

OmniCellTOSG: The First Cell Text-Omic Signaling Graphs Dataset for Joint LLM and GNN Modeling

Heming Zhang*
hemingzhang@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Tim Xu*
tianqi.x@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Dekang Cao*
c.dekang@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Shunning Liang
l.shunning@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Lars Schimmelpfennig
l.schimmelpfennig@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Levi Kaster
k.levi@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Di Huang
di.huang@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Carlos Cruchaga
cruchagac@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Guangfu Li
gli@uchc.edu
University of Connecticut
Storrs, CT, USA

Michael Province
mprovince@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Yixin Chen
ychen25@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Philip Payne
prpayne@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Fuhai Li[†]
fuhai.li@wustl.edu
Washington University in St. Louis
St. Louis, MO, USA

Abstract

The human body consists of approximately 37 trillion cells, all originating from a single embryonic cell and sharing the same copy of genome. The complex, robust and accurate cell signaling systems, regulated by varying abundance of proteins and their interactions, create diverse cell types with different functions at different organs. The cell signaling systems are evolved and altered by many factors, like age, sex, diet, environment exposures and diseases. However, it remains an open problem to decode cell signaling systems or patterns in normal development or diseases because the systems are rather complex consists of tens of thousands of genes/proteins and massive interactions among them. Recently, hundreds of millions of single cell omic data have been being generated by many research projects, which provide the solid basis for understanding

cell signaling systems, like the key genes/proteins and their signaling interactions, within diverse cell-subpopulations in different conditions. Moreover, inspired by the success of foundation models that are pre-trained on massive datasets, like large language models (LLMs) and large vision models (LVMs), in this study, we build the first dataset of cell text-omic signaling graphs (TOSGs), named OmniCellTOSG. Each TOSG represents the signaling graph/system of an individual cell or meta-cell, and associated with labels, like organ, disease, sex, age, cell-subtype. The unique contributions of the OmniCellTOSG are two-folds. First, the TOSGs represents a novel and ideal graph data model for decoding cell signaling systems via graph reasoning by incorporating human-understandable textual annotation/prior knowledge (like biological functions, cellular locations, related signaling pathways, related diseases and drugs), with numeric values (that represent the observed abundance levels genes/proteins) in different cells in different organs and conditions. Also new paradigm-shift data analysis models like the joint LLM and GNN models are needed to analyze the TOSGs. Secondly, OmniCellTOSG consists of large-scale cell text-omic signaling graphs, using single cell RNAseq (scRNAseq) data from 120 millions cells from diverse tissues/organs, health vs diseases. The OmniCellTOSG data are structured in a format that is fully compatible with PyTorch, and can facilitate the development of novel joint LLM and graph neural network (GNN) foundation cell signaling models for decoding the complex cell signaling systems via TOSG graph reasoning. It could shift the paradigm in life sciences, healthcare and

*Equal contribution

[†]Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD' 25, August 3-7, 2025, Toronto, Canada

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXXXXX>

precision medicine research. The number of cells in OmniCellTOSG keeps growing and will be updated regularly. Dataset and code are available at Github¹.

CCS Concepts

• **Applied computing** → **Bioinformatics**; • **Computing methodologies** → **Artificial intelligence**.

Keywords

Large Language Models, Graph Neural Networks, Cell Signaling Graphs, Single Cell, scRNAseq, snRNAseq

ACM Reference Format:

Heming Zhang, Tim Xu, Dekang Cao, Shunning Liang, Lars Schimmelpfennig, Levi Kaster, Di Huang, Carlos Cruchaga, Guangfu Li, Michael Province, Yixin Chen, Philip Payne, and Fuhai Li. 2025. OmniCellTOSG: The First Cell Text-Omic Signaling Graphs Dataset for Joint LLM and GNN Modeling. In *Proceedings of Proceedings of the 31th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD' 25)*. ACM, New York, NY, USA, 11 pages. <https://doi.org/XXXXXXXXXX>

Copy

1 Introduction

The human body consists of approximately 37.2 trillion cells, all originating from a single embryonic cell and sharing the same copy of genome. The complex, robust and accurate cell signaling systems, regulated by varying abundance of proteins and their interactions, create diverse cell types with different functions at different organs. The cell signaling systems are evolved and altered by many factors, like age, sex, diet, environment exposures and diseases. Though many biomarkers and knowledge have been uncovered in life science and healthcare studies, the cell signaling systems still remain mysterious. For example, what are the panoramic view of cell signaling systems (all the entities and their interactions) within the cells? How do the cell signaling systems evolve and altered by the factors, like age, sex and diseases? What are the disease related cell subtypes and the interactions among these cells? How can we perturb these cells' signaling systems and interactions to prevent and treat diseases. These questions (e.g., disease pathogenesis) are some of the major reasons of why there is no drug can cure complex diseases, like Alzheimer's disease (AD), cancer, heart disease, and many other chronic inflammation related diseases, like kidney failure and liver hepatitis and cirrhosis.

Recently, single-cell/nucleus RNA sequencing (sc/snRNAseq) has revolutionized our ability to measure transcriptomic abundance at individual cell level. With the sc/snRNAseq, it is feasible to identify the cell types/subtypes in (disease) tissues, and investigate the cell signaling systems and their signaling interactions within a niche or microenvironment (ME). For example, hundreds of millions of sc/snRNAseq data were generated by Human Cell Atlas (HCA)[24], Brain Cell Atlas[4], and many studies of diverse diseases[18, 20]. These datasets are valuable to systematically investigate and decode the cell signaling systems. In another word, it is crucial to understand which groups of genes/proteins with different abundance levels work together coordinately to achieve the specific biological

functions or tasks in the diverse cell sub-populations in a disease or organ niche.

Furthermore, the success of large-scale foundation models, like chatGPT, have revolutionized AI research and applications. The foundation models were pretrained on massive and diverse datasets via self-supervised learning (SSL), and thus can have the generalized understanding of the information patterns embedded within the massive and diverse training data, consequently, serving as a solid base upon which specialized adaptations can be developed to tackle specific tasks or challenges. Therefore, the disease specific data analysis only measures and observes a limited of number of cell signaling systems patterns/scenarios. Consequently, deep learning models trained on small-scale, disease-specific datasets are prone to bias and overfitting, often converging to local minima. This challenge is analogous to the limitations of training ChatGPT-scale foundation models on restricted language datasets. Thus, foundation models have emerged as a promising approach to address these issues. However, due to their inherent "black box" nature, most foundation models struggle to effectively integrate detailed biological information pertaining to cell signaling interactions. Herein, in this study, we build the OmniCellTOSG dataset. To the best of our knowledge, OmniCellTOSG is the first Text-Omic Signaling Graph (TOSG) dataset. It creates a new graph data type integrating both human-understandable text-attributed information and numerical omic features. The textual information annotates the accumulated knowledge of genes or proteins, like the biological functions, cellular locations, related diseases and drugs. The omic feature indicates the abundance level of the genes/proteins. In its current version, the human-interpretable annotations are derived from BioMedGraphica[34], a unified database compiling prior knowledge on genes, proteins, drugs, diseases, and phenotypes from diverse data sources. Future iterations will incorporate enriched knowledge from extensive literature, synthesized using advanced large language models such as ChatGPT-4. This unprecedented dataset is expected to facilitate the development of novel joint foundation models that integrate large language models (LLMs) with graph neural networks (GNNs) to decode complex cell signaling networks for interesting cell signaling patterns.

In the following sections, we detailed the construction and utilization of OmniCellTOSG, including data sourcing, preprocessing workflows, graph-generation protocols, and data loading using our developed package, CellTOSGDataset. By releasing OmniCellTOSG to the broader research community, we aim to foster collaboration between data scientists and biologists, ultimately accelerating breakthroughs in precision medicine and enhancing our understanding of disease pathogenesis.

