

1 **scMontage: Fast and Robust Gene Expression Similarity Search**  
2 **for Massive Single-cell Data**

3

4 **Authors' Names**

5 Tomoya Mori<sup>1,#,\*</sup>, Naila Shinwari<sup>1,\*</sup>, Wataru Fujibuchi<sup>1,†</sup>

6

7 **Authors' Address Information**

8 <sup>1</sup> Center for iPS Cell Research and Application (CiRA), Kyoto University, 53 Kawahara-cho, Sho-

9 goin, Sakyo-ku, Kyoto 606-8507, Japan

10 <sup>#</sup> Present address: Bioinformatics Center, Institute for Chemical Research, Kyoto University, Gokasho,

11 Uji, Kyoto 611-0011, Japan

12

13 <sup>\*</sup> These authors contributed equally to this work

14

15 <sup>†</sup> Corresponding author

16 E-mail: [fujibuchi-g@cira.kyoto-u.ac.jp](mailto:fujibuchi-g@cira.kyoto-u.ac.jp) (Fujibuchi W)

17 **Abstract**

18 Single-cell RNA-seq (scRNA-seq) analysis is widely used to characterize cell types or detect  
19 heterogeneity of cell states at much higher resolutions than ever before. Here we introduce scMontage  
20 (<https://scmontage.stemcellinformatics.org>), a gene expression similarity search server dedicated to  
21 scRNA-seq data, which can rapidly compare a query with thousands of samples within a few seconds.  
22 The scMontage search is based on Spearman's rank correlation coefficient and its robustness is  
23 ensured by introducing Fisher's Z-transformation and Z-test. Furthermore, search results are linked to  
24 a human cell database SHOGoiN (<http://shogoin.stemcellinformatics.org>), which enable users to fast  
25 access to additional cell-type specific information. The scMontage is available not only as a web server  
26 but also as a stand-alone application for user's own data, and thus it enhances the reliability and  
27 throughput of cell analysis and helps users gain new insights into their research.

28

29 **KEYWORDS:** Human Cell Atlas; Massive single-cell data; Gene expression profile similarity;  
30 Cell type analysis; Fisher Z-transformation

## 31 **Introduction**

32 Technology for single-cell analysis has evolved to reveal cell profiles at much higher resolutions than  
33 ever before. As an example, several studies have demonstrated that the computational analysis of  
34 single-cell RNA-seq (scRNA-seq) data can discover novel cells or cell subtypes. The recently  
35 launched Human Cell Atlas (HCA) project [1] is expected to further accelerate the production of  
36 single-cell data on an extraordinary scale. These unprecedented massive-scale data will be available  
37 to the public through International Nucleotide Sequence Database Collaboration (INSDC) sites such  
38 as the Gene Expression Omnibus (GEO) [2] and the Sequence Read Archive (SRA) [3]. Thus, data  
39 mining by very fast gene expression profile similarity searches has become increasingly important in  
40 terms of screening, clustering, and finding cells.

41 The concept of similarity searches for gene expression profiles was proposed nearly 20 years  
42 ago [4]. CellMontage [6] is the first practical and large-scale implementation that provides users quick  
43 searches against a large-scale microarray database for similar gene expression profiles based on  
44 Spearman's rank correlation.

45 Here, we propose scMontage, a renovated gene expression similarity search server, which  
46 is developed for analyzing massive-scale scRNA-seq data, based on the SHOGoin human cell type  
47 database (<http://shogoin.stemcellinformatics.org>) with statistically robust Fisher's Z-transformed  
48 correlation coefficient. Currently, the scMontage server provides human and mouse scRNA-seq data

49 and allows users to quickly access cell-type-specific biological information, such as cell taxonomy,  
50 lineage map, cell marker, and so on. The scMontage enhances the throughput and reliability of single-  
51 cell analysis and helps users gain new insights into massive scRNA-seq data.

