

OmniGeo: Towards a Multimodal Large Language Models for Geospatial Artificial Intelligence

Long Yuan*
Beijing Jiaotong University
Beijing, China
24120408@bjtu.edu.cn

Fengran Mo*
University of Montreal
Montreal, Quebec, Canada
fengran.mo@umontreal.ca

Kaiyu Huang†
Beijing Jiaotong University
Beijing, China
kyhuang@bjtu.edu.cn

Wenjie Wang
Beijing Jiaotong University
Beijing, China

Wangyuxuan Zhai
Beijing Jiaotong University
Beijing, China

Xiaoyu Zhu
Beijing Jiaotong University
Beijing, China

You Li
Beijing Jiaotong University
Beijing, China

Jinan Xu
Beijing Jiaotong University
Beijing, China

Jian-Yun Nie
University of Montreal
Montreal, Quebec, Canada

Abstract

The rapid advancement of multimodal large language models (LLMs) has opened new frontiers in artificial intelligence, enabling the integration of diverse large-scale data types such as text, images, and spatial information. In this paper, we explore the potential of multimodal LLMs (MLLM) for geospatial artificial intelligence (GeoAI), a field that leverages spatial data to address challenges in domains including Geospatial Semantics, Health Geography, Urban Geography, Urban Perception, and Remote Sensing. We propose a MLLM (OmniGeo) tailored to geospatial applications, capable of processing and analyzing heterogeneous data sources, including satellite imagery, geospatial metadata, and textual descriptions. By combining the strengths of natural language understanding and spatial reasoning, our model enhances the ability of instruction following and the accuracy of GeoAI systems. Results demonstrate that our model outperforms task-specific models and existing LLMs on diverse geospatial tasks, effectively addressing the multimodality nature while achieving competitive results on the zero-shot geospatial tasks. Our code will be released after publication.

Keywords

Multimodal large language model, Geospatial artificial intelligence

ACM Reference Format:

Long Yuan, Fengran Mo, Kaiyu Huang, Wenjie Wang, Wangyuxuan Zhai, Xiaoyu Zhu, You Li, Jinan Xu, and Jian-Yun Nie. 2025. OmniGeo: Towards a Multimodal Large Language Models for Geospatial Artificial Intelligence. In . ACM, New York, NY, USA, 15 pages. <https://doi.org/XXXXXXX.XXXXXXX>

*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.
Conference'17, Washington, DC, USA

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-XXXX-X/2018/06
<https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Recent advancements in machine learning (ML) and artificial intelligence (AI) underscore the immense potential of data and computational power [20, 49, 63]. Exceptionally large language models (LLMs), trained on large-scale datasets, have demonstrated remarkable effectiveness in a wide range of learning tasks [11, 47, 69]. In particular, the unprecedented success of LLMs has catalyzed a paradigm shift in the training of ML and AI models. Instead of developing task-specific models from scratch for each individual task, the generality and adaptability of LLMs are leveraged to perform multiple tasks simultaneously with a single model through few-shot/zero-shot learning techniques [40]. This enables deployment across a diverse range of domains [26, 44], such as healthcare [33, 34], education [16], law [29, 39], and finance [4, 58]. LLMs have encapsulated the knowledge embedded in their training corpus, which encompasses billions or even trillions of tokens sourced from the Internet [43]. Therefore, we aim to investigate approaches to extract the geospatial knowledge that LLMs possess, enhancing a variety of geospatial machine learning tasks.

However, as shown in Figure 1, geospatial tasks typically include sequence labeling, time series prediction forecasting, geospatial image classification, and tasks related to urban functions [40]. Thus, geographic science is an inherently complex discipline that encompasses a wide range of tasks, indicating that addressing geospatial tasks requires the integration of multiple modalities. The primary technical challenge in Geospatial Artificial Intelligence (GeoAI) lies in its inherently multimodal nature. The data modalities in GeoAI include text, images (e.g., RS images or street view images), trajectory data, knowledge graphs, and geospatial vector data, which encapsulate critical geospatial information [8, 38, 42, 51]. Each modality exhibits unique structural characteristics that demand distinct representations. As a result, the multimodal nature of GeoAI complicates the straightforward application of existing LLMs across the full spectrum of GeoAI tasks. Few existing studies integrate these diverse representations while incorporating suitable inductive biases into a single model requires careful and thoughtful design.

In this paper, we propose a unified MLLM that leverages diverse geographic spatial data to guide the implicit geospatial knowledge

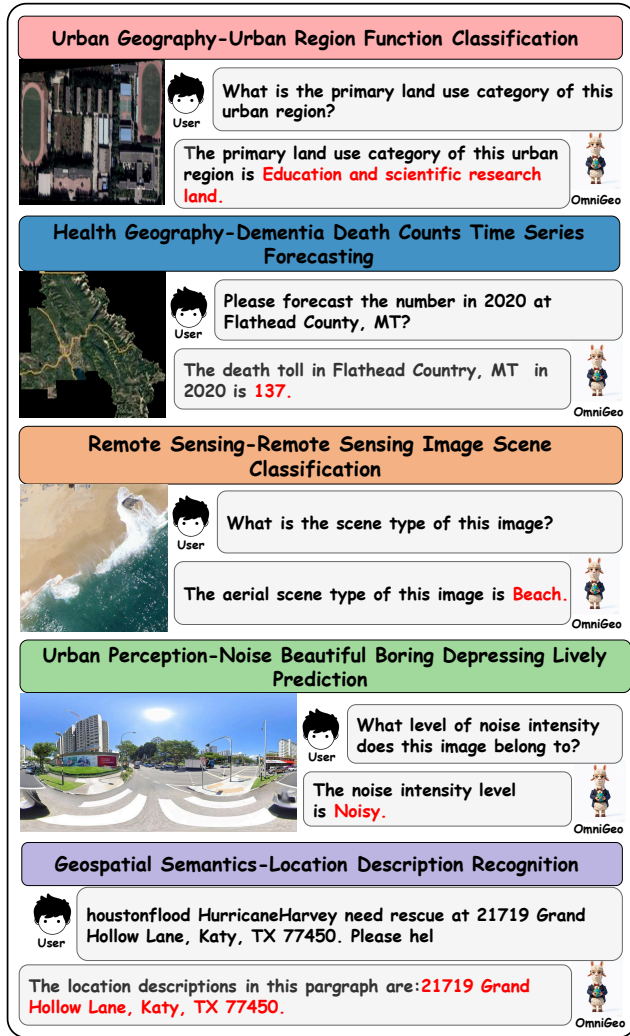


Figure 1: The Illustration of all the tasks covered by OmniGeo through an engaging dialogue.

learned by the foundation model during unsupervised pre-training, allowing it to better understand geoscience tasks and provide accurate responses. The data utilized include geographic textual data from tweets and webpages, spatial polygon vector data from CDC Wonder (dementia data at the state and county level), spatial point and polygon vector data from Gaode Maps (points of interest, POI), RS images from WorldView, street view images from Google Maps, and spatial polygon vector data from city government websites (urban planning layer data). Initially, because LLMs are unable to directly comprehend the spatial distribution of large-scale POI data of the region, we convert the geographic spatial vector data (POI and dementia data) into text paragraphs, crop RS images by region, and generate captions for both RS and street view images. Subsequently, the aforementioned geographic spatial data are aligned with geographic entities to create a multimodal instruction fine-tuning dataset. Finally, we obtain OmniGeo using these instruction data in LoRA and full-supervision fine-tuning manner. This

model achieves commendable performance across multiple geospatial tasks, proving that the model have successfully extracted the implicit geospatial knowledge from the base model, thereby enabling more accurate geospatial task inference.

Our contributions are summarized as (1) We propose OmniGeo, a MLLM to facilitate the integration of geospatial information from different modalities, mitigating interference between tasks, and enabling a single large model to simultaneously handle multiple heterogeneous geospatial tasks. (2) OmniGeo has constructed 12 multimodal instruction fine-tuning data for GeoAI and achieves competitive performance on geospatial tasks. (3) To the best of our knowledge, OmniGeo is the first large-scale multimodal language model that covers the five core tasks in the GeoAI field (health geography, urban geography, RS images, urban perception, and semantic analysis), taking the step towards the multimodal development of GeoAI.

2 Related Work

Geospatial Artificial Intelligence. For toponym recognition and location description recognition, it is regarded as a subtask of Named Entity Recognition (NER), with the goal of identifying named places or location descriptions from text. early methods utilized general NER tools such as Stanford NER [14] and spaCy NER [21] to uniformly identify geospatial semantics. Wang et al. [53] were the first to design the NeuroTPR, an RNN-based model which can extract location information from text.

For urban region function classification, it aims to classify the primary functional categories of urban region based on various geospatial data within the target region, which is beneficial for city planning and resource allocation. early methods [2, 64] primarily relied on researchers' subjective analysis of urban planning maps and human mobility data to determine urban region function categories. With the acquisition and application of multi-source heterogeneous geographic data, Jing et al. [30] developed an RS image semantic interpretation model to analyze the characteristics of specific urban functional regions, whereas Qi et al. [50] performed a qualitative analysis of local urban region functions using GPS-based taxi mobility data. Existing research, including Place2Vec [61] and HGI [27], has developed region-specific semantic embeddings, to effectively carry out various downstream tasks.

For RS image scene classification, Yao et al. [60] leverage stacked sparse autoencoders to classify scenes via learning a large number of discriminative image features. To transfer the pre-trained CNN to RS image classification tasks, He et al. [19] designed the MSCP algorithm to automatically select and combine multi-level feature maps extracted from the pre-trained CNN.

Multimodal Large Language Models for Domain Tasks. MLLM research in GeoAI is primarily focused on handling RS images. To understand objects with arbitrary sizes and orientations, Wang et al. [52] designed the Remote-Sensing-RVSA with a new rotation-invariant scalable attention mechanism, achieving SOTA performance in several RS image visual tasks. Pang et al. [48] developed a large-scale RS image-text dataset, Hnstd, to perform a variety of common RS image tasks and show clear improvements over previous VLMs. Kartik Kuckreja et al. [32] constructed a new RS

multimodal instruction-following dataset to train models capable of performing region-level reasoning, and named the model GeoChat.