2 Related Work

Recently, massive single cell/nucleus RNAseq datasets have been generated to study the diversity of cellular transcriptomics, e.g., the human cell atlas (HCA) project, CZ CellxGene Database, Allen Brain Cell Atlas, and massive datasets in AD knowledge portal, like SEA-AD and many AD studies, as well as many other diseases⁶, and datasets from Gene Expression Omnibus (GEO). In addition, a few exploratory foundation models (only using the gene expression value in random orders), have been developed based on the massive

¹<https://github.com/FuhaiLiAiLab/OmniCellTOSG>

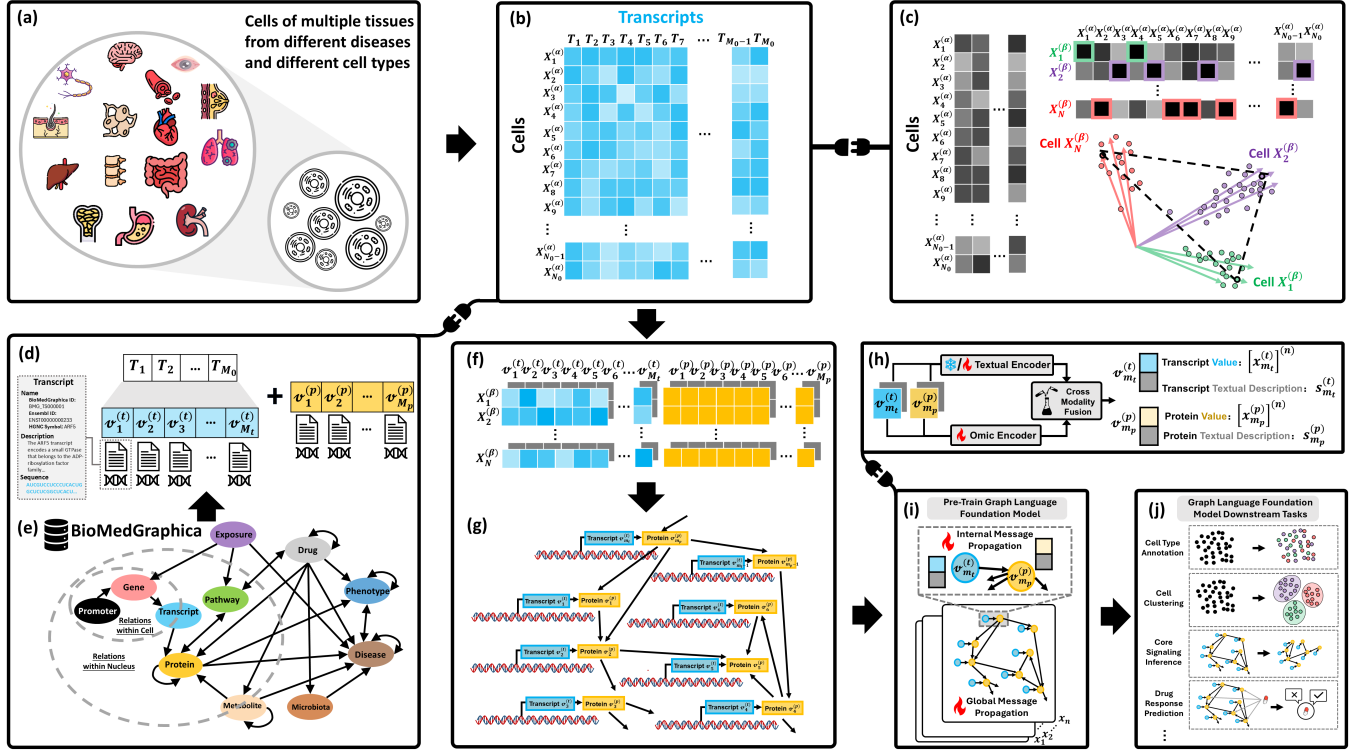


Figure 1: Overview of text-omic signaling graph (TOSG) generation. (a) Millions of single cells collected from multiple tissues, diseases, and cell types. **(b)** The values in the collected h5ad files for those N_0 single cells. **(c)** Archtypal analysis to aggregate N_0 cells into N meta-cells. **(d-e)** Integrating transcript entities into text-omic signaling network with M ($M = M_t + M_p$) matched entities by retrieving the knowledge base from BioMedGraphica. **(f-g)** Generate the text-omic signaling graphs for the matched and virtual entities. **(h)** Joint text-encoder and omic encoder with cross-modality fusion. **(i-j)** Message propagation on the generated text-omic signaling graphs, encapsulating the fused biological and textual information for foundation model training and downstream tasks.

single cell omic data, like A cell atlas foundation model for scalable search of similar human cells[15], geneFormer[26], scGPT[5], scFoundation[13], and GET (general expression transformer)[10]. However, none of them was built by incorporating cell signaling pathways/graphs for the purpose of inferring the disease dysfunctional signaling pathways, and decoding cell signaling graph patterns among diverse cell types under different conditions. Aside from that, recent investigations have revealed that off-the-shelf large language models can struggle with complex reasoning tasks in the biomedical domain, often producing inaccurate or hallucinated outputs [3]. Recent studies have demonstrated that integrating knowledge graphs can significantly enhance model reasoning capabilities, effectively mitigating issues such as hallucination in large language models. At the same time, GNN models are both with challenges of the expressive power of graph neural networks [1, 7]. By combining biologically meaningful knowledge graphs with quantitative omic features, one can more accurately capture complex cellular interactions. Motivated by these advances, we developed the OmniCellTOSG dataset, which incorporates both human-interpretable text annotations and numerical omic features. The textual annotations encode accumulated prior knowledge about genes and proteins—including their biological functions, cellular

localizations, and associations with diseases and drugs—while the numerical features represent their abundance levels. For example, it is well known that the apolipoprotein E4 (APOE4) gene is a significant genetic risk factor for Alzheimer’s disease. Whereas, the APOE2 gene is a variant of the APOE gene that may lower the risk of Alzheimer’s disease. It’s also associated with longevity and reduced cognitive decline. And APOE gene plays a role in cholesterol metabolism in the brain. Therefore, both the human-understandable accumulated prior knowledge and numerical omic features are crucial for decode dysfunctional signaling pathways.

Moreover, disease specific smaller-scale data analysis only measures and observes a limited of number of cell signaling systems scenarios. Consequently, models trained on the small-scale disease specific data might be biased and overfitting to the noisy signal reaching a local minimum. Just like it is limited and infeasible to train ChatGPT scale foundation models on a small language dataset. While using OmniCellTOSGs, hundreds of millions of single cell data, covering diverse tissues/organs, diseases, cell types, age, sex, diet, environmental exposures, will reflect diverse and comprehensive cell signaling patterns. Thus, the OmniCellTOSG data will be valuable to develop novel joint LLM + GNN cell signaling graph foundation AI models.

3 OmniCellTOSG Datasets

In this work, we introduce OmniCellTOSG, a comprehensive dataset that integrates single-cell transcriptomic data from multiple sources with detailed textual annotations. Data were collected from the CellxGene, GEO, Brain Cell Atlas, and SEA-AD repositories, yielding millions of cells across diverse tissues and disease conditions. Rigorous preprocessing ensured cross-dataset compatibility through quality control, normalization, and systematic grouping of organ/tissue and disease labels. Additionally, cell type annotations were refined through a combination of automated methods by CellTypist[6] and manual curation, reducing 910 initial cell types to 134 major categories, along with classifying 22 organ types and 21 disease types. The dataset was originally composed of 117,519,978 cells and was subsequently refined to 547,168 cells, all numerically encoded for downstream analysis. Finally, individual cells were aggregated into meta-cells using the SEACells algorithm, and these meta-cell data were integrated with gene-regulatory network information to construct text-omic signaling graphs, which serve as the foundation for training a joint LLM-GNN cell signaling graph model.

3.1 Data Collection

The dataset was assembled from four primary sources, detailed as follows, with the collection procedures described in the appendix.

CellxGene datasets. The data downloaded from CellxGene consists of over 42 million single cells in H5AD AnnData files, derived from 91 human tissues and encompassing 28 disease studies as well as general single-cell data[19, 23].

