

1 Abstract

The recent explosion of single cell transcriptomics has led to the challenge of developing data analysis pipelines that are both fully reproducible and modular while allowing interoperability across multiple systems and institutions.

We present scAN10, a processing pipeline of 10X single cell RNAseq data, that inherits the ability to be executed on most computational infrastructures, thanks to Nextflow DSL2. The modular nature of Nextflow pipelines allows to easily integrate and assess different blocks for a given analysis step.

We illustrate the benefit of using scAN10 by showing its ability to assess the impact of the mapping step on the resulting output using a clinical 10X scRNAseq analysis of a human pituitary gonadotroph tumour.

29 **2 Introduction**

30 The recent explosion of single cell transcriptomics, mostly through scR-
 31 Naseq, has led to the challenge of developing data analysis pipelines that are
 32 both fully reproducible and modular while allowing interoperability across
 33 multiple systems and institutions.

34 The initial step in scRNAseq data analysis consist in generating a count
 35 matrix from **fastq** sequence files. This step is often overlooked (see e.g. [1]),
 36 although it can represent a very critical step (see e.g. [2]). Therefore one
 37 needs freely available analysis pipeline that could allow to verify the impact
 38 of some early analysis step easily, such as the nature of the **GTF** file used [3]
 39 or downstream steps such as the normalization method on the generation of
 40 the count matrix.

41 The most widely used existing solutions like the Seurat suite [4] have been
 42 designed to be as much user-friendly as possible and therefore does not offer
 43 easy solution for incorporating alternative low-level analysis steps.

44 This is why in the present paper we approached this challenge by designing
 45 **scAN10**, a processing pipeline of 10X single cell RNAseq data, that inherits
 46 the ability to be executed on most computational infrastructures, thanks to
 47 Nextflow DSL2. The modular nature of Nextflow pipelines allows to easily
 48 integrate and assess different bricks blocks for a given analysis step. **scAN10**
 49 is available as an open source Gitlab repository. It takes raw paired-end

50 **fastq** files and genomic files (FASTA and GTF annotation files) as input,
51 and outputs a clustered dimensionally-reduced version of the dataset.

52 We illustrate the benefit of using **scAN10** by showing its ability to assess the
53 impact of the mapping step on the resulting UMAP projection as well as on
54 some specific gene identification using a clinical 10X scRNA-seq analysis of
55 one human pituitary gonadotroph tumour.

56 **3 Results**

57 **3.1 Pipeline description**

58 Figure 1 describes the overall processing of the sequencing files with the
59 ordering of all steps described in section 5.

60 **3.2 Raw dataset**

61 To showcase the applicability of **scAN10** we processed a clinical dataset from
62 a human pituitary gonadotroph tumour acquired from one male patient and
63 sequenced by 10Xgenomics. The dataset was processed using **scAN10** with
64 the following parameters :

- 65 • `max_feature_RNA` = 7000
- 66 • `max_percent_mito` = 25

67 The clustering (resolution =0.7) and UMAP embedding were done based on

the ten first principal components of the PCA, as determined by the rule of thumb heuristic and the broken stick method [5].

3.3 Version annotation effect

We first assess the impact of the Ensembl version of the GTF file on the final output. GTF files and their corresponding FASTA files were downloaded via ftp protocol using the --version parameter (see section 5). We set Cellranger as the default mapper and then assessed the impact of 4 different annotation releases (93, 98, 103 and 106) on the number of detected genes (figure 2A) and on the count per genes as assessed with the CHGA gene. (Figure 2B).

The overall impact of the GTF version seems relatively modest, especially in regard to the number of counts. The 106 version allowed to identify a larger number of genes and was kept for the next step.

3.4 Comparing filtered with unfiltered annotations

We then assessed the impact of filtrating the GTF file with the mkgtf Cellranger function. This filtration step is intended to remove unwanted genes classified by biotype. We used the default values of that function that removes biotypes such as gene_biotype:pseudogene from the GTF annotation file.

This step removes some ambiguity between reads location by allowing reads that would be flagged as multi mapped reads to be included in the quantifi-

88 cation.

89 We observed that this filtration step indeed had a major impact on both the
90 number of genes detected (Figure 3A) but also on gene counts (Figure 3B).

91 **3.5 Cellranger versus Kallisto-bustools**

92 We then compared the impact of two popular alignment tools for single-cell
93 RNA sequencing (Kallisto-bustools and Cellranger) using the 10x Genomics
94 pre-built Cellranger reference packages version 2020-A for human.

95 As seen in Figure 4A, 81 % of the genes were identified by both algorithms
96 whereas Kallisto-bustools identified more genes than Cellranger.

97 The impact of the mapper on counts for specific genes seemed to be negligible
98 (Figure 4B). Therefore this tends to favor Kallisto-bustools for downstream
99 analyses.

