# Supplementary Figure 1

**Cells**

Genes

| | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 |
|---|---|---|---|---|---|---|---|---|---|---|
| | 2 | 3 | 3 | 0 | 2 | 0 | 0 | 0 | 0 | 0 |
| | 3 | 3 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 1 |
| | 0 | 1 | 0 | 0 | 0 | 5 | 0 | 6 | 3 | 6 |
| | 3 | 2 | 2 | 2 | 2 | 0 | 0 | 0 | 1 | 0 |

**Cell Clustering**

( C1, C2, C3 C4, C5 )     ( C6, C7, C8 C9, C10 )

**Imputation**

Genes

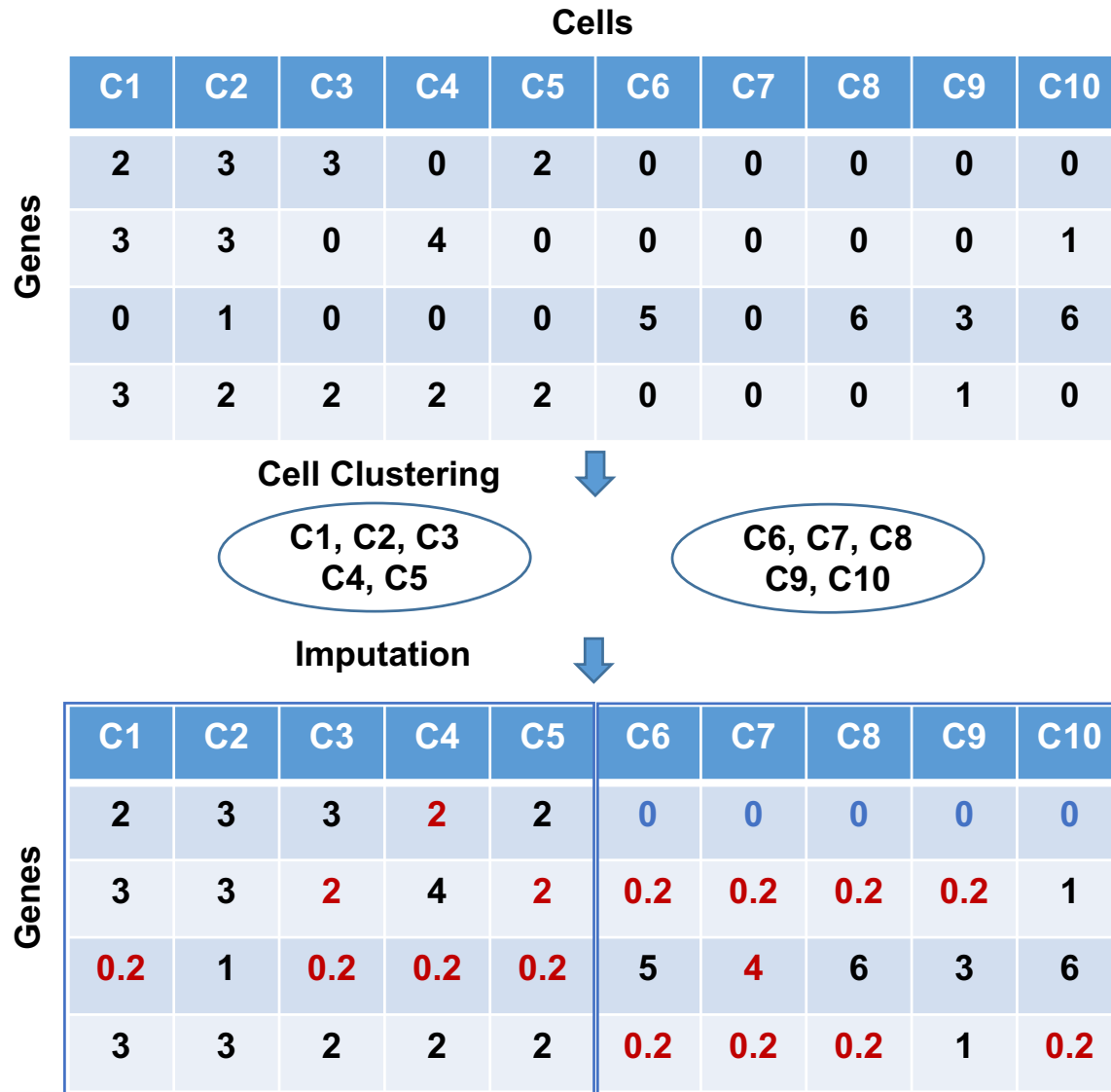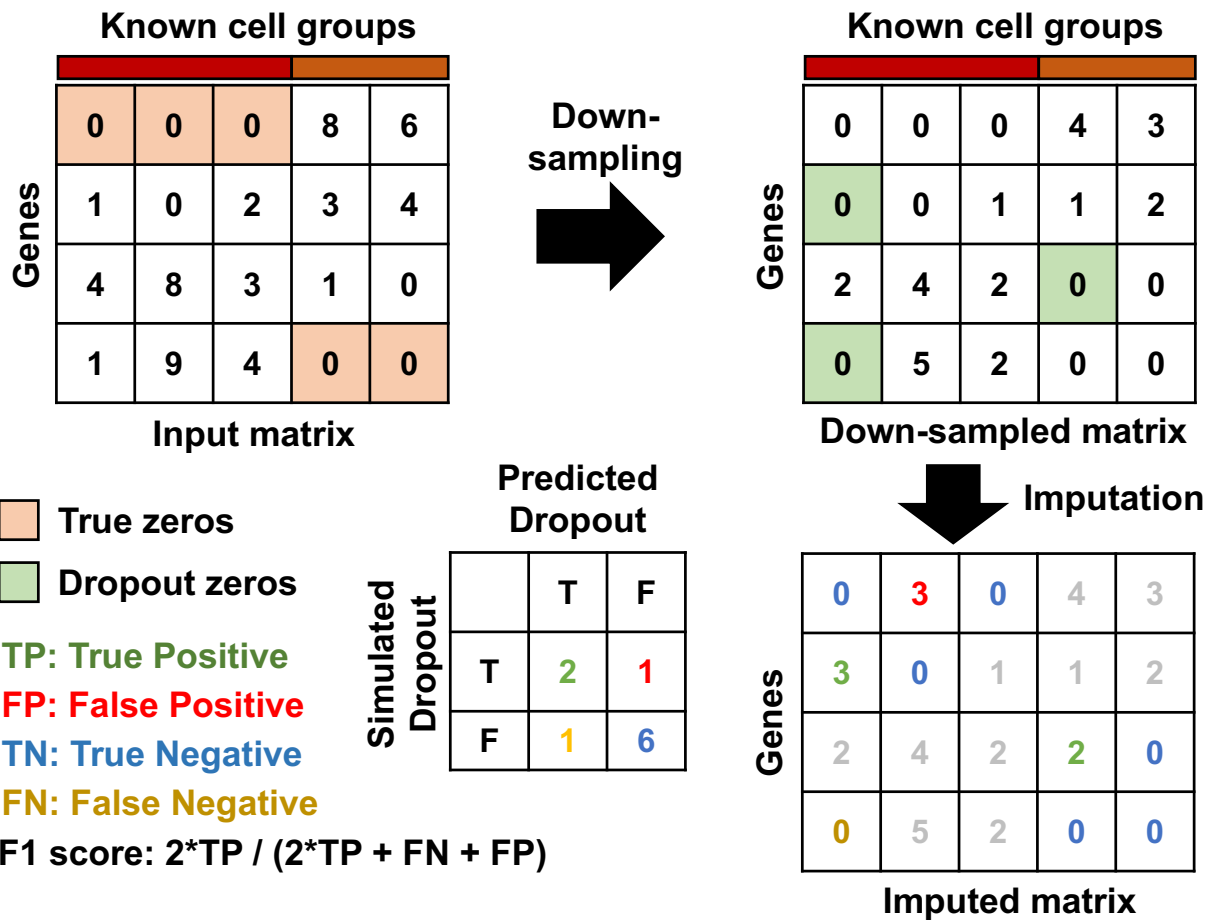| | C1 | C2 | C3 | C4 | C5 | C6 | C7 | C8 | C9 | C10 |
|---|---|---|---|---|---|---|---|---|---|---|
| | 2 | 3 | 3 | 2 | 2 | 0 | 0 | 0 | 0 | 0 |
| | 3 | 3 | 2 | 4 | 2 | 0.2 | 0.2 | 0.2 | 0.2 | 1 |
| | 0.2 | 1 | 0.2 | 0.2 | 0.2 | 5 | 4 | 6 | 3 | 6 |
| | 3 | 3 | 2 | 2 | 2 | 0.2 | 0.2 | 0.2 | 1 | 0.2 |

