# Towards Urban General Intelligence Through Urban Foundation Models

Hao Liu, Weijia Zhang, Jindong Han,
Zhao Xu, Hang Ni, Hui Xiong

The Hong Kong University of Science and Technology (Guangzhou) &
The Hong Kong University of Science and Technology

KDD2024
BARCELONA, SPAIN

Tutorial
Website

Survey
Paper

Github

# OUTLINE

**Introduction**

**Challenges of Building UFMs**

**Overview of UFMs**

**Prospects of UFMs**

**Summary**

# OUTLINE

**Introduction**

**Challenges of Building UFMs**

**Overview of UFMs**

**Prospects of UFMs**

**Summary**

# City Evolution – Challenge & Opportunity

- City, as a dynamically evolving entity, is driven by a multitude of factors, such as population growth, environmental changes, and technological advancements.

**Challenge**

Traffic jam

Energy crisis

Climate change

...

**Opportunity**

Urban big data

Artificial Intelligence

Smarter & sustainable city

United Nation Released the 17 Sustainable Development Goals (SDGs). https://sdgs.un.org/goals

# Smart City and Urban Intelligence

- Smart city has deeply penetrated into our daily life.


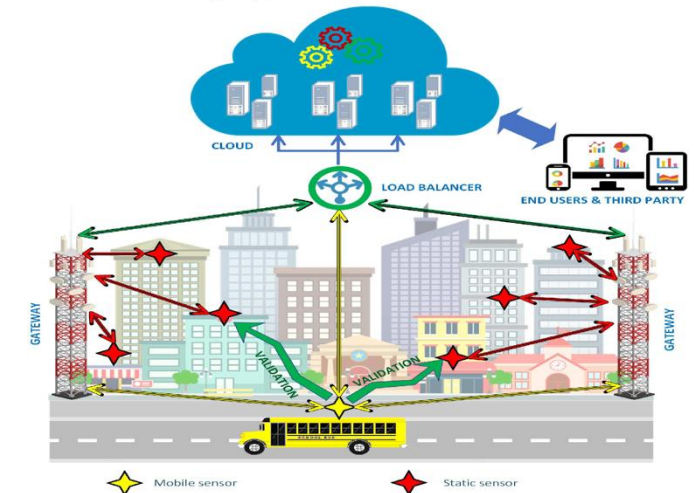**Intelligent transportation systems**


**Smart building**


**E-Government Services**


**Urban surveillance**


**Smart grid**


**Air pollution analysis**

# Urban Intelligence: Are We There Yet?

- Urban Intelligence (UI) refers to the strategic integration of various technologies and practices to enhance the management, operation, and development of urban areas.
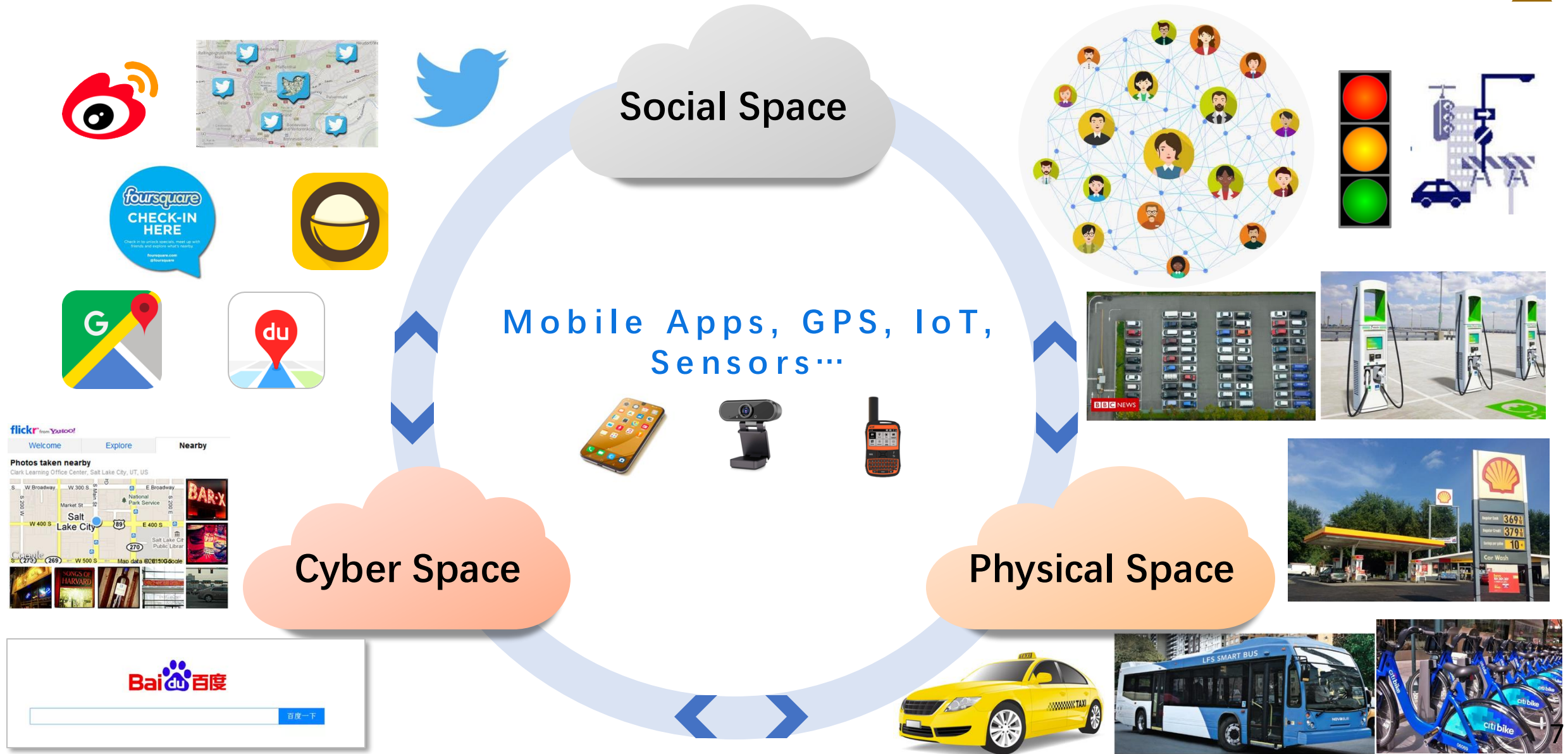


**Drainage System in Forbidden City**



**Beijing flood, 2023**

Social Space

Mobile Apps, GPS, IoT, Sensors…

Cyber Space

Physical Space

# Limitation 2: Task Specific Machine Learning Modeling

- Most data-driven Urban Intelligence solutions follow the machine learning development pipeline

