## Supplemental information

## Generative pretraining from large-scale

## transcriptomes for single-cell deciphering

**Hongru Shen, Jilei Liu, Jiani Hu, Xilin Shen, Chao Zhang, Dan Wu, Mengyao Feng, Meng Yang, Yang Li, Yichen Yang, Wei Wang, Qiang Zhang, Jilong Yang, Kexin Chen, and Xiangchun Li**

## Supplementary Materials



**A** (*HCA*, **MNN**, n = 282,558, *kBET* = 0.88)  **B** (*HCA*, **Combat**, n = 282,558, *kBET* = 0.78)  **C** (*HCA*, **Harmony**, n = 282,558, *kBET* = 0.88)  **D** (*HCA*, **Seurat**, n = 282,558, *kBET* = 0.86)

**E** (*HCA*, **Pegasus**, n = 282,558, *kBET* = 0.59)  **F** (*HCA*, **Scaronoma**, n = 282,558, *kBET* = 0.82)  **G** (*HCA*, **DESC**, n = 282,558, *kBET* = 0.94)  **H** (*HCA*, **iMAP**, n = 282,558, *kBET* = 0.93)

**I** (*HCA*, **ScVI**, n = 282,558, *kBET* = 0.85)  **J** (*HCA*, **scArches**, n = 282,558, *kBET* = 0.84)  **K** (*HCA*, **BBKNN**, n = 282,558, *kBET* = 0.86)  **L** *kBET* for different methods

Legend:
- CD14+ Monocytes
- CD16+ Monocytes
- CD4+ Naive T cells
- cDCs
- Cytotoxic T cells
- Erythroid cells
- HSCs
- Megakaryocytes
- Memory B cells
- MSCs
- Naive B cells
- Naive CD8+ T cells
- NK cells
- pDCs
- Plasma cells
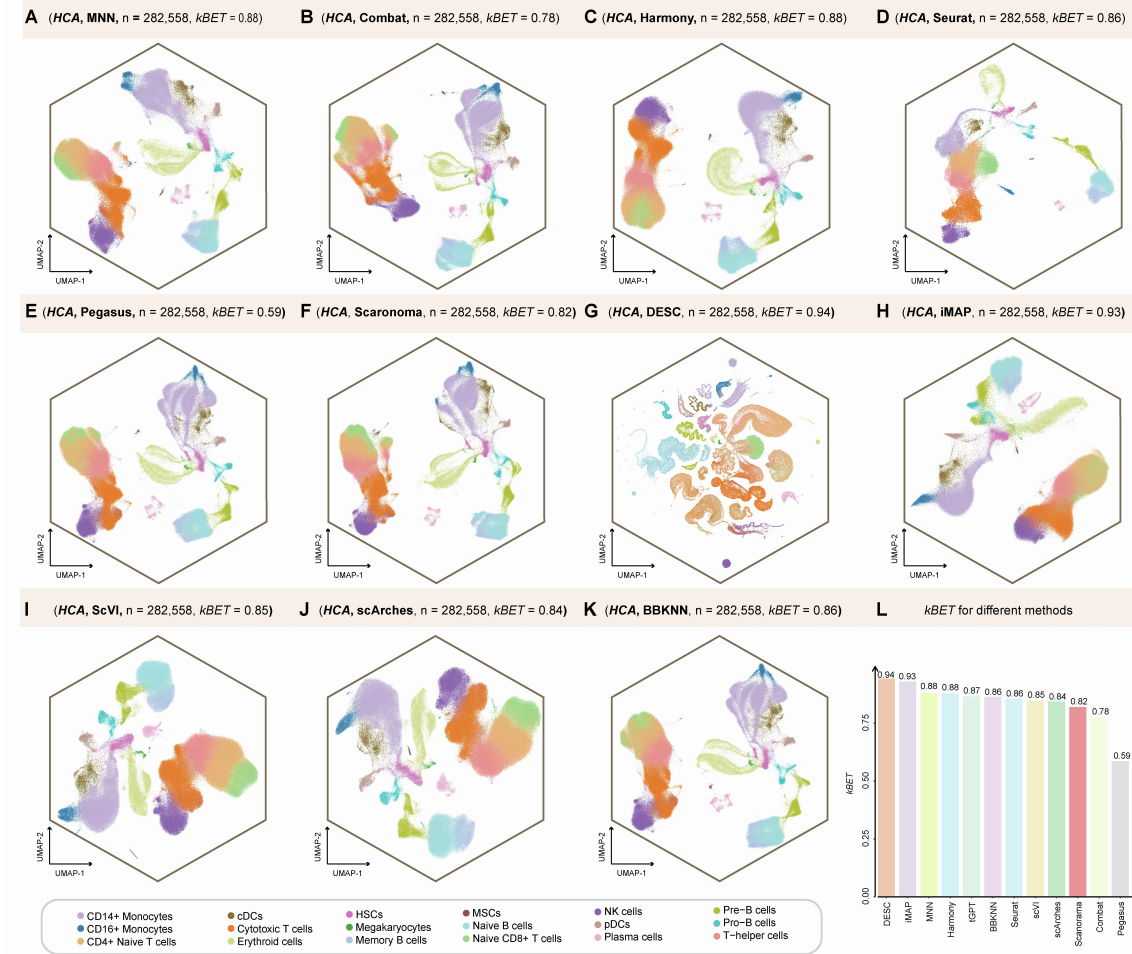- Pre−B cells
- Pro−B cells
- T−helper cells

