

Trustworthy Spatial Intelligence

Learning, Calibrating, and Reasoning Toward World Models of Cities

Dingyi Zhuang
Ph.D. Candidate @ MIT
Advisor: Prof. Jinhua Zhao

eMERGE seminar @ UC Berkeley
August 27, 2025

What's Up? 🙌

Not just the greeting... but also benchmark datasets!

What's "Up"?



A dog
on
a table

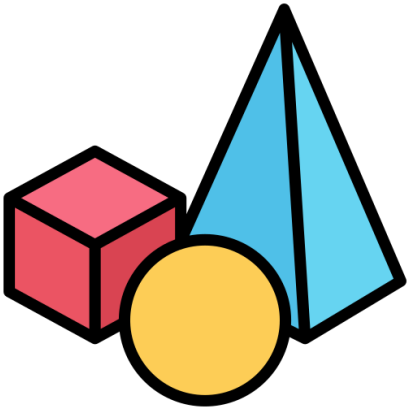


A dog
right of
a table

What does it really mean for a model to “understand space”?

What is Spatial Intelligence?

What things are



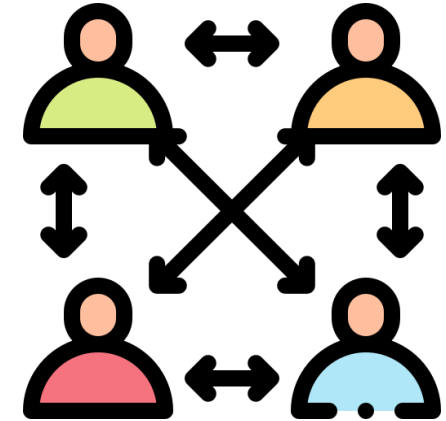
Object recognition

Where they are



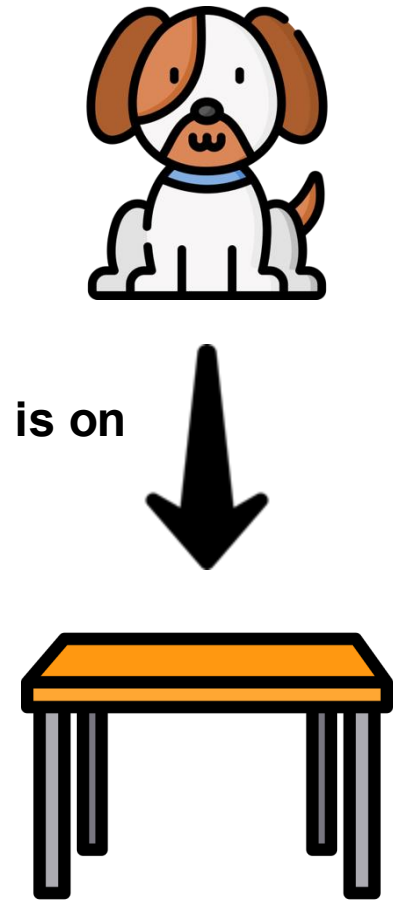
Spatial localization

How they relate



Relationships & interactions

Relationships Are the Foundations of Spatial Intelligence



Objects in isolation = limited meanings

Locations gain meaning through connections

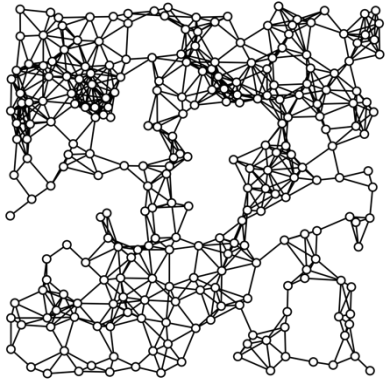
Relationship captures how entities influence or depend on one another

At its core, spatial intelligence is really about how well we can learn and represent relationships