1    **Supplementary Figure 1. Basic procedure for clustering-based imputation.** Upper

2    matrix is a gene by cell matrix. After clustering on gene by cell matrix, we observe C1–

3    C5 as one cluster and C6–C10 as the other cluster. Imputation is performed by

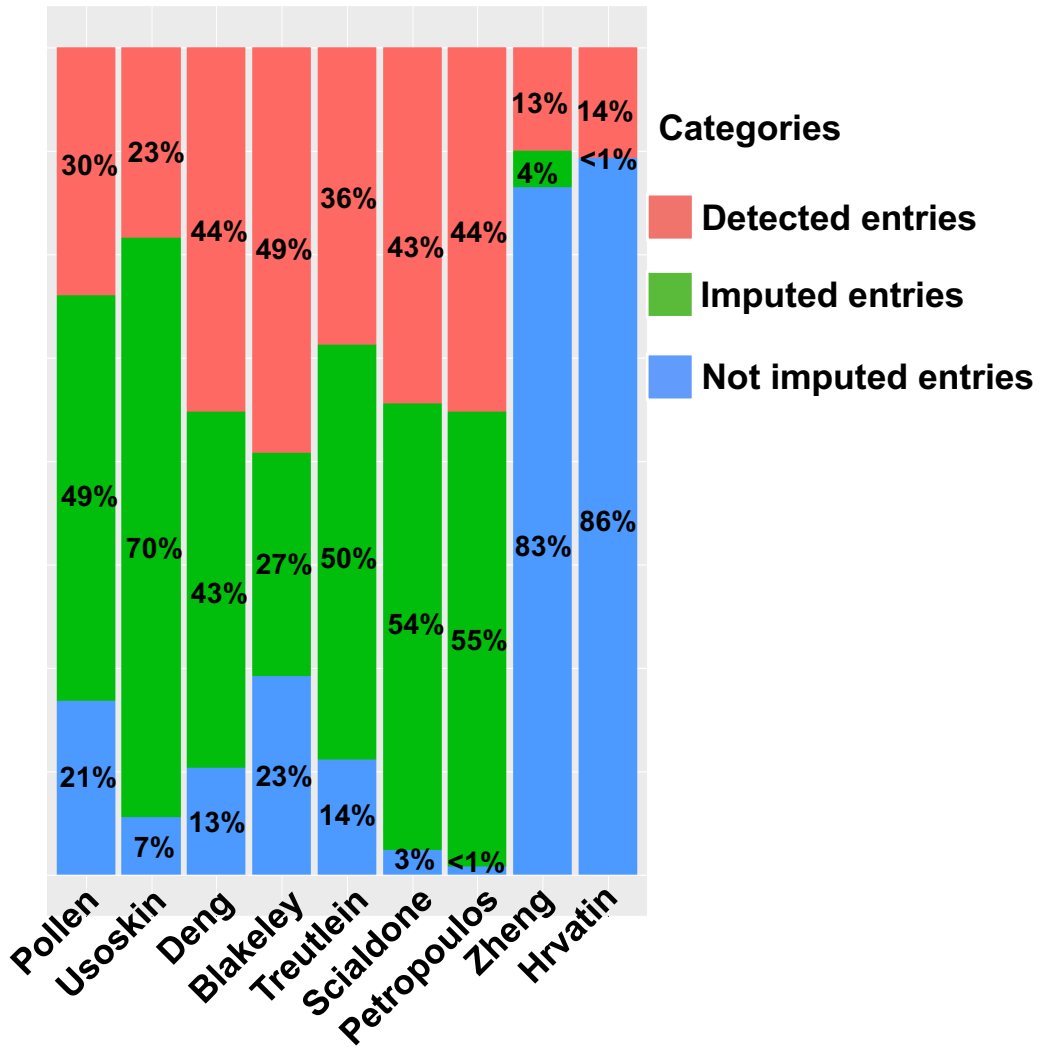4    averaging each cluster.