**Figure S1. The UMAP visualization plots of different batch-correction methods on the HCA dataset (A to K) and *kBET* acceptance rate (L), related to Figure 2.**
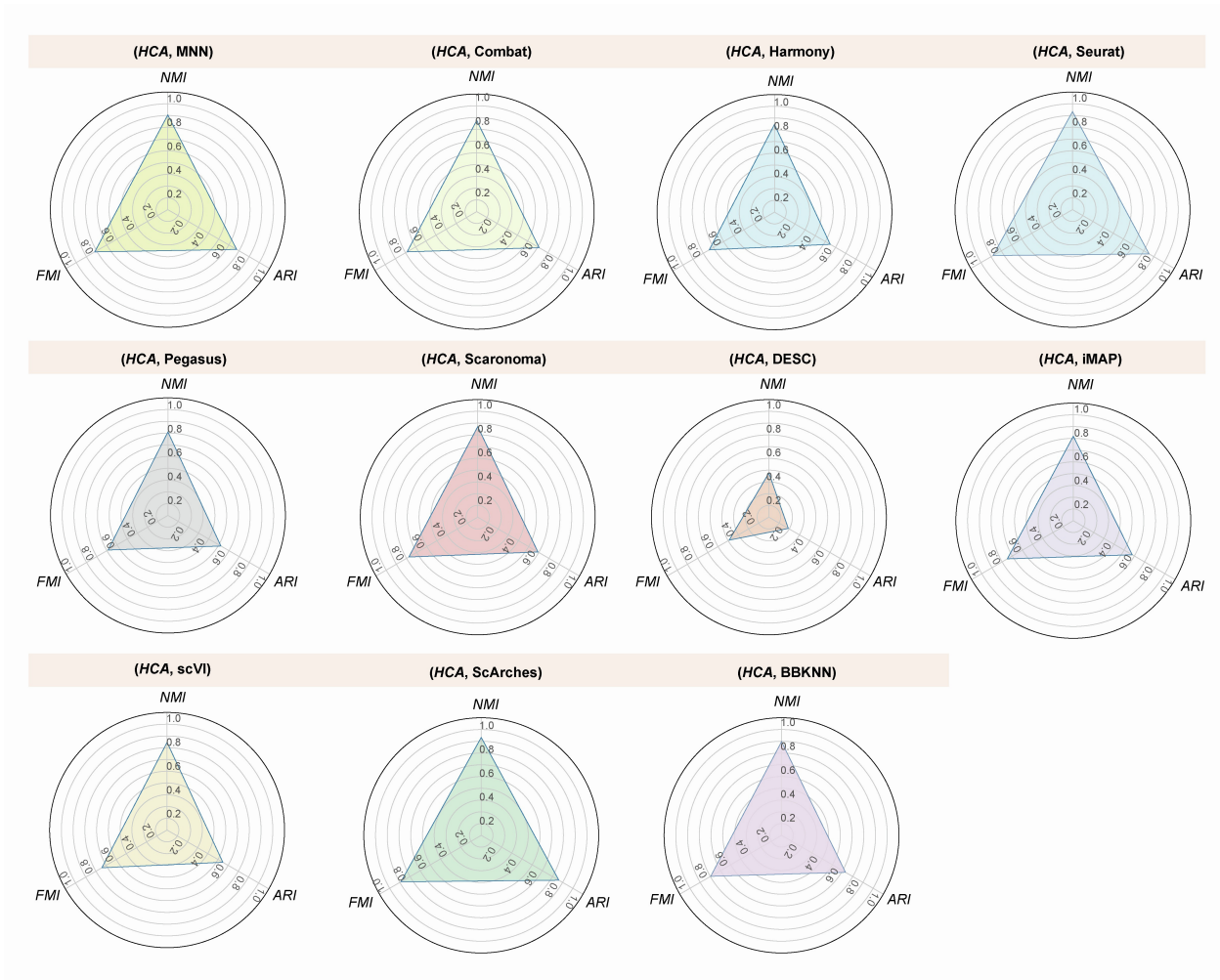
**Figure S 2 . Radar charts illustrating the best clustering metrics for different batch-correction methods obtained from grid search on the *HCA* dataset, related to Figure 2.** *ARI*, Adjusted Rand Index; *NMI*, Normalized Mutual information; *FMI*, Fowlkes-Mallows Index.
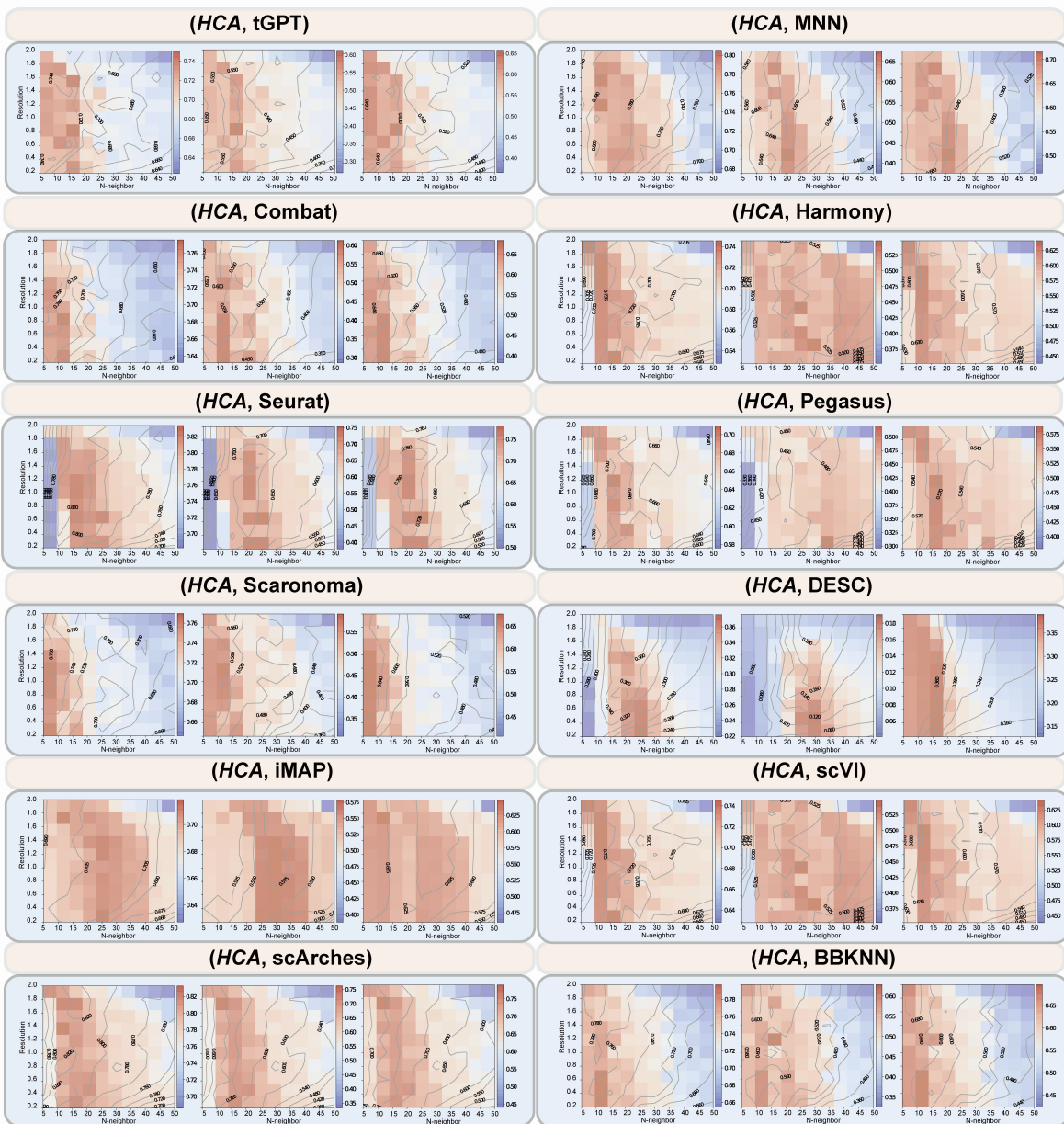
**Figure S3. The clustering performance with grid search for resolution and number of neighbors for different batch-correction methods on the *HCA* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.
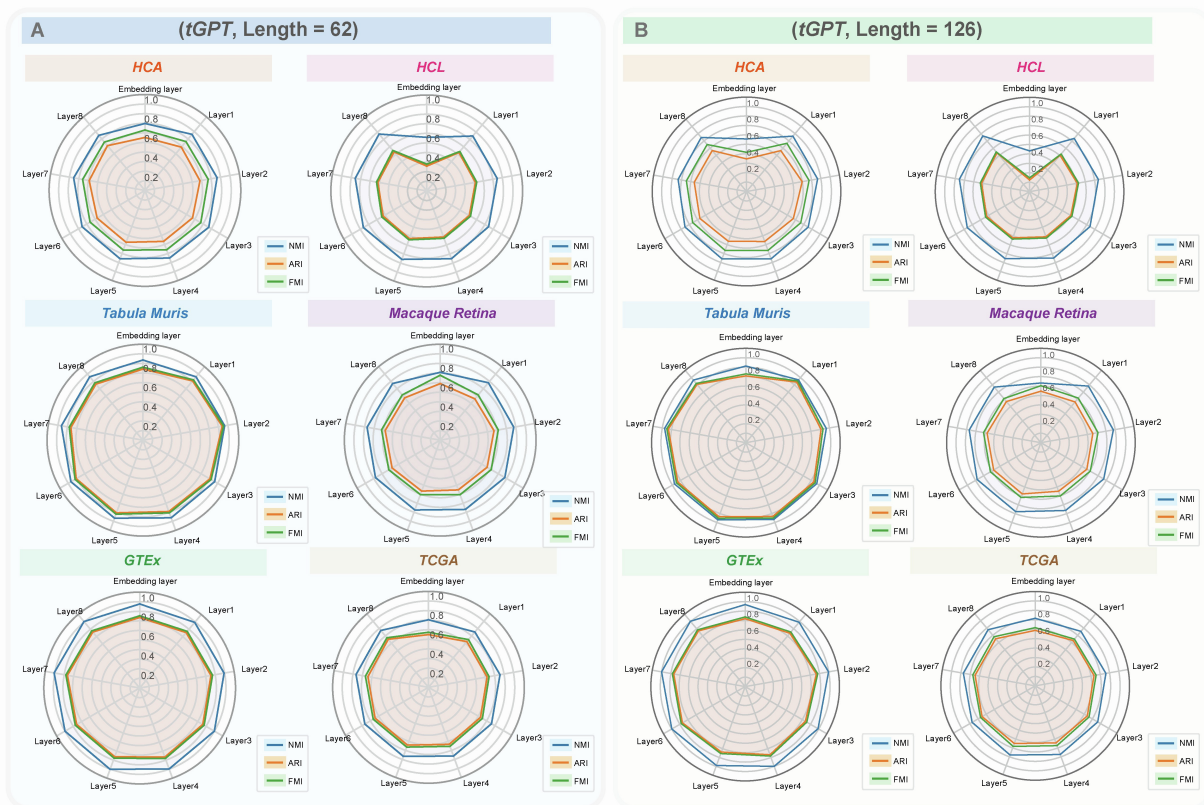
**Figure S4. Radar charts illustrating the clustering performance achieved by feature representations extracted from different layers of *tGPT* for the top 62 (A) and 126 (B) expressing genes on *HCA*, *HCL*, *Tabula Muris*, *Macaque Retina*, *GTEx*, and *TCGA* datasets, related to Figure 2.**
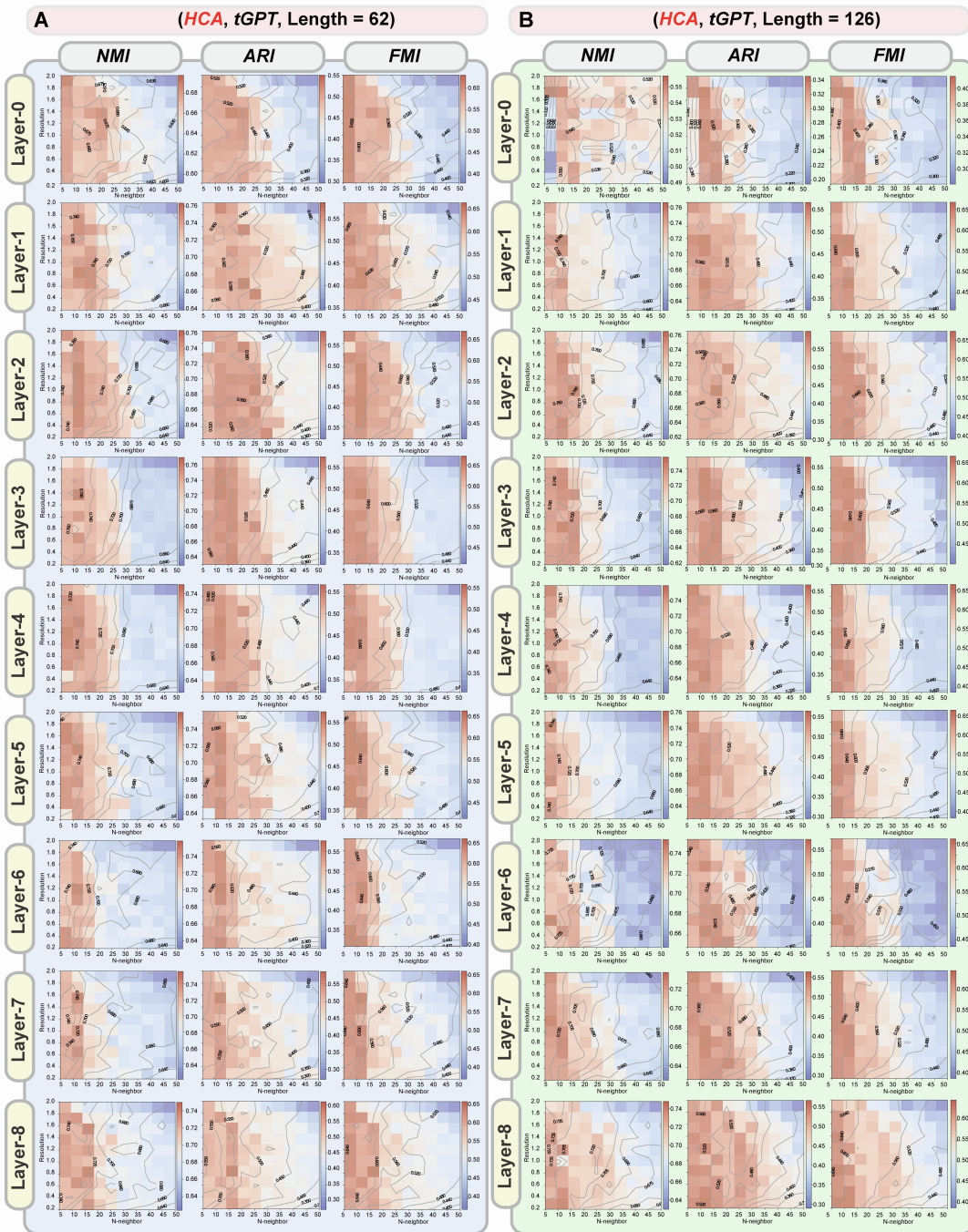
**Figure S 5. The clustering performance with grid search for resolution and number of neighbors for the top 62 (A) and 126 (B) expressing genes among feature representations extracted from different layers on the *HCA* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.

