

Check for



Citation: Soneson C, Srivastava A, Patro R, Stadler MB (2021) Preprocessing choices affect RNA velocity results for droplets cRNA-seq data. PLoS Comput Biol 17(1): e1008585. https://doi.org/10.1371/journal.pcbi.1008585

Editor: Min Li, Central South University, CHINA

Received: April 6, 2020

Accepted: November 30, 2020

Published: January 11, 2021

Copyright: 0 2021 Soneson et al. This is an open access article distributed under the terms of the Creative Commons Att ribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Avai lability Statement: The spermatogenesis data can be downloaded from GEO, accession number GSE109033. The pancreas data set was downloaded from GEO, accession number GSE132188. The dentate gyrus data set was downloaded from GEO, accession number GSE96315. The PFC data set is accessible from GEO under accession number GSE124952 (sample accession number GSI/N3569979). The OldB rain data set is accessible from GEO under accession number GSI/N3722108. The code used for our analyses is available on https://github.com/csoreson/ma_velocity_quant.

RESEARCH ARTICLE

Preprocessing choices affect RNA velocity results for droplet scRNA-seq data

Charlotte Soneson 51.2*, Avi Srivastava 53.4, Rob Patro 5, Michael B. Stadler 51.25

- 1 Friedrich Miescher Institute for Biomedical Research, Basel, Switzerland, 2 SIB Swiss Institute of Bioinformatics, Basel, Switzerland, 3 New York Genome Center, New York, United States of America, 4 Center for Genomics and Systems Biology, New York University, New York, United States of America, 5 Department of Computer Science, University of Maryland, College Park, Maryland, United States of America, 6 University of Basel, Basel, Switzerland
- * charlotte.soneson@fmi.ch

Abstract

Experimental single-cell approaches are becoming widely used for many purposes, including investigation of the dynamic behaviour of developing biological systems. Consequently, a large number of computational methods for extracting dynamic information from such data have been developed. One example is RNA velocity analysis, in which spliced and unspliced RNA abundances are jointly modeled in order to infer a 'direction of change' and thereby a future state for each cell in the gene expression space. Naturally, the accuracy and interpretability of the inferred RNA velocities depend crucially on the correctness of the estimated abundances. Here, we systematically compare five widely used quantification tools, in total yielding thirteen different quantification approaches, in terms of their estimates of spliced and unspliced RNA abundances in five experimental dropletscRNA-seg data sets. We show that there are substantial differences between the quantifications obtained from different tools, and identify typical genes for which such discrepancies are observed. We further show that these abundance differences propagate to the downstream analysis, and can have a large effect on estimated velocities as well as the biological interpretation. Our results highlight that abundance quantification is a crucial aspect of the RNA velocity analysis workflow, and that both the definition of the genomic features of interest and the quantification algorithm itself require careful consideration.

Author summary

Applied to single-cell RNA-seq data, RNA velocity analysis provides a way to estimate the rate of change of the gene expression levels in individual cells. This, in turn, enables estimation of what the gene expression profile of each cell will look like a short time into the future and lets researchers infer likely developmental relationships among different types of cells in a tissue. An important first step in this type of analysis consists of estimating the expression levels of unspliced pre-mRNA as well as mature mRNA in each cell. Several methods are available for this purpose, and in this study we perform a comparison of these tools and highlight respective advantages and disadvantages. We envision that the

Funding: This work has been supported by R01 HG009937 (US National Institutes of Health, https://www.min.gov/), NSF COF-1750472 (US National Science Foundation, https://www.nsf. gov/), and CNS-1763680 (US National Science Foundation, https://www.nsf.gov/) to R.P. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: Thave read the journal's policy and the authors of this manuscript have the following competing interests: R.P. is a co-founder of Ocean Genomics Inc. results will be informative for researchers performing RNA velocity analysis, as well as to guide future developments in the evaluated and newly proposed quantification methods.

This is a PLOS Computational Biology Benchmarking paper.

Introduction

Single-cell RNA-seq (scRNA-seq) enables high-throughput profiling of gene expression on a transcriptome-wide scale in individual cells [1–6]. The increased resolution compared to bulk RNA-seq, where only average expression profiles across populations of cells are obtained, provides vastly improved potential to study a variety of biological questions. One such question concerns the dynamics of biological systems, reflected in, e.g., cellular differentiation and development [7]. While such dynamical processes would ideally be studied via repeated transcriptome-wide expression profiling of the same cells over time, this is not possible with current scRNA-seq protocols. Existing analysis methods for so called trajectory inference are instead applied to one or several snapshots of a population of cells, assumed to comprise all stages of the trajectory of interest. Many computational methods for trajectory inference from scRNA-seq have been presented in the literature (reviewed by [8]). These methods typically use the similarity of the gene expression profiles between cells to construct a (possibly branching) path through the observed set of cells, representing the trajectory of interest. Projecting the cells onto this path then provides an ordering of the cells by so called pseudotime.

A different approach to the investigation of developmental processes in scRNA-seq data instead exploits the underlying molecular dynamics. The feasibility of such an approach is based on the observation that, with several commonly used library preparation protocols, not only exonic, but also intronic and exon/intron boundary-spanning reads are obtained [9], and the insight that considering these in combination with the exonic reads would allow for direct inference of developmental relationships among cells. Similar observations, coupled with a simple differential equation model of transcriptional dynamics, were previously used for investigation of pre-mRNA dynamics and transcriptional and post-transcriptional regulation of gene expression in exon arrays [10] and bulk RNA-seq [11], as well as estimation of transcription, processing and degradation rates in bulk RNA-seq [12]. For scRNA-seq, La Manno et al [9] used the differential equation model of [10], describing the rate of change of unspliced premRNA as well as spliced mRNA molecules, as a basis for their investigations. They defined the RNA velocity for a given gene in a given cell, at a given time point, as the instantaneous rate of change of the spliced mRNA abundance. Combining the RNA velocities with the estimated mRNA abundances enables reconstruction of the state of each cell at a timepoint in the near future. With the increased popularity of RNA velocity applications in scRNA-seq studies, several computational tools have been developed, both for the preprocessing of the reads and for the RNA velocity estimation. The original velocyto software [9] estimates velocities under a steady-state assumption, and provides both Python and R implementations. More recently, Bergen et al [13] relaxed the steady-state assumption and considered the full dynamical model, thereby enabling application of the RNA velocity framework to a broader set of biological systems and states. The dynamical model is implemented in their scVelo Python package, which also includes an efficient implementation of the steady-state model.

The input to both velocyto and scVelo effectively consists of two gene-by-cell count matrices; one representing mRNA ("spliced") abundances and one representing pre-mRNA ("unspliced") abundances. In practice, these two types of abundances are typically quantified using reads mapping to the exonic and intronic regions of genes, respectively. It is worth pointing out that this type of quantification differs from the more commonly used 'gene expression' estimation, for which a range of tools are available [4, 14–21]. There, a single matrix representing the abundance of mature mRNA is returned. Only a subset of the quantification tools provide explicit modules or instructions for joint estimation of both mature mRNA and pre-mRNA abundances, and those are the focus of the present study. Velocyto provides functions for parsing a BAM file of aligned reads (obtained by other tools such as Cell-Ranger [4]) and generating the two count matrices. Also dropEst [20] can be used to estimate exonic and intronic counts from an existing BAM file, e.g., from CellRanger. Dedicated end-to-end functionality for estimation of spliced and unspliced abundances from raw scRNA-seq reads is available within the kallisto|bustools software suite [14, 15] and in STARsolo, the single-cell mode of the STAR aligner [17]. Furthermore, assuming a properly specified set of reference sequences, the required counts can also be obtained using other general-purpose tools for quantification of droplet scRNA-seq data, such as alevin [16].

To our knowledge, no critical evaluation of the characteristics of the spliced and unspliced count matrices generated by these tools, and the effects on the downstream RNA velocity estimates, has been performed to date. In this study, we therefore used five public experimental droplet scRNA-seq data sets, generated with the 10x Genomics platform, to compare spliced and unspliced abundance estimates obtained by velocyto, kallisto | bustools, STARsolo, dropEst and alevin. Including alternative index definitions and parameter settings, we analyzed each of the five experimental data sets with a total of up to thirteen different quantification approaches (tool/parameter/index combinations). We illustrate that the estimates of spliced and unspliced abundances can be strongly affected by the choice of tool, as well as by the delineation of exonic and intronic regions; in particular, how genomic regions that are exonic in some annotated transcript isoforms and intronic in others are treated. Moreover, we show that abundance estimation is a critical step of the analysis workflow and that differences at this stage can have considerable effects on both the RNA velocity estimates and subsequent biological interpretation.

Methods

Data

In this study, we used five public single-cell RNA-seq data sets, generated by different laboratories using popular droplet-based protocols from 10x Genomics. Three of the five data sets (Pancreas, Dentate gyrus and Spermatogenesis) comprise cells from dynamically developing systems, while the remaining two (PFC and OldBrain) contain differentiated cells from adult mouse brain and were chosen as negative control data sets, assumed to not harbor a strong dynamic signal.

- The Pancreas data set [22] stems from a study of endocrine development in mouse, and was acquired with the 10x Genomics Chromium platform, using v2 chemistry. We downloaded the FASTQ files containing the reads from the cells at stage E15.5 from the Gene Expression Omnibus, accession number GSE132188. The RNA read length is 151 nt. A subset of this data set was used for illustration by [13], and is included as an example data set in the scVelo package. After the respective quantifications, we retained only the cells that are also included in the scVelo example data set, from which we also retrieved cell type labels. The final processed data set used for our analyses contains 3,696 cells.
- The Dentate gyrus data set [23] considers the developing mouse dentate gyrus, and was acquired with 10x Genomics v1 chemistry. The individual FASTQ files for cells from P12

and P35 were downloaded from the Gene Expression Omnibus, accession number GSE95315, and the reads from each time point were combined in a single pair of FASTQ files, with cell barcodes and UMI sequences in one file and the read sequence in the other. The RNA read length is 98 nt. Similarly to the Pancreas data set, the Dentate gyrus data set is also available as an example data set in scVelo, and was used for illustration by [13]. Only cells that were also studied in the scVelo paper were retained for our analysis, and cell type labels were obtained from the scVelo example data set. CellRanger, velocyto and dropEst were not run on this data set since the downloaded FASTQ files were not in the format expected by these tools. The final processed data set used for our analyses contains 2,914 cells.

