

# Journal Pre-proof

CDCP: a visualization and analyzing platform for single-cell datasets

Yuejiao Li, Tao Yang, Tingting Lai, Lijin You, Fan Yang, Jiaying Qiu, Lina Wang, Wensi Du, Cong Hua, Zhicheng Xu, Jia Cai, Zhiyong Li, Yiqun Liu, Ling Li, Minwen Zhang, Jing Chen, Lei Zhang, Dongsheng Chen, Weiwen Wang, Shiping Liu, Liang Wu, Wenjun Zeng, Bo Wang, Xiaofeng Wei, Longqi Liu, Fengzhen Chen



PII: S1673-8527(21)00374-X

DOI: <https://doi.org/10.1016/j.jgg.2021.12.004>

Reference: JGG 999

To appear in: *Journal of Genetics and Genomics*

Received Date: 1 September 2021

Revised Date: 9 December 2021

Accepted Date: 12 December 2021

Please cite this article as: Li, Y., Yang, T., Lai, T., You, L., Yang, F., Qiu, J., Wang, L., Du, W., Hua, C., Xu, Z., Cai, J., Li, Z., Liu, Y., Li, L., Zhang, M., Chen, J., Zhang, L., Chen, D., Wang, W., Liu, S., Wu, L., Zeng, W., Wang, B., Wei, X., Liu, L., Chen, F., CDCP: a visualization and analyzing platform for single-cell datasets, *Journal of Genetics and Genomics*, <https://doi.org/10.1016/j.jgg.2021.12.004>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Copyright © 2021, Institute of Genetics and Developmental Biology, Chinese Academy of Sciences, and Genetics Society of China. Published by Elsevier Limited and Science Press. All rights reserved.

## 1 **CDCP: a visualization and analyzing platform for single-cell datasets**

2  
3 Recently, the rapid advances of single-cell sequencing technologies, including  
4 sequencing of single-cell genomics, transcriptomes, epigenetics, and spatial  
5 transcriptomes, have empowered researchers to scrutinize cellular heterogeneity, gene  
6 expression, epigenetic modifications and spatial information of gene expression at  
7 single-cell resolution, which also leads to a continuous accumulation of massive single-  
8 cell datasets (Clark et al., 2016; Gawad et al., 2016). However, the scarcity of a suitable  
9 platform for analyzing single-cell datasets is a major impediment for in-depth mining  
10 of the existing data.

11  
12 Currently, there are few databases available that integrate only different aspects of  
13 datasets (Taverna et al., 2020; Longo et al., 2021; Zhao et al., 2021). For example,  
14 CancerSEA, a database covering single-cell transcriptomics of 41,900 cells from 25  
15 cancer types, aims to explore the different functional states of cancer cells at a single-  
16 cell resolution (Yuan et al., 2019). Abugessaisa et al. (2018) established SCPortalen, a  
17 database of publicly available human and mouse single-cell datasets. Consequently,  
18 regardless of these databases, there is still an immense need to build a comprehensive  
19 platform encompassing single-cell datasets covering different cell types of all species.  
20 Likewise, the continuous gathering of single-cell datasets has brought challenges for  
21 the research community, especially for bench scientists, to process and even conduct  
22 routine analyses of their own massive datasets. However, learning bioinformatics would  
23 definitely help to navigate and process these complex datasets, but it could be very  
24 time-consuming and costly. Therefore, it is imperative to establish a user-friendly  
25 dataset analysis platform that will meet the individual needs of bench researchers.

26  
27 Here, we develop a freely accessible (<https://db.cngb.org/cdcp/>) Cell-omics Data  
28 Coordinate Platform (CDCP) to share and integrate single-cell datasets. It allows  
29 visualization of each single-cell dataset with tSNE cell dimensionality reduction graphs,  
30 cluster analysis graphs of different cell types, and histograms displaying the number of

31 different cell types. In addition, the expression patterns of multiple genes in different  
32 cell types or clusters can be displayed with cluster and violin graphs. Moreover, CDCP  
33 contains detailed information on the biological samples of the datasets and allows  
34 downloading the raw sequences and expression matrices of each single-cell dataset.  
35 CDCP, unlike other platforms such as the UCSC single cell browser and ASAP,  
36 provides a platform to share and integrate single-cell transcriptomics datasets and  
37 allows users to upload their data to make it a real-time updated platform. Importantly,  
38 personalized analysis is also supported in the CDCP.