**Figure S 6. The clustering performance with grid search for resolution and number of neighbors for the top 62 (A) and 126 (B) expressing genes among feature representations extracted from different layers on the *HCL* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.
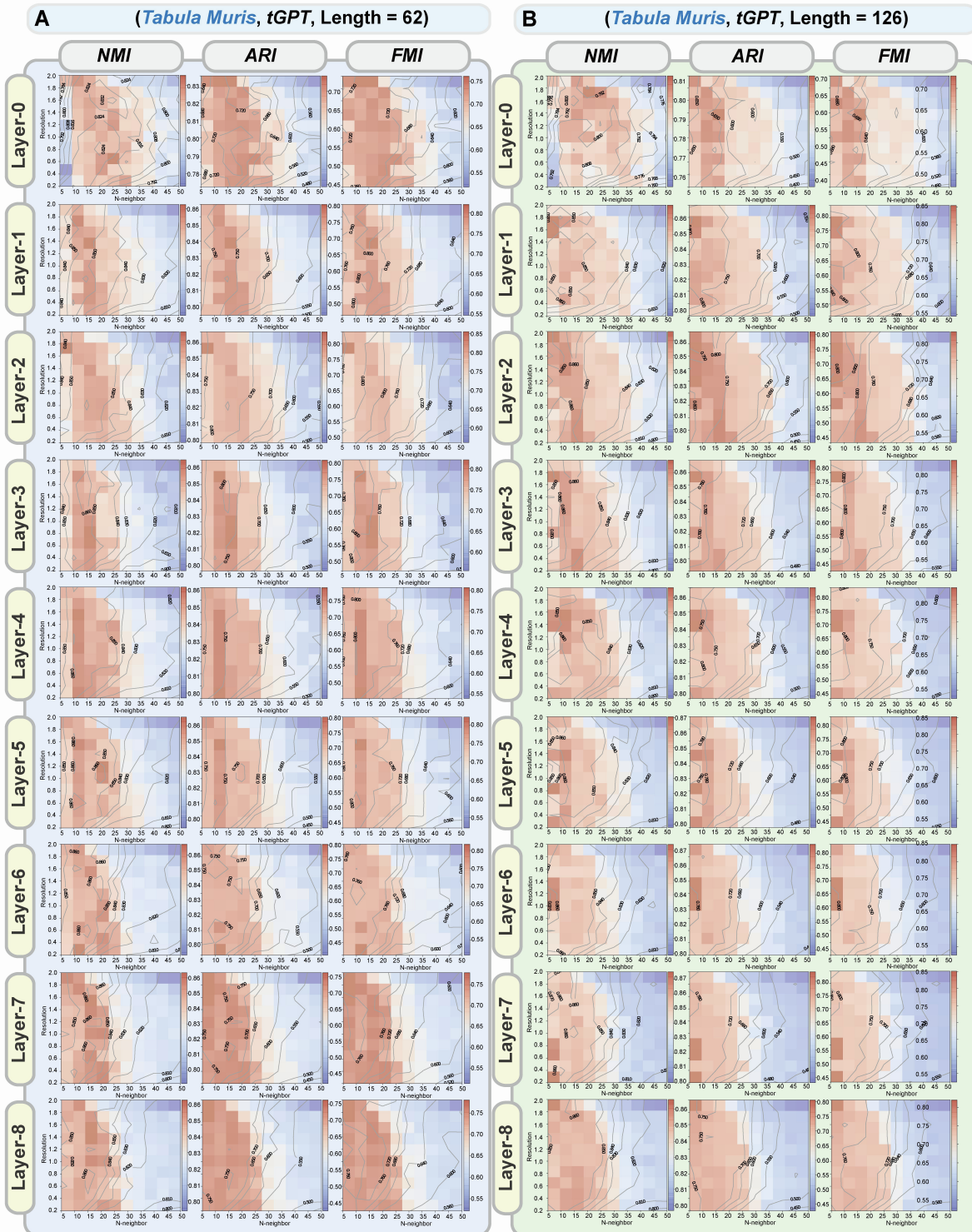
**Figure S 7. The clustering performance with grid search for resolution and number of neighbors for the top 62 (A) and 126 (B) expressing genes among feature representations extracted from different layers on the *Tabula Muris* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.

