



# SCSit: A high-efficiency preprocessing tool for single-cell sequencing data from SPLiT-seq



Mei-Wei Luan<sup>a,1</sup>, Jia-Lun Lin<sup>b,1</sup>, Ye-Fan Wang<sup>a</sup>, Yu-Xiao Liu<sup>b</sup>, Chuan-Le Xiao<sup>c</sup>, Rongling Wu<sup>d</sup>, Shang-Qian Xie<sup>a,\*</sup>

<sup>a</sup> Key Laboratory of Genetics and Germplasm Innovation of Tropical Special Forest Trees and Ornamental Plants (Ministry of Education), School of Life Science, Hainan University, Haikou 570228, China

<sup>b</sup> College of Biomedical Information and Engineering, Hainan Medical University, Haikou 571199, China

<sup>c</sup> State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-sen University, Guangzhou 510060, China

<sup>d</sup> Public Health Sciences and Statistics and Center for Statistical Genetics, Pennsylvania State University, Hershey, PA, USA

## ARTICLE INFO

### Article history:

Received 30 March 2021

Received in revised form 12 August 2021

Accepted 12 August 2021

Available online 14 August 2021

### Keywords:

SCSit

Single cell sequencing

SPLiT-seq

Preprocessing tool

Cell type identification

## ABSTRACT

SPLiT-seq provides a low-cost platform to generate single-cell data by labeling the cellular origin of RNA through four rounds of combinatorial barcoding. However, an automatic and rapid method for preprocessing and classifying single-cell sequencing (SCS) data from SPLiT-seq, which directly identified and labeled combinatorial barcoding reads and distinguished special cell sequencing data, is currently lacking. Here, we develop a high-efficiency preprocessing tool for single-cell sequencing data from SPLiT-seq (SCSit), which can directly identify combinatorial barcodes and UMI of cell types and obtain more labeled reads, and remarkably enhance the retained data from SCS due to the exact alignment of insertion and deletion. Compared with the original method used in SPLiT-seq, the consistency of identified reads from SCSit increases to 97%, and mapped reads are twice than the original. Furthermore, the runtime of SCSit is less than 10% of the original. It can accurately and rapidly analyze SPLiT-seq raw data and obtain labeled reads, as well as effectively improve the single-cell data from SPLiT-seq platform. The data and source of SCSit are available on the GitHub website <https://github.com/shang-qian/SCSit>.

© 2021 The Author(s). Published by Elsevier B.V. on behalf of Research Network of Computational and Structural Biotechnology. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 1. Introduction

Cellular heterogeneity is a widespread phenomenon in biology by which cells vary in genetic and genomic factors [1]. Cell heterogeneity may have dramatic impact on biological processes and diseases, such as bacterial sepsis, cell immunity, aneurysm and cardiomyocyte [2–5]. High-throughput single-cell sequencing (SCS) technologies, such as next generation sequencing (NGS) and third generation sequencing (TGS), have been developed and widely used to identify cell types from morphologically similar cell populations and multi-cellular tissues [6–10]. Compared with conventional sequencing technologies, SCS have an obvious advantage in cell type identification at the single-cell level, especially for low-abundance gene information that may be easily neglected in previous tissue-level studies [11]. SCS provide a cutting-edge technology to measure the real-time expression of genes in a single

cell [12–14] and reveal inter-cellular heterogeneity [15–17], which play an important role in understanding cellular features and functions in tumors [18], developmental biology [19,20], microbiology [21], and neuroscience [22,23].

At present, several single-cell sequencing technologies have been developed, including DroNC-seq [24], CROP-seq [25], LIANTI [26], and scCOOL-seq [27], scSLAM-seq [28], DART-seq [29] and TAP-seq [30]. DroNC-seq combines single nucleus RNA-seq (sNuc-seq) and Drop-seq using microfluidic beads marking up single-cell DNA, showing efficient and sensitive capabilities to identify single-cell types [24]. CROP-seq, called CRISPR droplet sequencing, enables pooled CRISPR-Cas9 screening with single-cell droplet, which facilitates high-throughput single-cell sequencing in a cost-effective way [25]. LIANTI linearly amplifies the whole genome DNA sequence by inserting the transposons in single cells, which significantly increases the depth and resolution of single-cell DNA sequencing [26]. scCOOL-seq is a single-cell complex sequencing technology that simultaneously characterizes the chromatin state, nucleosome location, DNA methylation, copy number variation and chromosome ploidy [27]. scSLAM-seq is a single-cell,

\* Corresponding author.

E-mail address: [sqianxie@foxmail.com](mailto:sqianxie@foxmail.com) (S.-Q. Xie).