Relationships in Actions

Macro & Micro Perspectives

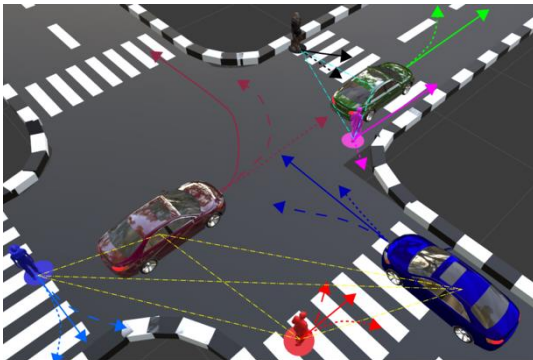
Geographic Information System



Urban value emerges from proximity and connectivity

(AAAI 2020)¹

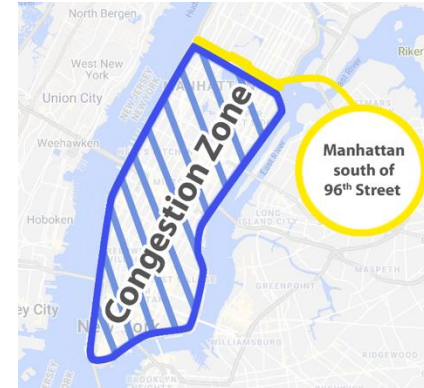
Autonomous Driving



Driving depends on relationships between vehicles, lanes, and pedestrians.

(On going work)

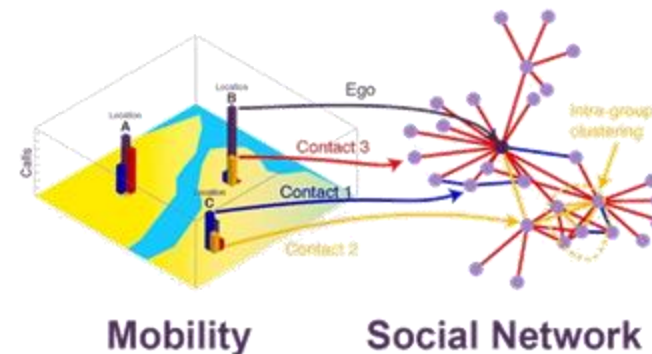
Transportation Policies



Policy relies on relationships between central and peripheral zones.

(On going work)

Social & Mobility Networks



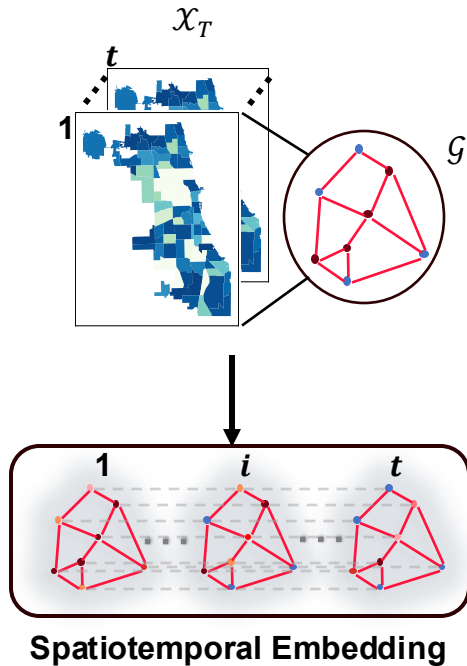
Mobility emerges from spatial relationships linking people, places, and opportunities

(On going work)

Research Agenda

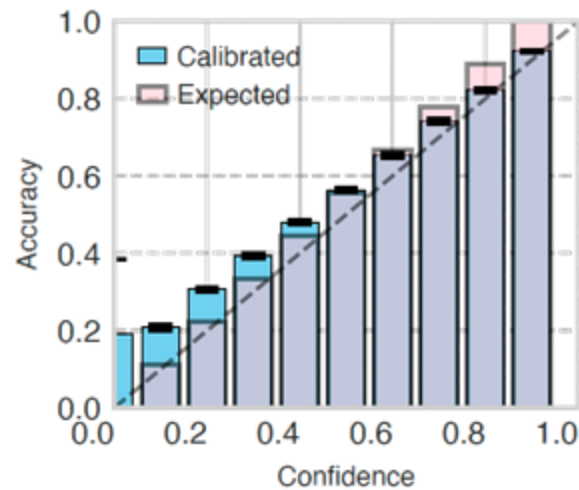
Learning, Calibrating, Reasoning about Relationships