5

# Supplementary Figure 2

8    **Supplementary Figure 2. Overview of the down-sampling studies on**

9    **discriminating true zeros and dropout zeros.**  We defined the *true zeros* as the

10   genes where expression levels are consistently zero across all cells belonging to one

11   cell cluster.  To generate the *dropout zero*, we randomly down-sampled the raw

12   sequencing reads to a certain percent (e.g. 25%) of the total number of reads, mapped

13   the sampled reads onto the genome and computed the corresponding gene-cell read

14   count matrices.  We defined *dropout zero* as the genes where expression levels are

15   zero in the down-sampled datasets, but are positive in the full dataset.  The imputed

16   zero events could be therefore grouped into four situations: (1) true positive (TP,

17   imputed dropout zeros), (2) true negative (TN, non-imputed true zeros), (3) false

18   positive (FP, imputed true zeros) and (4) false negative (FN, non-imputed dropout

19   zeros).  The F1 score (the harmonic mean of precision and recall) was used to evaluate

20   the imputation performance of each method on down-sampled datasets.
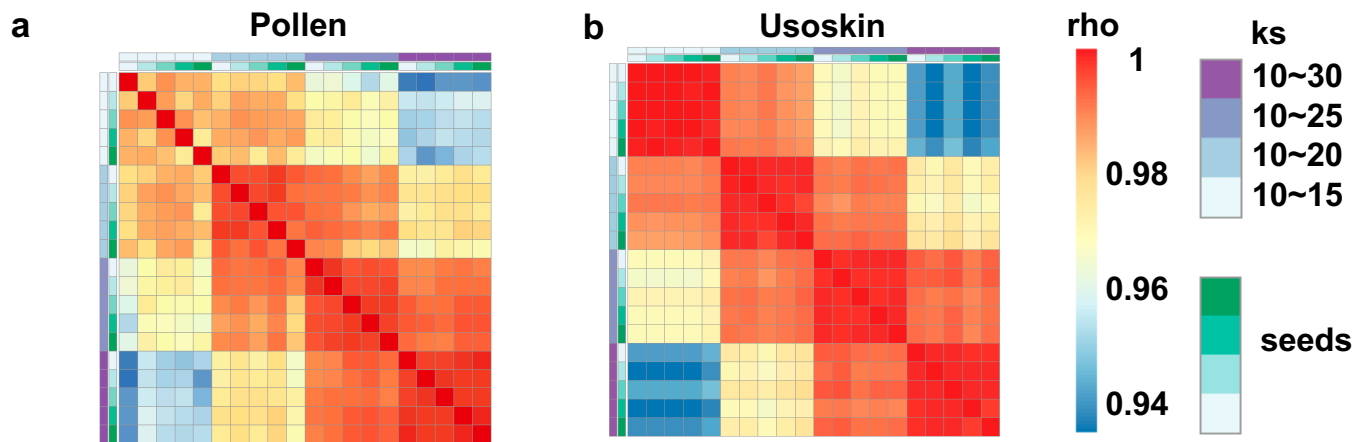
21

22

23

# Supplementary Figure 3

24    **Supplementary Figure 3. Overview of percent of detected, imputed and not**

25    **imputed entries in scRNA-seq datasets used in this study.** Percentage of detected

26    (input read count > 0), imputed (input read count is zero and the imputed read count is

27    positive), and not imputed entries (both input and imputed read count are zeros) for nine

28    different scRNA-seq datasets.  Genes that were expressed in less than 2 cells were

29    excluded before this analysis.