39

40 CDCP is established and maintained using the Django web framework  
41 (<https://www.djangoproject.com/>), the Python programming language and Pycharm  
42 (<https://www.jetbrains.com/pycharm/>). Centos7 (<https://www.centos.org/>) is chosen as  
43 the operating system of the server. In terms of system architecture, we use NGINX  
44 (<http://nginx.org/>) to provide static resource access, uWSGI ([https://uwsgi-](https://uwsgi-docs.readthedocs.io/en/latest/)  
45 [docs.readthedocs.io/en/latest/](https://uwsgi-docs.readthedocs.io/en/latest/)) to deploy query and download services, and Postgresql  
46 to store metadata. In terms of data security, the database is deployed in the CNGB  
47 (Wang et al., 2019), which has passed the three-level review of information security  
48 level protection and the protection capability review of trusted cloud services. Moreover,  
49 all services have been deployed with high availability. A user-friendly analysis pipeline  
50 is developed in the platform to facilitate the reanalysis of single-cell datasets of interest.  
51 By providing a single-cell expression matrix, users can calculate quality control metrics,  
52 annotate highly variable genes, perform principal component analysis, compute a  
53 neighborhood graph of observations, embed the neighborhood graph and find marker  
54 genes for characterizing groups.

55

56 **Database content and statistics.** CDCP aims to provide a platform for sharing single-  
57 cell datasets and allows query of the expression of different genes in specific cell types  
58 or clusters. By now CDCP contains expression profiles of 6467 samples and 474,573  
59 cells from public datasets in species including humans, monkeys and other animals.  
60 Although the existing datasets are collected from only 21 studies, the application of this

61 platform will attract a large number of researchers to upload and share their datasets on  
62 it.

63

64 **User interface.** A series of user-friendly interfaces are used to support various functions  
65 of CDCP. It provides five main functions: data exploration, associated databases,  
66 bioinformatics analysis, data visualization and data submission (Fig. 1A). On the “home”  
67 page, users can quickly have a view of the statistical data of the projects, samples,  
68 organisms and cells covered by CDCP (Fig. 1B). By clicking the buttons of “Project”,  
69 “Sample”, “Organism” or “Cell”, users can be directed to the “Explore” page for data  
70 retrieval. A summary of datasets for diverse organisms is presented below. Users can  
71 obtain the number of projects, samples, cells, and organs provided in different  
72 categories (Fig. S1A). In addition, a brief introduction of shared databases, including  
73 MBA, NHPCA, HCL and VThunter, is provided for users (Fig. S1A).

74

75 *Data exploration.* On the “Explore” page, one can retrieve a group of cells by selecting  
76 the attributes of the Project, Organism, Tissue, Disease, Library strategy or Release date.  
77 The query results that meet the search criteria will be displayed, including detailed  
78 information about the datasets (Fig. S1B). Actions for downloading raw data and  
79 visualizing results are also available. For data not available for download, users can  
80 apply to the submitter and download it with permission. CDCP supports visualization  
81 for data submitted with an expression matrix. If only the original sequencing data are  
82 submitted, the platform does not support visualization now. When clicking a specific  
83 project ID, users can be directed to the corresponding page that displays the detailed  
84 meta information of each single-cell sample. This includes the following information:  
85 the project name summarizing the type of the dataset, the abstract of the relevant  
86 publication, the submitter and affiliation of the dataset, and the metadata for  
87 downloading as a text file. Moreover, statistical information, including sample size,  
88 experiments and runs, is presented. By clicking on the number of samples, users can  
89 access detailed information about each sample, including age, gender, and origin.  
90 Similarly, click the number of experiments to obtain the detailed information of each

91 experiment, such as library strategy and library layout.