Learning



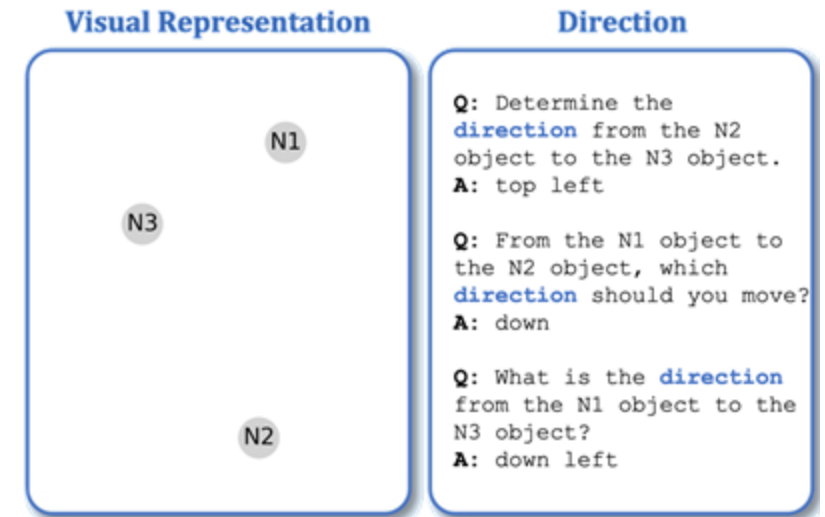
Learning patterns in spatiotemporal data

Calibrating



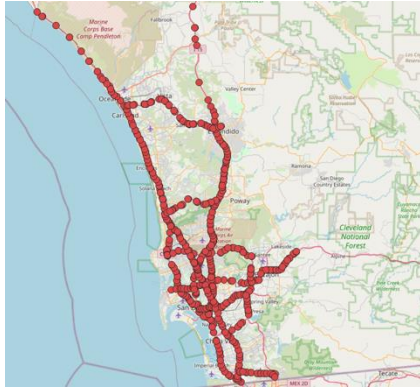
Calibrating model confidences and uncertainty

Reasoning



From pattern learning to reasoning with relationships

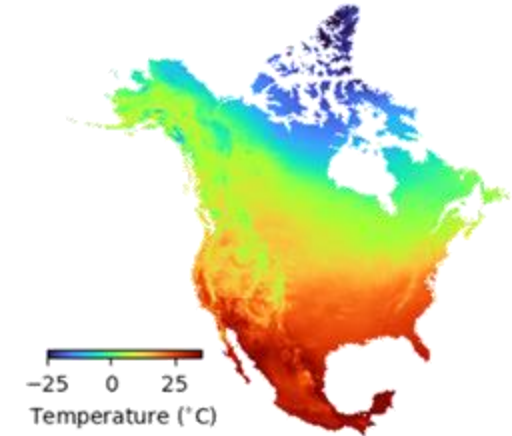
The Central Question: Learning Hidden Relationships in Cities



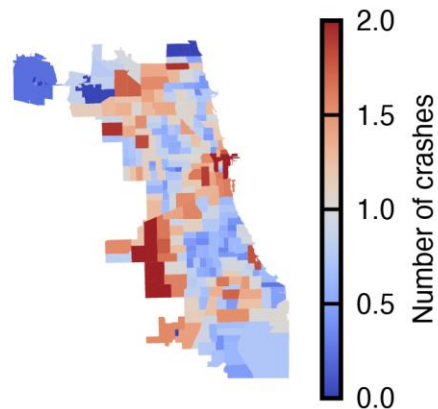
Highway Speed (San Diego)
(TRB 2025)²



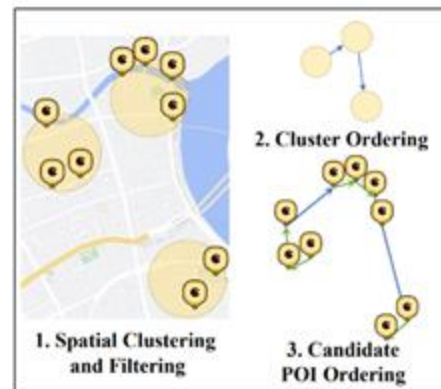
Transit Demand (Singapore)
(TR-C)³



Temperature (NA)