30

# Supplementary Figure 4

a         Pollen          b        Usoskin



rho

1

0.98

0.96

0.94

ks

10~30

10~25

10~20

10~15

seeds

31 **Supplementary Figure 4. DrImpute was robust on the different choices of the $k$**

32 **ranges and random seeds for the (a) Pollen and (b) Usoskin datasets.** The

33 robustness of imputation results were evaluated on different choices of number clusters:

34 $k = 10 - 15$ (default), $k = 10 - 20$, $k = 10 - 25$ and $k = 10 - 30$, as well as different

35 random number seeds (1 - 5) for k-means initialization. The robustness was

36 quantitatively measured as Pearson's correlation coefficient of imputed zero entries

37 between any two conditions (choices of $k$ ranges and random seeds). The color of the

38 heatmap indicates the Pearson's correlation coefficient.
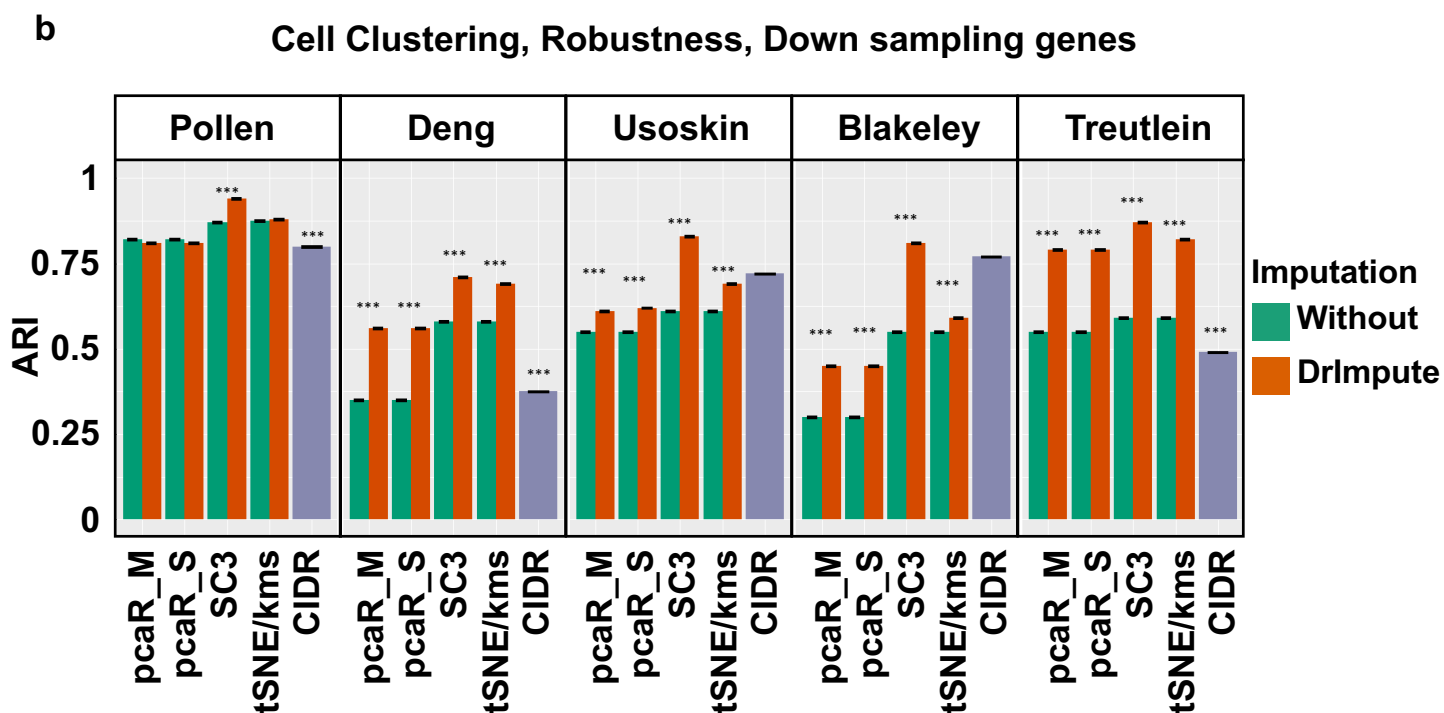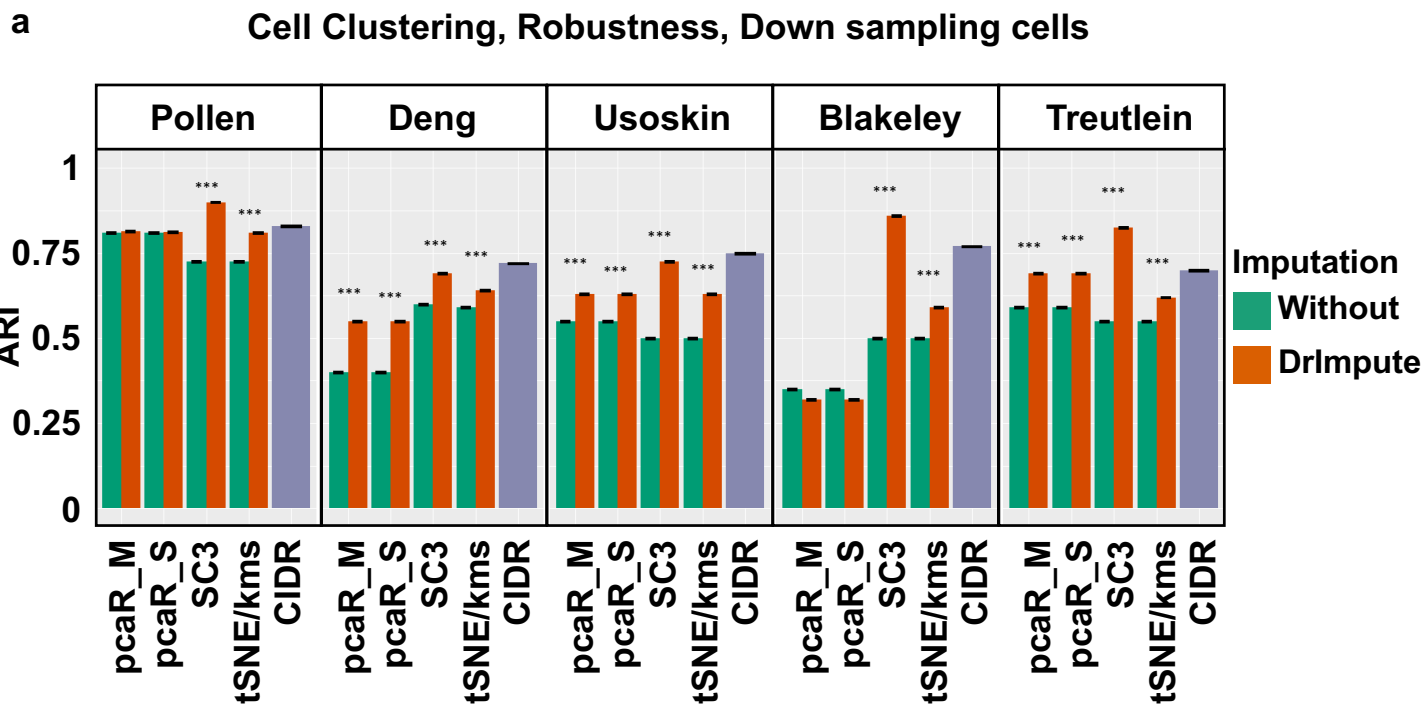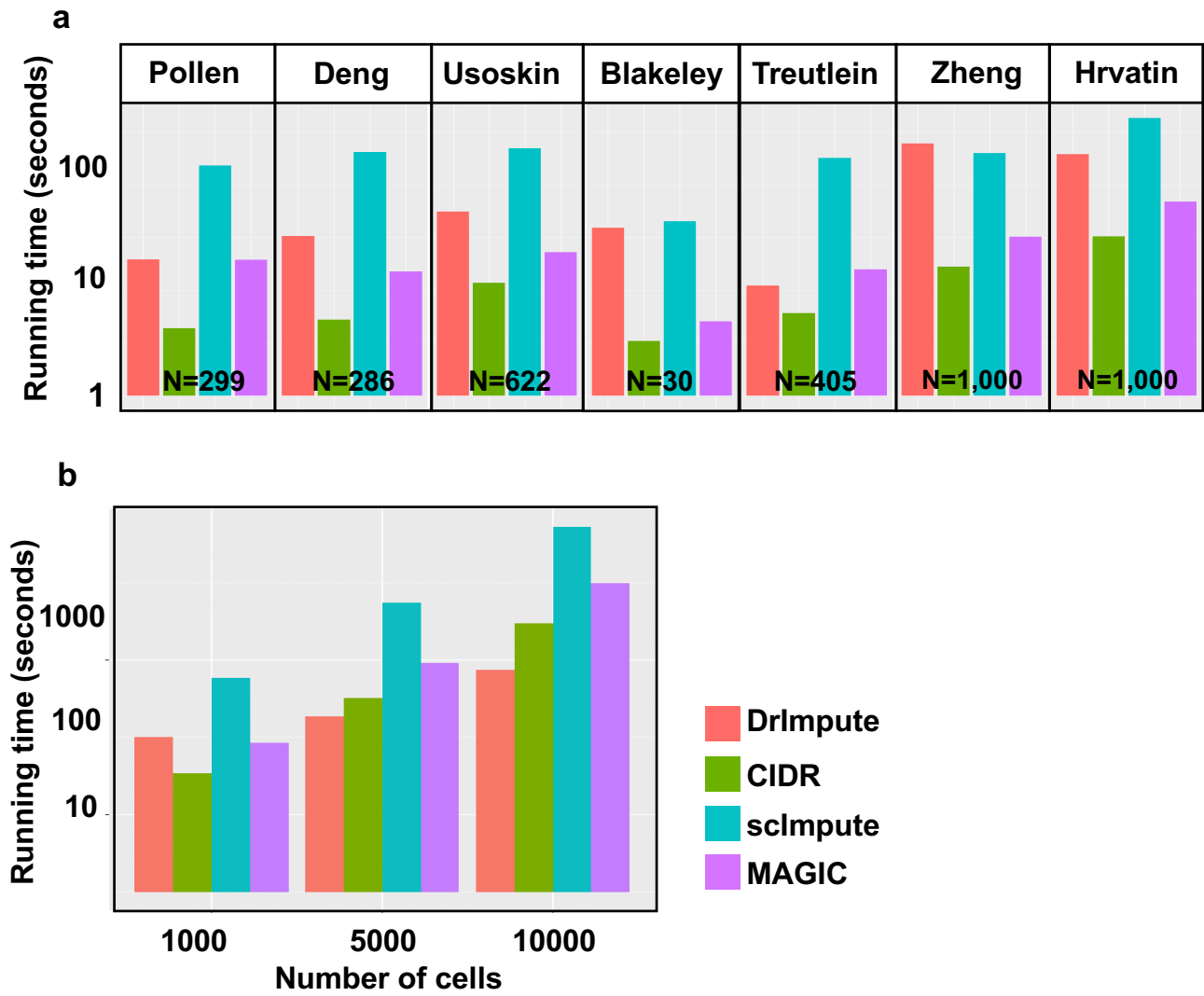
39

