

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

3

4 Xiaochuan Liu<sup>1</sup>, Jadwiga R Bienkowska<sup>2</sup>, and Wenyan Zhong<sup>1</sup>

5 <sup>1</sup>Oncology Research & Development, Pfizer Worldwide Research and Development, Pearl  
6 River, NY 10965, USA

7 <sup>2</sup>Oncology Research & Development, Pfizer Worldwide Research and Development, San Diego,  
8 CA 92121, USA

9

10 **Corresponding authors:**

11 [wenyan.zhong@pfizer.com](mailto:wenyan.zhong@pfizer.com)

12 [Jadwiga.R.Bienkowska@pfizer.com](mailto:Jadwiga.R.Bienkowska@pfizer.com)

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

## 29 **Abstract**

30

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

41

## 42 **Introduction**

43

44 Transposable elements (TEs) are mobile DNA sequences which have the capacity to  
45 move from one location to another on the genome[1]. TEs make up a considerable fraction of  
46 most eukaryotic genomes and can be classified into retrotransposons and DNA transposons  
47 according to their different mechanisms of transposition and chromosomal integration[2, 3].  
48 Retrotransposons are made of Long Terminal Repeats (LTRs) and non-LTRs that include long  
49 interspersed nuclear elements (LINEs) and short interspersed nuclear elements (SINEs) that  
50 mobilize via a RNA intermediate, while DNA transposons mobilize and function through a DNA  
51 intermediate[4-6]. TEs can be transcribed from the genome[7] and have been demonstrated to  
52 play important roles in the mammalian embryogenesis[8, 9], neurodevelopment[10, 11], and  
53 immune functions[12, 13]. Furthermore, aberrant expressions of TEs have been linked to

54 cancers[14-16], neurodegenerative disorders[17, 18], and immune-mediated inflammation[19,  
55 20]. Therefore, it has become increasingly important to explore biological roles of TEs  
56 expression. However, genome-wide analysis of TEs expression from high throughput RNA  
57 sequencing data has been a challenging computational problem. TEs contain highly repetitive  
58 sequence elements, making it arduous to unambiguously assign reads to the correct genomic  
59 location and accurately quantitate their expression level. Several bioinformatics tools have been  
60 developed to address this challenge with relatively good success [16, 21-23]. Recently, SQuIRE  
61 was reported to have the capability to quantify locus-specific expression of TEs from RNA-Seq  
62 data[23]. In addition, RNA-Seq data has long been used to detect dysregulated genes between  
63 different disease and/or drug treatment conditions to help understand disease mechanisms and/or  
64 drug response mechanisms. Therefore, it is of great interest to quantify both TEs and gene  
65 expression to elucidate contribution of both to disease mechanisms. Although many open source  
66 software and tools exist for analysing gene [24-26] and TEs expression, there are considerable  
67 challenges to efficiently apply these tools. In general, these multi-step data processing pipelines  
68 use many different tools. Correct versions of each tool need to be installed separately, and  
69 appropriate options, parameters, different reference genome and gene annotation files have to be  
70 set at each step. This can be quite tedious and challenging especially for non-computational  
71 users. Additionally, to ensure reproducibility of the analysis results, it is critical to capture  
72 analysis parameters from each step of the process. Equally important, to enable general use of  
73 the pipeline, the pipeline should be platform agnostic. Thus far, a one-stop computational  
74 framework for the comprehensive analysis of gene and locus-specific TEs expression from  
75 RNA-Seq data does not exist.

76           To address this need, we developed GeneTEFlow, a reproducible and platform-  
77 independent workflow, for the comprehensive analysis of gene and locus-specific TEs  
78 expression from RNA-Seq data using Nextflow[27] and Docker[28] technologies. GeneTEFlow  
79 provides several features and advantages for integrated gene and TEs transcriptomic analysis.  
80 First, by employing Docker technology, GeneTEFlow encapsulates bioinformatics tools and  
81 applications of specific versions into Docker containers enabling tracking, eliminating the need  
82 for software installation by users, and ensuring portability of the pipeline on multiple computing  
83 platforms including stand-alone workstations, high-performance computing (HPC) clusters, and  
84 cloud computing systems. Second, GeneTEflow uses Nextflow to define the computational  
85 workflows, not only enabling parallelization and complete automation of the analysis, but also  
86 providing capability to track analysis parameters. Thus, GeneTEFlow allows users to generate  
87 reproducible analysis results through utilization of both Docker and Nextflow in a platform  
88 independent manner. Lastly, GeneTEFlow has modular architecture, and modules in  
89 GeneTEFlow can be turned on or off, providing developers with flexibility to build extensions  
90 tailored to specific analysis needs.