<sup>1</sup> These authors contributed equally to this work

thiol-(SH)-linked alkylation of RNA for metabolic labelling sequencing which records transcriptional activity directly integrating metabolic RNA labelling and biochemical nucleoside conversion [28]. DART-seq enables multiplexed amplicon sequencing and transcriptome profiling in single cells [29]. TAP-seq, called targeted Perturb-seq, focuses single-cell RNA-seq coverage on genes of interest and permits a routine analysis of thousands of CRISPR-mediated perturbations within a single experiment [30]. Although the SCS technologies mentioned above have their own advantages and characteristics, they all require custom microfluidics or microwells for cell sorting to obtain single cells, resulting in the high cost of single-cell sequencing.

Recently, Rosenberg et al. developed a single-cell RNA-seq method, split-pool ligation-based transcriptome sequencing (SPLiT-seq), which labeled the cellular origin of RNA through four rounds of combinatorial barcoding and unique molecular identifier (UMI) (Fig. 1A) [12]. SPLiT-seq eliminated the need of single cells isolation because of the index information of DNA barcodes. The alignment of cell barcodes could be used to identify cell types from SPLiT-seq data, and this principle greatly reduced SCS cost and experimental requirements, making it to be widely used in single cell research. However, there is currently no automatic and rapid preprocessing method that enables the classification of single-cell sequencing data from SPLiT-seq. The existing methods simply based on ordinary alignment tools, such as BLAST or BWA, are time-consuming and fallibility for simultaneous determination of all three barcodes in different regions of each sequence. BLAST, an ordinary alignment tool, takes a lot of computation time in determination of all three barcodes by blastn-short and can only run serially on a single CPU core. BWA is faster, but the length of the barcodes (8 bp) is too short for BWA to take advantage for sequence alignment. Moreover, the alignment results of existing alignment tools considering the allowable 1–2 base mismatches require special screening, which are unable to automatically label barcodes with fault-tolerant matching. To date, there is a lack of an automatic and rapid method to identify and label combinatorial barcoding reads and distinguish special cell sequencing data from SPLiT-seq. Therefore, we develop a high-efficiency preprocessing tool for single-cell sequencing data from SPLiT-seq (SCSit), which automatically identifies three rounds of barcode and UMI and significant increase the clean SCS reads due to the accurate detection of insertion and deletion of barcodes in the alignment. SCSit effectively solves the classification and extraction of cell type labels from SPLiT-seq and achieves more accurate single-cell data.

## 2. Materials and methods

### 2.1. Feature of SCS data

Raw data of SPLiT-seq was sequenced on Illumina platform using 150 nucleotide (nt) kits and paired-end sequencing. Read1 included the transcript (cDNA) and R1 primer sequences, and Read2 covered UMI, three BC barcode combinations and cDNA (Fig. 1B). Thus, the identification of Read2 determines the accuracy and efficiency of cell type classification in SCS data, and it is a key step in SCSit.

### 2.2. Identification of index position of barcodes in Read2

Five contents contain UMI, three BC barcodes and cDNA, and the UMI and three barcodes in Read2 were used as specific tag to obtain labeled reads that identified different cell types (Fig. 1B). Each barcode is composed of indicator sequences of cell type (8 nt) and index sequences of barcodes (Table S1). The index sequences of barcode 1 and 2 (*index21*, 30 nt), barcode 2 and 3 (*index32*, 30 nt) were joined each end to end (Fig. 1B and Table S1), and the joint sequences (*index21* and *index32*) were used to identify each round barcode in Read2. The sequences of *index21* and *index32* were divided into 23 segments by 8 nt k-mer. Then the Read2 were mapped to 23 segments by sliding window (8 nt) and detected the position of index sequences of three barcodes. There were three situations of detection of index sequence in Read2:

(1) Complete match of index sequence: sliding window sequence (8 nt) of Read2 was completely and continuously matched to index sequence of barcode, the position of index sequence of Read2 was determined by the start of matched window1 ( $W_1$ ,  $M_s$ ) and the end of matched window23 ( $W_{23}$ ,  $M_e$ ), and  $M_e - M_s$  was equal to the length of index sequence (30 nt) (Fig. 2A).

