@pfizer.com (JRB)



OPEN ACCESS

Citation: Liu X, Bienkowska JR, Zhong W (2020) GeneTEFlow: A Nextflow-based pipeline for analysing gene and transposable elements expression from RNA-Seq data. PLoS ONE 15(8): e0232994. <https://doi.org/10.1371/journal.pone.0232994>

Editor: Min Zhao, University of the Sunshine Coast, AUSTRALIA

Received: April 24, 2020

Accepted: August 2, 2020

Published: August 31, 2020

Copyright: © 2020 Liu et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Data are available from GitHub (<https://github.com/zhongw2/GeneTEFlow>).

Funding: JRB is an employee of Pfizer Inc. WZ was an employee of Pfizer Inc. and XL was a contractor of Pfizer Inc. when the work was being conducted. The funder provided support in the form of salaries for authors XL, JRB and WZ., but did not have any additional role in the study design, data collection and analysis, decision to publish, or preparation of

Abstract

Transposable elements (TEs) are mobile genetic elements in eukaryotic genomes. Recent research highlights the important role of TEs in the embryogenesis, neurodevelopment, and immune functions. However, there is a lack of a one-stop and easy to use computational pipeline for expression analysis of both genes and locus-specific TEs from RNA-Seq data. Here, we present GeneTEFlow, a fully automated, reproducible and platform-independent workflow, for the comprehensive analysis of gene and locus-specific TEs expression from RNA-Seq data employing Nextflow and Docker technologies. This application will help researchers more easily perform integrated analysis of both gene and TEs expression, leading to a better understanding of roles of gene and TEs regulation in human diseases. GeneTEFlow is freely available at <https://github.com/zhongw2/GeneTEFlow>.

Introduction

Transposable elements (TEs) are mobile DNA sequences which have the capacity to move from one location to another on the genome [1]. TEs make up a considerable fraction of most eukaryotic genomes and can be classified into retrotransposons and DNA transposons according to their different mechanisms of transposition and chromosomal integration [2, 3]. Retrotransposons are made of Long Terminal Repeats (LTRs) and non-LTRs that include long interspersed nuclear elements (LINEs) and short interspersed nuclear elements (SINEs) that mobilize via a RNA intermediate, while DNA transposons mobilize and function through a DNA intermediate [4–6]. TEs can be transcribed from the genome [7] and have been demonstrated to play important roles in the mammalian embryogenesis [8, 9], neurodevelopment [10, 11], and immune functions [12, 13]. Furthermore, aberrant expressions of TEs have been linked to cancers [14–16], neurodegenerative disorders [17, 18], and immune-mediated inflammation [19, 12]. Therefore, it has become increasingly important to explore biological roles of TEs expression. However, genome-wide analysis of TEs expression from high throughput RNA sequencing data has been a challenging computational problem. TEs contain highly repetitive sequence elements, making it arduous to unambiguously assign reads to the correct

the manuscript. The specific roles of these authors are articulated in the 'author contributions' section.

Competing interests: JRB is an employee of Pfizer Inc. WZ was an employee of Pfizer Inc. and XL was a contractor of Pfizer Inc. when the work was being conducted. This does not alter our adherence to PLOS ONE policies on sharing data and materials.

genomic location and accurately quantitate their expression level. Several bioinformatics tools have been developed to address this challenge with relatively good success [16, 20–22].

Recently, SQuIRE was reported to have the capability to quantify locus-specific expression of TEs from RNA-Seq data [22]. In addition, RNA-Seq data has long been used to detect dysregulated genes between different disease and/or drug treatment conditions to help understand disease mechanisms and/or drug response mechanisms. Therefore, it is of great interest to quantify both TEs and gene expression to elucidate contribution of both to disease mechanisms. Although many open source software and tools exist for analysing gene [23–25] and TEs expression, there are considerable challenges to efficiently apply these tools. In general, these multi-step data processing pipelines use many different tools. Correct versions of each tool need to be installed separately, and appropriate options, parameters, different reference genome and gene annotation files have to be set at each step. This can be quite tedious and challenging especially for non-computational users. Additionally, to ensure reproducibility of the analysis results, it is critical to capture analysis parameters from each step of the process. Equally important, to enable general use of the pipeline, the pipeline should be platform agnostic. Thus far, a one-stop computational framework for the comprehensive analysis of gene and locus-specific TEs expression from RNA-Seq data does not exist.

