

Berke - Superbio.ai





Haotian Cui

scGPT: towards building a foundation model for single-cell multi-omics using generative AI

Authors: Haotian Cui*, Chloe Wang*, Hassaan Maan, Kuan Pang, Fengning Luo, Bo Wang

Present by: Haotian Cui





Haotian Cui

scGPT: towards building a foundation model for single-cell multi-omics using generative AI

Authors: [Haotian Cui*](#), [Chloe Wang*](#), [Hassaan Maan](#), [Kuan Pang](#),
[Fengning Luo](#), [Bo Wang](#)

Present by: [Haotian Cui](#)



Thank you for the introduction! I am so thrilled to have the opportunity to present our work scGPT here today. I will start with the introduction to the work and then let's try to see some application demo around zero-shot reference mapping. Hopefully these can give you some idea about how to easily use the model

So let me start and feel free to interrupt with any questions

scGPT is about building foundation models for single-cell omics

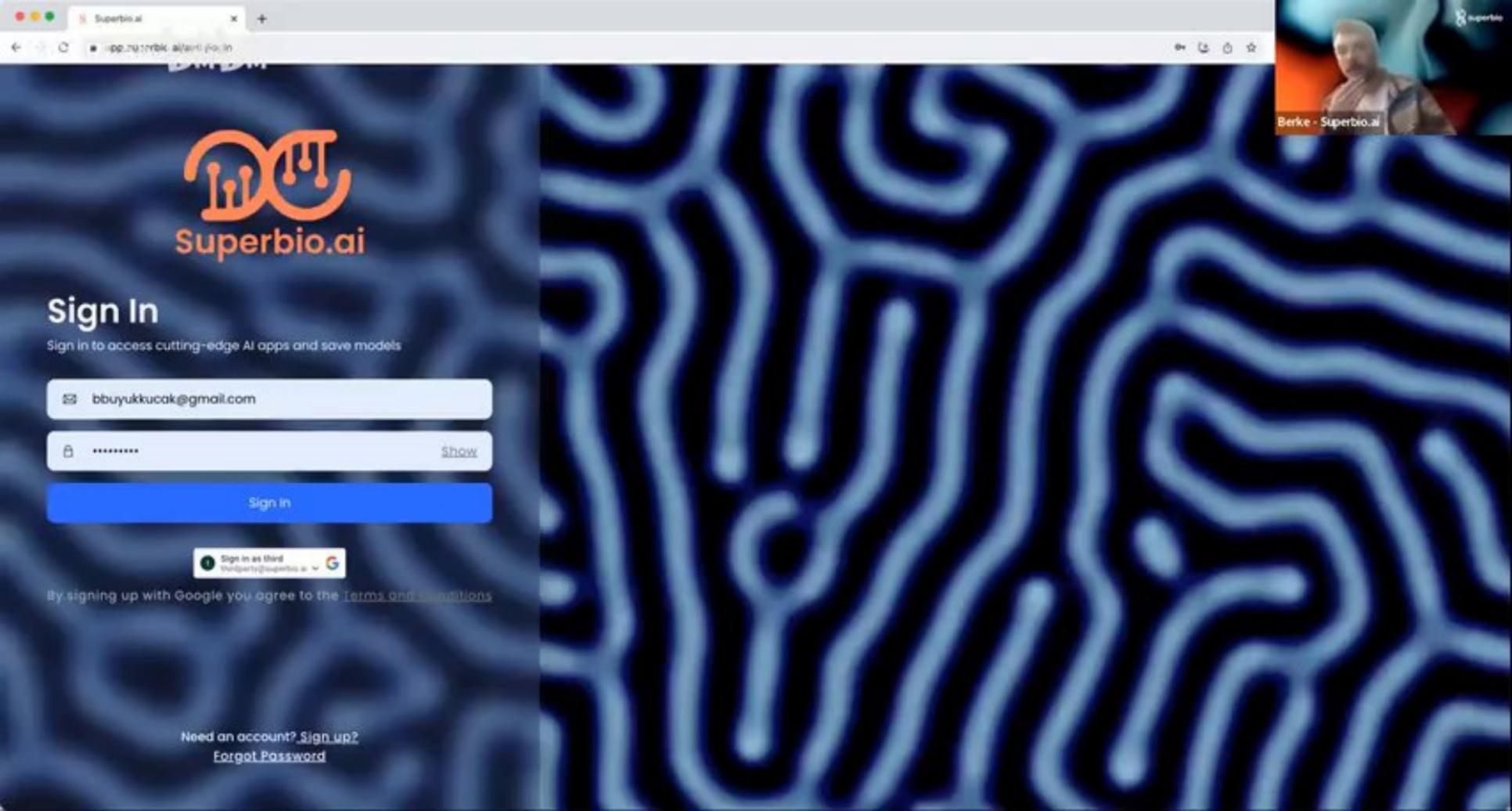
scGPT: towards building a foundation model for single-cell multi-omics using generative AI

Authors: Haotian Cui*, Chloe Wang*, Hassaan Maan, Kuan Pang, Fengning Luo, Bo Wang

Present by: Haotian Cui







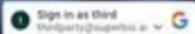
Sign In

Sign in to access cutting-edge AI apps and save models

bbuyukkucak@gmail.com

Show

Sign In



By signing up with Google you agree to the [Terms and Conditions](#)

Need an account? [Sign up?](#)
[Forgot Password](#)

Berke - Superbio.ai

Superbio

Feedback Tutorials Upgrade

Berke - Superbio.ai

App Store My Library

What are you looking for? SORT: Recently added APPS: 650

FOCUS Clear

Biomarker Discovery an... Biomedical Analysis
 Biomedical Image Analy... Biomedical Language Ma...
 Cheminformatics +10

METHOD Clear

3D Protein Structure P... AIMNET ANI AutoML
 Base Caller +42

DATA EXTENSION Clear

bai bam bed bim csv +17

[Submit Your Own App](#)
[Can't Find What You Need?](#)

scGPT: Reference Mapping Using Cell Embedding
 WangLab at University of Toronto | Cui et al **Single-Cell Bioinformatics** **Cell Type Annotation** ▶ 26 ☆ 1
[Example Results](#)

ProtGPT2: De Novo Protein Generation
 Noelia Ferruz, Steffen Schmidt & Birte Höcker **Biomedical Language Models** **Proteomics** **GPT** ▶ 10 ☆ 2
 De novo protein generator and protein sequence completer.
[Example Results](#)

ProtGPT2: Finetune
 Noelia Ferruz, Steffen Schmidt & Birte Höcker **Biomedical Language Models** **Proteomics** **GPT** ▶ 2 ☆ 1
[Example Results](#)

RFDiffusion: Accurate Protein Design
 RosettaCommons **Structural Bioinformatics** **3D Protein Structure Prediction** **Diffusion** ▶ 7 ☆ 3
[Example Results](#)

SearchBio: Find & Summarize Scientific Articles
 Superbio **Biomedical Language Models** **GPT** ▶ 29 ☆ 4

[Terms of Service](#) [Privacy Policy](#)



AlphaFold2: 3D Protein Structure Prediction/DeepMind v0.2

Structural Bioinformatics 3D Protein Structure Prediction Fasta [Example Results](#) [See Source](#)

Released: Jul-01-2022

Revolutionary 3D protein structure prediction model with high accuracy that allows for a close study of proteins without an expensive X-ray crystallography process

Example use case: Prediction of protein disorders, studying enzymes, and developing novel protein-based drugs

> [Show More](#)

Previous Job Parameters

AF-20231130-210551-cpu

AF-20231114-041757-cpu

AF-20231010-155238-cpu

AF-20230913-004251-cpu

Protein Name

MQRQNPNPYQ_4 MQRQNPNPYQNNQAQPNRPNQNPYGTQNNQ

Myoglobin MGLSDGEWQLVLNVWGKVEADIPGHGQEVLRIFK

Myoglobin MGLSDGEWQLVLNVWGKVEADIPGHGQEVLRIFK

Myoglobin MGLSDGEWQLVLNVWGKVEADIPGHGQEVLRIFK

RUNS ▶ 96 SAVE ↻ 15 [Share](#)

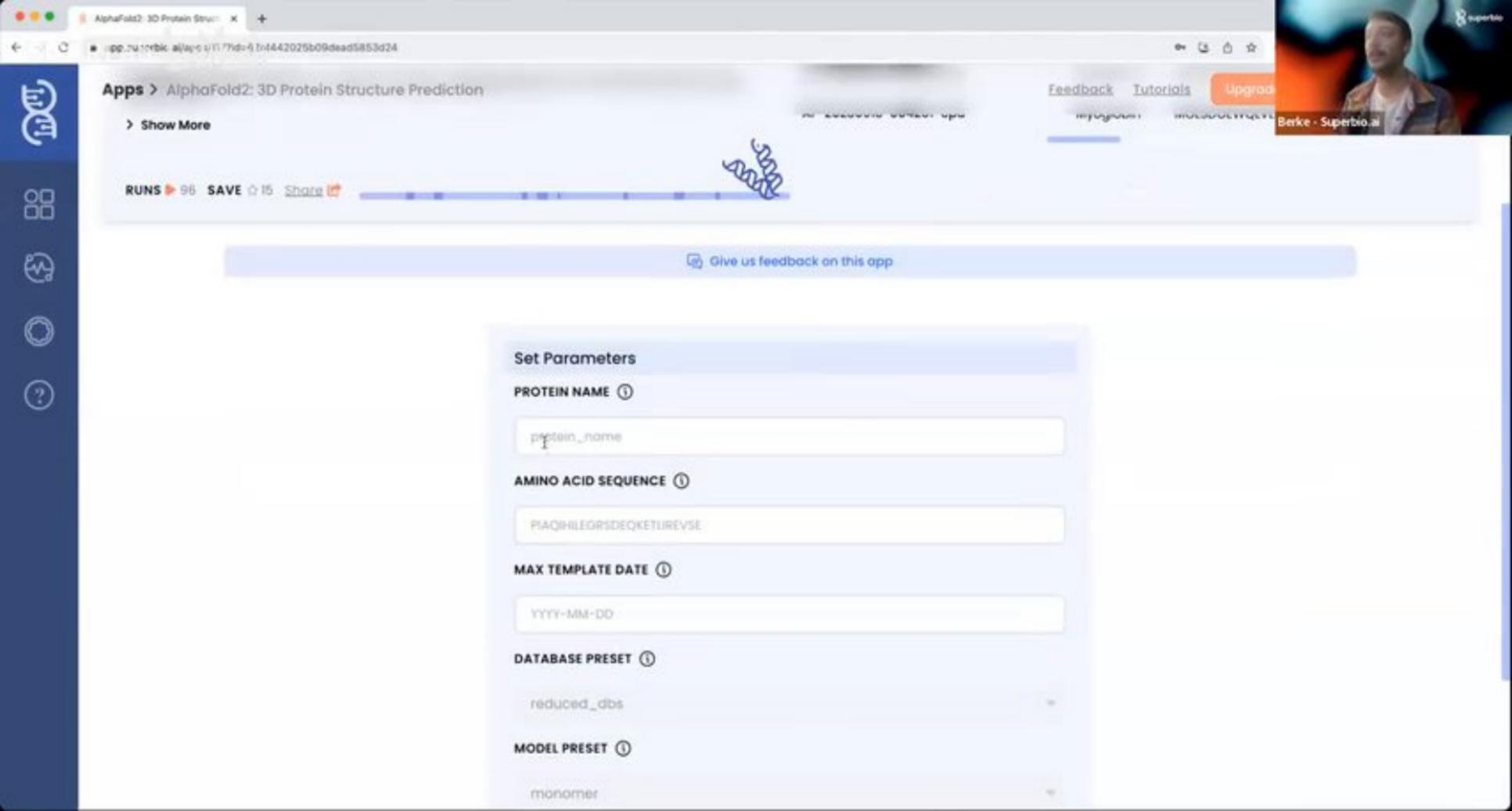
[Give us feedback on this app](#)

Set Parameters

PROTEIN NAME ⓘ

AMINO ACID SEQUENCE ⓘ

MAX TEMPLATE DATE ⓘ



Apps > AlphaFold2: 3D Protein Structure Prediction

> Show More

RUNS ▶ 96 SAVE ▶ 15 [Share](#)

[Give us feedback on this app](#)

Set Parameters

PROTEIN NAME ⓘ

AMINO ACID SEQUENCE ⓘ

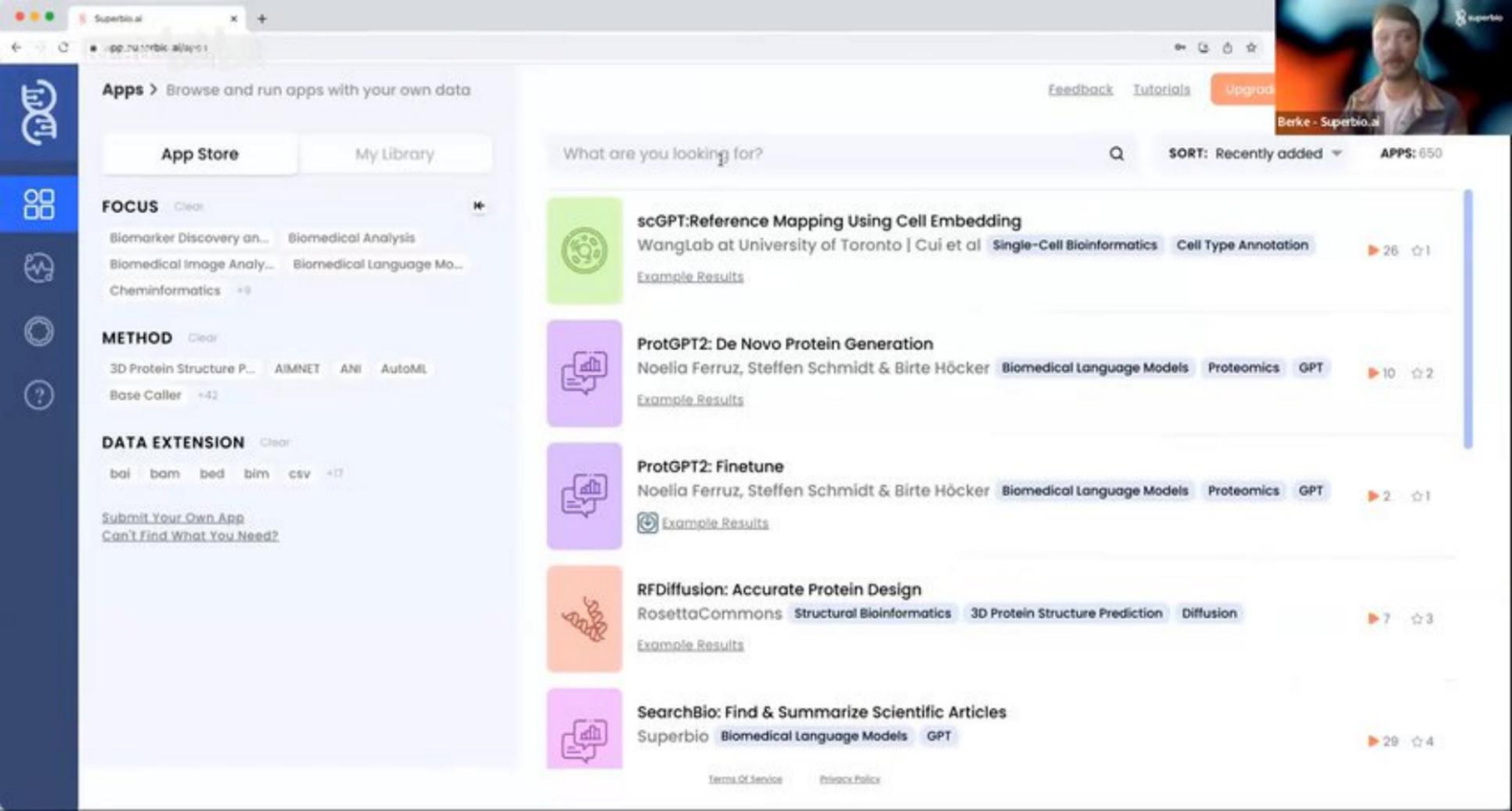
MAX TEMPLATE DATE ⓘ

DATABASE PRESET ⓘ

reduced_dbs

MODEL PRESET ⓘ

monomer



Apps > Browse and run apps with your own data

App Store

My Library

FOCUS [Clear](#)

[Biomarker Discovery an...](#) [Biomedical Analysis](#)
[Biomedical Image Analy...](#) [Biomedical Language Mo...](#)
[Cheminformatics](#) +9

METHOD [Clear](#)

[3D Protein Structure P...](#) [AIMNET](#) [ANI](#) [AutoML](#)
[Base Caller](#) +42

DATA EXTENSION [Clear](#)

[bai](#) [bam](#) [bed](#) [bim](#) [csv](#) +17

[Submit Your Own App](#)
[Can't Find What You Need?](#)

[Feedback](#) [Tutorials](#) [Upgrade](#)

What are you looking for?