- The Spermatogenesis data set [24] consists of steady-state spermatogenic cells from an adult mouse, and was processed with 10x Genomics Chromium v2 chemistry. The submitted BAM file was downloaded from the Gene Expression Omnibus, accession number GSE109033 (sample accession number GSM2928341), and converted to FASTQ format using the bamtofastq utility (v1.1.2) from 10x Genomics (https://support.10xgenomics.com/docs/bamtofastq). The RNA read length is 100 nt. Cell type labels were obtained from the corresponding loupe browser file provided by the authors, downloaded from https://data.mendeley.com/datasets/kxd5f8vpt4/1. Only cells that were also included in this file were retained for further analysis after quantification. The final processed data set used for our analyses contains 1,829 cells.
- The PFC data set [25] consists of cells from the prefrontal cortex of an adult mouse, and was generated using 10x Genomics Chromium v2 chemistry. Since only limited dynamics is expected in this data set, it is used here as a negative control. The submitted BAM file was downloaded from the Gene Expression Omníbus, accession number GSE124952 (sample accession GSM3559979), and converted to FASTQ format using bamtofastq. The RNA read length is 98 nt. Only cells annotated to 'PFC sample 2' were used, and cell type labels assigned by the data generators were obtained from the GEO record metadata. The final processed data set used for our analyses contains 1,267 cells.
- The OldBrain data set [26] consists of cells from the dissociated brain of a 21-22 months old mouse, and was generated using 10x Genomics Chromium v2 chemistry. Similarly to the PFC data set, it is used as a negative control. The submitted BAM file was downloaded from the Gene Expression Omnibus, accession number GSE129788 (sample accession GSM3722108), and converted to FASTQ format using bamtofastq. The RNA read length is 57 nt. Only cells annotated to 'sample 37' were used, and cell type labels assigned by the data generators were obtained from the GEO record metadata. The final processed data set used for our analyses contains 1,823 cells.

Feature sequence extraction

All analyses are based on reference files from Gencode, mouse release M21 [27]. The desired output from each quantification method is a pair of count matrices; one containing 'spliced' or 'exonic' counts, and the other containing 'unspliced' or 'intronic' counts for each gene in each cell. For simplicity, in the remainder of this paper, the terms 'spliced' and 'exonic' will be used interchangeably to refer to the counts representing the processed mRNA abundances, and 'unspliced' and 'intronic' counts will similarly refer to the counts representing the unprocessed pre-mRNA abundances.

To enable this type of quantification with alevin and kallisto bustools (as opposed to the more standard, gene-level expression quantification), the Gencode reference files were

Table 1. Summary of reference FASTA files containing exonic and intronic sequences. Only files generated by eisa R were used for creation of quantification indices.

FASTA file name	Sequence extraction	Description
BUSpaRse-separate	BUSpaRte	Introns extracted separately from each transcript isoform
BUSpaRse-collapse	BUSpaRte	Introns extracted after collapsing all transcripts of a gene
eisaR-separate	eisaR	Introns extracted separately from each transcript isoform
eisaR-collapse	eisaR	Introns extracted after collapsing all transcripts of a gene

https://doi.org/10.1371/journal.pcbi.1008585.t001

processed as follows (also summarized in Table 1). First, we used the BUSpaRse R/Bioconductor package (v1.0.0) [28] to extract transcript and intron sequences from the genome sequence and the Gencode annotation GTF file. While the definition of the transcripts is unambiguous, the BUSpaRse package supports two ways of defining the introns (see \$1 Fig for a schematic). The 'separate' approach considers each transcript separately when extracting the intronic regions (and thus, an intron can overlap with an exonic region of an alternate transcript), while the 'collapse' approach first collapses the isoforms of a gene (taking the union of all the exonic regions) before defining the introns as any regions of the gene locus that are not exonic in any of the annotated transcript isoforms of the gene. In effect, the 'separate' approach thus considers exonic and intronic regions on an equal footing, while the 'collapse' approach represents a prior belief that an ambiguous read is more likely to come from an exon than from an intron. A flanking sequence of length L-1 (where L is the RNA read length of the respective study) is added on each side of each intron to account for reads mapping across exon/intron boundaries. For comparison, we also reimplemented the extraction of transcript and intron sequences (for both the 'separate' and 'collapse' approaches) directly using functions from the GenomicFeatures and BSgenome R/Bioconductor packages [29, 30]. For convenience, code used to extract the features has been included in the eisaR R/Bioconductor package (https:// www.bio.conductor.org/packages/eisaR/). In each case, the extracted transcript and intron sequences were written to a joint FASTA file, summarized in Table 1. Upon comparison of the two implementations, we noticed that the BUSpaRse release version at the time (v1.0.0) returned erroneous feature sequences for multi-exonic transcripts and corresponding introns on the negative strand. This was partly corrected in a later release (v1.2.2) (S2 Fig). However, for this reason, the features extracted by BUSpaRse were not used for further analyses.

Reference index generation

The combined transcript and intron FASTA files were used to build the following quantification indices (summarized in Table 2):

- a joint transcript and intron index for Salmon (v1.0.0) [31]
- an index for Salmon, considering the transcripts as the features of interest and providing the introns as decoy sequences [32]
- an index for Salmon, considering the introns as the features of interest and providing the transcripts as decoy sequences
- a joint transcript and intron index for kallisto (v0.46.0) [33]

In addition to the indices based on transcripts and introns, we built one Salmon index from the original Gencode FASTA file with the annotated transcripts, and one Salmon index from a FASTA file combining the annotated transcripts and fully unspliced versions of all transcripts.

Index type Index building Input reference files (Table 1) Target sequences Decoy sequences spliced transcripts and introns salm on-joi nt Saknon e isaR-separate, eisaR-collapse genome salmon-spliced Sahnon transcriptome spliced transcripts ge nom e salmon-spliced-unspliced Salmon transcriptome + unspliced transcripts spliced and unspliced transcripts genome introns + genome salmon-spliced-decoy Salmon e isaR-separate, eisaR-collapse spliced transcripts salmon-introns-decoy Salmon e isaR-separate, eisaR-collapse introns spliced transcripts + genome kalli sto-ioin t kallicts e isaR-separate, eisaR-collapse spliced transcripts and introns N/A cellranger CellRanger genome and GTF file whole genome N/A star STAR genome and GTF file whole genome N/A

Table 2. Summary of reference indices. The full index name is constructed by concatenating the index type to the FASTA file name, e.g. salmon-joint-eisaR-separate refers to the salmon-joint index built for the sequences in the eisaR-separate FASTA file.

https://doi.org/10.1371/journal.pcbi.1008585.t002

For all Salman indices, the complete genome sequence was provided as a decoy sequence [32], with the aim to exclude reads coming from intergenic regions of the genome. Across data sets and reference specifications, this excluded between less than 1 and 2.5% of the reads from the quantification. For the quantification based on spliced-only transcripts, in which the genome decoy would also capture unambiguously intronic reads, this number was between 2 and 10%. Finally, we built an index for CellRanger (v3.0.2) [4] and one for STAR (v2.7.3a) [17] based on the reference genome and GTF file from Gencode. The splice junction database overhang in the STAR index was set to 150nt, which is at least as long as the read length minus one for all data sets considered here.

To investigate the effect of the choice of flank length in the intron definition, we further built Salmon indices (using the 'separate' intron definition) with flank lengths equal to the read length minus 21bp, and the read length minus 41bp. We also built an index for kallisto bustools using the kb-python wrapper [15, 33], which uses the 'separate' intron definition and sets the flank length to 30 bp.

Sequence uniqueness estimation

Finally, to aid in the interpretation, we estimated the sequence uniqueness for each gene, relative to all other genes, separately for each FASTA file generated as described above (Table 1). The gene uniqueness was defined as the fraction of unique k-mers in the gene, that is, the fraction of the constituent k-mers that are not found in any other gene. For each data set, the k-mer length was set to be equal to the RNA read length minus 1. The sequence uniqueness for a gene was calculated in two different ways. In each case, the full FASTA file with transcript and intron sequences was used as input. First, we estimated separate uniqueness values for the exonic and intronic parts of a gene, by assigning the transcripts and introns of a gene to distinct gene IDs in the uniqueness calculations. With this approach, k-mers that are exonic in one transcript and intronic in another (even if these are isoforms of the same gene) are classified as non-unique, and thus the uniqueness score of an exon sharing some of its sequence with an intron of another transcript from the same gene is reduced. Second, we estimated an overall gene uniqueness, by again considering all (exonic and intronic) sequences, but not penalizing shared sequence between introns and exons of the same gene.

Quantification

alevin. For each of the Salmon indices described above, we ran alevin (v1.0.0) [16] to estimate exonic and intronic abundances for each annotated gene. It is worth noting that for abundances obtained with the salmon-spliced-unspliced index, 'exonic' and 'intronic'

abundances refer directly to spliced mRNA and unspliced pre-mRNA abundance estimates, respectively, while for all other indices, 'exonic' abundances refer to exonic regions and 'intronic' abundances to intronic regions. This affects, for example, how reads aligning completely within an exon are used in the quantification: For the salmon-spliced-unspliced index, such a read could stem from either the spliced or the unspliced transcript (since also the latter contains the exons) and could contribute to the abundance of either (or both) of them, while with the other indices, it will be considered 'exonic'.