To address this need, we developed GeneTEFlow, a reproducible and platform-independent workflow, for the comprehensive analysis of gene and locus-specific TEs expression from RNA-Seq data using Nextflow [26] and Docker [27] technologies. GeneTEFlow provides several features and advantages for integrated gene and TEs transcriptomic analysis. First, by employing Docker technology, GeneTEFlow encapsulates bioinformatics tools and applications of specific versions into Docker containers enabling tracking, eliminating the need for software installation by users, and ensuring portability of the pipeline on multiple computing platforms including stand-alone workstations, high-performance computing (HPC) clusters, and cloud computing systems. Second, GeneTEFlow uses Nextflow to define the computational workflows, not only enabling parallelization and complete automation of the analysis, but also providing capability to track analysis parameters. Thus, GeneTEFlow allows users to generate reproducible analysis results through utilization of both Docker and Nextflow in a platform independent manner. Lastly, GeneTEFlow has modular architecture, and modules in GeneTEFlow can be turned on or off, providing developers with flexibility to build extensions tailored to specific analysis needs.

Implementation

The GeneTEFlow pipeline was developed using Nextflow, a portable, flexible, and reproducible workflow management system, and Docker technology, a solution to securely build and run applications on multiple platforms. To build the GeneTEFlow pipeline, a series of bioinformatics tools (S1 Table) were selected for QC, quantitation and differential expression analysis of genes and TEs from RNA-Seq data. These bioinformatics tools and custom scripts were built into four Docker containers to ensure portability of the workflow on different computational platforms. Data processing and analysis steps were implemented by modules using Nextflow. Modules are connected through channels and can be run in parallel. Each module in GeneTEFlow can include any executable Linux scripts such as Perl, R, or Python. Parameters for each module are defined in a configuration file.

A conceptual workflow of GeneTEFlow is illustrated in Fig 1. The workflow includes four major inputs: raw sequence files in fastq format, a sample meta data file in excel format, reference genome and gene annotation files, and a Nextflow configuration file. The sample meta data file contains detailed sample information and the design of group comparisons between

