# **UrbanLLaVA**: A Multi-modal Large Language Model for Urban Intelligence

**Jie Feng**

**Department of Electronic Engineering**

**Tsinghua University**

**https://vonfeng.github.io/**
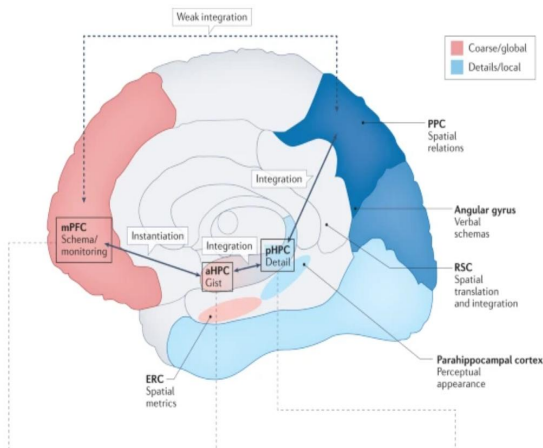
**2025.10**

Arxiv

Github

# Spatial Intelligence

- **Spatial Intelligence—the capacity to understand, reason, and act on spatial information—is the foundation of embodied AI and a key prerequisite for AGI.**



- **Spatial Intelligence is also for Cognitive Science, Urban Science, Complex Systems, Earth Science**

# Spatial Intelligence → Urban Intelligence

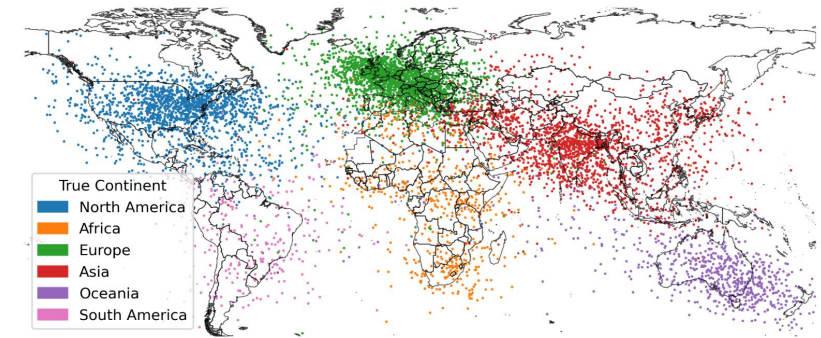- **General LLMs have significant limitations in spatial knowledge, understanding, and reasoning.**



World Level

Urban Level

Embodied Level

Xu, Wenrui*, Dalin Lyu*, Weihang Wang*, **Jie Feng**, Chen Gao, and Yong Li. "Defining and Evaluating VLMs' Basic Spatial Abilities: A Perspective from Psychometrics." ACL 2025 Main.
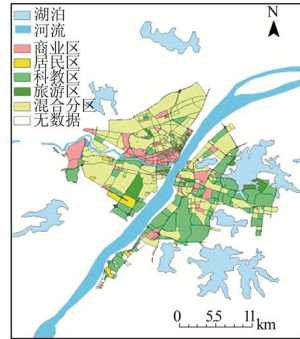
Gurnee, Wes, and Max Tegmark. "Language models represent space and time." ICLR 2024.

# Urban Intelligence

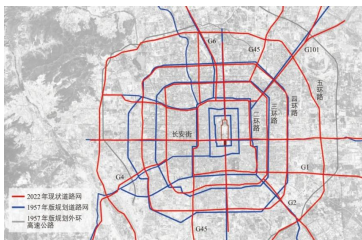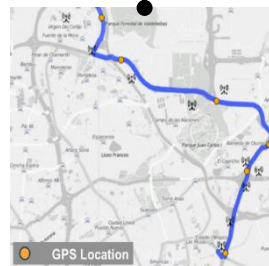- **We need strengthen LLMs spatial intelligence at Urban scale for embodied and generalized AI.**



**Remote Sensing**

**Regions**

**Street View**

**Road Network**

**Mobility Trajectory**

**POI**

# Before UrbanLLaVA: CityBench

- **A comprehensive urban spatial intelligence benchmark**



**8 tasks in 13 cities around the world**

Feng, Jie*, Jun Zhang*, Tianhui Liu*, Xin Zhang, Tianjian Ouyang, Junbo Yan, Yuwei Du, Siqi Guo, and Yong Li. "CityBench: Evaluating the capabilities of large language models for urban tasks." KDD 2025 Dataset & Benchmark.