Different from the above work, we hope to achieve a new paradigm in which a single model can cover all core tasks of GeoAI, in an attempt to alleviate the inherent challenges brought by multimodal complex data.

3 Preliminaries

3.1 Task Definition

Toponym Recognition and Location Description Recognition:

Toponym Recognition is a subtask of NER. The challenge of this task lies in correctly distinguishing geographically similar entities that appear as names of coarse-grained locations (such as cities, states, and countries). Formally, given a text $T = \{t_1, t_2, \dots, t_n\}$, the goal is to assign each token t_i a label $y_i \in \{B\text{-LOC}, I\text{-LOC}, O\}$, where

$$y_i \in \begin{cases} B\text{-LOC} & \text{if } t_i \text{ is the beginning of a location description,} \\ I\text{-LOC} & \text{if } t_i \text{ is inside a location description,} \\ O & \text{if } t_i \text{ is outside a location description.} \end{cases} \quad (1)$$

Moreover, the goal of Location Description Recognition is to identify the location description that appears within a given text. The challenge of this task lies in whether it can fully output more fine-grained location descriptions, such as home addresses, highway exits, and road intersections, rather than large-scale geographic entities like cities, states, and countries. Formally, given a text $T = \{t_1, t_2, \dots, t_n\}$, the task is to assign each token t_i a label $y_i \in \{B\text{-LDR}, I\text{-LDR}, O\}$, where

$$y_i \in \begin{cases} B\text{-LDR} & \text{if } t_i \text{ is the beginning of a location description,} \\ I\text{-LDR} & \text{if } t_i \text{ is inside a location description,} \\ O & \text{if } t_i \text{ is outside a location description.} \end{cases} \quad (2)$$

Dementia Death Counts Time Series Forecasting: Given historical dementia mortality time series data $\{Y_t^r, Y_{t-1}^r, \dots, Y_1^r\}$ for a specific region r (e.g., state, county) and the corresponding RS image RS_r for the region at a fixed time t , the goal is to predict the mortality at the next time point, Y_{t+1}^r . The challenge of this task lies in whether it can eliminate "geographic bias" [13] and reasonably estimate the death count, considering the geographical distribution differences [1] in death rate growth. The task can be formally defined as

$$Y_{t+1}^r = f(\{Y_t^r, Y_{t-1}^r, \dots, Y_1^r\}, RS_r) \quad (3)$$

where Y_t^r represents the number of deaths in region r at time t , RS_r is the RS image of the region r .

Urban Region Function Classification: The goal of this task is to determine the human activity patterns within the region based on the distribution of POI data and regional RS images, and correctly estimate the primary urban function category. Formally, this can be expressed as:

$$\hat{F}_r = f(\mathbf{P}_r, RS_r), \quad (4)$$

where \hat{F}_r is the predicted primary urban function of region r and \mathbf{P}_r is a set representing the counts of different types of POIs in region r , i.e., $\mathbf{P}_r = (p_{r1}, p_{r2}, \dots, p_{rn})$. The p_{ri} and RS_r denote the number of POIs of type i and the RS image of the region r .

Remote Sensing Image Scene Classification: The goal of this task is to correctly distinguish highly similar scene types based on

information such as the types, quantities, and layouts of geographic entities in the RS images. Formally, it can be expressed as $\hat{C}_i = f(RS_i)$. The \hat{C}_i and RS_i are the predicted scene and the original type of the RS image i .

Urban Perception Prediction: Based on street view images, this task involves judging the fine-grained perceptual features of urban neighborhoods from seven perceptual indicators (Noise, Beautiful, Boring, Depressing, Lively, Safe, Wealthy), with each indicator subdivided into four perceptual features. The challenge of this task lies in whether the model possesses advanced human perception knowledge [65, 66]. Formally, this can be expressed as:

$$\hat{P}_i = f(SVI_i), \quad (5)$$

where \hat{P}_i is the predicted perception feature vector of the urban region in image i . Specifically, \hat{P}_i consists of seven perception indicators (**Noise, Beautiful, Boring, Depressing, Lively, Safe, Wealthy**), and the prediction for each indicator is the maximum value among its four different levels of perceptual features:

$$\hat{P}_i = \max(P_i^1, P_i^2, P_i^3, P_i^4), \quad \text{for each } i \in \{1, 2, \dots, 7\}, \quad (6)$$

where P_i^k ($k = 1, 2, 3, 4$) represents the k -th level feature of the i -th perception indicator. SVI_i is the street view image i .

3.2 Exploration of LLMs for GeoAI

As a starting point for this study, we empirically demonstrate the potential of leveraging LLMs for addressing geospatial tasks. We aim to show that our investigation not only highlight the effectiveness of general-purpose, few-shot learners in the geospatial domain, but also challenge the prevailing paradigm of training task-specific models as a standard practice in GeoAI research. To this end, we compare the performance of specific ML models and LLMs. These models are evaluated against multiple supervised, task-specific baselines on two representative geospatial tasks: (1) toponym recognition and (2) RS image scene classification.

Table 1: Accuracy of the baseline models on two representative geospatial tasks.

Model	Toponym Recognition	Remote Sensing Image Classification
Stanford NER	0.757	-
AlexNet	-	0.812
LLaVA	0.708	0.648
GPT-4	0.731	0.794

As shown in Table 1, GPT-4 achieves promising results, demonstrating that a significant amount of geospatial knowledge embedded in LLMs. This suggests that using instruction data to guide foundation models in constructing specialized geospatial LLMs holds significant potential. However, extracting this knowledge from LLMs is a non-trivial task. LLMs are unable to perform spatial reasoning in a way that is grounded in the real world. Therefore, it is critical to align geospatial information from different modalities to uncover accurate geospatial knowledge during inference.

4 Methodology

In this section, we first describe the sources of geospatial data and the process of constructing multimodal instruction fine-tuning

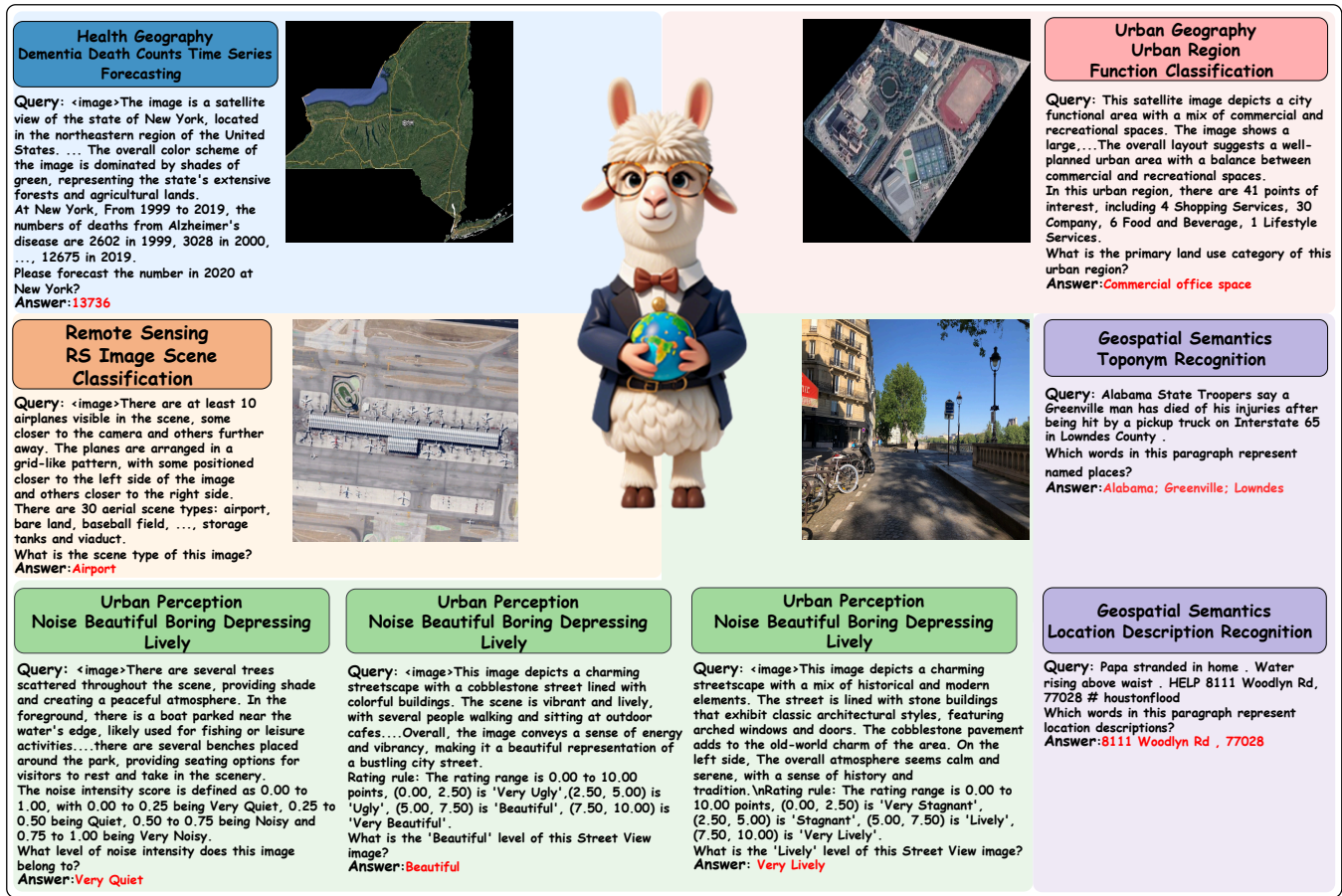


Figure 2: A detailed illustration of the image-text instruction data and geospatial tasks covered by OmniGeo.

datasets. We then propose the use of multi-task joint training to activate the inherent geospatial knowledge within MLLMs, thereby enhancing their instruction-following capabilities and improving accuracy in geospatial tasks

4.1 Multimodal Instruction Data Construction

In GeoAI, RS images provide a macro perspective of urban spatial structure or geographical locations, while SVI depict a more granular view of urban street blocks. Various tasks in GeoAI are closely associated with different geographic entities, and captions for remote sensing and street view images will further enhance this collaborative information. Therefore, the core motivation behind constructing a fine-tuning dataset is to equip each task with images (either RS or SVI) and complement them with detailed captions. This will assist in spatial reasoning and the 3D perception of geographic entities by the model.