## 91 **Implementation**

92  
93           The GeneTEFlow pipeline was developed using Nextflow, a portable, flexible, and  
94 reproducible workflow management system, and Docker technology, a solution to securely build  
95 and run applications on multiple platforms. To build the GeneTEFlow pipeline, a series of  
96 bioinformatics tools (S1 Table) were selected for QC, quantitation and differential expression  
97 analysis of genes and TEs from RNA-Seq data. These bioinformatics tools and custom scripts  
98 were built into four Docker containers to ensure portability of the workflow on different  
99 computational platforms. Data processing and analysis steps were implemented by modules

100 using Nextflow. Modules are connected through channels and can be run in parallel. Each  
101 module in GeneTEFlow can include any executable Linux scripts such as Perl, R, or Python.  
102 Parameters for each module are defined in a configuration file.

103 A conceptual workflow of GeneTEFlow is illustrated in Fig 1. The workflow includes  
104 four major inputs: raw sequence files in fastq format, a sample meta data file in excel format,  
105 reference genome and gene annotation files, and a Nextflow configuration file. The sample meta  
106 data file contains detailed sample information and the design of group comparisons between  
107 different experimental conditions. Human reference genome UCSC hg38 with the gene  
108 annotation (.gtf) was downloaded from Illumina iGenomes collections[29] and used by the  
109 bioinformatics tools included in GeneTEFlow. Scheduling of computational resources for each  
110 application module is defined in the configuration file.

111  
112 **Fig 1.** Illustration of GeneTEFlow: a Nexflow-based pipeline for identification of differentially  
113 expressed genes and locus specific transposable elements from RNA-Seq data.

114  
115 GeneTEFlow analysis is performed in following steps: QC, expression quantification,  
116 differential expression and down-stream analysis. First, adapter sequences are trimmed off from  
117 the Illumina raw reads using Trimmomatic(v0.36)[30] for single-end or paired-end reads, and  
118 low-quality reads are filtered out. Next, FastQC(v0.11.7)[31] is executed to survey the quality of  
119 sequencing reads, and report is generated to help identify any potential issues of the high  
120 throughput sequencing data. Reference genome index for mapping sequencing reads to mRNA  
121 genes is built using “rsem-prepare-reference” of RSEM (v.1.3.0). Reads remaining after the pre-  
122 processing step are mapped to the reference genome using STAR(v2.6.0c)[32]. Gene level

123 expression is quantitated as expected counts and transcripts per million (TPM) using “rsem-  
124 calculate-expression” of RSEM(v1.3.0) with default parameters [33]. Custom Perl scripts were  
125 developed to aggregate data from each sample into a single data matrix for expected counts and  
126 TPM values respectively. The expression quantification of locus-specific TEs is performed by  
127 SQuIRE[23].

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

135 Differential expression analysis of genes and transposable elements is performed using  
136 DESeq2(v1.18.1) package[34]. Significantly up-regulated and down-regulated genes and TEs are  
137 summarized in a table. To analyse overlap among significantly regulated genes and TEs from  
138 pair-wise comparisons between different sample groups we use Venn diagrams. We perform  
139 hierarchical clustering of significantly dysregulated genes or TEs using R package  
140 “ComplexHeatmap” [35] with euclidean distance and average linkage clustering parameters.  
141 Gene set enrichment analysis (GSEA, v3.0) [36] is conducted using collections from the  
142 Molecular Signatures Database (MSigDB) [37]. The outputs (S2 Table) from GeneTEFlow are  
143 organized into several folders predefined in a GeneTEFlow configuration file. A tutorial with  
144 detailed instructions on how to set up and run GeneTEFlow is provided at  
145 <https://github.com/zhongw2/GeneTEFlow>