# Before UrbanLLaVA: CityGPT

- **For urban geospatial knowledge: CityGPT, using daily behavior simulation to synthesize data**
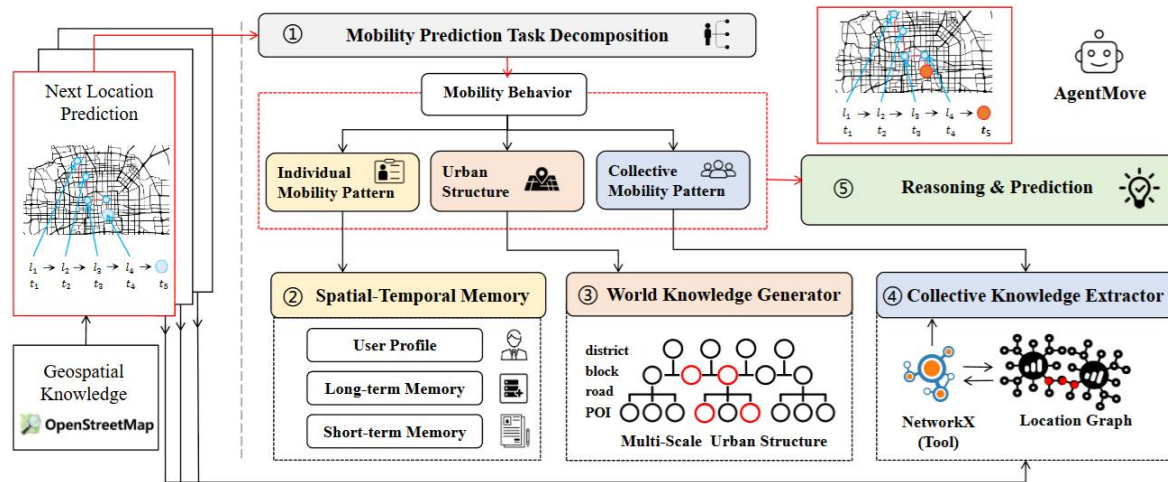


Large-scale human behavior simulation to synthesize high-quality, multi-source spatial data.

Domain-data mixed adaptive training to balance general and specialized capabilities.
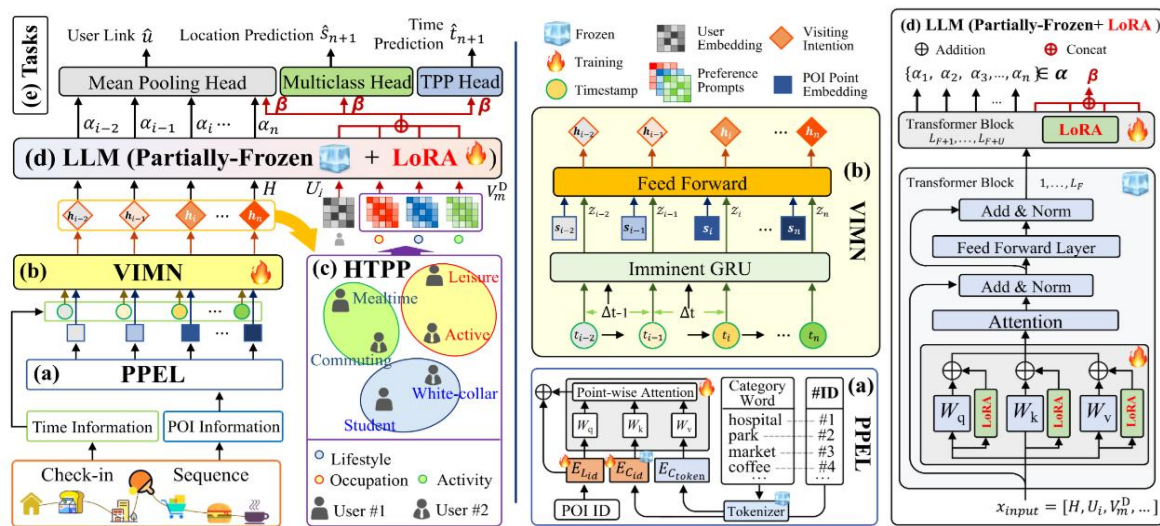
Systematic evaluation of language-based urban spatial cognition.