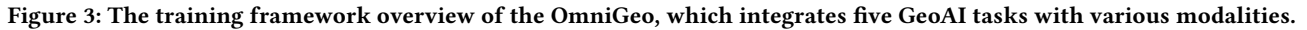
Data sources. In GeoAI, the primary data structures encompass text, images (e.g., remote sensing and street view images), geographic knowledge graphs, geospatial vector data, videos, and audio. This paper includes geographic text data from tweets and web pages, geospatial polygon vector data from CDC Wonder (dementia data at the state and county levels in the US), geospatial point and

polygon vector data from Amap (POI), remote sensing images from WorldView, street view images from Google Maps, and geospatial polygon vector data from city government websites (urban planning layer data).

Data construction. Based on the aforementioned data sources, we process all the data to construct instruction datasets that can be utilized by LLMs. The specific construction methods for different tasks and data modalities are outlined as follows:

For **Health Geography**, historical death statistics are first retrieved from CDC Wonder and converted into text. Data from 1999 to 2019 are used for model training, while data from 2020 are reserved for evaluation. Moreover, we add the image data (RS) as supplementary information for the Health Geography to provide modal alignment. In particular, the WorldView series of RS images are cropped region by region, and captions are generated using InstructBLIP [10], which are subsequently concatenated with the statistical data. This process results in two datasets: one at the US state level and the other at the US county level.

For **Urban Geography**, POI data for Beijing and Shenzhen are initially retrieved from Gaode Maps and classified into 11 categories. Functional region partition layer data for these cities are then obtained from government websites, and the POI data are converted



For **Geospatial Semantics**, textual resources containing geospatial information are abundant. Thus, existing geospatial benchmark datasets are directly employed. LGL, GeoVirus, and NEEL are used for toponym recognition, while HaveyTweet2017 and GeoCorpora are focused on location description recognition. The contents are mainly gathered from web pages, tweets, and other sources.

Urban Region Function Classification: As shown in the red block in Figure 2, the goal of this task is to model potential human activity patterns based on POI data from the target region. In particular, we integrate RS images and their associated captions to capture high-level semantic information about specific regions, such as the spatial arrangement of urban geographical entities, including shopping malls, schools, and residential districts. This approach enables OmniGeo to perceive the distribution of POIs within the region from a broader, more comprehensive perspective. The output of the model represents the primary functional category of the region.

RS Image Scene Classification: As shown in the orange block in Figure 2, the model processes Rs images and their corresponding captions for specific scenes, with the goal of capturing unified high-level semantics across scenes exhibiting the same geographic context but varying spatial and spectral resolutions. Additionally, the model is designed to differentiate challenging samples, where geographic entities remain consistent but scene types diverge.

Urban Perception Prediction: As shown in the green block in Figure 2, the urban perception prediction is different from typical image classification, which typically classifies images based solely on visual instances. Instead, the model is required to simulate high-level human perception abilities by analyzing street view images and their associated captions. It estimates various perceptual dimensions, including Noise, Beauty, Boredom, Depression, Livelihood, Safety, and Wealth, where high-level perceptual knowledge is more complex and challenging to assess. The output of OmniGeo consists of one of four distinct levels for each of the seven perceptual indicators.

Toponym Recognition: The language capabilities of advanced LLMs are sufficiently powerful, yet extracting frequently occurring and easily confused place names within the same text remains difficult. The output of OmniGeo is a semicolon-separated list of toponyms.

Location Description Recognition: This task presents greater challenges than toponym recognition, as real-world location descriptions exhibit considerable disarray due to factors such as cultural differences, varying educational levels, and diverse lifestyle habits, particularly in relation to the handling of special symbols. The output is a list of location descriptions separated by semicolons.

4.3 Training Objective

To further align and internalize geospatial knowledge with the intrinsic knowledge of MLLMs, we select LLaVA1.5 and Qwen2-VL from top open-source MLLMs as the base for OmniGeo. Multimodal instruction fine-tuning data from multiple tasks will be proportionally mixed for joint training, using both LoRA and full supervision fine-tuning approaches. The base models accept default system instructions and task-specific queries, with the training goal being to generate task-specific answers in an autoregressive manner.

Specifically, for tokenized training texts, $[x_1, \dots, x_{n_x}, y_1, \dots, y_{n_y}]$ where x and y denote queries and ground-truth answer, respectively. The training loss \mathcal{L}_Q is defined as:

$$\mathcal{L}_Q = -\log p(y|x) = \sum_{i=1}^{n_y} -\log p(y_i | y_{i-1}, \dots, y_1, x_{n_x}, \dots, x_1) \quad (7)$$

In terms of LoRA (Low-Rank Adaptation) fine-tuning, given the pre-trained weight matrix $W_0 \in \mathbb{R}^{d_{in} \times d_{out}}$, the fine-tuning introduces low-rank matrices $A \in \mathbb{R}^{d_{in} \times r}$ and $B \in \mathbb{R}^{r \times d_{out}}$, where r is the rank of the adaptation. The objective function is defined as

$$\mathcal{L}_{LoRA} = \mathbb{E}_{(x,y)} [\mathcal{L}(f(W_0 + AB, x), y)] \quad (8)$$

where $A \in \mathbb{R}^{d_{in} \times r}$ and $B \in \mathbb{R}^{r \times d_{out}}$ are the low-rank matrices learned during fine-tuning. W_0 is the pre-trained model weight matrix and $f(W_0 + AB, x)$ is the model output after LoRA fine-tuning. LoRA fine-tuning adjusts the low-rank matrices A and B to

adapt the model to various geospatial tasks, without requiring the retraining of the entire weight matrix W_0 .

5 Experiments

5.1 Datasets

We performed experiments on 19 datasets across 5 subtasks, including 14 multimodal datasets (1-4) and 5 text-only datasets (5), with a total of 128,060 samples, including (1) **Health Geography:** Self-constructed Dataset. US state-level and US country-level; (2) **Urban Geography:** Self-constructed Dataset. UG-Beijing and UG-Shenzhen; (3) **Remote Sensing:** geographic benchmark dataset. AID [57] and UC-Merced [59]; (4) **Urban Perception:** self-constructed dataset. For noise indicators: Noise-SG and Noise-SZ. For other indicators: GS-Beautiful, GS-Boring, GS-Depressing, GS-Lively, GS-Safe and GS-Wealthy; and (5) **Geospatial Semantics:** geographic benchmark dataset. LGL [36], GeoVirus [17], and NEEL [23] are toponym recognition datasets, whereas HaveyTweet2017 [24] and GeoCorpora [28] are location description recognition datasets.

The geospatial semantic analysis datasets is randomly divided into an 80% training set and a 20% test set. For the dementia death prediction datasets, it is randomly split into 80% training set and 20% test set at both the state and county levels. The urban region functional classification datasets is split using stratified random sampling into 80% for training and 20% for testing. For the RS image scene classification datasets, the AID is divided by stratified random sampling into 80% training and 20% testing, while UC-Merced uses the official test set. For the urban perception datasets, the 7 urban perception datasets are split using stratified random sampling into 80% training set and 20% test set.

5.2 Evaluation Metrics

For **Health Geography:** time series forecasting task. The evaluation metrics include: Mean Absolute Error (MAE), Mean Squared Error (MSE), Mean Absolute Percentage Error (MAPE), and R2-score. For **Urban Geography:** multimodal multi-class classification task. Considering the issue of class imbalance, the main evaluation metrics are: Precision (P), Recall (R), and Weighted-F1. For **Remote Sensing:** multimodal multi-class classification task. The evaluation metrics include: P, R, Weighted-F1. For **Urban Perception:** multimodal multi-class classification task. Considering the issue of class imbalance, the main evaluation metrics are: P, R, Weighted-F1. For **Geospatial Semantics:** Sequence labeling task. The evaluation metrics are: P, R, F1-score.

5.3 Baseline Methods

We use four groups of models as baselines: task-specific ML models, deep neural network models in general domains, top open-source LLMs, MLLMs, and a closed-source model GPT-4o. They are detailed as: (1) **ARIMA** [3]: An advanced time series forecasting ML model, for the temporal relationship of dementia mortality modeling; (2) **Place2Vec** and **HGI**: learn region-specific semantic embeddings to effectively perform various downstream tasks; (3) **AlexNet** [31], **ResNet18** [18], **ResNet50** [18], **DenseNet161** [25]: Four deep neural network models that have demonstrated strong classification

Table 2: Results of OmniGeo and the baselines on *Health Geography*, the dementia death counts time series forecasting task.

Model	US state-level				US country-level			
	MSE↓	MAE↓	MAPE↓	R^2 ↑	MSE↓	MAE↓	MAPE↓	R^2 ↑
ARIMA	562768.0000	462.0000	6.70%	0.9840	708935.4100	374.1000	4.71%	0.9800
GPT-4o	746173.8182	420.3636	0.07%	0.9792	1361.8132	17.0774	0.18%	0.9841
LLaVA1.5-7B	935866.3636	451.8182	0.06%	0.9739	1284.8832	16.0810	0.16%	0.9849
Qwen2-VL-7B	1234893.4444	688.5556	23.45%	0.9704	263919.3817	148.5340	383.21%	-
OmniGeo (LLaVA)	25329.8182	117.6364	0.03%	0.9993	462.3494	11.6622	0.15%	0.9947
OmniGeo (Qwen2)	192489.0000	290.8182	0.06%	0.9946	482.4748	12.3655	0.16%	0.9948

Table 3: Results of OmniGeo and the baselines on *Urban Geography* task with urban region function classification.