100 We finally assessed the impact of the mapper choice on the final clustering
101 step. As seen in Figure 5, the impact was modest but apparent (e.g. cluster
102 number 1 in the CR dataset was split in two in the KB dataset).

103 **4 Discussion**

104 We described **scAN10**, a Nextflow based processing pipeline of 10X single cell
105 RNAseq data.

106 By applying `scAN10` to a clinical dataset we showed that the impact of the
107 annotation version was relatively modest although using the latest Ensembl
108 release (106) of the `GTF` and `FASTA` allows to identify a larger number of genes.

109 As expected, filtrating the `GTF` files by removing unwanted genes based on
110 10X reference packages generation had a major impact both on the number
111 of genes but also on gene counts.

112 However, recent study observed differences in the mitochondrial content of
113 the resulting cells when comparing a filtered annotation to the full annota-
114 tion. Therefore removing processed pseudogene might lead to an enrichment
115 of the mitochondrial content [6].

116 Futhermore when using Kallisto-bustools instead of Cellranger the impact
117 of the count numbers for specific genes seemed to be small but meaningful.
118 Kallisto-bustools produced higher total number of genes detected which 5169
119 unique genes as compared to Cellranger.

120 The final combination that was found to be the most effective for our dataset
121 therefore was using Kallisto-bustools together with the filtered version of
122 the 106 `GTF`. There is no reason to believe that such parameters might be
123 universally applicable, and we therefore recommend the use of `scAN10` so
124 assess such an impact on any other dataset before proceeding with higher
125 level analysis.

126 With Nextflow each step is encapsulated in independent blocks called pro-

cesses. Each process communicates via channels. The orchestration of the workflow with DSL2 syntax allows to easily modify the pipeline, by adding new processes modules. In the future we expect to include some normalization procedures to the pipeline:

- The basic global-scaling normalization method from Seurat that divides the feature expression for each cell by the total expression and multiplies this by a scale factor and log-transforms the result.[7]
- **Sctransform** which uses regularized negative binomial regression and computes Pearson residuals that correspond to the normalized expression levels for each transcript.[8]
- **Scran** which uses pooling-based size factor estimation.[9]

The use of **scAN10** should be made straightforward to assess a combination of low-level steps together with the normalization step on the resulting output.

One major limitation using Nextflow is a lack of interactivity during pipeline running. Indeed, when Nextflow pipelines run, although it can output files, there is no interactive dialogue that could allow the user to modify parameters during the run. All the pipelines parameter need to be defined and set in the launching command.

We finally believe that **scAN10** will be a useful tool for the growing community of 10X scRNAseq *aficionados*.

147 **5 Material and methods**

148 **5.1 Study ethic approval**

149 This work is part of the SPACE-PIT study (MR004 n21-5439). It was ap-
 150 proved by the Hospices Civils de Lyon ethical committee and registered at
 151 the “Centre National Information et Liberté” (CNIL.fr) under the reference
 152 20_098. Informed consent was obtained from the patient.

153 **5.2 Single cells preparation and sequencing**

154 A tumor fragment from a gonadotroph surgically-resected adenoma was col-
 155 lected in Dulbecco’s Modified Eagle Medium (DMEM, cat 41965062; Life
 156 Technologies). Single-cell suspension of the resected fragment was obtained
 157 through mechanical enzymatic dissociation (Collagenase P, cat 11213865001)
 158 then passed through a 100 μ m mesh-strainer (#732-2759, VWR interna-
 159 tional).

160 Red blood cells were eliminated using a 10-minute incubation with a com-
 161 mercial red blood cell lysis buffer (eBioscience, cat #00-4300-54). The whole
 162 process was achieved within the 2 hours following the surgical resection,
 163 cell viability was evaluated to reach at least 70 percent prior encapsula-
 164 tion. Generation of the library was done using a Chromium controller from
 165 10xGenomics. The entire procedure was achieved as recommended by the
 166 manufacturer’s for the v3 reagent kit (10xGenomics). Single cell suspension

167 was loaded onto a Chromium Single Cell A Chip, aiming for 5,000 cells. The
 168 cDNA was amplified after a reverse transcription step prior a SPRIselect
 169 (Beckman Coulter), a cleaning, a quantification and an enzymatic fragmen-
 170 tation prior to the library sequencing on a NextSeq500 system (Illumina).