# Before UrbanLLaVA: AgentMove/Mobility-LLM

- **For urban mobility intelligence: language-driven methods (AgentMove) and modality-aligned modeling (Mobility-LLM).**
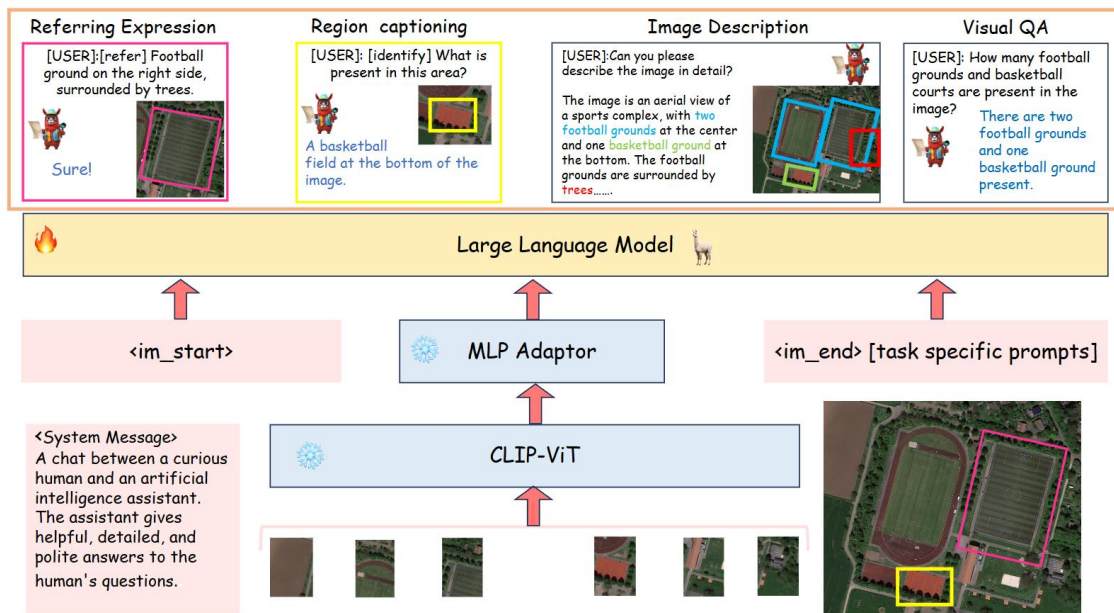


**AgentMove**: Models movement behavior by translating mobility trajectories into natural language descriptions and performing the modeling within the natural language space.
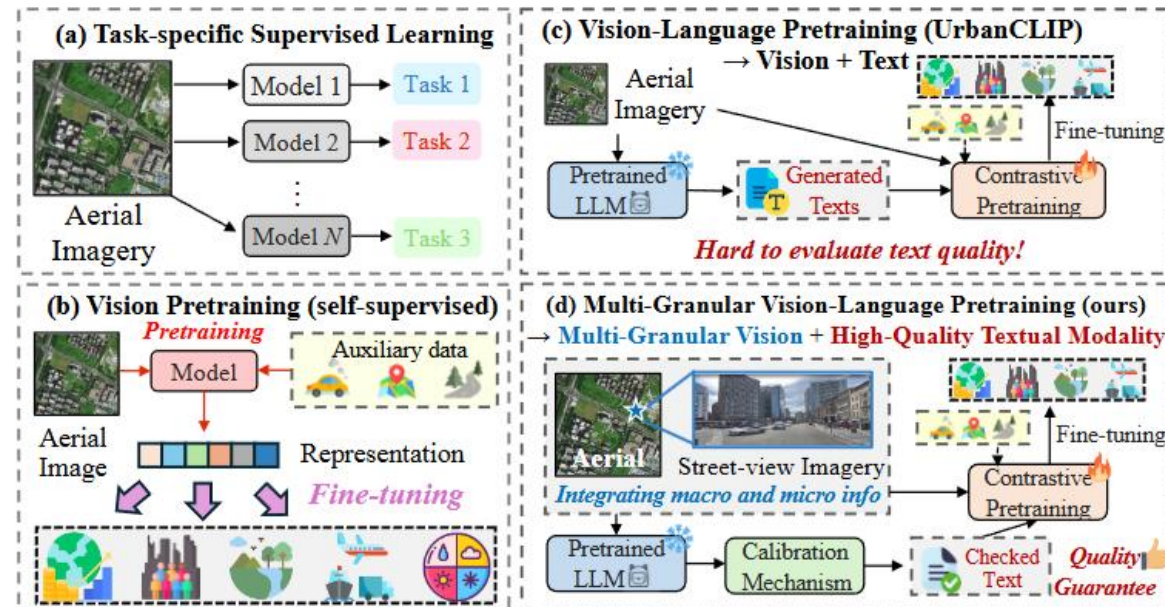
**Mobility-LLM**: Models movement by using classical models to process mobility trajectories, obtaining an embedded representation that is then integrated into the LLM.