Model	UG-Shenzhen			
	Acc↑	P↑	R↑	F1↑
Place2Vec	0.5674	0.4737	0.5674	0.5164
HGI	0.6620	0.6893	0.6620	0.6754
BLIP-2	0.3567	0.1992	0.3567	0.2556
GPT-4o	0.8277	0.8798	0.8277	0.8530
LLaVA1.5-7B	0.3914	0.2755	0.3914	0.3234
Qwen2-VL-7B	0.3670	0.2314	0.3670	0.2839
OmniGeo (LLaVA)	0.8268	0.8872	0.8268	0.8559
OmniGeo (Qwen2)	0.8277	0.8823	0.8277	0.8541

and generalization abilities on non-RS image datasets; (4) **Stanford NER**, **spaCy NER**, and **NeuroTPR**: General NER tools. NeuroTPR is a recurrent neural network-based model that can extract location information from text; (5) **BERT** [12], **OpenCLIP-B** [7], **BLIP2** [35], **LLaVA1.5** [37], **Qwen2-VL** [55], **MiniCPM** [62], **InternVL2** [5, 6]: Vanilla BERT and six top open-source vision-language models; and (6) **GPT-4o**: A typical closed-source visual model that outperforms most open-source models in the general domain.

5.4 Implementation Details

The base models of OmniGeo are LLaVA1.5-7B and Qwen2-VL-7B, with model parameters and pre-trained weights set according to the official configurations. OmniGeo (LLaVA) is fine-tuned with either full fine-tuning or LoRA fine-tuning on 8 A6000 GPUs with 48GB memory, training for 3 to 5 epochs, with a learning rate of $2e-6$ and a warmup ratio of 0.03. OmniGeo (Qwen2) is fine-tuned with either full fine-tuning or LoRA fine-tuning on 4 L20 GPUs with 48GB memory, training for 3 to 5 epochs, with a learning rate of $1e-5$ and a warmup ratio of 0.1.

6 Results

6.1 Main Results

Health Geography task performance is shown in Table 2 where LLaVA1.5 is on par with GPT-4o and both of them surpass the task-specific model ARIMA and the vision LLM Qwen2-VL. However, our

OmniGeo, leveraging rich geospatial knowledge and multimodal advantages, outperforms all these baselines in both regression metrics and model interpretability scores.

Urban Geography task results are reported in Table 3. After multimodal instruction fine-tuning, OmniGeo achieves absolute improvements of 0.3473 and 0.2227 in Weighted-F1 compared to Place2Vec and HGI, respectively. Besides, it also performs on par with the best closed-source model GPT-4o, which suggests that OmniGeo can utilize RS images and POI data to model regional-level semantic information and human activity patterns. Notably, most open-source MLLMs show very low classification performance, which implies that knowledge from the general domain cannot be directly transferred to this task due to the significant domain gap. Nevertheless, our designed strategies address such issues and thus enhance the classification accuracy of OmniGeo.

Remote Sensing and Urban Perception tasks comparison are presented in Table 4. Benefiting from the geospatial knowledge gained during fine-tuning, OmniGeo can understand high-level semantic information of each scene type and achieve the best classification performance. The typical MLLMs perform even worse than the four DL models, due to the lack of explicit geographic spatial knowledge and its application capabilities. The traditional DL models exhibit good classification attributing to the specific fine-tuning. These observations indicate the importance of grasping or arousing the domain knowledge parametrically.

Geospatial Semantics task results are shown in Table 5, where our OmniGeo still outperforms the other systems based on various types of backbone models. It even outperforms the strong baseline GPT-4o, which indicates that with proper optimization, open-source models can be viable alternatives for geographic knowledge enhancement, offering better performance than the best commercial models in the GeoAI task.

6.2 Ablation Studies

To investigate the effectiveness of leveraging multimodal information and the synergy between geospatial tasks for model training, we conduct ablation studies from two perspectives: i) integrating multimodal information via different fine-tuned configurations (e.g., LoRA or full parameters) and ii) the affect of jointly fine-tuning various tasks simultaneously. The results are shown in Table 6 and Table 7, respectively.

Impact of multimodal information. As shown in Table 6, fusing multimodal information (image in our cases) significantly improves

Table 4: Results of OmniGeo and the baseline models on the seven *Urban Perception* prediction tasks and *Remote Sensing* (RS) image scene classification task with Weighted-F1 scores.

Model	Noise	Beautiful	Boring	Depressing	Lively	Safe	Wealthy	RS
AlexNet	0.3642	0.4343	0.4998	0.4547	0.5144	0.4550	0.4581	0.7410
ResNet18	0.3802	0.4146	0.4464	0.4218	0.4819	0.4333	0.4278	0.7564
ResNet50	0.3359	0.4332	0.4565	0.4440	0.4977	0.4253	0.4428	0.7252
DenseNet161	0.2940	0.4501	0.4593	0.4587	0.4975	0.4833	0.4727	0.7882
OpenCLIP-B-9B	0.3063	0.2576	0.3426	0.1976	0.3428	0.2906	0.2554	0.5766
BLIP2	0.2506	0.2865	0.1797	0.1151	0.2818	0.1271	0.1538	0.3722
GPT-4o	0.3133	0.2686	0.3262	0.3485	0.1934	0.3313	0.2862	0.7125
MiniCPM	0.1927	0.1336	0.2874	0.2418	0.2049	0.1102	0.2287	0.0141
InternVL2-8B	0.2811	0.3051	0.1025	0.2832	0.2671	0.2413	0.3135	0.6792
LLaVA1.5-7B	0.2412	0.4728	0.0747	0.1153	0.2135	0.2342	0.2412	0.2394
Qwen2-VL-7B	0.1209	0.2365	0.2884	0.1604	0.2917	0.1853	0.2993	0.5634
OmniGeo (LLaVA)	0.4407	0.5438	0.5413	0.5159	0.6174	0.5362	0.5047	0.9303
OmniGeo (Qwen2)	0.4397	0.5645	0.5578	0.5740	0.6188	0.6062	0.4830	0.9104

Table 5: Results of OmniGeo and the baseline models on *Geospatial Semantics* tasks in terms of the toponym recognition and location description recognition.

Model	Toponym Recognition						Location Description Recognition		
	LGL			GeoVirus			HaveyTweet2017		
	P↑	R↑	F1-score↑	P↑	R↑	F1-score↑	P↑	R↑	F1-score↑
Stanford NER	0.5872	0.5329	0.5588	0.9043	0.8351	0.8683	0.5140	0.3816	0.4380
spaCy NER	0.2390	0.5371	0.3308	0.5248	0.8028	0.6349	0.4261	0.4145	0.4203
NeuroTPR	0.3701	0.6578	0.4737	0.7059	0.8929	0.7884	0.5428	0.5686	0.5554
BERT	0.5248	0.6647	0.5865	0.5670	0.7520	0.6483	0.6376	0.7549	0.6913
BLIP-2	0.6667	0.0042	0.0083	0.5000	0.0039	0.0078	0.4500	0.0501	0.0902
GPT-4o	0.5285	0.8204	0.6429	0.8841	0.8207	0.8512	0.4434	0.6657	0.5323
LLaVA1.5-7B	0.4160	0.5710	0.4814	0.8081	0.5494	0.6541	0.3705	0.5181	0.4361
Qwen2-VL-7B	0.2673	0.3512	0.3036	0.6164	0.3889	0.4769	0.4237	0.0696	0.1196
OmniGeo (LLaVA)	0.7140	0.6722	0.6925	0.8917	0.8425	0.8664	0.7366	0.8022	0.7680
OmniGeo (Qwen2)	0.7806	0.8204	0.8000	0.9240	0.9167	0.9203	0.7884	0.7577	0.7727

Table 6: Results of the modal ablation experiments on the dementia death counts time series forecasting and urban region function classification tasks.

Model	Config.	w/ image	US state-level				US country-level				Urban Geography(UG-Shenzhen)			
			MSE↓	MAE↓	MAPE↓	R ² ↑	MSE↓	MAE↓	MAPE↓	R ² ↑	P↑	R↑	Weighted-F1↑	
Qwen2-VL-7B	LoRA	✓	198724.8182	250.4545	0.05%	0.9945	509.0063	12.3903	0.17%	0.9946	0.8823	0.8277	0.8541	
	LoRA	✗	275918.5455	281.4545	0.05%	0.9923	482.4748	12.3655	0.16%	0.9948	0.3509	0.4316	0.3871	
	Full	✓	192489.0000	290.8182	0.06%	0.9946	3391.0498	15.9542	0.17%	0.9620	0.8821	0.8212	0.8506	
	Full	✗	1043293.1429	757.1429	0.22%	0.9793	5355.6794	28.9660	0.21%	0.9431	0.6049	0.4682	0.5278	
LLaVA1.5-7B	LoRA	✓	25329.8182	117.6364	0.03%	0.9993	706.5149	13.2911	0.16%	0.9920	0.8827	0.8305	0.8558	
	LoRA	✗	135622.0000	196.5455	0.04%	0.9962	668.5937	13.3393	0.16%	0.9925	0.6786	0.4026	0.5054	
	Full	✓	143825.5455	225.3636	0.05%	0.9960	462.3494	11.6622	0.15%	0.9947	0.8872	0.8268	0.8559	
	Full	✗	175533.9091	275.1818	0.06%	0.9951	498.8000	12.3462	0.17%	0.9943	0.3868	0.4345	0.4093	

the performance on top of two backbone models compared to using single modal information (text) only. The improvements are across various datasets and fine-tuning settings (both LoRA and full parameters). The slight performance decrease is only observed on the US country-level dataset, which might be because of the low resolution of RS images and noise from atypical country data.