  - Requires abundant labeled data for training -> limited budget/ authority

  - Possess specific skill for designated tasks -> limited generalizability

  - The limit of the model approx. to the quality of the data -> limited knowledge

# Urban General Intelligence (UGI)

- Analogous to Artificial General Intelligence (AGI), UGI is envisioned to autonomously perform any intellectual task related to urban contexts, rivaling or even surpassing human capabilities, thereby transforming cities into more livable, resilient, and adaptive spaces.
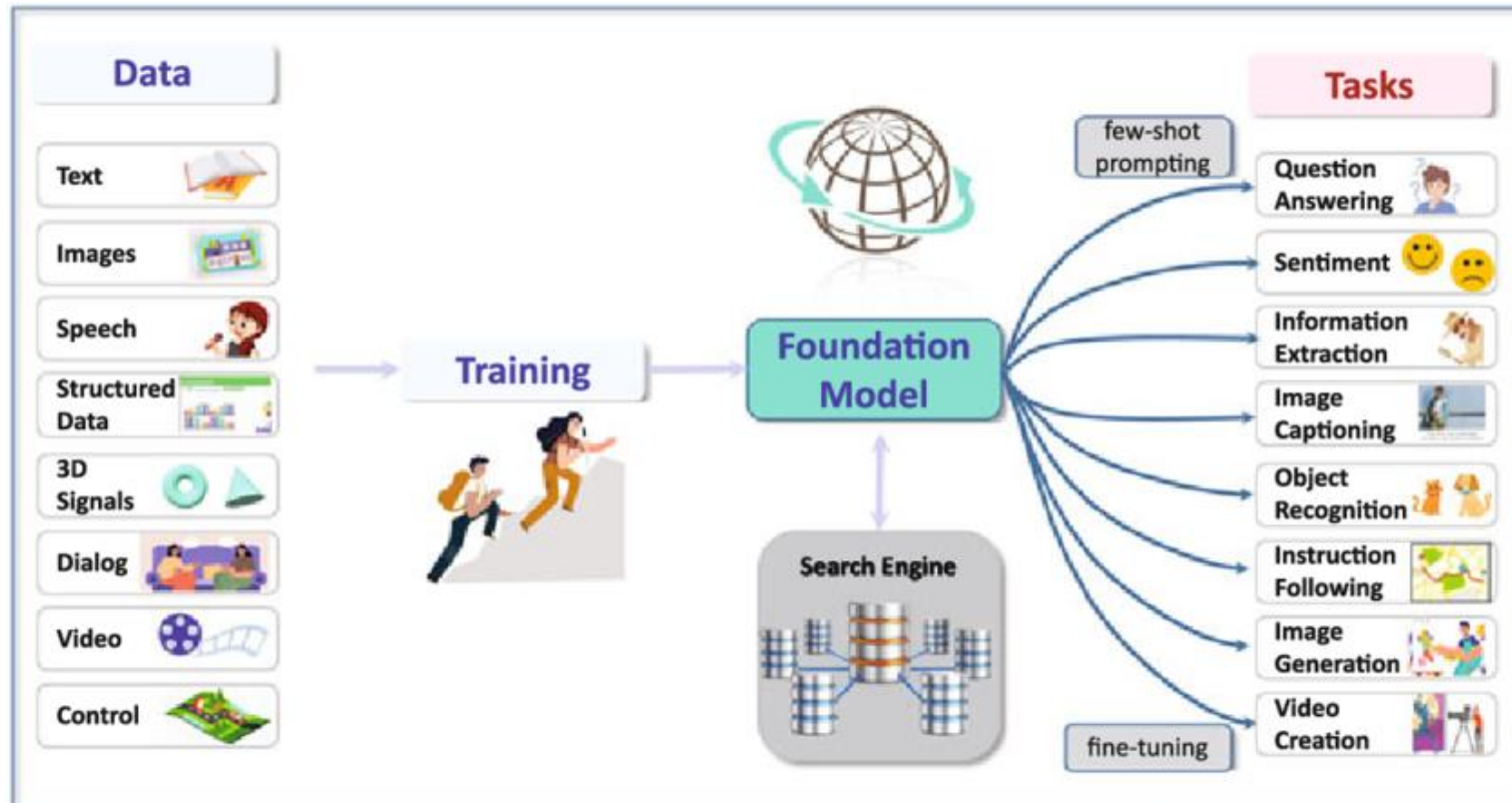


AGI



UGI

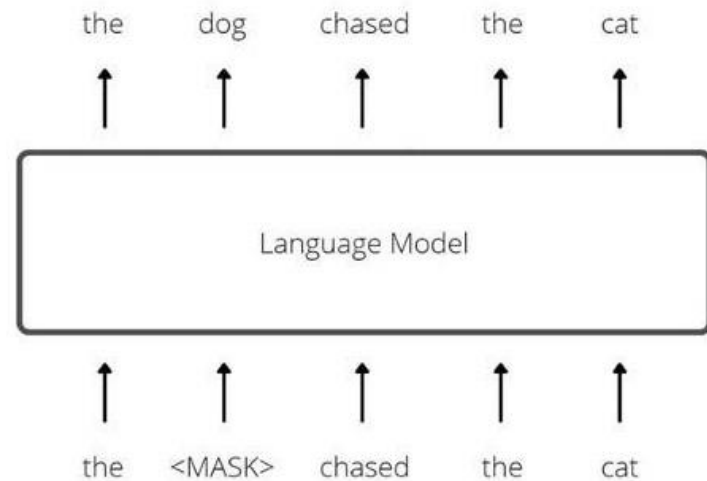# Towards UGI: Are Foundation Models a Potential Solution?