For the alevin quantifications, the transcripts and introns (or unspliced transcripts) from the same gene were manually annotated with different gene IDs, in order to obtain separate exonic and intronic gene-level abundances despite estimating them jointly. For the indices with decoys, the exonic gene counts were defined as the counts obtained when quantifying against the transcript index (with introns as decoys) and the intronic gene counts were similarly obtained by quantifying against the intron index (with transcripts as decoys). Hence, for these approaches, it is possible for a read that maps equally well to an exonic and an intronic sequence to be included in both the exonic and the intronic count matrices, and thus be counted twice.

kallisto|bustools. For each of the kallisto indices, we applied kallisto|bustools (v0.46.0)
[14, 15] to generate a BUS file. Barcodes were corrected using the list of available cell barcodes from 10x Genomics for the appropriate chemistry version, and the BUS file was sorted using kallisto|bustools. Next, the BUS file was subset with the capture command of kallisto|bustools to generate separate BUS files to be used for the quantification of exonic and intronic features, respectively. The gene-level exonic and intronic counts were subsequently obtained using the count command. The capture was performed using two different approaches:

- "include", where the features of interest for the quantification at hand are provided to kallistol
 bustools. In other words, the transcript IDs are provided as the -c argument to quantify the
 exonic abundances, and the intron IDs are provided to quantify the intronic abundances. In
 practice, this means that reads in equivalence classes containing at least one transcript are
 retained for the exonic quantification, and reads in equivalence classes containing at least
 one intron are retained for the intronic quantification. Hence, equivalence classes containing
 both exonic and intronic features will be provided to both the exonic and intronic quantification steps.
- 'exclude', where the features that are not of interest for the quantification at hand are provided to kallisto|bustools, and subsequently excluded. In other words, intron IDs are provided as the -c argument to quantify the exonic abundances, and the transcript IDs are provided to quantify the intronic abundances, and in addition the -x flag is used to indicate that the provided IDs represent sequences to be excluded. In practice, this means that only reads in equivalence classes that don't contain any introns will be retained for the exonic quantification, and only reads in equivalence classes that don't contain any transcripts will be retained for the intronic quantification. Hence, equivalence classes containing both exonic and intronic features will be excluded in both steps.

In addition to the manual application of kallisto|bustools as described above, we applied the kb-python wrapper for quantification based on the corresponding index. With its default settings, it calls kallisto|bustools with the 'exclude' capture approach.

Velocyto, drop Est and STARsolo. CellRanger (v3.0.2) [4] and velocyto (v0.17) [9] were run with default settings to generate exonic and intronic counts based on the CellRanger index. Also drop Est (v0.8.6) [20] was run on the BAM files output by CellRanger, specifying the -V flag to instruct drop Est to return exonic, intronic and exon-intron spanning UMI

Table 3. Summary of quantification approaches. In addition to the strategies included in this table, we applied a levin_sep_gtr after using different flank length when defining the intronic reference sequences, alevin_sep_gtr in unstranded mode, and the kb-python wrap per around kall isto]bustools. The effects of these modifications are evaluated in S3 Fig.

Method ID Quant method		In dex Exonic/	Exonic/intronic quantification	
alevin_coll_gtr	alevin	salmon-jo int-eisaR-collapse joint		N/A
alevin_coll_decoy_gtr	alevin	exonic: salmon-spliced-decoy-eisaR-collapse, intronic: salmon- introns-decoy-eisaR-collapse		N/A
alevin_sep_gtr	alevin	salmon-jo int-e isaR-separate joint		N/A
ale vin_sep_dec oy_gtr	alevin	exonic: salmon-spli ced-decoy-eisaR-separate, intronic: salmon- introns-decoy-eisaR-separate		N/A
alevin_spliced_unspliced_gtr	alevin	salmon-spliced-unspliced joint		N/A
kallisto bus_coll_eucl	kallisto bustools	kallisto-joint-eisaR-collapse separate		exclude
kallisto bus_coll_incl	kallisto bustools	kallisto-joint-eisaR-collapse separate		include
kallisto bus_sep_excl	kallisto bustools	kallisto-joint-eisaR-separate separate		eschide
kaili sto bus_sep_incl	kallisto bustools	kallisto-joint-eisaR-separate separate		include
starsolo	STAR	star joint		N/A
starsolo_diff	STAR	star exonic: G —Gene	iene, intronic: GeneFull	N/A
velocyto	velocyto	CellRanger joint		N/A
dropest	drapEst	CellRanger joint		N/A

https://doi.org/10.1371/journal.pcbi.1008585.t003

count matrices. The intronic and exon-intron spanning UMI counts were summed and used to represent the unspliced UMI count. STARsolo (v2.7.3a) [17] was run using the STAR index, specifying the SOLOfeatures argument to generate 'Velocity' (exonic and intronic), 'Gene' (regular exonic gene expression) and 'GeneFull' (reads with any overlap with the gene locus) counts. Based on these count matrices, we obtained exonic and intronic count matrices in two different ways. First, we directly used the 'Velocity' count matrices as exonic and intronic counts (below referred to as starsolo). Second, we used the 'Gene' count matrix as the exonic counts, and the difference between the 'GeneFull' and 'Gene' counts as the intronic counts (below referred to as starsolo_diff). For genes where the 'Gene' counts were higher than the 'GeneFull' counts, the intronic count was set to zero. This can happen, for example, for a gene located in the intron of another gene. In the 'GeneFull' quantification, reads mapping to such a gene are considered ambiguous and therefore discarded. However, they may be assigned in the 'Gene' quantification, if they are compatible with the annotated gene model. An overview of the evaluated quantification approaches is provided in Table 3.

Cell filtering and data processing. For each quantification, we generated a Single Cell Experiment object [34] containing the exonic and intronic counts. Only cells and genes included
by all methods were retained for further analysis. For the Pancreas and Dentate gyrus data sets,
we further subset the objects to only the cells analyzed by [13], for the PFC data we kept only
cells annotated to PFC sample 2, for the OldBrain data only cells annotated to sample 37, and
for the Spermatogenesis data set only cells with an assigned cell type label provided by the data
generators were retained. For visualization purposes, we calculated a single low-dimensional
representation based on the alevin quantification of only the spliced mRNAs (using the original transcriptome FASTA file from Gencode). After normalization with scater (v1.14.6) [35],
using the library sizes to define size factors, we extracted 30 principal components from the

log-transformed normalized count values. The scater package was then used to apply UMAP [36] with default parameters to the PCA output to obtain a two-dimensional representation that was used for visualization of estimated velocities.

Visualization

In order to visualize the coverage pattern of reads within genomic regions, we subset the BAM file generated by CellRanger to only the reads assigned to the retained cell barcodes, using the subset_bam tool (v1.0) from 10x Genomics (https://github.com/10XGenomics/subset-bam). Next, we used BEDTools (v2.27.1) [37] to calculate the coverage along the genome, separately for reads on the positive and negative strand. The bedGnaphToBigWig script from Kent Tools (v20190212) [38] was used to convert the resulting bedGraph file to bigwig format.

Coverage patterns, together with annotated gene models, were visualized using the Gviz R/Bioconductor package (v1.30.3) [39]. In these figures, the annotated gene models are visualized by their genomic coordinates, together with coverage tracks of reads aligned to the positive and negative strand of the genome. The alignments are aggregated across all the retained cells in the data set. All alignments contained in the BAM file are included; hence, multimapping reads are represented in all the reported mapping locations. Moreover, no UMI deduplication is performed and thus the number of reads reported in the coverage tracks are often higher than the total UMI count returned by any of the counting methods. It is also important to note that while the gene models and coverage tracks are represented with respect to a genomic reference for ease of interpretation, both alevin and kallisto |bustools perform the quantification based on mapping to transcriptomic features, not alignment to the genome. Thus, these plots are not intended to provide an exact correspondence between mapped reads and estimated UMI counts, but rather serve as illustrations to aid in the understanding of the causes of differences between the counts from the various methods.

RNA velocity estimation

SingleCellExperiment objects with exonic and intronic gene-level UMI counts were converted to AnnData objects [40] using the anndata2ri package v1.0 (https://github.com/theislab/anndata2ri). The scVelo package (v0.1.24) [13] was then used to normalize the counts and select the 2,000 most highly variable genes separately for each quantification approach, after excluding all genes with less than 20 assigned reads across the exonic and intronic components (only summing across cells with nonzero exonic and intronic count). Note that by default, scVelo selects highly variable genes based on the spliced counts only. RNA velocity estimates were obtained using the dynamical model implemented in scVelo. For comparison, we also performed downstream analysis and visualization of the RNA velocity using only the genes that were selected (and for which scVelo returned a finite velocity value) by scVelo with all the quantification approaches.

The scVelo analysis returns a gene-by-cell matrix of estimated velocities, as well as corresponding matrices of normalized (spliced and unspliced) abundances. Based on these matrices, we estimated both gene- and cell-wise Spearman correlations between the different types of abundances, as well as between the abundances and the velocity estimates. For each gene, we also fit a linear model without intercept with the normalized spliced abundance as the predictor and the normalized unspliced abundance as the response, and estimated the R^2 for the model, similarly to the fit_r^2 value currently returned by scVelo. A large R^2 value indicates that most cells lie close to a diagonal line in the spliced-unspliced abundance space, suggesting that the gene is in a steady state across the majority of the cells. We used this to contrast the

data sets where we do not expect much systematic dynamics to those where we do, expecting a shift towards larger \mathbb{R}^2 values in the former. For each cell, scVelo further provides an estimate of the velocity confidence (representing the average correlation of the velocity vector of the cell and those of its neighbors). It is worth noting that the velocity calculations are performed separately for each input gene, and the resulting values are therefore independent of which other genes are included in the data set (under the assumption that the normalized abundance values stay unchanged).