52

## 53 **Results**

54 A profile search in scMontage can be implemented by selecting a database and inputting a query  
55 profile. After selecting the database by specifying the organism and the platform, the user can limit  
56 the genes for calculation to particular types according to Gene Ontology [7]. As a query, it is possible  
57 to either upload a gene expression profile or directly paste gene expression data in CM format.

58 **Figure 1** shows an example screen shot when human pancreatic alpha cell (GEO id:  
59 GSM1901473) is queried to the database, where ‘H. sapiens’, ‘HiS-eq2000/2500’, and  
60 ‘MF:transcription factor activity, protein binding’ are selected. The results show that the first hit is the  
61 query itself, as expected, and the top hits come from the pancreas alpha cells (**Table 1**, Table S1). The  
62 description column contains SHOGoiN Cell IDs (in parentheses) from which a user can access  
63 integrated cell type information by the SHOGoiN database. Similarly, when human pancreatic islet  
64 cell (GEO id: GSM1901455) is queried to the database under the same database setting as the previous  
65 search, the pancreatic islet cell sample is found in the top hits with high statistical significance though  
66 less number of pancreatic islet cell samples are contained in the database than the other pancreatic

67 cells (Table S2). The reliability of the scMontage search is not limited to human cell samples. Table  
68 S3 shows the search result when mouse Reg4-positive intestinal cell is queried to the database of “M.  
69 musculus, SINGLECELL: all” with “MF:transcription factor activity protein binding” genes. The top  
70 hit is the Reg4-positive intestinal cell and most of the top hits are small intestinal cells. Therefore, the  
71 scMontage search is robust not only for cell types but also for species.

72 In addition, **Figure 2** shows a comparison of statistical evaluation between CellMontage and  
73 scMontage for a mouse lung cell sample (GEO id: GSM1271921) under the database setting of  
74 “SINGLECELL: all” and “MF:transcription factor activity, protein binding”. The histograms indicate  
75 the distributions of the Spearman’s rank correlation coefficient  $r$ , the t-statistic  $t_r =$   
76  $r\sqrt{(n-2)/(1-r^2)}$ , the Z-transformed sample correlation coefficient  $z_r$ , and the standardized Z-  
77 transformed sample correlation coefficient  $z$ , respectively. In CellMontage, the distribution of  $t_r$   
78 does not follow t-distribution when the population correlation coefficient between query and database  
79 profiles is non-zero. In scMontage, however, the distribution of  $z_r$  approximately follows the normal  
80 distribution whose mean is  $z_\rho = 0.42$ . Consequently, the standardized Z-transformed sample  
81 correlation coefficient  $z$  follows the standard normal distribution.

82

## 83 **Discussion**

84 We developed scMontage that can be used for the validation and functional prediction of unknown

85 cell types obtained from tissues or derived from stem cells at the single-cell level. The scMontage also  
86 provides quick access to additional information of various cell types in the SHOGoin database from  
87 the search results. It is highly expected that a vast amount of single-cell gene expression profiles will  
88 be produced from the HCA projects or other research groups in the future. Therefore, scMontage will  
89 become an important tool for providing a very fast and powerful environment that can accelerate  
90 massive single-cell data analysis by extracting information on gene expression similarity relationships  
91 between known and unknown as well as within known/unknown cell types.

92

### 93 **Materials and methods**

94 The scMontage basically runs on Spearman's rank correlation coefficient as a similarity metric of gene  
95 expression profiles using a very fast algorithm, RaPiDS [8], for vast calculation, which enables a linear  
96 time search with a small constant for the size of the database. As a result, scMontage can compare a  
97 query with tens of thousands of samples in the database within a minute. The Spearman's rank  
98 correlation coefficient  $r$  between two rank numbers is defined as

$$r = 1 - \frac{6 \sum_{i=1}^n D_i^2}{n(n+1)(n-1)}, \quad (1)$$

99 where  $D_i$  and  $n$  indicate the rank difference between gene  $i$  and the number of genes to be used for  
100 calculation, respectively. As the output of scMontage, profiles with the highest similarity to the query  
101 are ranked by their statistical significance on the basis of the Fisher's Z-transformation of the rank