SORT: Recently added

APPS: 650



scGPT:Reference Mapping Using Cell Embedding

WangLab at University of Toronto | Cui et al [Single-Cell Bioinformatics](#) [Cell Type Annotation](#)

▶ 26 ☆ 1

[Example Results](#)



ProtGPT2: De Novo Protein Generation

Noelia Ferruz, Steffen Schmidt & Birte Höcker [Biomedical Language Models](#) [Proteomics](#) [GPT](#)

▶ 10 ☆ 2

[Example Results](#)



ProtGPT2: Finetune

Noelia Ferruz, Steffen Schmidt & Birte Höcker [Biomedical Language Models](#) [Proteomics](#) [GPT](#)

▶ 2 ☆ 1

[Example Results](#)



RFDiffusion: Accurate Protein Design

RosettaCommons [Structural Bioinformatics](#) [3D Protein Structure Prediction](#) [Diffusion](#)

▶ 7 ☆ 3

[Example Results](#)



SearchBio: Find & Summarize Scientific Articles

Superbio [Biomedical Language Models](#) [GPT](#)

▶ 29 ☆ 4



scGPT: towards building a foundation model for single-cell multi-omics using generative AI

Authors: Haotian Cui*, Chloe Wang*, Hassaan Maan, Kuan Pang, Fengning Luo, Bo Wang

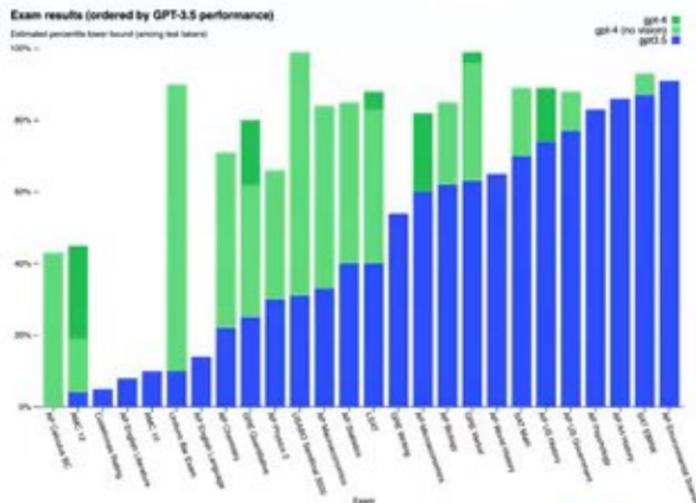
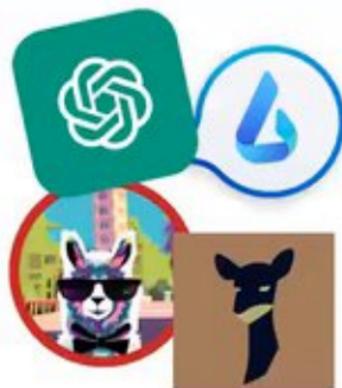
Present by: Haotian Cui





Haotian Cui

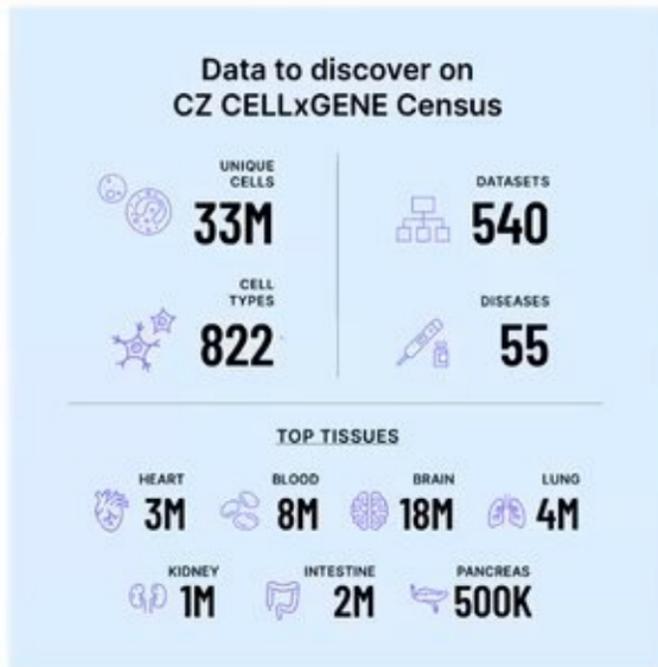
- Generative pretraining has achieved remarkable success in various domains such as natural language processing and computer vision
- Specifically, the combination of large-scale diverse datasets and pre-trained transformers has emerged as a promising approach for developing foundation models.



GPT-4 performances on standard exams. Credit to [OpenAI](#)



- While texts are made up of words, cells can be characterized by genes.
- This analogy inspires us to explore foundation models for single cell. In the long run, we expect such model has the potential to gain general understanding of molecular biology.
- This is feasible based on two considerations:
 - Transformers would be suitable to learn interactions between gene tokens.
 - The vast scale of sequencing data, of diverse cell types and conditions, and it is growing exponentially.

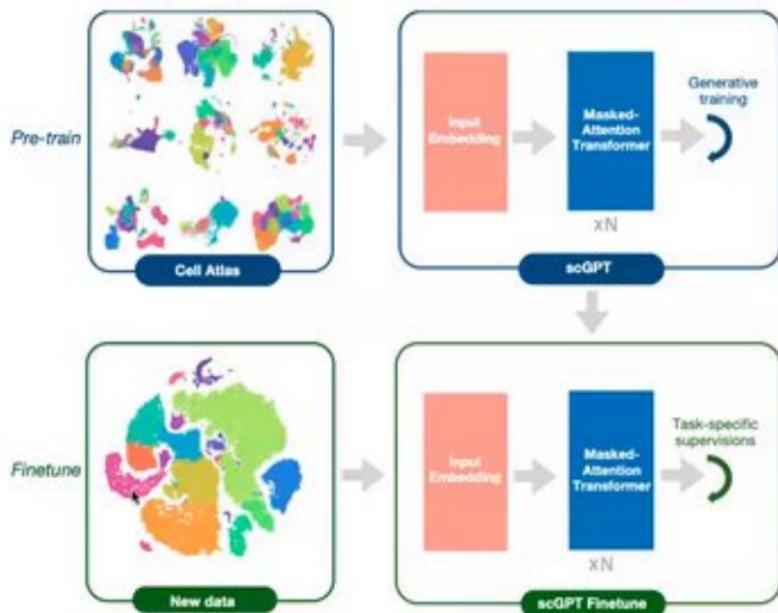


We are particularly excited about the recent launch of CellxGene census. Figure credit to [@JCoolScience](#)



Haotian Cui

- We present scGPT to construct a single-cell foundation model. The workflow contains two steps:
 1. Generative **pre-training** on over 33 million normal human cells
 2. **Fine-tuning** for downstream tasks



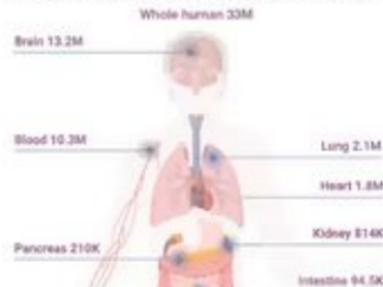


Haotian Cui

- In the pre-training:
 - Large-scale expression matrices were collected and preprocessed
 - We selected around 33 million normal human cells of 51 organs/tissues from the CZ CellxGene collection
 - Numbers of cells from major tissues shown on the right



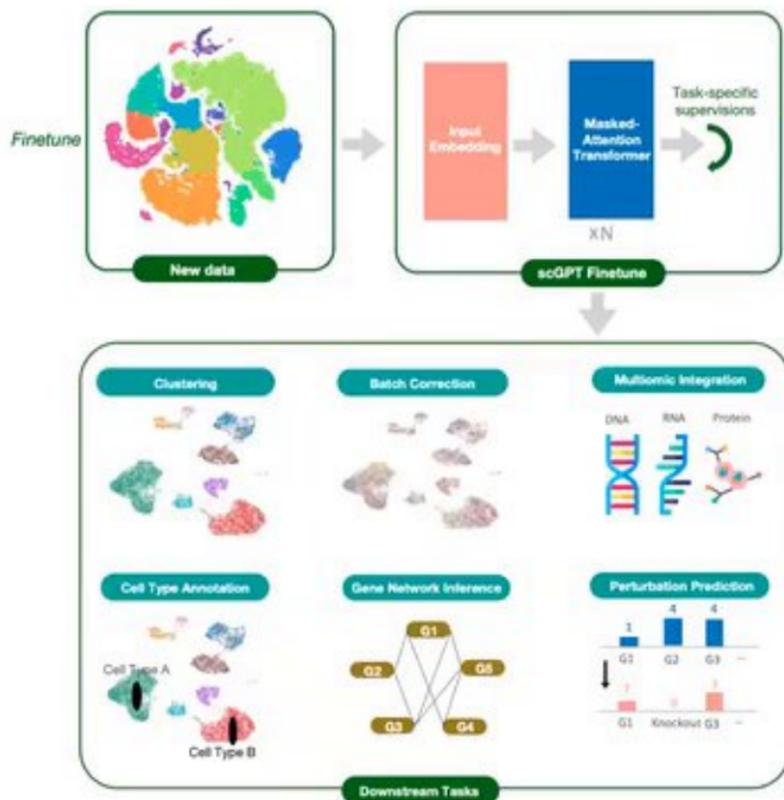
Cell numbers and origin tissues included in the pretraining





- In fine-tuning:

- A set of versatile learning objectives are adopted to readily support various applications
- Self-supervised objectives:
 - Gene Expression Prediction (GEP)
 - Gene Expression Prediction for Cell Modelling (GEPC)
 - Elastic Cell Similarity (ECS)
- Supervised objectives:
 - Domain Adaptation via Reverse Back-propagation (DAR)
 - Cell Type Classification (CLS)



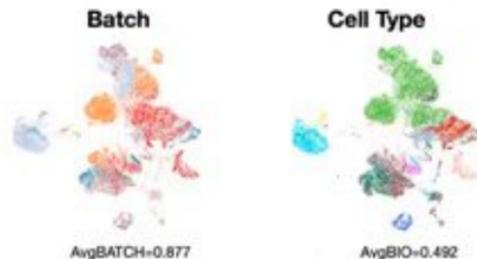


Haotian Cui

- **Zero-shot**, meaning using the pretrained model to generate embeddings for new data, without any further training, **much faster, more accessible**
- UMAP of zero-shot embeddings on two disease datasets, compared to HVGs
 - Considerable ability to distinguish cell types
 - Note the pretraining process is not designed for mitigating technical batch effects

COVID-19

scGPT zero-shot

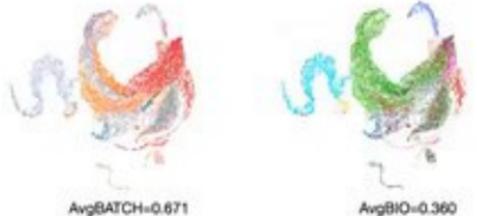


- Sample Batch**
- 16X
 - COVID 19 (spiral)
 - Fraxipag
 - H3L
 - Krasnow_batch 1a
 - Krasnow_batch 2
 - Krasnow_batch 3
 - Krasnow_batch 2
 - Krasnow_pretrain 3
 - Northwestern_Moharik_2018Reynhan
 - Origan_A
 - Origan_P
 - Origan_SJ
 - Sample_Meyer_2019Madison
 - Sun_sample1_CS
 - Sun_sample2_AC
 - Sun_sample3_TB
 - Sun_sample4_TC

Cell Type

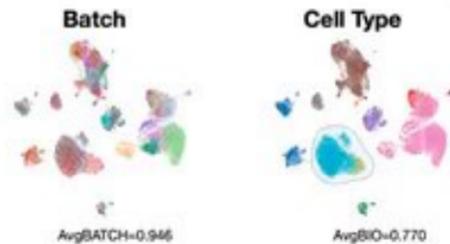
- A01
- A02
- B cell
- CCR7+ T
- CD4+ T cells
- CD8 T
- CD8+ T cells
- CD28+ B cells
- CD14+ Monocytes
- CD18+ Monocytes
- CD28+ B cells
- Clonal
- DC_activated
- Dendritic cell
- Erythromyelin
- Erythroid progenitors
- HSPCs
- IGM21+ Dendritic
- M2 Macrophage
- Macrophages
- Mast cells
- Megakaryocyte progenitors
- Negakaryocytes
- Monocyte progenitors
- Monocytes
- NK

HVG+PCs



Lung-Kim

scGPT zero-shot



- Patient**
- P0006
 - P0008
 - P0018
 - P0019
 - P0020
 - P0025
 - P0026
 - P0028
 - P0030
 - P0031
 - P0034
 - P1006
 - P1008
 - P1009
 - P1016
- Cell Type**
- B_cell
 - Dendritic
 - Endothelial
 - Erythelial
 - Fibroblast
 - Macrophage
 - Malignant
 - Mast
 - MC_cell
 - T_cell

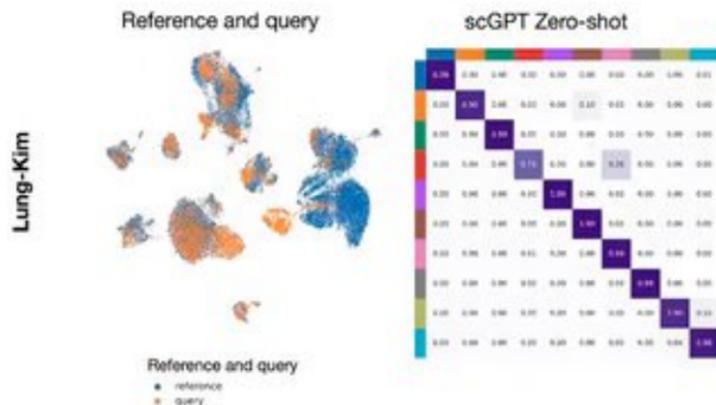
HVG+PCs





Haotian Cui

- We also use the embeddings in a reference mapping manner
- We split the data to simulate a setting that some of the patient cells, query set, were not annotated. Cell types were propagated from reference data set to label the query set.
- The zero-shot workflow demonstrates competitive accuracy across comparisons.



Dataset	Method	Evaluation Metrics			
		<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>MacroF1</i>
COVID-19	scGPT (fine-tuned)	0.891	0.577	0.567	0.547
	scGPT (zero-shot)	0.867	0.513	0.476	0.468
	Azimuth	0.878	0.515	0.418	0.444
	expiMap	0.730	0.441	0.325	0.335
Lung-Kim	scGPT (fine-tuned)	0.974	0.959	0.967	0.962
	scGPT (zero-shot)	0.968	0.970	0.933	0.948
	Azimuth	0.970	0.957	0.962	0.959
	expiMap	0.920	0.940	0.846	0.878



Hao Tian Cui

- In the pre-training:
 - Large-scale expression matrices were collected and preprocessed
 - We selected around 33 million normal human cells of 51 organs/tissues from the CZ CellxGene collection
 - Numbers of cells from major tissues shown on the right



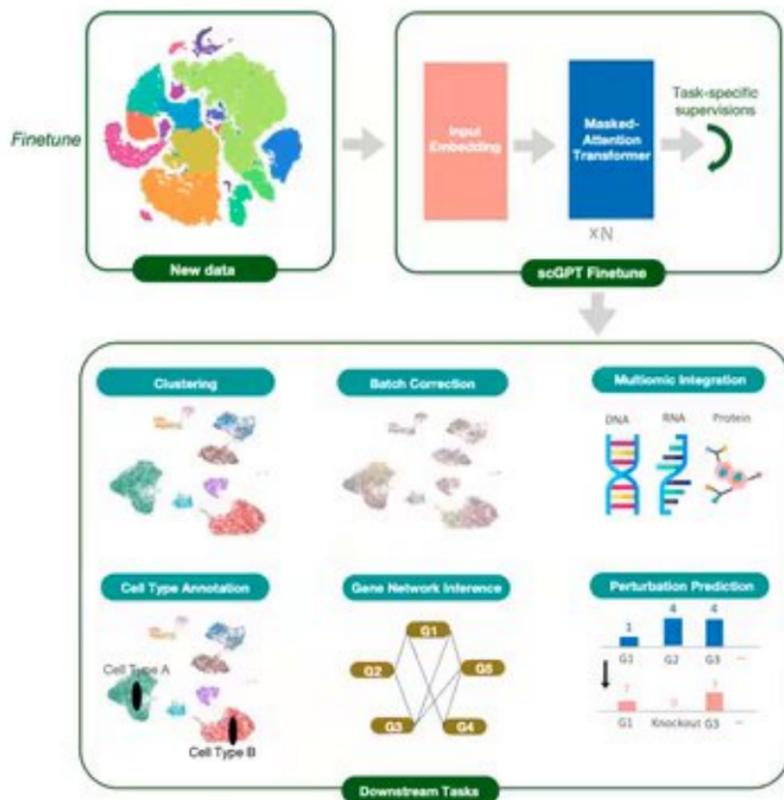
Cell numbers and origin tissues included in the pretraining





- In fine-tuning:

- A set of versatile learning objectives are adopted to readily support various applications
- Self-supervised objectives:
 - Gene Expression Prediction (GEP)
 - Gene Expression Prediction for Cell Modelling (GEPC)
 - Elastic Cell Similarity (ECS)
- Supervised objectives:
 - Domain Adaptation via Reverse Back-propagation (DAR)
 - Cell Type Classification (CLS)