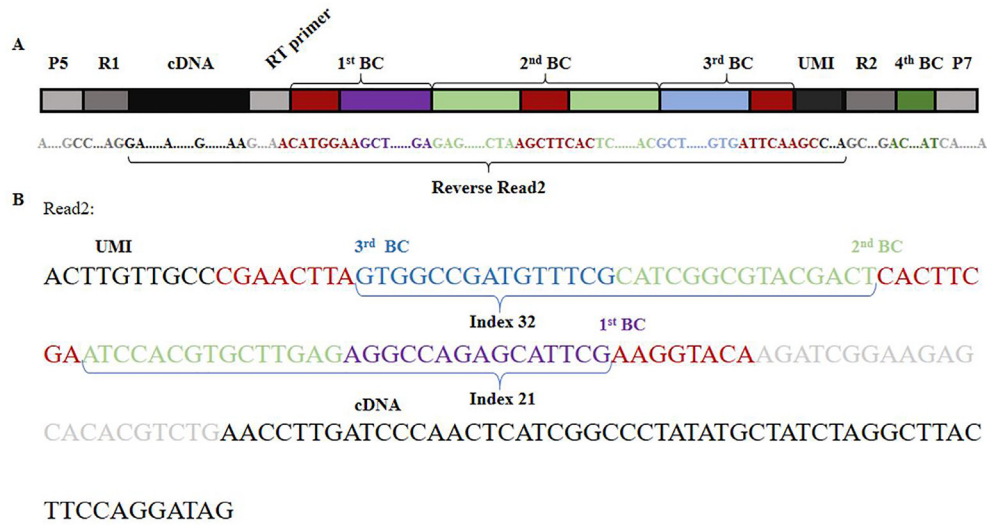
- (2) Considerate one mismatch: i) If one mismatch exists in the position  $m$  of sequences ( $9 \leq m \leq 21$ ), and  $W_1$  and  $W_{23}$  are all completely matched. The index sequence of Read2 could be identified while  $M_e - M_s = 30$ , and the start and end position of index sequences are determined from  $M_s$  to  $M_e$  (Fig. 2B). ii) If the mismatch occurs in the position  $m$  ( $1 \leq m \leq 8$ ) of matched  $W_1$ , the index sequence of Read2 could be identified while  $M_e - M_s = 30 - m$ , and the start and end position of index sequences are determined from  $M_s - m$  to  $M_e$  (Fig. 2B). iii) If the mismatch occurs in the position  $m$  ( $22 \leq m \leq 30$ ) of matched  $W_{23}$ , the index sequence of Read2 could be identified while  $M_e - M_s = 30 - m$ , and the start and end position of index sequences were determined from  $M_s$  to  $M_e + m$  (Fig. 2B).
- (3) Considerate one INDEL (length of index sequence in Read2  $\neq$  30 nt): i) If INDEL is present in the position  $m$  ( $9 \leq m \leq 21$ ) of between matched  $W_1$  and  $W_{23}$ . The index sequence of Read2 could be identified while  $M_e - M_s$  is equal to the length of matched index sequence containing insertion (31 nt) or deletion (29 nt), and the start and end position of index sequences are determined from  $M_s$  to  $M_e$  (Fig. 2C). ii) If the INDEL exists in the position  $m$  ( $1 \leq m \leq 8$ ) of matched  $W_1$ , the index sequence of Read2 could be identified while  $M_e - M_s = 30 - m$ , and the start and end position of index sequences with insertion are determined from  $M_s - m - 1$  to  $M_e$  and with deletion from  $M_s - m + 1$  to  $M_e$  (Fig. 2C). iii) If the INDEL exists in the position  $m$  ( $22 \leq m \leq 30$ ) of matched  $W_{23}$ , the index sequence of Read2 could be identified while  $M_e - M_s = 30 - m$ , and the start and end position of index sequences with insertion are determined from  $M_s$  to  $M_e + m + 1$  and with deletion from  $M_s$  to  $M_e + m$  (Fig. 2C).

### 2.3. Identification of indicator sequences of barcodes in Read2

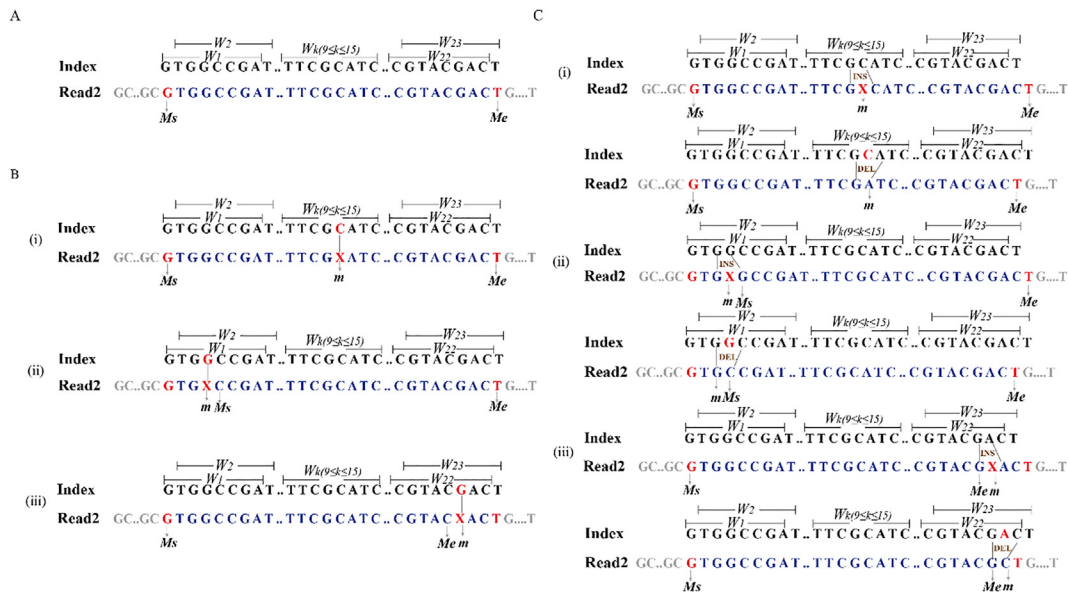
Based on the identification of round index of barcodes, the labeled reads could be obtained, which cell type could be classified by the alignment of indicator sequences of barcodes. For the identification of indicator sequences, we executed a quad to decimal conversion, and used 0, 1, 2, and 3 to present A, T, G, and C. Then the fragment sequences were converted to decimal number (*intSeq*), which was calculated by