92

93 *Associated single-cell databases.* CDCP also provides external links to databases of  
94 MBA, NHPCA, HCL and VThunter in the “Databases” module. The Macaque Brain  
95 Atlas (MBA) depicts the macaque brain atlas and allows searching for gene expression  
96 in diverse brain cell types at single-cell resolution. The Non-Human Primate Cell Atlas  
97 (NHPCA) is a single-cell transcriptomics database that provides visualization and  
98 analysis of transcriptomics and epigenetic single-cell datasets sampled from non-  
99 human primate organs or tissues. Human Cell Landscape (HCL) illustrates a basic  
100 scheme of the human cell landscape to identify the cell type composition of human  
101 organs using single-cell RNA-seq. VThunter allows single-cell screening of virus  
102 receptor expression in the whole animal kingdom and predicts the host tropism of  
103 various viruses. These available databases will provide more comprehensive  
104 information for users interested in brain research, human, non-human primate models,  
105 and virus research.

106

107 *Data analysis.* On the “Analysis” page, the description and workflow of Codeplot, the  
108 analyzing tool, is introduced (Fig. S1C). To facilitate the analysis and sharing of single-  
109 cell transcriptome sequencing datasets, a single-cell workspace was implemented to  
110 manage relevant datasets and perform comprehensive bioinformatics analysis and  
111 visualization based on the expression matrix datasets archived in CDCP. By clicking  
112 the “single-cell workspace”, users can be directed to the Codeplot page, which provides  
113 a computing platform with a trusted execution environment for bioinformatics analysis.  
114 Users only need to provide a single cell expression matrix (csv/tsv) for data analysis.  
115 Sufficient optional adjustment parameters are provided to support user-defined  
116 parameters. Each step will generate the hdf5 file as input for the next step. We also  
117 provide each step as a workflow for users to debug.

118

119 *Data visualization.* Visualization of the selected dataset using tSNE analysis is  
120 displayed on the “Visualization” page, with dot colors representing different cell

121 clusters or cell types (Fig. 1C). Users can browse the cell clusters and cell numbers of  
122 each cell type in different projects by selecting a project and a dataset in the  
123 corresponding drop-down box presented in the hierarchical navigation menu (Fig. 1D  
124 and 1E). Moreover, by selecting the gene of interest in the “Primary Gene List” drop-  
125 down box, users can analyze the expression pattern of this gene in all clusters and cell  
126 types (Figs. 1F, 1G, S1D). In addition, the EXHIBIT Double-positive function was  
127 established to describe the co-expression patterns of two genes in different cell types  
128 and clusters with the exhibition of a tSNE plot and violin plot (Fig. S1E), which can be  
129 used to reveal functional gene annotation and transcriptional regulatory programs.

130

131 *Data submission.* CDCP allows scientists to submit their own single-cell  
132 transcriptomics datasets for sharing with the research community. Single-cell data  
133 submission supports five main entities: project, sample, experiment, run, and analysis  
134 data. The submission of single-cell datasets is mainly divided into two processes:  
135 submission of original sequencing data and submission of analysis results. By clicking  
136 “The submission of single-cell data” on the “Submission” page, users will be directed  
137 to the submission portal page (Fig. S1F). Users need to fill detailed information of each  
138 sample. Before submitting data files, project and sample information should be  
139 submitted. The detailed steps for submitting datasets and the data type and format for  
140 submission are available on the “Submission” page to guide users. For large files, it is  
141 recommended to submit data through aspera, which supports resumable upload.

142

143 **Example application.** To demonstrate the potential application of CDCP, we take the  
144 query of the human breast immune cell dataset as an example. By selecting “homo  
145 sapiens” in the “Organism” drop-down box and “breast” in the “Tissue” drop-down box,  
146 a project that met the search criteria was presented (Fig. S1B). This project depicts a  
147 single-cell atlas of 9683 immune cells from 14 triple-negative breast cancer samples  
148 (Qiu et al., 2019). By clicking the “Visualization” button on the right, users can be  
149 directed to the Visualization page. As shown by the tSNE plot, 22 cell clusters were  
150 obtained and annotated as T cells (9 clusters), macrophages (6 clusters), B cells (3