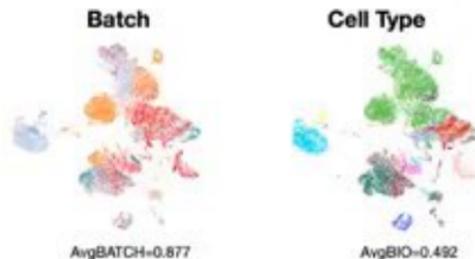




- **Zero-shot**, meaning using the pretrained model to generate embeddings for new data, without any further training, **much faster, more accessible**
- UMAP of zero-shot embeddings on two disease datasets, compared to HVGs
 - Considerable ability to distinguish cell types
 - Note the pretraining process is not designed for mitigating technical batch effects

COVID-19

scGPT zero-shot



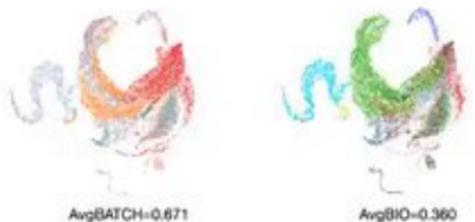
Sample Batch

- 18X
- COVID-19 (spatial)
- Freytag
- H3C
- Krasnow_0001_04
- Krasnow_0001_02
- Krasnow_0001_03
- Krasnow_0001_02
- Krasnow_0001_03
- Northwestern_Michael_2020Mayhew
- Oelgem_A
- Oelgem_T
- Oelgem_U
- Sample_Meyer_2021Madison
- Sun_sample1_CS
- Sun_sample1_AC
- Sun_sample1_TB
- Sun_sample1_TC

Cell Type

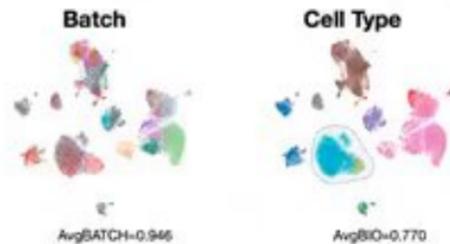
- A17
- A12
- B cell
- CD8+ T
- CD4+ T cells
- CD8 T
- CD8+ T cells
- CD20+ B cells
- CD14+ Monocytes
- CD15+ Monocytes
- CD20+ B cells
- Cluster
- DC_activated
- Dendritic cell
- Erythrocytes
- Erythroid progenitors
- HSPCs
- K562+ dendritic
- M2 Macrophage
- Macrophages
- Mast cells
- Megakaryocyte progenitors
- Megakaryocytes
- Monocyte progenitors
- Monocytes
- NK

HVG+PCs



Lung-Kim

scGPT zero-shot



- Patient
- P0006
 - P0008
 - P0018
 - P0019
 - P0020
 - P0025
 - P0029
 - P0030
 - P0032
 - P0034
 - P1006
 - P1018
 - P1049
 - P1058
- Cell Type
- E_cell
 - Dendritic
 - Endothelial
 - Epithelial
 - Fibroblast
 - Macrophage
 - Mast
 - NK cell
 - T_cell

HVG+PCs

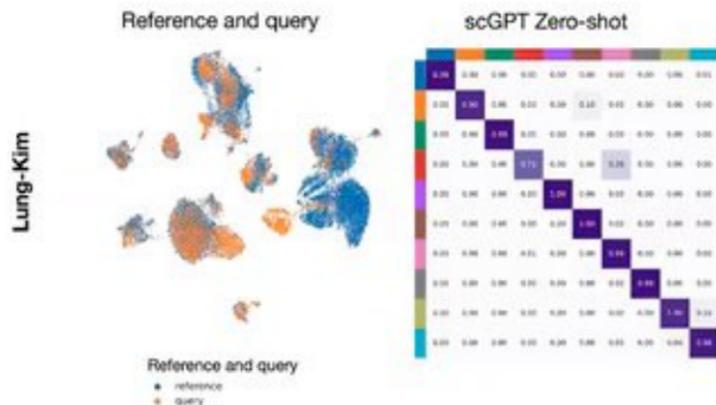


MCBRLab Zero-shot reference mapping



Hao Tian Cui

- We also use the embeddings in a reference mapping manner
- We split the data to simulate a setting that some of the patient cells, query set, were not annotated. Cell types were propagated from reference data set to label the query set.
- The zero-shot workflow demonstrates competitive accuracy across comparisons.



Dataset	Method	Evaluation Metrics			
		<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>MacroF1</i>
COVID-19	scGPT (fine-tuned)	0.891	0.577	0.567	0.547
	scGPT (zero-shot)	0.867	0.513	0.476	0.468
	Azimuth	0.878	0.515	0.418	0.444
	expiMap	0.730	0.441	0.325	0.335
Lung-Kim	scGPT (fine-tuned)	0.974	0.959	0.967	0.962
	scGPT (zero-shot)	0.968	0.970	0.933	0.948
	Azimuth	0.970	0.957	0.962	0.959
	expiMap	0.920	0.940	0.846	0.878



- Zero-shot embeddings can further be used to map cells and propagate meta-info from large cell atlas
- For this new feature, we generated the cell embeddings for millions of cells on CellXGene, and build an index for efficient similarity search using the faiss toolbox.
- This can be used to query new cells super fast, [https://github.com/bowang-lab/scGPT/blob/main/tutorials/Tutorial Reference Mapping.ipynb](https://github.com/bowang-lab/scGPT/blob/main/tutorials/Tutorial%20Reference%20Mapping.ipynb)

The search runs super fast, especially on GPU. Here the similarity search for 4,000 query cells within the whole reference of millions should take around 7 second on CPU and 0.1 second on GPU.

```
In [10]: %time
          k = 50
          # test with the first 100 cells
          distances, idx = index.search(test_embd, k)
```

```
CPU times: user 252 ms, sys: 493 ms, total: 745 ms
Wall time: 133 ms
```

Demo of zero-shot application



Hao Tian Cui

- On google colab



Demo of zero-shot application

- On google colab



Microsoft PowerPoint | File | Edit | View | Insert | Format | Layout | Tools | Slide Show | Window | Help

Tutorial_ZeroShot_Reference_Mapping.ipynb | New Tab | scGPT/tutorials at main · bowang-lab/scGPT

colab:usurdf:google.com/github/bowang-lab/scGPT/biob/main/tutorials/zero-shot/Tutorial_ZeroShot_Reference_Mapping.ipynb#scrollTo=q-SKW5BGecq

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

14 Colab AI



Zero-shot reference mapping tutorial with scGPT

Introduction

This tutorial covers the zero-shot reference mapping with scGPT. This workflow achieves accurate and fast reference mapping for scRNA-seq datasets without fine-tuning (or any extensive training) of scGPT. To further boost the performance, we recommend fine-tuning scGPT.

We will use COVID-19 dataset to demonstrate the zero-shot reference mapping. You can download the processed reference and query datasets from [here](#). The COVID-19 dataset is derived from the work by [Lotfollahi et al.](#), which contains 18 distinct batches and diverse samples from lung tissues. The reference dataset consists of 15,997 cells and the query dataset contains 4,003 cells. You may place the dataset under `data` directory in the outer level.

Particularly, we use the `scGPT_human` model to provide embeddings out of the box. You may download it from [here](#).

The zero-shot reference mapping workflow is as follows:

1. Load and pre-process the dataset
2. Generate scGPT embeddings for each cell in reference and query datasets
3. Transfer the annotations from reference to query dataset

At the [appendix](#) of this tutorial, we will also showcase the zero-shot reference mapping on Lung dataset. You can find the dataset [here](#).

We use a similarity-based method for transferring the annotation, which involves comparing the similarity between the query cell embedding and the reference cell embeddings. We use [FAISS](#) to perform the similarity search.

[Open in Colab](#)

```
[1] # Specifically for Google Colab, install dependencies and download data
```

```
import os
import sys

if "google.colab" in sys.modules:
    print("Running on Google Colab")
    print("Installing dependencies...")
    !pip install -U scgpt
```

0s completed at 12:03 PM

Files

main + Q

Go to file

- vscode
- data
- docs
- examples
- scgpt
- tests
- tutorials
 - zero-shot
 - Tutorial_Annotation.ipynb
 - Tutorial_Attention_GRN.ipynb
 - Tutorial_GRN.ipynb
 - Tutorial_Integration.ipynb
 - Tutorial_Multiomics.ipynb
 - Tutorial_Perturbation.ipynb
 - Tutorial_Reference_Mapping.L...
 - build_atlas_index_faiss.py
 - .gitignore
 - .readthedocs.yaml
 - LICENSE
 - README.md
 - poetry.lock
 - pyproject.toml

scGPT / tutorials /

Add file ...

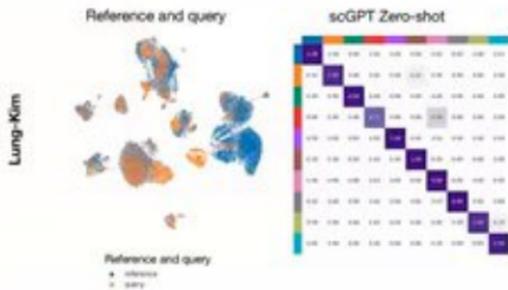
subercul add colab badges 18 hours ago History

Name	Last commit message	Last commit date
..		
zero-shot	add colab badges	18 hours ago
Tutorial_Annotation.ipynb	update annotation tutorial	5 months ago
Tutorial_Attention_GRN.ipynb	add attention-grn tutorial	5 months ago
Tutorial_GRN.ipynb	update doc with tutorial notebooks	6 months ago
Tutorial_Integration.ipynb	update load model in notebook	last month
Tutorial_Multiomics.ipynb	update multi-omic tutorial	2 months ago
Tutorial_Perturbation.ipynb	update docs	3 months ago
Tutorial_Reference_Mapping.ipynb	update value counts for predictions, fix #135	20 hours ago
build_atlas_index_faiss.py	update build atlas index faiss	3 months ago



Zero-shot reference mapping

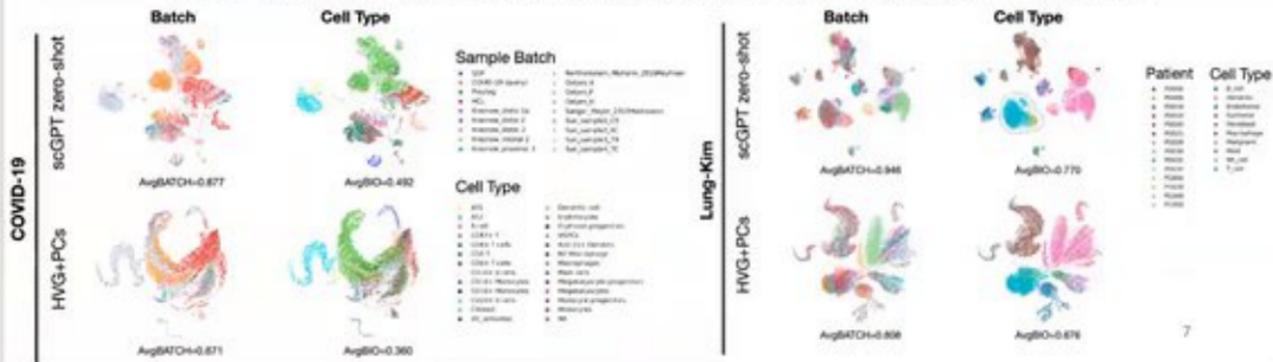
- We also use the embeddings in a reference mapping manner
- We split the data to simulate a setting that some of the patient cells, query set, were not annotated. Cell types were propagated from reference data set to label the query set.
- The zero-shot workflow demonstrates competitive accuracy across comparisons.



Dataset	Method	Evaluation Metrics			
		Accuracy	Precision	Recall	MacroF1
COVID-19	scGPT (fine-tuned)	0.891	0.577	0.567	0.547
	scGPT (zero-shot)	0.867	0.513	0.476	0.468
	Azimuth	0.878	0.515	0.418	0.444
	expMap	0.730	0.441	0.325	0.335
Lung-Kim	scGPT (fine-tuned)	0.974	0.959	0.967	0.962
	scGPT (zero-shot)	0.968	0.970	0.933	0.948
	Azimuth	0.970	0.957	0.962	0.959
	expMap	0.920	0.940	0.846	0.878

Zero-shot applications

- **Zero-shot**, meaning using the pretrained model to generate embeddings for new data, without any further training, **much faster, more accessible**
- UMAP of zero-shot embeddings on two disease datasets, compared to HVGs
 - Considerable ability to distinguish cell types
 - Note the pretraining process is not designed for mitigating technical batch effects



Have been working new extensions exploring the possibility of zero-shot applications. Here, I'd like to share some new results

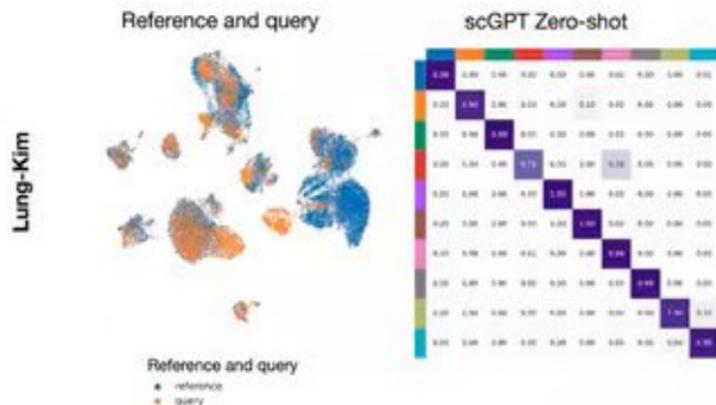
Particularly want to mention the influence of pretraining, so of the technical info is actually signals

Haotian Cui

MCBRLab ~~Zero~~shot reference mapping



- We also use the embeddings in a reference mapping manner
- We split the data to simulate a setting that some of the patient cells, query set, were not annotated. Cell types were propagated from reference data set to label the query set.
- The zero-shot workflow demonstrates competitive accuracy across comparisons.

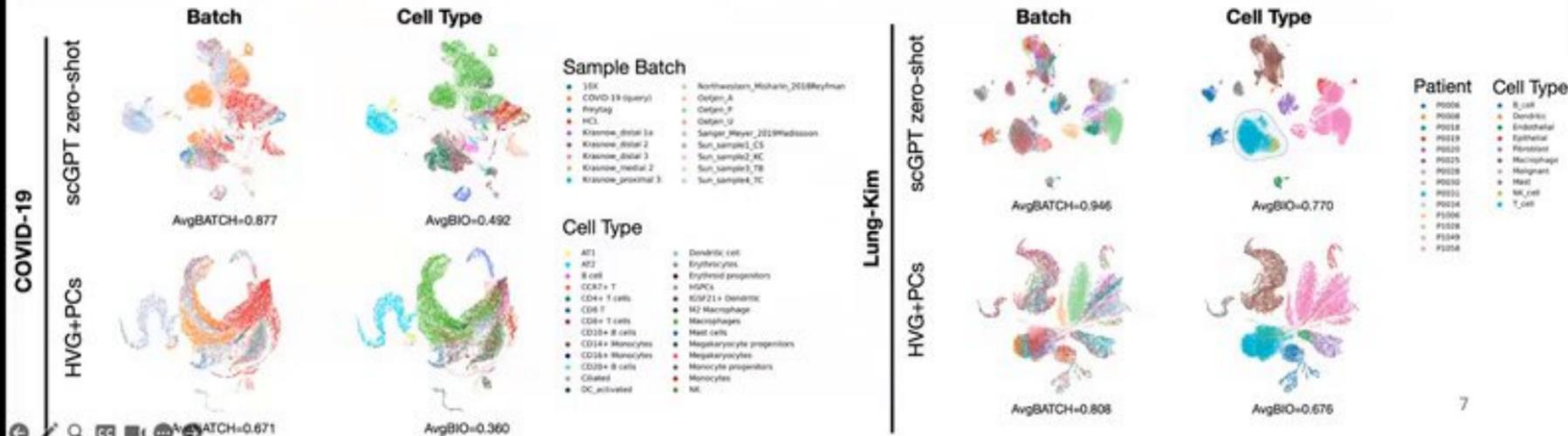


Dataset	Method	Evaluation Metrics			
		<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>MacroF1</i>
COVID-19	scGPT (fine-tuned)	0.891	0.577	0.567	0.547
	scGPT (zero-shot)	0.867	0.513	0.476	0.468
	Azimuth	0.878	0.515	0.418	0.444
	expiMap	0.730	0.441	0.325	0.335
Lung-Kim	scGPT (fine-tuned)	0.974	0.959	0.967	0.962
	scGPT (zero-shot)	0.968	0.970	0.933	0.948
	Azimuth	0.970	0.957	0.962	0.959
	expiMap	0.920	0.940	0.846	0.878



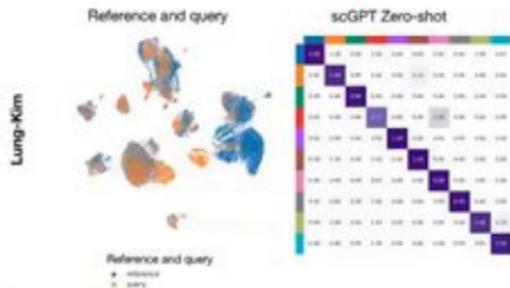
Haotian Cui

- **Zero-shot**, meaning using the pretrained model to generate embeddings for new data, without any further training, **much faster, more accessible**
- UMAP of zero-shot embeddings on two disease datasets, compared to HVGs
 - Considerable ability to distinguish cell types
 - Note the pretraining process is not designed for mitigating technical batch effects



Zero-shot reference mapping

- We also use the embeddings in a reference mapping manner
- We split the data to simulate a setting that some of the patient cells, query set, were not annotated. Cell types were propagated from reference data set to label the query set.
- The zero-shot workflow demonstrates competitive accuracy across comparisons.