Brain Cell Atlas datasets. The data obtained from the Brain Cell Atlas comes from human brain single-cell studies, encompassing 23 disease types and over 7 million single cells in H5AD AnnData files[4, 29].

SEA-AD datasets. Many AD studies deposit AD related datasets into the AD Knowledge Portal[12, 18]. For example, the SEA-AD consortium data contributes a significant portion focused specifically on Alzheimer’s Disease, providing over 68 million cells from brain tissue samples and the original format is h5 containing raw feature-barcode matrices [11, 17, 21].

GEO datasets. Many studies deposited the datasets in the Gene Expression Omnibus (GEO). The data downloaded from GEO contains a wide variety of organ and disease types, covering 15 major tissue types and over 20 disease conditions. The original format of the data was available in two types: 1) Matrix Market format consisting of gene expression matrices, cell barcodes, and feature annotations, and 2) Compressed CSV format containing gene-cell expression matrices. We are collecting and adding more sc/snRNAseq datasets into OmniCellTOSG.

3.2 Data Preprocessing

In order to proceed further with the experiments, the dataset needs to be processed to ensure quality control and standardization. The gene list from[23] the scFoundation[13] article was used to process the data to enable cross-dataset compatibility. All the data underwent rigorous filtering with a minimum threshold of 100 genes per cell to ensure reliable expression profiles. The final pre-processed dataset contains 118 million high-quality cells of 894 different cell types.

3.2.1 Grouping Organs/Tissues The data were grouped based on platform, organ/tissue, and disease. For H5AD files where the raw data already contained organ/tissue and disease information, such as datasets from CellxGene[23] and the Brain Cell Atlas[4], the grouping was directly derived from the dataset’s original classification. For H5AD files that only contained the gene expression matrix without additional label information, such as raw data from GEO and SEA-AD, manual classification was performed by consulting the dataset’s original description. After all datasets were categorized, organ/tissue and disease classifications were systematically reviewed, finer-grained tissue and disease classifications were merged into broader organ and disease categories. Tissue regions were grouped according to anatomical structures into larger organ categories. Similarly, disease subtypes with related pathological characteristics were consolidated into broader disease categories. For instance, various subtypes of gliomas, such as glioblastoma, mixed gliomas, and oligodendroglioma, were combined into the general gliomas category.

3.2.2 Cell Type Annotations To specifically address the lack of cell type annotation in the GEO and SEA-AD datasets, CellTypist, a cell type classification tool, was used for annotation. This process included data normalization (scale factor setting as $1e4$), log-transformation, and majority voting for prediction confidence. For these two datasets, the selection of the appropriate CellTypist pre-trained model was guided by the grouping results of organ/tissue and disease, ensuring that the model used was best suited for the corresponding organ and disease classification. After processing with CellTypist, the labeled H5AD files were enriched with observations for tissues and diseases, addressing the lack of metadata in the original datasets that contained only the gene expression matrix. This ensured that the annotated data could be easily tracked during downstream analysis, eliminating the need for separate files to map H5AD files to their corresponding organs/tissues and diseases.

3.3 Converting Single Cells to Meta-cells

To mitigate the inherent sparsity and noise in single-cell RNA sequencing (scRNAseq) data, we adopt a meta-cell strategy based on the SEACells algorithm[21]. Our approach is designed to ensure consistency across datasets from diverse sources by employing uniform preprocessing, feature selection, and dimensionality reduction procedures before meta-cell aggregation.

3.3.1 Data Preprocessing and Normalization Let the raw single-cell seq data be represented by $\mathcal{X}^{(\alpha)} = \{X_1^{(\alpha)}, X_2^{(\alpha)}, \dots, X_{N_0}^{(\alpha)}, \dots, X_{N_0}^{(\alpha)}\}$, where $X_{n_0}^{(\alpha)} \in \mathbb{R}^{M_0}$ denotes the cell, and N_0 is the number of cells collected from various data resources and M_0 is the number of elements in transcript entity set $\mathcal{T} = \{T_1, T_2, \dots, T_{m_0}, \dots, T_{M_0}\}$. Furthermore, cell annotations are collected or inferred as $\mathcal{Y}^{(\alpha)} = \{Y_{ct}^{(\alpha)}, Y_{og}^{(\alpha)}, Y_{ds}^{(\alpha)}\}$, where $Y_{ct}^{(\alpha)}$ represents cell type names, $Y_{og}^{(\alpha)}$ for tissue names, $Y_{ds}^{(\alpha)}$ for disease names in N_0 cells. To alleviate computational demands, raw data files (stored in H5AD format) are partitioned into subsets of no more than 50,000 cells. For datasets requiring normalization, we first apply total count normalization by scaling UMI counts of each cell to a fixed total of 10,000, followed by a log1p transformation to stabilize variance. Pre-normalized

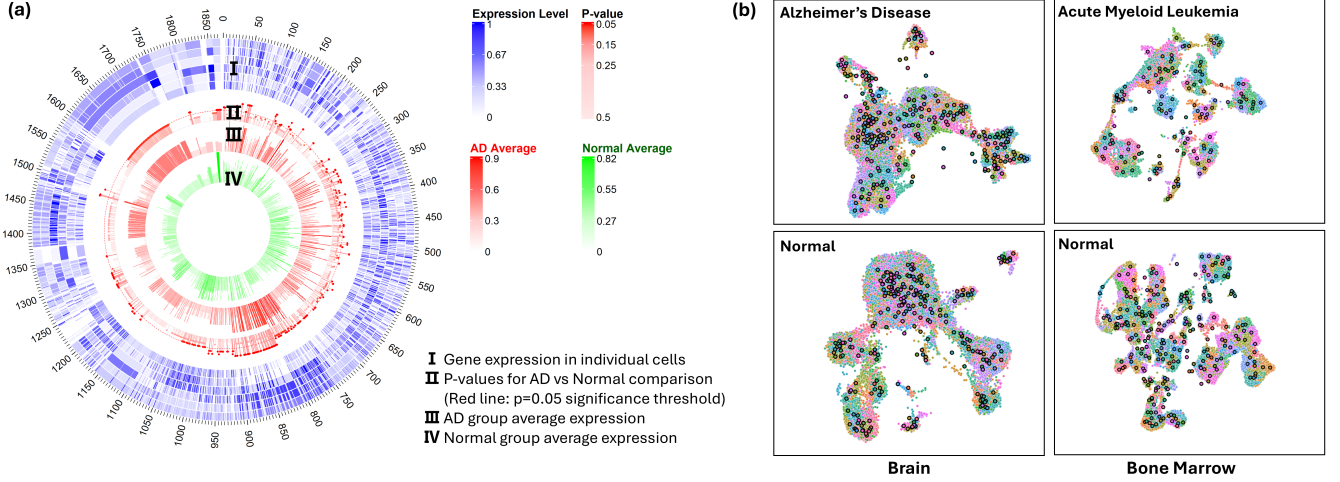


Figure 2: Observation of Meta-Cell Gene Expression Distributions and Clustering Patterns. (a) Circular visualization of differential gene expression between Alzheimer’s Disease (AD) and normal brain samples. The concentric rings represent: (I) Gene expression profiles in individual cells, with the outer three rings corresponding to AD samples and the inner three rings to normal samples, randomly selected from the dataset; (II) P-values derived from a t-test comparing AD and normal cells, with the red line indicating the $p < 0.05$ significance threshold; and (III–IV) Mean gene expression levels for AD and normal groups, respectively. (b) UMAP visualization of meta-cell clustering results for brain and bone marrow tissues. The first column presents AD and corresponding normal samples from the brain, while the second column shows Acute Myeloid Leukemia and normal samples from the bone marrow. Each color represents a cluster corresponding to a distinct cell type, with black circles indicating clusters consolidated into a single meta-cell.

datasets are used as provided, ensuring uniformity across different data sources. In addition, Uniform feature selection is performed by identifying the top 1,500 highly variable genes from each dataset. We then apply Principal Component Analysis (PCA[2]) with 50 components to reduce dimensionality while preserving essential variance. Based on the PCA-reduced features, a K-Nearest Neighbor (KNN[22]) graph is constructed to maintain the underlying structural relationships among cells.