$$\text{intSeq} = \sum_{i=0}^n \text{trans}(\text{seq}[i]) * 4^{(n-i-1)}$$

where  $\text{seq}[i]$  denotes the  $i$ th base of a fragment,  $n$  is the length of one fragment. The three round barcodes were converted and stored in decimal number lists, and the *intSeq* values were used as the index of three lists (Fig. 3A). The indicator sequences of barcodes



**Fig. 1.** Schematic overview of barcoded cDNA molecules from SPLiT-seq data (A: labeling transcriptome with SPLiT-seq, P5, P7, R1 and R2 refer PCR step of library preparation, 1st BC, 2nd BC, 3rd BC and 4th BC refer to 4 rounds barcodes, RT primer refer to reverse transcription primer, UMI refer to unique molecular identifier, B: Components of sequenced Read containing three barcodes, UMI and cDNA. Index32 is the joint sequences of 3rd BC and 2nd BC, and Index21 is the joint sequences of 2nd BC and 1st BC).



**Fig. 2.** Index sequence identification of barcode in Read2 (INS: insertion, DEL: deletion, A: complete match of index sequence. B: considerate one mismatch. C: considerate one INDEL).

in Read2 were mapped to the barcode lists by using the *intSeq* index, thus labeled reads could rapidly be marked that cell types could be rapidly identified in SCSit.

To further make the SCSit more applicable, we used the distance function to allow mismatch between indicator sequence of Read2 and three round barcodes. The distance function is defined by

$$Dist(x, y) = \sum_{i=0}^{n-1} D(x[i], y[i]) * \left(1 + \frac{i}{n}\right),$$

$$D(x, y) = \begin{cases} 0 & x = y \\ 1 & x \neq y \end{cases},$$

where  $x[i]$  and  $y[i]$  are referred to as the  $i$ th base of sequence  $x$  and  $y$ ,  $n$  referred to as the length of them. For the indicator sequence of

barcode,  $n = 8$ . If the value of *Dist* is less than threshold (2), the barcode identification of Read2 is valid, otherwise Read2 is abandoned.

#### 2.4. Identification and preprocessing of Read1

Raw Read1 sequences were retained with the valid Read2, and trim the reads by PrimerList (forward and reverse, Table S2) used in the SPLiT-seq literature. The sequences of PrimerList were divided into  $n$  segments ( $r_1, \dots, r_n$ ) by k-mer (8 nt), and similarly converted to the decimal number (*intSeq*) (Fig. 3A). Read1 sequence was divided into  $l$  ( $q_1, \dots, q_l$ ) segments. If the last segment is less than 8 bases, then it achieves the 8 bases backwards ( $q_l$ ) (Fig. 3B). Then the fragment was direct inquiry by *intSeq* index of PrimerList. The first and last matched segment of Read1 were recorded and used to trim the primer sequences of Read1.

A

Barcode round1 list			Barcode round2 list			Barcode round1 list		
No.	Sequence	intSeq	No.	Sequence	intSeq	No.	Sequence	intSeq
1	AACGTGAT	3681	1	ACAGATTC	12823	1	AACCGAGA	3976
2	AAACATCG	798	2	ATTGGCTC	5815	2	AAGACGGA	2280
3	ATGCCTAA	7120	3	CAAGGAGC	49803	3	ACACAGAA	13088
4	AGTGGTCA	9884	4	CACCTTAC	53075	4	ACGTATCA	14620
5	ACCACTGT	15577	5	CCATCCTC	61943	5	AGCAGGAA	11424
6	ACATTGGC	12651	6	CCGACAAC	63683	6	ATCCTGTA	8036
...	...	...	...	...	...	...	...	...
96	CCATCCTC	61943	96	AAGACGGA	2280	96	AGTGGTCA	9884