Low-dimensional embedding of velocities

All velocity estimates were embedded into the same low-dimensional representation, calculated from the spliced-only abundance quantification by alevin as described above. Based on the estimated velocity vectors and the differences between the expression profiles of different cells, scVelo calculates a cosine correlation (π_{ij}) for each pair of 'neighboring' cells. A high value of π_{ij} indicates that the velocity vector of cell i points in the direction from cell i to cell jin gene expression space. The low-dimensional embedded velocity vector for cell i, calculated by scVelo, is given by

$$\tilde{v}_i = \sum_{j \neq i} \left(\tilde{\pi}_{ij} - \frac{1}{n} \right) \tilde{\delta}_{ij}$$

where $\hat{\pi}_{ij}$ is the transition probability from cell i to cell j (derived from the cosine correlation π_{ii}), n is the number of cells, and

$$\tilde{\delta}_{ij} = \frac{\tilde{s}_j - \tilde{s}_i}{\|\tilde{s}_i - \tilde{s}_i\|}$$

is the normalized difference of the coordinates of cells i and j in the low-dimensional embedding. The fact that $\tilde{\delta}_{ij}$ is normalized implies that the length of \tilde{v}_i indicates to what extent the cells to which cell i has high transition probabilities are all located in the same direction from cell i in the low-dimensional representation. It further implies that the embedded velocities are potentially more comparable across methods than the original velocity vectors, since the magnitudes of the latter depend on the normalized abundance levels of the genes, and since the velocity vectors will only be directly comparable between methods if they are based on the same set of input genes. In order to compare the velocity embeddings across methods, we calculate a concordance score for each cell. The score for cell i is defined as the ratio between the length of the sum of the embedded velocity vectors for cell i across all quantification methods, and the sum of the lengths of the individual embedded velocity vectors. In other words, the score for cell i is given by

$$c_{i} = \frac{\|\sum_{m} \bar{v}_{i}^{(m)}\|}{\sum_{m} \|\bar{v}_{i}^{(m)}\|}$$

where the sum is taken over all methods m, and $\tilde{v}_i^{(n)}$ is the embedded velocity vector for cell iwith method m. If all embedded velocity vectors for cell i point in the same direction in the low-dimensional representation, this ratio will be close to 1, while if there is less concordance between the different quantification methods, the ratio will be lower than 1.

Results

The total UMI count varies between methods

With the aim of characterizing global differences between the counting methods, we first directly compared the total UMI count assigned to exonic and intronic targets by each of the methods. We summed up the counts across all cells, either across all genes or within gene subsets stratified by sequence uniqueness (Fig. 1, Pancreas data; S4 Fig., OldBrain data). There are considerable differences in the assigned UMI counts between the methods. Moreover, these differences are not confined to a small number of susceptible genes, but can be seen across a large fraction of the expressed genes (S5 Fig.).

Overall, dropest gave the highest total UMI counts, mostly resulting from higher counts for genes with a low degree of sequence uniqueness. This effect is stronger in the data sets with a higher degree of multi-mapping, such as the OldBrain data set (\$4 Fig). We hypothesized that this effect is due to insufficient UMI collapsing, resulting from the use of only the alignments designated as primary alignments in the BAM files output by CellRanger. This reported primary alignment can vary between different reads stemming from the same underlying transcript molecule. In such cases, the corresponding UMI could contribute to the count for multiple different genes. Among the other approaches, starsalo_diff and the alevin-based quantification approaches based on transcript/intron annotations gave the highest total UMI counts, mainly driven by higher counts for the exonic targets. Also here, the effect is predominantly driven by genes with a low fraction of unique k-mers. alevin includes multi-mapping reads in the quantification, thus increasing the counts for these targets. This is in contrast to kallisto | bustools, velocyto and starsalo, which by default exclude ambiguous reads that map to multiple genes from the quantification.

Velocyto and starsolo_diff gave the highest UMI counts for genes whose transcripts are all shorter than the read length (genes in the 'NA' category in Fig.1, for which no uniqueness could be calculated since all transcripts were shorter than the employed k-mer length). However, velocyto considered most of these reads to be exonic, while starsolo_diff assigned them to the intronic features. This behaviour of starsolo_diff is likely due to the generation of the intronic counts by subtraction of the exonic count from the full gene locus count. A read which partly overlaps the gene locus but is not consistent with the annotated gene model would be included in the 'GeneFull' count but not in the 'Gene' count, and thus considered an intronic read, regardless of whether or not the gene actually contains any introns. Similarly, starsolo_diff assigned higher counts than both velocyto and the default 'Velocity' counting of starsolo, for both exonic and intronic features.

As expected, the alevin-based approaches where exonic and intronic features are quantified separately, as well as the 'include' capture mode of the kallisto|bustools approaches, tend to give higher total UMI count than quantifying exonic and intronic features jointly with alevin or running kallisto|bustools in 'exclude' capture mode, especially for the 'separate' intron definition. This is likely due to double-counting of some reads that map equally well to an exon and an intron. The difference between the 'include' and 'exclude' capture approaches is smaller for the 'collapse' annotation, since in that case, no genomic regions are annotated as both intronic and exonic for the same gene. The same is true for the difference between the alevin quantifications employing joint and separate quantification. It is worth noting that the length of the flanking region chosen when constructing the intronic features for the quantification also influences the counts (S3 Fig.). A shorter flank length typically leads to a lower unspliced count, since a larger fraction of the read must overlap the intron for the read to be considered potentially intronic.

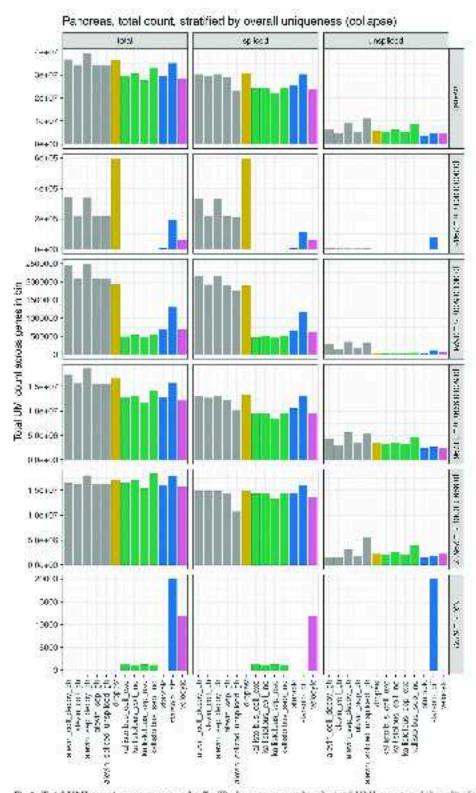


Fig.1. Total UMI count across genes and cells. The bars correspond to the total UMI count, and the split of these into counts for exonic and intronic targets, for each quantification method in the Pancreas data set. Largely similar patterns were seen in the other four data sets. In addition to the overall count (top row), the figure shows the total count after stratifying genes by the oresult fraction of unique k-mers (using the 'collapse' annotation), indicated in the vertical panel headers together with the number of genes in the category. The genes for which no uniqueness information

could be calculated (the "NA" category) are those for which all transcripts are shorter than the chosen k-mer length (which was set to the read length minus one; here 190nt).

https://doi.org/10.1371/journal.pcbi.1008585.g001

In addition to the absolute counts, also the fraction of UMIs assigned to unspliced targets varies between methods, with the largest fraction of intronic counts obtained by alevin_spliced_unspliced_gtr. This was expected, given that for this method, the 'intronic' features are the full pre-mRNA molecules and thus contain both exonic and intronic sequences. Hence, also reads falling in exons may be assigned to the unspliced features.

Individual genes exemplify methodological differences

Next, we aimed to find individual genes whose count patterns exemplify the main methodological differences among the counting strategies. First, we restricted the set of genes to those that were selected as highly variable by scVelo, and thus used for velocity estimation, for at least one counting method. Notably, while a large fraction of these genes were selected across all quantification methods, there were non-negligible differences between the sets of selected genes (S6 Fig). In particular, alevin_spliced_unspliced_gtr gave the largest number of unique highly variable genes, followed by the starsalo variants and alevin_sep_gtr depending on the data set. For each of the retained genes, we calculated the fraction of the total counts that were assigned to unspliced features (summarized across all cells). Next, we calculated the standard deviation, across the quantification methods, of these intronic count fractions and selected the top 10% of the genes based on this measure. These genes were partitioned into 10 clusters based on the Pearson correlation dissimilarity $(\sqrt{2(1-\rho)})$ where ρ is the Pearson correlation) between the fractions of intron-assigned counts across methods, using hierarchical clustering with complete linkage (Fig 2). The gene clusters reveal typical cases where the methods yield different exonic and/or intronic counts. Representative genes for each cluster, selected among the genes with the highest correlation with the cluster centroid, are discussed below and illustrated in \$7 and \$8 Figs.

Genes with ambiguous regions. (clusters 1, 2, 3, exemplified by Tspan3, Ssr1 and Alg2 in S7 Fig). For genes in which many of the base positions are annotated to both exons and introns (in different isoforms), the choice of how to define introns ('separate' vs 'collapse' approaches) has a major effect on the quantifications. If exons are collapsed before the introns are defined, reads falling in ambiguous regions are considered exonic, leading to a higher exonic and a lower intronic count than with the 'separate' intron definition. This effect can be seen in approximately half of the genes with the highest variability in the fraction of unspliced counts (Fig.2), and manifests itself via a low fraction of unspliced reads for the alevin and kallisto|bustoals approaches based on annotations explicitly obtained with the 'collapse' intron definition approach, as well as dropest. Also starsalo_diff falls in this category, since the intronic regions are not considered when the exonic counts are estimated (via the 'Gene' count), and thus any read that is compatible with at least one transcript model is considered exonic.

With the 'separate' intron definition, running alevin with decoys or kallisto|bustools with the 'include' capture double-counts many reads falling completely in ambiguous regions, giving high values of both exonic and intronic counts. Conversely, running kallisto|bustools with 'exclude' capture discards many reads in ambiguous regions, since they will typically be assigned to equivalence classes containing both exonic and intronic targets, and this counting strategy therefore often returns low counts for both types of features. While these effects can be seen in the absolute counts, they do not necessarily affect the ratio of spliced and unspliced counts (S7 Fig).

