# Self-supervised contrastive learning for integrative single cell RNA-seq data analysis

Wenkai Han[1,2,#], Yuqi Cheng[2,3,#], Jiayang Chen[2,#], Huawen Zhong[4], Zhihang Hu[2], Siyuan Chen[1], Licheng Zong[2], Irwin King[2], Xin Gao[1,*], Yu Li[2,5,*]

[1]Computer, Electrical and Mathematical Sciences and Engineering Division (CEMSE), Computational Bioscience Research Center (CBRC), King Abdullah University of Science and Technology (KAUST), Thuwal, 23955, Saudi Arabia

[2]Department of Computer Science and Engineering (CSE), The Chinese University of Hong Kong (CUHK), Hong Kong SAR, China

[3]Weill Cornell Graduate School of Medical Sciences, Weill Cornell Medicine, New York, NY, 10065, USA

[4]Biological and Environmental Sciences & Engineering Division (BESE), Red Sea Research Center (RSRC), King Abdullah University of Science and Technology (KAUST), Thuwal, 23955, Saudi Arabia

[5]The CUHK Shenzhen Research Institute, Hi-Tech Park, Nanshan, Shenzhen, 518057, China.

[#]The first three authors contributed equally to this paper.

[*]To whom all correspondence should be addressed. E-mail: xin.gao@kaust.edu.sa, Tel: +966-12-8080323. E-mail: liyu@cse.cuhk.edu.hk, Tel: +852-39438397

1

## 1    **Abstract**

2    Single-cell RNA-sequencing (scRNA-seq) has become a powerful tool to reveal the complex

3    biological diversity and heterogeneity among cell populations. However, the technical noise

4    and bias of the technology still have negative impacts on the downstream analysis. Here, we

5    present a self-supervised Contrastive LEArning framework for scRNA-seq (CLEAR) profile

6    representation and the downstream analysis. CLEAR overcomes the heterogeneity of the

7    experimental data with a specifically designed representation learning task and thus can handle

8    batch effects and dropout events. In the task, the deep learning model learns to pull together

9    the representations of similar cells while pushing apart distinct cells, without manual labeling.

10   It achieves superior performance on a broad range of fundamental tasks, including clustering,

11   visualization, dropout correction, batch effect removal, and pseudo-time inference. The

12   proposed method successfully identifies and illustrates inflammatory-related mechanisms in a

13   COVID-19 disease study with 43,695 single cells from peripheral blood mononuclear cells.

14   Further experiments to process a million-scale single-cell dataset demonstrate the scalability

15   of CLEAR. This scalable method generates effective scRNA-seq data representation while

16   eliminating technical noise, and it will serve as a general computational framework for single-

17   cell data analysis.

1 **Introduction**

2      Single-cell RNA sequencing (scRNA-seq) has been a powerful tool for measuring the transcriptome-

3      wide gene expression in individual cells and understanding the heterogeneity among cell populations[1,]

4      [2]. It has been facilitating researchers to investigate several critical biomedical topics, such as cancer[3]

5      and autoimmunity[4]. Despite its promises, the unique properties of the scRNA-seq data, such as extreme

6      sparsity and high variability[5], have posed a number of computational challenges to researchers[6, 7]. To

7      analyze the data, among all the steps[7], the key processing is to obtain a reliable low-dimension

8      representation for each cell, which can preserve the biological signature of the cell for downstream

9      analysis while eliminating technical noise[8, 9].

10      The existing commonly used methods to perform the above processing are based on different backbone

11      algorithms and assumptions. The earliest methods utilize the traditional dimension reduction algorithms,

12      such as Principal Component Analysis (PCA), followed by $k$-means or hierarchical clustering to group

13      cells[5, 10-15]. Although these methods are widely used, their assumption, that is, the complex single-cell

14      transcriptomics can be accurately mapped onto a low-dimensional space by a generalized linear model,

15      may not be necessarily justified[8]. Considering the complexity of the data, researchers have developed

16      multiple kernel-based spectral clustering methods to learn more robust similarity matrices for cells[16, 17].

17      However, the time and space complexity of such methods impede the broad applications of the methods[5].

18      In contrast, the graph-based methods enjoy high speed and scalability[14, 15, 18, 19]. But such methods are

19      hyper-parameter sensitive. The choice of $k$ for the widely used $k$-nearest-neighbors graph affects the

20      size and number of final clusters[7]. Because of the model capacity and scalability of deep learning

21      methods, almost all the recently developed methods are based on antoencoder[5, 9, 20-25] (AE) or variational

22      autoencoder[8, 26, 27] (VAE), which can also incorporate the biostatistical models[28, 29] seamlessly. However,

23      as AE and VAE methods are unsupervised learning methods, it is very difficult to control and decide

24      what the deep learning models will learn, although some very recent studies try to impose constraints

25      and our prior knowledge about the problem onto the low-dimensional space[5, 27]. Researchers have also

26      tried to utilize manual labeling as supervision for training the models, accompanied by transfer

27      learning[22] or meta-learning[30], but such methods encounter scalability issues and have strong

28      assumptions on the homogeneity of different datasets, making them less popular than the above methods.

1    As discussed above, almost all the existing methods are based on unsupervised learning[7], regardless of

2    the specific algorithm. Without accessible supervision, for the deep learning-based methods, it is hard

3    to guide the training process of the model and explain why a particular transformation is learned,

4    although the model may work well. To promote the scRNA-seq data analysis, we indeed have some

5    specific requirements for the model. For example, the functionally similar cells should be close in the

6    transformed space, while distinct cells should be distant[7]; the model should overcome the batch effect

7    and map the cells of the same type but from different experiments into the same region[8]. Unsupervised

8    learning methods may have difficulty in incorporating these requirements explicitly. Here, we propose

9    a novel method, CLEAR, for integrative single-cell RNA-seq data analysis, based on a new machine

10   learning scenario, self-supervised learning, which can model all the above requirements explicitly.

11   More specifically, we design our method based on self-supervised contrastive learning[31], where we

12   construct the training labels from the unlabeled data. For the gene expression profile of each cell, we

13   distort the data slightly by adding noise to the raw data, which mimics the technical noise in the

14   biological experiments. During training, we force the model to produce similar low-dimension

15   representations for the raw data and the corresponding distorted profile (positive pairs). Meanwhile, we

16   train the model to output distant representations for cells of different types (negative pairs). Intuitively,

17   the deep learning model learns to pull together the representations of similar cells while pushing apart

18   different cells, only utilizing labels constructed from the data without manual labeling.

19   Based on self-supervised contrastive learning, CLEAR achieves superior performance on a broad range

20   of fundamental tasks for single-cell RNA-seq data analysis, including clustering, visualization, dropout

21   correction, batch effect removal, and pseudo-time inference. As for clustering, CLEAR can outperform

22   the popular tools and recently proposed tools on diverse datasets from different organisms. Applied on

23   a dataset from a COVID-19 disease study with 43,695 single cells from peripheral blood mononuclear

24   cells, CLEAR successfully identifies and illustrates inflammatory-related mechanisms. Further

25   experiments to process a million-scale single-cell dataset demonstrate the scalability and potential of

26   CLEAR to handle the emerging large-scale cell atlases. With the capability of generating effective

27   scRNA-seq data representation while eliminating technical noise, the proposed method can serve as a

28   general computational framework for single-cell data analysis.

1    **Results**

2    **Overview of CLEAR**

3    Unlike most existing methods, which are based on unsupervised learning to map the single-cell gene

4    expression profile to the low-dimension space, we develop CLEAR base on self-supervised learning.

5    That is, although we do not have the golden standard supervised information, such as the cell type, we

6    train the deep learning model using the supervision constructed from the unlabeled data themselves.

7    Notice that we can incorporate our prior knowledge about single-cell RNA-seq data, such as noise and

8    dropout events, into the model training process implicitly and seamlessly when we build the label from

9    the unlabeled data. More specifically, we design CLEAR based on self-supervised contrastive learning[31].

10   As shown in **Fig. 1**, eventually, we also want to train a deep learning encoder to map the gene expression

11   profile into the low-dimension space. However, in addition to that, we further want the trained model

12   to force functionally similar cells close in the transformed space while distinct cells being distant. Here,

13   the model should also be robust to technical noise, such as dropout events. That is, the profiles from the

14   same cell, no matter with or without dropout events, should be mapped into the same place in the low-

15   dimension space. Although it is difficult to estimate the noise level of the real dataset, we can add

16   simulated noise to the data and force the trained to be robust to them. Based on the above idea, we

17   design CLEAR as shown in **Fig. 1**. Given the single-cell gene expression profile, we add different

18   simulated noise, such as Gaussian noise and simulated dropout events, to it (data augmentation),

19   resulting in distorted profiles (augmented data). The raw profile and the corresponding distorted profiles

20   from the same cell are positive pairs, while the profiles from different cells are negative pairs. When

21   training the model, we force the model to produce similar representations for the positive pairs while

22   distinct for the negative pairs (contrastive learning). Intuitively, we pull together the representations of

23   functionally close cells in the low-dimension space while pushing apart the embeddings of the

24   dissimilar ones. CLEAR does not have any assumptions on the data distribution or the encoder

25   architecture. It can eliminate technical noise and generate effective scRNA-seq data representation,

26   which is suitable for a range of downstream applications, such as clustering, batch effect correction,

27   and time-trajectory inference, as discussed below.