B

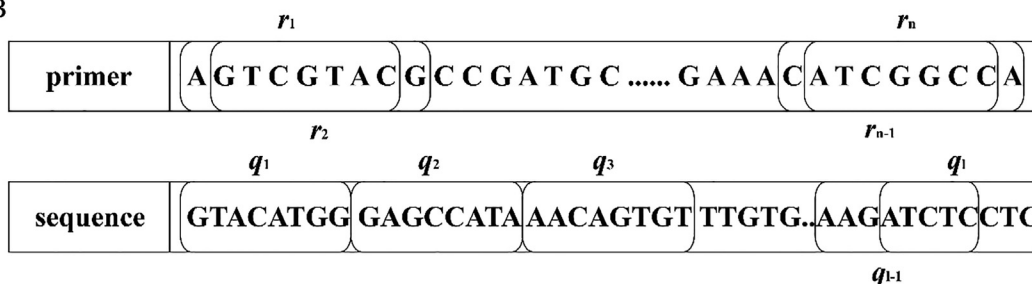


Fig. 3. Decimal conversion of barcodes and segmentation of sequences (A: decimal conversion, B: the sequence was divided into  $n$  segments with 8 nt).

### 2.5. Implementation and validation

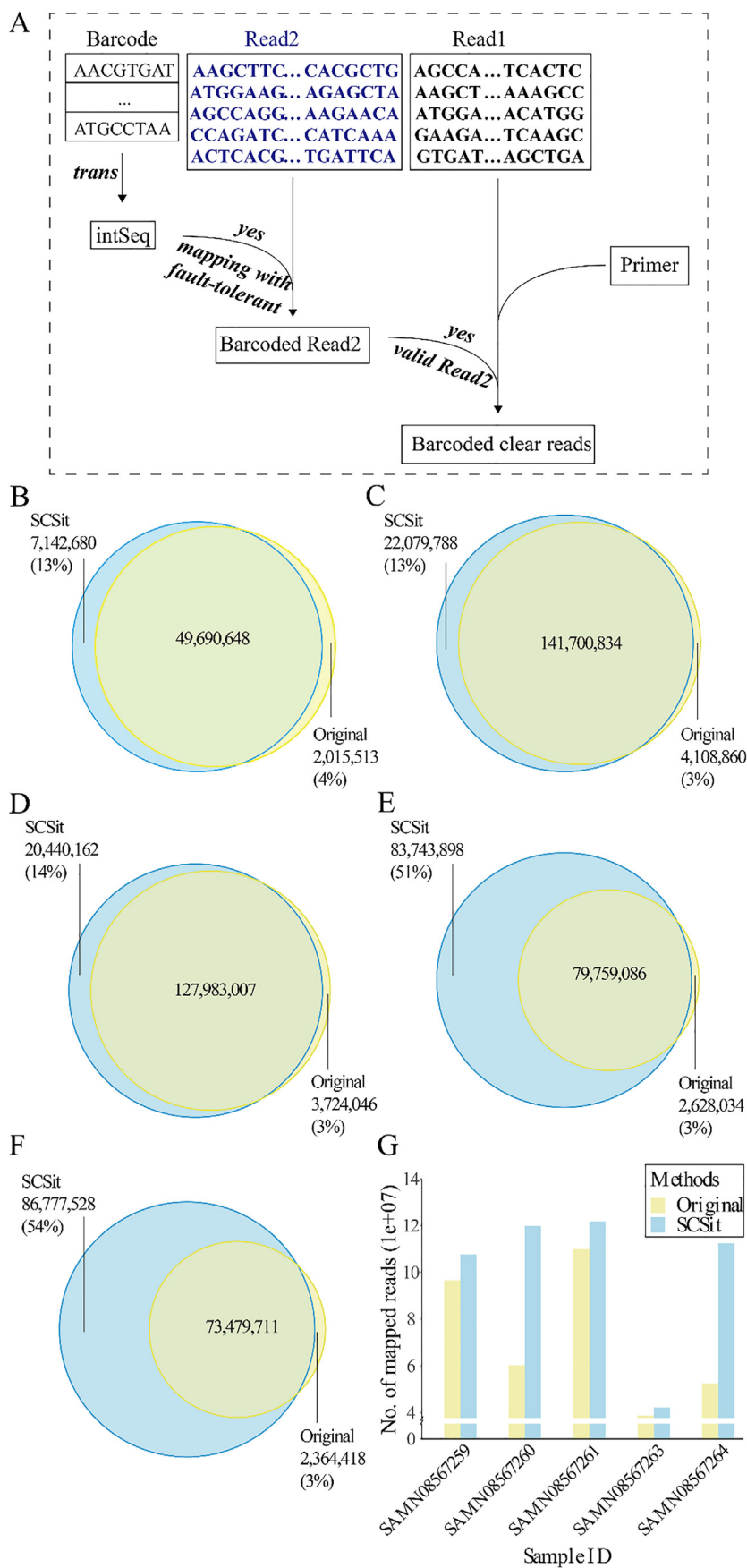
SCSit was developed by C program and parallel compute by multithread to quickly identify the cell type of SPLiT-seq reads. The raw reads FASTQ format of SPLiT-seq data as input executed SCSit program and output labeled reads with combinatorial barcodes and UMI. The output Read1 FASTQ format file was composed of combinatorial barcodes and UMI, and Read2 file was corresponding sequencing data (Fig. 4A). To validate the accuracy and reliability of SCSit, we collected and compared the treated five samples from mouse (SAMN08567263) and mixed human and mouse cells containing HEK293, HeLaS3, and NIH/3T3 (SAMN08567259, SAMN08567260, SAMN08567261, and SAMN08567264) used in the SPLiT-seq literature [12]. The original method used in SPLiT-seq discarded reads that last 6 bases of them did not match barcode sequence, UMI region were then filtered based on quality score that read with greater than 1 low-quality base ( $\text{phred} \leq 10$ ) and three 8 nt cell barcodes with more than one mismatch. We compared the SCSit and original method by using the filtered clean data. The source of SCSit and validation data were available on the GitHub website <https://github.com/shang-qian/SCSit>.