Dataset	Method	Evaluation Metrics			
		Accuracy	Precision	Recall	MacroF1
COVID-19	scGPT (fine-tuned)	0.891	0.577	0.567	0.547
	scGPT (zero-shot)	0.867	0.513	0.476	0.468
	Azimuth	0.878	0.515	0.418	0.444
	expMap	0.730	0.441	0.325	0.335
Lung-Kim	scGPT (fine-tuned)	0.974	0.959	0.967	0.962
	scGPT (zero-shot)	0.968	0.970	0.933	0.948
	Azimuth	0.970	0.957	0.962	0.959
	expMap	0.920	0.940	0.846	0.878

Files

main + Q

Go to file 1

- > .vscode
- > data
- docs
- examples
- scgpt
- tests
- tutorials
 - zero-shot
 - Tutorial_Annotation.ipynb
 - Tutorial_Attention_GRN.ipynb
 - Tutorial_GRN.ipynb
 - Tutorial_Integration.ipynb
 - Tutorial_Multiomics.ipynb
 - Tutorial_Perturbation.ipynb
 - Tutorial_Reference_Mapping_L...
 - build_atlas_index_faiss.py
 - .gitignore
 - .readthedocs.yaml
 - LICENSE
 - README.md
 - poetry.lock
 - pyproject.toml

scGPT / tutorials /

Add file ...

subercul add colab badges 10:45:18 - 18 hours ago History

Name	Last commit message	Last commit date
..		
zero-shot	add colab badges	18 hours ago
Tutorial_Annotation.ipynb	update annotation tutorial	5 months ago
Tutorial_Attention_GRN.ipynb	add attention grn tutorial	5 months ago
Tutorial_GRN.ipynb	update doc with tutorial notebooks	6 months ago
Tutorial_Integration.ipynb	update load model in notebook	last month
Tutorial_Multiomics.ipynb	update multi omic tutorial	2 months ago
Tutorial_Perturbation.ipynb	update docs	3 months ago
Tutorial_Reference_Mapping.ipynb	update value counts for predictions, fix #135	20 hours ago
build_atlas_index_faiss.py	update build atlas index faiss	3 months ago



Zero-shot reference mapping tutorial with scGPT

Introduction

This tutorial covers the zero-shot reference mapping with scGPT. This workflow achieves accurate and fast reference mapping for scRNA-seq datasets without fine-tuning (or any extensive training) of scGPT. To further boost the performance, we recommend fine-tuning scGPT.

We will use COVID-19 dataset to demonstrate the zero-shot reference mapping. You can download the processed reference and query datasets from [here](#). The COVID-19 dataset is derived from the work by [Lotfollahi et al.](#), which contains 18 distinct batches and diverse samples from lung tissues. The reference dataset consists of 15,997 cells and the query dataset contains 4,003 cells. You may place the dataset under `data` directory in the outer level.

Particularly, we use the `scGPT_human` model to provide embeddings out of the box. You may download it from [here](#).

The zero-shot reference mapping workflow is as follows:

1. Load and pre-process the dataset
2. Generate scGPT embeddings for each cell in reference and query datasets
3. Transfer the annotations from reference to query dataset

At the [appendix](#) of this tutorial, we will also showcase the zero-shot reference mapping on Lung dataset. You can find the dataset [here](#).

We use a similarity-based method for transferring the annotation, which involves comparing the similarity between the query cell embedding and the reference cell embeddings. We use [FAISS](#) to perform the similarity search.

[Open in Colab](#)

```
In [1]: # Specifically for Google Colab, install dependencies and download data
import os
import sys

if "google.colab" in sys.modules:
    print("Running on Google Colab")
    print("Installing dependencies...")
    !pip install -U scgpt
    # the optional dependency of flash-attention is skipped on colab
    !pip install wandb louvain faiss-cpu

# NOTE: May need to restart runtime after the installation
```

Haotian Cui

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

14 Colab AI

```
Building directory structure completed
Retrieving folder list completed
Building directory structure
Downloading...
From: https://drive.google.com/uc?id=tp0tycRWjw0b70w001d08dxUre80Ma-
To: /data/lung/sample_proc_lung_test.h5ad
100% [██████████] 15.4M/15.4M [00:00<00:00, 66.4MB/s]
Downloading...
From: https://drive.google.com/uc?id=1M4wni18qvE4MAoanTSC@u131tth
To: /data/lung/sample_proc_lung_train.h5ad
100% [██████████] 60.5M/60.5M [00:01<00:00, 42.3MB/s] Downloading model ckpt...

Download completed
Retrieving folder list
Processing file 1hh2zGkYwAx3Dyov03RG5tZ301zmSgdk1 args.json
Processing file 14AebJfG0UF847Eg40hK57HC1rb0fy0Tm best_model.pt
Processing file 1H3E_MJ-0136A0V6jLbna2EdvgPaqvcC vocab.json
Building directory structure completed
Retrieving folder list completed
Building directory structure
Downloading...
From: https://drive.google.com/uc?id=1hh2zGkYwAx3Dyov03RG5tZ301zmSgdk1
To: /save/scGPT_human/args.json
100% [██████████] 1.30k/1.30k [00:00<00:00, 5.00MB/s]
Downloading...
From (original): https://drive.google.com/uc?id=14AebJfG0UF847Eg40hK57HC1rb0fy0Tm
From (redirected): https://drive.google.com/uc?id=14AebJfG0UF847Eg40hK57HC1rb0fy0Tm&confirm=t&uiid=9d1338e6-98ff-4466-b4fe-c32caeb306ef
To: /save/scGPT_human/best_model.pt
100% [██████████] 205M/205M [00:09<00:00, 22.3MB/s]
Downloading...
From: https://drive.google.com/uc?id=1H3E\_MJ-0136A0V6jLbna2EdvgPaqvcC
To: /save/scGPT_human/vocab.json
100% [██████████] 1.32M/1.32M [00:00<00:00, 145MB/s]
Download completed
```

Import scGPT and dependencies

```
[2] from pathlib import Path
import numpy as np
from scipy.stats import mode
import scanpy as sc
import warnings
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.metrics import confusion_matrix
```

0s completed at 12:03 PM



Haotian Cui

Google Chrome | Tutorial_ZeroShot_Reference_Mapping.ipynb | Colab AI

Zero-shot reference mapping tutorial with scGPT

Introduction

This tutorial covers the zero-shot reference mapping with scGPT. This workflow achieves accurate and fast reference mapping for scRNA-seq datasets without fine-tuning (or any extensive training) of scGPT. To further boost the performance, we recommend fine-tuning scGPT.

We will use COVID-19 dataset to demonstrate the zero-shot reference mapping. You can download the processed reference and query datasets from [here](#). The COVID-19 dataset is derived from the work by [Lotfollahi et al.](#), which contains 18 distinct batches and diverse samples from lung tissues. The reference dataset consists of 15,997 cells and the query dataset contains 4,003 cells. You may place the dataset under `data` directory in the outer level.

Particularly, we use the `scGPT_human` model to provide embeddings out of the box. You may download it from [here](#).

The zero-shot reference mapping workflow is as follows:

1. Load and pre-process the dataset
2. Generate scGPT embeddings for each cell in reference and query datasets
3. Transfer the annotations from reference to query dataset

At the [appendix](#) of this tutorial, we will also showcase the zero-shot reference mapping on Lung dataset. You can find the dataset [here](#).

We use a similarity-based method for transferring the annotation, which involves comparing the similarity between the query cell embedding and the reference cell embeddings. We use [FAISS](#) to perform the similarity search.

[Open in Colab](#)

```
# Specifically for Google Colab, install dependencies and download data
import os
import sys

if "google.colab" in sys.modules:
    print("Running on Google Colab")
    print("Installing dependencies...")
    !pip install -U scgpt
```

Os completed at 12:03 PM



Zero-shot reference mapping tutorial with scGPT

Introduction

This tutorial covers the zero-shot reference mapping with scGPT. This workflow achieves reference mapping on new query datasets without fine-tuning (or any extensive training) of scGPT. To further boost the

We will use COVID-19 dataset to demonstrate the zero-shot reference mapping. You can find the dataset from [here](#). The COVID-19 dataset is derived from the work by [Lotfollahi et al.](#), which covers various human tissues. The reference dataset consists of 15,997 cells and the query dataset contains 100 cells. Both datasets are located in the directory in the outer level.

Particularly, we use the `scGPT_human` model to provide embeddings out of the box. You can find the model in the `scGPT` directory.

The zero-shot reference mapping workflow is as follows:

1. Load and pre-process the dataset
2. Generate scGPT embeddings for each cell in reference and query datasets
3. Transfer the annotations from reference to query dataset

At the [appendix](#) of this tutorial, we will also showcase the zero-shot reference mapping on a new query dataset. We use a similarity-based method for transferring the annotation, which involves computing the similarity between the reference cell embeddings and the reference cell embeddings. We use [FAISS](#) to perform the similarity search.

[Open in Colab](#)

```
# Specifically for Google Colab, install dependencies and download data
```

```
import os
import sys
```

```
if "google.colab" in sys.modules:
    print("Running on Google Colab")
    print("Installing dependencies...")
    !pip install -U scgpt
```

Change runtime type

Runtime type

Python 3

Hardware accelerator



CPU



T4 GPU



A100 GPU



V100 GPU



TPU

Want access to premium GPUs? [Purchase additional compute units](#)

Cancel

Save



Haotian Cui

Zero-shot reference mapping tutorial with scGPT

Introduction

This tutorial covers the zero-shot reference mapping with scGPT. This workflow achieves accurate and fast reference mapping for scRNA-seq datasets without fine-tuning (or any extensive training) of scGPT. To further boost the performance, we recommend fine-tuning scGPT.

We will use COVID-19 dataset to demonstrate the zero-shot reference mapping. You can download the processed reference and query datasets from [here](#). The COVID-19 dataset is derived from the work by [Lotfollahi et al.](#) which contains 18 distinct batches and diverse samples from lung tissues. The reference dataset consists of 15,997 cells and the query dataset contains 4,003 cells. You may place the dataset under `data` directory in the outer level.

Particularly, we use the `scGPT_human` model to provide embeddings out of the box. You may download it from [here](#).

The zero-shot reference mapping workflow is as follows:

1. Load and pre-process the dataset
2. Generate scGPT embeddings for each cell in reference and query datasets
3. Transfer the annotations from reference to query dataset

At the [appendix](#) of this tutorial, we will also showcase the zero-shot reference mapping on Lung dataset. You can find the dataset [here](#).

We use a similarity-based method for transferring the annotation, which involves comparing the similarity between the query cell embedding and the reference cell embeddings. We use [FAISS](#) to perform the similarity search.

[Open in Colab](#)

Specifically for Google Colab, install dependencies and download data

```
import os
import sys

if "google.colab" in sys.modules:
    print("Running on Google Colab")
    print("Installing dependencies...")
    !pip install -U scgpt
    # the optional dependency of flash-attention is skipped on colab
    !pip install wandb louvain faiss-cpu
```



Haotian Cu

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

3. Transfer the annotations from reference to query dataset

At the [appendix](#) of this tutorial, we will also showcase the zero-shot reference mapping on Lung dataset. You can find the dataset [here](#).

We use a similarity-based method for transferring the annotation, which involves comparing the similarity between the query cell embedding and the reference cell embeddings. We use [FAISS](#) to perform the similarity search.

[Open in Colab](#)

Specifically for Google Colab, install dependencies and download data

```
import os
import sys

if "google.colab" in sys.modules:
    print("Running on Google Colab")
    print("Installing dependencies...")
    !pip install -U scgpt
    # the optional dependency of flash-attention is skipped on colab
    !pip install wandb louvain faiss-cpu

    # NOTE: May need to restart runtime after the installation

    print("Downloading data and model ckpt...")
    !pip install -q -U gdown
    import gdown

    data_dir = "../../data"
    if not os.path.exists(data_dir):
        os.mkdir(data_dir)
    if not os.path.exists(os.path.join(data_dir, "covid")):
        gdown.download_folder(
            "https://drive.google.com/drive/folders/1jSpPunG00ed71vDsk8FS7UvndHGshQ5",
            output=os.path.join(data_dir),
        )
    if not os.path.exists(os.path.join(data_dir, "lung")):
        gdown.download_folder(
            "https://drive.google.com/drive/folders/1gbf07VqxCOkfzqHA1h6h088zFv6p8u0",
            output=os.path.join(data_dir),
        )

    print("Downloading model ckpt...")
    model_dir = "../../save/scGPT_human"
```



Haotian Cui

Colab AI



Haotian Cui

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive

```
print("Downloading data and model ckpt...")
!pip install -q -U gdown
!import gdown

data_dir = "../data"
if not os.path.exists(data_dir):
    os.mkdir(data_dir)
if not os.path.exists(os.path.join(data_dir, "covid")):
    gdown.download_folder(
        "https://drive.google.com/drive/folders/1jSPaPun5Q0md71vDsK8F5TjvndHGshQ5",
        output=os.path.join(data_dir,
    )
if not os.path.exists(os.path.join(data_dir, "lung")):
    gdown.download_folder(
        "https://drive.google.com/drive/folders/1gbF07VqxCKkfzghAljh088zfv6pd8wD",
        output=os.path.join(data_dir,
    )

print("Downloading model ckpt...")
model_dir = "../save/scGPT_human"
if not os.path.exists(model_dir):
    !mkdir -p $model_dir
    gdown.download_folder(
        "https://drive.google.com/drive/folders/1oWh_-ZR8htoG02Fw24HP41f9LoomVo-y",
        output=model_dir,
    )
```

```
Requirement already satisfied: aiosignal==1.2 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets<3.0.0,=>2.3.0->scgpt) (1.3.1)
Requirement already satisfied: async-timeout<5.0,=>4.0 in /usr/local/lib/python3.10/dist-packages (from aiohttp->datasets<3.0.0,=>2.3.0->scgpt) (4.0.3)
Requirement already satisfied: texttable==1.6.2 in /usr/local/lib/python3.10/dist-packages (from igraph<0.11,=>0.10.0->leidenalg==0.8.10->scgpt) (1.7.0)
Requirement already satisfied: ml-types==0.2.0 in /usr/local/lib/python3.10/dist-packages (from jax==0.4.6->orbx<0.1.0->scgpt) (0.2.0)
Requirement already satisfied: opt-einsum in /usr/local/lib/python3.10/dist-packages (from jax==0.4.6->orbx<0.1.0->scgpt) (3.3.0)
Requirement already satisfied: contourpy==1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib==3.4->scanpy<2.0.0,=>1.9.1->scgpt) (1.2.0)
Requirement already satisfied: cyclus==0.10 in /usr/local/lib/python3.10/dist-packages (from matplotlib==3.4->scanpy<2.0.0,=>1.9.1->scgpt) (0.12.1)
Requirement already satisfied: fonttools==4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib==3.4->scanpy<2.0.0,=>1.9.1->scgpt) (4.46.0)
Requirement already satisfied: kiwisolver==1.0.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib==3.4->scanpy<2.0.0,=>1.9.1->scgpt) (1.4.5)
Requirement already satisfied: pillow==6.2.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib==3.4->scanpy<2.0.0,=>1.9.1->scgpt) (9.4.0)
Requirement already satisfied: pyparsing==2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib==3.4->scanpy<2.0.0,=>1.9.1->scgpt) (3.1.1)
Requirement already satisfied: contextlib2 in /usr/local/lib/python3.10/dist-packages (from ml-collections==0.1.1->scvi-tools<1.0,=>0.16.0->scgpt) (21.6.0)
Requirement already satisfied: et-xmlfile in /usr/local/lib/python3.10/dist-packages (from openpyxl==3.0->scvi-tools<1.0,=>0.16.0->scgpt) (1.1.0)
Requirement already satisfied: pyro-api==0.1.1 in /usr/local/lib/python3.10/dist-packages (from pyro-ppl==1.6.0->scvi-tools<1.0,=>0.16.0->scgpt) (0.1.2)
Requirement already satisfied: lightning-utilities==0.6.0.post0 in /usr/local/lib/python3.10/dist-packages (from pytorch-lightning<1.10.0,=>1.9.0->scvi-tools<1.0,=>0.16.0->scgpt) (0.10.0)
Requirement already satisfied: charset-normalizer<4,=>2 in /usr/local/lib/python3.10/dist-packages (from requests==2.19.0->datasets<3.0.0,=>2.3.0->scgpt) (3.3.2)
Requirement already satisfied: idna==2.5 in /usr/local/lib/python3.10/dist-packages (from requests==2.19.0->datasets<3.0.0,=>2.3.0->scgpt) (3.6)
Requirement already satisfied: certifi==2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests==2.19.0->datasets<3.0.0,=>2.3.0->scgpt) (2023.11.17)
Requirement already satisfied: markdown-it-py==2.2.0 in /usr/local/lib/python3.10/dist-packages (from rich==12.0.0->scvi-tools<1.0,=>0.16.0->scgpt) (3.0.0)
Requirement already satisfied: oyaml==3.0.0,=>2.13.0 in /usr/local/lib/python3.10/dist-packages (from rich==12.0.0->scvi-tools<1.0,=>0.16.0->scgpt) (2.16.1)
```