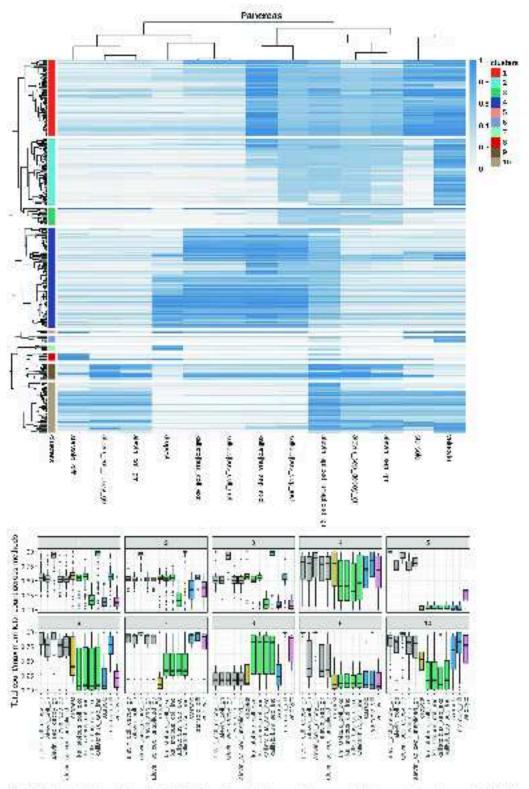


Fig 2. Clustering of genes based on the fraction of unspliced counts. Top panel: Hearmap showing the unspliced UMI count fraction, across all cells, in the Pancreas data. The dendrogram was cut into 10 clusters, indicated by colors to the left of the hearmap. Bottom panel: Relative total count assigned to each gene by the different methods, stratified by gene cluster.

https://doi.org/10.13/1/journal.pdbi.1008585.g002

For genes with many ambiguous regions, velocyto often returns a relatively low number of spliced counts, and consequently a large fraction of unspliced counts (clusters 1, 2). This likely follows from the default 'permissive' counting logic where, essentially, a read contributes to the spliced count only if it is consistent with the exonic region of a transcript model, but does not map to an intronic region or an exon/intron boundary of any other transcript model. Also starsolo often assigns a low total count for this group of genes.

Genes overlapping (introns of) other genes. (clusters 4, 5, 6, 9, 10, exemplified by Rassf1, Gm21983, Chkb, Cnot6 and Tmem120b in S8 Fig). Overlaps between genes can take many different shapes; exons of one gene can overlap either introns or exons of other genes, and the overlap can be on the same or on opposite strands. Reads falling in an exon of one gene and in an intron of another gene on the same strand are considered ambiguous, and are therefore discarded, by kallisto|bustools, velocyto and starsolo. In contrast, dropest and alevin do not discard such reads. dropest preferentially assigns the reads to the exonic region, while alevin instead often distributes the reads between the two features. Thus, alevin tends to give a higher fraction of unspliced counts, and in many cases also a higher total count, for genes with other genes in their introns (cluster 9, exemplified by Cnot6 in S8 Fig). It is worth observing that alevin will often also assign reads to the gene located within the intron (Gm12191 in S8 Fig, which is the gene in the intron of Cnot6). As for the previous category of genes, alevin with the decoy approach double-counts reads mapping equally well to an exon and an intron, regardless of whether or not the exon and the intron belong to the same gene.

In cases of exonic overlap between genes on the same strand (clusters 5 and 6, exemplified by Gm21983 and Chkb in S8 Fig), all methods except alevin consider the corresponding reads ambiguous and discard them, leading to a large difference in the total counts (Fig.2). The main difference between many genes in clusters 5 and 6 is that in cluster 6, the gene of interest overlaps partly with an intron of the other gene, and the reads in the overlapping region may be counted by the 'Gene' approach of starsolo_diff and by dropest. In cluster 5, the gene of interest overlaps only with exons of the other gene in the region.

Reads falling in overlapping regions of genes on apposite strands are not considered ambiguous if the strandedness of the reads is taken into account. However, not accounting for the strandedness (illustrated here by the behaviour of kallisto bustools and dropest) implies treating such overlaps similarly to same-strand overlaps. For example, for a gene located inside the intron of another gene on the opposite strand, all reads not mapping across an exon-exon junction will map equally well to the two genes and thus be discarded (cluster 10, exemplified by Tmem120b in \$8 Fig.). As before, dropest tends to give preference to exonic overlaps, while the reads falling in the intronic regions in the two overlapping genes are often discarded. Exonic overlaps between genes on opposite strands also leads to a decrease in the number of assigned reads if the strandedness is not taken into account (exemplified by Rassf1 in S8 Fig). The assigned count can also increase by performing the quantification in an unstranded manner. For example, reads mapping to the negative strand, not overlapping any feature there but overlapping a gene on the positive strand, can be assigned to the latter. Intronic reads resulting from discordant priming from poly-T sequences, as observed by [9], would also be incorporated with an strand-agnostic counting approach. The observation of [9] could be reproduced in our data sets (\$9 Fig), and the incorporation of opposite-strand reads by kallisto bustools and dropest is exemplified by Ssr1 and Brsk2 in \$7 Fig.

Genes with reads only partly overlapping the gene body. (cluster 8, exemplified by 1810019D21Rik in S7 Fig). In some cases, reads extend outside the annotated gene body. The starsolo_diff "GeneFull" count incorporates these reads, while the "Gene" count from the same method does not, since they are not compatible with the annotated gene model. As an effect, the difference between them (which is used as the intronic count) can be high, even in cases

where the gene does not contain introns, and the total count assigned by starsolo_diff is higher than most of the other methods. Velocyto and, for some genes in the cluster, the kallisto |bustools methods, yield a high spliced count, while the alevin methods, dropest and starsolo return lower values.

Reads falling in purely exonic regions. Across the clusters, alevin_spliced_unspliced_gtr
often shows a different read assignment compared to the other methods. As previously noted,
there is a fundamental difference between the reference used for alevin_spliced_unspliced_gtr
(which considers full spliced and unspliced transcripts) and those used for the other methods
(which consider transcripts and introns). This implies that reads falling completely in exonic
regions, from the point of view of alevin_spliced_unspliced_gtr, are equally likely to come from
a spliced as an unspliced molecule, while the other methods consider such reads unambiguously exonic. In practice, this can lead to a lower exonic count and a more even exonic/intronic count ratio for alevin_spliced_unspliced_gtr than for the other methods. This is exemplified
by Map1b (cluster 10) in S7 Fig.

Large differences in inferred velocities between quantification methods

In the previous sections we showed that there are noticeable differences between the quantification methods, in terms of the total number of UMI counts as well as the distribution of these between spliced and unspliced targets. Next, we asked whether these differences could be seen also in the velocity estimates from scVelo, and in the embedding of these in a low-dimensional representation of the cells, which is arguably the most widely used way of interpreting RNA velocity estimates. For our analyses, we provided scVelo with raw spliced and unspliced UMI count matrices. These were then filtered and normalized by scVelo, and the RNA velocity was estimated for each input gene and each cell. Velocities were estimated for either the individual sets of 2,000 highly variable genes from each quantification method, or the set of genes that were selected by scVelo (and obtained a valid velocity value) with all the quantifications.

Interestingly, the estimated velocities consistently showed a lower correlation between methods than the normalized (spliced, unspliced or aggregated) abundances, when calculated across either cells or genes (\$10 Fig). Within a cell, there was also a relatively strong correlation between the total gene abundance and the absolute value of the velocity (\$11 Fig). This should be factored in when comparing absolute velocities across genes, and it may also suggest that velocity estimates are not directly comparable across quantification methods if the number of assigned reads are very different. For a given gene, the fraction of unspliced counts was also moderately positively correlated with the estimated velocity. The spliced and unspliced abundances were positively correlated for all quantification methods, suggesting that the intronic signal is indeed real and of potential biological relevance, rather than just the result of, e.g., contamination by genomic DNA. Finally, we noticed a moderate positive correlation between the abundance of a gene and the likelihood of the velocity fit in four of the five data sets (exemplified by the Pancreas data set in \$12 Fig), while it was substantially lower in the Dentate gyrus data set (\$13 Fig).

The velocity estimates were visualized by embedding them into a UMAP representation based on the alevin_spliced_gtr quantification (note that the alevin_spliced_gtr counts are not used to estimate the velocity, since no intronic counts are estimated). The UMAP embedding was compared to other types of embeddings (PCA, tSNE, UMAP based on aggregated abundances, unspliced abundances only or spliced and unspliced abundances concatenated), in terms of the length of the embedded velocity vectors as well as the average distance between the velocity vector of each cell and its 10 nearest neighbors. These comparisons suggested that the differences between embeddings were relatively minor, but that UMAP often provided a

slightly more interpretable representation (S14 and S15 Figs). Embeddings based solely on unspliced abundances were the least interpretable from a velocity perspective.

Differences in velocity estimates directly affect biological interpretation

From the UMAP visualizations it is immediately apparent that the differences in the estimated abundances between the quantification methods directly influence interpretation, e.g., indicated by streamlines pointing in different directions in certain regions of the plots (S16–S18 Figs). These differences are not captured by the velocity confidence estimates returned by scVelo for an individual method, which are often high for all the cells (S19–S21 Figs), suggesting that the differences between methods are systematic rather than just the result of random fluctuations or uncertainty in the velocity estimation. The similarity among the low-dimensional velocity embeddings based on different quantification methods increased somewhat when they were derived from the set of shared genes (Fig 3). However, considerable differences were still seen, indicating that the quantification does not only influence the velocity interpretation via the selection of genes.

The lack of unambiguous ground truth complicates a direct evaluation of the accuracy of velocity estimates from the different quantification methods. In addition, the typical way of interpreting velocity estimates by means of embedded stream lines in a low-dimensional space provides a relatively coarse-grained measure. Nevertheless, Fig 3 and S16 Fig suggest that for the Pancreas data, the largest differences between methods appear in the differentiated Alpha, Beta, Delta and Epsilon cell types (left). Here, alevin_sep_decoy_gtr, dropest and the kallisto| bus methods induce (partly or fully) a 'back-flow', with streamline arrows pointing from the differentiated cells back towards the pre-endocrine cells. A similar observation can be made for alevin_spliced_unspliced_gtr and kallisto|bus_sep_ind for the pre-endocrine cells. The cycling nature of the ductal cells is visible in the embeddings of velocities from most quantification methods, with the exception of alevin_spliced_unspliced_gtr, alevin_sep_decoy_gtr, kallisto|bus_coll_incl_and kallisto|bus_sep_ind.