28

1  **Overall clustering performance**

2  To access how the representation from CLEAR helps to cluster, we evaluate the proposed method,

3  combined with the $k$-means clustering algorithm, on ten published datasets with expert-annotated

4  labels[32-38]. The label information is only available during testing. We compare our model with several

5  state-of-the-art methods that are widely used for scRNA-seq data and belong to different categories,

6  including PCA-based tools (Seurat[10], SC3[11], CIDR[12], SINCERA[13]), graph-based methods (Seurat[18],

7  scGNN[24]), deep generative models (scVI[8], scDHA[21], scGNN[24], ItClust[22]), and transfer learning

8  approach (ItClust[22]). Evaluated on the same datasets with five-fold cross validation, CLEAR achieves

9  substantially better performance in clustering adjusted Rand index (ARI) score than all the other

10  methods on most datasets (**Fig. 2a**, **Supplementary Table 6**). In particular, on average, CLEAR

11  improves over the second-best method, scDHA, by 4.56% regarding the score. To evaluate the

12  performance more comprehensively, we also use other metrics, such as normalized mutual

13  information(NMI), where the compared methods show similar results. To better understand the

14  representation produced by each method, we use uniform manifold approximation (UMAP) to project

15  the internal representations into a two-dimensional space and visualize them. (**Fig. 2b**, **Supp Fig 1-9**)

16  As shown in the figure, CLEAR learns to embed similar cells within the same clusters while separate

17  dissimilar cells well among different clusters. Compared to the other methods, it produces more similar

18  clustering results as the ground truth cell annotation. Furthermore, as illustrated in **Fig. 2c,** the river plot

19  of the Hravtin dataset, comparing the CLEAR clustering and expert annotation, suggests that they are

20  nearly perfectly matched. On the other hand, scDHA tends to under-cluster the dataset, *e.g.*, the

21  interneurons are mixed up with the Exicitory cell, while Seurat is likely to over-cluster the cells, *e.g*.,

22  oligodendrocytes and Excitatory cells are split into many subclusters. Although CLEAR does not access

23  any human supervision on marker genes, it can recover the ground truth directly for this dataset,

24  suggesting that  the proposed framework can implicitly capture the data's biological features.

25  Furthermore, to demonstrate the effect of the proposed self-supervised contrastive learning settings, we

26  perform an ablation study on the data augmentation operations, removing each operation one by one

27  and recording the performance change. As shown in **Supplementary Table 8,10**, removing either

28  augmentation step will lead to the decreased ARI performance of CLEAR, which strongly indicates

1    that the introduced noise and instance discrimination task help the model to capture the real cell-cell

2    relationships. In addition, we further check the influence of using highly variable genes, discovering

3    that those low variable genes may reduce the signal ratio and harm the model performance, which is

4    consistent with the previous study[24, 39]. The design and comprehensive results of the ablation studies,

5    together with the hyper-parameter selection, are detailed in the **Supplementary Table 8-12**.

6

7    **CLEAR corrects dropout events and batch effects effectively**

8    Dropout events and batch effects are notorious in scRNA-seq data analysis, which should be handled

9    properly. We next evaluate the robustness of CLEAR when encountering dropout events. Although it

10   is impossible to recover the actual gene expression levels and determine how dropouts impact the data,

11   we simulate the dropout effects by randomly masking non-zero entries into zero with a hypergeometric

12   distribution. Given the additional artificial dropouts, clustering becomes much more difficult. We test

13   the eight competing approaches together with CLEAR on the Hravtin dataset,  containing 48,266 single

14   cells with 25,187 genes and thus 1.2 billion read counts. Among these reads, 94.2% of them are zeros.

15   We set 10%, 30%, 60%, 80% dropout rates for the non-zero entries, respectively, resulting in a masked

16   dataset with up to 98.8% zeros (**Supp Fig. 16**). CLEAR achieves the best performance in handling

17   dropout events in terms of clustering, even when 80% of the non-zero entries are masked, suggesting

18   that it is robust and has the potential to extract important features in some extreme cases. Although the

19   performance of scDHA is similar to that of CLEAR when no dropouts are introduced, it becomes worse

20   when the dropout rate is 80%. (**Fig. 2d**).

21   Although several methods have been proposed to correct batch effects, which are undesirable variability

22   in the scRNA-seq datasets from technical and biological noise, most of them work as separate modules,

23   focusing on one variable, and thus cannot generalize to the large complex atlas projects. CLEAR,

24   however, has the potential to model multiple batch effects in an end-to-end fashion. Here, we assess

25   CLEAR on correcting batch effects. Specifically, we first evaluate CLEAR on a dataset[40], consisting of

26   batches with shared cell types and biologically similar but unshared cell types. The goal of the batch

27   effect removal algorithms is to integrate common cell types while maintaining separation between

28   highly similar cells in different batches (**Methods**). As shown in **Fig. 3a**, CLEAR can separate

1   difference cell types while mix up DoubleNeg and pDC cells from different batches. The biological

2   similarity between CD141 cells and CD1C cells is also represented on the figure: the distance between

3   CD141 cell cluster and CD1C cell cluster is closer than the other two clusters. On the other hand, scVI

4   and SIMLR bring DoubleNeg and pDC cells closer but do not mix the batches well. Seurat can mitigate

5   the batch effects in DoubleNeg and pDC cells but split CD141 cells into 2 clusters. ItCLUST mixed up

6   all cells, regardless of batch and cell type, suggesting that it could not handle the dataset.

7   We further quantify the performance of different methods regarding batch effect removal with two

8   metrics, average silhouette width (ASW) and adjusted rand index(ARI), on six datasets (**Datasets**). We

9   further calculate each metric in three aspects: cell type (cARI, cASW), batch mixing ($1 - bARI$, $1 -$

10  $bASW$), and the Harmonic mean of the two (f1_ARI, f1_ASW). As shown in **Fig. 3b**, CLEAR achieves

11  the best balance between cell type identification and batch mixing. Futhermore, CLEAR outperforms

12  all the other baselines under various complex batch effects settings, even though it was not designed to

13  do so (**Supp Fig. 10-15**). In particular, on the Tabula Muris Senis cell atlas, which covers the life span

14  of a mouse and contains many batches, including cells from several mouses with different identities,

15  ages, genders, and from different chips, CLEAR mixes all the cells of the same type from different

16  batches while separating distinct cell types well.

17

18  **Pseudo-time inference**

19  Another thriving topic in single-cell RNA-seq data analysis is pseudo-time inference, also known as

20  trajectory inference. It aims to infer the ordering of cells along a one-dimensional manifold (pseudo-

21  time) from the gene expression profiles. Usually, the inferring algorithms will benefit much from better

22  data representations. Here, we evaluate whether the representation produced by CLEAR can facilitate

23  the downstream pseudo-time inference. We use the CLEAR embeddeings and the PAGA[41] algorithm

24  to generate the pseudo-time. We compare it with two other popular methods, SCANPY[14] and

25  Monocle3[42], using two mouse embryo development datasets: Yan[33] and Deng[32]. We show the cells

26  ordered by pseudo-time in **Fig. 4** . Ideally, the points should fall on the diagonal, indicating the relative

27  relationship among the cells. The time inferred with CLEAR is strongly correlated with the true

28  development stages, where Monocle3 mixes the cells from different development stages. We also use

1    the R-squared value to quantify the performance. CLEAR achieves the highest value ($R^2 = 0.957$),

2    compared with SCANPY ($R^2 = -0.014$) and Monocle3 ($R^2 = 0.884$). We further illustrate the cell

3    embeddings in the 2D space with UMAP, as shown in **Fig. 4**. The smooth lines indicate the time-

4    trajectory from different methods. The trajectories inferred by CLEAR follow the development stages

5    precisely. It starts at the zygote, goes through two cells, four cells, eight cells, and 16 cells, and finally

6    stops at the blast cells. However, for Monocle3 and SCANPY, there is no clear trajectory among the

7    cells. The cells in the early stages tend to mix, while cells in the late stages form another big group. The

8    above experiments suggest that cell embeddings from CLEAR can facilitate the downstream algorithms

9    in producing better biologically meaningful trajectories.