**Figure S 8. The clustering performance with grid search for resolution and number of neighbors for the top 62 (A) and 126 (B) expressing genes among feature representations extracted from different layers on the *Macaque Retina* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.
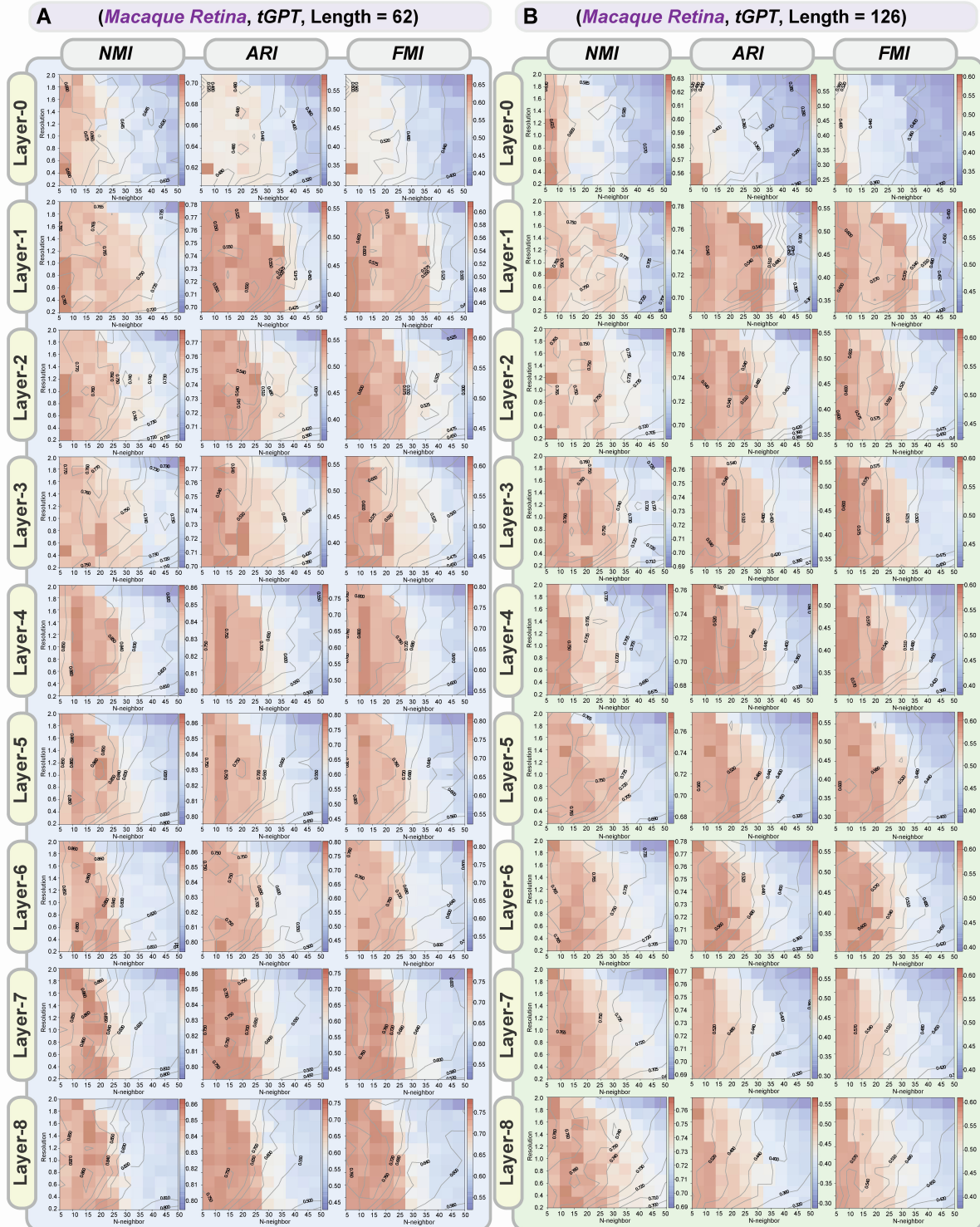
**Figure S9. The clustering performance with grid search for resolution and number of neighbors for the top 62 (A) and 126 (B) expressing genes among feature representations extracted from different layers on the *GTEx* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.