- Recent advances on Large Language Model (e.g., ChatGPT) and Vision Transformer (e.g., SAM) demonstrate remarkable generalization and emergent abilities on a wide spectrum of tasks.
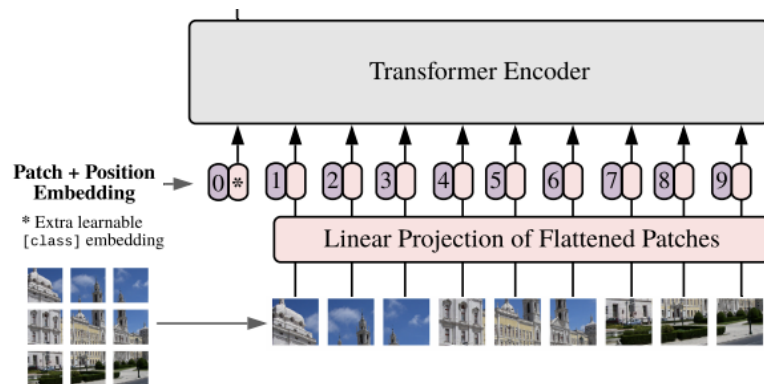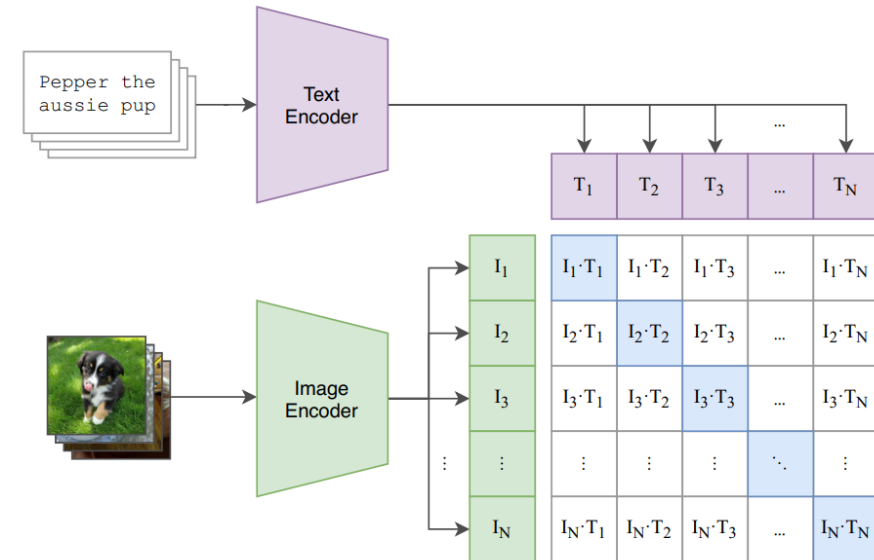
# Basics of Foundation Model

- Pretraining of language, vision, and multi-modal foundation models.