10

11   **CLEAR illustrates peripheral immune cells atlas and inflammatory-related mechanisms in**

12   **COVID-19.**

13   To demonstrate the application potential of CLEAR on real-world biology research, we apply it to

14   analyze a newly published COVID-19 dataset[43] (GEO accession number GSE150728), containing

15   44721 cells (43695 cells after quality control) collected from six healthy controls and seven COVID-

16   19 samples. Four of the seven COVID samples are collected from patients with acute respiratory distress

17   syndrome (ARDS) in clinical (**Fig. 5a**, **Supplementary Table 1**). We perform dimensionality reduction

18   by CLEAR and graph-based clustering, identifying 32 clusters and visualizing them via uniform

19   manifold approximation and projection (UMAP). We calculate the differential expressed genes (DEGs)

20   of each cluster to annotate cell types manually. The cell types of monocytes (CD14+ and CD16+), T

21   cells (CD4+ and CD8+), natural killer (NK) cells, B cells, plasmablasts, conventional dendritic cells

22   (DCs), plasmacytoid dendritic cells (pDC), stem cell (SC) and eosinophil, neutrophil, platelets, and red

23   blood cells (RBCs) are identified (**Fig. 5b,d**, **Supplementary Table 2**).

24   To assess the general atlas of immune responses and perturbation during different COVID-19 statuses,

25   we quantify the proportions of immune cell subsets in health donors (HDs), moderate (without ARDS),

26   or severe COVID-19 (with ARDS) individuals (**Fig. 5c**). Consistent with previous reports[43-45], several

27   immune cell subsets vary between healthy donors and COVID-19 samples, and we observe a significant

28   depletion of NK cells, DC, pDC, and CD16+ monocytes. We also note an elevated frequency of

1    plasmablasts, especially in patients with ARDS, which indicates that, together with the published

2    clinical observations[46], acute COVID-19 response may be associated with a severe humoral immune

3    response.

4    Several previous studies have shown that severe COVID-19 has been associated with dysregulated

5    immune responses, which may be induced by the abnormal activation or suppression of inflammatory

6    reaction[47-50]. To reveal inflammatory-related mechanisms in COVID-19, we perform transcription level

7    analysis on monocytes in more granularity. We first examine the expression of 'COVID cytokine storm'

8    marker genes which encode pro-inflammatory cytokines reported before produced by monocytes,

9    including *IL1B, IL2, IL6, IL10, TNF*[51, 52]. Interestingly, we do not find significant expression of these

10    pro-inflammatory genes in monocytes (**Fig. 5e**), consistent with recent research with deeper profiling

11    of immune cells[43, 49] , suggesting that COVID-19 may also present an immune suppression status. To

12    further analyze transcription changes driving monocyte response remodeling in COVID-19, we conduct

13    differential expression (DE) analysis and cellular pathway analysis by comparing COVID samples to

14    HDs. Given that the dysregulation of CD14+ monocyte plays a more dominant role in COVID-19

15    progress[53], we especially investigate the transcription profile changes in CD14+ monocytes. An

16    increased IFN-stimulated gene (ISG) set and decreased major histocompatibility complex (MHC)

17    molecules in CD14+ monocyte compared to HDs are observed (**Fig. 5f**). Scoring the samples with

18    published MHC-related genes and ISGs respectively also reveal that downregulation of MHC gene

19    expression and upregulation of ISGs are significant in CD14+ monocytes across all the COVID patients

20    (**Fig. 5g, h**). The dominant effect of the IFN response is consistent with the acute viral infection. But

21    the suppression of MHC molecules may hinder the ability to activate lymphocytes and raise an effective

22    anti-viral response. We then apply Gene Ontology (GO) analysis, combined with GSEA, to study the

23    biological pathway changes in CD14+ monocytes with different COVID statuses. Significant ISG

24    upregulation in CD14+ monocyte in moderate samples is also reflected in the pathway analysis, such

25    as Type I interferon response (**Fig. 5i**), which may indicate a more active interferon level in moderate

26    COVID patients and have the potential to become a clinical blood test marker to monitor COVID

27    progress. Interestingly, we also find a secretion pathway and myeloid leukocyte activation upregulation

1    in severe samples (**Fig. 5j**). This may suggest a dysregulated CD14+ monocytes activation in patients

2    with ARDS.

3

4    **CLEAR handles million-scale scRNA-seq datasets**

5    With the unprecedented increase in sequencing scale of the recent scRNA-seq experiment platform, the

6    ability to process million-scale single-cell sequencing datasets is increasingly essential. However, many

7    published tools require complicated parameter setting tunning and cause burdens on the users with the

8    split-merge process[54]. This has become a big challenge. The proposed method, CLEAR, is a robust and

9    scalable framework, which can resolve the problem naturally. It can perform million-level dataset

10   dimension reduction in parallel while getting rid of the tedious parameter tunning process. To test the

11   scalability of CLEAR, we apply it on a newly published million-level COVID PBMC scRNA-seq

12   dataset (GEO accession number GSE158055), which contains around 1.5 million cells from COVID

13   samples. We use CLEAR with the default parameters to conduct dimension reduction, visualizing the

14   produced representations of the dataset with UMAP. CLEAR identifies 40 clusters, which are then

15   annotated manually according to each cluster's top 100 differential expressed genes (**Fig. 6a**). Among

16   them, we find 14 subtypes and then plot selected marker genes for each cell type. Satisfyingly, a

17   significant expression track of these marker genes is obtained under the higher level (*e.g*., CD4+ T cells

18   and CD8+ T cells are combined as T cells) of these subtypes (**Fig. 6b**). Performing sensitive feature

19   extraction while eliminating technical noise on the million-scale dataset, CLEAR is an easy-to-use and

20   well-performed large-scale scRNA-seq data analysis tool, which has the potential to assist the

21   construction and refinement of cell atlases.

22

23   **Discussion and Conclusion**

24   scRNA-seq has become a powerful and essential tool in biological research. With the accumulated data

25   and the emerging cell atlases, the demand for practical computational tools to process and analyze such

26   data has never been fully satisfied. Based on the current situation, the newly developed tools to process

27   the single-cell data should, first of all, learn effective representations for the profile while eliminating

28   the technical noise within the data. Secondly, they should have sufficient scalability to handle the

11

1    million-scale unlabeled datasets in the field. Here, we introduce such a framework, CLEAR, based on

2    self-supervised contrastive learning. By introducing noise during training and forcing the model to pull

3    together the representation of functionally similar cells while pushing apart dissimilar cells with a

4    carefully designed task, we managed to train the model to produce effective representations for the

5    single-cell profile. CLEAR achieves superior performance on a broad range of fundamental tasks,

6    including clustering, visualization, dropout correction, batch effect removal, and pseudo-time inference.

7    Furthermore, it scalable enough to handle a million-scale dataset, which suggests its potential to handle

8    the emerging cell atlases.

9    In the future, CLEAR can be further developed from both the biological aspect and the machine learning

10   aspect. Regarding the biological application, obviously, CLEAR is a very flexible framework to

11   perform data integration, no matter the single-omics, multi-dataset integration (cell atlases

12   construction), or multi-omics integration (e.g., the integration of scRNA-seq and scATAC-seq data). In

13   terms of the machine learning technical details, more advanced methods to handle data imbalance and

14   incorporate prior knowledge, such as partially labeled data, should be developed. We believe that our

15   framework, CLEAR, will become an alternative approach for single-cell data analysis.

16

1 **Methods**

2 **Datasets**

3 Below we describe how we obtained and preprocessed each dataset. Notice that, to show the

4 generalization property of CLEAR, we used various datasets with different sequencing protocols

5 (Smart-Seq, Smart-Seq2, DropSeq, CEL-Seq2, *etc*.), from different tissues, and with diverse data sizes

6 (from 90 cells to 1.46 million cells). Unless otherwise noted, we obtained the cell type annotation

7 information of each dataset from the original data paper. For the Hrvatin, Kolodziejczyk, Muraro,

8 Pollen, and Tabula Muris Senis datasets, to improve the data quality, we filtered out low-quality cells

9 with fewer than 200 genes and genes expressed in less than three cells.

10 Yan dataset. The Yan dataset refers to the human preimplantation embryos and embryonic stem cells.

11 In this dataset, 90 cells were sequenced with the Tang protocol. We downloaded the dataset from

12 Hemberg Group's website. We used Scanpy to log-transform the dataset, and then each cell was

13 normalized to 10,000 read counts. After that, highly variable genes were selected. Finally, we scaled

14 the dataset to unit variance and zero mean.

15 Deng dataset. The Deng dataset refers to the mouse preimplantation embryos and embryonic stem cells

16 of mixed background. 268 cells were sequenced via two protocols, Smart-Seq and Smart-Seq2. We

17 downloaded the dataset from Hemberg Group's website and performed a log-normalization

18 transformation on the RPKM expression values with SCANPY.

19 Hrvatin dataset. The Hrvatin dataset refers to the mouse primary visual cortex cells under different

20 simulation conditions, which was downloaded from GEO database (accession number: GSE102827). It

21 contains 48,266 cells from 6-8 week-old mice, which were sequenced by DropSeq. After low-quality

22 data filtering, we performed the log transformation, per-cell count normalization, and highly variable

23 gene selection steps as mentioned above.

24 Kolodziejczyk dataset. This dataset, which was downloaded from the Hemberg Group's website,

25 contains 704 embryonic stem cells. They were sequenced with three batches under SMARTer protocol.

26 We removed the low-quality data and the spiked-in cells. The top 2,000 most variable genes were

27 selected for downstream analysis after we normalized each cell to 10,000 read counts.