### 3. Results

Five datasets of different species from SPLiT-seq were used to perform the assessment of SCSit. The identified labeled reads and runtime were compared between SCSit and original methods (Table 1). The results indicate the identified labeled reads of SCSit in SAMN08567263 (56,833,343), SAMN08567261 (163,780,622),

SAMN08567260 (163,502,984), SAMN08567259 (148,423,169) and SAMN08567264 (160,257,239) were more than those in the original method (Table 1). And the rate of identified reads in SCSit were all more than 65% that were distinctly higher than those by the original method. The consistency of SCSit identified reads was 96 ~ 97% in the original (Table 1). The reads uniquely identified by SCSit are all more than 13 percent in five samples (Fig. 4B–F). Especially, almost double increase rates of identified reads are found in SAMN08567260 and SAMN08567264 (Table 1). The main reason for the obvious improvement of SCSit in labeled reads identification is the consideration of indel and mismatch of barcodes alignment and UMIs (Table S3), which further illustrates the necessity of developing proper method to obtain labeled reads from SPLiT-seq data. Besides, we assessed the runtime of SCSit for five datasets under 4 cores of CPU. Results demonstrate that the runtime of five samples in SCSit is less than the original method with blastn-short (Table 1). The runtime of SCSit was mainly used to identify the indicator of barcodes and trimming primer, while the original took two part time that contains barcodes alignment with blastn-short and UMI identification.

To further validate the accuracy of obtain labeled reads from SCSit, we mapped the identified reads to either the combined mm10-hg19 genome or mm10 genome using STAR [31]. The mapped reads number of SCSit are more than these in the original method in five samples, SAMN08567264 has the most incremental mapped reads (114,732,417) and twice than the original method (Table 1 and Fig. 4G). The 93 ~ 98% of uniquely mapped reads by SCSit are consistent with the reads in the original, which directly enhances the number of mapped reads (Table S3). The above results illustrate that SCSit is an accurate and efficiency tool to obtain labeled reads from SPLiT-seq.



**Fig. 4.** Pretreatment of SCS reads and comparison clear reads from SCSit and original method (A: work flow of SCSit, B-F: SAMN08567263, SAMN08567261, SAMN08567259, SAMN08567260, and SAMN08567264, G: comparison of mapped reads between SCSit and original).

**Table 1**

The statistical assessment of SCSit in five samples.

Sample ID	Method	Raw reads No.	Identified reads No. (ratio*)	Consistency rate in the original (%)	Runtime (h)	Mapped reads No.
SAMN08567263	SCSit	77,621,181	56,833,328 (73.22%)	96.10	0.62	44,541,508
	Original		51,706,161 (66.61%)		17.68	41,067,548
SAMN08567261	SCSit	218,683,580	163,780,622 (74.89%)	97.18	1.90	123,956,003
	Original		145,809,694 (66.68%)		50.34	112,226,501
SAMN08567259	SCSit	221,577,898	148,423,169 (66.98%)	97.17	1.98	109,834,181
	Original		131,707,053 (59.44%)		51.01	98,810,811
SAMN08567260	SCSit	215,597,675	163,502,984 (75.84%)	96.81	1.87	122,062,984
	Original		82,387,120 (38.21%)		49.58	62,516,718
SAMN08567264	SCSit	241,868,411	160,257,239 (66.26%)	96.88	2.15	114,732,417
	Original		75,844,129 (31.36%)		55.67	54,887,436

ratio\*: the percentage of identified reads in raw reads. Original referred to original method.

## 4. Discussion

SCSit, an automatic, rapid and accurate preprocessing tool for single-cell sequencing data, and obtain labeled reads for SPLiT-seq data which can directly identifies cell types. SPLiT-seq labels cell types by four rounds of combinatorial barcoding and UMI [12]. K-mer alignment algorithm that completely considers the mismatch and indel in barcode sequences and UMI are used to obtain labeled reads and classify cell types in SCSit, and conversion index of decimal conversion greatly improves the efficiency of alignment. The comparison of identified reads and consistency ratio with the original illustrates that SCSit has a high-efficiency preprocessing performance for cell type's identification of SCS data from SPLiT-seq.