3.3.2 Meta-cells Generation via SEACells Meta-cell generation is performed using the SEACells algorithm. With a fixed aggregation size of N cells per meta cell, SEACells first measures cell-to-cell similarity and then decomposes the resulting structure via archetypal analysis. Cells near the convex hulls of the data distribution are grouped together, yielding a new set of meta cells denoted by $\mathcal{X}^{(\beta)} = \{X_1^{(\beta)}, X_2^{(\beta)}, \dots, X_n^{(\beta)}, \dots, X_N^{(\beta)}\}$, where $X_n^{(\beta)} \in \mathbb{R}^{M_0}$ represents a meta-cell. The associated labels for the meta cells are computed by aggregating the raw cell labels through majority voting, resulting with $\mathcal{Y}^{(\beta)} = \{Y_{ct}^{(\beta)}, Y_{og}^{(\beta)}, Y_{ds}^{(\beta)}\}$, where $Y_{ct}^{(\beta)}, Y_{og}^{(\beta)}, Y_{ds}^{(\beta)}$ are the corresponded cell type names, tissue names, disease names of meta-cells (see **Figure 2**). Nevertheless, due to the diverse origins of the datasets, cell type nomenclature varies in both naming conventions and granularity, necessitating cluster grouping for 894 identified cell types. We first apply keyword-based matching to cluster cell types with distinct naming patterns, such as grouping all T cell variants together. For the remaining types, TF-IDF vectorization followed by hierarchical clustering (Ward’s method with a 1.5 threshold) is used, while unannotated types are isolated into a separate category. This approach consolidates the 910 cell types into 135 major categories, with manual validation ensuring

cluster accuracy. Based on the strategy for grouping organs and diseases as aforementioned, the corresponding annotations are derived accordingly. For downstream classification tasks, cell types, tissue types, and disease types are numerically encoded. Consequently, we construct the label set $\mathcal{Y} = \{Y_{ct}, Y_{og}, Y_{ds}, Y\}$, where $Y_{ct} \in \mathbb{R}^N, Y_{og} \in \mathbb{R}^N, Y_{ds} \in \mathbb{R}^N$ denote the cell type, tissue type and disease type labels of each meta-cell, respectively. In addition, an extra label $Y \in \mathbb{R}^N$ is introduced to indicate the disease status of each cell by marking normal cells as zero and others as one.

3.4 Text-Omic Signaling Graphs Generation

In this stage, the pre-processed single-cell transcriptomic data are integrated with gene-regulatory network information to enable omic analysis and signaling graph construction.

3.4.1 Entity Matching With the pre-processed single-cell transcriptomic dataset, denoted as $\mathcal{X}^{(\beta)} \in \mathbb{R}^{N \times M_0}$, is systematically integrated into the BioMedGraphica framework, incorporating the gene-regulatory network. Using the mapping match table, those M_0 transcript features will be mapped into the M_t transcript entities. In details, each transcript element in set \mathcal{T} will be mapped and extended to the transcript entities set $\mathcal{V}^{(t)} = \{v_1^{(t)}, v_2^{(t)}, \dots, v_{m_t}^{(t)}, \dots, v_{M_t}^{(t)}\}$. By linking transcript nodes within the network to the protein-protein interaction (PPI) graph, proteins are treated as virtual nodes with adding the new entity set $\mathcal{V}^{(p)}$. This integration yields the entity set $\mathcal{V} = \{\mathcal{V}^{(t)}, \mathcal{V}^{(p)}\}$, where $|\mathcal{V}| = |\mathcal{V}^{(t)}| + |\mathcal{V}^{(p)}| = M_t + M_p = M$. Also the feature set $\mathcal{X} = \{\mathcal{X}^{(t)}, \mathcal{X}^{(p)}\}$ are also generated, where $\mathcal{X} \in \mathbb{R}^{N \times M}$, $\mathcal{X}^{(t)} \in \mathbb{R}^{N \times M_t}$ and $\mathcal{X}^{(p)} \in \mathbb{R}^{N \times M_p}$ correspond to the transcriptomic and proteomic feature sets, respectively.

Table 1: OmniCellTOSG Dataset Overview: Detailed Statistics for Cases with Over 800 Meta-Cells

Diseases	Organ/Tissue Types	# of Original Cells	# of Result Cells
Alzheimer's Disease	Brain	69,834,238	315,611
Amyotrophic Lateral Sclerosis	Brain	163,883	819
Gastrointestinal Cancers	Stomach	236,207	1154
General	Multiple organs*	41,742,976	202,306
Gliomas	Brain	1,822,859	9003
Kidney Cancer	Blood, Kidney	191,169	954
Lung Cancer	Adrenal Gland, Brain, Liver, Lung, Lymph Node	2,189,381	10,925
Lymphoma	Bone Marrow, Lymph Node	182,448	911
Nasopharyngeal Carcinoma	Blood, Nasopharynx	176,447	871
... For datasets with fewer than 800 meta-cells, please refer to Table 4 ...			
Total	-	117,519,978	547,168

* Multiple organs include: Adrenal Gland, Blood, Bone Marrow, Brain, Breast, Cervical Spinal Cord, Esophagus, Eye, Gonad, Heart, Intestine, Kidney, Liver, Lung, Lymph Node, Mouth, Pancreas, Skin, Stomach, Uterus.

3.4.2 TOSGs Construction From the perspective of single cell side, the multi-omics \mathcal{X} can be decomposed as $\{X_1, X_2, \dots, X_n, \dots, X_N\}$, where each sample X_n resides in \mathbb{R}^M . Additionally, the cell label matrices set \mathcal{Y} , and given that the cell label set are consistent with label for meta cells, $\mathcal{Y}^{(\beta)}$. Beyond transcriptomic features and virtual proteomic features, an auxiliary node textual information dataset, $\mathcal{S} = \{S^{(\varphi)}, S^{(\chi)}, S^{(\psi)}\}$, is incorporated. Each of those entity textual information corresponds to the node in entity set \mathcal{V} . The $S^{(\varphi)} = [s_1^{(\varphi)}, s_2^{(\varphi)}, \dots, s_m^{(\varphi)}, \dots, s_M^{(\varphi)}]$, representing the entity names (e.g., HGNC symbol, Ensembl ID), $S^{(\chi)} = [s_1^{(\chi)}, s_2^{(\chi)}, \dots, s_m^{(\chi)}, \dots, s_M^{(\chi)}]$, representing the entity textual descriptions (e.g., Uniprot protein description), and $S^{(\psi)} = [s_1^{(\psi)}, s_2^{(\psi)}, \dots, s_m^{(\psi)}, \dots, s_M^{(\psi)}]$, representing biochemical information (e.g., biosequences or chemical structures, such as IChIKey). Therefore, for any entity, v_m , it has the textual information set $s_m = \{s_m^{(\varphi)}, s_m^{(\chi)}, s_m^{(\psi)}\}$. And the entity textual information dataset, \mathcal{S} , enhances the graph's expressivity, facilitating the generation of a textual-attributed transcriptomic signaling knowledge graph.