1    Muraro dataset. This dataset, sequenced with the CEL-Seq2 protocol, contains 2126 cells from the

2    human pancreas. We downloaded the data from the GEO database (accession number GSE85241), and

3    removed the low-quality data and cells containing a higher number of mitochondrial genes and spike-

4    in RNAs. The highly variable genes were then selected after log transformation.

5    Pollen dataset. The dataset, sequenced with the SMARTer protocol, contains 301 cells in the developing

6    cerebral cortex from 11 populations. We downloaded it from the Hemberg Group's website, removing

7    the low-quality data and cells having a higher number of mitochondrial genes and spike-in RNAs. The

8    highly variable genes were then selected after log transformation.

9    Tabula Muris Senis dataset. The entire dataset, sequenced with 10X, contains more than 100,000 cells.

10   It was generated across the lifespan of mice, including 23 tissues and organs, but here we only focus on

11   four tissues: bladder, mammary gland, limb muscle, and diaphragm. These datasets allow us to examine

12   the batch effect and cell clustering. After downloading the raw data from

13   https://figshare.com/projects/Tabula_Muris_Senis/64982, we removed the low-quality data and cells

14   containing a higher number of mitochondrial genes and spike-in RNAs. Each cell was normalized to

15   10,000 read counts. The highly variable genes were then selected after log transformation. Finally, we

16   scaled the dataset to unit variance and zero mean.

17   Human dendritic cells dataset. It consists of human blood dendritic cell (DC) data from Villani *et al* [55].

18   We downloaded the data from https://github.com/JinmiaoChenLab/Batch-effect-removal-

19   benchmarking/tree/master/Data. The dataset is composed of two batches. Each batch contains three cell

20   types. Both batches share two cell types (pDC and double negative), while remaining one unshared

21   biologically similar cell type (CD141 and CD1C, respectively).

22   COVID PBMC dataset. This dataset (GEO accession number: GSE150728), generated by Wilk *et al* [43],

23   contains 44,271 cells sequenced with the Seq-Well platform. We have eight peripheral blood samples

24   from seven SARS-COV-2 patients and six healthy controls. We removed the cells with a higher number

25   of mitochondrial genes and spike-in RNAs. Each cell was normalized to 10,000 read counts. The highly

26   variable genes were then selected after log transformation. Finally, we scaled the dataset to unit variance

27   and zero mean. The cell type information was annotated using marker genes by experts.

1    COVID large-scale dataset. This dataset (GEO accession number: GSE158055), generated by Ren *et al*

2    [56], contains more than 1.46 million cells generated through 10X Genomics. 171 COVID-19 patients and

3    25 healthy individuals were enrolled with PBMC, BALF, PFMC, and sputum samples. We removed

4    the cells containing a higher number of mitochondrial genes and spike-in RNAs. Each cell was

5    normalized to 10,000 read counts. The highly variable genes were then selected after log transformation.

6    Finally, we scaled the dataset to unit variance and zero mean. The cell type information was annotated

7    using marker genes by experts.

8

9    **The CLEAR framework**

10   The key idea of CLEAR is to learn effective cell representations, considering noise in the data, and to

11   pull together the representation of functionally similar cells, while pushing apart dissimilar cells. We

12   achieve the goal with self-supervised contrastive learning. Given the single-cell gene expression profile,

13   we add different simulated noise, such as Gaussian noise and simulated dropout events, to it (data

14   augmentation), resulting in distorted profiles (augmented data). The raw profile and the corresponding

15   distorted profiles from the same cell are positive pairs, while the profiles from different cells are

16   negative pairs. When training the model, we force the model to produce similar representations for the

17   positive pairs while distinct for the negative pairs (contrastive learning). More specifically, by

18   discriminating the positive pairs from a large number of negatives, CLEAR learns a locally smooth

19   nonlinear mapping function $f_\theta$ that pulls together multiple distortions of a cell in the embedding space

20   and pushes away the other samples. The locally smooth function is also helpful for the global

21   embeddings. In the transformed space, cells with similar expression patterns form clusters, which are

22   likely to be cells of the same cell types. The function $f_\theta$ is parameterized by a deep neural network,

23   whose parameters can be optimized in an end-to-end manner. The detailed workflows are as below.

24   1. Data augmentation. We first perform data augmentation to generate training pairs. Each cell will

25   have two augmented versions, and thus a minibatch of $N$ cells is augmented to $2N$ cells. This step will

26   be discussed in detail in **Data augmentation**.

27   2. Constructing negative labels with data from multiple minibatches. For data in one minibatch, we can

28   consider the two data points generated from the same gene expression profile as a positive pair while

1    the other combinations as negatives. However, if we only consider the negatives within a minibatch,

2    the learned mapping function is less likely to be effective for global clustering. To make the locally

3    smooth function $f_\theta$ have a global effect, we should consider negatives from other minibatches. We

4    achieve that by maintaining a queue with data from multiple minibatches. When the current minibatch

5    is enqueued, the oldest minibatch will be dequeued. Within the queue, a specific distorted profile only

6    has one positive pair match, while all the other profiles are negatives for it. Notice that the conceptual

7    difference between minibatch and the queue arises from the hardware limitation. If the GPU memory

8    is large enough and we can feed all the data in one minibatch, we can discard the queue maintenance.

9    3. Loss function. Let $X = \{x_k \in R^G\}_{k=1}^{2MN}$ be the queue consisting of a number of gene expression

10    profiles, where $G$ denotes for the number of genes; $N$ stands for the batch size; $M$ stands for the number

11    of batches stored in the queue. In one batch, $N$ samples are augmented into $2N$ samples. Consequently,

12    the queue consisting of $M$ minibatches contains $2MN$ augmented samples. $x_k$ denotes for the $k$-th

13    (distorted) gene expression profile in the queue. For a pair of positive samples $x_i$ and $x_j$ (derived from

14    one original sample), the other $2MN - 2$ samples are treated as negatives. To distinguish the positive

15    pair from the negatives, we use the following pairwise contrastive InfoNCE loss:

16
$$L_{i,j} = -log \frac{e^{\left(x_i \cdot \frac{x_j}{\tau}\right)}}{\sum_{k=1,k\neq i}^{2MN} e^{\left(x_i \cdot \frac{x_k}{\tau}\right)}}. \tag{1}$$

17    Note that $L_{i,j}$ is asymmetrical. Suppose we put all the pairs in an order, such that $2i-1$ and $2i$ denote

18    for the paired augmentations, then the summed-up loss is:

19
$$\mathcal{L} = \frac{1}{2N} \sum_{i=1}^{N} (L_{2i-1,2i} + L_{2i,2i-1}). \tag{2}$$

20    4. Momentum update. As suggested by He *et al.* [57], a rapidly changing encoder network will reduce the

21    representations' consistency, resulting in poor performance. To deal with the problem, we utilize two

22    encoders, a slow-evolving key encoder $f_k$ and a fast-evolving query encoder $f_q$. Denoting the

23    parameters of $f_q$ as $\theta_q$ and those of $f_k$ as $\theta_k$, we update the query encoder by the normal back-

24    propagation. However, for the key encoder, we update it with momentum:

25
$$\theta_k \leftarrow m\theta_k + (1-m)\theta_q. \tag{3}$$

Here $m \in [0,1)$ is a momentum coefficient. The momentum update makes the encoder network evolve smoothly.

5. Inference. After we train the model, the query encoder network $f_q$ is the final productive network, which outputs the representation of a single cell gene expression profile. After obtaining the representations of all the cells in a dataset, we cluster the cells with the common clustering algorithms (*e.g.*, k-means algorithm, Louvain algorithm, and Leiden algorithm). Finally, cell types are assigned to the discovered clusters based on the differential expression genes in the cluster.

**Data augmentation**

Data augmentation is critical to the success of self-supervised contrastive learning. We use the following ways of data augmentation, considering noise during real experiments. Note that the augmentations are performed in a specific order (as shown below). Not all the steps will be certainly conducted, with each step having a probability of being chosen or dropped.

1. Random mask. We randomly replace some gene expression values with zero in the profile of the target cell. The mask percentage is 0.2, while the probability of executing the step is 0.5. Notice that this synthetic noise is similar to the dropout events in the single-cell sequencing experiments.

2. Gaussian noise. We randomly replace some gene expression values in the target cell profile with numbers drawn from a predefined Gaussian distribution. The noise percentage is 0.8. The mean of Gaussian distribution is 0, while the standard deviation is 0.2. The probability of executing this step is 0.5.

3. Random swap. For a gene expression profile, we randomly choose an even number of gene expression values and construct pairs from the subset, then swapping the gene expression values inside each pair. The total percentage that performs swapping is 0.1. The probability of execution is 0.5.

4. Crossover with another cell. We randomly choose another cell in the dataset as the crossover source and then select some genes from the target gene expression profile, swapping the gene expression value between the two cells. 25% of the gene expression data in one cell will be exchanged with the other cell. The probability of executing this step is 0.5. This exchanging step will not influence the next batch or the next training epoch.