✓ 0s completed at 12:29 PM

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive

T4 Colab AI

```
[18] Requirement already satisfied: urllib3<3,=>1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests<3,=>2.0.0->wandb) (2.0.7)
Requirement already satisfied: certifi<=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests<3,=>2.0.0->wandb) (2023.11.17)
Requirement already satisfied: smmap<6,=>3.0.1 in /usr/local/lib/python3.10/dist-packages (from gitdb<5,=>4.0.1->GitPython!<3.1.29,=>1.0.0->wandb) (5.0.1)
Downloading data and model ckpt...
Downloading model ckpt...
```

Import scGPT and dependencies

```
from pathlib import Path
import numpy as np
from scipy.stats import mode
import scanpy as sc
import warnings
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
from sklearn.metrics import confusion_matrix
import seaborn as sns
import pandas as pd
import sys

sys.path.insert(0, "..")

import scgpt as scg

# extra dependency for similarity search
try:
    import faiss

    faiss_imported = True
except ImportError:
    faiss_imported = False
    print(
        "faiss not installed! We highly recommend installing it for fast similarity search."
    )
print("To install it, see https://github.com/facebookresearch/faiss/wiki/Installing-Faiss")

warnings.filterwarnings("ignore", category=ResourceWarning)
```



Google Chrome File Edit View History Bookmarks Profiles Tab Window Help

colinrusso.github.io/scGPT/tutorial/zero-shot/tutorial_zero_shot_reference_mapping.ipynb

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

[19]

Load the dataset, you may download the dataset from [here](#). We set the columns storing gene name columns, batch key and cell type key.

```
model_dir = Path("../save/scGPT_human")
adata = sc.read_h5ad("../data/covid/batch_covid_subsampled_train.h5ad")
test_adata = sc.read_h5ad("../data/covid/batch_covid_subsampled_test.h5ad")

cell_type_key = "celltype"
gene_col = "gene_name"
```

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' or should_run_async(code)

Embed the reference dataset

```
!ref_embed_adata = scg.tasks.embed_data(
    adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)
```

scGPT - INFO - match 1173/1200 genes in vocabulary of size 60697.
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' or should_run_async(code)

0s completed at 12:29 PM





```
[19]
Load the dataset, you may download the dataset from here. We set the columns storing gene name columns, batch key and cell type key.

model_dir = Path("../data/covid/batch_covid_subsamped_train.h5ad")
adata = sc.read_h5ad(model_dir)
test_adata = sc.read_h5ad("../data/covid/batch_covid_subsamped_test.h5ad")

cell_type_key = "celltype"
gene_col = "gene_name"

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' or and should_run_async(code)

Embed the reference dataset

!ref_embed_adata = scg.tasks.embed_data(
    adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)

scGPT - INFO - match 1173/1200 genes in vocabulary of size 60697.
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' or and should_run_async(code)
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-warnings.warn[
Embedding cells: 94% |██████████| 236/250 [00:32<00:01, 7.38it/s]
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive

[19]

Load the dataset, you may download the dataset from [here](#). We set the columns storing gene name columns, batch key and cell type key.

```
model_dir = Path("../save/scGPT_human")
adata = sc.read_h5ad("../data/covid/batch_covid_subsampled_train.h5ad")
test_adata = sc.read_h5ad("../data/covid/batch_covid_subsampled_test.h5ad")

cell_type_key = "celltype"
gene_col = "gene_name"
```

`/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)`

Embed the reference dataset

```
[21] ref_embed_adata = scg.tasks.embed_data(
    adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)
```

```
scGPT - INFO - match 1173/1200 genes in vocabulary of size 60697.
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-warnings.warn
Embedding cells: 100% [██████████] 250/250 [00:34<00:00, 7.29it/s]
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element '.obsn['X_scGPT']' of view, Initializing view as actual.
adata.obsn["X_scGPT"] = cell_embeddings
```



Haotian Cui

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

Colab AI

```
model_dir = Path("../save/scGPT_human")
adata = sc.read_h5ad("../data/covid/batch_covid_subsampled_train.h5ad")
test_adata = sc.read_h5ad("../data/covid/batch_covid_subsampled_test.h5ad")

cell_type_key = "celltype"
gene_col = "gene_name"
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' or and should_run_async(code)
```

Embed the reference dataset

```
ref_embed_adata = scg.tasks.embed_data(
    adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)
```

```
scGPT - INFO - match 1173/1200 genes in vocabulary of size 68697.
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' or and should_run_async(code)
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-warnings.warn()
Embedding cells: 100% |██████████| 250/250 [00:34<00:00, 7.29it/s]
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element '.obsm["X_scGPT"]' of view, initializing view as actual.
adata.obsm["X_scGPT"] = cell_embeddings
```

Embed the query dataset

```
[22] test_embed_adata = scg.tasks.embed_data(
    test_adata,
    model_dir
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

```
[20] cell_type_key = "celltype"  
     gene_col = "gene_name"  
  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)
```

Embed the reference dataset

```
ref_embed_adata = scg.tasks.embed_data(  
    adata,  
    model_dir,  
    gene_col=gene_col,  
    batch_size=64,  
)  
  
scgPT - INFO - match 1173/1200 genes in vocabulary of size 60697.  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)  
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-attn.  
Embedding cells: 100% [██████████] 250/250 [00:34<00:00, 7.29it/s]  
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element '.obsn['X_scgPT']' of view, initializing view as actual.  
adata.obsn["X_scgPT"] = cell_embeddings
```

Embed the query dataset

```
[22] test_embed_adata = scg.tasks.embed_data(  
    test_adata,  
    model_dir,  
    gene_col=gene_col,  
    batch_size=64,  
)
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

Colab AI

```
scGPT - INFO - match 1173/1200 genes in vocabulary of size 60697.  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar  
and should_run_async(code)  
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-  
warnings.warn(  
Embedding cells: 100% [██████████] 250/250 [00:34<00:00, 7.29it/s]  
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element '.obsn['X_scGPT']' of view, Initializing view as actual.  
adata.obsn["X_scGPT"] = cell_embeddings
```

Embed the query dataset

```
test_embed_adata = scg.tasks.embed_data(  
    test_adata,  
    model_dir,  
    gene_col=gene_col,  
    batch_size=64,  
)  
  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar  
and should_run_async(code)  
scGPT - INFO - match 1173/1200 genes in vocabulary of size 60697.  
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-  
warnings.warn(  
Embedding cells: 100% [██████████] 63/63 [00:09<00:00, 6.86it/s]  
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element '.obsn['X_scGPT']' of view, Initializing view as actual.  
adata.obsn["X_scGPT"] = cell_embeddings
```

Mark the reference vs. query cells and mask the cell types on query cells

```
[23] # concatenate the two datasets  
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")  
# mark the reference vs. query dataset  
adata_concat.obs["dataset"] = np.zeros(adata_concat.n_obs) + 1  
adata_concat.obs["dataset"][test_embed_adata.n_obs:] = 2
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

```
+ Code + Text Copy to Drive  
[22] and should_run_async(code)  
scgpt - INFO - match 1173/1200 genes in vocabulary of size 60697.  
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-warnings.warn  
Embedding cells: 100% [██████████] 63/63 [00:09<00:00, 5.86it/s]  
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element 'obs['X_scgpt']' of view, initializing view as actual.  
adata.obs['X_scgpt'] = cell_embeddings
```

Mark the reference vs. query cells and mask the cell types on query cells

```
# concatenate the two datasets  
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")  
# mark the reference vs. query dataset  
adata_concat.obs["is_ref"] = ["Query"] * len(test_embed_adata) + ["Reference"] * len(ref_embed_adata)  
adata_concat.obs["is_ref"] = adata_concat.obs["is_ref"].astype("category")  
# mask the query dataset cell types  
adata_concat.obs[cell_type_key] = adata_concat.obs[cell_type_key].astype("category")  
adata_concat.obs[cell_type_key] = adata_concat.obs[cell_type_key].cat.add_categories(["To be predicted"])  
adata_concat.obs[cell_type_key][: len(test_embed_adata)] = "To be predicted"  
  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)  
<ipython-input-23-e57926d27f57>:2: FutureWarning: Use adata_concat instead of AnnData.concatenate, AnnData.concatenate is deprecated and will be removed in the future. See the tutorial for concat at: https://anndata.readthedocs.io/en/latest/concat.html  
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")
```

+ Code + Text

Visualize the embeddings

We visualize the embeddings from query and reference datasets using UMAP

```
[24] sc.pp.neighbors(adata_concat, use_rep="X_scgpt")  
sc.tl.umap(adata_concat)  
sc.pl.umap(adata_concat, color="is_ref", cell_type_key, wspace=0.4, frameon=False, ncols=1)  
  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

Colab AI

```
[22] /usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-
warnings.warn(
Embedding cells: 100% [██████████] 63/63 [00:09<00:00, 6.86it/s]
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element '.obsn['X_scGPT']' of view, initializing view as actual.
adata.obsn['X_scGPT'] = cell_embeddings
```

Mark the reference vs. query cells and mask the cell types on query cells

```
# concatenate the two datasets
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")
# mark the reference vs. query dataset
adata_concat.obs["is_ref"] = ["Query"] * len(test_embed_adata) + ["Reference"] * len(
    ref_embed_adata
)
adata_concat.obs["is_ref"] = adata_concat.obs["is_ref"].astype("category")
# mask the query dataset cell types
adata_concat.obs[cell_type_key] = adata_concat.obs[cell_type_key].astype("category")
adata_concat.obs[cell_type_key] = adata_concat.obs[cell_type_key].cat.add_categories(["To be predicted"])
adata_concat.obs[cell_type_key][: len(test_embed_adata)] = "To be predicted"
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar
and should_run_async(code)
<ipython-input-23-e57926427f57>:2: FutureWarning: Use adata.concatenate instead of AnnData.concatenate, AnnData.concatenate is deprecated and will be removed in the future. See the tutorial for concat at: https://ann
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")
```

Visualize the embeddings

We visualize the embeddings from query and reference datasets using UMAP

```
[24] sc.pp.neighbors(adata_concat, use_rep="X_scGPT")
sc.tl.umap(adata_concat)
sc.pl.umap(
    adata_concat, color=["is_ref", cell_type_key], wspace=0.4, frameon=False, ncols=1
)
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar
and should_run_async(code)
/usr/local/lib/python3.10/dist-packages/scanpy/plotting/_tools/scatterplots.py:394: UserWarning: No data for colormapping provided via 'c'. Parameters 'cmap' will be ignored
rax = scatterpl
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

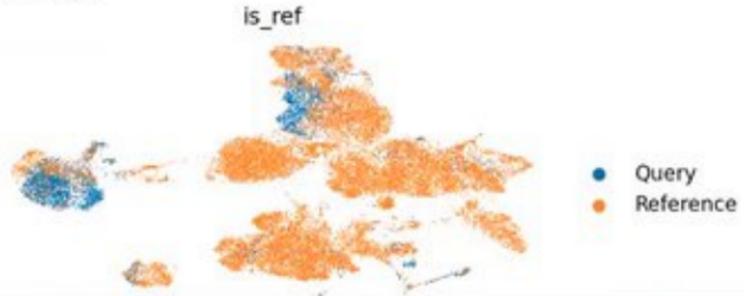
File Edit View Insert Runtime Tools Help Cannot save changes

```
+ Code + Text Copy to Drive  
adata_concat.obs[cell_type_key] = adata_concat.obs[cell_type_key].cat.add_categories(["To be predicted"])  
adata_concat.obs[cell_type_key][: len(test_embed_adata)] = "To be predicted"  
  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)  
<ipython-input-23-e57926d27f57>:2: FutureWarning: Use adata_concat instead of AnnData.concatenate, AnnData.concatenate is deprecated and will be removed in the future. See the tutorial for concat at: https://an  
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")
```

Visualize the embeddings

We visualize the embeddings from query and reference datasets using UMAP

```
[24] sc.pp.neighbors(adata_concat, use_rep="X_scGPT")  
sc.tl.umap(adata_concat)  
sc.pl.umap(  
    adata_concat, color="is_ref", cell_type_key, wspace=0.4, frameon=False, ncols=1  
)  
  
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)  
/usr/local/lib/python3.10/dist-packages/scanpy/plotting/_tools/scatterplots.py:394: UserWarning: No data for colormapping provided via 'c'. Parameters 'cmap' will be ignored  
cax = scatter(  
/usr/local/lib/python3.10/dist-packages/scanpy/plotting/_tools/scatterplots.py:394: UserWarning: No data for colormapping provided via 'c'. Parameters 'cmap' will be ignored  
cax = scatter(  
  
is_ref
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive



We visualize the embeddings from query and reference datasets using UMAP

```

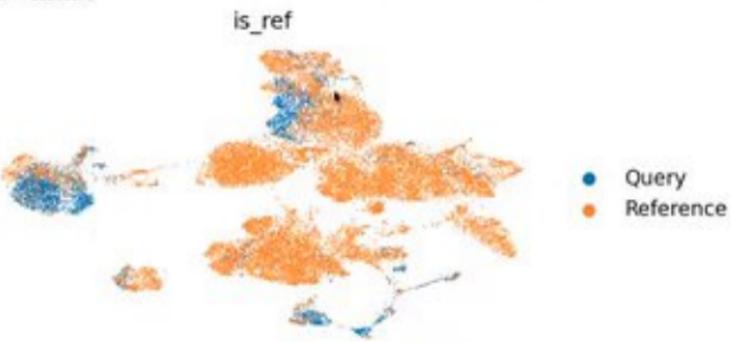
sc.pp.neighbors(adata_concat, use_rep="X_scGPT")
sc.tl.umap(adata_concat)
sc.pl.umap(
    adata_concat, color="is_ref", cell_type_key, wspace=0.4, frameon=False, ncols=1
)

```

```

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar
and should_run_async(code)
/usr/local/lib/python3.10/dist-packages/scrapy/plotting/_tools/scatterplots.py:394: UserWarning: No data for colormapping provided via 'c'. Parameters 'cmap' will be ignored
cax = scatter(
/usr/local/lib/python3.10/dist-packages/scrapy/plotting/_tools/scatterplots.py:394: UserWarning: No data for colormapping provided via 'c'. Parameters 'cmap' will be ignored
cax = scatter(

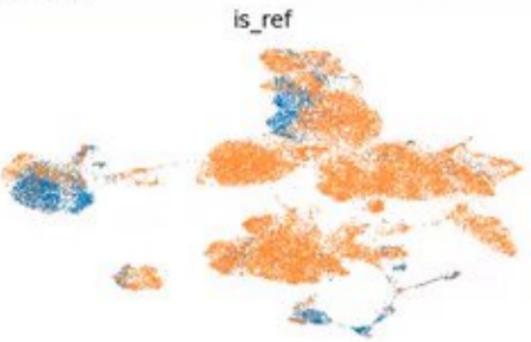
```



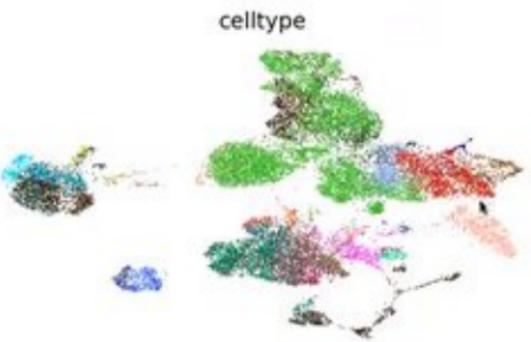


Haotian Cui