SCSit identifies more labeled reads in five samples, the uniquely identified reads were 7,142,680, 22,079,788, 20,440,162, 83,743,898 and 86,777,528 in sample SAMN08567263, SAMN08567261, SAMN08567259, SAMN08567260 and SAMN08567264, respectively (Table S3). The unique additional reads of SCSit contain the barcodes with indel, unidentified UMI and barcodes in the original, and the barcodes absence that one of the three barcodes was missing (Table S3). The UMI and barcodes unidentified reads in original are 90 ~ 97% of uniquely identified reads in SCSit (Table S3). In order to facilitate comparison of SCSit and original method that discard the barcode sequence with more than one mismatch [12], we used one mismatch and one indel for barcodes alignment in this study. Actually, SCSit is also suitable for the case of more than one mismatch and indel in the identification of Read2. However, evaluating results from SCSit with more than one mismatch and indel show that the identified reads only increase 4 ~ 5% but takes 20 ~ 40% more runtime (Table S4). Considering the optimal balance of efficiency, we use one mismatch and one indel as default setting in SCSit.

SCSit is a high-efficiency preprocessing tool for single-cell sequencing data and provides accurately and rapidly processing of SPLiT-seq raw sequencing data for biologists and bioinformatician. In order to make researchers better understand and use SCSit, we shared the source code under the MIT license on GitHub (<https://github.com/shang-qian/SCSit>) and integrated it into Conda environment for convenient use by researchers. The alignment of barcodes with fault-tolerant and indel is the main reason for the high-efficiency and rapid preprocessing in SCSit, which are the basis of sequence alignment algorithm, so it ensures the reliability and accuracy of SCSit. Although the fault-tolerant and indel alignment were proposed for SPLiT-seq data in this study, the core principle could be widely used in other single-cell sequencing data similar to SPLiT-seq that using barcode sequence information. Recently, single cell sequencing is a hot topic in the field of life science, SCSit will be updated and improved in time to accommodate more single-cell sequencing platform data.

## 5. Conclusions

SCSit, a rapid and high-efficiency preprocessing tool for single-cell sequencing data was developed in this study. It could accurately analyze SPLiT-seq raw data and labeled reads, and effectively improved the single-cell data from SPLiT-seq platform.

## 6. Contributions

SQX conceived the project and designed the experiments, MWL and YXL collected datasets and performed the bioinformatics analysis, JLL validated the method, YFW containerized SCSit in a Conda environment, MWL, CLX, and SQX wrote the manuscript. RW critically read and revised the manuscript. All authors read and approved the final manuscript.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

This work was supported by grants from the National Natural Science Foundation of China (31760316, 32060149 and 31871326), Hainan Provincial Natural Science Foundation of China (320RC500 and ZDKJ201815), Priming Scientific Research Foundation of Hainan University (KYQD(ZR)1721).

## Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.csbj.2021.08.021>.

## References

- [1] Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, et al. Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* 2015;161(5):1202–14.
- [2] Reyes M, Filbin MR, Bhattacharyya RP, Billman K, Eisenhaure T, et al. An immune-cell signature of bacterial sepsis. *Nat Med* 2020;26(3):333–40.
- [3] Papalexi E, Satija R. Single-cell RNA sequencing to explore immune cell heterogeneity. *Nat Rev Immunol* 2018;18(1):35–45.
- [4] Liu Z, Wang L, Welch JD, Ma H, Zhou Y, et al. Single-cell transcriptomics reconstructs fate conversion from fibroblast to cardiomyocyte. *Nature* 2017;551(7678):100–4.
- [5] Zhao G, Lu H, Chang Z, Zhao Y, Zhu T, et al. Single-cell RNA sequencing reveals the cellular heterogeneity of aneurysmal infrarenal abdominal aorta. *Cardiovasc Res* 2021;117(5):1402–16.
- [6] Tung N, Battelli C, Allen B, Kaldate R, Bhatnagar S, et al. Frequency of mutations in individuals with breast cancer referred for BRCA1 and BRCA2 testing using next-generation sequencing with a 25-gene panel. *Cancer* 2015;121(1):25–33.