**Masked and next token prediction**

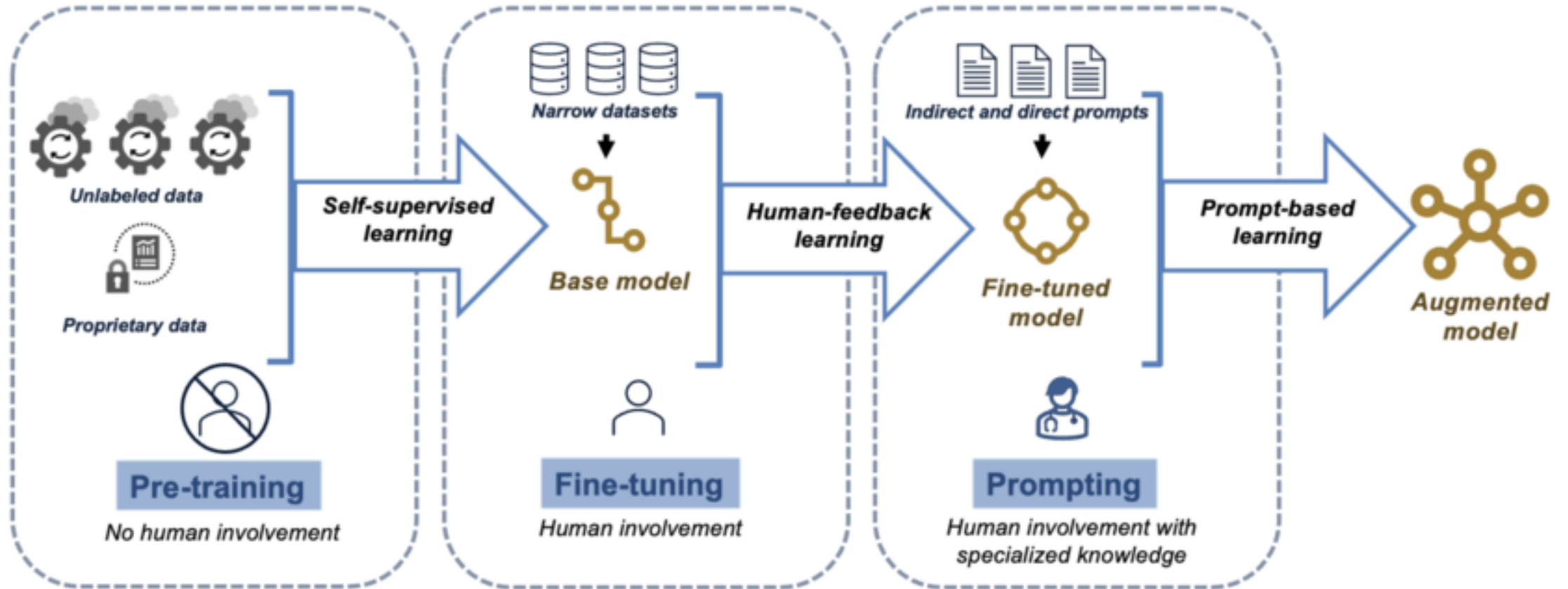

**Masked patch prediction**



**Cross-modal alignment**

# Basics of Foundation Model

- Typical foundation model training and serving stages.

# Can LLMs Memorize Spatial Knowledge?

- LLM has compressed world knowledge in training corpus obtained from the Internet

  - Does LLM memorized geospatial knowledge? If yes, how to extract them?



**User:** Describe the address "6th Avenue, Manhattan, New York" in one detailed paragraph.

**GPT-3.5:** 6th Avenue, Manhattan, New York, is a bustling thoroughfare that cuts through the heart of this iconic borough. It stretches from the southern tip of Manhattan, where it kisses the edge of the verdant Washington Square Park, and winds its way northward, passing through diverse neighborhoods like Greenwich Village, Chelsea, and Midtown. The avenue is a symphony of

...

**GPT-3.5 can describe addresses**



**Basic Prompt:** Coordinates: (40.76208, -73.98042)

Population Density (On a Scale from 0.0 to 9.9):