**Figure S 10. The clustering performance with grid search for resolution and number of neighbors for the top 62 (A) and 126 (B) expressing genes among feature representations extracted from different layers on the *TCGA* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.
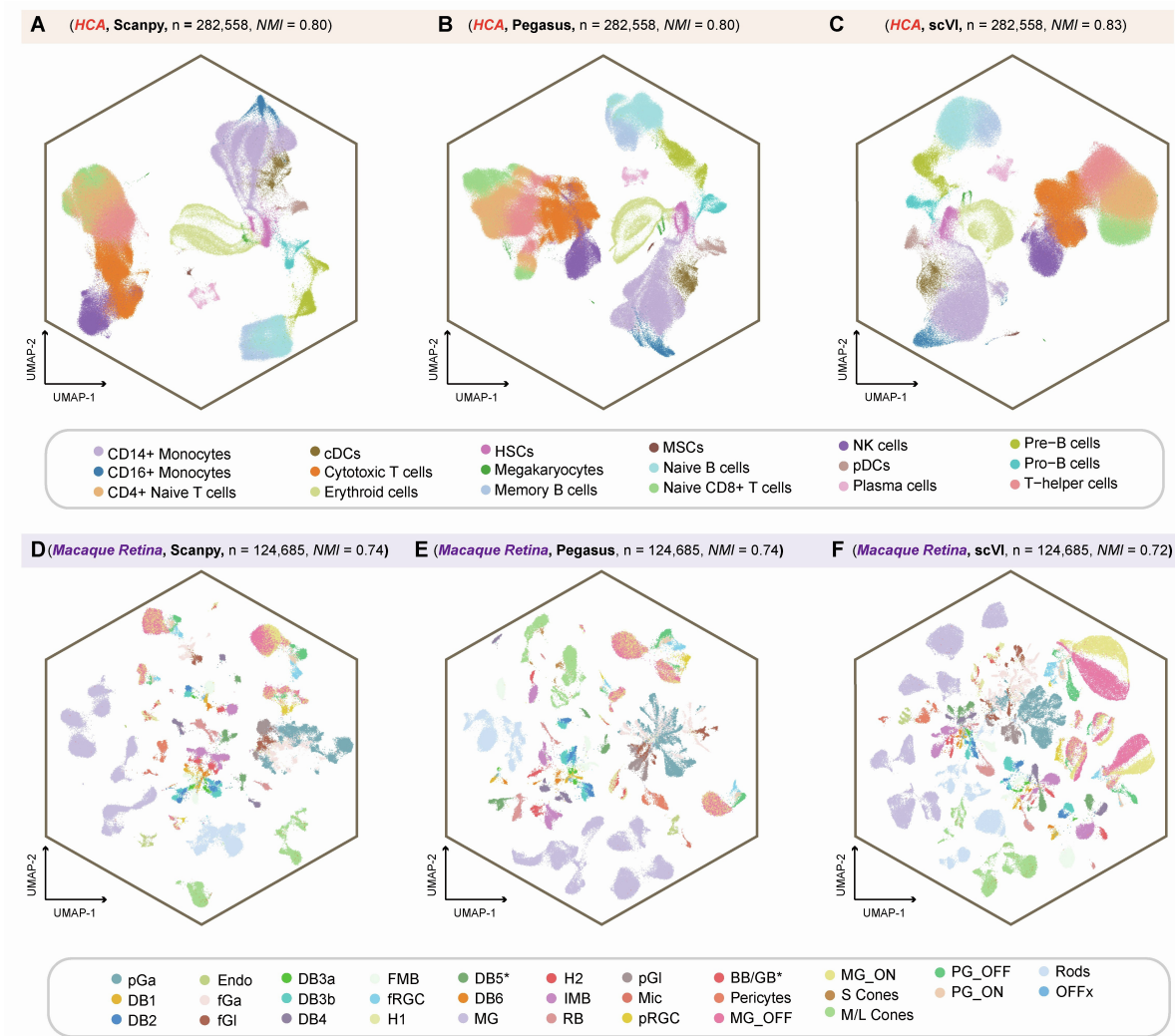
**Figure S11. UMAP visualization on different datasets obtained from different methods, related to Figure 2.**
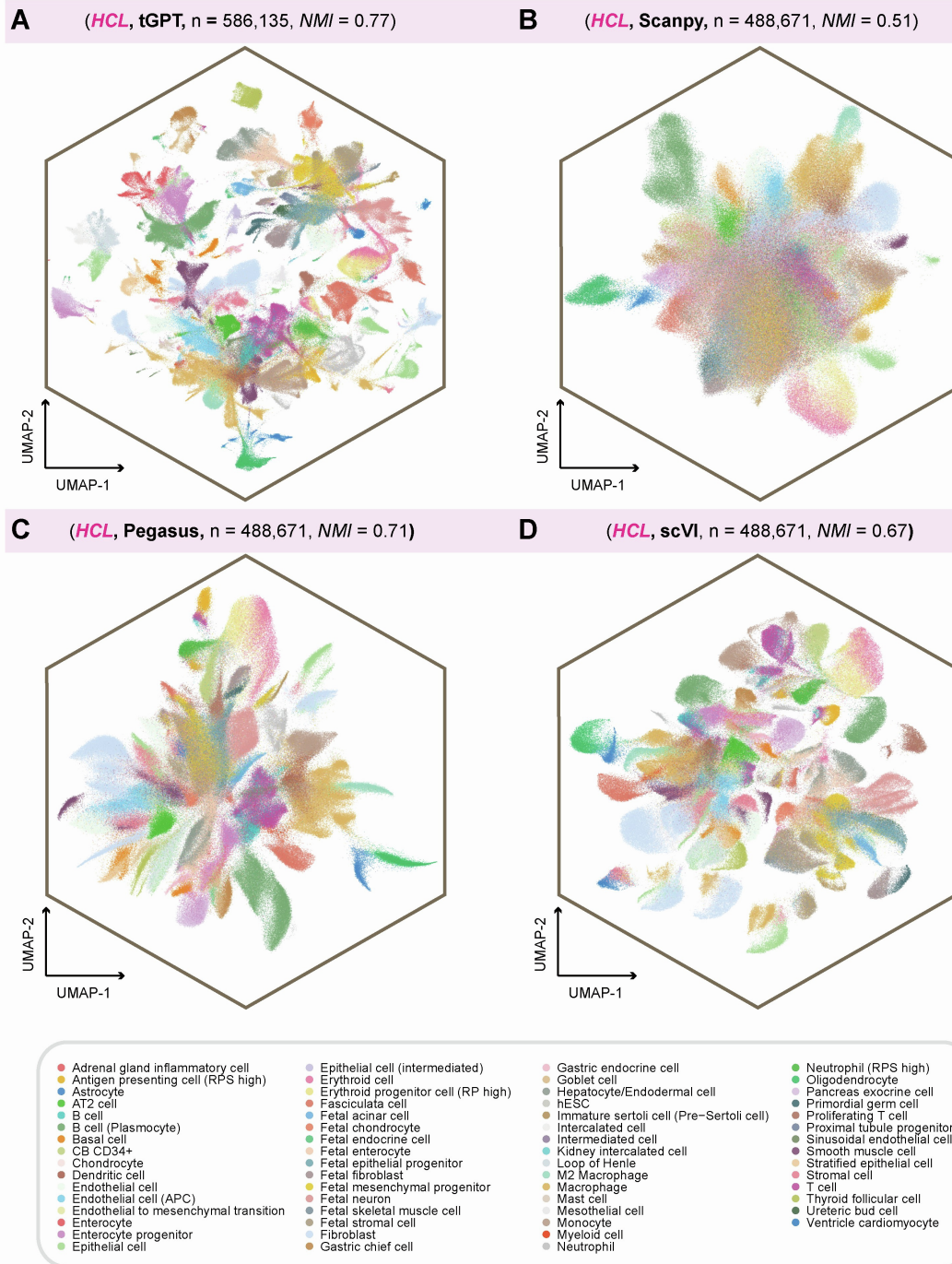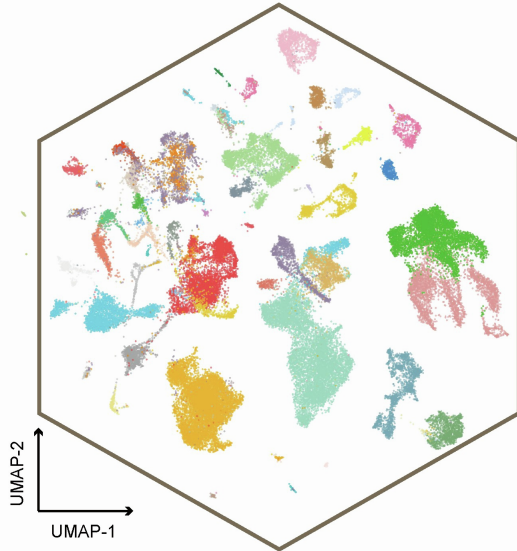
**A** (*HCL*, **tGPT,** n = 586,135, *NMI* = 0.77)

**B** (*HCL*, **Scanpy,** n = 488,671, *NMI* = 0.51)

**C** (*HCL*, **Pegasus,** n = 488,671, *NMI* = 0.71)

**D** (*HCL*, **scVI**, n = 488,671, *NMI* = 0.67)

Legend:
- Adrenal gland inflammatory cell
- Antigen presenting cell (RPS high)
- Astrocyte
- AT2 cell
- B cell
- B cell (Plasmocyte)
- Basal cell
- CB CD34+
- Chondrocyte
- Dendritic cell
- Endothelial cell
- Endothelial cell (APC)
- Endothelial to mesenchymal transition
- Enterocyte
- Enterocyte progenitor
- Epithelial cell
- Epithelial cell (intermediated)
- Erythroid cell
- Erythroid progenitor cell (RP high)
- Fasciculata cell
- Fetal acinar cell
- Fetal chondrocyte
- Fetal endocrine cell
- Fetal enterocyte
- Fetal epithelial progenitor
- Fetal fibroblast
- Fetal mesenchymal progenitor
- Fetal neuron
- Fetal skeletal muscle cell
- Fetal stromal cell
- Fibroblast
- Gastric chief cell
- Gastric endocrine cell
- Goblet cell
- Hepatocyte/Endodermal cell
- hESC
- Immature sertoli cell (Pre−Sertoli cell)
- Intercalated cell
- Intermediated cell
- Kidney intercalated cell
- Loop of Henle
- M2 Macrophage
- Macrophage
- Mast cell
- Mesothelial cell
- Monocyte
- Myeloid cell
- Neutrophil
- Neutrophil (RPS high)
- Oligodendrocyte
- Pancreas exocrine cell
- Primordial germ cell
- Proliferating T cell
- Proximal tubule progenitor
- Sinusoidal endothelial cell
- Smooth muscle cell
- Stratified epithelial cell
- Stromal cell
- T cell
- Thyroid follicular cell
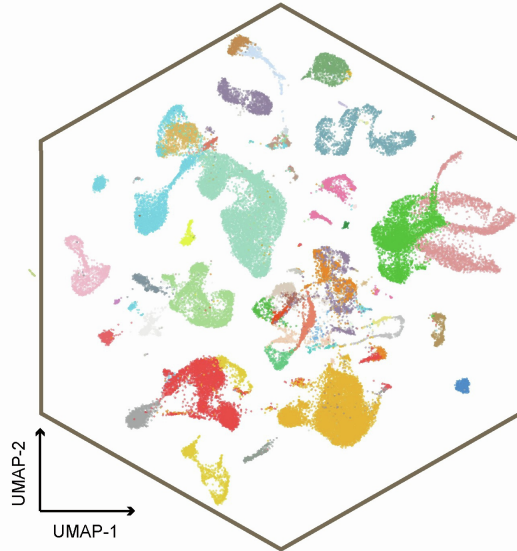- Ureteric bud cell
- Ventricle cardiomyocyte