- [7] Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, et al. A map of human genome variation from population-scale sequencing. *Nature* 2010;467(7319):1061–73.
- [8] Montgomery SB, Sammeth M, Gutierrez-Arcelus M, Lach RP, Ingle C, et al. Transcriptome genetics using second generation sequencing in a Caucasian population. *Nature* 2010;464(7289):773–7.
- [9] Ha G, Roth A, Khattri J, Ho J, Yap D, et al. TITAN: inference of copy number architectures in clonal cell populations from tumor whole-genome sequence data. *Genome Res* 2014;24(11):1881–93.
- [10] Gawad C, Koh W, Quake SR. Single-cell genome sequencing: current state of the science. *Nat Rev Genet* 2016;17(3):175–88.
- [11] Ramskold D, Luo S, Wang YC, Li R, Deng Q, et al. Full-length mRNA-Seq from single-cell levels of RNA and individual circulating tumor cells. *Nat Biotechnol* 2012;30(8):777–82.
- [12] Rosenberg AB, Roco CM, Muscat RA, Kuchina A, Sample P, et al. Single-cell profiling of the developing mouse brain and spinal cord with split-pool barcoding. *Science* 2018;360(6385):176–82.
- [13] Xue Z, Huang K, Cai C, Cai L, Jiang CY, et al. Genetic programs in human and mouse early embryos revealed by single-cell RNA sequencing. *Nature* 2013;500(7464):593–7.
- [14] Voet T, Kumar P, Van Loo P, Cooke SL, Marshall J, et al. Single-cell paired-end genome sequencing reveals structural variation per cell cycle. *Nucleic Acids Res* 2013;41(12):6119–38.
- [15] wDavis-Marcisak EF, Sherman TD, Orugunta P, Stein-O'Brien GL, Puram SV, et al. Differential variation analysis enables detection of tumor heterogeneity using single-cell RNA-sequencing data. *Cancer Res* 2019;79(19):5102–12.
- [16] Angerer P, Simon L, Tritschler S, Wolf FA, Fischer D, et al. Single cells make big data: New challenges and opportunities in transcriptomics. *Curr Opin Syst Biol* 2017;4(1):85–91.
- [17] Vieira Braga FA, Teichmann SA, Chen X. Genetics and immunity in the era of single-cell genomics. *Hum Mol Genet* 2016;25(R2):141–8.
- [18] Ting DT, Wittner BS, Ligorio M, Vincent Jordan N, Shah AM, et al. Single-cell RNA sequencing identifies extracellular matrix gene expression by pancreatic circulating tumor cells. *Cell Rep* 2014;8(6):1905–18.
- [19] Filbin MG, Tirosh I, Hovestadt V, Shaw ML, Escalante LE, et al. Developmental and oncogenic programs in H3K27M gliomas dissected by single-cell RNA-seq. *Science* 2018;360(6386):331–5.
- [20] Marioni JC, Arendt D. How Single-Cell Genomics Is Changing Evolutionary and Developmental Biology. *Annu Rev Cell Dev Biol* 2017;33(1):537–53.
- [21] Hess M, Sczyrba A, Egan R, Kim TW, Chokhawala H, et al. Metagenomic discovery of biomass-degrading genes and genomes from cow rumen. *Science* 2011;331(6016):463–7.
- [22] Jordao MJC, Sankowski R, Brendecke SM, Sagar GL, et al. Single-cell profiling identifies myeloid cell subsets with distinct fates during neuroinflammation. *Science* 2019;363(6425):7554–71.
- [23] Tirosh I, Venteicher AS, Hebert C, Escalante LE, Patel AP, et al. Single-cell RNA-seq supports a developmental hierarchy in human oligodendroglioma. *Nature* 2016;539(7628):309–13.
- [24] Habib N, Avraham-Davidi I, Basu A, Burks T, Shekhar K, et al. Massively parallel single-nucleus RNA-seq with DroNc-seq. *Nat Methods* 2017;14(10):955–8.
- [25] Datlinger P, Rendeiro AF, Schmidl C, Krausgruber T, Traxler P, et al. Pooled CRISPR screening with single-cell transcriptome readout. *Nat Methods* 2017;14(3):297–301.
- [26] Chen C, Xing D, Tan L, Li H, Zhou G, et al. Single-cell whole-genome analyses by Linear Amplification via Transposon Insertion (LIANTI). *Science* 2017;356(6334):189–94.
- [27] Guo F, Li L, Li J, Wu X, Hu B, et al. Single-cell multi-omics sequencing of mouse early embryos and embryonic stem cells. *Cell Res* 2017;27(8):967–88.
- [28] Erhard F, Baptista MAP, Krammer T, Hennig T, Lange M, et al. scSLAM-seq reveals core features of transcription dynamics in single cells. *Nature* 2019;571(7765):419–23.
- [29] Saikia M, Burnham P, Keshavjee SH, Wang MFZ, Heyang M, et al. Simultaneous multiplexed amplicon sequencing and transcriptome profiling in single cells. *Nat Methods* 2019;16(1):59–62.
- [30] Schraivogel D, Gschwind AR, Milbank JH, Leonce DR, Jakob P, et al. Targeted Perturb-seq enables genome-scale genetic screens in single cells. *Nat Methods* 2020;17(6):629–35.
- [31] Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;29(1):15–21.