102 correlation coefficient, which is drastically improved from the CellMontage approach. The  
103 distribution of Fisher's Z-transformed sample correlation coefficient  $z_r$  approximately follows the  
104 normal distribution with a mean  $z_p$  and a standard deviation  $1/\sqrt{n-3}$  regardless of the size of  $n$ ,  
105 where  $z_p$  is approximated as the mean of  $z_r$  when it appears in standardization as the following  
106 equations:

$$z_r = \frac{1}{2} \ln \left( \frac{1+r}{1-r} \right), \quad (2)$$

$$z = \frac{z_r - z_p}{\frac{1}{\sqrt{n-3}}}. \quad (3)$$

107 Thus, scMontage can correct the significance that the population correlation coefficient between query  
108 and database profiles is non-zero, which often occurs due to common cell properties such as cell cycle  
109 states observed at single-cell level regardless of cell types.

110 The scMontage server currently provides 5,035 single-cell transcriptome data (1,424 human  
111 and 3,611 mouse cell samples on 23 August 2018) whose cell types are available by original submitters.  
112 Raw sequence data are acquired from SRA, and their read counts are computed by mapping them to  
113 human/mouse reference genome sequences downloaded from Ensembl [9] using Bowtie2 [10] and  
114 counting the mapped reads by HTSeq [11].

115 Furthermore, scMontage results are linked to the SHOGoiN database, a repository for  
116 accumulating, integrating, and providing cell information of human and other model organisms. This  
117 allows users fast access to additional cell-type-specific information, such as cell taxonomy, lineage

118 map, cell marker, DNA methylation, and morphological image.



119 **Authors' contributions**

120 WF conceptualized and designed the study. TM, NS, and WF developed the server and drafted the  
121 paper. All authors have read and approved the final manuscript.

122

123 **Competing Interests**

124 The authors declare that they have no competing interests.

125

126 **Acknowledgements**

127 This work was partially supported by the Core Center for iPS Cell Research, Research Center Network  
128 for Realization of Regenerative Medicine, Japan Agency for Medical Research and Development,  
129 Grant-in-Aid for Scientific Research on Innovative Areas, The Ministry of Education, Culture, Sports,  
130 Science and Technology, and the iPS Cell Research Fund. The authors deeply appreciate Dr. Peter  
131 Karagiannis for kindly reviewing the manuscript.

## 132 **References**

- 133 [1] Rozenblatt-Rosen O, Stubbington MJT, Regev A, Teichmann SA. The Human Cell Atlas: from  
134 vision to reality. *Nature* 2017;550:451–3.
- 135 [2] Barrett T, Troup DB, Wilhite SE, Ledoux P, Rudnev D, Evangelista C, et al. NCBI GEO: mining  
136 tens of millions of expression profiles - database and tools update. *Nucleic Acids Res* 2007;35:  
137 D760–5.
- 138 [3] Leinonen R, Sugawara H, Shumway M on behalf of the International Nucleotide Sequence  
139 Database Collaboration. The Sequence Read Archive. *Nucleic Acids Res* 2011;39:D19–21.
- 140 [4] Bassett DE Jr, Eisen MB, Boguski MS. Gene expression informatics - it's all in your time. *Nat*  
141 *Genet* 1999;21:51–5.
- 142 [5] Hunter L, Taylor RC, Leach SM, Simon R. GEST: a gene expression search tool based on a novel  
143 Bayesian similarity metric. *Bioinformatics* 2001;17:S115–22.
- 144 [6] Fujibuchi W, Kiseleva L, Taniguchi T, Harada H, Horton P. CellMontage: similar expression  
145 profile search server. *Bioinformatics* 2007;23:3103–4.
- 146 [7] The Gene Ontology Consortium. Expansion of the Gene Ontology knowledgebase and resources.  
147 *Nucleic Acids Res* 2017;45:D331–8.
- 148 [8] Horton PB, Kiseleva L, Fujibuchi W. RaPiDS: an algorithm for rapid expression profile database  
149 search. *Genome Inform* 2006;17:67–76.