```
!cax = scatter!  
/usr/local/lib/python3.10/dist-packages/scanpy/plotting/_tools/scatterplots.py:394: UserWarning: No data for colormapping provided via 'c'. Parameters 'cmap' will be ignored  
!cax = scatter!
```



- Query
- Reference



- AT1
- AT2
- B cell
- CCR7+ T
- CD4+ T cells
- CD8 T
- CD8+ T cells
- CD10+ B cells
- CD14+ Monocytes
- CD16+ Monocytes
- CD20+ B cells
- Ciliated
- DC activated
- Erythrocytes
- Erythroid progenitors
- HSPCs
- IGSF21+ Dendritic
- M2 Macrophage
- Macrophages
- Mast cells
- Megakaryocyte progenitors
- Megakaryocytes
- Monocyte progenitors
- Monocytes
- NK
- NK cells
- Neutrophil
- Plasma
- Proliferating Macrophage
- Proliferating T
- Secretory
- Signaling Alveolar Epithelial Type 2
- T cell
- Treg
- Tregs
- innate T
- mDC
- pDC
- To be predicted

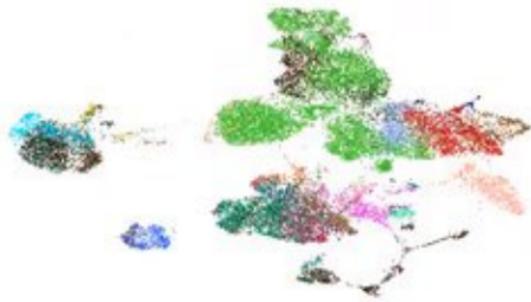
```
/usr/local/lib/python3.10/dist-packages/scannpy/plotting/_tools/scatterplots.py:394: UserWarning: No data for colormapping provided via "c". Parameters "cmap" will be ignored  
cax = scatter()
```

is_ref



- Query
- Reference

celltype

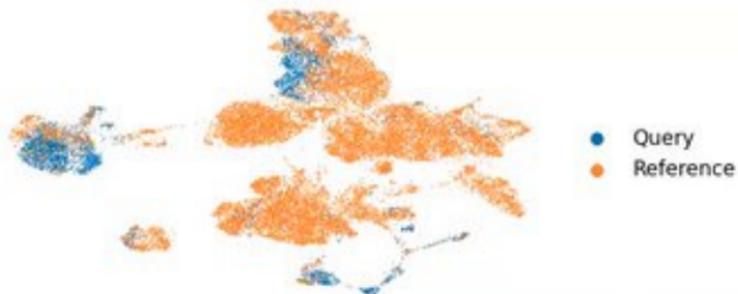


- | | | |
|-------------------|-----------------------------|--|
| ● AT1 | ● Erythrocytes | ● Neutrophil |
| ● AT2 | ● Erythroid progenitors | ● Plasma |
| ● B cell | ● HSPCs | ● Proliferating Macrophage |
| ● CCR7+ T | ● IGSF21+ Dendritic | ● Proliferating T |
| ● CD4+ T cells | ● M2 Macrophage | ● Secretory |
| ● CD8 T | ● Macrophages | ● Signaling Alveolar Epithelial Type 2 |
| ● CD8+ T cells | ● Mast cells | ● T cell |
| ● CD10+ B cells | ● Megakaryocyte progenitors | ● Treg |
| ● CD14+ Monocytes | ● Megakaryocytes | ● Tregs |
| ● CD16+ Monocytes | ● Monocyte progenitors | ● innate T |
| ● CD20+ B cells | ● Monocytes | ● mDC |
| ● Ciliated | ● NK | ● pDC |
| ● DC_activated | ● NKT cells | ● To be predicted |

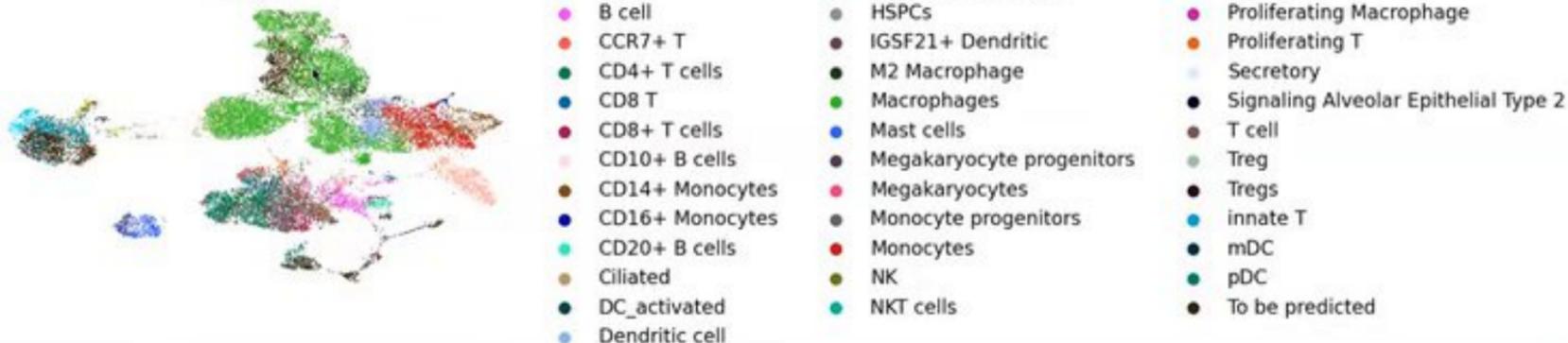


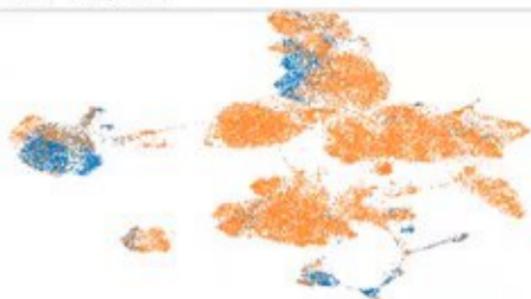
Haotian Cui

is_ref

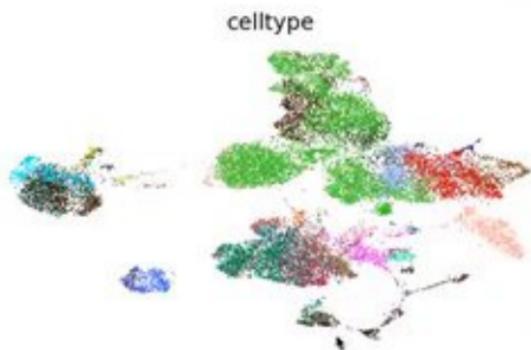


celltype



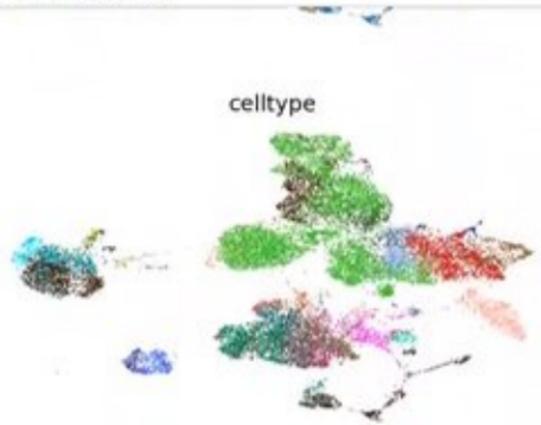


● Query
● Reference



celltype

- AT1
- AT2
- B cell
- CCR7+ T
- CD4+ T cells
- CD8 T
- CD8+ T cells
- CD10+ B cells
- CD14+ Monocytes
- CD16+ Monocytes
- CD20+ B cells
- Ciliated
- DC_activated
- Dendritic cell
- Erythrocytes
- Erythroid progenitors
- HSPCs
- IGSF21+ Dendritic
- M2 Macrophage
- Macrophages
- Mast cells
- Megakaryocyte progenitors
- Megakaryocytes
- Monocyte progenitors
- Monocytes
- NK
- NKT cells
- Neutrophil
- Plasma
- Proliferating Macrophage
- Proliferating T
- Secretory
- Signaling Alveolar Epithelial Type 2
- T cell
- Treg
- Tregs
- innate T
- mDC
- pDC
- To be predicted



- | | | |
|-------------------|-----------------------------|--|
| ● AT1 | ● Erythrocytes | ● Neutrophil |
| ● AT2 | ● Erythroid progenitors | ● Plasma |
| ● B cell | ● HSPCs | ● Proliferating Macrophage |
| ● CCR7+ T | ● IGSF21+ Dendritic | ● Proliferating T |
| ● CD4+ T cells | ● M2 Macrophage | ● Secretory |
| ● CD8 T | ● Macrophages | ● Signaling Alveolar Epithelial Type 2 |
| ● CD8+ T cells | ● Mast cells | ● T cell |
| ● CD10+ B cells | ● Megakaryocyte progenitors | ● Treg |
| ● CD14+ Monocytes | ● Megakaryocytes | ● Tregs |
| ● CD16+ Monocytes | ● Monocyte progenitors | ● innate T |
| ● CD20+ B cells | ● Monocytes | ● mDC |
| ● Ciliated | ● NK | ● pDC |
| ● DC_activated | ● NKT cells | ● To be predicted |
| ● Dendritic cell | | |

Reference mapping and transfer the annotations

We run the reference mapping using cell-level majority voting. You may adjust the `k` parameter to control the number of nearest neighbors to consider for voting.

```
[25] # Those functions are only used when faiss is not installed
def l2_sim(a, b):
    sims = -np.linalg.norm(a - b, axis=1)
    return sims

def get_similar_vectors(vector, ref, top_k=10):
```

Google Chrome File Edit View History Bookmarks Files Tab Window Help

colab:usar2f...@github.com/github/bowang-lab/scGPT/blob/main/tutorials/zero-shot/Tutorial_ZeroShot_Reference_Mapping.ipynb#scrollTo=scLk5jPKKGecw

Tutorial_ZeroShot_Reference_Mapping.ipynb
File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive

Haotian Cui

- Ciliated
- DC_activated
- Dendritic cell
- NK
- NKT cells
- pDC
- To be predicted

Reference mapping and transfer the annotations

We run the reference mapping using cell-level majority voting. You may adjust the `k` parameter to control the number of nearest neighbors to consider for voting.

Those functions are only used when faiss is not installed

```
def l2_sim(a, b):
    sims = -np.linalg.norm(a - b, axis=1)
    return sims

def get_similar_vectors(vector, ref, top_k=10):
    # sims = cos_sim(vector, ref)
    sims = l2_sim(vector, ref)

    top_k_idx = np.argsort(sims)[::-1][:top_k]
    return top_k_idx, sims[top_k_idx]
```

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar and should_run_async(code)

```
[26] ref_cell_embeddings = ref_embed_adata.obsn["X_scGPT"]
test_embd = test_embed_adata.obsn["X_scGPT"]

k = 10 # number of neighbors

index = faiss.IndexFlatL2(ref_cell_embeddings.shape[1])
index.add(ref_cell_embeddings)

# Query dataset, k - number of closest elements (returns 2 numpy arrays)
distances, labels = index.query(test_embd, k)
```

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

T4 Colab AI

```
[25] sims = -np.linalg.norm(a - b, axis=1)
      return sims
```

```
def get_similar_vectors(vector, ref, top_k=10):
    # sims = cos_sim(vector, ref)
    sims = l2_sim(vector, ref)

    top_k_idx = np.argsort(sims)[::-1][:top_k]
    return top_k_idx, sims[top_k_idx]
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' or should_run_async(code)
```

```
ref_cell_embeddings = ref_embed_adata.obs["X_scGPT"]
test_embd = test_embed_adata.obs["X_scGPT"]
```

```
k = 10 # number of neighbors
```

```
index = faiss.IndexFlatL2(ref_cell_embeddings.shape[1])
index.add(ref_cell_embeddings)
```

```
# Query dataset, k - number of closest elements (returns 2 numpy arrays)
distances, labels = index.search(test_embd, k)
```

```
idx_list=[i for i in range(test_embd.shape[0])]
```

```
preds = []
```

```
sim_list = distances
```

```
for k in idx_list:
```

```
    if faiss_imported:
```

```
        idx = labels[k]
```

```
    else:
```

```
        idx, sim = get_similar_vectors(test_embd[k][np.newaxis, ...], ref_cell_embeddings, k)
```

```
        pred = ref_embed_adata.obs[cell_type_key][idx].value_counts()
```

```
        preds.append(pred.index[0])
```

```
gt = test_adata.obs[cell_type_key].to_numpy()
```

Evaluate the performance

```
[27] res_dict = {
```



Haotian Cui



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

+ Code + Text Copy to Drive

```
[25] top_k_idx = np.argsort(sims)[::-1][:top_k]
      return top_k_idx, sims[top_k_idx]

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar
and should_run_async(code)

ref_cell_embeddings = ref_embed_adata.obsn["X_scGPT"]
test_embd = test_embed_adata.obsn["X_scGPT"]

k = 10 # number of neighbors

index = faiss.IndexFlatL2(ref_cell_embeddings.shape[1])
index.add(ref_cell_embeddings)

# Query dataset, k - number of closest elements (returns 2 numpy arrays)
distances, labels = index.search(test_embd, k)

idx_list=[i for i in range(test_embd.shape[0])]
preds = []
sim_list = distances
for k in idx_list:
    if faiss_imported:
        idx = labels[k]
    else:
        idx, sim = get_similar_vectors(test_embd[k][np.newaxis, ...], ref_cell_embeddings, k)
    pred = ref_embed_adata.obs[cell_type_key][idx].value_counts()
    preds.append(pred.index[0])
gt = test_adata.obs[cell_type_key].to_numpy()
```

▼ Evaluate the performance

+ Code + Text

```
[27] res_dict = {
      "accuracy": accuracy_score(gt, preds),
      "precision": precision_score(gt, preds, average="macro"),
      "recall": recall_score(gt, preds, average="macro"),
      "macro_f1": f1_score(gt, preds, average="macro"),
    }

res_dict
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

```
Code + Text Copy to Drive
preds = []
sim_list = distances
for k in idx_list:
    if faiss_imported:
        idx = labels[k]
    else:
        idx, sim = get_similar_vectors(test_embd[k][np.newaxis, ...], ref_cell_embeddings, k)
    pred = ref_embed_adata.obs[cell_type_key][idx].value_counts()
    preds.append(pred.index[0])
gt = test_adata.obs[cell_type_key].to_numpy()
```

Evaluate the performance

```
res_dict = {
    "accuracy": accuracy_score(gt, preds),
    "precision": precision_score(gt, preds, average="macro"),
    "recall": recall_score(gt, preds, average="macro"),
    "macro_f1": f1_score(gt, preds, average="macro"),
}
res_dict
```

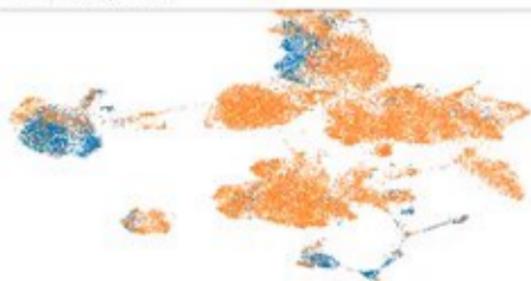
Warning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control.

Warning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to control.

```
{'accuracy': 0.8748944291781164,
 'precision': 0.5274756111472559,
 'recall': 0.48119880193882796,
 'macro_f1': 0.46823887224181864}
```

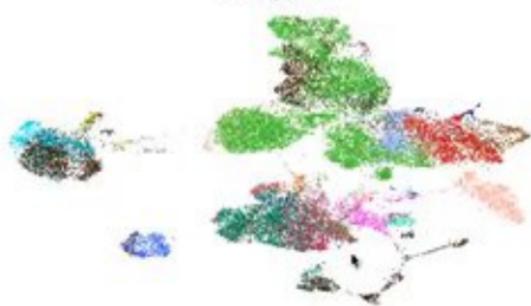
Confusion matrix:

```
[28] y_true = gt
y_pred = preds
cell_type_list = np.unique(y_true)
matrix = confusion_matrix(y_true, y_pred, labels=cell_type_list)
matrix = matrix.astype("float") / matrix.sum(axis=1)[:, np.newaxis]
```



● Query
● Reference

celltype



- | | | |
|-------------------|-----------------------------|--|
| ● AT1 | ● Erythrocytes | ● Neutrophil |
| ● AT2 | ● Erythroid progenitors | ● Plasma |
| ● B cell | ● HSPCs | ● Proliferating Macrophage |
| ● CCR7+ T | ● IGSF21+ Dendritic | ● Proliferating T |
| ● CD4+ T cells | ● M2 Macrophage | ● Secretory |
| ● CD8 T | ● Macrophages | ● Signaling Alveolar Epithelial Type 2 |
| ● CD8+ T cells | ● Mast cells | ● T cell |
| ● CD10+ B cells | ● Megakaryocyte progenitors | ● Treg |
| ● CD14+ Monocytes | ● Megakaryocytes | ● Tregs |
| ● CD16+ Monocytes | ● Monocyte progenitors | ● innate T |
| ● CD20+ B cells | ● Monocytes | ● mDC |
| ● Ciliated | ● NK | ● pDC |
| ● DC_activated | ● NKT cells | ● To be predicted |
| ● Dendritic cell | | |



Haotian Cui

Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

```
Code Text Copy to Drive  
distances, labels = index.search(test_embd, k)  
idx_list=[i for i in range(test_embd.shape[0])]  
preds = []  
sim_list = distances  
for k in idx_list:  
    if faiss_imported:  
        idx = labels[k]  
    else:  
        idx, sim = get_similar_vectors(test_embd[k][np.newaxis, ...], ref_cell_embeddings, k)  
        pred = ref_embed_adata.obs[cell_type_key][idx].value_counts()  
        preds.append(pred.index[0])  
gt = test_adata.obs[cell_type_key].to_numpy()
```

Evaluate the performance

```
res_dict = {  
    "accuracy": accuracy_score(gt, preds),  
    "precision": precision_score(gt, preds, average="macro"),  
    "recall": recall_score(gt, preds, average="macro"),  
    "macro_f1": f1_score(gt, preds, average="macro"),  
}
```

res_dict

```
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' param  
_warn_prf(average, modifier, msg_start, len(result))  
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to  
_warn_prf(average, modifier, msg_start, len(result))  
{  
  "accuracy": 0.8748944291781164,  
  "precision": 0.527475611472559,  
  "recall": 0.48119888193882796,  
  "macro_f1": 0.46823887224181864}
```

Confusion matrix:

```
[28]: y_true = gt  
y_pred = preds  
cell_type_list = np.unique(y_true)  
matrix = confusion_matrix(y_true, y_pred, labels=cell_type_list)
```



Haotian Cu

```
else:
    idx, sim = get_similar_vectors(test_embd[k][np.newaxis, ...], ref_cell_embeddings, k)
    pred = ref_embed_adata.obs[cell_type_key][idx].value_counts()
    preds.append(pred.index[0])
gt = test_adata.obs[cell_type_key].to_numpy()
```

Evaluate the performance

```
res_dict = {
    "accuracy": accuracy_score(gt, preds),
    "precision": precision_score(gt, preds, average="macro"),
    "recall": recall_score(gt, preds, average="macro"),
    "macro_f1": f1_score(gt, preds, average="macro"),
}

res_dict

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control.
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to control.
{'accuracy': 0.8748944291781164,
 'precision': 0.5274756111472559,
 'recall': 0.48119888193882796,
 'macro_f1': 0.46823887224181864}
```