171 **5.3 Implementation**

172 The `scAN10` pipeline was powered and supported by the reactive workflow
 173 manager Nextflow [10]. In addition, the pipeline was coded with the DSL2
 174 syntax extension. Nextflow simplifies the writing of computational pipelines
 175 by making them portable, scalable, parallelizable and ensuring a high level of
 176 reproducibility. Nextflow provides native support for container technologies
 177 such as Docker or Singularity. Each process in the pipeline will be run in
 178 a container. A reproducible container environment is built for each process
 179 from Docker images stored on the DockerHub. The analyses can be run
 180 on the user's preferred computing platform. Using the configuration file and
 181 corresponding profile, the pipeline can be run on a local computer via Docker
 182 or Singularity, as well as on a high-performance computing (HPC) cluster
 183 or in cloud-based environments. Nextflow includes a cache-based pipeline
 184 resume feature, no matter what the reason was for its stopping.

185 **5.4 Input**

186 **scAN10** takes three mandatory parameters as input: the paired-end **FASTQ** R1
187 and R2 from 10X Chromium sequencing and two genomics files (one **FASTA**
188 an one **GTF**).

189 **FASTQ** files store the nucleotide sequence and the associated sequencing qual-
190 ity scores. Those files must be provided through the `--fastq` pipeline pa-
191 rameter and needs to be compressed in **gzip** format. Reads are sequenced
192 in paired-ends thus 2 reads will be produced for each sequencing lane. In
193 the case a sample has been sequenced on several lanes, all reads R1 can be
194 concatenate together and all reads R2 together. The mapping step require
195 the input of two additional files corresponding to the species of interest:

- 196 • A **FASTA** genomic file which stores the raw genome sequence.
- 197 • A Gene transfer format (**GTF**) file which stores genome annotation in-
198 cluding gene positions.

199
200 One should note that for human datasets, **FASTA** and **GTF** files can be down-
201 loaded automatically by specifying a version number as an entry parameter
202 (`--version`) available on the ENSEMBL database.

203 5.5 Preprocessing & Mapping

204 First, **FASTQ** files are processed and trimmed by using **Fastp** v0.20.1 , an
 205 ultra-fast **FASTQ** preprocessor with useful quality control and data-filtering
 206 features [11]. Reads with phred quality ≥ 30 (-q 30) is qualified to the quan-
 207 tification step. Length filtering is disabled (-L) while adapter sequence auto-
 208 detection is enabled (--detect_adapter_for_pe).

209 To reduce overlapping annotation, we recommend and provide an optional
 210 parameter --filtergtf allow the **GTF** filtration with Cellranger's **mkgtf** function
 211 based on the same biotype attribute used to generate the **GTF** file for the hu-
 212 man Cell Ranger reference package. (see : <https://support.10xgenomics.com/single-cell-gene-expression/software/release-notes/build>)
 213

214 The processed files are then mapped of a reference genome to quantify gene
 215 expression. In **scAN10**, the user can specify two different mappers (see below):
 216 **Kallisto-bustools** v0.26.0 [12] or **Cellranger** v5.0.1 [13].

- 217 • **Kallisto-bustools** is used through the Python wrapper : **kb_python**.
 218 Starting with **FASTA** and **GTF** files an index of the reference can be built
 219 as a colored De Bruij graph with **Kallisto** via **kb ref.** with default pa-
 220 rameter. Once an index has been generated or downloaded, **kb count**
 221 uses **Kallisto** to pseudoalign reads and **bustools** to quantify the data.
- 222 • **Cellranger** created and prepared reference package with **Cellranger's**
 223 **mkref** function. The alignment was run via **Cellranger's** **count** with

224 default parameters as described on `10xgenomics.com`.

225 5.6 Quality control

226 **Empty Droplets** With Droplet-based data most of the barcodes in the
 227 matrix correspond to empty droplets (eg barcode with sum expression over
 228 gene of 0). While Cellranger takes care of the empty droplet filtering, empty
 229 droplets from Kallisto-bustools gene expression matrix need to be removed.
 230 The Kallisto-bustools outputs were imported into R with customized R func-
 231 tion. The UMI total counts were ranked using `DropletUtils::barcodeRanks`
 232 function from `DropletUtils` v1.14.2 . Empty droplets were removed by se-
 233 lecting the inflection point value on the resulting knee plot (lower cutoff =
 234 10). The Cellranger matrix is imported from the standard filtered barcode
 235 output. After importation, either gene expression matrices were converted as
 236 Seurat object (Seurat v4.0.4) with `Seurat::CreateSeuratObject` function,
 237 including features detected in at least 3 cells (`min.cells = 3`).

238 **Low quality cells** After removing empty droplets the pipeline computes 3
 239 QC metrics per sample: the number of unique features per cells, the number
 240 of UMIs by cells and the percent of mitochondrial counts by cells. These QC
 241 metrics are then used to discard three main types of low quality cells [14] :

- 242 1. Cells in apoptosis may exhibit high % mitochondrial and low number
 243 of UMIs per cells

- 244 2. Cells that failed during library preparation exhibit low number of unique
245 gene counts and low number of UMIs per cells
- 246 3. The pipeline uses the R package `DoubletFinder` v2.0.3 to detect and
247 remove the potential doublets from the dataset. See [15] and [16].