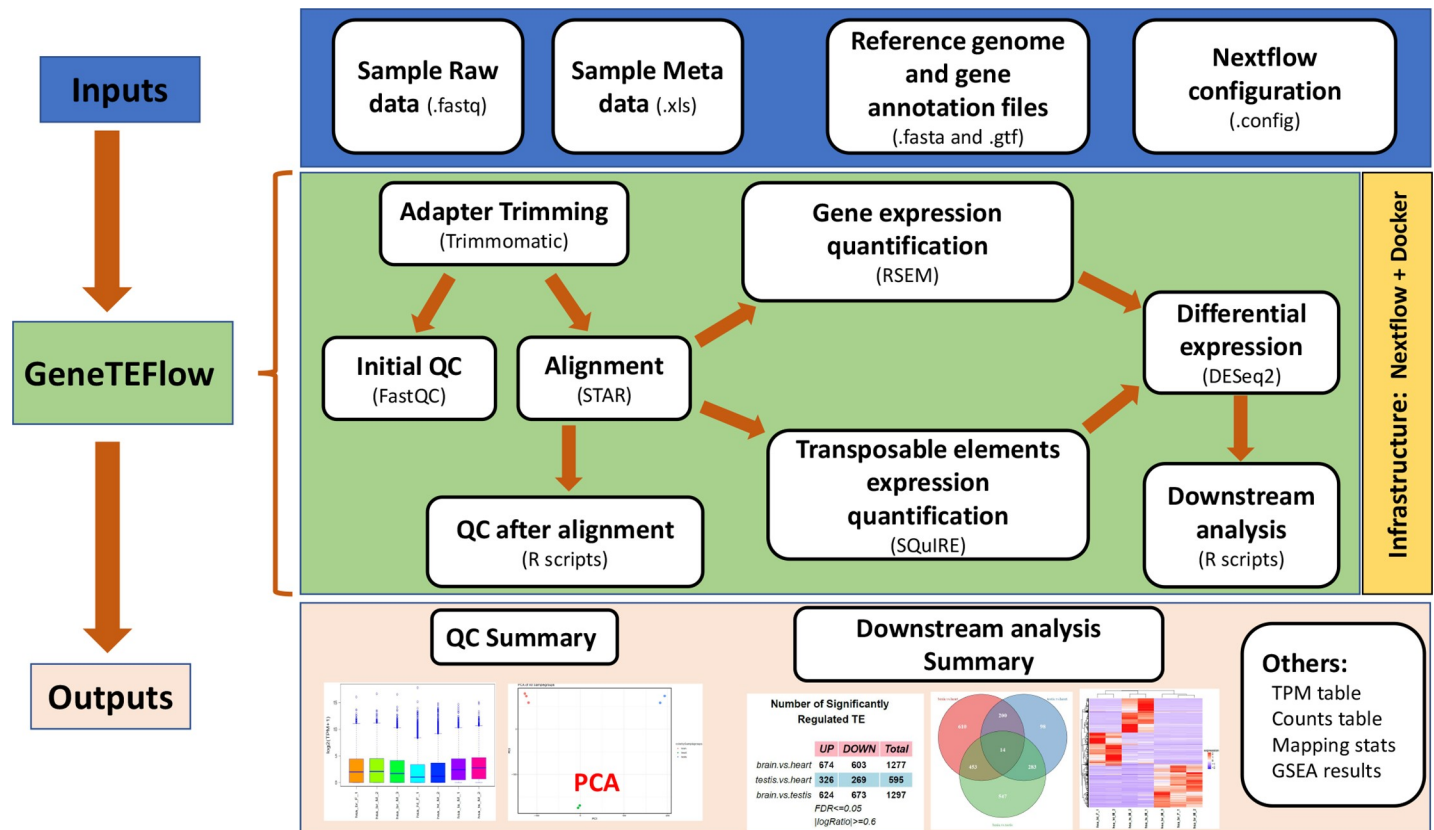


Fig 1. Illustration of GeneTEFlow: a Nexflow-based pipeline for identification of differentially expressed genes and locus specific transposable elements from RNA-Seq data.

<https://doi.org/10.1371/journal.pone.0232994.g001>

different experimental conditions. Human reference genome UCSC hg38 with the gene annotation (.gtf) was downloaded from Illumina iGenomes collections [28] and used by the bioinformatics tools included in GeneTEFlow. Scheduling of computational resources for each application module is defined in the configuration file.

GeneTEFlow analysis is performed in following steps: QC, expression quantification, differential expression and down-stream analysis. First, adapter sequences are trimmed off from the Illumina raw reads using Trimmomatic(v0.36) [29] for single-end or paired-end reads, and low-quality reads are filtered out. Next, FastQC(v0.11.7) [30] is executed to survey the quality of sequencing reads, and report is generated to help identify any potential issues of the high throughput sequencing data. Reference genome index for mapping sequencing reads to mRNA genes is built using "rsem-prepare-reference" of RSEM (v.1.3.0). Reads remaining after the pre-processing step are mapped to the reference genome using STAR(v2.6.0c) [31]. Gene level expression is quantitated as expected counts and transcripts per million (TPM) using "rsem-calculate-expression" of RSEM(v1.3.0) with default parameters [32]. Custom Perl scripts were developed to aggregate data from each sample into a single data matrix for expected counts and TPM values respectively. The expression quantification of locus-specific TEs is performed by SQuIRE [22]. After aligning reads to reference genome, SQuIRE classifies reads into unique reads (reads mapped to a single locus) and multi-mapping reads (reads mapped to multiple locations). Then SQuIRE calculates unique read expression of each annotated TE and assigns the fractions of multi-mapping reads based on the normalized unique read expression of each TE. Finally, expectation-maximization (EM) algorithm was

implemented in SQuIRE to recursively calculate the fractions of multi-mapping reads (E-step) of each TE, and re-estimate total read counts (M-step) until convergence.

In addition, we also implemented quality control measures after reads alignment step to detect potential outlier samples resulted from experimental errors. Boxplot and density plot are used to evaluate the overall consistency of the expression distribution for each sample. Sample correlation analysis is performed with Pearson method using TPM values to assess the correlation between biological replicates from each sample group. Principal component analysis (PCA) is employed to identify potential outlier samples and to investigate relationships among sample groups.