Gong, Letian, et al. "Mobility-LLM: Learning visiting intentions and travel preference from human mobility data with large language models." NeurIPS 2024
Feng, Jie, et al. "AgentMove: A large language model based agentic framework for zero-shot next location prediction." NAACL 2025.

# Before UrbanLLaVA: GeoChat/UrbanVLP

- **For urban vision intelligence, e.g., GeoChat and UrbanVLP.**



**GeoChat**: Addresses urban remote sensing recognition and detection tasks using a classical multi-modal architecture.
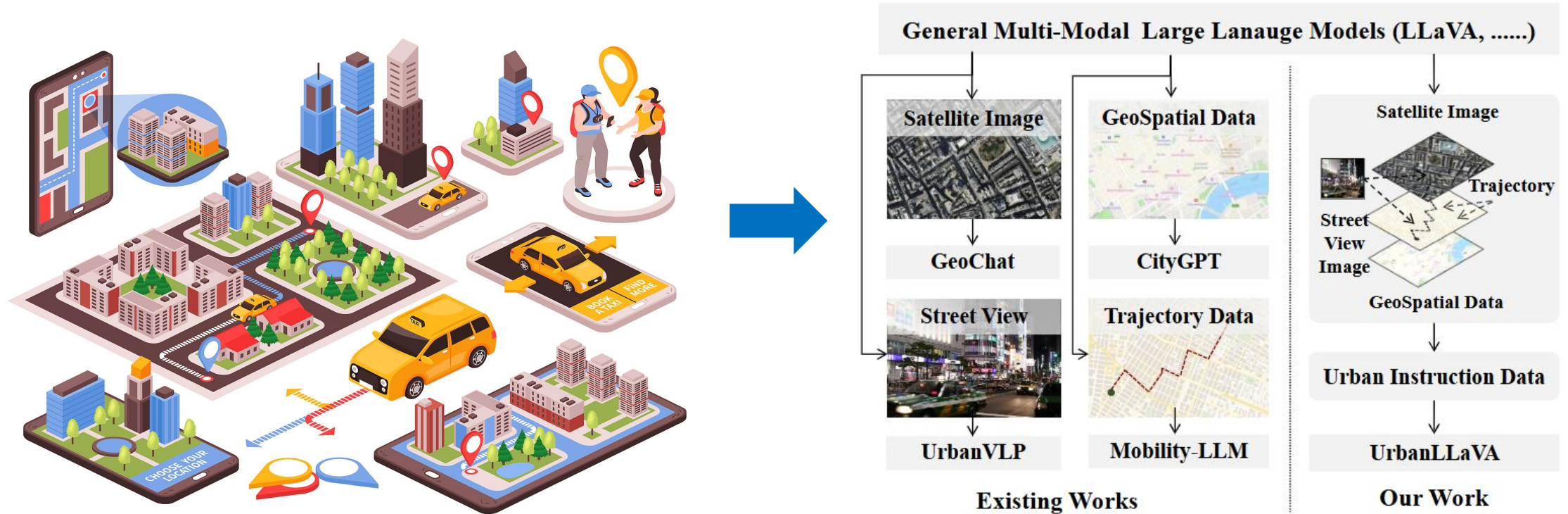
**UrbanVLP**: Focuses on predicting socio-economic indicators from street view imagery, leveraging contrastive learning and LLM annotation.