248 **Thresholding** Default thresholding parameters used by the `scAN10` pipeline
249 are:

- 250 1. `min_feature_RNA=500`
- 251 2. `min_ncount_RNA=0`
- 252 3. `max_percent_mito="adaptive"`
- 253 4. `max_feature_RNA="adaptive"`
- 254 5. `max_ncount_RNA="adaptive"`

To define maximum values of threshold `scAN10` allows to use an adaptive filter defined by a certain number of median absolute deviations (MADs) away from the median [17]:

$$median + 3 * mad$$

255 **Non-expressed genes** Genes with sum count along cells equal to 0 (eg
256 not-expressed genes) are removed.

257 5.7 Normalization

258 To normalize the filtered gene expression matrix, the user can use **Sanity**,
 259 a Bayesian algorithm to infer gene-expression state [18]. We are fully aware
 260 that the normalization of single cell transcriptomic data is a research field on
 261 its own (see e.g. [2] and citations therein), and we expect this block in the
 262 pipeline to be susceptible to be modified in future versions of the pipeline.
 263 The modularity of the Nextflow syntax makes it ideal for such additions.

264 5.8 Clustering and two-dimensional visualization

265 A final step consists in variable features selection with **Seurat::FindVariableFeatures**
 266 using the **vst** method (`selection.method = "vst"`) and selecting the 2000 first
 267 highly variable genes (`nfeatures = 2000`), followed by a first linear dimension-
 268 ality reduction using PCA (**Seurat::RunPCA** with default parameter). The
 269 M first axis of the PCA are then used for the nearest-neighbor graph con-
 270 struction with **Seurat::FindNeighbors** function (`dims=1:M`).

271 Cluster determination was performed using the Louvain algorithm [19] run
 272 with **Seurat::FindClusters** function (`algorithm = 1`). The resolution pa-
 273 rameter set by default to 0.7 is use for increasing (values >1) or decreasing
 274 (values <1) the number of clusters obtained. The quality of the clustering
 275 was assessed using the Silhouette score [20]. Finally, non-linear dimensional
 276 representation (using either t-SNE [21] or UMAP [22]) is then performed
 277 using **Seurat::RunTSNE** or **Seurat::RunUMAP** with default parameters using

the same number of dimensions than the nearest-neighbor graph building.

The resolution parameters used for clustering and the number of principal components kept for both clustering and dimension reduction embeddings can be respectively modified by the user at the start of the pipeline by the `--resolution_clustering` and the `--principal_component` parameters.

Alternatively this last step can be skipped allowing the user to use their own clustering method. Similarly, to avoid the introduction of layers of complexity and simplify the pipeline usage, the automatic annotation of clusters was not introduced. Users can annotate their dataset manually.

5.9 Output

Exhaustive list of processes outputs is available on the Readme of the gitlab repository (section 6).

6 Availability

scAN10 is freely available at : <https://gitbio.ens-lyon.fr/LBMC/sbdlm/>
`scan10`

7 Acknowledgements

We thank the Institut Convergence Plascan (Grant Number ANR-17-CONV-0002) for their support.

296 Mirela Diana Ilie has been supported by the Fondation ARC pour la recherche
297 sur le cancer (ARCMD-DOC12020020001361).

298 We thank (i) Emmanuel Jouanneau and Alexandre Vasiljevic from the Neu-
299 rosurgery and Pathology Departments, Reference Center for Rare Pituitary
300 Diseases HYPO, “Groupement Hospitalier Est ”Hospices Civils de Lyon,
301 Bron, France for providing diagnosed Gonadotroph tumor material; (ii) Hec-
302 tor Hernandez-Vargas for his advices and support with setting the single cell
303 dissociation and subsequent bioinformatic analysis; and (iii) Laurent Modolo
304 from the LBMC for his help using the Nextflow syntax.

305 We thank the Pôle Scientifique de Modélisation Numérique (PSMN, Ecole
306 Normale Supérieure de Lyon) where computations were performed.

307 We finally thank the BioSyL Federation (<http://www.biosyl.org>), the LabEx
308 Ecofect (ANR-11-LABX-0048) and the LabEx Milyon of the University of
309 Lyon for inspiring scientific events.