Impact of the jointly fine-tuning. From Table 7, we observe that our OmniGeo (LLaVA) fine-tuned with multiple tasks jointly shows more robust results compared to the baselines where each task was fine-tuned individually based on LLaVA-1.5-7B. This is mainly attributed to the cross-modal and cross-task knowledge sharing and transfer via multiple task resources jointly fine-tuning.

Table 7: Results of task ablation on the base model LLaVA1.5-7B.

Model	Health Geography								Urban Geography		Urban Perception(Noise)		RS	
	US state-level				US country-level				UG-Shenzhen	UG-Beijing(zs)	UP-Singapore	UP-Shenzhen(zs)	AID	UC-Merced(zs)
	MSE↓	MAE↓	MAPE↓	R ² ↑	MSE↓	MAE↓	MAPE↓	R ² ↑	F1↑	F1↑	F1↑	F1↑	F1↑	F1↑
w/. HG-s & HG-c	424105.4545	302.3636	0.04%	0.9881	519.6601	12.1513	0.15%	0.9941	-	-	-	-	-	-
w/. UG-SZ	-	-	-	-	-	-	-	-	0.8637	0.3460	-	-	-	-
w/. UP-SG	-	-	-	-	-	-	-	-	-	-	0.2632	0.2410	-	-
w/. RS-AID	-	-	-	-	-	-	-	-	-	-	-	-	0.8012	0.5832
OmniGeo (Ours)	25329.8182	117.6364	0.03%	0.9993	462.3494	11.6622	0.15%	0.9947	0.8559	0.3557	0.4407	0.4683	0.9303	0.6025

A performance slightly decreased is observed on the Shenzhen dataset, which might imply the inherent noise in multi-task datasets and minor information loss due to domain transfer.

7 Conclusion and Future Work

This study contributes a large-scale, high-quality visual instruction-following dataset to the GeoAI community. Besides, we explore the potential of MLLMs in GeoAI, facilitating cross-modal geographic knowledge sharing and transfer, and develop the first MLLM, OmniGeo, covering all core tasks in GeoAI. It achieves performance surpassing or on par with GPT-4o, effectively addressing the inherent multimodal nature of data in GeoAI and taking the first step towards the deployment of MLLMs in GeoAI. In the future, we can further develop more robust and effective systems for GeoAI, which provide flexible interaction [45, 46] with the domain users [54] and support more geographic tasks [41].

References

- [1] Igor Akushevich, Arseniy P Yashkin, Anatoliy I Yashin, and Julia Kravchenko. 2021. Geographic disparities in mortality from Alzheimer’s disease and related dementias. *Journal of the American Geriatrics Society* 69, 8 (2021), 2306–2315.
- [2] Reza Askarizad and Jinliao He. 2022. Perception of spatial legibility and its association with human mobility patterns: An empirical assessment of the historical districts in rasht, iran. *International Journal of Environmental Research and Public Health* 19, 22 (2022), 15258.
- [3] George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. 2015. *Time series analysis: forecasting and control*. John Wiley & Sons.
- [4] Zhiyu Chen, Shiyang Li, Charese Smiley, Zhiqiang Ma, Sameena Shah, and William Yang Wang. 2022. ConvFinQA: Exploring the Chain of Numerical Reasoning in Conversational Finance Question Answering. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 6279–6292. doi:10.18653/v1/2022.emnlp-main.421
- [5] Zhe Chen, Weiyun Wang, Hao Tian, Shenglong Ye, Zhangwei Gao, Erfei Cui, Wenwen Tong, Kongzhi Hu, Jiapeng Luo, Zheng Ma, et al. 2024. How Far Are We to GPT-4V? Closing the Gap to Commercial Multimodal Models with Open-Source Suites. *arXiv preprint arXiv:2404.16821* (2024).
- [6] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. 2024. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 24185–24198.
- [7] Mehdi Cherti, Romain Beaumont, Ross Wightman, Mitchell Wortsman, Gabriel Ilharco, Cade Gordon, Christoph Schuhmann, Ludwig Schmidt, and Jenia Jitsev. 2023. Reproducible scaling laws for contrastive language-image learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2818–2829.
- [8] Seongjin Choi, Jiwon Kim, and Hwasoo Yeo. 2021. TrajGAIL: Generating urban vehicle trajectories using generative adversarial imitation learning. *Transportation Research Part C: Emerging Technologies* 128 (2021), 103091. doi:10.1016/j.trc.2021.103091
- [9] Philippa J. Clarke, Jennifer Weuve, Lisa Barnes, Denis A. Evans, and Carlos F. Mendes de Leon. 2015. Cognitive decline and the neighborhood environment. *Annals of Epidemiology* 25, 11 (2015), 849–854. doi:10.1016/j.annepidem.2015.07.001
- [10] Wenliang Dai, Junnan Li, DONGXU LI, Anthony Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale N Fung, and Steven Hoi. 2023. InstructBLIP: Towards General-purpose Vision-Language Models with Instruction Tuning. In *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (Eds.), Vol. 36. Curran Associates, Inc., 49250–49267. https://proceedings.neurips.cc/paper_files/paper/2023/file/9a6a435e75419a836fe47ab6793623e6-Paper-Conference.pdf
- [11] DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucang Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaoqin Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuan Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shutong Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanbiao Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yudian Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yuxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. 2024. DeepSeek-V3 Technical Report. arXiv:2412.19437 [cs.CL] <https://arxiv.org/abs/2412.19437>
- [12] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. doi:10.18653/v1/N19-1423
- [13] Fahim Faisal and Antonios Anastasopoulos. 2022. Geographic and geopolitical biases of language models. *arXiv preprint arXiv:2212.10408* (2022).
- [14] Jenny Rose Finkel, Trond Grenager, and Christopher D Manning. 2005. Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the 43rd annual meeting of the association for computational linguistics (ACL ’05)*. 363–370.
- [15] Jason Fletcher, Katie Jajtner, and Jinho Kim. 2024. Geographic disparities in Alzheimer’s disease and related dementia mortality in the US: Comparing impacts of place of birth and place of residence. *SSM - Population Health* 27 (2024), 101708. doi:10.1016/j.ssmph.2024.101708
- [16] Wensheng Gan, Zhenlian Qi, Jiayang Wu, and Jerry Chun-Wei Lin. 2023. Large Language Models in Education: Vision and Opportunities. arXiv:2311.13160 [cs.AI] <https://arxiv.org/abs/2311.13160>
- [17] Milan Gritta, Mohammad Taher Pilehvar, and Nigel Collier. 2018. Which Melbourne? Augmenting Geocoding with Maps. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Iryna Gurevych and Yusuke Miyao (Eds.). Association for Computational