Afterwards, to construct the text-omic signaling graph, expressed as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, relations / edges between entities should be identified. As aforementioned, in this signaling graph, the vertex set is defined as $\mathcal{V} = \{\mathcal{V}^{(t)}, \mathcal{V}^{(p)}\}$. And two types of relations / edges, internal signaling and gene regulatory signaling, are selected. In details, the constructed signaling graph can be decomposed into two distinct subgraphs: the internal signaling subgraph, $\mathcal{G}^{(in)} = (\mathcal{V}^{(in)}, \mathcal{E}^{(in)})$, which encapsulates the molecular mechanisms governing protein translation, and the PPI-based gene regulatory subgraph, $\mathcal{G}^{(PPI)} = (\mathcal{V}^{(PPI)}, \mathcal{E}^{(PPI)})$, capturing protein-protein interactions, jointly composing the edge set $\mathcal{E} = \{\mathcal{E}^{(in)}, \mathcal{E}^{(PPI)}\}$. Specifically, $\mathcal{G}^{(in)}$ consists of all vertices such that $\mathcal{V}^{(in)} = \mathcal{V}$, with cardinality $|\mathcal{V}^{(in)}| = M = M_t + M_p$, while $\mathcal{G}^{(PPI)}$ is constrained to the protein nodes, i.e., $\mathcal{V}^{(PPI)} = \mathcal{V}^{(p)}$. Along with this, the dataset will be intergated as $\mathcal{D} = \{\mathcal{X}, \mathcal{S}, \mathcal{E}\}$ to be released.

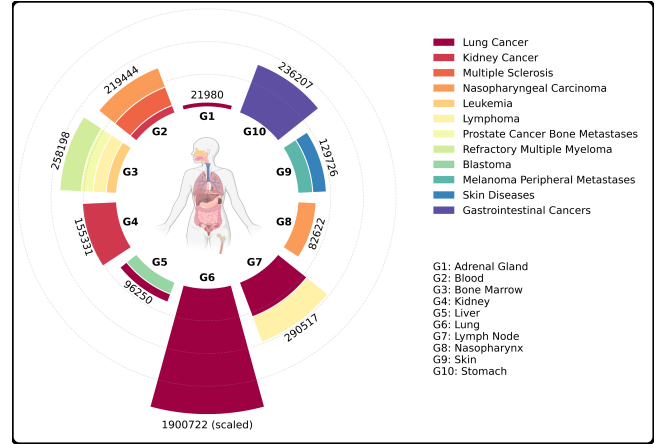


Figure 3: Overview of the filtered dataset, highlighting diseased cells from various organ groups after excluding normal cells and brain cells due to their high abundance. Each colored segment (G1 to G10) represents a distinct organ category, with numeric labels indicating the total number of cells retained in each group.

3.5 CellTOSGDataset Package

In the CellTOSGDataset package, the data matrix $\mathcal{X} \in \mathbb{R}^{N \times M}$ is formatted as a NumPy file. To optimize memory usage during processing, the input H5AD files are partitioned into 1024 MB chunks. Each partition within this size limit is then used to construct the corresponding x . npy and y . npy files. We extensively collected and curated multiple human single-cell datasets and employed meta-cell analysis to extract the core biological characteristics of cell groups. Starting with 117,519,978 raw cells, we distilled the data into a final set of 547,168 meta-cells (see **Figure 3** for the distribution of cells across organs). Detailed information regarding the organs and diseases is provided in **Table 1**. Due to the large volume of cells, the datasets are organized hierarchically by organ and disease. The final

datasets are available online via a Box folder². After downloading the files locally, users can load the data using the following Python code:

```

1 from dataset import CellTOSGDataset
2
3 CellTOSG = CellTOSGDataset(
4     root="./CellTOSG_dataset",
5     categories="get_organ_disease",
6     name="brain-AD",
7     label_type="status",
8     seed=2025,
9     ratio=0.01,
10    shuffle=True
11 )
12
13 x = CellTOSG.data
14 y = CellTOSG.labels
15 edge_index = CellTOSG.edge_index
16 internal_edge_index = CellTOSG.internal_edge_index
17 ppi_edge_index = CellTOSG.ppi_edge_index
18 s_name = CellTOSG.s_name
19 s_desc = CellTOSG.s_desc
20 s_bio = CellTOSG.s_bio
21
22 print(f"Data Load Finished, {len(xAll)} samples in
    total.")

```

This API extracts data from the specified root directory, where the full dataset is stored. The parameter categories determines the dataset subset to be retrieved. For instance, "get_organ_disease" indicates that the user wishes to obtain disease-specific cells from a given organ (e.g., brain-AD for Alzheimer's Disease cells from the brain). The label_type parameter accepts four options, ct, og, ds and status, which correspond to the four types to labels in the \mathcal{Y} . As to the ratio, this parameters will extracted this ratio of samples from whole candidate cells, since some files are pretty large and it will burst the memory storage. By using this ratio, we will sampling this ratio of cells from whole candidate cells. Aside from this, user can also shffule the data. To rebalance the dataset, we also implicitly intergrate the method for sampling the equal number of normal cells from same organs to serve as the control group, given that disease cell numbers are less than normal cells in OmniCellTOSG.

3.6 Joint LLM-GNN Cell Signaling Graph Foundation Model

3.6.1 CellTOSG Foundation Model Given the intergated text-omic signaling graph dataset \mathcal{D} , which contains the single cell text-omic signaling graph, denoted as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and its text-omic feature sets \mathcal{X}, \mathcal{S} , we can pretrain our foundation model by generating the node mask set $\mathcal{E}_{\text{mask}} \sim \text{Bernoulli}(p)$, where $p < 1$ is the ratio of the masked edges for set $\mathcal{E}^{(\text{PPI})}$ to mask out the signaling flows in protein-protein interactions, $\mathcal{E}^{(\text{PPI})}$. Then, the model will be pretrained by

$$\mathcal{H} = f_{\text{pre}}(\mathcal{X}, \mathcal{S}, \mathcal{E}, \mathcal{E}_{\text{mask}}) \quad (1)$$

²<https://wustl.box.com/s/6hr0yprwrmrkykldslw76etw6nluyw8>

, where $\mathcal{H} \in \mathbb{R}^{N \times M \times d}$ is the entity embeddings, and $f_{\text{pre}}(\cdot)$ is the pre-trained foundation model.

In details, to merge the text-omics feature sets \mathcal{X}, \mathcal{S} into unified entity embeddings, bi-encoder framework was leveraged by

$$\mathcal{X}' = \text{OmicEncoder}(\mathcal{X}) \quad (2)$$

$$\mathcal{S}' = \text{TextEncoder}(\mathcal{S}) \quad (3)$$

$$\mathcal{H}' = \text{CrossModalityEncoder}(\mathcal{X}', \mathcal{S}') \quad (4)$$

, where the OmicEncoder is the linear transformation and $\mathcal{X}' \in \mathbb{R}^{N \times M \times d'}$ the TextEncoder can be BERT-based or other LLMs and $\mathcal{S}' = \{S^{(\gamma)}, S^{(\theta)}, S^{(\rho)}\}$, where $S^{(\gamma)} \in \mathbb{R}^{M \times d'}$, $S^{(\theta)} \in \mathbb{R}^{M \times d'}$, $S^{(\rho)} \in \mathbb{R}^{M \times d'}$ are encoded as entity name, textual description and biochemical embeddings. The CrossModalityEncoder will fuse the omic embeddings and the textual embeddings with $\mathcal{H}' \in \mathbb{R}^{N \times M \times d'}$.

Afterwards, the internal signaling will be propagated by using graph encoder with

$$\mathcal{H}^{(\text{in})} = \text{GNN}_{\text{in}}(\mathcal{H}', \mathcal{E}^{(\text{in})}) \quad (5)$$

, where $\mathcal{H}^{(\text{in})} \in \mathbb{R}^{N \times M \times d}$. Finally, with the prepared entity embedding, the foundation model will be pretrained by masking nodes with

$$\mathcal{H} = \text{GNN}_{\text{pre}}(\mathcal{H}^{(\text{in})}, \mathcal{E}^{(\text{PPI})}, \mathcal{E}_{\text{mask}}) \quad (6)$$

3.6.2 Model Downstream Tasks Ultimately, the objective is to use the pretrained foundation model, $f_{\text{pre}}(\cdot)$, that synergistically integrates the incoming feature set $\mathcal{X}_0 \in \mathbb{R}^{N_0 \times M}$, node descriptions \mathcal{S} , and graph topology \mathcal{E} to predict cell-specific outcomes. As to the unsupervised task, the latent embedding for the incoming feature set will be generated by

$$\mathcal{H}^{(0)} = f_{\text{pre}}(\mathcal{X}_0, \mathcal{S}, \mathcal{E}) \quad (7)$$

, where $\mathcal{H}^{(0)} \in \mathbb{R}^{N_0 \times M \times d}$. With this latent embeddings, those N_0 will be clustered into K clusters.