Kuckreja, Kartik, et al. "Geochat: Grounded large vision-language model for remote sensing. " CVPR 2024
Hao, Xixuan, et al. "Urbanvlp: Multi-granularity vision-language pretraining for urban socioeconomic indicator prediction." AAAI 2025
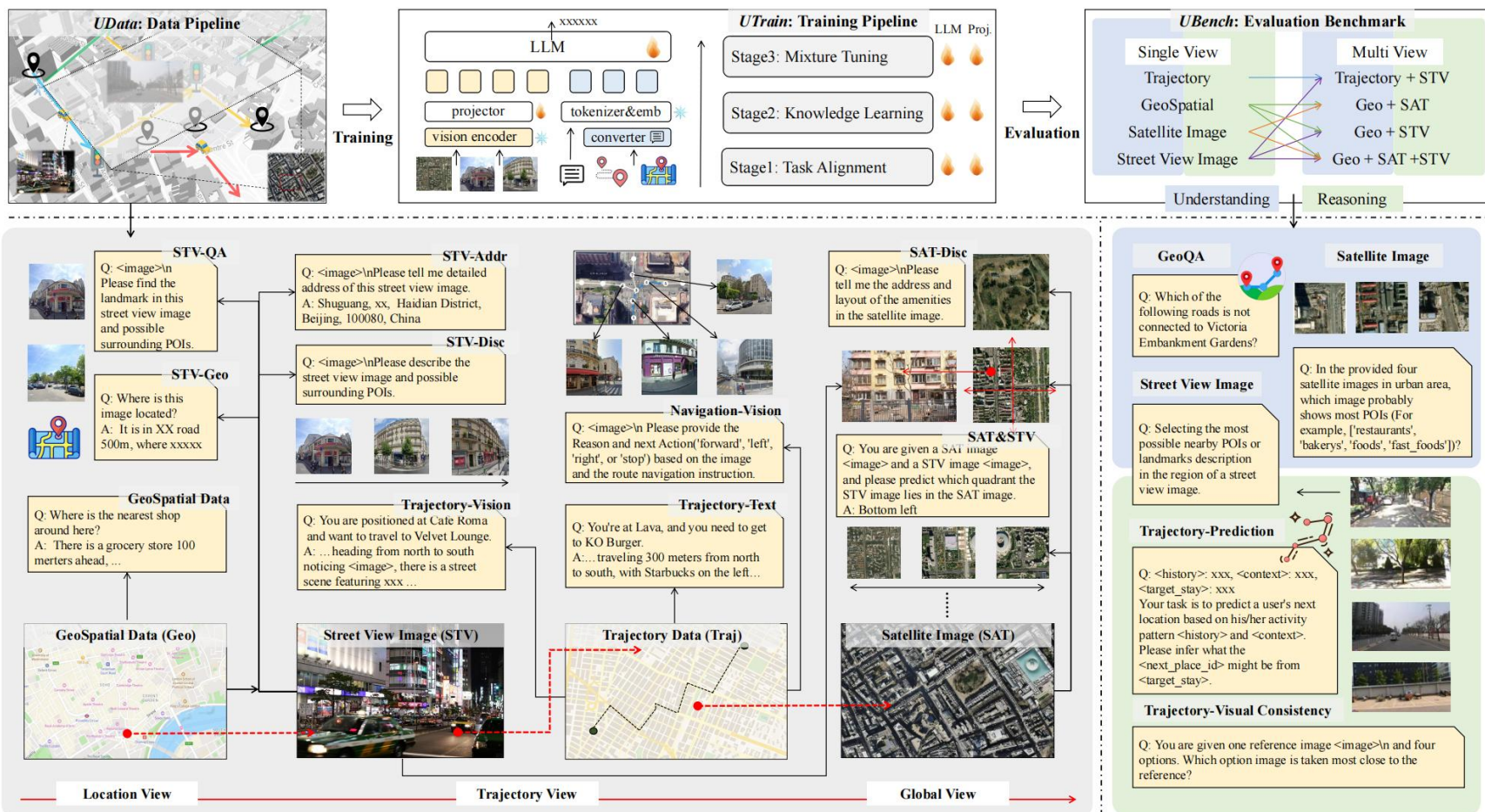
# UrbanLLaVA: Building Unified Foundation Model for Urban Intelligence

- **UrbanLLaVA:** Achieves unified modeling of **remote sensing, street view, geographical data,** and **mobility trajectories** by deeply aligning and fusing **visual information** with **geospatial knowledge**.