- Linguistics, Melbourne, Australia, 1285–1296. doi:10.18653/v1/P18-1119
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [19] Nanjun He, Leyuan Fang, Shutao Li, Antonio Plaza, and Javier Plaza. 2018. Remote sensing scene classification using multilayer stacked covariance pooling. *IEEE Transactions on Geoscience and Remote Sensing* 56, 12 (2018), 6899–6910.
- [20] Noah Hollmann, Samuel Müller, Lennart Purucker, Arjun Krishnakumar, Max Körfer, Shi Bin Hoo, Robin Tibor Schirrmeyer, and Frank Hutter. 2025. Accurate predictions on small data with a tabular foundation model. *Nature* 637, 8045 (jan 2025), 319–326. doi:10.1038/s41586-024-08328-6
- [21] Matthew Honnibal. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. (No Title) (2017).
- [22] Yujun Hou, Matias Quintana, Maxim Khomiakov, Winston Yap, Jiani Ouyang, Koichi Ito, Zeyu Wang, Tianhong Zhao, and Filip Biljecki. 2024. Global Streetscapes – A comprehensive dataset of 10 million street-level images across 688 cities for urban science and analytics. *ISPRS Journal of Photogrammetry and Remote Sensing* 215 (2024), 216–238. doi:10.1016/j.isprsjprs.2024.06.023
- [23] Xuke Hu, Yeran Sun, Jens Kersten, Zhiyong Zhou, Friederike Klan, and Hongchao Fan. 2023. How can voting mechanisms improve the robustness and generalizability of toponym disambiguation? *International Journal of Applied Earth Observation and Geoinformation* 117 (2023), 103191. doi:10.1016/j.jag.2023.103191
- [24] Yingjie Hu and Jimin Wang. 2020. How do people describe locations during a natural disaster: an analysis of tweets from Hurricane Harvey. doi:10.48550/arXiv.2009.12914
- [25] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4700–4708.
- [26] Kaiyu Huang, Fengran Mo, Xinyu Zhang, Hongliang Li, You Li, Yuanchi Zhang, Weijian Yi, Yulong Mao, Jinchun Liu, Yuzhuang Xu, et al. 2024. A survey on large language models with multilingualism: Recent advances and new frontiers. *arXiv preprint arXiv:2405.10936* (2024).
- [27] Weiming Huang, Daokun Zhang, Gengchen Mai, Xu Guo, and Lizhen Cui. 2023. Learning urban region representations with POIs and hierarchical graph infomax. *ISPRS Journal of Photogrammetry and Remote Sensing* 196 (2023), 134–145.
- [28] Alan M. MacEachren Jan Oliver Wallgrün, Morteza Karimzadeh and Scott Pezanowski. 2018. GeoCorpora: building a corpus to test and train microblog geoparsers. *International Journal of Geographical Information Science* 32, 1 (2018), 1–29. doi:10.1080/13658816.2017.1368523 arXiv:https://doi.org/10.1080/13658816.2017.1368523
- [29] Hang Jiang, Xijie Zhang, Robert Mahari, Daniel Kessler, Eric Ma, Tal August, Irene Li, Alex Pentland, Yoon Kim, Deb Roy, and Jad Kabbara. 2024. Leveraging Large Language Models for Learning Complex Legal Concepts through Storytelling. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 7194–7219. doi:10.18653/v1/2024.acl-long.388
- [30] Yu Jing. 2008. Remote sensing semantic model for city planning. *Computer Applications*, 51 (2008), 348–435.
- [31] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* 25 (2012).
- [32] Kartik Kuckreja, Muhammad Sohail Danish, Muzammal Naseer, Abhijit Das, Salman Khan, and Fahad Shahbaz Khan. 2024. Geochat: Grounded large vision-language model for remote sensing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 27831–27840.
- [33] Suhyeon Lee, Won Jun Kim, Jinho Chang, and Jong Chul Ye. 2024. LLM-CXR: Instruction-Finetuned LLM for CXR Image Understanding and Generation. arXiv:2305.11490 [cs.CV] https://arxiv.org/abs/2305.11490
- [34] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. 2023. LLaVA-Med: Training a Large Language and-Vision Assistant for Biomedicine in One Day. In *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (Eds.), Vol. 36. Curran Associates, Inc., 28541–28564. https://proceedings.neurips.cc/paper_files/paper/2023/file/5abdcf8edccacba028c6662789194572-Paper-Datasets_and_Benchmarks.pdf
- [35] Junning Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *Proceedings of the 40th International Conference on Machine Learning (Honolulu, Hawaii, USA) (ICML '23)*. JMLR.org, Article 814, 13 pages.
- [36] Michael D. Lieberman, Hanan Samet, and Jagan Sankaranarayanan. 2010. Geo-tagging with local lexicons to build indexes for textually-specified spatial data. In *2010 IEEE 26th International Conference on Data Engineering (ICDE 2010)*. 201–212. doi:10.1109/ICDE.2010.5447903
- [37] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning. In *Advances in Neural Information Processing Systems*, A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (Eds.), Vol. 36. Curran Associates, Inc., 34892–34916. https://proceedings.neurips.cc/paper_files/paper/2023/file/6dcf277ea32ce3288914fa369fede0-Paper-Conference.pdf
- [38] Kang Liu, Song Gao, Peiyuan Qiu, Xiliang Liu, Bo Yan, and Feng Lu. 2017. Road2Vec: Measuring Traffic Interactions in Urban Road System from Massive Travel Routes. *International Journal of Geo-Information* 6 (10 2017), 321. doi:10.3390/ijgi6110321
- [39] Robert Mahari, Dominik Stambach, Elliott Ash, and Alex Pentland. 2023. The Law and NLP: Bridging Disciplinary Disconnects. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, Singapore, 3445–3454. doi:10.18653/v1/2023.findings-emnlp.224
- [40] Gengchen Mai, Weiming Huang, Jin Sun, Suhang Song, Deepak Mishra, Ninghao Liu, Song Gao, Tianming Liu, Gao Cong, Yingjie Hu, Chris Cundy, Ziyuan Li, Rui Zhu, and Ni Lao. 2024. On the Opportunities and Challenges of Foundation Models for GeoAI (Vision Paper). *ACM Trans. Spatial Algorithms Syst.* 10, 2, Article 11 (July 2024), 46 pages. doi:10.1145/3653070
- [41] Gengchen Mai, Yiqun Xie, Xiaowei Jia, Ni Lao, Jimmeng Rao, Qing Zhu, Zeping Liu, Yao-Yi Chiang, and Junfeng Jiao. 2025. Towards the next generation of Geospatial Artificial Intelligence. *International Journal of Applied Earth Observation and Geoinformation* 136 (2025), 104368.
- [42] Gengchen Mai, Bo Yan, Krzysztof Janowicz, and Rui Zhu. 2020. Relaxing Unanswerable Geographic Questions Using A Spatially Explicit Knowledge Graph Embedding Model. In *Geospatial Technologies for Local and Regional Development*, Phaedon Kyriakidis, Diofantos Hadjimitsis, Dimitrios Skarlatos, and Ali Mansourian (Eds.). Springer International Publishing, Cham, 21–39.
- [43] Rohin Manvi, Samar Khanna, Marshall Burke, David Lobell, and Stefano Ermon. 2024. Large language models are geographically biased. In *Proceedings of the 41st International Conference on Machine Learning (Vienna, Austria) (ICML '24)*. JMLR.org, Article 1409, 16 pages.
- [44] Fengran Mo, Kelong Mao, Ziliang Zhao, Hongjin Qian, Haonan Chen, Yiruo Cheng, Xiaoxi Li, Yutao Zhu, Zhicheng Dou, and Jian-Yun Nie. 2024. A survey of conversational search. *arXiv preprint arXiv:2410.15576* (2024).
- [45] Fengran Mo, Kelong Mao, Yutao Zhu, Yihong Wu, Kaiyu Huang, and Jian-Yun Nie. 2023. ConvGQR: Generative Query Reformulation for Conversational Search. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 4998–5012.
- [46] Fengran Mo, Jian-Yun Nie, Kaiyu Huang, Kelong Mao, Yutao Zhu, Peng Li, and Yang Liu. 2023. Learning to relate to previous turns in conversational search. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1722–1732.
- [47] OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aidan Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpouras, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gilda Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hesham Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Makin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukas Kaiser, Luke Metz, Madeline Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kauffer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glase, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renney

- Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Petersen, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vittor Pong, Vlad Fomenko, Weiyei Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang, Yunyun Wang, Zheng Shao, and Zhuohan Li. 2024. OpenAI o1 System Card. arXiv:2412.16720 [cs.AI] <https://arxiv.org/abs/2412.16720>
- [48] Chao Pang, Xingxing Weng, Jiang Wu, Jiayu Li, Yi Liu, Jiaxing Sun, Weijia Li, Shuai Wang, Litong Feng, Gui-Song Xia, and Conghui He. 2024. VHM: Versatile and Honest Vision Language Model for Remote Sensing Image Analysis. arXiv:2403.20213 [cs.CV] <https://arxiv.org/abs/2403.20213>
- [49] Andrew Pannone, Aditya Raj, Hari Krishnan Ravichandran, Sarbajit Das, Ziheng Chen, Collin A. Price, Mahmooda Sultana, and Saptarshi Das. 2024. Robust chemical analysis with graphene chemosensors and machine learning. *Nature* 634, 8034 (oct 2024), 572–578. doi:10.1038/s41586-024-08003-w
- [50] Guande Qi, Xiaolong Li, Shijian Li, Gang Pan, Zonghui Wang, and Daqing Zhang. 2011. Measuring social functions of city regions from large-scale taxi behaviors. In *2011 IEEE International Conference on Pervasive Computing and Communications Workshops (PERCOM Workshops)*. IEEE, 384–388.
- [51] Jimeng Rao, Song Gao, and Sijia Zhu. 2023. CATS: Conditional Adversarial Trajectory Synthesis for Privacy-Preserving Trajectory Data Publication Using Deep Learning Approaches. arXiv:2309.11587 [cs.LG] <https://arxiv.org/abs/2309.11587>
- [52] Di Wang, Qiming Zhang, Yufei Xu, Jing Zhang, Bo Du, Dacheng Tao, and Liangpei Zhang. 2022. Advancing plain vision transformer toward remote sensing foundation model. *IEEE Transactions on Geoscience and Remote Sensing* 61 (2022), 1–15.
- [53] Jimin Wang, Yingjie Hu, and Kenneth Joseph. 2020. NeuroTPR: A neuro-net toponym recognition model for extracting locations from social media messages. *Transactions in GIS* 24, 3 (2020), 719–735.
- [54] Jiayin Wang, Fengran Mo, Weizhi Ma, Peijie Sun, Min Zhang, and Jian-Yun Nie. 2024. A User-Centric Multi-Intent Benchmark for Evaluating Large Language Models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. 3588–3612.
- [55] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Ke-Yang Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. *ArXiv abs/2409.12191* (2024). <https://api.semanticscholar.org/CorpusID:272704132>
- [56] Yu-Tzu Wu, A. Matthew Prina, and Carol Brayne. 2015. The association between community environment and cognitive function: a systematic review. *Social Psychiatry and Psychiatric Epidemiology* 50, 3 (mar 2015), 351–362. doi:10.1007/s00127-014-0945-6
- [57] Gui-Song Xia, Jingwen Hu, Fan Hu, Baoguang Shi, Xiang Bai, Yanfei Zhong, Liangpei Zhang, and Xiaoqiang Lu. 2017. AID: A Benchmark Data Set for Performance Evaluation of Aerial Scene Classification. *IEEE Transactions on Geoscience and Remote Sensing* 55, 7 (2017), 3965–3981. doi:10.1109/TGRS.2017.2685945
- [58] Haochong Xia, Shuo Sun, Xinrun Wang, and Bo An. 2024. Market-GAN: Adding Control to Financial Market Data Generation with Semantic Context. *Proceedings of the AAAI Conference on Artificial Intelligence* 38, 14 (Mar. 2024), 15996–16004. doi:10.1609/aaai.v38i14.29531
- [59] Yi Yang and Shawn Newsam. 2010. Bag-Of-Visual-Words and Spatial Extensions for Land-Use Classification. In *ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems (ACM GIS)*.
- [60] Xiwen Yao, Junwei Han, Gong Cheng, Xueming Qian, and Lei Guo. 2016. Semantic annotation of high-resolution satellite images via weakly supervised learning. *IEEE Transactions on Geoscience and Remote Sensing* 54, 6 (2016), 3660–3671.
- [61] Yao Yao, Xia Li, Xiaoping Liu, Penghua Liu, Zhaotang Liang, Jinbao Zhang, and Ke Mai. 2017. Sensing spatial distribution of urban land use by integrating points-of-interest and Google Word2Vec model. *International Journal of Geographical Information Science* 31, 4 (2017), 825–848.
- [62] Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng Zou, Haoye Zhang, Shengding Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. MiniCPM-V: A GPT-4V Level MLLM on Your Phone. arXiv:2408.01800 [cs.CV] <https://arxiv.org/abs/2408.01800>
- [63] Seong-Keun Yoo, Conall W. Fitzgerald, Byuri Angela Cho, Bailey G. Fitzgerald, Catherine Han, Elizabeth S. Koh, Abhinav Pandey, Hannah Sfreddo, Fionnuala Crowley, Michelle Rudsteyn Korostin, Neha Debnath, Yan Leyfman, Cristina Valero, Mark Lee, Joris L. Vos, Andrew Sangho Lee, Karena Zhao, Stanley Lam, Ezekiel Oluomuyide, Fengshen Kuo, Eric A. Wilson, Pauline Hamon, Clotilde Hennequin, Miriam Saffern, Lynda Vuong, A. Ari Hakimi, Brian Brown, Miriam Merad, Sacha Gnjatich, Nina Bhardwaj, Matthew D. Galsky, Eric E. Schadt, Robert M. Samstein, Thomas U. Marron, Mithat Gönen, Luc G. T. Morris, and Diego Chowell. 2025. Prediction of checkpoint inhibitor immunotherapy efficacy for cancer using routine blood tests and clinical data. *Nature Medicine* (2025). doi:10.1038/s41591-024-03398-5 Published: 2025/01/06.
- [64] Yihong Yuan and Martin Raubal. 2016. Analyzing the distribution of human activity space from mobile phone usage: an individual and urban-oriented study. *International Journal of Geographical Information Science* 30, 8 (2016), 1594–1621.
- [65] Fan Zhang, Zhuangyuan Fan, Yuhao Kang, Yujie Hu, and Carlo Ratti. 2021. "Perception bias": Deciphering a mismatch between urban crime and perception of safety. *Landscape and Urban Planning* 207 (2021), 104003.
- [66] Fan Zhang, Bolei Zhou, Liu Liu, Yu Liu, Helene H Fung, Hui Lin, and Carlo Ratti. 2018. Measuring human perceptions of a large-scale urban region using machine learning. *Landscape and Urban Planning* 180 (2018), 148–160.
- [67] Ziwei Zhang, Liang Wu, Liufeng Tao, Sheng Hu, Hui Long, Yongyang Xu, Jinquan Li, Jingjing Zhang, Zhijun Zhou, Jing Liu, Cheng Cai, Hong Zhang, Dan Liu, Yan Zeng, and Wei Luo. 2024. Geospatial Applications in Alzheimer's Disease Research and Beyond: A Systematic Review. *Annals of the American Association of Geographers* (2024). <https://api.semanticscholar.org/CorpusID:271724904>
- [68] Tianhong Zhao, Xiucheng Liang, Wei Tu, Zhengdong Huang, and Filip Biljecki. 2023. Sensing urban soundscapes from street view imagery. *Computers, Environment and Urban Systems* 99 (2023), 101915. doi:10.1016/j.compenvurbysys.2022.101915
- [69] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. MiniGPT-4: Enhancing Vision-Language Understanding with Advanced Large Language Models. *arXiv preprint arXiv:2304.10592* (2023).