**Finetuned GPT-3.5:** 5.2

**Our Prompt:** Coordinates: (40.76208, -73.98042)

Address: "Calyon Building, 6th Avenue, Manhattan Community Board 5, Manhattan, New York County, City of New York, New York, United States"
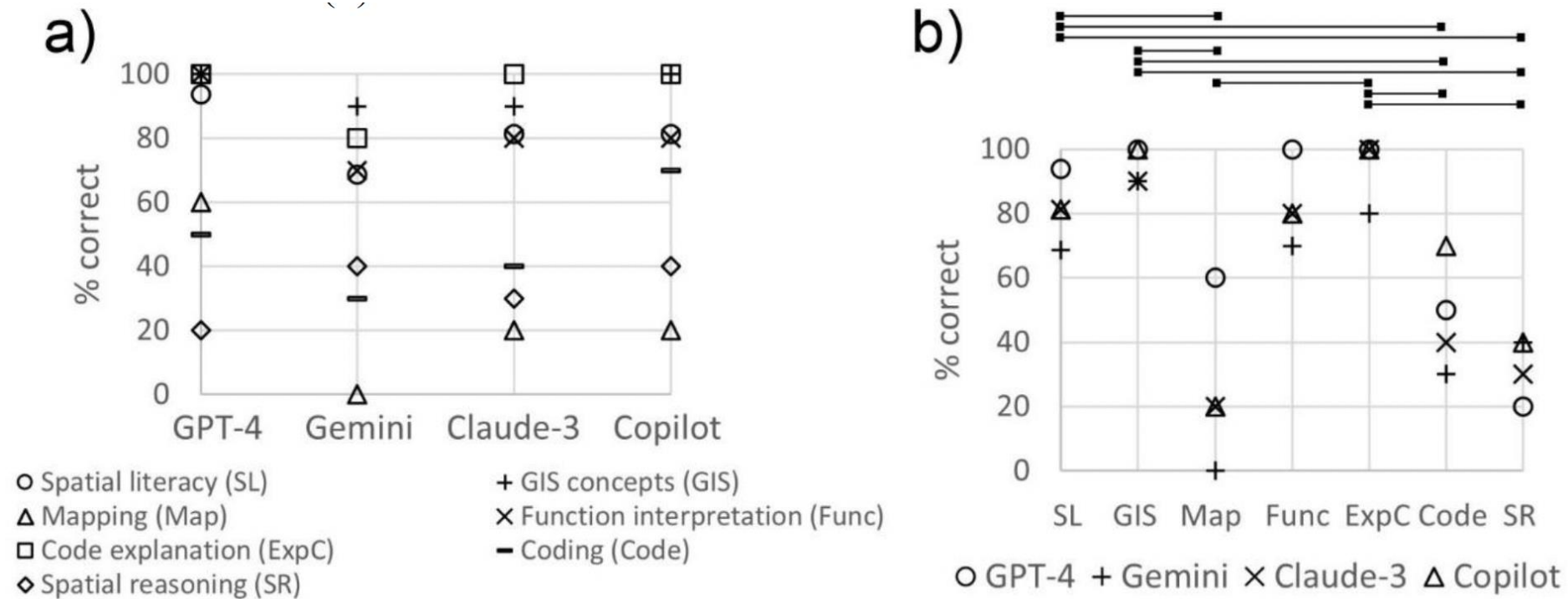
...

**Finetuned GPT-3.5:** 9.0

**Prompting with additional map info**

Rohin Manvi, et al. GeoLLM: Extracting Geospatial Knowledge from Large Language Models. ICLR 2024.