310 **8 Author contributions**

311 Maxime Lepetit wrote the scAN10 script, performed the analysis, generated
312 the figures and wrote the paper. Mirela Diana Ilie participated in generating
313 and analyzing the 10X data. Marie Chanal participated in generating the
314 10X data. Gerald Raverot participated in generating and analyzing the 10X
315 data. Philippe Bertolino helped designing the study, analyzing the results,

316 reviewed and edited the manuscript, and secured the funding. Franck Picard
317 helped designing the study and analyzing the results, reviewed and edited the
318 manuscript, and secured the funding. Olivier Gandrillon helped designing the
319 study and analyzing the results, wrote the paper, and secured the funding.

320 9 References

- 321 [1] M. D. Luecken and F. J. Theis. “Current best practices in single-cell
322 RNA-seq analysis: a tutorial”. In: Mol Syst Biol 15.6 (2019), e8746.
- 323 [2] R. Zhang, G. S. Atwal, and W. K. Lim. “Noise regularization removes
324 correlation artifacts in single-cell RNA-seq data preprocessing”. In:
325 Patterns (N Y) 2.3 (2021), p. 100211.
- 326 [3] A.-H. Pool, H. Poldsam, S. Chen, M. Thomson, and Y. Oka. “En-
327 hanced recovery of single-cell RNA-sequencing reads for missing gene
328 expression data”. In: bioRxiv (2022), p. 2022.04.26.489449.
- 329 [4] Y. Hao, S. Hao, E. Andersen-Nissen, 3rd Mauck W. M., S. Zheng,
330 A. Butler, M. J. Lee, A. J. Wilk, C. Darby, M. Zager, P. Hoffman,
331 M. Stoeckius, E. Papalexi, E. P. Mimitou, J. Jain, A. Srivastava, T.
332 Stuart, L. M. Fleming, B. Yeung, A. J. Rogers, J. M. McElrath, C. A.
333 Blish, R. Gottardo, P. Smibert, and R. Satija. “Integrated analysis of
334 multimodal single-cell data”. In: Cell 184.13 (2021), 3573–3587 e29.

- 335 [5] D. A. Jackson. “Stopping Rules in Principal Components Analysis: A
336 Comparison of Heuristical and Statistical Approaches”. In: Ecology 74
337 (1993), pp. 2204–2214.
- 338 [6] Ralf Schulze Brüning, Lukas Tombor, Marcel H Schulz, Stefanie Dim-
339 meler, and David John. “Comparative analysis of common alignment
340 tools for single-cell RNA sequencing”. In: GigaScience 11 (Jan. 2022).
341 giac001. eprint: [https://academic.oup.com/gigascience/article-](https://academic.oup.com/gigascience/article-pdf/doi/10.1093/gigascience/giac001/42297447/giac001.pdf)
342 [pdf/doi/10.1093/gigascience/giac001/42297447/giac001.pdf](https://academic.oup.com/gigascience/article-pdf/doi/10.1093/gigascience/giac001/42297447/giac001.pdf).
- 343 [7] R. Satija, J. Farrell, and D. Gennert. “Spatial reconstruction of single-
344 cell gene expression data.” In: Nat Biotechnol 33.5 (2015), pp. 495–
345 502.
- 346 [8] C. Hafemeister and R. Satija. “Normalization and variance stabiliza-
347 tion of single-cell RNA-seq data using regularized negative binomial
348 regression.” In: Genome Biol 20.1 (2019), p. 296.
- 349 [9] A.T. L. Lun, K. Bach, and J.C. Marioni. “Pooling across cells to nor-
350 malize single-cell RNA sequencing data with many zero counts.” In:
351 Genome Biol 17.1 (2019), p. 75.
- 352 [10] P. Di Tommaso, M. Chatzou, E. W. Floden, P. P. Barja, E. Palumbo,
353 and C Notredame. “Nextflow enables reproducible computational work-
354 flows”. In: Nature biotechnology 35.4 (2017), pp. 316–319.
- 355 [11] S. Chen, Y. Zhou, Y. Chen, and J. Gu. “fastp: an ultra-fast all-in-one
356 FASTQ preprocessor”. In: Bioinformatics 34.17 (2018), pp. i884–i890.