## 146 **Application of GeneTEFlow**

147  
148 We applied GeneTEFlow to a public dataset from Brawan's study [38] investigating  
149 tissue-specific expression changes of genes and transposable elements. Human RNA-Seq data  
150 from brain, heart and testis tissues were downloaded from GEO (accession number: GSE30352)  
151 (S3 Table). Expression analysis of genes and TEs were performed using GeneTEFlow and  
152 results are shown in Fig 2. Gene expression analysis was performed using RSEM and DESeq2  
153 modules while TEs expression analysis was conducted using SQuIRE and DESeq2 modules  
154 within GeneTEFlow. Significantly regulated genes were identified with FDR less than 0.05 and  
155 fold change greater than 2. Significantly regulated locus-specific transposable elements were  
156 identified with FDR less than 0.05 and fold change greater than 1.5. The number of significantly  
157 regulated genes and transposable elements were summarized into two tables respectively (Fig 2,  
158 top panels). Using GeneTEFlow, we detected genes and TE differentially expressed between  
159 different tissue types (brain vs heart tissues: 6,264 genes and 1,277 TEs; testis vs heart tissues:  
160 7,066 genes and 595 TE; brain vs testis tissues: 8,125 genes and 1,297 TEs) with most  
161 significant gene and TE expression differences observed being between brain and testis tissues.  
162 Our analysis identified large number of both genes and TEs with tissue specific patterns (Fig 2,  
163 middle panels and bottom panels). More in depth analysis to include additional tissue types  
164 would be required to fully understand the tissue specific gene and TEs expression and their  
165 relationship. GeneTEFlow is a computational solution to facilitate such studies.

166  
167 **Fig 2.** Differential expression analysis results of genes and transposable elements from  
168 GeneTEFlow. Left panels: gene results; right panels: TEs results. Top panels: number of  
169 significantly regulated genes or TEs in each sample group comparison. Significance was defined

170 as following:  $FDR \leq 0.05$  and fold change  $\geq 2$  for gene expression analysis;  $FDR \leq 0.05$  and fold  
171 change  $\geq 1.5$  for TEs expression analysis. Middle panels: overlaps of significantly regulated  
172 genes or TEs amongst sample group comparisons. Bottom panels: hierarchical clustering of  
173 significantly regulated genes or TEs.

174

175 In addition to quantification of TEs expression, SQuIRE provides quantification of gene  
176 expression. Therefore, we compared gene level expression quantification between RSEM and  
177 SQuIRE (S1 Fig). The results showed high concordance (correlation coefficient:  $\sim 97\%$ ) of the  
178 gene level expression quantification between the two methods (S1 Fig, highlighted in red box)  
179 suggesting a robust measurement for both gene and TEs expression by SQuIRE.

180

## 181 **Conclusions**

182

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

191

## 192 **Competing interests**

193 WZ and JRB are employees of Pfizer Inc.

194 XL was contractor of Pfizer Inc. when the work was being conducted.



195  
196

## 197 **Funding**

198 Not applicable.

199

## 200 **Authors' contributions**

201  
202 WZ conceptualized the work. XL and WZ designed and implemented the pipeline. XL, JRB and

203 WZ drafted and revised the manuscript. All authors read and approved the final manuscript.

204

## 205 **Acknowledgements**

206  
207 We gratefully acknowledge inputs and support from our colleagues: Jeremy Myers, Keith Ching,  
208 Corey Dasilva and Da Tse.

209

## 210 **References**

- 211
- 212 1. Biémont C, Vieira C. Junk DNA as an evolutionary force. *Nature*. 2006;443(7111):521-4.  
213 doi: 10.1038/443521a.
  - 214 2. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, et al. A unified  
215 classification system for eukaryotic transposable elements. *Nature Reviews Genetics*.  
216 2007;8(12):973-82. doi: 10.1038/nrg2165.
  - 217 3. Feschotte C. Transposable elements and the evolution of regulatory networks. *Nature*  
218 *Reviews Genetics*. 2008;9(5):397-405. doi: 10.1038/nrg2337.
  - 219 4. Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M, et al. Ten  
220 things you should know about transposable elements. *Genome Biology*. 2018;19(1):199. doi:  
221 10.1186/s13059-018-1577-z.
  - 222 5. Lanciano S, Mirouze M. Transposable elements: all mobile, all different, some stress  
223 responsive, some adaptive? *Current Opinion in Genetics & Development*. 2018;49:106-14. doi:  
224 <https://doi.org/10.1016/j.gde.2018.04.002>.
  - 225 6. Chuong EB, Elde NC, Feschotte C. Regulatory activities of transposable elements: from  
226 conflicts to benefits. *Nature Reviews Genetics*. 2017;18(2):71-86. doi: 10.1038/nrg.2016.139.
  - 227 7. Rebollo R, Romanish MT, Mager DL. Transposable Elements: An Abundant and Natural  
228 Source of Regulatory Sequences for Host Genes. *Annual Review of Genetics*. 2012;46(1):21-42.  
229 doi: 10.1146/annurev-genet-110711-155621. PubMed PMID: 22905872.