Also for the Spermatogenesis data (Fig. 3 and S17 Fig.), the largest differences between methods are seen towards the end of the developmental trajectory. Again, many methods (with the exception of alevin_sep_gtr and kallisto|bus_sep_excl) induce a back-flow, with streamline arrows pointing from the late round spermatids towards the mid round spermatid cluster. In most cases, this back-flow continues through (part of) the mid round spermatid cluster as well.

In the Dentate gyrus data (\$18 Fig), the lowest concordance between velocities based on different quantifications is seen for the cells in the granule cell lineage (middle part). While some quantifications indicate a direction largely from neuroblasts to granule cells (e.g., alevin_ sep_gtr, kallisto|bus_coll_ind), others indicate a strong movement in the opposite direction (e.g., alevin_sep_decoy_gtr, kallisto|bus_sep_incl). All methods except kallisto|bus_coll_excl, kallisto|bus_sep_excl and kallisto|bus_coll_incl show a strong dynamic flow within the mature granule cells, and there is further disagreement within the astrocyte cell cluster. Overall, the results from these three data sets highlight that the biological interpretation can be strongly affected by the choice of quantification method.

Negative control data

The PFC and OldBrain data sets were used to compare the methods in terms of their performance on 'negative control' data sets, that is, data sets where no strong systematic dynamics are expected. While it is intrinsically hard to find a living biological system without dynamics, we expected to find a larger fraction of genes in a steady state, in which the expression levels are maintained at a constant level over time. Under the assumption of a constant degradation

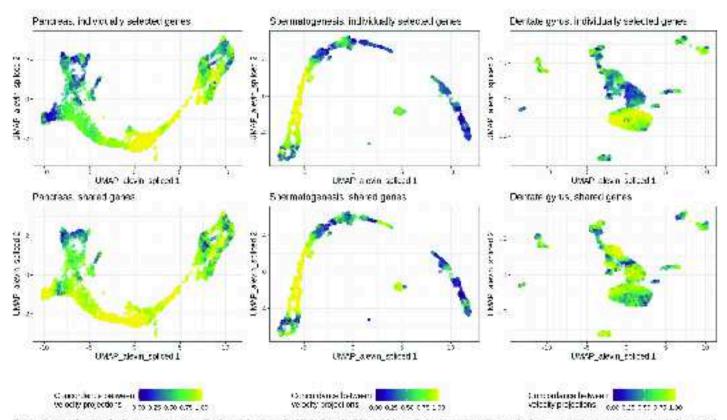


Fig.3. Concordance of velocity projections. Each panel colors the UMAP embedding of the cells by the concordance of velocity projections across methods for one of the three data sets expected to contain systematic dynamics. The concordance for a given cell was defined as the ratio between the length of the sum of the velocity embeddings for the different quantification methods, and the sum of the lengths of the same embedding vectors. A value close to 1 indicates that the embedded velocity vectors all point in the same direction, while a low value indicates larger deviations among them. For each data set, the top panel shows the concordance between the velocity embeddings derived based on all the 2,000 genes selected by at Velo for each method, while the bottom panel shows the concordance between the embeddings obtained when only the genes selected for all methods were considered (see also 56 Fig.).

https://doi.org/10.1371/journal.pcbi.1008585.g003

and splicing rate for a gene across cells, we would thus expect a linear relationship between the observed spliced and unspliced abundance for a gene, across all cells. Similarly to scVelo, we therefore fit a zero-intercept line for the unspliced vs spliced normalized abundance for each gene across all cells, and estimate the R² for this line. A high value suggests that the gene is in a steady state across a large fraction of the cells in the data set. Indeed, all methods showed generally higher R² values for the negative control data sets (Fig.1), with a smaller difference seen for alevin_spliced_gtr, kallisto|bus_coll_incl and kallisto|bus_sep_ind. The three quantification methods that detected the strongest differences between positive and negative control datasets were alevin_coll_gtr, alevin_sep_gtr and starsolo.

Discussion and conclusions

In this study, we have compared different counting strategies for obtaining the spliced and unspliced count matrices required for RNA velocity analysis. Using five experimental droplet scRNA-seq data sets, we have shown that there are considerable differences between the count matrices obtained by different methods that are widely used in the field, and that these differences directly influence the downstream analysis and interpretation of the estimated velocities. This effect is mediated partly by an impact on the genes that are selected for inclusion by scVelo, but differences affecting the interpretation remain even when the same set of genes is used across all methods.

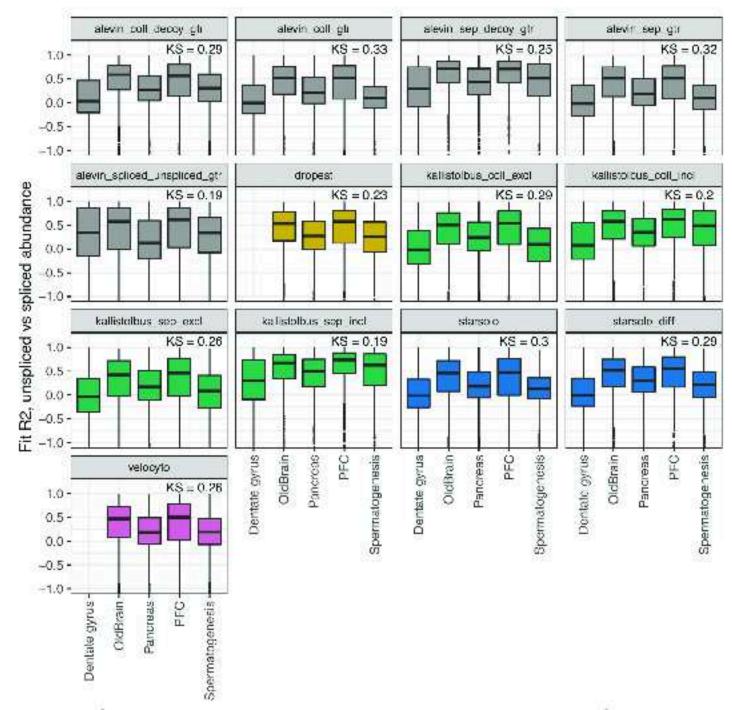


Fig 4. Estimated R² for linear fit of unspliced vs spliced abundance. The box plots show the distribution, across genes, of the estimated R² value for a linear fit of the unspliced vs spliced normalized abundances estimated by scVola. The Kolmogorov-Smimov statistic indicated within each panel is obtained by comparing the estimated R² values between the three data sets expected to contain a dynamical signal (Dentate gyrus, Pancreas, Spermatogenesis) and the two data sets where a dynamical signal is not expected (Old Brain, PPC). The y-axis is truncated at a value of -1, but the box p to is represent the full set of values.

https://doi.org/10.1371/journal.pcbi.1008585.g004

Given the relative immaturity of the RNA velocity field, and the lack of a generally accepted method for generating realistic, simulated data with known ground truth for this application, it is challenging to rank the quantification methods in terms of absolute performance. However, some clear themes emerge from our analysis, and are summarized below and in Table 4.

Table 4. Summary of the evaluated quantification methods. The first group of methods are recommended for use with droplet scRNA-seq data, owing to their good performance on the experimental data sets used in this study. The last group of methods are not recommended for use with droplet scRNA-seq data. We emphasize that this table provides a snapshot of the capabilities of the evaluated versions of the respective methods. In addition, the relative performance may be different for data from full-length scRNA-seq protocols.

Metho d ID	Agreement	Strand- aware	Counting of multi-gene reads	Comments	
alevin_sep_gtr	good	yes	yes		
slevin_coll_gtr	good	yes	yes		
kallisto bus_sep_exel	good	no	no	Reads in ambiguous regions typically discarded	
starsalo	good	yes	no		
alevin_coll_decoy_gtr	variable	yes	yes		
velocyto	variable	yes	no		
kallisto bus_coll_e xcl	variable	no	no		
alevin_sep_decoy_gtr	poor	yes	yes	Reads in ambiguous regions double-counted.	
alevin_spliced_un splic ed_gt r	variable	yes	yes	3' tag protocols does not provide enough information for spliced/unspliced resolution.	
dropest	poor	no	partially	Insufficient UMI collapsing	
kallisto bus_coll_i ncl	poor	no	no		
kallisto bus_sep_i nel	poor	no	no	Reads in ambiguous regions double-counted.	
starsolo_diff	variable	yes	partially	Non-zero unspliced count possible even without introns; nominally negative unspliced count possible.	

Agreement with expected velocity field for the experimental data sets included in this study.

https://doi.org/10.1371/journal.pcbi.1008585.t004

First, counting exonic and intronic reads separately, without consideration of whether the read could have resulted from the other type of feature (exemplified here mainly by alevin_sep_decoy_gtr and kallisto|bus_sep_incl) leads to double-counting of reads, and velocities that agree less well with expectations.

Second, not considering the strandedness of the reads from 10x Genomics (here exemplified by the kallisto|bustools variants, dropest, and by explicitly running alevin in unstranded mode) implies that many reads in regions where genes on different strands overlap each other (exonically or intronically) are considered ambiguous. Depending on the method, these reads may consequently be excluded from the quantifications. However, at the same time it provides the ability to include reads resulting from discordant internal priming.

Third, deriving the intronic reads by subtracting the 'Gene' count from STARsalo from the corresponding 'GeneFull' count (here denoted starsalo_diff) sometimes has unexpected consequences, since the 'GeneFull' counting considers all reads that overlap the gene locus, while the 'Gene' counting requires that the reads are consistent with the transcript model. Thus, a gene can obtain a nonzero 'intronic' count despite completely lacking annotated introns. Moreover, genes located within introns of other genes will often obtain a nominally negative intronic UMI count since non-junction-spanning reads mapping to the former will be considered ambiguous and discarded in the 'GeneFull' counting.