# Quantitative Study of LLM's Spatial Capability?

- LLM generally performed well on tasks related to spatial literacy, GIS theory, but revealed weaknesses in mapping and spatial reasoning
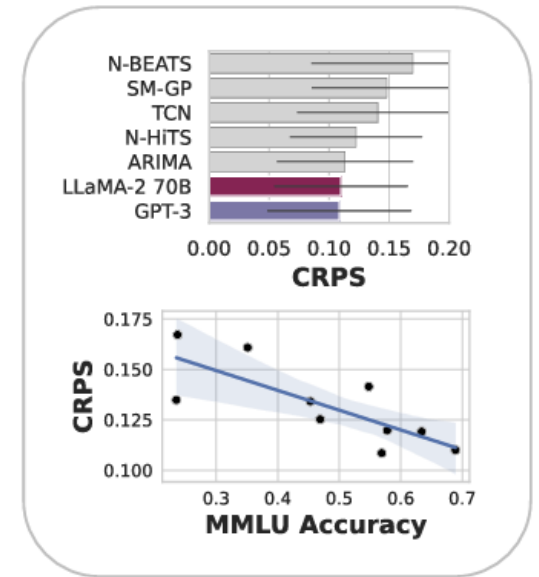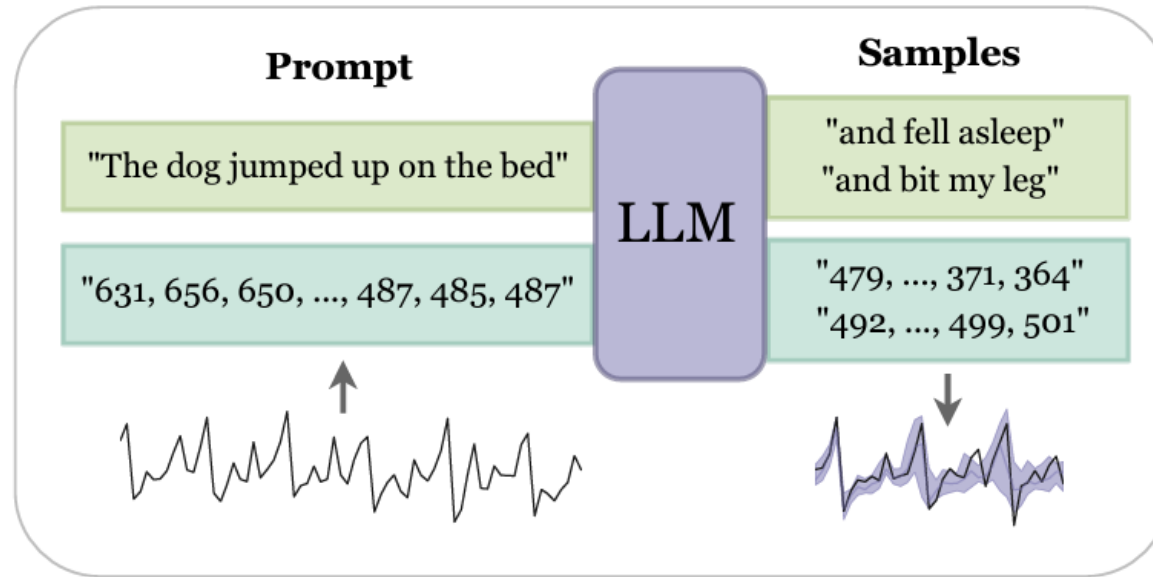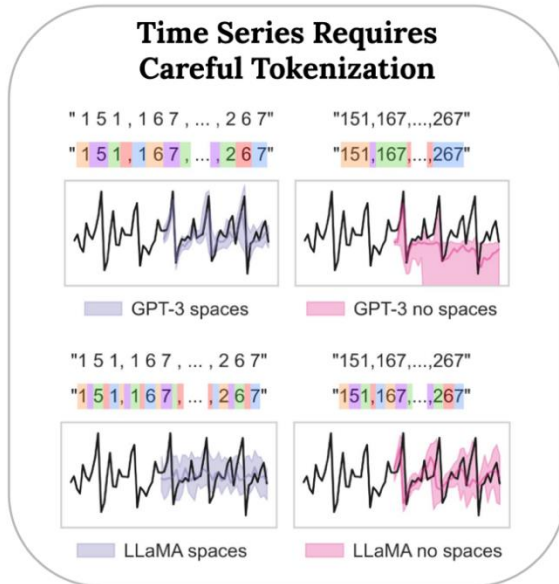


**Performance of 4 mainstream LLMs on 7 different spatial tasks.**

# Can We Stimulate LLMs for Temporal Analysis?

- LLM can be adopted for solving time series forecasting by translating timeseries to nature language

  - GPT-3 and LLaMA-2 can zero-shot extrapolate time series comparable to or exceeding the performance of task-specific time series models.



Nate Gruver, et al. Large Language Models Are Zero-Shot Time Series Forecasters. NeurIPS 2023.

# Can We Teach LLMs for Solving Urban Tasks?

- LLM can be finetuned to accomplish transportation tasks

  - Knowledge injected: traffic engineering documents, research papers, examination documents, etc.

  - Supported Tasks: synthetic traffic scenarios, explaining traffic phenomena, answering traffic-related questions, providing traffic recommendations, and generating traffic reports.



**LLM on solving traffic sign recognition, driving test, landmark recognition, captioning. LLM can be armed with multi-modal capability.**

Peng Wang, et al. TransGPT: Multi-modal Generative Pre-trained Transformer for Transportation. Arxiv 2024.

# Urban Foundation Models (UFMs): A Way to UGI



Urban Foundation Models (UFMs) are a family of large-scale models pre-trained on vast amounts of multi-source, multi-granularity, and multimodal urban data. They acquire notable general-purpose capabilities in the pre-training phase, exhibiting remarkable emergent abilities and adaptability dedicated to a range of urban application domains.

# OUTLINE

Introduction

**Challenges of Building UFMs**

Overview of UFMs

Prospects of UFMs

Summary

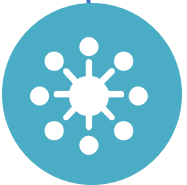# Challenges of Existing Foundation Models for UGI

## Data Misalignment

- Multi-domain heterogeneous urban data distributed in different parties/stakeholders.
- A bunch of well-developed urban tools with rich urban knowledge encoded.