151 clusters), dendritic cells (2 clusters) and 2 unassigned clusters. T cells are the dominant  
152 group of immune cells, followed by macrophages, B cells and dendritic cells. The  
153 different characteristics of T cell clusters can be revealed by known markers and  
154 differentially expressed genes (DEGs) among clusters. For instance, among the nine  
155 clusters of T cells, clusters T7 and T9 exhibited features of regulatory T cells with high  
156 expression of *FOXP3* (ENSG00000049768) and *IL2RA* (ENSG00000134460) (Figs.  
157 S1D and S2A). The DEGs were identified among clusters, such as the expression of  
158 *CD200* (ENSG00000091972) in cluster T8, the expression of *SOX4*  
159 (ENSG00000124766) in cluster T9 and the expression of *ZNF683* (ENSG00000176083)  
160 in T6 (Fig. S2B-S2D). These results indicate that CDCP is a user-friendly, reliable and  
161 practical platform for in-depth mining of existing single-cell datasets.

162  
163 CDCP aims to provide users with a comprehensive resource of published single-cell  
164 datasets to make them easily accessible to biology researchers. As of July 2021, CDCP  
165 has collected single-cell databases from 6,467 samples and 474,573 cells. This enables  
166 retrieval and visualization of single-cell datasets and the expression profiles of a  
167 specific gene in different cell clusters or cell types. In addition, CDCP also allows  
168 visualization of the co-expression patterns of two genes for scientists interested in  
169 studying gene correlations. CDCP contains single-cell RNA-seq and ATAC-seq of  
170 different cells from humans, monkeys and other animals. This may help scientists in  
171 diverse research fields retrieve datasets of interest in a more convenient way. In addition  
172 to providing a comprehensive resource for single-cell datasets, an additional strength  
173 of CDCP is the supply of the network analysis toolkit (CodePlot), which facilitates  
174 customized data processing by users. Additional work will focus on integrating more  
175 data sources. Our future goal is to update the database with more single-cell datasets to  
176 cover more species, including animals, plants, and microbes. Despite the advantages,  
177 there are some limitations of the CDCP database; for example, not all single-cell  
178 datasets are available for download. This depends on whether the submitter has enabled  
179 downloadable permissions. Additionally, only a few datasets support visualization  
180 because some of the submitted datasets lack the expression matrix necessary for



181 automatic visualization. Further improvements will be made to ensure the submission  
182 of the expression matrix.

183

184 In summary, CDCP provides a substantial source for the research community to access  
185 and analyze single-cell datasets of interest. The diverse functions of CDCP, such as  
186 visualization and analyzing tools, may help biologists attain valuable information from  
187 the available public as well as their own datasets.

188

### 189 **Conflict of interest**

190 The authors declare no competing interests.

191

### 192 **Acknowledgments**

193 This work was supported by the China National GeneBank (CNGB).

194

### 195 **References**

196 Clark, S.J., Lee, H.J., Smallwood, S.A., Kelsey, G., Reik, W., 2016. Single-cell  
197 epigenomics: powerful new methods for understanding gene regulation and cell  
198 identity. *Genome Biol.* 17, 72.

199 Gawad, C., Koh, W., Quake, S.R., 2016. Single-cell genome sequencing: current state  
200 of the science. *Nat Rev Genet.* 17, 175-188.

201 Longo, S.K., Guo, M.G., Ji, A.L., Khavari, P.A., 2021. Integrating single-cell and spatial  
202 transcriptomics to elucidate intercellular tissue dynamics. *Nat Rev Genet.* 22,  
203 627-644.

204 Qiu, S., Hong, R., Zhuang, Z., Li, Y., Zhu, L., Lin, X., Zheng, Q., Liu, S., Zhang, K.,  
205 Huang, M., Lee, K., Lu, Q., Xia, W., Xu, F., Wang, X., Tang, J., Xiao, X., Wei,  
206 W., Yuan, Z., Shi, Y., Hou, Y., Zhang, X., Wang, J., Yang, H., Zhan, Q., Li,  
207 B., Wang, S., 2019. A Single-Cell Immune Atlas of Triple Negative Breast  
208 Cancer Reveals Novel Immune Cell Subsets. *bioRxiv.* 566968.