Differential expression analysis of genes and transposable elements is performed using DESeq2(v1.18.1) package [33]. Significantly up-regulated and down-regulated genes and TEs are summarized in a table. To analyse overlap among significantly regulated genes and TEs from pair-wise comparisons between different sample groups we use Venn diagrams. We perform hierarchical clustering of significantly dysregulated genes or TEs using R package “ComplexHeatmap” [34] with euclidean distance and average linkage clustering parameters. Gene set enrichment analysis (GSEA, v3.0) [35] is conducted using collections from the Molecular Signatures Database (MSigDB) [36]. The outputs (S2 Table) from GeneTEFlow are organized into several folders predefined in a GeneTEFlow configuration file.

In addition, GeneTEFlow can be run in step-by-step mode, which allows users to explore their RNA-Seq data. On GitHub, we include an example on how to run GeneTEFlow in a flexible manner where user has an option to remove some low-quality samples before proceeding with the differential expression analysis. A tutorial with detailed instructions on how to set up and run GeneTEFlow is provided at <https://github.com/zhongw2/GeneTEFlow>.

Application of GeneTEFlow

We applied GeneTEFlow to a public dataset from Brawan’s study [37] investigating tissue-specific expression changes of genes and transposable elements. Human RNA-Seq data from brain, heart and testis tissues were downloaded from GEO (accession number: GSE30352) (S3 Table). Expression analysis of genes and TEs were performed using GeneTEFlow and results are shown in Fig 2. Gene expression analysis was performed using RSEM and DESeq2 modules while TEs expression analysis was conducted using SQuIRE and DESeq2 modules within GeneTEFlow. Significantly regulated genes were identified with FDR less than 0.05 and fold change greater than 2. Significantly regulated locus-specific transposable elements were identified with FDR less than 0.05 and fold change greater than 1.5. The number of significantly regulated genes and transposable elements were summarized into two tables respectively (Fig 2, top panels). Using GeneTEFlow, we detected genes and TE differentially expressed between different tissue types (brain vs heart tissues: 6,264 genes and 1,277 TEs; testis vs heart tissues: 7,066 genes and 595 TE; brain vs testis tissues: 8,125 genes and 1,297 TEs) with most significant gene and TE expression differences observed being between brain and testis tissues. Our analysis identified large number of both genes and TEs with tissue specific expression patterns (Fig 2, middle panels and bottom panels). More in depth analysis to include additional tissue types would be required to fully understand the tissue specific gene and TEs expression and their relationship. GeneTEFlow is a computational solution to facilitate such studies.

Although SQuIRE provides both gene and TEs expression quantification, we also implemented the widely used RSEM method to provide users with alternative approaches when only gene expression quantification is desired. We compared gene level expression quantification between RSEM and SQuIRE (S1 Fig). The results showed high concordance (correlation coefficient: ~97%) of the gene level expression quantification between the two methods (S1