Traffic Crash (Chicago)
(SIGSPATIAL 2024)⁴



Point of Interests (Shanghai)
(EMNLP 2024)⁵

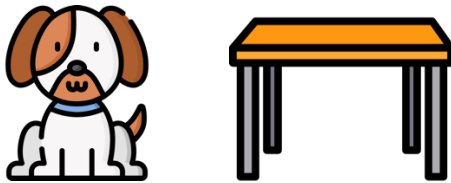
Challenges:

- Sparsity (high-resolution and missing information)
- High-dimensionality (city scale)
- Multi-modality (various data structures)

Graphs as Lenses for Learning Hidden Patterns

A graph G is the combination of

*Nodes \mathcal{V}
(entities)*



*Edges \mathcal{E}
(relationships)*



Adjacency Matrix A

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

Road network



Scene graph



Transit network



- *Origin-destination (OD) networks*
- *Social networks*
- *Land-use interaction graphs*
- *Accessibility graphs*
- *Urban knowledge graphs*
- *Distance-based graphs*
- ...

Graph Neural Networks (GNNs): Learning from Structure



- A GNN learns from graphs via **message passing**
- Each node **aggregates** neighbor info and **updates** its state
- Intuition: *“I adjust my choice based on what neighbors tell me”*



- Link prediction: the simplest way GNNs model relationships — decide if an edge should exist between two nodes.
- How: Based on updated node information
 - **Edge exists:** If two nodes become similar after exchanging neighbor info
 - **No edge:** If they remain very different

Using GNNs to Learn Travel Demand from Relationships

*Uncertainty Quantification of Sparse Travel Demand
Prediction with Spatial-temporal GNNs*



Dingyi Zhuang



Shenhao Wang



Haris N. Koutsopoulos



Jinhua Zhao

*ACM SIGKDD Conference on Knowledge Discovery and Data Mining 2022
Oral Presentation, <10%*

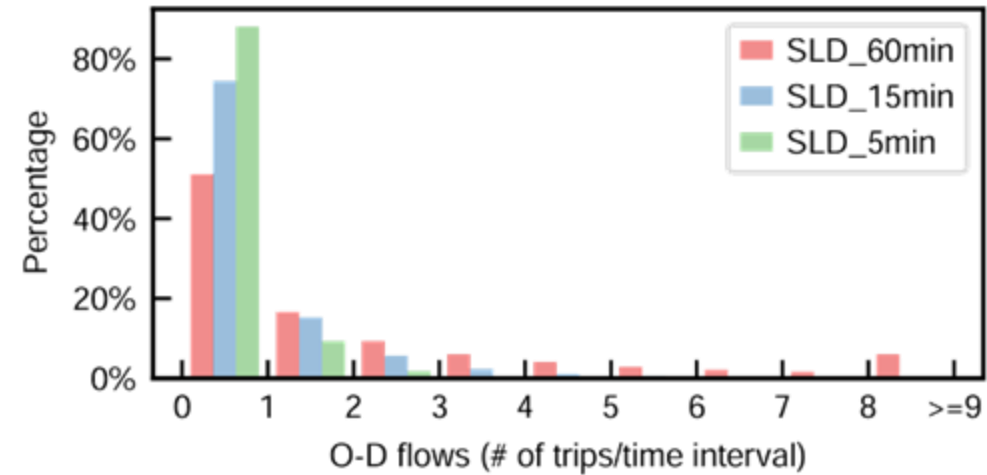
Sparsity in Spatiotemporal Transportation Data