- 230 8. Percharde M, Lin C-J, Yin Y, Guan J, Peixoto GA, Bulut-Karslioglu A, et al. A LINE1-  
231 Nucleolin Partnership Regulates Early Development and ESC Identity. *Cell*. 2018;174(2):391-  
232 405.e19. doi: <https://doi.org/10.1016/j.cell.2018.05.043>.
- 233 9. Garcia-Perez JL, Widmann TJ, Adams IR. The impact of transposable elements on  
234 mammalian development. *Development*. 2016;143(22):4101-14. doi: 10.1242/dev.132639.
- 235 10. Sun W, Samimi H, Gamez M, Zare H, Frost B. Pathogenic tau-induced piRNA depletion  
236 promotes neuronal death through transposable element dysregulation in neurodegenerative  
237 tauopathies. *Nature Neuroscience*. 2018;21(8):1038-48. doi: 10.1038/s41593-018-0194-1.
- 238 11. Guo C, Jeong H-H, Hsieh Y-C, Klein H-U, Bennett DA, De Jager PL, et al. Tau  
239 Activates Transposable Elements in Alzheimer's Disease. *Cell Reports*. 2018;23(10):2874-80.  
240 doi: <https://doi.org/10.1016/j.celrep.2018.05.004>.
- 241 12. Colombo AR, Elias HK, Ramsingh G. Senescence induction universally activates  
242 transposable element expression. *Cell Cycle*. 2018;17(14):1846-57. doi:  
243 10.1080/15384101.2018.1502576.
- 244 13. Koonin EV, Krupovic M. Evolution of adaptive immunity from transposable elements  
245 combined with innate immune systems. *Nature Reviews Genetics*. 2015;16(3):184-92. doi:  
246 10.1038/nrg3859.
- 247 14. Colombo AR, Triche T, Ramsingh G. Transposable Element Expression in Acute  
248 Myeloid Leukemia Transcriptome and Prognosis. *Scientific Reports*. 2018;8(1):16449. doi:  
249 10.1038/s41598-018-34189-x.
- 250 15. Burns KH. Transposable elements in cancer. *Nature Reviews Cancer*. 2017;17(7):415-24.  
251 doi: 10.1038/nrc.2017.35.
- 252 16. Criscione SW, Zhang Y, Thompson W, Sedivy JM, Neretti N. Transcriptional landscape  
253 of repetitive elements in normal and cancer human cells. *BMC Genomics*. 2014;15(1):583. doi:  
254 10.1186/1471-2164-15-583.
- 255 17. Krug L, Chatterjee N, Borges-Monroy R, Hearn S, Liao W-W, Morrill K, et al.  
256 Retrotransposon activation contributes to neurodegeneration in a Drosophila TDP-43 model of  
257 ALS. *PLOS Genetics*. 2017;13(3):e1006635. doi: 10.1371/journal.pgen.1006635.
- 258 18. Tam OH, Ostrow LW, Gale Hammell M. Diseases of the nERVous system:  
259 retrotransposon activity in neurodegenerative disease. *Mobile DNA*. 2019;10(1):32. doi:  
260 10.1186/s13100-019-0176-1.
- 261 19. De Cecco M, Ito T, Petrashen AP, Elias AE, Skvir NJ, Criscione SW, et al. L1 drives IFN  
262 in senescent cells and promotes age-associated inflammation. *Nature*. 2019;566(7742):73-8. doi:  
263 10.1038/s41586-018-0784-9.
- 264 20. Colombo AR, Elias HK, Ramsingh G. Senescence induction universally activates  
265 transposable element expression. *Cell cycle (Georgetown, Tex)*. 2018;17(14):1846-57. Epub  
266 08/16. doi: 10.1080/15384101.2018.1502576. PubMed PMID: 30080431.
- 267 21. Jin Y, Tam OH, Paniagua E, Hammell M. TETranscripts: a package for including  
268 transposable elements in differential expression analysis of RNA-seq datasets. *Bioinformatics*.  
269 2015;31(22):3593-9. doi: 10.1093/bioinformatics/btv422.
- 270 22. Lerat E, Fablet M, Modolo L, Lopez-Maestre H, Vieira C. TETools facilitates big data  
271 expression analysis of transposable elements and reveals an antagonism between their activity  
272 and that of piRNA genes. *Nucleic Acids Research*. 2016;45(4):e17-e. doi: 10.1093/nar/gkw953.
- 273 23. Yang WR, Ardeljan D, Pacyna CN, Payer LM, Burns KH. SQUIRE reveals locus-specific  
274 regulation of interspersed repeat expression. *Nucleic Acids Research*. 2019;47(5):e27-e. doi:  
275 10.1093/nar/gky1301.