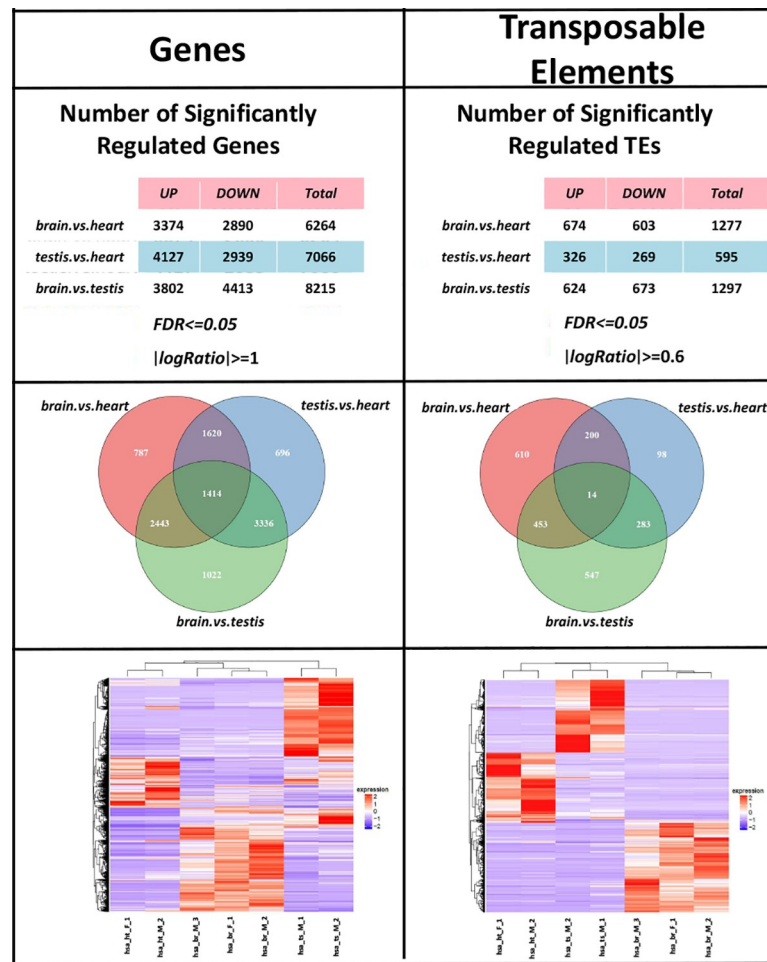


Fig 2. Differential expression analysis results of genes and transposable elements from GeneTEFlow. Left panels: gene results; right panels: TEs results. Top panels: number of significantly regulated genes or TEs in each sample group comparison. Significance was defined as following: $FDR \leq 0.05$ and fold change ≥ 2 for gene expression analysis; $FDR \leq 0.05$ and fold change ≥ 1.5 for TEs expression analysis. Middle panels: overlaps of significantly regulated genes or TEs amongst sample group comparisons. Bottom panels: hierarchical clustering of significantly regulated genes or TEs.

<https://doi.org/10.1371/journal.pone.0232994.g002>

[Fig](#), highlighted in red box) suggesting a robust measurement for both gene and TEs expression by SQuIRE.

Conclusions

In conclusion, we have developed and made available an automated pipeline to comprehensively analyse both gene and locus-specific TEs expression from RNA-Seq data. Taking advantage of the advanced functionalities provided by Nextflow and Docker, GeneTEFlow allows users to run analysis reproducibly on different computing platforms without the need for individual tool installation and manual version tracking. We believe this pipeline will be of great help to further our understanding of roles of both gene and TEs regulation in human diseases. This pipeline is flexible and can be easily extended to include additional types of analysis such as alternative splicing, fusion genes, and so on.

Supporting information

S1 Fig. Comparison of gene expression quantification by RSEM and SQuIRE. Gene expression (total 22,955 genes) of samples from brain tissues (left), heart tissues (middle), and testis tissues (right) was calculated by both RSEM and SQuIRE. Lower diagonal panels: pairwise comparisons using $\log_2(\text{TPM} + 1)$ of 22,955 genes. Upper diagonal panels: *Pearson* correlation coefficient of each comparison. Panels highlighted in red: *Pearson* correlation coefficient of comparisons between RSEM and SQuIRE gene expression quantification of the same sample. Rep_: replicate, _RSEM: quantification performed by RSEM, _SQuIRE: quantification performed by SQuIRE.

(TIF)

S1 Table. Major bioinformatics tools installed in GeneTEFlow.

(DOCX)

S2 Table. Major outputs from GeneTEFlow.

(DOCX)

S3 Table. Human RNA-Seq data used in the example application of GeneTEFlow.

(DOCX)

Acknowledgments

We gratefully acknowledge inputs and support from our colleagues: Jeremy Myers, Keith Ching, Corey Dasilva and Da Tse.

Author Contributions

Conceptualization: Wenyan Zhong.

Formal analysis: Xiaochuan Liu.

Methodology: Xiaochuan Liu, Wenyan Zhong.

Resources: Jadwiga R. Bienkowska.

Software: Xiaochuan Liu, Wenyan Zhong.

Supervision: Jadwiga R. Bienkowska, Wenyan Zhong.

Visualization: Xiaochuan Liu, Wenyan Zhong.

Writing – original draft: Xiaochuan Liu.

Writing – review & editing: Jadwiga R. Bienkowska, Wenyan Zhong.