- 150 [9] Zerbino DR, Achuthan P, Akanni W, Amode MR, Barrell D, Bhai J, et al. Ensembl 2018. *Nucleic*  
151 *Acids Res* 2017;46:D754–61.
- 152 [10] Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*  
153 2012;9:357–9.
- 154 [11] Anders S, Pyi PT, Huber W. HTSeq - a Python framework to work with high-throughput  
155 sequencing data. *Bioinformatics* 2015;31:166–9.

156 **Figures**

**Example gene expression profile of CM format**

```
>GSM1901473 |GPL11154 (HiSeq
2000)|H.sapiens|311000101000000000000-020 (Pancreas, Alpha
cell (Pancreatic islet))
ENSG00000000003:0 ENSG00000000005:0 ENSG000000000419:0
ENSG000000000457:347 ENSG000000000460:1132
```

**Query Input screen**

**Database settings**

Specify database (databases are constructed by platform):

**CMDB**  SHOGiN (Single-cell)

**Subset**  Ensembl

**Species**  H. sapiens  M. musculus

SINGLECELL: HiSeq2000/2500 (948)

Specify genes used to search:

- Selected by GO term (for Human/Mouse only): MF:transcription factor activity, protein binding (625)
- Upload file: Choose File no file selected **format**

**Query settings**

Specify query:

- Example query: H.sapiens: alpha cell (GPL11154:HiSeq 2000; Ensembl)
- Paste Query (enter a set of unigene ids or supported gene ids with values in CM format):
 

```
>GSM1901473 |GPL11154 (HiSeq
2000)|H.sapiens|311000101000000000000-020 (Pancreas, Alpha
cell (Pancreatic islet))
ENSG00000000003:0 ENSG00000000005:0 ENSG000000000419:0
ENSG000000000457:347 ENSG000000000460:1132
ENSG000000000938:0 ENSG000000000971:0
ENSG000000001036:1722 ENSG000000001084:0
ENSG000000001167:0
ENSG000000001460:1 ENSG000000001461:1 ENSG000000001497:0
ENSG000000001561:0 ENSG000000001617:1 ENSG000000001626:1
```
- Upload CM file: Choose File no file selected

**Query ID conversion (mandatory!)**

Select your query ID: Ensembl

Query IDs are converted into IDs of the database specified above (Ensembl).