209 Taverna, F., Goveia, J., Karakach, T.K., Khan, S., Rohlenova, K., Treps, L.,  
210 Subramanian, A., Schoonjans, L., Dewerchin, M., Eelen, G., Carmeliet, P., 2020.  
211 BIOMEX: an interactive workflow for (single cell) omics data interpretation  
212 and visualization. *Nucleic Acids Res.* 48, W385-W394.

213 Wang, B., Liu, F., Zhang, E.C., Wo, C.L., Chen, J., Qian, P.Y., Lu, H.R., Zeng, W.J.,  
214 Chen, T., Wei, J.P., Wan, Q., Wang, R., Xu, X., 2019. [The China National  
215 GeneBank horizontal line owned by all, completed by all and shared by all]. *Yi*  
216 *Chuan.* 41, 761-772.

217 Yuan, H., Yan, M., Zhang, G., Liu, W., Deng, C., Liao, G., Xu, L., Luo, T., Yan, H.,



218 Long, Z., Shi, A., Zhao, T., Xiao, Y., Li, X., 2019. CancerSEA: a cancer single-  
219 cell state atlas. *Nucleic Acids Res.* 47, D900-D908.  
220 Zhao, T., Lyu, S., Lu, G., Juan, L., Zeng, X., Wei, Z., Hao, J., Peng, J., 2021. SC2disease:  
221 a manually curated database of single-cell transcriptome for human diseases.  
222 *Nucleic Acids Res.* 49, D1413-D1419.

223

224 *Yuejiao Li<sup>1</sup>,*  
225 *BGI-Shenzhen, Shenzhen 518083, China*

226

227 *Tao Yang<sup>1</sup>, Tingting Lai<sup>1</sup>, Lijin You<sup>1</sup>, Fan Yang<sup>1</sup>,*  
228 *China National GeneBank, Shenzhen 518120, China*

229

230 *Jiaying Qiu<sup>1</sup>,*  
231 *BGI-Shenzhen, Shenzhen 518083, China*  
232 *College of Life Sciences, University of Chinese Academy of Sciences, Beijing 100049,*  
233 *China*

234

235 *Lina Wang, Wensi Du, Cong Hua, Zhicheng Xu, Jia Cai, Zhiyong Li, Yiqun Liu, Ling*  
236 *Li, Minwen Zhang, Jing Chen, Lei Zhang,*  
237 *China National GeneBank, Shenzhen 518120, China*

238

239 *Dongsheng Chen,*  
240 *BGI-Shenzhen, Shenzhen 518083, China*  
241 *College of Life Sciences, University of Chinese Academy of Sciences, Beijing 100049,*  
242 *China*

243

244 *Weiwen Wang,*  
245 *China National GeneBank, Shenzhen 518120, China*

246

247 *Shiping Liu, Liang Wu,*  
248 *BGI-Shenzhen, Shenzhen 518083, China*

249  
250 *Wenjun Zeng, Bo Wang, Xiaofeng Wei\*,*  
251 *China National GeneBank, Shenzhen 518120, China*

252  
253 *Longqi Liu\*,*  
254 *BGI-Shenzhen, Shenzhen 518083, China*

255  
256 *Fengzhen Chen\**  
257 *China National GeneBank, Shenzhen 518120, China*

258  
259 <sup>1</sup> *These authors contributed equally*

260 *\* Corresponding authors.*

261 *Email addresses: [weixiaofeng@cngb.org](mailto:weixiaofeng@cngb.org) (X. Wei), [liulongqi@genomics.cn](mailto:liulongqi@genomics.cn) (L. Liu),*  
262 *[chenfengzhen@cngb.org](mailto:chenfengzhen@cngb.org) (F. Chen)*

263

## 264 **Figure legend**

265 **Fig. 1.** Overview and application of CDCP. **A:** Overview of the CDCP platform. **B:** The  
266 homepage of CDCP illustrating the covered single-cell datasets and five main functions.  
267 **C:** Visualization of the dataset with a tSNE cell dimensionality reduction graph. **D:** The  
268 distribution of cell clusters in different cell types. **E:** Cell numbers in different clusters  
269 are displayed. **F:** tSNE graph showing the expression profile of *FOXP3*  
270 (ENSG00000049768) in different clusters. **G:** The expression of *FOXP3*  
271 (ENSG00000049768) in different cell types.