References

1. Biémont C, Vieira C. Junk DNA as an evolutionary force. *Nature*. 2006; 443(7111):521–4. <https://doi.org/10.1038/443521a> PMID: 17024082
2. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, et al. A unified classification system for eukaryotic transposable elements. *Nature Reviews Genetics*. 2007; 8(12):973–82. <https://doi.org/10.1038/nrg2165> PMID: 17984973
3. Feschotte C. Transposable elements and the evolution of regulatory networks. *Nature Reviews Genetics*. 2008; 9(5):397–405. <https://doi.org/10.1038/nrg2337> PMID: 18368054
4. Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M, et al. Ten things you should know about transposable elements. *Genome Biology*. 2018; 19(1):199. <https://doi.org/10.1186/s13059-018-1577-z> PMID: 30454069

5. Lanciano S, Mirouze M. Transposable elements: all mobile, all different, some stress responsive, some adaptive? *Current Opinion in Genetics & Development*. 2018; 49:106–14. <https://doi.org/10.1016/j.gde.2018.04.002>.
6. Chuong EB, Elde NC, Feschotte C. Regulatory activities of transposable elements: from conflicts to benefits. *Nature Reviews Genetics*. 2017; 18(2):71–86. <https://doi.org/10.1038/nrg.2016.139> PMID: 27867194
7. Rebollo R, Romanish MT, Mager DL. Transposable Elements: An Abundant and Natural Source of Regulatory Sequences for Host Genes. *Annual Review of Genetics*. 2012; 46(1):21–42. <https://doi.org/10.1146/annurev-genet-110711-155621> PMID: 22905872.
8. Percharde M, Lin C-J, Yin Y, Guan J, Peixoto GA, Bulut-Karslioglu A, et al. A LINE1-Nucleolin Partnership Regulates Early Development and ESC Identity. *Cell*. 2018; 174(2):391–405.e19. <https://doi.org/10.1016/j.cell.2018.05.043> PMID: 29937225
9. Garcia-Perez JL, Widmann TJ, Adams IR. The impact of transposable elements on mammalian development. *Development*. 2016; 143(22):4101–14. <https://doi.org/10.1242/dev.132639> PMID: 27875251
10. Sun W, Samimi H, Gamez M, Zare H, Frost B. Pathogenic tau-induced piRNA depletion promotes neuronal death through transposable element dysregulation in neurodegenerative tauopathies. *Nature Neuroscience*. 2018; 21(8):1038–48. <https://doi.org/10.1038/s41593-018-0194-1> PMID: 30038280
11. Guo C, Jeong H-H, Hsieh Y-C, Klein H-U, Bennett DA, De Jager PL, et al. Tau Activates Transposable Elements in Alzheimer's Disease. *Cell Reports*. 2018; 23(10):2874–80. <https://doi.org/10.1016/j.celrep.2018.05.004> PMID: 29874575
12. Colombo AR, Elias HK, Ramsingh G. Senescence induction universally activates transposable element expression. *Cell Cycle*. 2018; 17(14):1846–57. <https://doi.org/10.1080/15384101.2018.1502576> PMID: 30080431
13. Koonin EV, Krupovic M. Evolution of adaptive immunity from transposable elements combined with innate immune systems. *Nature Reviews Genetics*. 2015; 16(3):184–92. <https://doi.org/10.1038/nrg3859> PMID: 25488578
14. Colombo AR, Triche T, Ramsingh G. Transposable Element Expression in Acute Myeloid Leukemia Transcriptome and Prognosis. *Scientific Reports*. 2018; 8(1):16449. <https://doi.org/10.1038/s41598-018-34189-x> PMID: 30401833
15. Burns KH. Transposable elements in cancer. *Nature Reviews Cancer*. 2017; 17(7):415–24. <https://doi.org/10.1038/nrc.2017.35> PMID: 28642606
16. Criscione SW, Zhang Y, Thompson W, Sedivy JM, Neretti N. Transcriptional landscape of repetitive elements in normal and cancer human cells. *BMC Genomics*. 2014; 15(1):583. <https://doi.org/10.1186/1471-2164-15-583> PMID: 25012247
17. Krug L, Chatterjee N, Borges-Monroy R, Hearn S, Liao W-W, Morrill K, et al. Retrotransposon activation contributes to neurodegeneration in a *Drosophila* TDP-43 model of ALS. *PLOS Genetics*. 2017; 13(3):e1006635. <https://doi.org/10.1371/journal.pgen.1006635> PMID: 28301478
18. Tam OH, Ostrow LW, Gale Hammell M. Diseases of the nERVous system: retrotransposon activity in neurodegenerative disease. *Mobile DNA*. 2019; 10(1):32. <https://doi.org/10.1186/s13100-019-0176-1> PMID: 31372185
19. De Cecco M, Ito T, Petrashen AP, Elias AE, Skvir NJ, Criscione SW, et al. L1 drives IFN in senescent cells and promotes age-associated inflammation. *Nature*. 2019; 566(7742):73–8. <https://doi.org/10.1038/s41586-018-0784-9> PMID: 30728521
20. Jin Y, Tam OH, Paniagua E, Hammell M. TETranscripts: a package for including transposable elements in differential expression analysis of RNA-seq datasets. *Bioinformatics*. 2015; 31(22):3593–9. <https://doi.org/10.1093/bioinformatics/btv422> PMID: 26206304
21. Lerat E, Fablet M, Modolo L, Lopez-Maestre H, Vieira C. TETools facilitates big data expression analysis of transposable elements and reveals an antagonism between their activity and that of piRNA genes. *Nucleic Acids Research*. 2016; 45(4):e17–e. <https://doi.org/10.1093/nar/gkw953> PMID: 28204592
22. Yang WR, Ardeljan D, Pacyna CN, Payer LM, Burns KH. SQUIRE reveals locus-specific regulation of interspersed repeat expression. *Nucleic Acids Research*. 2019; 47(5):e27–e. <https://doi.org/10.1093/nar/gky1301> PMID: 30624635
23. Varet H, Brillat-Guéguen L, Coppée J-Y, Dillies M-A. SARTools: A DESeq2- and EdgeR-Based R Pipeline for Comprehensive Differential Analysis of RNA-Seq Data. *PLOS ONE*. 2016; 11(6):e0157022. <https://doi.org/10.1371/journal.pone.0157022> PMID: 27280887
24. Finotello F, Di Camillo B. Measuring differential gene expression with RNA-seq: challenges and strategies for data analysis. *Briefings in Functional Genomics*. 2014; 14(2):130–42. <https://doi.org/10.1093/bfgp/elu035> PMID: 25240000