- 276 24. Varet H, Brillet-Guéguen L, Coppée J-Y, Dillies M-A. SARTools: A DESeq2- and  
277 EdgeR-Based R Pipeline for Comprehensive Differential Analysis of RNA-Seq Data. PLOS  
278 ONE. 2016;11(6):e0157022. doi: 10.1371/journal.pone.0157022.
- 279 25. Finotello F, Di Camillo B. Measuring differential gene expression with RNA-seq:  
280 challenges and strategies for data analysis. Briefings in Functional Genomics. 2014;14(2):130-  
281 42. doi: 10.1093/bfgp/elu035.
- 282 26. Ewels PA, Peltzer A, Fillinger S, Patel H, Alneberg J, Wilm A, et al. The nf-core  
283 framework for community-curated bioinformatics pipelines. Nature Biotechnology.  
284 2020;38(3):276-8. doi: 10.1038/s41587-020-0439-x.
- 285 27. Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow  
286 enables reproducible computational workflows. Nature Biotechnology. 2017;35(4):316-9. doi:  
287 10.1038/nbt.3820.
- 288 28. Merkel D. Docker: lightweight Linux containers for consistent development and  
289 deployment. Linux J. 2014;2014(239):Article 2.
- 290 29. iGenomes : [https://support.illumina.com/sequencing/sequencing\\_software/igenome.html](https://support.illumina.com/sequencing/sequencing_software/igenome.html).
- 291 30. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence  
292 data. Bioinformatics. 2014;30(15):2114-20. Epub 04/01. doi: 10.1093/bioinformatics/btu170.  
293 PubMed PMID: 24695404.
- 294 31. FastQC : <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
- 295 32. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast  
296 universal RNA-seq aligner. Bioinformatics. 2012;29(1):15-21. doi:  
297 10.1093/bioinformatics/bts635.
- 298 33. Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with or  
299 without a reference genome. BMC Bioinformatics. 2011;12(1):323. doi: 10.1186/1471-2105-12-  
300 323.
- 301 34. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for  
302 RNA-seq data with DESeq2. Genome Biology. 2014;15(12):550. doi: 10.1186/s13059-014-  
303 0550-8.
- 304 35. Gu Z, Eils R, Schlesner M. Complex heatmaps reveal patterns and correlations in  
305 multidimensional genomic data. Bioinformatics. 2016;32(18):2847-9. doi:  
306 10.1093/bioinformatics/btw313.
- 307 36. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al.  
308 Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide  
309 expression profiles. Proceedings of the National Academy of Sciences. 2005;102(43):15545-50.  
310 doi: 10.1073/pnas.0506580102.
- 311 37. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP.  
312 Molecular signatures database (MSigDB) 3.0. Bioinformatics. 2011;27(12):1739-40. doi:  
313 10.1093/bioinformatics/btr260.
- 314 38. Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, et al. The  
315 evolution of gene expression levels in mammalian organs. Nature. 2011;478(7369):343-8. doi:  
316 10.1038/nature10532.