**Feng, Jie**, Shengyuan Wang, Tianhui Liu, Yanxin Xi, and Yong Li. "UrbanLLaVA: A Multi-modal Large Language Model for Urban Intelligence with Spatial Reasoning." ICCV 2025

# UrbanLLaVA: Building Unified Foundation Model for Urban Intelligence
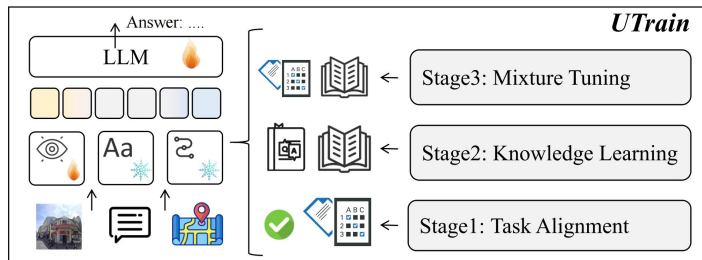
- **UrbanLLaVA: UData + Utrain + UBench**



## Training Data

Table 1. Detailed information about *UBench* for Beijing, 'STV' refers to street view image, and 'SAT' refers to satellite image.
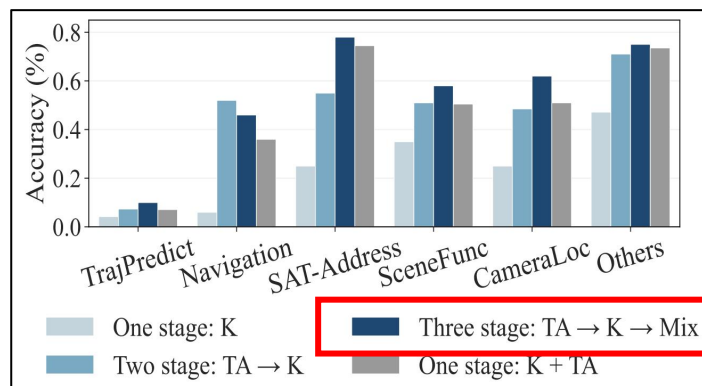
| Tasks | Data | Category | Metrics | Samples | Source |
|---|---|---|---|---|---|
| GeoQA | Geospatial Data | GeoQA | Avg. Accuracy | 1450 | CityBench |
| TrajPredict | Trajectory Data | Geo+Traj | Top-1 | 500 | CityBench |
| Navigation | Single STV | Geo+Traj | Success Rate | 50 | CityBench |
| SceneComp | Multi SAT | Geo+SAT | Accuracy | 200 | UrBench |
| ImgRetrieval | Multi STV & SAT | Geo+SS | Accuracy | 200 | UrBench |
| CameraLoc | Multi STV & SAT | Geo+SS | Accuracy | 200 | UrBench |
| STV-Address | Single STV | Geo+STV | Accuracy | 200 | *UBench* |
| STV-Landmark | Single STV | Geo+STV | Accuracy | 200 | *UBench* |
| SAT-Address | Single SAT | Geo+SAT | Accuracy | 200 | *UBench* |
| SAT-Landuse | Single SAT | Geo+SAT | Accuracy | 200 | *UBench* |
| STV-Outlier | Multi STV | Geo+STV | Accuracy | 200 | *UBench* |
| SceneFunc | Multi SAT | Geo+SAT | Accuracy | 200 | *UBench* |

**Feng, Jie**, Shengyuan Wang, Tianhui Liu, Yanxin Xi, and Yong Li. "UrbanLLaVA: A Multi-modal Large Language Model for Urban Intelligence with Spatial Reasoning." ICCV 2025

- **Uses a multi-stage training pipeline for urban domain multi-modal and multi-scenario tasks, leading to significant performance gains on critical benchmarks.**
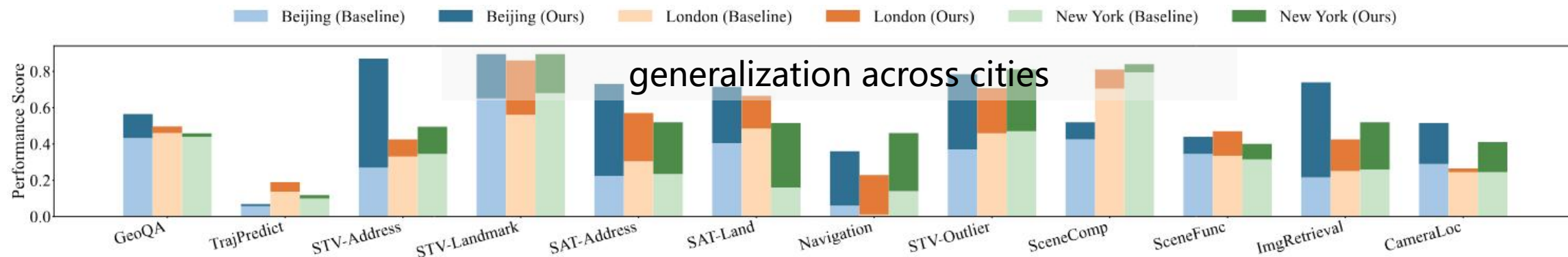