5. Crossover with many cells. We randomly choose several cells in the dataset as the crossover source and some genes from the target gene expression profile, swap the expression values between the source cell and the target cells. 25% of the gene expression data in the cell will be exchanged with the selected cells. The probability of execution is 0.5. This step would not influence the next batch or the next training epoch.

**Architecture and hyperparameters**

The encoder neural network in CLEAR consists of two fully connected layers. The query encoder and the key encoder share the same architecture. The first layer has 1024 nodes, while the second layer has 128 nodes. The $ReLU$ function, defined as $ReLU(x) = max(0, x)$, is used as the nonlinear activation function after the linear transformation. We use Adam optimizer with the learning rate as 1 and the cosine learning schedule. We train the paired neural networks for 200 epochs. Temperature, $\tau$, in the CLEAR's objective function, is set to be 0.2. The momentum coefficient $m$ is 0.999. The hyper-parameters are determined using grid search with cross-validation.

**Performance evaluation**

To evaluate the standard clustering performance of the proposed method, we use the adjusted Rand index (ARI) and Normalized Mutual Information (NMI). On the other hand, to benchmark different methods' performance on batch effect removal, we utilize ARI and Average Silhouette Width (ASW). In addition, we also used cell ARI (cARI) and batch ARI (bARI), as well as cell ASW (cASW) and batch ASW (bASW). Their definitions are as follows. Note that, during evaluation, we use the default parameters for all the criteria. More quantitative measurement also shown in **Supplementary Method 3**.

ARI measures the similarity between two partitions by comparing all the pairs of the samples adjusted by random permutation.

$$ARI = \frac{RI - E[RI]}{max(RI) - E[RI]}. \qquad (4)$$

Where Rand index (RI) is defined as

1
$$RI = \frac{a + b}{C_n^2}. \qquad (5)$$

2 Where $a$ is the number of pairs that are correctly labelled in the same set of the two partitions, and $b$ is

3 the number of pairs that are correctly labelled but not in the same set of the two partitions. $C_n^2$ is the

4 total number of pairs. $E[RI]$ is the expected $RI$ from a random model. ARI ranges from -1 to 1. A value

5 close to 0 suggests random labeling, while close to 1 means the nearly perfect cell type purity. To

6 evaluate batch effect removal, we also calculate three specific kinds of ARI, cell ARI (cARI) and batch

7 ARI (bARI), and the Harmonic mean of the two f1_ARI. A higher cARI corresponds to higher cell type

8 purity, while a bARI close to zero suggests superior batch effects removal.

9 NMI measures the amount of information obtained about one partition through observing the other

10 partition, ignoring the permutations:

11
$$NMI = \frac{2I(Y;C)}{[H(Y) + H(C)]}. \qquad (6)$$

12 Where $Y$ is the class labels, and $C$ is the cluster label. $H(.)$ is the entropy, and $I(Y;C)$ measures the

13 mutual information between $Y$ and $C$.

14 ASW measures the relative distance between inter-clusters and intra-clusters. The Silhouette width (SW)

15 is defined as:

16
$$SW = \frac{b - a}{\max(a, b)}. \qquad (7)$$

17 Where $a$ is the mean distance between a sample and all the inter-cluster points, while $b$ is the mean

18 distance between a sample and all the other points in the next nearest cluster. The ASW is defined as

19 the average of all the cells' Silhouette width within the entire dataset. The range of ASW is $[-1,1]$,

20 where 1 suggests the best clustering result and -1 suggests the worst clustering result. To evaluate batch

21 effect removal, we calculate three kinds of ASW, cell ASW (cASW), batch ASW (bASW), and the

22 Harmonic mean of the two f1_ASW. A higher cASW suggests better cell type purity, while a lower

23 bASW suggests better batch mixing.

24

25 **Software comparison and settings**

19

1    To evaluate the performance of CLEAR compared with other methods, we select the below several

2    software packages for comparison. All the evaluation codes and input data follow the instruction and

3    tutorials provided by each package (**Code Availability**).

4    For baseline clustering, we compare CLEAR with R-based tools including Seurat, SC3, CIDR,

5    SINCERA, scDHA, and SIMLR, and Python-based packages, such as ItClust, scVI, and scGNN. The

6    details of the software we used are: (i) Seurat version 4.0.1 from CRAN. The parameters are set as the

7    default value provided by the tutorial. (ii) SC3 version 1.2 from Bioconductor. The key parameter,

8    svm_num_cells, which means the number of randomly selected training cells to be used for SVM

9    prediction, is set as 5000. All other parameters follow the SC3 function instruction. (iii) CIDR version

10    0.1.5 from Github (https://github.com/VCCRI/CIDR). We set all the parameters as default following

11    the README file on Github. (iv) SINCERA from Github (https://github.com/xu-lab/SINCERA). The

12    parameters follow the pipeline of the demonstrations listed on their Github. (v) scDHA from Github

13    (https://github.com/duct317/scDHA). Parameters and data input format are set following the running

14    example. We also utilize its built-in function to generate 2d visualization of the representations. (vi)

15    ItClust from Github (https://github.com/jianhuupenn/ItClust), running with the default parameters. For

16    the sake of simplicity and convenience, we set the required source dataset in the code as "baron-human"

17    across all the experiments. (vii) SIMLR version 1.18.0 from Bioconductor, running with the default

18    parameters. We select the small-scale version of code because the large-scale version cannot run in our

19    environment. (viii) scVI from Github (https://github.com/YosefLab/scvi-tools). We use K-means for

20    clustering based on the embeddings generated from the trained VAE. (ix) scGNN version 1.0.2 from

21    Github (https://github.com/juexinwang/scGNN). We run the GPU version and set the hyper-parameters

22    following their example. We include LTMG inferring in preprocessing with the corresponding given

23    option of the code. All the hyper-parameters are set following the tutorial.

24    For the dropout clustering experiment, we create a dropout_sampling function for random sampling.

25    Random seeds are set to ensure each sampling is unique. Software parameters are the same as the

26    baseline clustering experiment.

27    Seurat, SC3, CIDR, and SINCERA are run on the PC with Intel(R) Core i7-8750H CPU, Window 10

28    operation system, 32GB physical memory. The virtual memory limitation of our working environment

1   is set as 100GB RAM, R version 4.0.3. We run scDHA, SIMLR, ItClust, scVI, and scGNN on a

2   workstation with Intel(R) Xeon(R) CPU E5-2667 v4, CentOS Linux release 7.7.1908 operation system,

3   Nvidia TITAN X GPU, 503GB physical memory.

4

5   **Case study of the COVID dataset**

6   We first apply CLEAR on the published COVID PBMC scRNA-seq dataset, the parameters of CLEAR

7   are set as in the **Architecture and hyperparameters**. Based on the 128 features generated by CLEAR,

8   we run Seurat (V 4.0.1) with the parameter $Resolution = 1.2$ to cluster all the cells, by which 32

9   clusters are identified. The cell type of each cluster is annotated by the top differentially expressed

10  genes found by the FindAllMarkers function. The statistical method is Wilcoxon Rank-Sum Test, and

11  the $LogFC\ threshold = 0.25$. All the gene expression level plots are generated by the FeaturePlot

12  function with the default parameters. To conduct different expression gene (DEG) analysis, we use the

13  Wilcoxon Rank-Sum Test to search for the DEGs between each pair of monocytes obtained from the

14  three groups (i.e., the health donors (HDs), moderate and severe (ARDS) groups). We put

15  $LogFC\ threshold = 0.25$ and show negative (downregulated) genes as well. We obtain two groups of

16  DEGs for each monocyte subtype and show the result in the Supplementary material (**Supplementary**

17  **Table 4-5**). Given that gene expression score can be calculated by the AddModuleScore function, we

18  use the function to calculate ISG score and MHC score for monocytes with the pre-determined

19  interferon-stimulated gene set and MHC-related gene set (**Supplementary Table 6**). The significant

20  test is also the Wilcoxon test in **Fig 5. g** and **h**. Finally, we use the DEGs we got before to perform Gene

21  Ontology (GO) analysis for each COVID stage and run GSEA analysis on the GO results. All the

22  parameters are set as default during GO and GSEA analysis.

23

24  **Data Availability**

25  We used 10 datasets for evaluating the performance of clustering and dropouts, one dataset for

26  benchmarking the batch effects removal. Two COVID-PBMC dataset for case study. The

1     details information and the links to the publicly available sources of the 13 datasets can be

2     found in the Method part.

3     **Code Availability**

4     An open-source implementation of CLEAR is available at GitHub:

5     https://github.com/ml4bio/CLEAR, under the MIT license.