Confusion matrix:

```
[28] y_true = gt
y_pred = preds
cell_type_list = np.unique(y_true)
matrix = confusion_matrix(y_true, y_pred, labels=cell_type_list)
matrix = matrix.astype("float") / matrix.sum(axis=1)[:, np.newaxis]

df = pd.DataFrame(matrix, index=cell_type_list[matrix.shape[0]], columns=cell_type_list[matrix.shape[1]])

ax = sns.clustermap(df,
                    cmap='Purples',
                    annot=True, font="2f",
                    annot_kws={'size': 8},
```



Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

Evaluate the performance

```
[x] [27] res_dict = {
    "accuracy": accuracy_score(gt, preds),
    "precision": precision_score(gt, preds, average="macro"),
    "recall": recall_score(gt, preds, average="macro"),
    "macro_f1": f1_score(gt, preds, average="macro"),
}

res_dict

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples. Use 'zero_division' parameter to control.
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1344: UndefinedMetricWarning: Recall is ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to control.
{'accuracy': 0.8748944291781164,
 'precision': 0.5274756111472559,
 'recall': 0.48119888193882796,
 'macro_f1': 0.46823887224181864}
```

Confusion matrix:

```
y_true = gt
y_pred = preds
cell_type_list = np.unique(y_true)
matrix = confusion_matrix(y_true, y_pred, labels=cell_type_list)
matrix = matrix.astype("float") / matrix.sum(axis=1)[:, np.newaxis]

df = pd.DataFrame(matrix, index=cell_type_list[:matrix.shape[0]], columns=cell_type_list[:matrix.shape[1]])

ax = sns.clustermap(df,
                    cmap='Purples',
                    annot=True, fmt=".2f",
                    annot_kws={'size': 8},
                    vmin=0,
                    vmax=1,
                    row_cluster=False,
                    col_cluster=False,
                    figsize=(14, 14))
```





Tutorial_ZeroShot_Reference_Mapping.ipynb

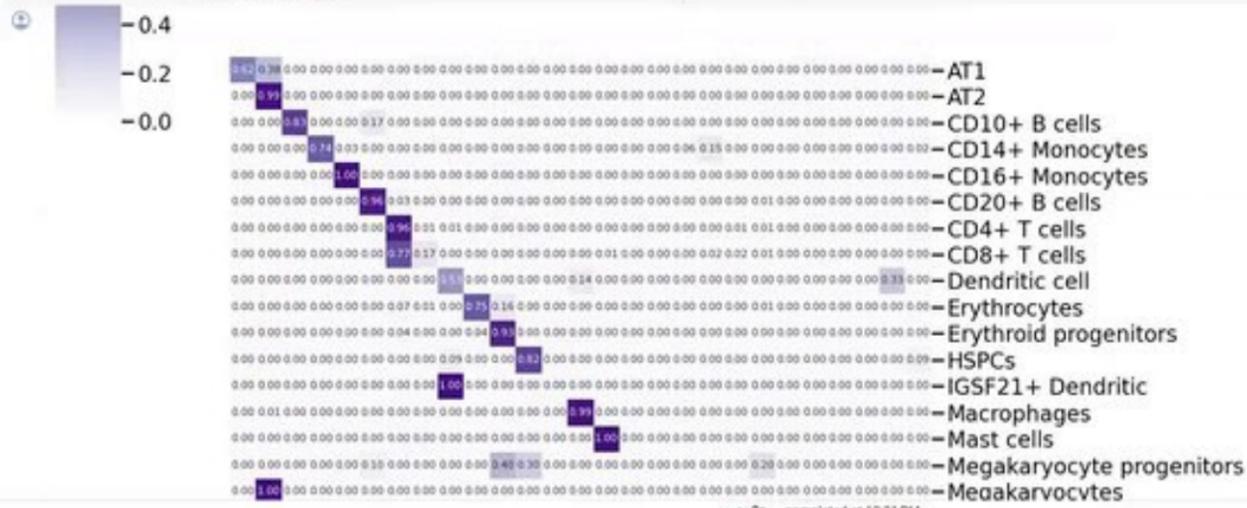
File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

```
y_true = gt
y_pred = preds
cell_type_list = np.unique(y_true)
matrix = confusion_matrix(y_true, y_pred, labels=cell_type_list)
matrix = matrix.astype("float") / matrix.sum(axis=1)[:, np.newaxis]

df = pd.DataFrame(matrix, index=cell_type_list[:matrix.shape[0]], columns=cell_type_list[:matrix.shape[1]])

ax = sns.clustermap(df,
                    cmap='Purples',
                    annot=True, font="2f",
                    annot_kws={"size": 8},
                    vmin=0,
                    vmax=1,
                    row_cluster=False,
                    col_cluster=False,
                    figsize=(14, 14))
```



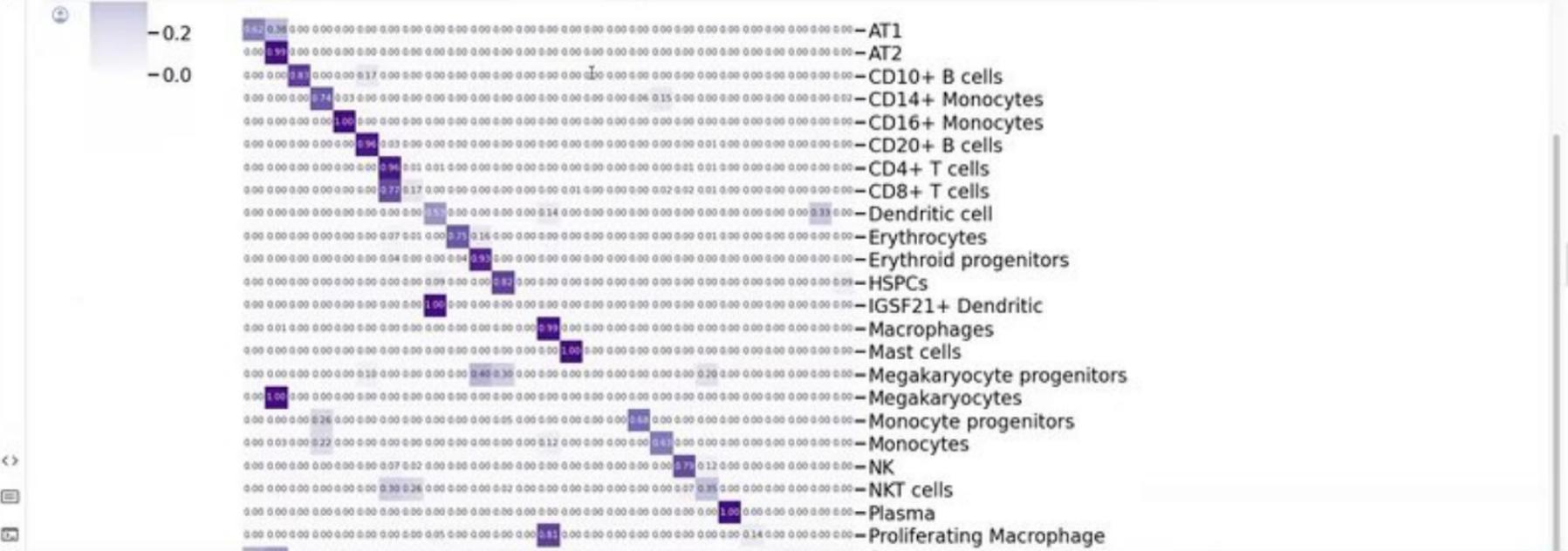


Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

```
ax = sns.clustermap(df,
                    cmap='Purples',
                    annot=True, font="21",
                    annot_kws={'size': 8},
                    vmin=0,
                    vmax=1,
                    row_cluster=False,
                    col_cluster=False,
                    figsize=(14, 14))
```



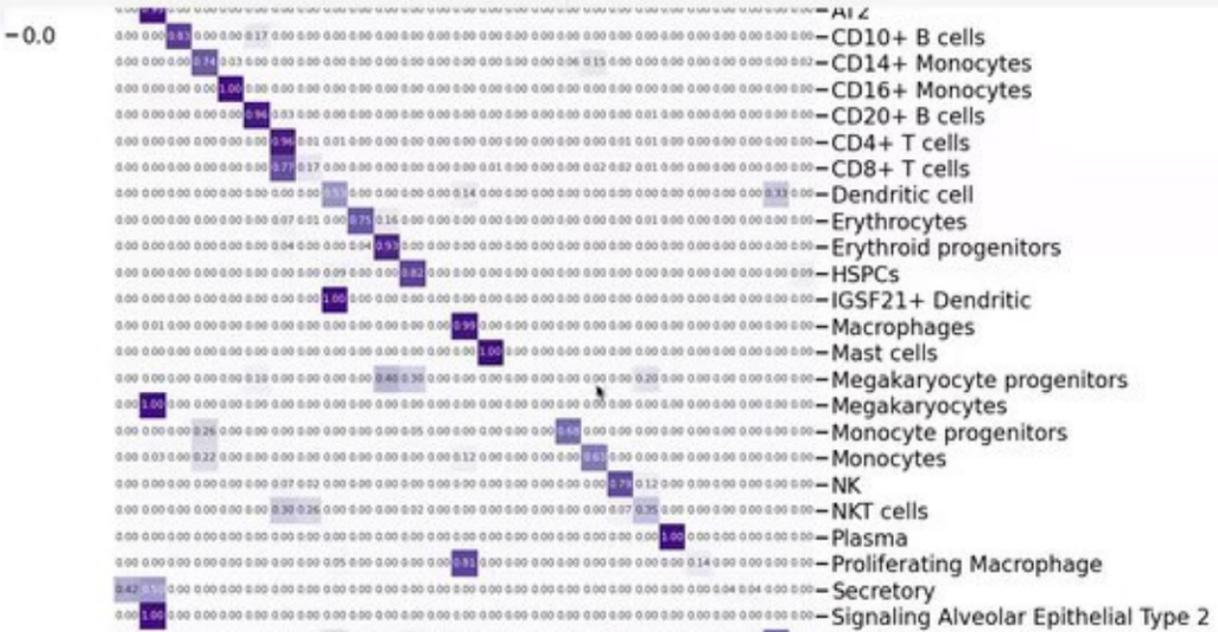


Tutorial_ZeroShot_Reference_Mapping.ipynb

File Edit View Insert Runtime Tools Help Cannot save changes

Code Text Copy to Drive

```
ax = sns.clustermap(df,
                    cmap='Purples',
                    annot=True, font="2f",
                    annot_kws={'size': 8},
                    vmin=0,
                    vmax=1,
                    row_cluster=False,
                    col_cluster=False,
                    figsize=(14, 14))
```





Haotian Cui

```
[29] sc.pp.normalize_total(adata, target_sum=1e4)
sc.pp.log1p(adata)
sc.pp.normalize_total(test_adata, target_sum=1e4)
sc.pp.log1p(test_adata)

gene_col = "gene_name"
cell_type_key = "cell_type"

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar
and should_run_async(code)

ref_embed_adata = scg.tasks.embed_data(
    adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)
test_embed_adata = scg.tasks.embed_data(
    test_adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)
# concatenate the two datasets
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")
# mark the reference vs. query dataset
adata_concat.obs["is_ref"] = ["Query"] * len(test_embed_adata) + ["Reference"] * len(
    ref_embed_adata
)
adata_concat.obs["is_ref"] = adata_concat.obs["is_ref"].astype("category")
# mask the query dataset cell types
adata_concat.obs[cell_type_key] = adata_concat.obs[cell_type_key].cat.add_categories(["To be predicted"])
adata_concat.obs[cell_type_key][: len(test_embed_adata)] = "To be predicted"

--- scGPT - INFO - match 2978/3000 genes in vocabulary of size 60697.
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' ar
and should_run_async(code)
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-
warnings.warn(
Embedding cells: 100% [██████████] 363/363 [03:00:00:00, 2.02it/s]
/usr/local/lib/python3.10/dist-packages/scgpt/tasks/cell_emb.py:279: ImplicitModificationWarning: Setting element '.obsn['X_scGPT']' of view, initializing view as actual.
adata.obsn["X_scGPT"] = cell_embeddings
scGPT - INFO - match 2978/3000 genes in vocabulary of size 60697.
/usr/local/lib/python3.10/dist-packages/scgpt/model/model.py:77: UserWarning: flash-attn is not installed, using pytorch transformer instead. Set use_fast_transformer=False to avoid this warning. Installing flash-
✓ Os completed at 12:31 PM
```

Appendix: Reference mapping on Lung-Kim dataset

Appendix: Reference mapping on Lung-Kim dataset

The dataset can be accessed from [here](#).

```
adata = sc.read_h5ad('../data/lung/sample_proc_lung_train.h5ad')
test_adata = sc.read_h5ad('../data/lung/sample_proc_lung_test.h5ad')
```

```
sc.pp.normalize_total(adata, target_sum=1e4)
sc.pp.log1p(adata)
sc.pp.normalize_total(test_adata, target_sum=1e4)
sc.pp.log1p(test_adata)
```

```
gene_col = "gene_name"
cell_type_key = "cell_type"
```

```
/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: 'should_run_async' will not call 'transform_cell' automatically in the future. Please pass the result to 'transformed_cell' and should_run_async(code)
```

```
ref_embed_adata = scg.tasks.embed_data(
    adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)
```

```
test_embed_adata = scg.tasks.embed_data(
    test_adata,
    model_dir,
    gene_col=gene_col,
    batch_size=64,
)
```

```
# concatenate the two datasets
adata_concat = test_embed_adata.concatenate(ref_embed_adata, batch_key="dataset")
# mark the reference vs. query dataset
adata_concat.obs["is_ref"] = ["Query"] * len(test_embed_adata) + ["Reference"] * len(
    ref_embed_adata
)
```

```
adata_concat.obs["cell"] = adata_concat.obs["cell"].astype("category")
```





Apps > scGP: Reference Mapping Using Cell Embedding

Feedback

Tutorials

Upgrade

setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



scgptRefMap-20231212-223918-...

Celltype

index

Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

Example

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

Preview

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

index

Need help?

Submit Job



Haotian Cui



Apps > scGPT: Reference Mapping Using Cell Embedding

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



Feedback

Tutorials

Upgrade

Haotian Cui

scgptxemap-20231212-225044-...

scgptRefMap-20231212-223918-...

Celltype

Index

Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

Example

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

Preview

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

Index

Need help?

Submit Job



Apps > scGPT:Reference Mapping Using Cell Embedding

[Feedback](#)[Tutorials](#)[Upgrade](#)

Haotian Cui

scGPT:Reference Mapping Using Cell Embedding/WangLab at University of Toronto | Cui et al. v0.1

[Single-Cell Bioinformatics](#) [Cell Type Annotation](#) [H5ad](#) [Csv](#) [Example Results](#) [See Source](#)

Released: Nov-06-2023

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.[Show More](#)RUNS 26 SAVE 1 [Share](#) 

Previous Job Parameters

scgptRefMap-20231212-225834-...

Cell Type Column Name (Optional)

celltype

Gene Column Name

gene_name

scgptRefMap-20231212-225544-...

index

scgptRefMap-20231212-223918-...

Celltype

index

[Give us feedback on this app](#)

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

[Example](#)

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

index



Haotian Cui

Apps > scGPT: Reference Mapping Using Cell Embedding

Feedback

Tutorials

Upgrade

Released: Nov-06-2023

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



scgptRefMap-20231212-225544-...

index

scgptRefMap-20231212-223918-...

Celltype

index

Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

[Example](#)

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsamped.h5ad

Preview

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL) ⓘ

Celltype

GENE COLUMN NAME ⓘ

Index

Need help?

Submit Job



Apps > scGPT:Reference Mapping Using Cell Embedding

[Feedback](#) [Tutorials](#)[Upgrade](#)

scGPT:Reference Mapping Using Cell Embedding/WangLab at University of Toronto |

Cui et al. 2023

[Single-Cell Bioinformatics](#) [Cell Type Annotation](#) [H5ad](#) [Csv](#) [Example Results](#) [See Source](#)

Released: Nov-06-2023

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.[Show More](#)RUNS ▶ 20 SAVE ★ 1 [Share](#)

Previous Job Parameters

scgptRefMap-20231212-225834-...

celltype

gene_name

scgptRefMap-20231212-225544-...

index

scgptRefMap-20231212-223918-...

Celltype

index

[Give us feedback on this app](#)

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

[Example](#)

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

index

Haotian Cui





Haotian Cu

Apps > scGPT:Reference Mapping Using Cell Embedding

Feedback

Tutorials

Upgrade

scGPT:Reference Mapping Using Cell Embedding/WangLab at University of Toronto | Cui et al. 2023

Single-Cell Bioinformatics Cell Type Annotation h5ad Csv Example Results See Source

Released: Nov-06-2023

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



Previous Job Parameters

scgptRefMap-20231212-225834-...

celltype

gene_name

scgptRefMap-20231212-225644-...

index

scgptRefMap-20231212-223918-...

Celltype

index

Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

Example

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

index



Apps > scGPT: Reference Mapping Using Cell Embedding

Efficiently map your data with the scGPT-Transcript model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

[Example](#)

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

[Preview](#)

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

index

[Need help?](#)

[Submit Job](#)

[Feedback](#)

[Tutorials](#)

[Upgrade](#)

scgptRefMap-20231212-223918-...