25. Ewels PA, Peltzer A, Fillinger S, Patel H, Alneberg J, Wilm A, et al. The nf-core framework for community-curated bioinformatics pipelines. *Nature Biotechnology*. 2020; 38(3):276–8. <https://doi.org/10.1038/s41587-020-0439-x> PMID: 32055031
26. Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow enables reproducible computational workflows. *Nature Biotechnology*. 2017; 35(4):316–9. <https://doi.org/10.1038/nbt.3820> PMID: 28398311
27. Merkel D. Docker: lightweight Linux containers for consistent development and deployment. *Linux J*. 2014; 2014(239):Article 2.
28. iGenomes: https://support.illumina.com/sequencing/sequencing_software/igenome.html.
29. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*. 2014; 30(15):2114–20. Epub 04/01. <https://doi.org/10.1093/bioinformatics/btu170> PMID: 24695404.
30. FastQC: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
31. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2012; 29(1):15–21. <https://doi.org/10.1093/bioinformatics/bts635> PMID: 23104886
32. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics*. 2011; 12(1):323. <https://doi.org/10.1186/1471-2105-12-323> PMID: 21816040
33. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology*. 2014; 15(12):550. <https://doi.org/10.1186/s13059-014-0550-8> PMID: 25516281
34. Gu Z, Eils R, Schlesner M. Complex heatmaps reveal patterns and correlations in multidimensional genomic data. *Bioinformatics*. 2016; 32(18):2847–9. <https://doi.org/10.1093/bioinformatics/btw313> PMID: 27207943
35. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*. 2005; 102(43):15545–50. <https://doi.org/10.1073/pnas.0506580102> PMID: 16199517
36. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. Molecular signatures database (MSigDB) 3.0. *Bioinformatics*. 2011; 27(12):1739–40. <https://doi.org/10.1093/bioinformatics/btr260> PMID: 21546393
37. Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, et al. The evolution of gene expression levels in mammalian organs. *Nature*. 2011; 478(7369):343–8. <https://doi.org/10.1038/nature10532> PMID: 22012392