A More Details of Experimental Setup

A.1 The Source of the Data and Model

You can easily obtain the datasets or models involved in this paper from the following resource list:

- (1) CDC Wonder: <https://wonder.cdc.gov/ucd-icd10.html>. The dementia mortality data are obtained from the US Centers for Disease Control and Prevention Wide-ranging Online Data for Epidemiologic Research (CDC WONDER).
- (2) Gaode Maps: <https://lbs.amap.com/>. Obtain POI data of Beijing and Shenzhen via API.
- (3) WorldView: <https://worldview.earthdata.nasa.gov/>. Through this website, you can obtain WorldView series remote sensing image data products of Beijing, Shenzhen and the United States.
- (4) Google Maps : <https://www.google.com/maps>. The source of the street view images used in the article. Alternatively, you can directly use the Visual-Soundscapes and Global Streetscapes datasets.
- (5) city government websites: <https://opendata.sz.gov.cn/> and https://www.beijing.gov.cn/gongkai/guihua/wngh/cqgh/201907/t20190701_100008.html.
- (6) LLaVA1.5-7B and Qwen2-VL-7B: The base model weights from the Huggingface website: <https://huggingface.co/>.

B Additional Experimental Results

B.1 Results in Zero-Shot Settings

As shown in Table 8, we demonstrate the comparison of the generalization performance of OmniGeo and baseline models on five datasets in zero-shot settings.

Geospatial Semantics task performance is shown in Table 8 where OmniGeo achieved competitive performance with GPT-4o on both the NEEL and GeoCorpora datasets.

Urban Geography task results are reported in Table 8, GPT-4o performs the best, achieving optimal results on all three metrics. OmniGeo (LLaVA) achieved results that are competitive with GPT-4o. It is worth noting that LLaVA1.5 performs slightly better than OmniGeo (LLaVA) on the Beijing dataset, which is understandable for two reasons. First, OmniGeo (LLaVA) acquired substantial knowledge about the urban spatial structure of Shenzhen during fine-tuning, and transferring this knowledge directly to Beijing is challenging due to the substantial differences in urban planning and structure between Shenzhen and Beijing. Second, the Beijing dataset contains a large amount of noise. Compared to Shenzhen, which is predominantly industrial, Beijing has many regions with small sizes and unclear land-use types, or mixed land-use types.

Urban Perception task comparison is presented in Table 8. OmniGeo (LLaVA) achieved the best performance on the Shenzhen dataset, with a particularly higher Weighted-F1 score of 0.1548 compared to GPT-4o. This suggests that OmniGeo (LLaVA) has the capability to transfer advanced perceptual knowledge across cities for urban perception tasks, which is precisely what we aimed for.

Remote Sensing task results are shown in Table 8. GPT-4o maintains a relatively high classification ability on the UC-Merced dataset, while OmniGeo (LLaVA) consistently demonstrates competitive RS image classification ability with GPT-4o. This indicates

that OmniGeo (LLaVA) correctly understands the semantics of each scene and transfers geographical spatial knowledge across RS images with varying spatial and spectral resolutions.

B.2 Discussion

For the dementia death counts time series fore-casting, Figure 4 visualizes the prediction capabilities of various baselines and OmniGeo (LLaVA) on the US state-level dataset based on PE values. Significant geographic distribution differences in prediction abilities exist across different states, possibly due to "geographical bias". Based on the position 0 in the legend, it can be seen that ARIMA, GPT-4o, and LLaVA1.5 tend to overestimate the death count across the entire dataset. Regarding the degree of overestimation, LLaVA1.5 tends to overestimate more significantly, while ARIMA and GPT-4o show similar behaviors, with overestimation tendencies in Florida and Maine, respectively. In contrast, the baseline models severely underestimate the death toll in New Mexico, while OmniGeo (LLaVA) shows balanced predictive performance. The results at the US country-level are similar, but Qwen2-VL performs poorly, which may be attributed to its poor instruction-following ability in this task and its lack of relevant geographical knowledge.

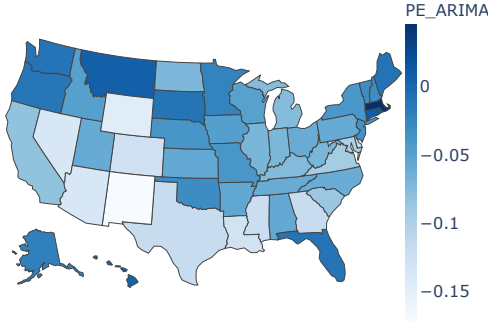
For the urban region function classification, The embedding models Place2Vec and HGI exhibit different confusion patterns: Place2Vec tends to predict "Commercial office space" as "Commercial service land," while HGI tends to predict "Commercial service land" as "Sports and cultural land". Blip2 failed in this task, as it predicted nearly all regional functional types as "Residential land," while the GPT-4o model performed well. In contrast, OmniGeo (LLaVA) demonstrated the most stable performance. The above results suggest that "Commercial service land," "Residential land," "Residential land," and "Commercial office space" are all challenging samples to classify, and different models exhibit varying confusion patterns across different land-use types, which may be related to the models' geographic knowledge capacity and the difficulty of transferring ontology knowledge into the GeoAI domain.

For the urban perception prediction, In the prediction of the Noise indicator, AlexNet tends to classify noise intensity as "Quiet" and "Very Quiet," while GPT-4o and InternVL2 exhibit similar behavior, with the only difference being that GPT-4o is more inclined to predict "Quiet," while InternVL2 tends to predict "Very Quiet". Interestingly, OmniGeo (LLaVA) is more likely to categorize it as "Noisy." In the prediction of the Wealthy metric, the first three models still exhibit different preference patterns, while OmniGeo (LLaVA) shows relatively balanced performance.

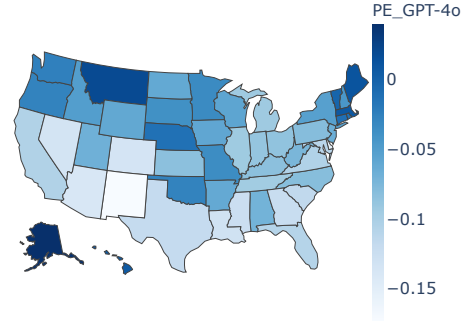
For the geospatial semantic, we provide two cases where OmniGeo (Qwen2) gave correct answers, but GPT-4o failed. In toponym recognition, given the text: "Romans, Sakai and the parolee died later in a shootout when the city's SWAT team stormed an apartment where the man was hiding.", GPT-4o mistakenly identifies "apartment" as a named place, which is evidently just a generic entity name. In location description recognition, given the tweet "HurricaneHarvey My mother needs help! Please send a boat to 4254 Geronimo Lake Dr. Houston, TX 77047 rescue Pleasehel", the correct answer is "4254 Geronimo Lake Dr. Houston, TX 77047", but GPT-4o erroneously splits it as "4254 Geronimo Lake Dr; Houston; TX; 77047", which is obviously an incomplete location description.

Table 8: Results of OmniGeo and the baseline models on geospatial related tasks under the zero-shot setting.