## Skill Misalignment

- Tailored skill requirements to accomplish urban tasks, e.g., spatial & temporal calculation, multi-modal signal consolidation capabilities.
- Do way really need chat ability in all urban tasks?

## Task Misalignment

- Misaligned objectives between pretrain/fine-tune of general purpose foundation model and urban tasks.
- Alignment of human preferences may limit the urban problem-solving capability.

# Characteristics of Urban Big Data



**Urban big data**

**Multi-source**

**Multi-granularity**

**Multi-modal**

Sensor networks,
Mobile devices,
Satellites,
…

**Macro-level**:
Regional demand,
Population movement,
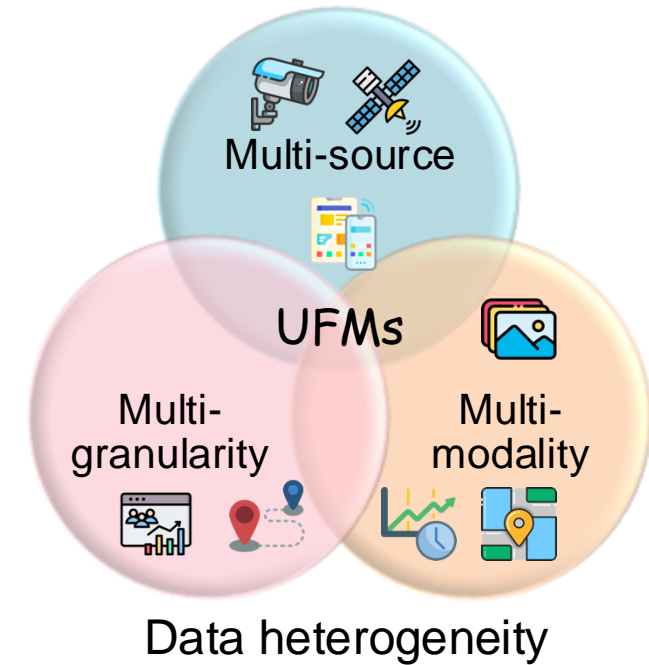Traffic flow
**Micro-level**:
Trajectory,
Social media

Text,
Image,
Trajectory,
Time series,
POI,
…

# Challenges of Data Misalignment

➢ **Multi-source, multi-granularity, and multi-modal data integration**
- Integrating diverse data sources at various granularities, from city-wide patterns to specific local details, poses a significant challenge due to the varied data modalities including text, images, and sensor readings.

➢ **Privacy and security concerns**
- Privacy and security are paramount, necessitating adherence to laws like GDPR, strong data anonymization, and defenses against data breaches and adversarial attacks to build trust and ensure ethical use.



Data heterogeneity



Privacy protection
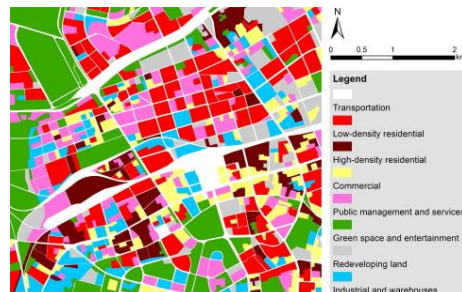
# Challenges of Skill & Task Misalignment

➢ **Spatio-temporal reasoning skills**

- Spatio-temporal reasoning aims to understand dynamic urban phenomena, which evolve over space and time through sophisticated modeling of temporal sequences and spatial distributions.



Spatial understanding     Spatial planning     Dynamic forecasting     Decision making
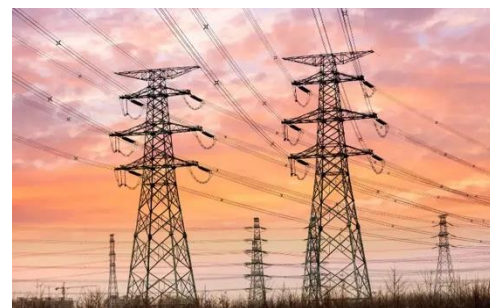
➢ **Versatility to diverse urban task domains**

- UFMs need to adapt to diverse urban task domains like transportation, energy, and environmental monitoring, ensuring versatility and effective generalization across different urban scenarios.



Traffic                    Environment                    Energy                    Climate