Celltype

Index

Haotian Cui



Haotian Cui

Apps > scgpt:Reference Mapping Using Cell Embedding

[Feedback](#) [Tutorials](#)[Upgrade](#)

Efficiently map your data with the scgpt framework model coverage profiles and embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

scgptRefMap-20231212-223918-...

Celltype

[Index](#)

Example use case: Understand the cell composition of newly collected samples.

[Show More](#)RUNS 20 SAVE 1 [Share](#)

Upload Single Cell RNA-Seq

DATA STRUCTURE

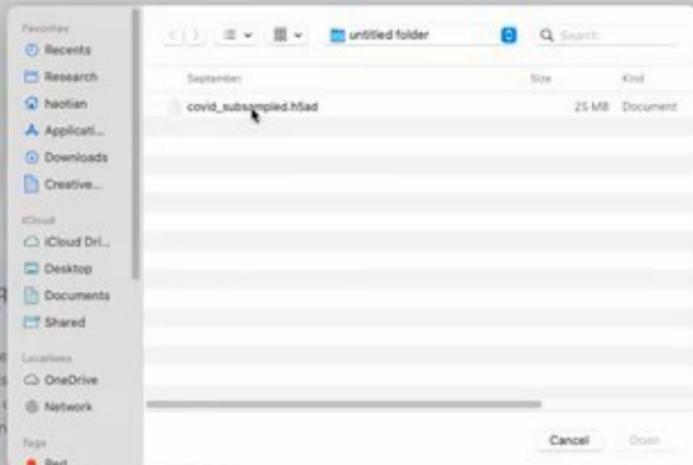
The h5ad format in scRNA-seq reference (HDF5) Annotated Data format. It is used for expression data. The h5ad format also includes a raw gene expression matrix and additional annotations.

📁 **UPLOAD FROM** Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

[Preview](#)[Need help?](#)[Submit Job](#)



Apps > scGPT: Reference Mapping Using Cell Embedding

Efficiently map your data with the scRNA-seq training model coverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



Feedback Tutorial

scgptRefMap-20231212-223918-...

Celltype

Index

You can't see this video

Haotian Cui

Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

Example

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

Upload from Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

Preview

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

Index

Need help?

Submit Job



Apps > scGPT:Reference Mapping Using Cell Embedding

[Feedback](#) [Tutorials](#)

Haotian Cui

scGPT:Reference Mapping Using Cell Embedding/WangLab at University of Toronto | [Previous Job Parameters](#)

Cui et al. v0.1

[Single-Cell Bioinformatics](#) [Cell Type Annotation](#) [H5ad](#) [Csv](#) [Example Results](#) [See Source](#)

Released Nov-06-2023

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

[Show More](#)

RUNS 28 SAVE 1 Share

Your previous job parameters will show up here so you can keep track of your jobs

[Give us feedback on this app](#)

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

[Example](#)

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

No Files Selected!

formats: h5ad

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL) ⓘ

Celltype

GENE COLUMN NAME ⓘ

index

Apps > scGPT:Reference Mapping Using Cell-Embedding

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS 2% SAVE 1 Share

Upload Single Cell RNA-Seq

DATA STRUCTURE

The h5ad format in scRNA-seq refers to the (.HDF5) Annotated Data format. It is used to store gene expression data. The h5ad format includes a raw gene expression matrix and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

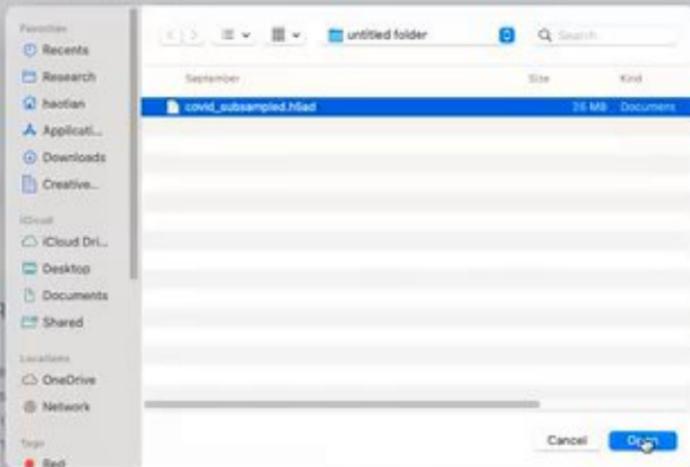
Remote

No Files Selected!
Allowed file formats: h5ad

Need help?

Need help?

Submit job



Feedback Tutorials

Upgrade

scgptRefMap-20231212-223918-...

scgptRefMap-20231212-223918-...

Celltype

Index

Haotian Cui



Loading Covid_subsamped.H5ad Obs Preview

Loading Covid_subsamped.H5ad Var Preview



Uploading your data, please wait...



Haotian Cui

Covid_subsampled.H5ad Vqr Preview

Covid_subsampled.H5ad Obs Preview

Feature Name	feature_types-1-0-1	gene_ids-1-0-1	gene_name	genome-1-0-1	highly_variable	highly_variable_nbatches	highly_variable_rank	means
HES4	Gene Expression	ENSG00000188290	HES4	GRCh38	true	13.00	286.00	0.20
ISG15	Gene Expression	ENSG00000187608	ISG15	GRCh38	true	15.00	247.00	5.55
TNFRSF4	Gene Expression	ENSG00000186827	TNFRSF4	GRCh38	true	13.00	438.00	0.06
MXRA8	Gene Expression	ENSG00000162576	MXRA8	GRCh38	true	4.00	421.00	0.01
ATAD3A	Gene Expression	ENSG00000197785	ATAD3A	GRCh38	true	3.00	435.00	0.05
PRDM16	Gene Expression	ENSG00000142611	PRDM16	GRCh38	true	1.00	376.00	0.01
UTS2	Gene Expression	ENSG00000049247	UTS2	GRCh38	true	3.00	87.00	0.02
TNFRSF9	Gene Expression	ENSG00000049249	TNFRSF9	GRCh38	true	5.00	392.00	0.02
ERRF1	Gene Expression	ENSG00000116285	ERRF1	GRCh38	true	9.00	342.00	0.21



Apps > scGP: Reference Mapping Using Cell Embedding

Feedback

Tutorials

Upgrade

setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

scgpRefMap-20231212-223918-...

Celltype

index

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

Example

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

Preview

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

index

Need help?

Submit Job



Haotian Cui



Haotian Cui

Covid_subsampled.H5ad Var Preview

Covid_subsampled.H5ad Obs Preview

Feature Name	feature_types-1-0-1	gene_ids-1-0-1	gene_name	genome-1-0-1	highly_variable	highly_variable_nbatches	highly_variable_rank	means
HES4	Gene Expression	ENSG00000188290	HES4	GRCh38	true	13.00	288.00	0.20
ISG15	Gene Expression	ENSG00000187608	ISG15	GRCh38	true	15.00	247.00	5.55
TNFRSF4	Gene Expression	ENSG00000186827	TNFRSF4	GRCh38	true	13.00	438.00	0.06
MXRA8	Gene Expression	ENSG00000162576	MXRA8	GRCh38	true	4.00	421.00	0.01
ATAD3A	Gene Expression	ENSG00000197785	ATAD3A	GRCh38	true	3.00	435.00	0.05
PRDM6	Gene Expression	ENSG00000142611	PRDM6	GRCh38	true	1.00	376.00	0.01
UTS2	Gene Expression	ENSG00000049247	UTS2	GRCh38	true	3.00	87.00	0.02
TNFRSF9	Gene Expression	ENSG00000049249	TNFRSF9	GRCh38	true	5.00	392.00	0.02
ERRFI1	Gene Expression	ENSG00000116285	ERRFI1	GRCh38	true	9.00	342.00	0.21



Haotian Cui



Apps > scGPT:Reference Mapping Using Cell Embedding

Feedback

Tutorials

Upgrade

setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

scgptRefMap-20231212-223918-...

Celltype

index

Example use case: Understand the cell composition of newly collected samples.

> Show More

RUNS ▶ 26 SAVE ★ 1 Share



Give us feedback on this app

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

Example

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsampled.h5ad

Preview

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL)

Celltype

GENE COLUMN NAME

index

This field is required

index

gene_name

Email & Symbols

Undo

Redo

Cut

Copy

Paste

Paste and Match Style

Select All

Language Settings

Writing Direction

Open in Reading Mode View

Paperpile

Inspect

Autofill

Submit Job



Apps > scGPT: Reference Mapping Using Cell Embedding

[Feedback](#) [Tutorials](#)[Upgrade](#)

setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

scgptRefMap-20231212-223918-...

Celltype

[index](#)

Example use case: Understand the cell composition of newly collected samples.

[Show More](#)RUNS ▶ 26 SAVE ★ 1 [Share](#)[Give us feedback on this app](#)

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

covid_subsamped.h5ad

[Preview](#)

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL) ⓘ

Celltype

GENE COLUMN NAME ⓘ

gene_name

[Run on GPU](#)[Need help?](#)[Submit Job](#)

Uploading your data, please wait...



Haotian Cui



Apps > scGPT:Reference Mapping Using Cell Embedding

[Feedback](#)[Tutorials](#)[Upgrade](#)

scGPT:Reference Mapping Using Cell Embedding/WangLab at University of Toronto | Cui et al. 2023

Single-Cell Bioinformatics Cell Type Annotation H5ad Csv [Example Results](#) [See Source](#)

Released Nov-06-2023

Efficiently map your data with the scGPT_human model. Leverage pre-trained embeddings in a zero-shot setting for quick, accurate cell analysis. Ideal for transferring cell type and disease condition metadata, this tool simplifies understanding cell compositions in new samples.

Example use case: Understand the cell composition of newly collected samples.

[Show More](#)

RUNS ▶ 26 SAVE ★ 1 [Share](#)



Uploading your data, please wait...

[Give us feedback on this app](#)

Upload Single Cell RNA-Seq Data File

DATA STRUCTURE

[Example](#)

The h5ad format in scRNA-seq refers to the Hierarchical Data Format 5 (HDF5) Annotated Data format. It is commonly used to store single-cell gene expression data. The h5ad format organizes data in a hierarchical structure, including a raw gene expression matrix, cell and gene metadata, and additional annotations.

UPLOAD FROM Size limits - Local: 100Mb / Remote: no limit

Local

Remote

Set Parameters

CELL TYPE COLUMN NAME (OPTIONAL) ⓘ

Celltype

GENE COLUMN NAME ⓘ

gene_name

Run on GPU

Haotian Cui

Jobs > Analyze and download your results

JOBS WORKSPACE

All Jobs

SORT BY

Recent

DATE Clear

From: dd/mm/yyyy

To: dd/mm/yyyy

PARAMETERS

Hover over a job and view the job's parameters here

Search by job name or app name



JOBS: 4

TITLE	STATUS	START TIME	DURATION
scgptRefMap-20231213-174000-gpu	Preparing	12/13/2023, 12:40:00 PM	0h:0m:0s
scgptRefMap-20231212-225834-gpu	Completed	12/12/2023, 5:58:34 PM	0h:2m:0s
scgptRefMap-20231212-225544-gpu	Completed	12/12/2023, 5:55:44 PM	0h:2m:5s
scgptRefMap-20231212-223918-gpu	Completed	12/12/2023, 5:39:18 PM	0h:1m:11s



Haotian Cui

Results

Below, you'll find a preview of the comparison table showcasing the original labels alongside the predicted labels if the original cell type column is provided in the parameters before submitting the job. For a comprehensive view, you can download the complete



Parameters

File H5ad: Covid_subsampled.H5ad
Cell Type Column Name (Optional): Celltype
Gene Column Name: Gene_name
Workflow_name: Workflow.Yml

Actions

- Download Output Files
- Share This Job
- Reopen App
- Save App To My Library
- Give us feedback on this app

Propagated CellXGene labels (Preview)

Original Labels	CellXGene Labels
Monocytes	classical monocyte
Mast cells	mast cell
AT2	type II pneumocyte
T cell	effector memory CD8-positive, alpha-beta T cell
CD4+ T cells	CD4-positive, alpha-beta memory T cell
AT2	type II pneumocyte
Macrophages	macrophage

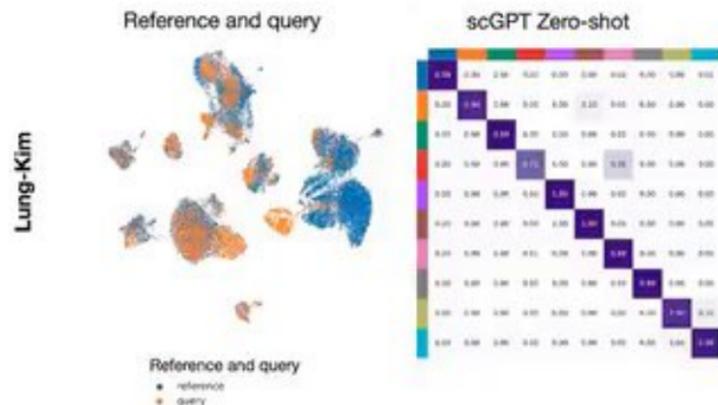
Contents

Propagated_CellXGene_Labels_(Preview)

MCBRLab Zero-shot reference mapping



- We also use the embeddings in a reference mapping manner
- We split the data to simulate a setting that some of the patient cells, query set, were not annotated. Cell types were propagated from reference data set to label the query set.
- The zero-shot workflow demonstrates competitive accuracy across comparisons.



Dataset	Method	Evaluation Metrics			
		<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>MacroF1</i>
COVID-19	scGPT (fine-tuned)	0.891	0.577	0.567	0.547
	scGPT (zero-shot)	0.867	0.513	0.476	0.468
	Azimuth	0.878	0.515	0.418	0.444
	expiMap	0.730	0.441	0.325	0.335
Lung-Kim	scGPT (fine-tuned)	0.974	0.959	0.967	0.962
	scGPT (zero-shot)	0.968	0.970	0.933	0.948
	Azimuth	0.970	0.957	0.962	0.959
	expiMap	0.920	0.940	0.846	0.878



Haotian Cui

- Zero-shot can be promising, we are going to provide further updates for the scenario:
 - Hierarchical cell type labeling of reference mapping from atlas
 - Going to provide new checkpoint specifically for zero-shot cell embeddings
 - Zero-shot applications using gene embeddings
 - Others



Thank you!

scGPT codebase: <https://github.com/bowang-lab/scGPT>,

Sincere gratitude to collaborators: Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Dr. Lin Zhang, Dr. Nan Duan, and Dr. Bo Wang



WangLab@UofT



Results

Below, you'll find a preview of the comparison table showcasing the original labels alongside the predicted labels if the original cell type column is provided in the parameters before submitting the job. For a comprehensive view, you can download the complete



Parameters

File H5ad: Covid_subsampled.H5ad
Cell Type Column Name (Optional): Celltype
Gene Column Name: Gene_name
Workflow_name: Workflow.Yml

Actions

- Download Output Files
- Share This Job
- Reopen App
- Save App To My Library
- Give us feedback on this app

Propagated CellXGene labels (Preview)

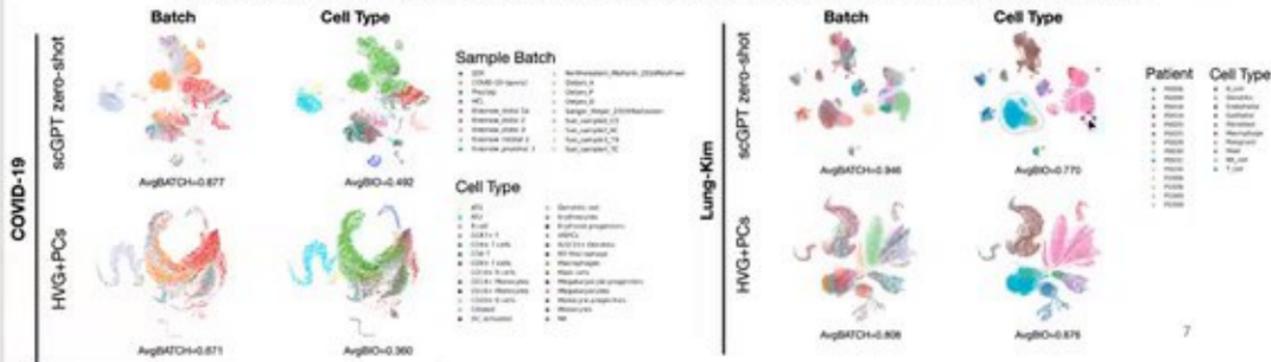
Original Labels	CellXGene Labels
Monocytes	classical monocyte
Mast cells	mast cell
AT2	type II pneumocyte
T cell	effector memory CD8-positive, alpha-beta T cell
CD4+ T cells	CD4-positive, alpha-beta memory T cell
AT2	type II pneumocyte
Macrophages	macrophage

Contents

Propagated_CellXGene_Labels_(Preview)

Zero-shot applications

- **Zero-shot**, meaning using the pretrained model to generate embeddings for new data, without any further training, **much faster, more accessible**
- UMAP of zero-shot embeddings on two disease datasets, compared to HVGs
 - Considerable ability to distinguish cell types
 - Note the pretraining process is not designed for mitigating technical batch effects



Have been working new extensions exploring the possibility of zero-shot applications. Here, I'd like to share some new results

Particularly want to mention the influence of pretraining, so of the technical info is actually signals



Haotian Cui