**Figure S12. The full annotation of UMAP visualization of different methods on the *HCL* dataset, related to Figure 2.** The *NMI* metric and annotation of cells are shown.
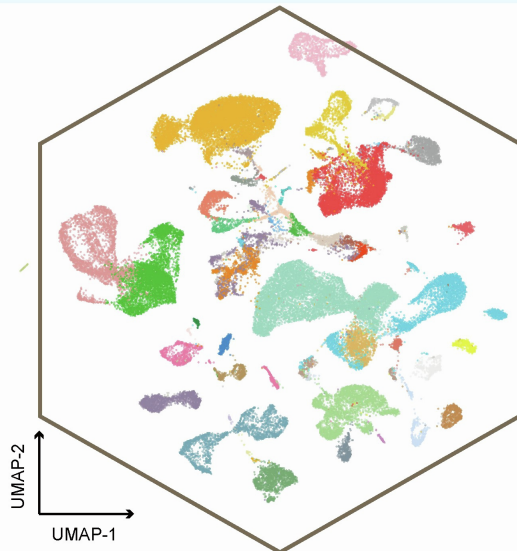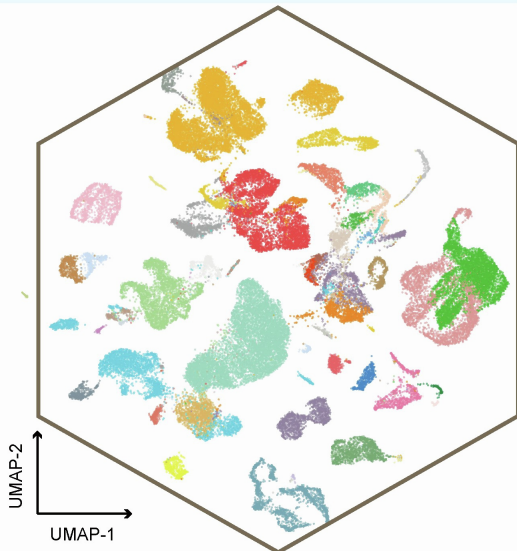
**A** (*Tabula Muris*, **tGPT**, n = 54,862, *NMI* = 0.87)

**B** (*Tabula Muris*, **Scanpy**, n = 54,779, *NMI* = 0.88)

**C** (*Tabula Muris*, **Pegasus**, n = 54,779, *NMI* = 0.86)

**D** (*Tabula Muris*, **scVI**, n = 54,779, *NMI* = 0.86)

- alveolar macrophage
- B cell
- basal cell
- basal cell of epidermis
- basophil
- bladder cell
- bladder urothelial cell
- blood cell
- cardiac muscle cell
- ciliated columnar cell of tracheobronchial tree
- classical monocyte
- dendritic cell
- DN1 thymic pro−T cell
- duct epithelial cell
- early pro−B cell
- endocardial cell
- endothelial cell
- endothelial cell of hepatic sinusoid
- epithelial cell
- erythroblast
- fibroblast
- Fraction A pre−pro B cell
- granulocyte
- granulocytopoietic cell
- hematopoietic precursor cell
- hepatocyte
- immature B cell
- immature T cell
- keratinocyte
- kidney capillary endothelial cell
- kidney cell
- kidney collecting duct epithelial cell
- kidney loop of Henle ascending limb epithelial cell
- kidney proximal straight tubule epithelial cell
- Langerhans cell
- late pro−B cell
- leukocyte
- luminal epithelial cell of mammary gland
- lung endothelial cell
- macrophage
- mast cell
- mesangial cell
- mesenchymal cell
- mesenchymal stem cell
- monocyte
- myeloid cell
- natural killer cell
- neuroendocrine cell
- non−classical monocyte
- proerythroblast
- promonocyte
- skeletal muscle satellite cell
- stromal cell
- T cell
- type II pneumocyte