For supervised learning, the foundation model will predict the cell outcomes by

$$\hat{\mathcal{Y}}_0 = \text{MLP}(f_{\text{pre}}(\mathcal{X}_0, \mathcal{S}, \mathcal{E})) \quad (8)$$

, where MLP is the linear classifier and $\hat{\mathcal{Y}}_0 \in \mathbb{R}^N$ represents the predicted cellular states, which depends on specific downstream tasks (e.g., cell type annotations or cellular condition (normal vs. disease)). Furthermore, the correponding inferenced core cell-specific signaling network will be generated based on the latent embeddings $\mathcal{H}^{(0)}$ by

$$\mathcal{A}^{(0)} = \text{ATT}(\mathcal{H}^{(0)}) \quad (9)$$

$$\mathcal{G}^{(0)} = f_{\text{core}}(\mathcal{A}^{(0)}, \delta) \quad (10)$$

, where ATT is the attention-based function to generate the entity similarity matrix $\mathcal{A}^{(0)} \in \mathbb{R}^{M \times M}$ and f_{core} is the core signaling inferring function for filtering out the edge lower than the threshold δ , resulting with the core signaling subgraph $\mathcal{G}^{(0)} = \{\mathcal{V}^{(0)}, \mathcal{E}^{(0)}\}$.

Table 2: Overall performance for cell types classification (CT) and cell status (Status) prediction for graph-based methods and CellTOSG-Class.

Model	Alzheimer’s Disease		Acute Myeloid Leukemia		Small Cell Lung Carcinoma		Clear Cell Renal Carcinoma	
	CT	Status	CT	Status	CT	Status	CT	Status
GCN	0.3225	0.5833	0.0714	0.9643	0.2745	0.5665	0.2667	0.4667
GIN	0.4365	0.3333	0.4286	0.9643	0.2331	0.8333	0.3222	0.4667
CellTOSG-Class	0.4365	0.5967	0.4286	0.9643	0.4526	0.5844	0.375	0.5333

4 Experiments and Results

Our OmniCellTOSG dataset serves as a robust benchmark for integrative text-omic analysis, illustrating that the incorporation of text-omic features and a graph language framework can significantly improve model performance over conventional approaches. By unifying textual embeddings (from LLM-based encoders) with GNN-based modeling of single-cell transcriptomes, OmniCellTOSG provides a systematic platform for the research community to develop, evaluate, and compare novel methodologies. Furthermore, users can easily load and customize these datasets by employing our Python API, where the categories, label_type, and ratio parameters facilitate flexible data selection, labeling, and sampling. To rebalance normal and disease cells within each organ, OmniCellTOSG also integrates an implicit sampling strategy that provides a balanced, scalable, and memory-efficient resource for advancing text-omic research.

As previously noted, the dataset is organized hierarchically by organ and disease, enabling users to generate novel benchmark subsets for various analytical objectives. In this work, we focus on Alzheimer’s disease in the brain (AD), acute myeloid leukemia (AML) in the bone marrow, renal cell carcinoma (RCC) in the kidney, and small cell lung carcinoma (SCLC) in the lung. The sampling ratios for these subsets are set to 0.01 for AD, 1.0 for AML, 0.2 (lung), and 0.1 (RCC), resulting in 120, 278, 252, and 296 samples, respectively. Meanwhile, we split the dataset into training and test dataset by ratio of 0.9. And we employ two label types—cell type and cell status—for downstream classification tasks with accuracy as the metric for model evaluation. Building on our pretrained model, we develop a downstream classifier, CellTOSG-Class, and compare its performance against state-of-the-art graph neural network (GNN) architectures (GCN [16], GAT [27], GIN [28], and UniMP [25]). For the textual encoder, we adopt DeBERTa [14] for entity names and descriptions, and leverage DNAGPT [31] and ProtGPT2 [9] for DNA/RNA and protein sequences, respectively—substituting thymine (T) with uracil (U) for RNA. **Table 2** illustrates that our model—pretrained using OmniCellTOSG—consistently outperforms competing models across a majority of evaluated tasks. This superior performance emphasizes the critical role of the data generation and training strategies inherent to OmniCellTOSG, demonstrating that these processes are pivotal in enhancing model predictive capabilities.

5 Discussions

Tissue-level and single-cell omic datasets are being generated to investigate disease pathogenesis, a cornerstone of precision medicine.

Graph neural networks (GNNs) have been widely applied for signaling network analysis—integrating omic data with signaling interactions—to identify key disease targets and infer pathways [8, 30, 32, 33]. Although these models have achieved superior predictive performance, current graph-based reasoning approaches for numeric omic cell signaling graphs only capture part of the scientific discovery process. Therefore, the prior knowledge and numeric data should be integrated during scientific discovery and knowledge reasoning. Therefore, in this study, we introduced the first large-scale single cell text-omic signaling graphs, OmniCellTOSG, by incorporating the human-understandable prior knowledge. Thus, the TOSGs represent a novel graph data model incorporating both text-attributed prior knowledge with numerical omic gene/protein abundance levels, which can facilitate the decoding of complex cell signaling systems. Moreover, novel paradigm-shift data analysis models like the joint LLM and GNN models are needed to analyze the TOSGs. In addition, the OmniCellTOSG consists of large-scale cell text-omic signaling graphs, using scRNAseq data, health vs diseases, and the number of cells in OmniCellTOSG keeps growing and will be updated regularly. Developing large-scale cell signaling foundation models is crucial. By pretraining on massive, diverse TOSG datasets from OmniCellTOSGs using self-supervised learning (SSL), these models can acquire a generalized understanding of complex signaling patterns and serve as robust bases for specialized adaptations. This approach outperforms disease-specific analyses—which capture only limited signaling scenarios and risk bias and overfitting—much like training ChatGPT-scale models on a small language dataset.

The OmniCellTOSG data are open-access, and organized in a Pytorch friendly format, and can facilitate the development of novel joint LLM and graph neural network (GNN) foundation cell signaling models for decoding the complex cell signaling systems via TOSG graph reasoning. It could shift the paradigm in life sciences, healthcare and precision medicine research. We are adding more TOSGs into the OmniCellTOSG to cover more essential factors, e.g., diseases, sex, age and other important factors, to facilitating the understanding of complex cell signaling systems and predicting potentially effective drugs and cocktails perturbing the dysfunctional cell signaling targets and pathways.

Acknowledgments

This research was partially supported by NLM 1R01LM013902-01A1, NIA R56AG065352, NIA 1R21AG078799-01A1 and NINDS 1RM1NS132962-01.