multi-stage training pipeline



multi-stage training approach yields superior performance on critical tasks compared to direct mixed training.

| City | Beijing | | | | | London | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Task Group | GeoQA | Geo+Traj | Geo+STV | Geo+SAT | Geo+SS | GeoQA | Geo+Traj | Geo+STV | Geo+SAT | Geo+SS |
| VILA1.5-3B | 0.3873 | 0.0200 | 0.3967 | 0.3200 | 0.2575 | 0.4362 | 0.0400 | 0.2557 | 0.2850 | 0.2725 |
| VILA1.5-8B | 0.4322 | 0.0589 | 0.4300 | 0.3488 | 0.2425 | 0.4841 | 0.0884 | 0.4495 | 0.4575 | 0.2575 |
| VILA1.5-13B | 0.4410 | 0.1156 | 0.5167 | 0.3638 | 0.2400 | 0.4592 | 0.1298 | 0.4991 | 0.4538 | 0.2625 |
| InternVL2-8B | 0.4709 | 0.1578 | 0.4667 | 0.3313 | 0.2325 | 0.4973 | 0.1347 | 0.4477 | 0.4763 | 0.2400 |
| InternVL2-26B | 0.4877 | 0.1478 | 0.4550 | 0.3825 | 0.2275 | 0.5168 | 0.1288 | 0.4923 | 0.5138 | 0.2425 |
| Qwen2VL-7B | 0.4950 | 0.1389 | 0.4383 | 0.3638 | 0.2675 | 0.4991 | 0.1560 | 0.4381 | 0.4863 | 0.2775 |
| Qwen2VL-72B | 0.5491 | 0.1611 | 0.5817 | 0.3588 | 0.2975 | 0.5802 | 0.2322 | 0.6375 | 0.4375 | 0.3250 |
| LLama3.2-11B | 0.4229 | 0.0756 | 0.4375 | 0.3075 | / | 0.4804 | 0.1180 | 0.4000 | 0.3800 | / |
| LLama3.2-90B | 0.4502 | 0.1056 | 0.5325 | 0.2925 | / | 0.5659 | 0.2010 | 0.5450 | 0.4700 | / |
| GPT4o-mini | 0.4542 | 0.1622 | 0.4350 | 0.3800 | 0.2475 | 0.5357 | 0.1278 | 0.4752 | 0.5388 | 0.2675 |
| GPT4o | 0.5479 | 0.1522 | 0.4300 | 0.4125 | 0.3025 | 0.6446 | 0.1300 | 0.5469 | 0.6050 | 0.2850 |
| *UrbanLLaVA*-VILA1.5-8B | 0.5682 | 0.2800 | 0.8650 | 0.6663 | 0.7025 | 0.6399 | 0.2680 | 0.7500 | 0.7100 | 0.4325 |
| vs. VILA1.5-8B | +31.47% | +375.38% | +101.16% | +91.03% | +189.69% | +32.18% | +203.17% | +66.85% | +55.19% | +67.96% |
| vs. Best Baseline | +3.48% | +72.63% | +48.70% | +61.53% | +132.23% | -0.73% | +15.42% | +17.65% | +17.36% | +33.08% |

achieves significant performance improvements across various tasks in multiple cities.



generalization across cities

• **Task Examples**



**STV-Outlier Task**



**STV-Landmark Task**



**SceneFunc Task with Satellite Image**

**STV-Description**

**CameraLoc Task with STV and Satellite Image**

**UrbanLLaVA@Github**

# Thanks!

**Arxiv**

**Contact Jie Feng via:**

https://vonfeng.github.io/

fengj12ee@hotmail.com