6     **Acknowledgement**

7     **Contributions**

8     **Competing Interests**

9     The authors declare no competing interests.

10     **Additional Information**

11

12   1.    Shapiro, E., Biezuner, T. & Linnarsson, S. Single-cell sequencing-based technologies
13        will revolutionize whole-organism science. *Nat Rev Genet* **14**, 618-630 (2013).
14   2.    Shalek, A.K. et al. Single-cell RNA-seq reveals dynamic paracrine control of cellular
15        variation. *Nature* **510**, 363-369 (2014).
16   3.    Maynard, A. et al. Therapy-Induced Evolution of Human Lung Cancer Revealed by
17        Single-Cell RNA Sequencing. *Cell* **182**, 1232-1251 e1222 (2020).
18   4.    van Galen, P. et al. Single-Cell RNA-Seq Reveals AML Hierarchies Relevant to Disease
19        Progression and Immunity. *Cell* **176**, 1265-1281 e1224 (2019).
20   5.    Tian, T., Zhang, J., Lin, X., Wei, Z. & Hakonarson, H. Model-based deep embedding for
21        constrained clustering analysis of single cell RNA-seq data. *Nat Commun* **12**, 1873
22        (2021).
23   6.    Stegle, O., Teichmann, S.A. & Marioni, J.C. Computational and analytical challenges
24        in single-cell transcriptomics. *Nat Rev Genet* **16**, 133-145 (2015).
25   7.    Kiselev, V.Y., Andrews, T.S. & Hemberg, M. Challenges in unsupervised clustering of
26        single-cell RNA-seq data. *Nature Reviews Genetics* **20**, 273-282 (2019).
27   8.    Lopez, R., Regier, J., Cole, M.B., Jordan, M.I. & Yosef, N. Deep generative modeling
28        for single-cell transcriptomics. *Nature Methods* **15**, 1053-+ (2018).
29   9.    Deng, Y., Bao, F., Dai, Q.H., Wu, L.F. & Altschuler, S.J. Scalable analysis of cell-type
30        composition from single-cell transcriptomics using deep recurrent learning. *Nature*
31        *Methods* **16**, 311-+ (2019).
32   10.   Satija, R., Farrell, J.A., Gennert, D., Schier, A.F. & Regev, A. Spatial reconstruction of
33        single-cell gene expression data. *Nat Biotechnol* **33**, 495-502 (2015).
34   11.   Kiselev, V.Y. et al. SC3: consensus clustering of single-cell RNA-seq data. *Nat*
35        *Methods* **14**, 483-486 (2017).
36   12.   Lin, P.J., Troup, M. & Ho, J.W.K. CIDR: Ultrafast and accurate clustering through
37        imputation for single-cell RNA-seq data. *Genome Biol* **18**, 59 (2017).

13.  Guo, M., Wang, H., Potter, S.S., Whitsett, J.A. & Xu, Y. SINCERA: A Pipeline for Single-Cell RNA-Seq Profiling Analysis. *PLoS Comput Biol* **11**, e1004575 (2015).

14.  Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol* **19**, 15 (2018).

15.  Levine, J.H. et al. Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. *Cell* **162**, 184-197 (2015).

16.  Wang, B., Zhu, J., Pierson, E., Ramazzotti, D. & Batzoglou, S. Visualization and analysis of single-cell RNA-seq data by kernel-based similarity learning. *Nat Methods* **14**, 414-416 (2017).

17.  Park, S. & Zhao, H.Y. Spectral clustering based on learning similarity matrix. *Bioinformatics* **34**, 2069-2076 (2018).

18.  Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol* **36**, 411-420 (2018).

19.  Dijk, D.v. et al. Recovering Gene Interactions from Single-Cell Data Using Data Diffusion. *Cell* **174**, 716-729.e727 (2018).

20.  Eraslan, G., Simon, L.M., Mircea, M., Mueller, N.S. & Theis, F.J. Single-cell RNA-seq denoising using a deep count autoencoder. *Nature Communications* **10**, 390 (2019).

21.  Tran, D. et al. Fast and precise single-cell data analysis using a hierarchical autoencoder. *Nature Communications* **12**, 1029 (2021).

22.  Hu, J. et al. Iterative transfer learning with neural network for clustering and cell type classification in single-cell RNA-seq analysis. *Nat Mach Intell* **2**, 607-618 (2020).

23.  Wang, J. et al. Data denoising with transfer learning in single-cell transcriptomics. *Nat Methods* **16**, 875-878 (2019).

24.  Wang, J.X. et al. scGNN is a novel graph neural network framework for single-cell RNA-Seq analyses. *Nature Communications* **12**, 1882 (2021).

25.  Li, X.J. et al. Deep learning enables accurate clustering with batch effect removal in single-cell RNA-seq analysis. *Nature Communications* **11** (2020).

26.  Ding, J.R., Condon, A. & Shah, S.P. Interpretable dimensionality reduction of single cell transcriptome data with deep generative models. *Nature Communications* **9** (2018).

27.  Ding, J. & Regev, A. Deep generative model embedding of single-cell RNA-Seq profiles on hyperspheres and hyperbolic spaces. *Nat Commun* **12**, 2554 (2021).

28.  Pierson, E. & Yau, C. ZIFA: Dimensionality reduction for zero-inflated single-cell gene expression analysis. *Genome Biol* **16** (2015).

29.  Risso, D., Perraudeau, F., Gribkova, S., Dudoit, S. & Vert, J.P. A general and flexible method for signal extraction from single-cell RNA-seq data. *Nat Commun* **9**, 284 (2018).

30.  Brbic, M. et al. MARS: discovering novel cell types across heterogeneous single-cell experiments. *Nat Methods* **17**, 1200-1206 (2020).

31.  Chen, T., Kornblith, S., M., N. & Hinton, G. A Simple Framework for Contrastive Learning of Visual Representations. *ICML-2020* (2020).

32.  Deng, Q., Ramsköld, D., Reinius, B. & Sandberg, R. Single-cell RNA-seq reveals dynamic, random monoallelic gene expression in mammalian cells. *Science* **343**, 193-196 (2014).

33.  Yan, L. et al. Single-cell RNA-Seq profiling of human preimplantation embryos and embryonic stem cells. *Nature structural & molecular biology* **20**, 1131-1139 (2013).

23

34. Pollen, A.A. et al. Low-coverage single-cell mRNA sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex. *Nature biotechnology* **32**, 1053-1058 (2014).

35. Kolodziejczyk, A.A. et al. Single cell RNA-sequencing of pluripotent states unlocks modular transcriptional variation. *Cell stem cell* **17**, 471-485 (2015).

36. Muraro, M.J. et al. A single-cell transcriptome atlas of the human pancreas. *Cell systems* **3**, 385-394. e383 (2016).

37. Hrvatin, S. et al. Single-cell analysis of experience-dependent transcriptomic states in the mouse visual cortex. *Nature neuroscience* **21**, 120-129 (2018).

38. Consortium, T.M. A single cell transcriptomic atlas characterizes aging tissues in the mouse. *Nature* **583**, 590 (2020).

39. Luecken, M.D. et al. Benchmarking atlas-level data integration in single-cell genomics. *BioRxiv* (2020).

40. Tran, H.T.N. et al. A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome biology* **21**, 1-32 (2020).

41. Wolf, F.A. et al. PAGA: graph abstraction reconciles clustering with trajectory inference through a topology preserving map of single cells. *Genome biology* **20**, 1-9 (2019).

42. Trapnell, C. et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology* **32**, 381-386 (2014).

43. Wilk, A.J. et al. A single-cell atlas of the peripheral immune response in patients with severe COVID-19. *Nat Med* **26**, 1070-1076 (2020).

44. Kuri-Cervantes, L. et al. Immunologic perturbations in severe COVID-19/SARS-CoV-2 infection. *bioRxiv* (2020).

45. Kuri-Cervantes, L. et al. Comprehensive mapping of immune perturbations associated with severe COVID-19. *Sci Immunol* **5** (2020).

46. Zhao, J. et al. Antibody Responses to SARS-CoV-2 in Patients With Novel Coronavirus Disease 2019. *Clinical Infectious Diseases* **71**, 2027-2034 (2020).

47. Choudhary, S., Sharma, K. & Silakari, O. The interplay between inflammatory pathways and COVID-19: A critical review on pathogenesis and therapeutic options. *Microb Pathog* **150**, 104673 (2021).

48. Hu, B., Huang, S. & Yin, L. The cytokine storm and COVID-19. *Journal of Medical Virology* **93**, 250-256 (2021).

49. Schulte-Schrepping, J. et al. Suppressive myeloid cells are a hallmark of severe COVID-19. *medRxiv*, 2020.2006.2003.20119818 (2020).

50. Unterman, A. et al. Single-Cell Omics Reveals Dyssynchrony of the Innate and Adaptive Immune System in Progressive COVID-19. *medRxiv*, 2020.2007.2016.20153437 (2020).

51. Guo, C. et al. Single-cell analysis of two severe COVID-19 patients reveals a monocyte-associated and tocilizumab-responding cytokine storm. *Nat Commun* **11**, 3924 (2020).

52. Ragab, D., Salah Eldin, H., Taeimah, M., Khattab, R. & Salem, R. The COVID-19 Cytokine Storm; What We Know So Far. *Frontiers in Immunology* **11** (2020).