**STEP 1** Select database, gene assignment resource, species, and platform

**STEP 2** Specify gene groups to be used when searching

**STEP 3-1** Input example query

**STEP 3-2** or paste your query

**STEP 3-3** or upload your query

**STEP 4** Select your gene assignment resource

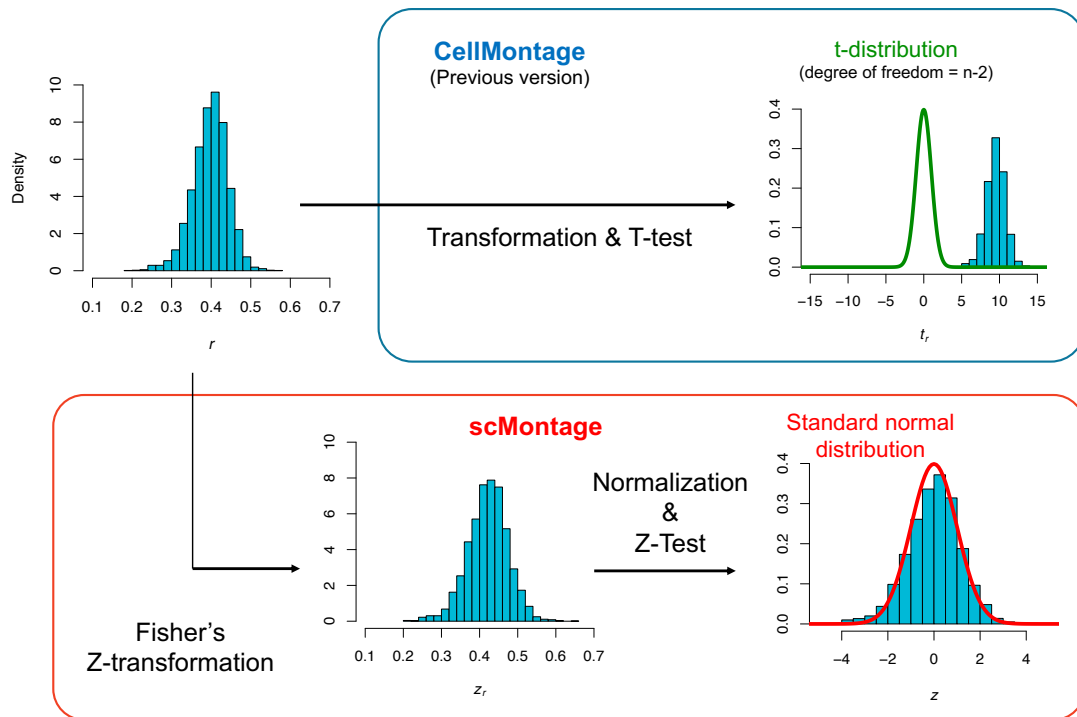
157

158 **Figure 1 Example of CM format and input screen of scMontage**

159 This example screen shot shows the case that human pancreatic alpha cell is queried to the database.

160 ‘H. sapiens’, ‘HiSeq2000/2500’, and ‘MF:transcription factor activity, protein binding’ are selected as

161 database settings. The search results are shown in Table 1.



162

163 **Figure 2 Statistical evaluations of search results from CellMontage and scMontage**  
164 **approaches**

165 The histograms indicate the distributions of the sample correlation coefficient  $z_r$ , the t-statistic  $t_r$   
166 the Z-transformed sample correlation coefficient  $z_r$ , and the standardized Z-transformed sample  
167 correlation coefficient  $z$  when a mouse lung cell sample (GSM1271921) is queried to the database  
168 under the database setting of “SINGLECELL: all” and “MF:transcription factor activity, protein  
169 binding”.

170 **Table**

171 **Table 1 Example search result**

<b>Sample id</b>	<b>Description</b>	<b>Correlation (P-value of Fisher's Z-transformed rank correlation coefficient)</b>
GSM1901473	Pancreas, Alpha cell (3110001010000000000000-020)	1.0 (0.0)
GSM1901487	Pancreas, Alpha cell (3110001010000000000000-020)	0.56 (3.6e-13)
GSM1901493	Pancreas, Alpha cell (3110001010000000000000-020)	0.54 (7.2e-12)
GSM1901488	Pancreas, Alpha cell (3110001010000000000000-020)	0.53 (1.9e-10)
GSM1901458	Pancreas, PP cell (3110001010000000000000-212)	0.53 (3.7e-10)
GSM1901497	Pancreas, Alpha cell (3110001010000000000000-020)	0.52 (1.1e-09)
GSM1901464	Pancreas, Duct cell (3110002050000000000000-090)	0.52 (1.6e-09)

172 *Note:* Search result when human pancreatic alpha cell (GEO id: GSM1901473) is queried. Numbers

173 in “Description” indicate SHOGoiN Cell IDs.

174 **Supplementary material**

175 **Supplementary Table S1 Search result when human pancreatic alpha cell (GEO id:**

176 **GSM1901473) is queried**

177

178 **Supplementary Table S2 Search result when human pancreatic islet cell (GEO id:**

179 **GSM1901455) is queried**

180

181 **Supplementary Table S3 Search result when mouse Reg4-positive intestinal cell (GEO id:**

182 **GSM1524296) is queried**