40

41

42

# Supplementary Figure 5

**a**



Cell Clustering, Robustness, Down sampling cells

**b**



Cell Clustering, Robustness, Down sampling genes

43    **Supplementary Figure 5. DrImpute significantly improved the performance of the**

44    **existing tools for cell type identification in robustness criteria.** To account for

45    robustness, original datasets were down-sampled by cells (a) or by genes (b); we

46    recorded clustering results for each data subset. ARIs are calculated for each pair of

47    data subsets. Barplot represents averaged ARIs. Blue interval represents one plus or

48    minus standard deviation of the data. Black interval represents one plus or minus

49    standard error of the data. Wilcoxon rank sum test is performed to compare before and

50    after imputation. For down-sampled cells, 16 out of 20 cases are improved. For down-

51    sampled genes, 18 out of 20 cases are improved ($***$ $p$ value $< 0.001$).

52

# Supplementary Figure 6

**Supplementary Figure 6. DrImpute is efficient on imputing large-scale scRNA-seq datasets. (a)** The running time of DrImpute, CIDR, scImpute and MAGIC on nine tested datasets are presented. The y-axis indicates the running time in seconds. **(b)** The running time of DrImpute, CIDR, scImpute and MAGIC on randomly sampled 1,000, 5,000 and 10,000 cells from the Zheng dataset are presented. All the analysis were performed on Intel Xeon 2.4GHz CPU. For both DrImpute and scImpute, 4 CPU cores were used for the analysis.