53. Schulte-Schrepping, J. et al. Severe COVID-19 Is Marked by a Dysregulated Myeloid Cell Compartment. *Cell* **182**, 1419-1440 e1423 (2020).

54.  Deng, Y., Bao, F., Dai, Q., Wu, L.F. & Altschuler, S.J. Scalable analysis of cell-type composition from single-cell transcriptomics using deep recurrent learning. *Nat Methods* **16**, 311-314 (2019).
55.  Villani, A.-C. et al. Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors. *Science* **356** (2017).
56.  Ren, X. et al. COVID-19 immune features revealed by a large-scale single-cell transcriptome atlas. *Cell* **184**, 1895-1913. e1819 (2021).
57.  Chen, X., Fan, H., Girshick, R. & He, K. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297* (2020).

**Fig. 1 | Overview of the proposed framework, CLEAR.** The proposed method is based on self-supervised contrastive learning. For the gene expression profile of each cell, we distort the data slightly by adding noise to the raw data, which mimics the technical noise in the biological experiments. When training the deep encoder model, we force the model to produce similar low-dimension representations for the raw data and the corresponding distorted profile while distant representations for cells of different types. Intuitively, the deep learning model learns to pull together the representations of similar cells while pushing apart different cells. By considering noise during training, CLEAR can produce effective representations while eliminating technical noise for the scRNA-seq profiles. Such representations have a broad range of applications, including clustering and classification, dropout event and batch effect correction, pseudo-time inference. CLEAR is also scalable to million-scale datasets without any overhead.
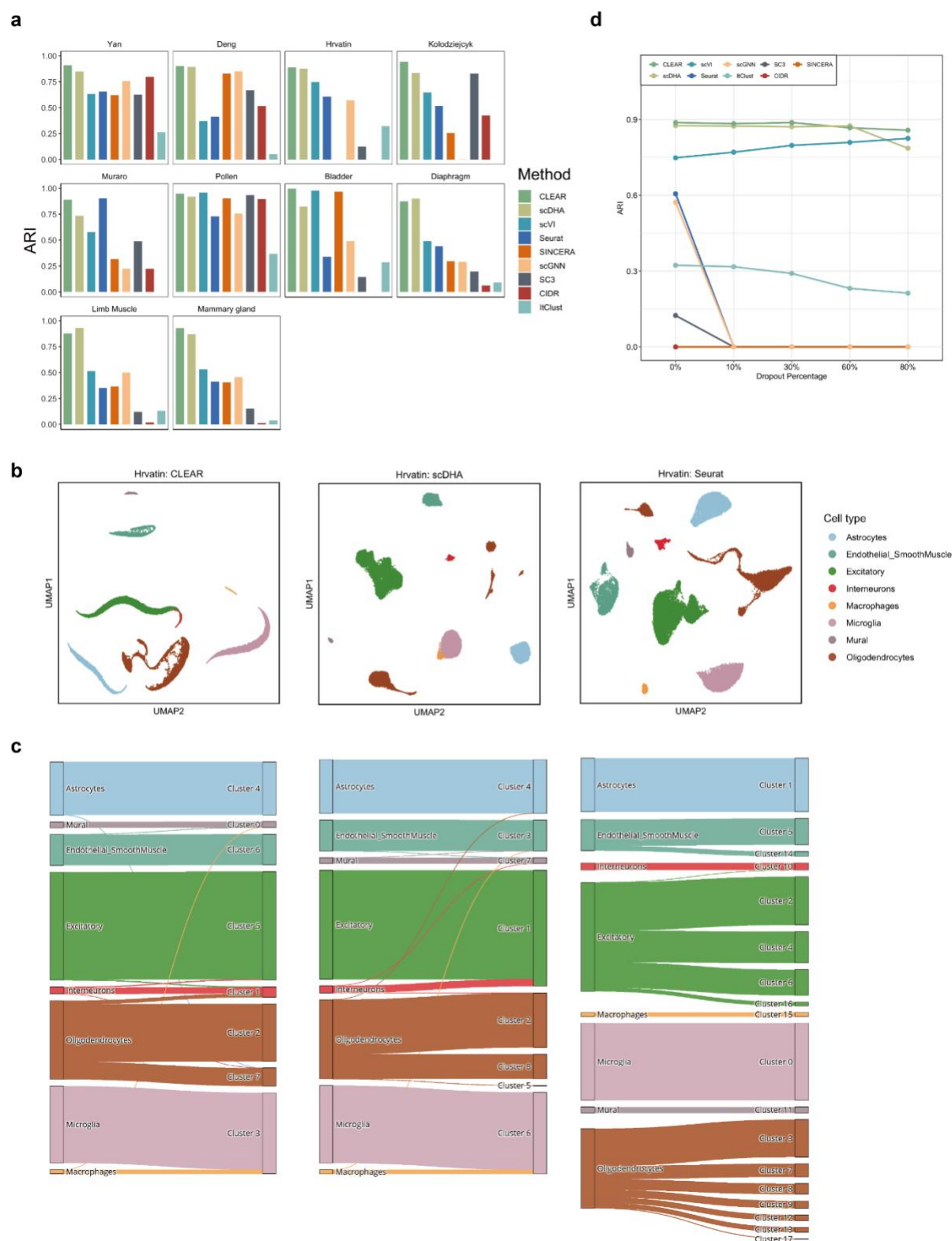
**Fig. 2 | The representation from CLEAR benefits clustering and dropout event correction. a** Clustering performance comparison of different methods on diverse datasets. On average, CLEAR improves over the second-best method, scDHA, by 4.56%, regarding ARI. **b** UMAP visualization of representations produced by CLEAR, scDHA, and Seruat on the Hrvatin dataset. **c** River plots of the Hrvatin dataset. CLEAR clustering matches almost perfectly with the expert annotation, without over-clustering or under-clustering. **d** Clustering performance change of different methods against different artificial dropout percentages in terms of ARI.