317  
318

## 319 **Supporting information**

320

321 **S1 Fig.** Comparison of gene expression quantification by RSEM and SQuIRE. Gene expression  
322 (total 22,955 genes) of samples from brain tissues (left), heart tissues (middle), and testis tissues  
323 (right) was calculated by both RSEM and SQuIRE. Lower diagonal panels: pairwise  
324 comparisons using  $\log_2(\text{TPM} + 1)$  of 22,955 genes. Upper diagonal panels: correlation  
325 coefficient of each comparison. Panels highlighted in red: correlation coefficient of comparisons  
326 between RSEM and SQuIRE gene expression quantification of the same sample. Rep\_ : replicate,  
327 \_RSEM: quantification performed by RSEM, \_SQuIRE: quantification performed by SQuIRE.  
328

329 **S1 Table.** Major bioinformatics tools installed in GeneTEFlow

330 **S2 Table.** Major outputs from GeneTEFlow

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

332 **S1\_File.** Supplemental tables: S1-S3 Tables.

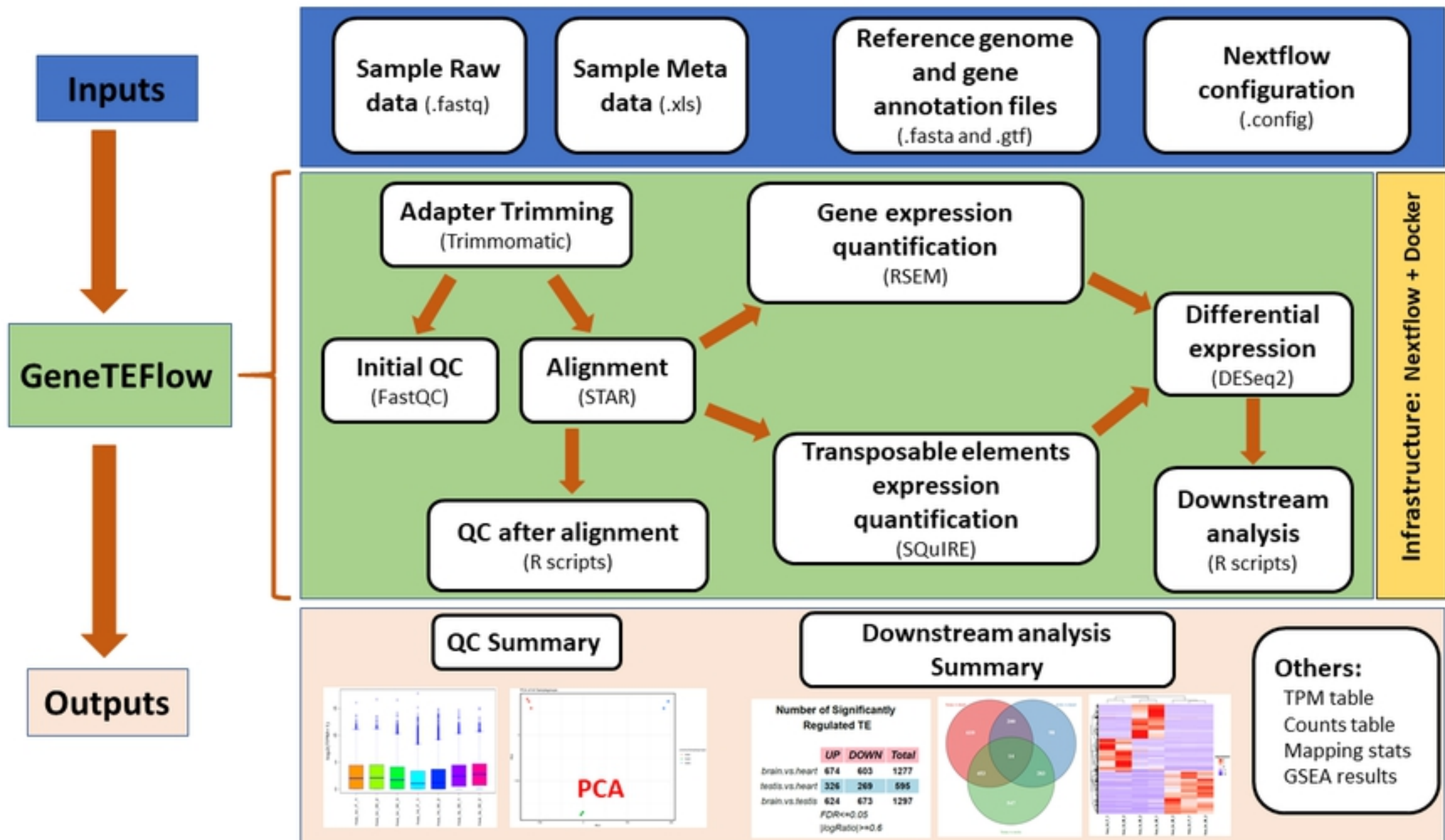


Fig 1

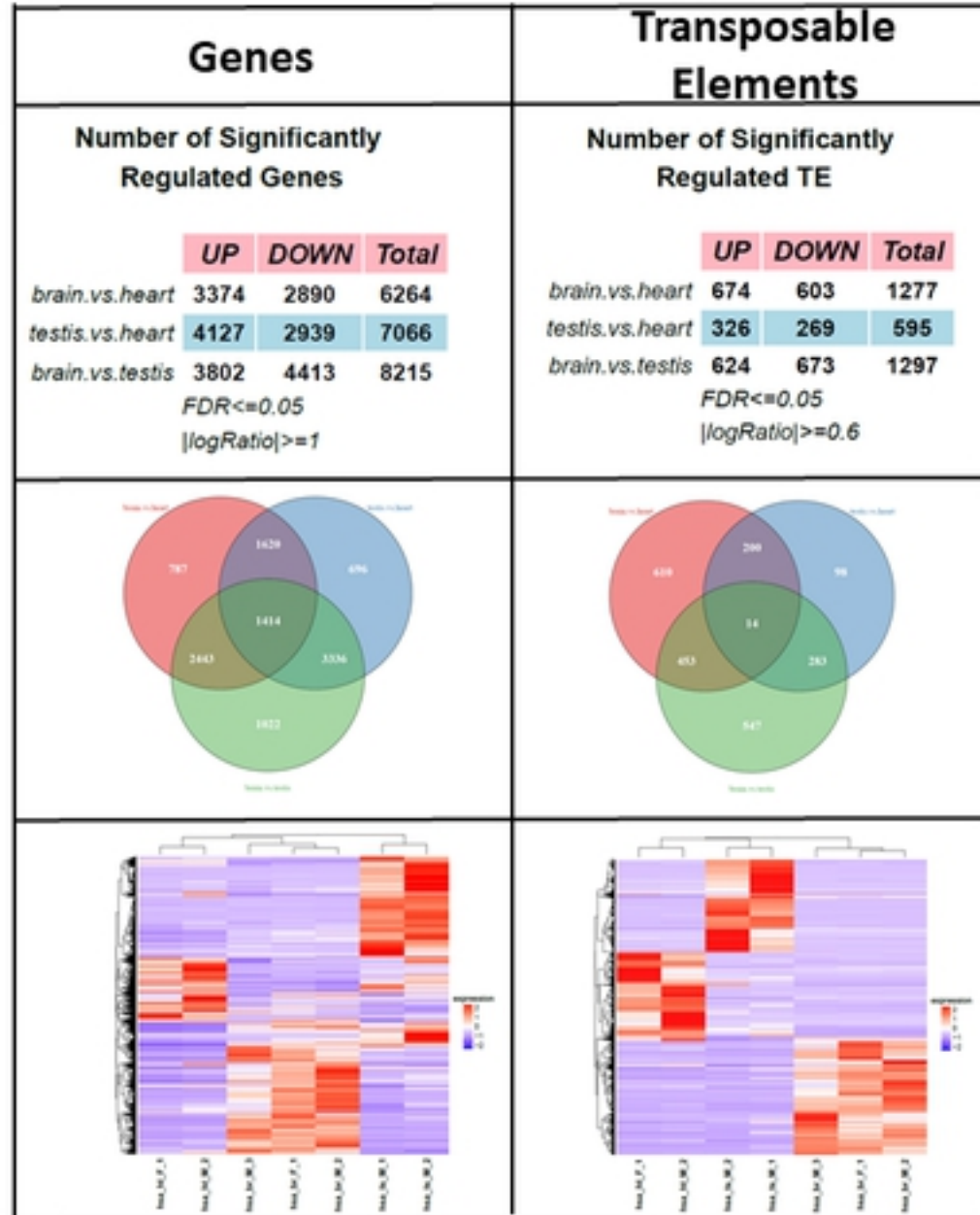


Fig 2