Fourth, for 3' tag data such as the 10x Genomics data we have considered in this study, quantifying the spliced and unspliced transcripts (rather than spliced transcripts and introns only) implies that for a large fraction of the reads, it is difficult to resolve whether they stem from the spliced or unspliced target. Thus, this type of reference may be more suitable for fulllength scRNA-seq protocols, where reads are sampled across the entire length of the transcript.

Among the counting strategies contrasted in this manuscript, alevin_sep_gtr, kallisto bus_ sep_excl, starsalo and alevin_coll_gtr provided velocity embeddings most in line with expectations in the three real data sets. The methods are also relatively similar in terms of run time, and considerably faster than running CellRanger+velocyto. However, even among these well-performing methods, there are large differences in the assigned counts as well as in the handling of ambiguous reads and genomic regions.

Going forward, we expect that improvements in counting strategies for scRNA-seq data, specifically tailored for RNA velocity preprocessing, will likely come alongside an increased understanding of the read generation process and the biases underlying specific scRNA-seq library preparation protocols, and that different counting strategies may be optimal for different types of scRNA-seq data. An increased understanding of the read generation process will also enable realistic simulation of sets of spliced and unspliced scRNA-seq reads, which in turn will provide an improved platform for objective evaluation of the performance of counting strategies.

Supporting information

S1 Fig. Schematic of the intron definition with the 'collapse' and 'separate' approaches. With the 'collapse' approach, the annotated isoforms of a gene are first collapsed before the introns are defined as any non-exonic region of the gene locus. With the 'separate' approach, introns are separately defined for each isoform. After extracting the intronic regions, a flanking region is added to each side of the introns to accommodate reads overlapping exon/intron boundaries.

(PDF)

S2 Fig. Comparison of feature extraction tools. A. The number of transcripts extracted by BUSpaRse and eisaR that have the correct or incorrect sequence, respectively, as determined by comparing to the transcriptome fasta file downloaded from Gencode. Transcripts are stratified by whether or not they are multi-exonic, and by the strand. B. The number of transcripts for which the lengths of the introns inferred by BUSpaRse and eisaR are correct or incorrect, respectively, as determined by comparing to the introns extracted by the intronsByTranscript function from the GenomicFeatures Bioconductor package ('separate' intron definition). C. The number of single-transcript genes for which the lengths of the introns inferred by BUSpaRse and eisaR are correct or incorrect, respectively, as determined by comparing to the introns extracted by the intronsByTranscript function from the GenomicFeatures Bioconductor package ('collapse' intron definition).

(PDF)

S3 Fig. Evaluation of the impact of changing the flank length in the intron extraction, as well as running alevin_sep_gtr in unstranded mode, in the Spermatogenesis data set. In each case, the methods are compared to the methodologically most similar among the methods discussed in the main text. alevin_sep_flankLXX_gtr, with XX set to either 20 or 40, corresponds to running alevin with introns defined using the 'separate' approach, and using a flank length equal to the read length minus (XX+1). The alevin_sep_gtr method corresponds to setting XX = 0. Further, alevin_sep_gtr_unstranded corresponds to running alevin in unstranded mode. The kb-python wrapper uses the same type of intron definition and capture approach as kallisto|bus_sep_excl, but fixes the flank length to 30bp, whereas for kallisto|bus_sep_excl, the read length minus 1 was used. Top row: scatter plot of the total spliced and unspliced count assigned to genes with the different methods. Bottom row, left: Spearman correlation between abundances and velocities for each pair of methods. Bottom row, right: A classical multidimensional scaling (MDS) plot based on the Euclidean distances among velocity values for the set of shared genes. The various modifications to the methods have an impact on the derived

velocities; however, the modified methods still cluster close to the corresponding base method. (PDF)

S4 Fig. Total UMI count across genes and cells. The bars correspond to the total UMI count, and the split of these into counts for exonic and intronic targets, for each quantification method in the OldBrain data set. In addition to the overall count (top row), the figure shows the total count after stratifying genes by the overall fraction of unique k-mers (using the 'collapse' annotation), indicated in the vertical panel headers together with the number of genes in the category. The genes for which no uniqueness information could be calculated (the 'NA' category) are those for which all transcripts are shorter than the chosen k-mer length (which was set to the read length minus one; here 56nt).

(PDF)

S5 Fig. Stratification of genes based on whether or not they obtain the same (spliced/ unspliced, respectively) count with different counting methods. For each pair of methods, the figure shows the fraction of genes that obtain the same count with the two methods (both overall and after excluding the genes assigned a count of 0 with both methods) as well as the fraction of genes where the difference between the assigned counts is non-zero, greater than 5% (of the average count assigned by the two methods) or greater than 10%. (PDF)

S6 Fig. Overlaps among the sets of highly variable genes selected by scVelo based on the different quantification approaches. Each column corresponds to the number of genes shared by a particular set of methods (indicated by black dots).
(PDF)

S7 Fig. Examples of genes illustrating major differences among the quantification methods. The coverage tracks show the number of reads overlapping each base position, on the positive (blue) or negative (red) strand. The isoforms of the main gene in focus of each panel are
similarly shown in blue (positive strand) or red (negative strand). Any overlapping features
from other genes are shown in the bottom annotation track, colored in muted blue or red,
depending on the annotated strand. The bottom panels show the total exonic and intronic
UMI count assigned to the displayed gene by the different quantification methods, as well as
the fraction of unspliced counts, and the fraction of unique k-mers in the gene overall, as well
as in the exons and introns.

(PDF)

S8 Fig. Additional examples of genes illustrating major differences among the quantification methods. The coverage tracks show the number of reads overlapping each base position,
on the positive (blue) or negative (red) strand. The isoforms of the main gene in focus of each
panel are similarly shown in blue (positive strand) or red (negative strand). Any overlapping
features from other genes are shown in the bottom annotation track, colored in muted blue or
red, depending on the annotated strand. The bottom panels show the total exonic and intronic
UMI count assigned to the displayed gene by the different quantification methods, as well as
the fraction of unspliced counts, and the fraction of unique k-mers in the gene overall, as well
as in the exons and introns.

(PDF)

S9 Fig. Average coverage of genomic regions around polyA/polyT stretches. Considered polyA/polyT stretches are at least 15 nt, with at most one mismatch per 15 nt in introns of genes on the forward and reverse strands. Similarly to [2], we observe consistent coverage

around discordant internal priming regions. (PDF)

S10 Fig. Cell- and gene-wise Spearman correlations of spliced, unspliced and total normalized abundance, as well as velocity estimates, between each pair of quantification methods. (PDF)

S11 Fig. Spearman correlation between total normalized abundance (from scVelo) and estimated velocity, and between spliced and unspliced normalized abundances, either by gene or by cell, for all quantification methods.

(PDF)

S12 Fig. Spearman correlation between average total normalized abundance across cells (from scVelo) and estimated likelihood of the velocity fit, for all quantification methods (Pancreas data).

(PDF)

S13 Fig. Spearman correlation between average total normalized abundance across cells (from scVelo) and estimated likelihood of the velocity fit, for all quantification methods (Dentate gyrus data).

(PDF)

S14 Fig. Distribution of the length of the embedded velocity vectors based on different quantification methods (panels), in different reduced dimension representations, for the Pancreas data set. Since the displacement vectors in the low-dimensional representation are normalized before a weighted average is taken to determine the velocity embedding for a given cell, the length of the velocity embedding indicates to what extent the cells to which a cell i has a high transition probability are all located in the same direction from cell i in the low-dimensional representation. Consequently, the length of the embedded velocity vector provides one way of measuring how interpretable the low-dimensional representation is, from a velocity point of view. While overall only small differences are seen between dimension reductions, UMAP generally showed a small improvement in interpretability compared to PCA and TSNE, and using only the unspliced counts as the basis for the embedding was typically less informative. (PDF)

S15 Fig. Distribution of the average similarity (dot product) between embedded velocity vectors among neighboring cells across a range of reduced dimension representations, for the Pancreas data set. A high similarity suggests that velocity streamlines in the reduced dimension representation are more easily interpretable, since they are designed to summarize the dynamics across neighbouring cells.

(PDE)

S16 Fig. Inferred velocities in the Pancreas data set, visualized on top of the same UMAP representation.

(PDF)

S17 Fig. Inferred velocities in the Spermatogenesis data set, visualized on top of the same UMAP representation.

(PDF)

S18 Fig. Inferred velocities in the Dentate gyrus data set, visualized on top of the same UMAP representation.

(PDF)

S19 Fig. Velocity confidence estimates (similarity between the estimated velocities for a cell and those of its neighbors) obtained by scVelo, for the Pancreas data set. (PDF)

S20 Fig. Velocity confidence estimates (similarity between the estimated velocities for a cell and those of its neighbors) obtained by scVelo, for the Spermatogenesis data set. (PDF)

S21 Fig. Velocity confidence estimates (similarity between the estimated velocities for a cell and those of its neighbors) obtained by scVelo, for the Dentate gyrus data set. (PDF)

Acknowledgments

The authors would like to thank Dania Machlab, Panagiotis Papasaikas and Michael I Love for helpful discussions and comments on the manuscript.

Author Contributions

Conceptualization: Charlotte Soneson, Avi Srivastava, Rob Patro, Michael B. Stadler.

Data curation: Charlotte Soneson, Avi Srivastava.

Formal analysis: Charlotte Soneson.

Funding acquisition: Rob Patro.

Methodology: Charlotte Soneson, Avi Srivastava, Rob Patro, Michael B. Stadler.

Resources: Michael B. Stadler.

Software: Charlotte Soneson, Avi Srivastava.

Supervision: Rob Patro.

Visualization: Charlotte Soneson.

Writing - original draft: Charlotte Soneson, Avi Srivastava, Michael B. Stadler.

Writing - review & editing: Charlotte Soneson, Avi Srivastava, Rob Patro, Michael B. Stadler.