- 357 [12] P. Melsted, A. S. Booesbaghi, F. Gao, E. Beltrame, L. Lu, K. E.
358 Hjorleifsson, J. Gehring, and L. Pachter. “Modular and efficient pre-
359 processing of single-cell RNA-seq”. In: Nat. Biotechnol. 39 (2021), pp. 813–
360 818.
- 361 [13] G. X. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R.
362 Wilson, S. B. Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu, M. T.
363 Gregory, J. Shuga, L. Montesclaros, J. G. Underwood, D. A. Masque-
364 lier, S. Y. Nishimura, M. Schnall-Levin, P. W. Wyatt, C. M. Hindson,
365 R. Bharadwaj, A. Wong, K. D. Ness, L. W. Beppu, H. J. Deeg, C. Mc-
366 Farland, K. R. Loeb, W. J. Valente, N. G. Ericson, E. A. Stevens, J. P.
367 Radich, T. S. Mikkelsen, B. J. Hindson, and J. H. Bielas. “Massively
368 parallel digital transcriptional profiling of single cells”. In: Nat Commun
369 8 (2017), p. 14049.
- 370 [14] T. Ilicic, J.K. Kim, A.A. Kolodziejczyk, F.O. Bagger, D.J. McCarthy,
371 J.C. Marioni, and Teichmann S.A. “Classification of low quality cells
372 from single-cell RNA-seq data.” In: Genome Biology 17.1 (2016), p. 29.
- 373 [15] C.S. McGinnis, L.M. Murrow, and Z.J. Gartne. “DoubletFinder: dou-
374 blet detection in single-cell RNA sequencing data using artificial nearest
375 neighbors.” In: Cell systems 8.4 (2019), 329–337. e4.
- 376 [16] N.M. Xi and J.J. Li. “Benchmarking Computational Doublet-Detection
377 Methods for Single-Cell RNA Sequencing Data”. In: Cell Systems 12.2
378 (2021), 176–194.e6.

- [17] D. J. McCarthy, K. R. Campbell, A. T. L. Lun, and Q. F. Wills. “Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R”. In: Bioinformatics 33.8 (Jan. 2017), pp. 1179–1186. eprint: <https://academic.oup.com/bioinformatics/article-pdf/33/8/1179/25150420/btw777.pdf>.
- [18] J. Breda, M. Zavolan, and E. van Nimwegen. “Bayesian inference of gene expression states from single-cell RNA-seq data”. In: Nat Biotechnol (2021).
- [19] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. “Fast unfolding of communities in large networks”. In: J. Stat. Mech. P10008 (2008), p. 12.
- [20] Peter J. Rousseeuw. “Silhouettes: A graphical aid to the interpretation and validation of cluster analysis”. In: Journal of Computational and Applied Mathematics 20 (1987), pp. 53–65.
- [21] L. Van der Maaten and G. Hinton. “Visualizing data using t-SNE”. In: Journal of Machine Learning Research 9 (2008), pp. 2579–2605.
- [22] L. McInnes, J. Healy, N. Saul, and L. Großberger. “UMAP: uniform manifold approximation and projection”. In: J. Open Source Softw 3 (2018), p. 861.

10 Figures caption

Figure 1 : A metromap view of the scAN10 pipeline.

400 Figure 2 : A. Number of detected genes when using different versions of the
401 GTF file. B. Violin plot representation of the impact of the GTF version on
402 the UMI counts for the CHGA gene.

403 Figure 3 : A. Number of detected genes when using either an unfiltered (106)
404 of filtered (filter106) version of the GTF file. B. Violin plot representation of
405 the filtration impact on the UMI counts for the CD68 gene.

406 Figure 4 : A. Number of detected genes when using either Cellranger (CR)
407 or Kallisto-bustools (KB) as an alignment tool. B. Violin plot representation
408 on the UMI counts for two genes, CHGA and RBP4.

409 Figure 5 : UMAP representation (A and B) and Silhouette scores (C and
410 D) of the clusters obtained on data processed with CellRanger (A and C)
411 or Kallisto-bustools (B and D). In E is shown an alluvial plot highlighting
412 the conservation and differences in cluster composition depending upon the
413 initial mapping method.