Fig. 3 | **CLEAR corrects batch effects effectively. a Upper panel:** UMAP visualization showing different methods' performance on integrating DoubleNeg and pDC cells from two batches. **Bottom panel:** UMAP visualization showing different methods' performance on separating four cell types. Notice that CLEAR's representations also preserve the biological similarity between CD141 cells and CD1C cells. **b** The quantitative performance of different methods on batch effect removal, measured by adjusted rand index(ARI) and average silhouette width (ASW).
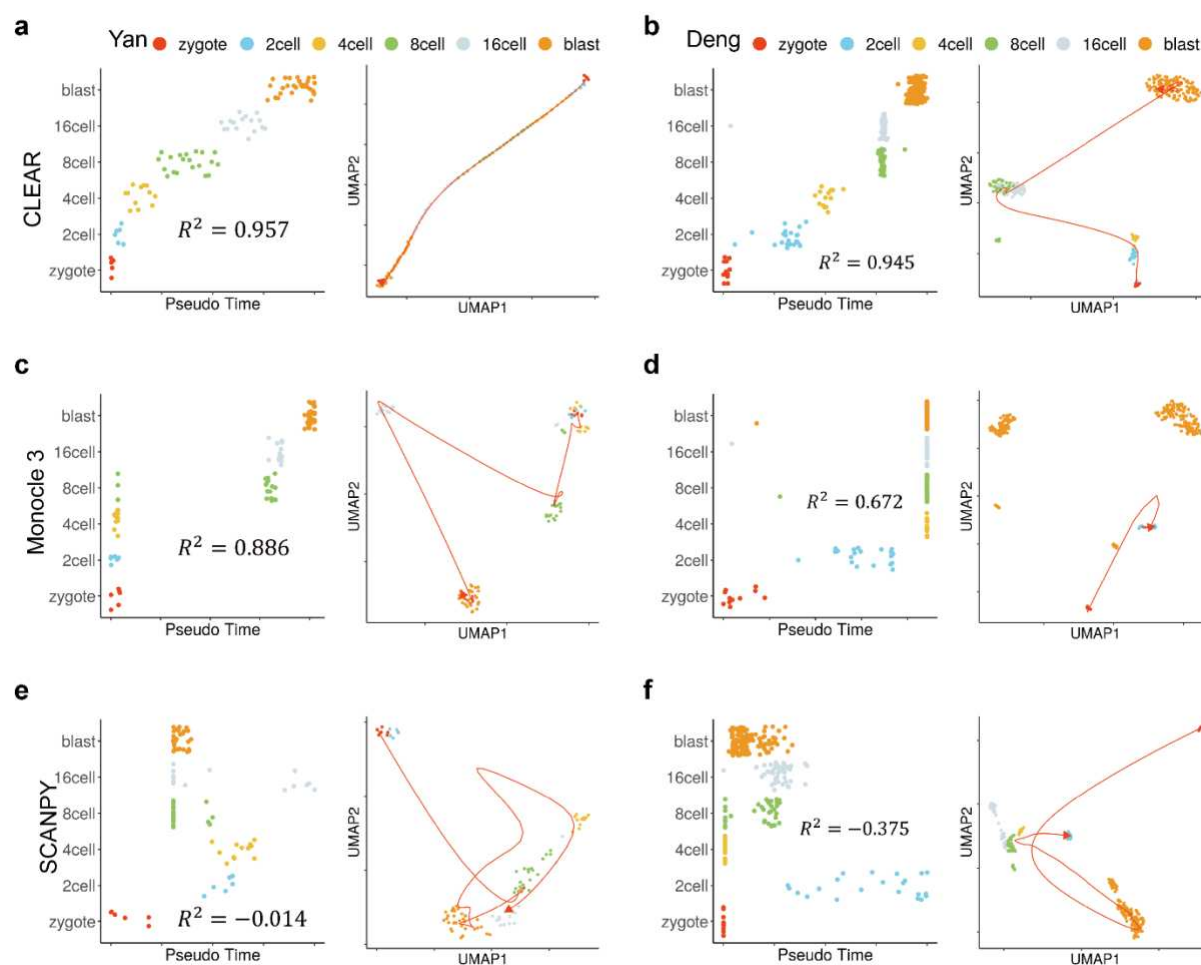
1
2 **Fig. 4 | CLEAR is helpful for pseudo-time inference. a** CLEAR's performance on the pseudo-time inference for
3 the Yan dataset. **Left figure:** Cells from the Yan dataset ordered by pseudo-time inferred from CLEAR. Ideally, the
4 points should fall on the diagonal. **Right figure:** UMAP visualization of time trajectory inferred from CLEAR. **b**
5 CLEAR's performance on the pseudo-time inference for the Deng dataset. **c,d** Monocle3's performance on the
6 pseudo-time inference for the Yan and Deng dataset. **e,f** SCANPY's performance on the pseudo-time inference
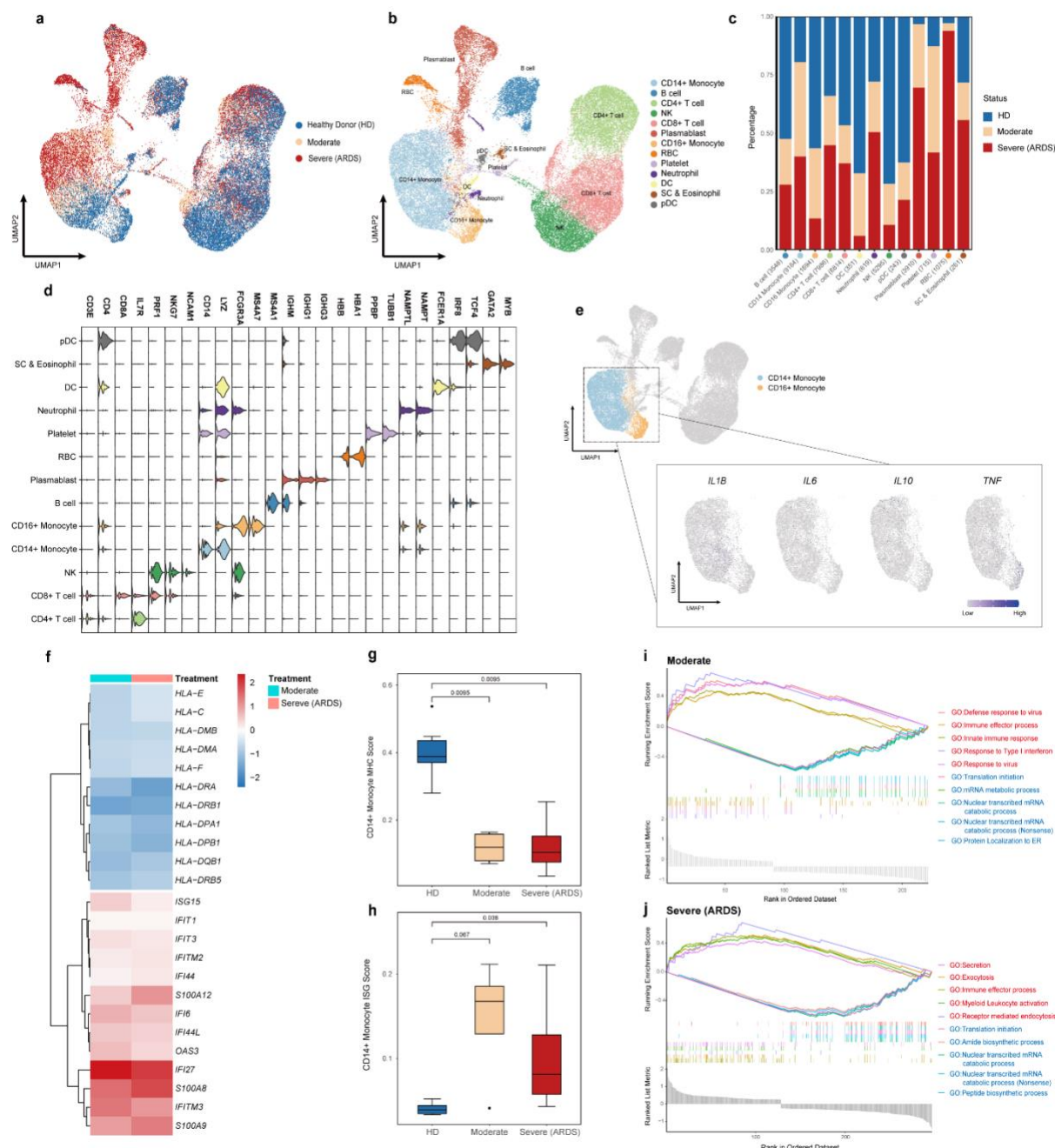7 for the Yan and Deng dataset.

1
2
3
4 **Fig. 5 | Peripheral immune cells atlas and inflammatory-related mechanisms in COVID-19 revealed by**
5 **CLEAR. a,b** UMAP visualization of the COVID-19 cell atlas **(a)** colored by COVID status and **(b)** colored by 13 cell
6 type clusters (n=43695 cells). **c** Bar plot showing the relative percentage of different cell types comparing three
7 COVID-19 statuses (HDs, Moderate status, and Severe status). **d** Stacked violin plot overview of the top-important
8 marker genes expression for each cell type. **e** UMAP visualization of the key pro-inflammatory cytokines expression
9 in both CD14+ and CD16+ monocytes. **f** Heatmap of IFN-stimulated genes and MHC-related genes in CD14+
10 monocyte. **g,h** Boxplots showing the mean **(g)** MHC-related score and **(h)** ISG score in CD14+ monocyte colored
11 by different COVID statuses (HDs--blue, Moderate--Oranger and Severe (ARDS)--Red). **i,j** Gene set enrichment
12 analysis (GSEA) of differential expressed gene (LogFC > 0.25) sets between **(i)** moderate CD14+ monocyte and
13 healthy donor CD14+ monocyte and **(j)** severe CD14+ monocyte and healthy donor CD14+ monocyte. Red
14 represents upregulated GO biological pathway, and blue represents downregulated GO biological pathway.
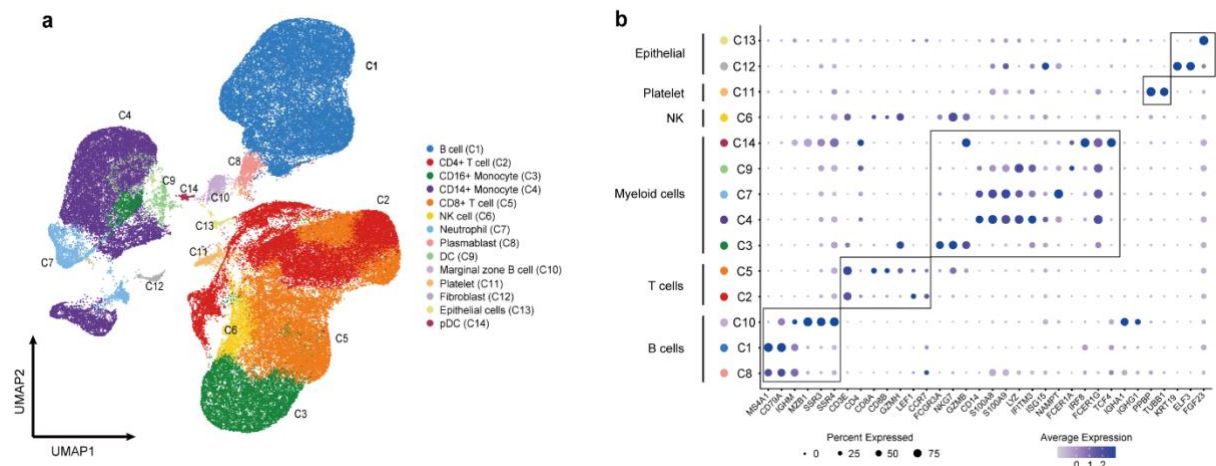
15

16

**Fig. 6 | COVID-19 PBMC cell atlas based on million-scale scRNA-seq dataset. a** UMAP embedding of peripheral blood mononuclear cells (PBMCs) from all samples (n=1.46 million cells) colored by manually-added cell types. **b** Dot plot showing percent expression and average expression of the selected marker genes for each cell type.