References

- [1] Ralph Abboud, Ismail Ilkan Ceylan, Martin Grohe, and Thomas Lukasiewicz. 2020. The surprising power of graph neural networks with random node initialization. *arXiv preprint arXiv:2010.01179* (2020).
- [2] Hervé Abdi and Lynne J Williams. 2010. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics* 2, 4 (2010), 433–459.
- [3] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big?. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*. 610–623.
- [4] Xinyue Chen, Yin Huang, Liangfeng Huang, Ziliang Huang, Zhao-Zhe Hao, Lahong Xu, Nana Xu, Zhi Li, Yonggao Mou, Mingli Ye, et al. 2024. A brain cell atlas integrating single-cell transcriptomes across human brain regions. *Nature Medicine* 30, 9 (2024), 2679–2691.
- [5] Haotian Cui, Chloe Wang, Hassaan Maan, Kuan Pang, Fengning Luo, Nan Duan, and Bo Wang. 2024. scGPT: toward building a foundation model for single-cell multi-omics using generative AI. *Nature Methods* 21, 8 (2024), 1470–1480.
- [6] C Domínguez Conde, C Xu, LB Jarvis, DB Rainbow, SB Wells, T Gomes, SK Howlett, O Suchanek, K Polanski, HW King, et al. 2022. Cross-tissue immune cell analysis reveals tissue-specific features in humans. *Science* 376, 6594 (2022), eabl5197.
- [7] Zehao Dong, Muhan Zhang, Philip RO Payne, Michael A Province, Carlos Cruchaga, Tianyu Zhao, Fuhai Li, and Yixin Chen. 2023. Rethinking the power of graph canonicalization in graph representation learning with stability. *arXiv preprint arXiv:2309.00738* (2023).
- [8] Zehao Dong, Qihang Zhao, Philip RO Payne, Michael A Province, Carlos Cruchaga, Muhan Zhang, Tianyu Zhao, Yixin Chen, and Fuhai Li. 2023. Highly accurate disease diagnosis and highly reproducible biomarker identification with PathFormer. *Research Square* (2023), rs–3.
- [9] Noelia Ferruz, Steffen Schmidt, and Birte Höcker. 2022. ProtGPT2 is a deep unsupervised language model for protein design. *Nature communications* 13, 1 (2022), 4348.
- [10] Xi Fu, Shentong Mo, Alejandro Buendia, Anouchka P Laurent, Anqi Shao, Maria del Mar Alvarez-Torres, Tianji Yu, Jimin Tan, Jiayu Su, Romella Sagatelian, et al. 2025. A foundation model of transcription across human cell types. *Nature* (2025), 1–9.
- [11] Mariano I Gabbito, Kyle J Travaglini, Victoria M Rachleff, Eitan S Kaplan, Brian Long, Jeanelle Ariza, Yi Ding, Joseph T Mahoney, Nick Dee, Jeff Goldy, et al. 2024. Integrated multimodal cell atlas of Alzheimer’s disease. *Nature Neuroscience* 27, 12 (2024), 2366–2383.
- [12] Anna K Greenwood, Kelsey S Montgomery, Nicole Kauer, Kara H Woo, Zoe J Leanza, William L Poehlman, Jake Gockley, Solveig K Sieberts, Ljubomir Bradic, Benjamin A Logsdon, et al. 2020. The AD knowledge portal: a repository for multi-omic data on Alzheimer’s disease and aging. *Current protocols in human genetics* 108, 1 (2020), e105.
- [13] Minsheng Hao, Jing Gong, Xin Zeng, Chiming Liu, Yucheng Guo, Xingyi Cheng, Taifeng Wang, Jianzhu Ma, Xuegong Zhang, and Le Song. 2024. Large-scale foundation model on single-cell transcriptomics. *Nature methods* 21, 8 (2024), 1481–1491.
- [14] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. DeBERTa: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654* (2020).
- [15] Graham Heimberg, Tony Kuo, Daryle J DePianto, Omar Salem, Tobias Heigl, Nathaniel Diamant, Gabriele Scalia, Tommaso Biancalani, Shannon J Turley, Jason R Rock, et al. 2024. A cell atlas foundation model for scalable search of similar human cells. *Nature* (2024), 1–3.
- [16] Thomas N Kipf and Max Welling. 2016. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907* (2016).
- [17] Ed S Lein and Erin E Gray. 2024. Seattle Alzheimer’s Disease Brain Cell Atlas (SEA-AD): A multi-faceted platform for discovery of cellular and molecular perturbations underlying Alzheimer’s disease. In *Alzheimer’s Association International Conference*. ALZ.
- [18] Hansruedi Mathys, Zhuyu Peng, Carles A Boix, Matheus B Victor, Noelle Leary, Sudhagar Babu, Ghada Abdelhady, Xueqiao Jiang, Ayesha P Ng, Kimia Ghafari, et al. 2023. Single-cell atlas reveals correlates of high cognitive function, dementia, and resilience to Alzheimer’s disease pathology. *Cell* 186, 20 (2023), 4365–4385.
- [19] Colin McGill, Bruce Martin, Charlotte Weaver, Sidney Bell, Lia Prins, Seve Badajoz, Brian McCandless, Angela Oliveira Pisco, Marcus Kinsella, Fiona Griffin, et al. 2021. Cellxgene: a performant, scalable exploration platform for high dimensional sparse matrices. *BioRxiv* (2021), 2021–04.
- [20] Jeremy A Miller, Michael J Hawrylycz, Matthew Aitken, Jeannelle Ariza, Rushil Chakrabarty, Song-Lin Ding, Yi Ding, Rebecca Ferrer, Jeff Goldy, Sergey Gratiy, et al. 2023. SEA-AD: Scientific analysis and open access resources targeting early changes in Alzheimer’s disease. *Alzheimer’s & Dementia* 19 (2023), e063478.
- [21] Sitara Persad, Zi-Ning Choo, Christine Dien, Noor Sohail, Ignas Masilionis, Ronan Chaligné, Tal Nawy, Chrysothemis C Brown, Roshan Sharma, Itzik Pe’er, et al. 2023. SEACells infers transcriptional and epigenomic cellular states from single-cell genomics data. *Nature Biotechnology* 41, 12 (2023), 1746–1757.
- [22] Leif E Peterson. 2009. K-nearest neighbor. *Scholarpedia* 4, 2 (2009), 1883.
- [23] CZI Cell Science Program, Shihla Abdulla, Brian Aevermann, Pedro Assis, Seve Badajoz, Sidney M Bell, Emanuele Bezzi, Batuhan Cakir, Jim Chaffer, Signe Chambers, et al. 2025. CZ CELLxGENE Discover: a single-cell data platform for scalable exploration, analysis and modeling of aggregated data. *Nucleic Acids Research* 53, D1 (2025), D886–D900.
- [24] Jennifer E Rood, Samantha Wynne, Lucia Robson, Anna Hupalowska, John Randlell, Sarah A Teichmann, and Aviv Regev. 2024. The Human Cell Atlas from a cell census to a unified foundation model. *Nature* (2024), 1–2.
- [25] Yunsheng Shi, Zhengjie Huang, Shikun Feng, Hui Zhong, Wenjin Wang, and Yu Sun. 2020. Masked label prediction: Unified message passing model for semi-supervised classification. *arXiv preprint arXiv:2009.03509* (2020).
- [26] Christina V Theodoris, Ling Xiao, Anant Chopra, Mark D Chaffin, Zeina R Al Sayed, Matthew C Hill, Helene Mantineo, Elizabeth M Brydon, Zexian Zeng, X Shirley Liu, et al. 2023. Transfer learning enables predictions in network biology. *Nature* 618, 7965 (2023), 616–624.
- [27] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903* (2017).
- [28] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2018. How powerful are graph neural networks? *arXiv preprint arXiv:1810.00826* (2018).
- [29] Zizhen Yao, Cindy TJ van Velthoven, Michael Kunst, Meng Zhang, Delissa McMillen, Changkyu Lee, Won Jung, Jeff Goldy, Aliya Abdelhak, Matthew Aitken, et al. 2023. A high-resolution transcriptomic and spatial atlas of cell types in the whole mouse brain. *Nature* 624, 7991 (2023), 317–332.
- [30] Lei Yu, Lei Liu, Zixian Zhang, Hengchang Zhang, and Xing Chen. 2021. MOGONET integrates multi-omics data using graph convolutional networks allowing patient classification and biomarker identification. *Nature Communications* 12, 1 (October 2021), 6510. <https://doi.org/10.1038/s41467-021-26624-7>
- [31] Daoan Zhang, Weitong Zhang, Bing He, Jianguo Zhang, Chenchen Qin, and Jianhua Yao. 2023. DNAGPT: a generalized pretrained tool for multiple DNA sequence analysis tasks. *bioRxiv* (2023), 2023–07.
- [32] Heming Zhang, Yixin Chen, Philip Payne, and Fuhai Li. 2024. Using DeepSignalFlow to mine signaling flows interpreting mechanism of synergy of cocktails. *npj Systems Biology and Applications* 10, 1 (2024), 92.
- [33] Heming Zhang, S Peter Goedegebuure, Li Ding, David DeNardo, Ryan C Fields, Michael Province, Yixin Chen, Philip Payne, and Fuhai Li. [n.d.]. M3NetFlow: A multi-scale multi-hop graph AI model for integrative multi-omic data analysis. *iScience* ([n.d.]).
- [34] Heming Zhang, Shunning Liang, Tim Xu, Wenyu Li, Di Huang, Yuhan Dong, Guangfu Li, J Philip Miller, S Peter Goedegebuure, Marco Sardiello, et al. 2024. BioMedGraphica: An All-in-One Platform for Biomedical Prior Knowledge and Omic Signaling Graph Generation. *bioRxiv* (2024), 2024–12.