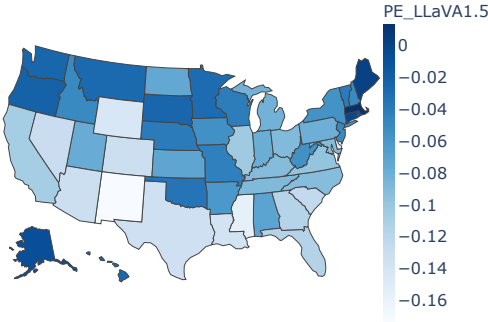
Model	Semantic Analysis						Urban Geography			Urban Perception(Noise)			RS		
	NEEL			GeoCorpora			UG-Beijing			UP-Shenzhen			UC-Merced		
	P↑	R↑	F1-score↑	P↑	R↑	F1-score↑	P↑	R↑	F1↑	P↑	R↑	F1↑	P↑	R↑	F1↑
Stanford NER	0.7331	0.5184	0.6074	0.8118	0.5118	0.6278	-	-	-	-	-	-	-	-	-
spaCy NER	0.5372	0.5488	0.5429	0.5562	0.5183	0.5366	-	-	-	-	-	-	-	-	-
NeuroTPR	0.7326	0.7582	0.7452	0.8199	0.7214	0.7675	-	-	-	-	-	-	-	-	-
BERT	0.5667	0.2213	0.3183	0.5292	0.3193	0.3983	-	-	-	-	-	-	-	-	-
BLIP-2	0.2000	0.0022	0.0043	0.8000	0.0014	0.0027	0.3790	0.1750	0.2395	0.3073	0.2724	0.2888	0.3088	0.2643	0.2848
GPT-4o	0.6980	0.8221	0.7550	0.7875	0.8597	0.8220	0.3884	0.3910	0.3897	0.2858	0.3471	0.3135	0.6726	0.6405	0.6562
LLaVA1.5-7B	0.4818	0.6876	0.5666	0.5711	0.6965	0.6276	0.3551	0.3643	0.3597	0.1704	0.4126	0.2412	0.4455	0.2738	0.3391
Qwen2-VL-7B	0.5954	0.4534	0.5148	0.2999	0.0681	0.1110	0.3025	0.1777	0.2239	0.1415	0.2951	0.1913	0.5047	0.2381	0.3235
OmniGeo (LLaVA)	0.7919	0.8039	0.7585	0.7791	0.7751	0.7771	0.3756	0.3378	0.3557	0.4640	0.4726	0.4683	0.6149	0.5905	0.6025
OmniGeo (Qwen2)	0.7177	0.7115	0.7146	0.7830	0.6636	0.7184	0.3666	0.2807	0.3180	0.5513	0.3858	0.4540	0.3837	0.2500	0.3072



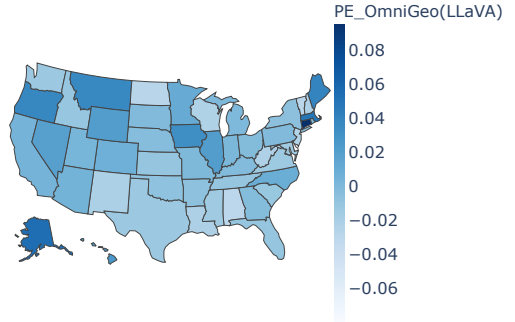
(a) ARIMA



(b) GPT-4o



(c) LLaVA1.5



(d) OmniGeo (LLaVA)

Figure 4: Prediction error plot for each baseline and OmniGeo (LLaVA) on the dementia deaths time series forecasting task. The color of each US region indicates the percentage error of each model in predicting that state $PE = (Prediction - True)/True$.

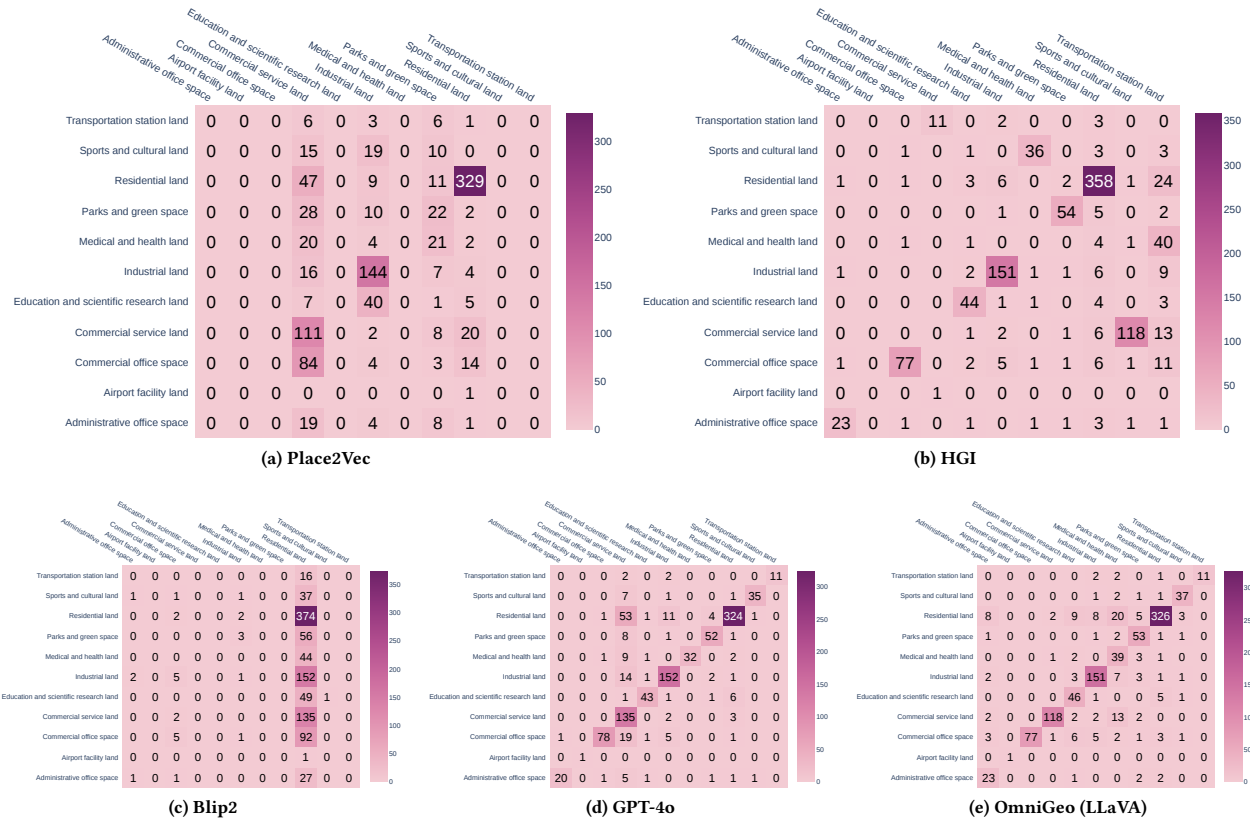


Figure 5: Confusion matrix comparison of Place2Vec, HGI, Blip2, GPT-4o and OmniGeo (LLaVA) models on the Shenzhen urban region function classification dataset

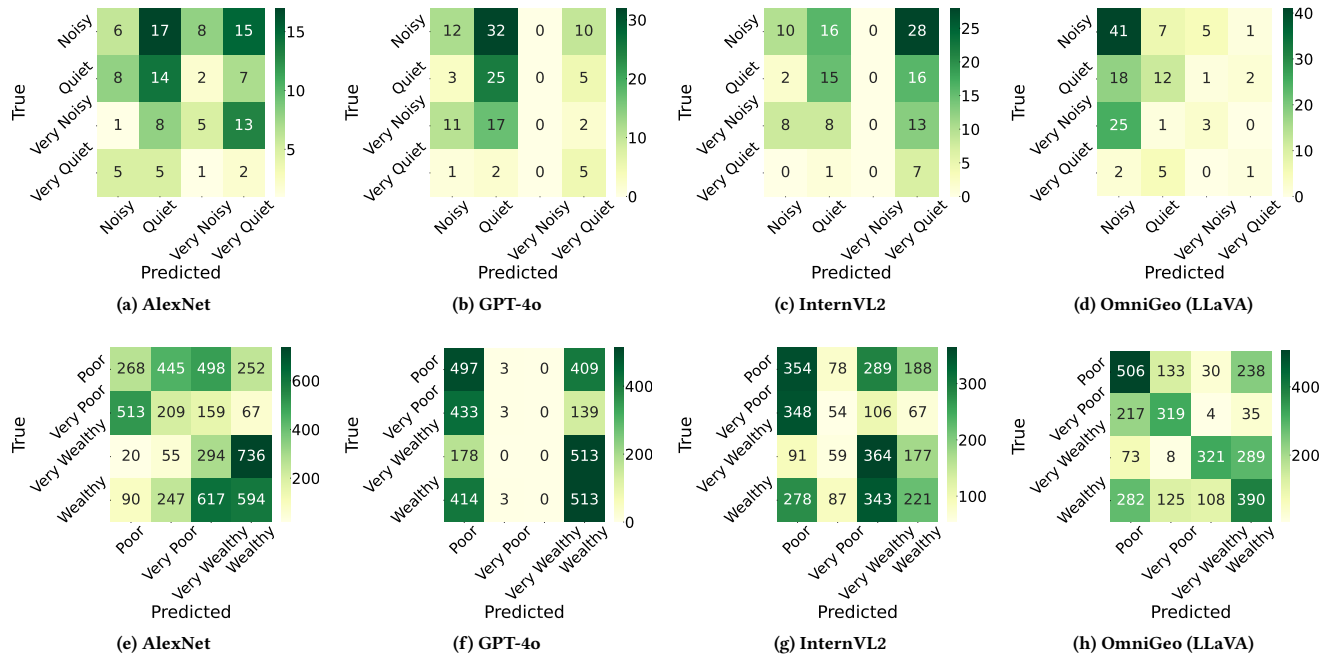


Figure 6: Confusion matrix comparison of AlexNet, InternVL2, GPT-4o and OmniGeo (LLaVA) in predicting Noise and Wealthy indicators of urban perception