11 Figures

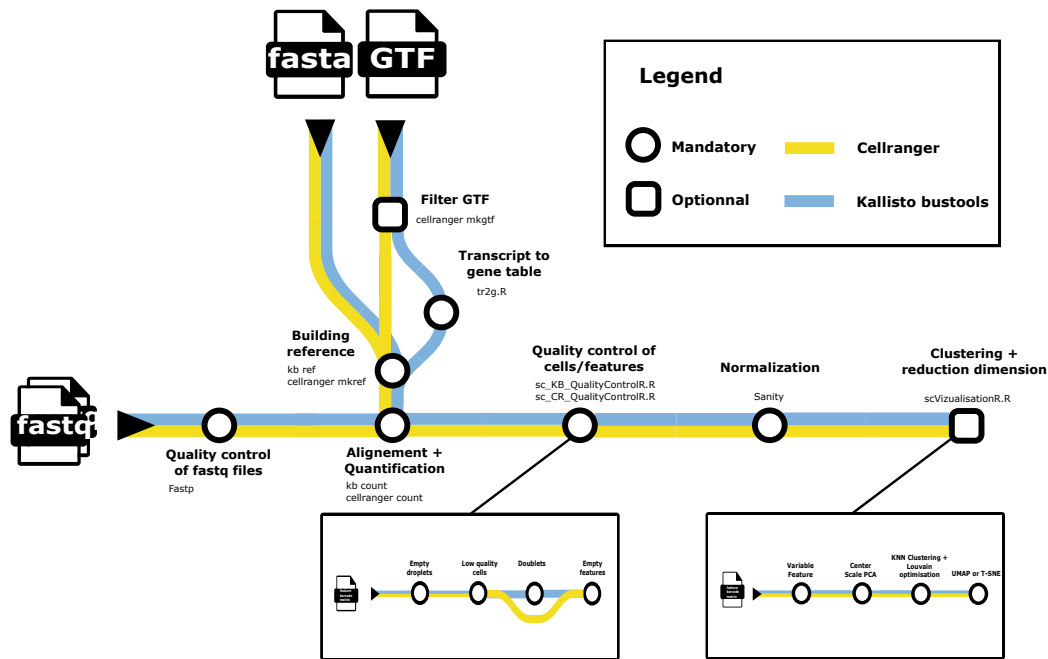


Figure 1

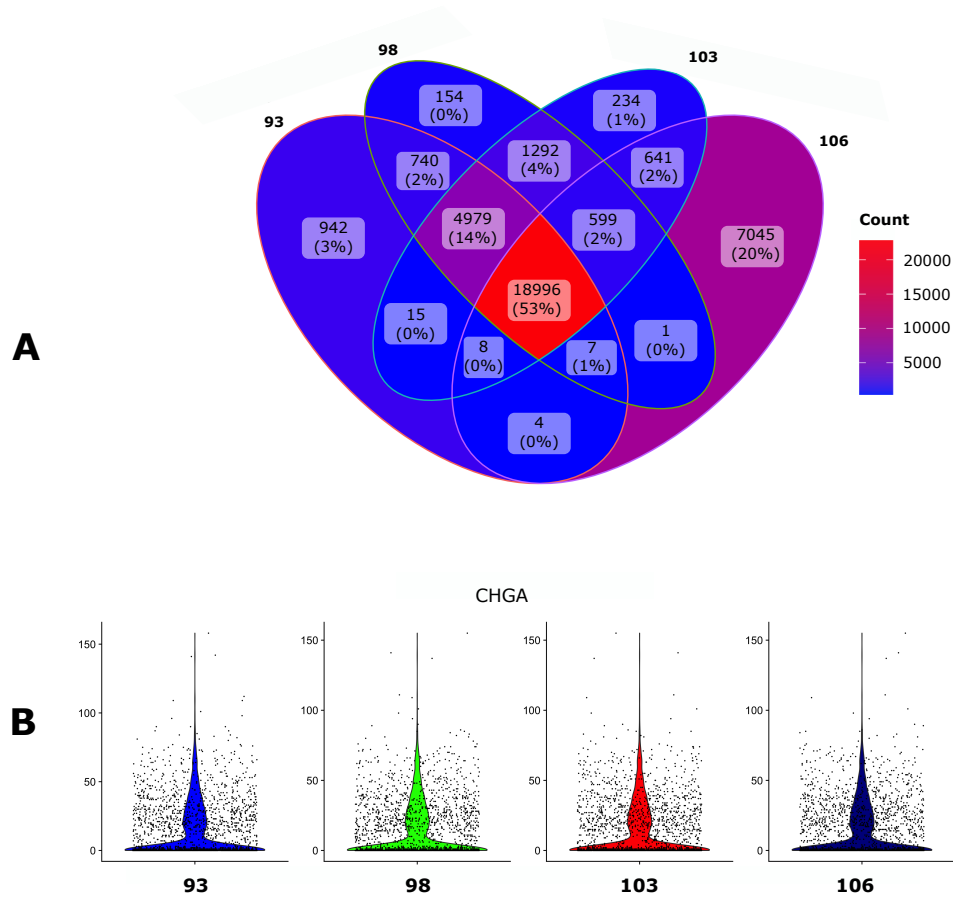
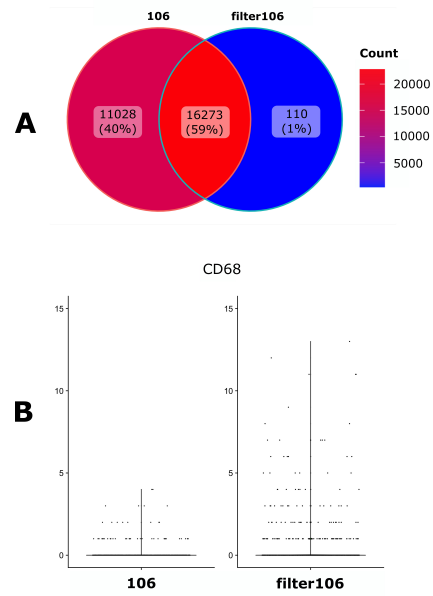