**Figure S13. The full annotation of UMAP visualization of different methods on the** ***Tabula Muris* dataset, related to Figure 2.** The *NMI* metric and annotation of cells are shown.
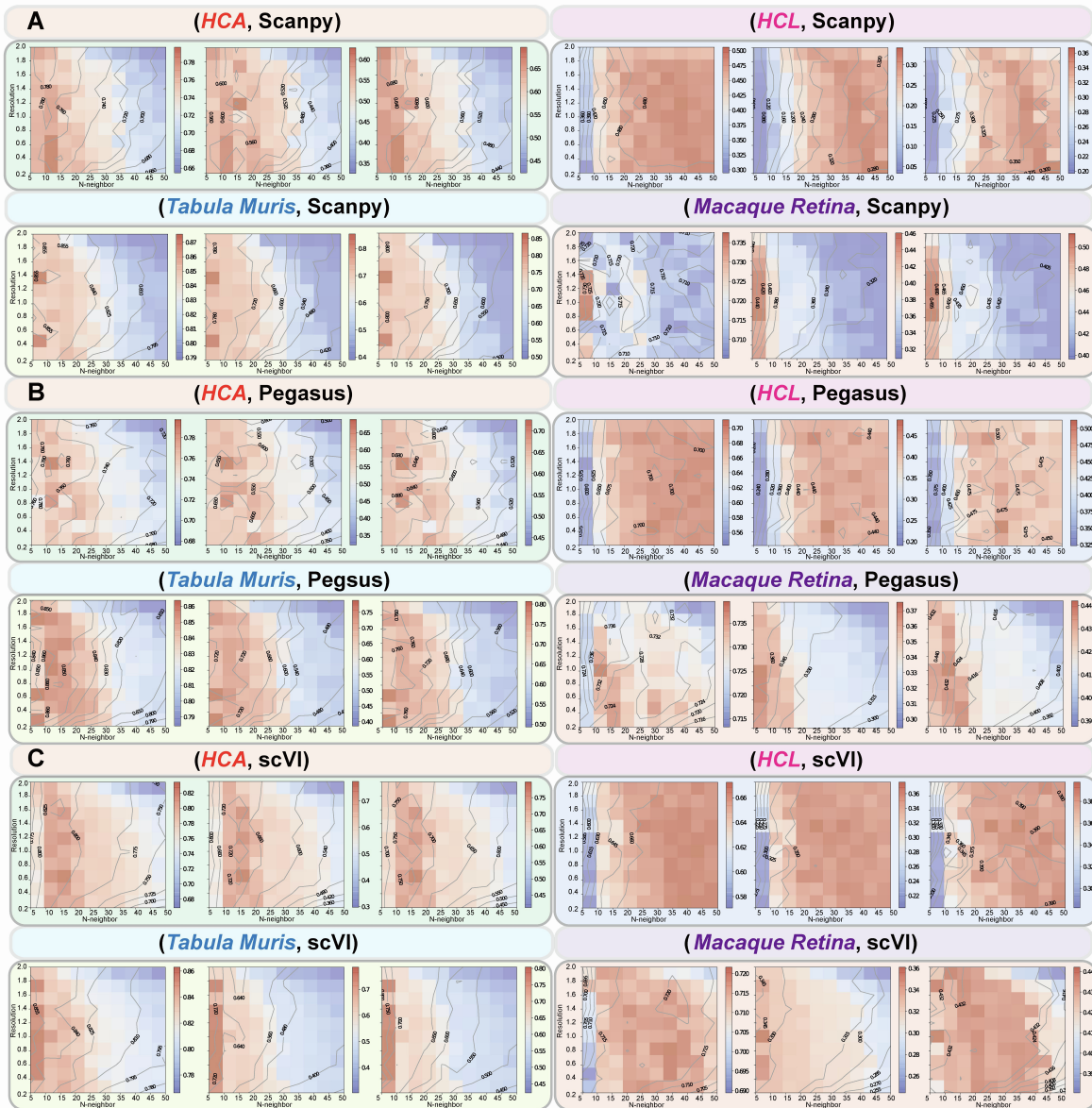
**Figure S14. The clustering performance with grid search for resolution and number of neighbors for *Scanpy (A), Pegasus (B),* and *scVI (C)* on the *HCA, HCL, Tabula Muris* and *Macaque Retina* dataset, related to Figure 2.** Contour maps depict different cluster metrics (i.e. *NMI*, *ARI* and *FMI*) with respect to different values of *Resolution* and *N-neighbors*.
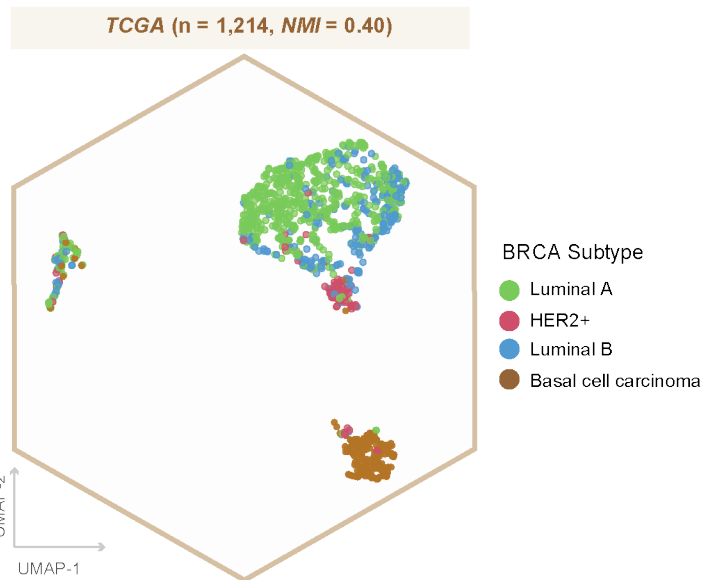
**Figure S15. The UMAP visualization plots of *tGPT* for molecular subtypes of BRCA from the *TCGA* datasets, related to Figure 2.** The NMI metric and annotation of cells are shown.
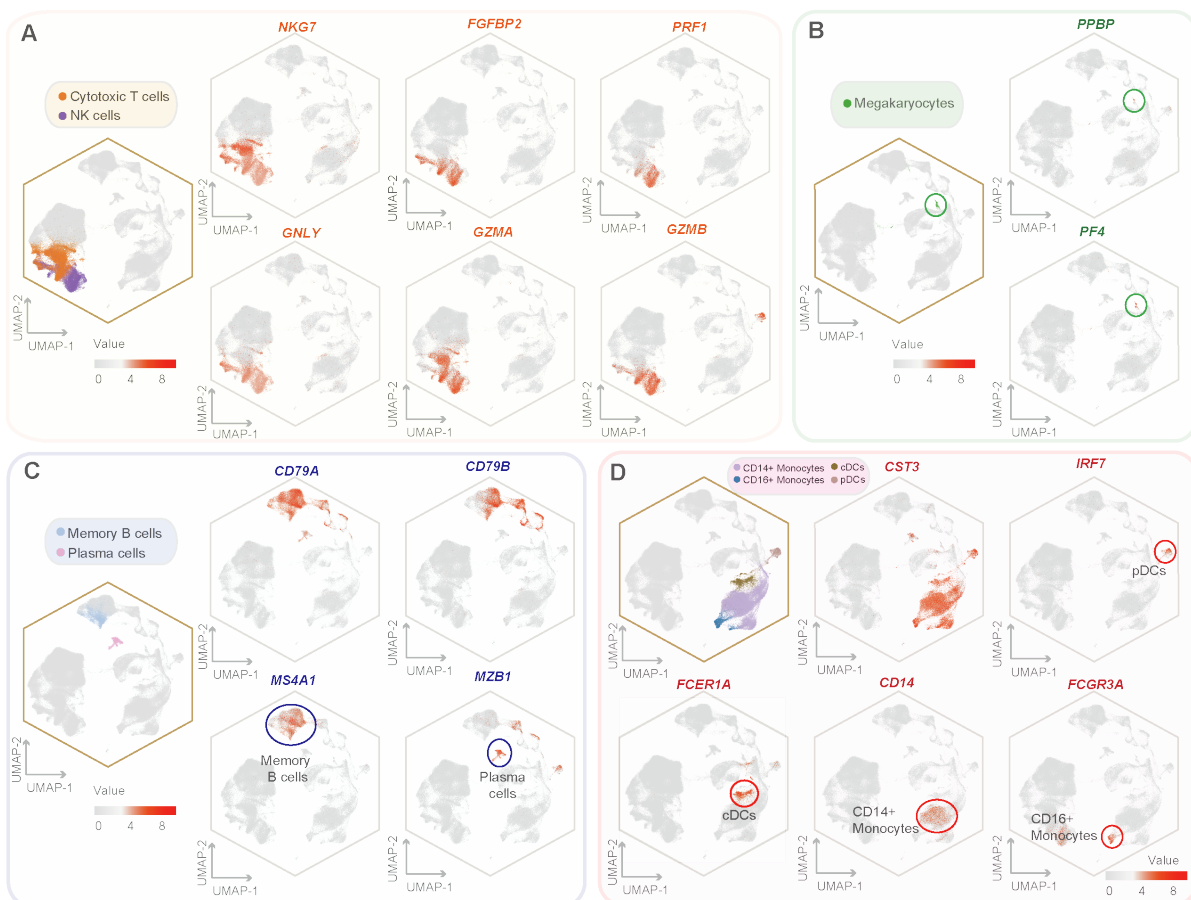