A Datasets collection

A.1 Data Sources and Download

CZ CellxGene Database Data from the CellxGene database was obtained using the CZI Science CELLxGENE Census Python API, with the census version set to '2023-05-15'. The data was downloaded in H5AD AnnData format. To minimize duplicate entries, SEA-AD-related data from this dataset was removed.

GEO Database. Data was downloaded from the Gene Expression Omnibus (GEO) database using automated shell scripts. The download links were structured in a standardized format: [GSM_ID] [Directory]/[Filename] [FTP_URL]. The data was available in two formats:

- Matrix Market format: consisting of barcodes.tsv.gz, features.tsv.gz, and matrix.mtx.gz files
- Compressed CSV format (csv.gz)

The download process was automated through a script that processes a links.txt file containing the download information (full version available at our GitHub repository³). This systematic approach ensured reliable data collection across all GEO datasets with automated retry mechanisms and SSL verification.

³<https://github.com/FuhaiLiAiLab/OmniCellTOSG>

Table 3: Example Cell Type Prediction Output

Cell ID	Genes	Predicted Label	Majority Vote	Conf.
AAACCCAAGATTGACA-1	215	L2-3 CUX2 NTNG1 PALMD	Oligo MOG OPALIN	0.297
AAACCCAAGCCTGCCA-1	236	Oligo MOG OPALIN	Endo CLDN5 SLC7A5	0.703
AAACCCAAGCGATCGA-1	3071	L6 OPRK1 THEMIS RGS6	L6 OPRK1 THEMIS RGS6	0.997
...				
TTTGTGTGTCGTCAGAT-1	4557	InN SST FREM1	InN SST FREM1	0.998
TTTGTGTGCTGGAGAG-1	4482	InN SST THSD7B	InN SST FREM1	1.000
TTTGTGTGCTTGGCTC-1	231	Oligo MOG OPALIN	Oligo MOG OPALIN	0.077

```

1 download_file() {
2     local project=$1
3     local file_name=$2
4     local file_url=$3
5
6     # Create directory
7     local dir=$(dirname "$file_name")
8     [[ -d $dir ]] || mkdir -p "$dir"
9
10    # Attempt download with retries
11    local curl_times=3
12    while [ $curl_times -gt 0 ]; do
13        curl --cacert /etc/ssl/certs/ca-
14            certificates.crt \
15            -C - -L -o "$file_name" "$file_url"
16        if [[ $? -eq 0 ]]; then
17            return
18        fi
19        curl_times=$((curl_times - 1))
20        sleep 1
21    done
22 }
23
24 # Process links.txt file
25 while read -r line || [[ -n "$line" ]]; do
26     set -- $line
27     download_file "$1" "$2" "$3"
28 done < "links.txt"
29
30 # Example links.txt format:
31 # GSM5706853 GSM5706853_P1_barcode.tsv.gz https
32   //ftp.ncbi.nlm.nih.gov/geo/samples/GSM5706nnn
33   /GSM5706853/suppl/GSM5706853_P1_barcode.tsv.
34   gz

```

Brain Cell Atlas. Data was manually downloaded from the dataset page of the Brain Cell Atlas project after setting the species filter to "Human". Since the processed data lacks unique identifiers, the original source dataset project IDs retained in the processed H5AD AnnData files were recorded to document the data sources used.

SEA-AD Database. Data from the Seattle Alzheimer's Disease Brain Cell Atlas (SEA-AD) consortium was accessed through their API using the Synapse client. Authentication was handled via personal access tokens, with data retrieved through the synID based on different folders, and data was downloaded in H5 format files containing raw feature-barcode matrices. Data access was managed through the Synapse Python client, which handled authentication

and maintained data provenance. Each downloaded dataset was organized in a consistent directory structure to facilitate subsequent processing steps.

```

1 import synapseclient
2 import synapseutils
3 syn = synapseclient.Synapse()
4 syn.login(authToken=token)
5 folder_ids = ["syn26273710", "syn51792375", "
6   syn52314491", "syn52314469", "syn61680896", "
7   syn52314488", "syn52314472"]
8 for folder_id in folder_ids:
9     synapseutils.syncFromSynapse(syn, folder_id)

```

A.2 Data Preprocessing

We developed a standardized preprocessing process to convert all datasets into a unified H5AD or H5 format. These formats, which efficiently store large-scale genomic data, are directly compatible with CellTypist, our chosen tool for cell type annotation. Both the H5AD format and the H5 format are hierarchical data formats designed to optimize the storage of large scientific datasets, with H5AD specifically designed for annotation matrices in single-cell genomics. These formats work well at efficiently processing large-scale data and maintaining relationships between genes, cells, and their annotations, making them ideal for work.

GEO Data Processing. For Matrix Market files, we utilized Scanpy's `read_10x_mtx` function to combine the three separate files (barcodes, features, matrix) into a single AnnData object. This was then saved as an H5AD file to preserve the complete data structure and annotations. For CSV files, we transformed the expression matrices into AnnData objects while preserving the gene-cell relationships, followed by conversion to H5AD format for consistent data handling. Gene name standardization was performed using reference gene list from scFoundation[13]. Quality filtering was applied with a minimum threshold of 100 genes per cell.

SEA-AD Data Processing. We directly loaded H5 files using Scanpy's `read_10x_h5` function. Duplicate genes were handled through unique name generation. Quality filtering was applied with a minimum threshold of 100 genes per cell.

CellxGene and Brain Cell Atlas Data Processing. Unlike GEO and SEA-AD datasets, the CellxGene and Brain Cell Atlas datasets were already provided in H5AD AnnData format. Therefore, no additional conversion was needed.

Table 4: Organ/Tissue Types and Disease Details (Fewer than 800 Meta-Cells)

... Continuation of **Table 1** ...

Diseases	Organ/Tissue Types	# of Original Cells	# of Result Cells
Autism Spectrum Disorder	Brain	52,003	259
Blastoma	Liver	57,445	287
Epilepsy	Brain	126,587	626
Frontotemporal Dementia	Brain	138,395	691
Leukemia	Bone Marrow	46,918	213
Major Depressive Disorder	Brain	41,944	209
Melanoma Brain Metastases	Brain	76,643	339
Melanoma Peripheral Metastases	Skin	67,581	334
Multiple Sclerosis	Blood, Brain	159,175	787
Prostate Cancer Bone Metastases	Bone Marrow	53,658	190
Refractory Multiple Myeloma	Bone Marrow	97,876	371
Skin Diseases	Skin	62,145	308

A.3 Cell Type Annotation

We employed CellTypist, a supervised cell type classification tool, for automated cell type annotation. A mapping strategy was developed to match each dataset with the most appropriate CellTypist model based on tissue origin and disease context. For instance, brain tissue samples were processed using the "Developing_Human_Brain" model, while immune-related samples utilized immune cell-specific models. After that, data preprocessing steps include normalization (scale factor: $1e4$) and log-transformation while maintaining matrix sparsity. Majority voting was implemented to resolve any conflicting predictions.

The processed data was stored in the H5AD format, which preserves raw expression counts, cell type annotations, quality metrics, and confidence scores for cell type predictions. Each dataset contains a standard set of fields including predicted labels, majority voting results, and confidence scores for reproducibility (see example in **Table 3**).

A.4 OmniCellTOSG Dataset Overview (Continuation)

See **Table 4** for details.