References

- Nature Methods. Method of the year 2013. Nat Methods. 2014; 11(1):1. https://doi.org/10.1038/nmeth. 2801 PMID:24524124
- Eberwine J, Sul JY, Barttai T, Kim J. The promise of single-cell sequencing. Nat Methods. 2014;11.
 25-27. https://doi.org/10.1039/rmeth.2769 PMID: 24524134
- Picelli S, Faridani OR, Björklund AK, Winberg G, Sagasser S, Sandberg R. Fulf-length RNA-seq from single cells using Smart-seq2. Nat Protoc. 2014;9(1):171–181. https://doi.org/10.1038/nprot.2014.008. PMID: 24385147
- Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital transcriptional profiling of single cells. Nat Commun. 2017; 8:14049. https://doi.org/10.1038/ncomms14049. PMID: 28091601
- Vanail, Hashimshony T. CEL-Seq2-Single-Cell RNA Sequencing by Multiplexed Linear Amplification. Methods Mol Biol. 2019; 1979:45–56. https://doi.org/10.1007/978-1-4939-9240-9_4 PMID: 31028831.
- Kulkarni A, Anderson AG, Merulio DP, Konopka G. Beyond bulk: a review of single cell transcriptomics methodologies and applications. Curr Opin Biotechnol. 2019; 58:129–136. https://doi.org/10.1016/j. cophia 2019.03.001 PMID: 30978843
- Lederer AR, La Manno G. The emergence and promise of single-cell temporal-omics approaches. Curr Opin Biotechnol. 2020; 63:70–78. https://doi.org/10.1016/j.copbio.2019.12.005 PMID: 31918114

- Saetens W, Cannoodt R, Todorov H, Saeys Y. A comparison of single-cell trajectory interence methods. Nat Biotechnot. 2019; 37:547–554. https://doi.org/10.1038/s41587-019-0071-9 PMID: 30936559
- La Manno G, Soldatov R, Zeisel A, Braun E, Hochgerner H, Petukhov V, et al. RNA velocity of single cells. Nature. 2018; 560:494–498. https://doi.org/10.1038/s41586-018-0414-6 PMID: 30089906.
- Zeisel A, Köstler WJ, Molotski N, Tsai JM, Krauthgamer R, Jacob-Hirsch J, et al. Coupled pre-mRNA and mRNA dynamics unveil operational strategies underlying transcriptional responses to stimuli. Mol Syst Biol. 2011; 7:529. https://doi.org/10.1038/msb.2011.62 PMID: 21915116
- Gaidatzis D, Burger L, Florescu M, Stadier MB. Analysis of intronic and exonic reads in RNA-seq data characterizes transcriptional and post-transcriptional regulation. Nat Biotechnol. 2015; 33:722–729. https://doi.org/10.1038/nbt.3269.PMID: 26098447
- Gray JM, Harmin DA, Boswell SA, Cloonan N, Mullen TE, Ling JJ, et al. SnapShot-Seq: A method for extracting genome-wide, in Vivo mRNA dynamics from a single total RNA sample. PLoS One. 2014; 9: e89673. https://doi.org/10.1371/journal.pone.0089673.PMID:24586954
- Bergen V, Lange M, Peidli S, Wolf FA, Theis FJ. Generalizing RNA velocity to transient cell states through dynamical modeling. Nat Biotechnol. 2020; https://doi.org/10.1038/s41587-020-0591-3 PMID: 307477-69
- Melsted P, Ntranos V, Pachter L. The Barcode, UMI, Settlormat and BUStools. Bioinformatics. 2019; 35:4472-4473. https://doi.org/10.1093bioinformatics/btz279 PMID: 31073610
- Meisted P, Booeshaghi AS, Gao F, da Veiga Beltrame E, Lu L, Hjorfeitsson KE, et al. Modular and efficient pre-processing of single-cell RNA-seq. bio Rxiv. 2019;
- Srivastava A, Malik L, Smith T, Sudbery I, Patro R. Alevin efficiently estimates accurate gene abundances from dscRNA-seq data. Genome Biol. 2019; 20:65. https://doi.org/10.1186/s13059-019-1670-y. PMID: 30917859
- Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNAseq aligner. Bioinformatics. 2013; 29:15–21. https://doi.org/10.1093/bioinformatics.bts635 PMID: 231.04986
- Smith T, Heger A, Sudbery I. UMI-tools: modeling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy. Genome Res. 2017; 27(3):491–499. https://doi.org/10.1101/gr. 209601.116 PMID:28100584
- Parekh S, Ziegenhain C, Vierh B, Enard W, Heilmann I. zUMIs—A tast and flexible pipeline to process RNA sequencing data with UMIs. Giga Science. 2018; 7: giy059. https://doi.org/10.1033/gig/ascience/ giy059 PMID: 29846586
- Petukhov V, Guo J, Baryawno N, Severe N, Scadden DT, Samsonova MG, et al. dropEst: pipeline for accurate estimation of molecular counts in droplet-based single-cell RNA-seq experiments. Genome Biol. 2018; 19(1):78. https://doi.org/10.1186/s13059-018-1449-6 PMID: 29921301
- Tian L, Su S, Dong X, Amann-Zaicenstein D, Biben C, Seidi A, et al. scPipe: A flexible R/Bibconductor preprocessing pipeline for single-cell RNA-sequencing data. PLoS Comput Biol. 2018; 14(8):e1008361. https://doi.org/10.1371/journal.pcbi.1008361 PMID: 30096162
- Bastidas-Ponce A, Tritschler S, Dony L, Scheibner K, Tarquis-Medina M, Salinno C, et al. Comprehensive single cell mRNA profiling reveals a detailed roadmap for pancreatic endocrinogenesis. Development. 2019; 146:dev173849. https://doi.org/10.1242/dev.173849.PMID: 31180421
- Hochgemer H, Zeisel A, Lönnerberg P, Linnarsson S. Conserved properties of dentate gyrus neurogenesis across postnatal development revealed by single-cell RNA sequencing. Nat Neurosci. 2018; 21:290–299. https://doi.org/10.1038/s41593-017-0056-2 PMID: 29333608
- Hermann BP, Cheng K, Singh A, Roa-De La Cruz L, Mutoji KN, Chen IC, et al. The Mammalian Spermatogenesis Single-Cell Transcriptome, from Spermatogonial Stem Cells to Spermatids. Cell Rep. 2018; 25:1650–1867.e8. https://doi.org/10.1016/j.celrep.2018.10.026 PMID: 30404016
- Bhattacherjee A, Djekidel MN, Chen R, Chen W, Tuesta LM, Zhang Y. Cell type-specific transcriptional programs in mouse prefrontal cortex during adolescence and addiction. Nat Commun. 2019; 10:4169. https://doi.org/10.1038/s41487-019-12054-8 PMID: 31519873
- Ximerakis M, Lipnick SL, Innes BT, Simmons SK, Adiconis X, Dionne D, et al. Single-cell transcriptomic profiling of the aging mouse brain. Nat Neurosci. 2019; 22: 1696–1708. https://doi.org/10.1038/s4.1593-019-0491-3 PMID: 31551601
- Frankish A, Diekhans M, Ferreira AM, Johnson R, Jungreis I, Loveland J, et al. G ENCO DE reference annotation for the human and mouse genomes. Nucleic Acids Res. 2019; 47 (D1):D766–D773. https:// doi.org/10.1093/nar/gky955 PMID: 30357393
- Moses L, Pachter L. BUSpaRse: kallisto | bustodis R utilities; 2019. Available from http://bioconductor.org/packages/BUSpaRse/

- Pages H. B.Sgenome: Software intrastructure for efficient representation of full genomes and their SNPs; 2019. Available from http://bioconductor.org/packages/BSgenome/
- Lawrence M, Huber W, Pages H, Aboyoun P, Carlson M, Gentleman R, et al. Software for Computing and Annotating Genomic Ranges. PLoS Comput Biol. 2013; 9. https://doi.org/10.1371/journal.pcbi. 1003118 PMID: 23950896
- Patro R, Duggal G, Love MI, Irizanty RA, Kingsford C. Salmon provides fast and bias-aware quantification of transcript expression. Nat Methods. 2017; 14:417–419. https://doi.org/10.1038/mmeth.4197. PMID: 28263959
- Srivastava A, Malik L, Sarkar H, Zakeri M, Almodaresi F, Sone son C, et al. Alignment and mapping methodology influence transcript abundance estimation. Genome Biol. 2020; 21(1):239. https://doi.org/ 10.1186/s13059-020-02151-8 PMID:32894187
- Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. Nat Biotechnol. 2016; 34:525. https://doi.org/10.1038/nbt0816-888d PMID: 27043002
- Lun A, Risso D. Single Cell Experiment: S4 Classes for Single Cell Data; 2019. Available from http:// bioconductor.org/packages/Single Cell Experiment/.
- McCarthy DJ, Campbell KR, Lun ATL, Wills QF. Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. Bioinformatics. 2017; p. btw777. https://doi.org/10. 1093/bioinformatics.btw777 PMID: 29088763
- Molnnes L, Healy J, Saul N, Großberger L. UMAP: Uniform Manifold Approximation and Projection. J. Open Source Softw. 2018; 3(29):961. https://doi.org/10.21105/joss.00961
- Quintan A.P., Hall I.M. B.EDTools: a flexible suite of utilities for comparing genomic features. Bioinformatics. 2010;26:841–842. https://doi.org/10.1093/bioinformatics.btq033.PMID: 20110279.
- KentWJ, Zweig AS, Barber G, Hinrichs AS, Karolchik D. BigWig and BigBed: enabling browsing of large distributed datasets. Bioinformatics. 2010; 26:2204–2207. https://doi.org/10.1093/bioinformatics/ btg351 PMID: 20639541
- Hahne F, Ivanek R. Visualizing Genomic Data Using Gviz and Bioconductor. In: Mathé E, Davis S, editors. Statistical Genomics: Methods and Protocols. New York; NY: Springer New York; 2016. p. 335–351
- Wolf FA, Angerer P, Theis FJ. SCANPY: large-scale single-cell gene expression data analysis. Genome Biol. 2018; 19:15. https://doi.org/10.1186/s13059-017-1382-0 PMID: 29409532