**Figure S16. Distinct features of different cell types from the *HCA* dataset learned by *tGPT*, related to Figure 3.** Scatter plots illustrating the distribution of attribution scores for different cell type specific genes across different cell types.
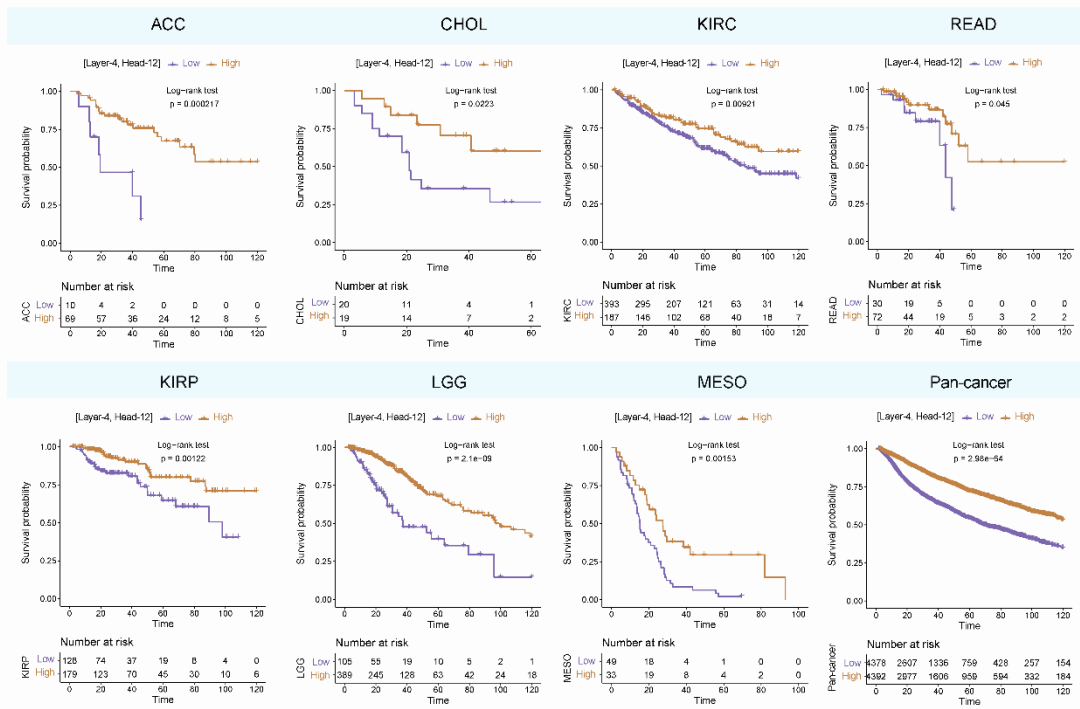
**Figure S17. The survival curves of the attention head related to overall survival across multiple cancer types, related to Figure 5.** ACC, Adrenocortical carcinoma; CHOL, Cholangiocarcinoma; KIRC, Kidney renal clear cell carcinoma; READ, Rectum adenocarcinoma; KIRP, Kidney renal papillary cell carcinoma; LGG, Brain Lower Grade Glioma; MESO, Mesothelioma.

**Table S2. The annotated cell labels on the different datasets, related to Figure 2.**

| HCA | HCL | Tabula Muris | Macaque Retina | GTEx | TCGA |
|---|---|---|---|---|---|
| Naive B cell | Neutrophil | myeloid cell | fGa | Adipose | ACC |
| CD14+ Monocyte | Stromal cell | alveolar macrophage | fGl | Tissue | BLCA |
| T-helper cell | Fibroblast | B cell | DB3b | Adrenal | NA |
| Pre-B cells | Monocyte | natural killer cell | FMB | Gland | DLBC |
| Naive CD8+ T cell | Macrophage | T cell | IMB | Blood | UCEC |
| Cytotoxic T cells | Antigen presenting cell (RPS high) | lung endothelial cell | DB5* | Vessel | SKCM |
| Pro-B cell | Mast cell | stromal cell | DB4 | Bladder | HNSC |
| CD4+ naive T cell | Sinusoidal endothelial cell | non-classical monocyte | DB2 | Brain | PRAD |
| NK cells | T cell | leukocyte | DB1 | Breast | KIRP |
| Erythroid cells | B cell | classical monocyte | BB/GB* | Blood | PAAD |
| cDCs | Dendritic cell | ciliated columnar cell of tracheobronchial tree | RB | Skin | SARC |
| Megakaryocyte | M2 Macrophage | type II pneumocyte | DB6 | Cervix Uteri | CESC |
| Memory B cell | Epithelial cell | mast cell | OFFx | Colon | COAD |
| Plasma cell | B cell (Plasmocyte) | monocyte | DB3a | Esophagus | LUSC |
| pDCs | Intercalated cell | granulocytopoietic cell | H1 | Fallopian Tube | READ |
| CD16+ Monocyte | Loop of Henle | promonocyte | H2 | Heart | KIRC |
| HSCs | Erythroid progenitor cell (RP high) | | MG | Kidney | LIHC |
| MSCs | Fetal epithelial progenitor | | Pericytes | Liver | BRCA |
| | | | Endo | Lung | OV |
| | | | Mic | | UCS |

| | | | | |
|---|---|---|---|---|
| Ureteric bud cell | granulocyte | M/L Cones | Salivary | GBM |
| Endothelial cell | erythroblast | S Cones | Gland | KICH |
| Endothelial cell (APC) | hematopoietic | Rods | Muscle | THCA |
| Smooth muscle cell | precursor cell | MG_OFF | Nerve | LGG |
| hESC | proerythroblast | PG_OFF | Ovary | LUAD |
| Stratified epithelial cell | late pro-B cell | PG_ON | Pancreas | MESO |
| Proximal tubule progenitor | basophil | MG_ON | Pituitary | PCPG |
| Fetal enterocyte | macrophage | fRGC | Prostate | TGCT |
| Myeloid cell | early pro-B cell | pGa | Small | UVM |
| Proliferating T cell | immature B cell | pGl | Intestine | THYM |
| Endothelial cell | Fraction A pre-pro B | pRGC | Spleen | CHOL |
| (endothelial to | cell | | Stomach | ESCA |
| mesenchymal transition) | basal cell of | | Testis | STAD |
| Enterocyte progenitor | epidermis | | Thyroid | |
| Enterocyte | keratinocyte | | Uterus | |
| Fetal stromal cell | Langerhans cell | | Vagina | |
| Erythroid cell | dendritic cell | | | |
| Hepatocyte/Endodermal | endothelial cell | | | |
| cell | fibroblast | | | |
| Fetal mesenchymal | endocardial cell | | | |
| progenitor | cardiac muscle cell | | | |
| Neutrophil (RPS high) | mesenchymal cell | | | |
| Fetal neuron | epithelial cell | | | |
| Fetal Neuron | blood cell | | | |
| Fetal endocrine cell | neuroendocrine cell | | | |
| AT2 cell | bladder cell | | | |
| Basal cell | bladder urothelial | | | |
| Epithelial cell | cell | | | |
| (intermediated) | luminal epithelial | | | |
| Chondrocyte | cell of mammary | | | |
| CB CD34+ | gland | | | |
| Fetal chondrocyte | basal cell | | | |
| Intermediated cell | kidney capillary | | | |
| Gastric endocrine cell | endothelial cell | | | |
| Primordial germ cell | mesangial cell | | | |
| Oligodendrocyte | kidney cell | | | |
| Astrocyte | kidney collecting | | | |
| Fasciculata cell | duct epithelial cell | | | |
| Immature sertoli cell (Pre- | kidney proximal | | | |
| Sertoli cell) | straight tubule | | | |
| Fetal fibroblast | epithelial cell | | | |
| Fetal skeletal muscle cell | kidney loop of Henle | | | |
| Fetal acinar cell | ascending limb | | | |
| Mesothelial cell | epithelial cell | | | |
| Goblet cell | immature T cell | | | |
| Ventricle cardiomyocyte | DN1 thymic pro-T | | | |
| Kidney intercalated cell | cell | | | |
| Thyroid follicular cell | hepatocyte | | | |
| Adrenal gland | duct epithelial cell | | | |

| | | | | |
|---|---|---|---|---|
| | inflammatory cell<br>Pancreas exocrine cell<br>Gastric chief cell | endothelial cell of<br>hepatic sinusoid<br>mesenchymal stem<br>cell<br>skeletal muscle<br>satellite cell | | | |

**Table S3. Running time of different methods on the four datasets, related to Figure 2.**

| Dataset<br><br>Method | Runtime (Seconds) | | | |
|---|---|---|---|---|
| | *HCA*<br>(n = 282,558) | *HCL*<br>(n = 586,135) | *Tabula Muris*<br>(n = 54,862) | *Macaque Retina* (n =<br>124, 965) |
| *tGPT* | 2513.5 | 2996.9 | 326.5 | 576.4 |
| *Scanpy* | 2352.2 | 3318.2 | 237.4 | 494.2 |
| *Pegasus* | 1575.8 | 1794.1 | 198.1 | 398.8 |
| *scVI* | 2849.5 | 3170.0 | 2084.0 | 2613.0 |

**Table S4. The generative metrics of *tGPT* on the four datasets, related to Figure 2.**

| Metrics<br><br>Dataset | *BLEU* |
|---|---|
| *HCA* | 0.77 |
| *HCL* | 0.69 |
| *Tabula Muris* | 0.76 |
| *Macaque Retina* | 0.75 |