Figure 2



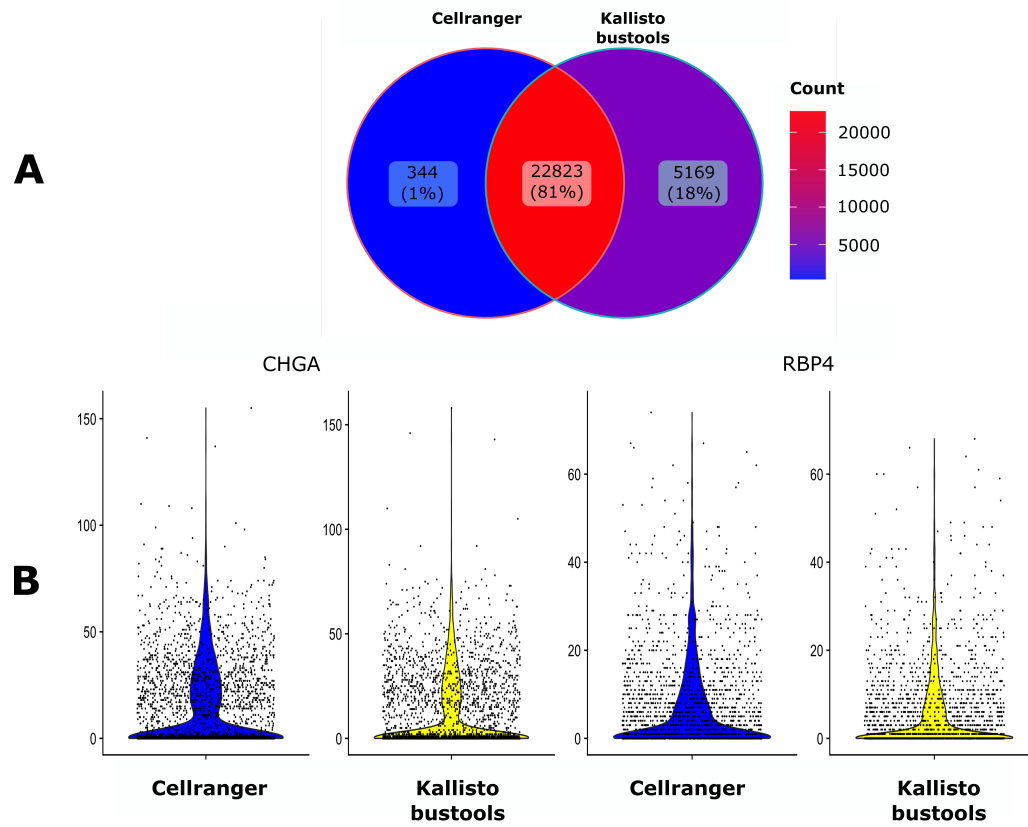


Figure 4

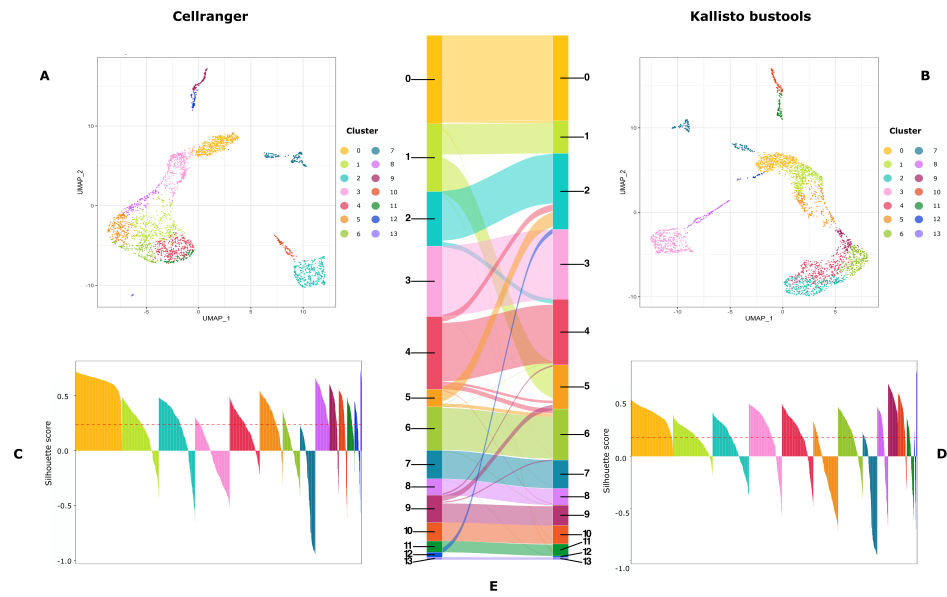


Figure 5