For-Hire Vehicles in NYC



High-resolution OD demand are highly sparse, with many zero entries

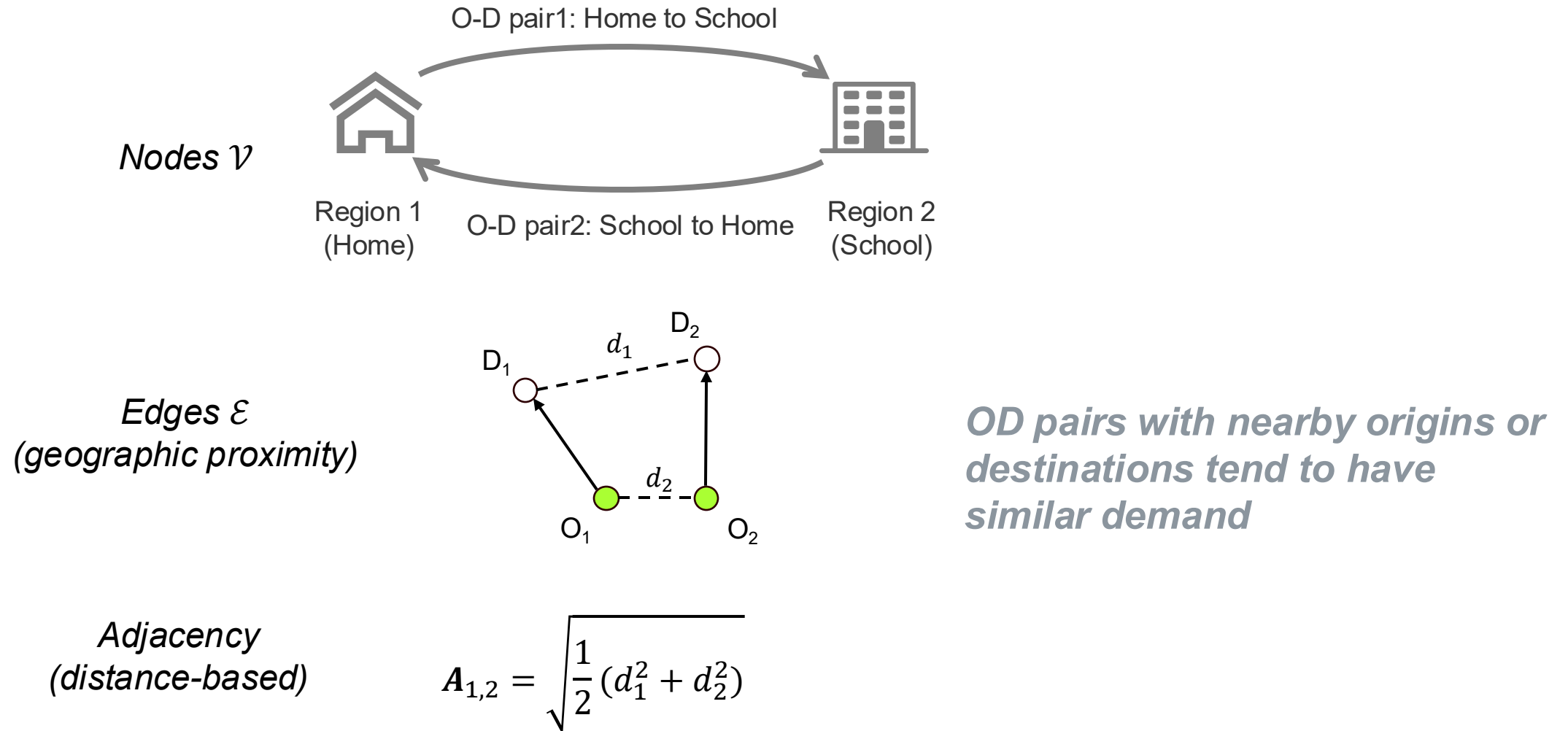
Distribution of travel demand



67 pick-up/drop-off zones, 67×67 OD pairs in total

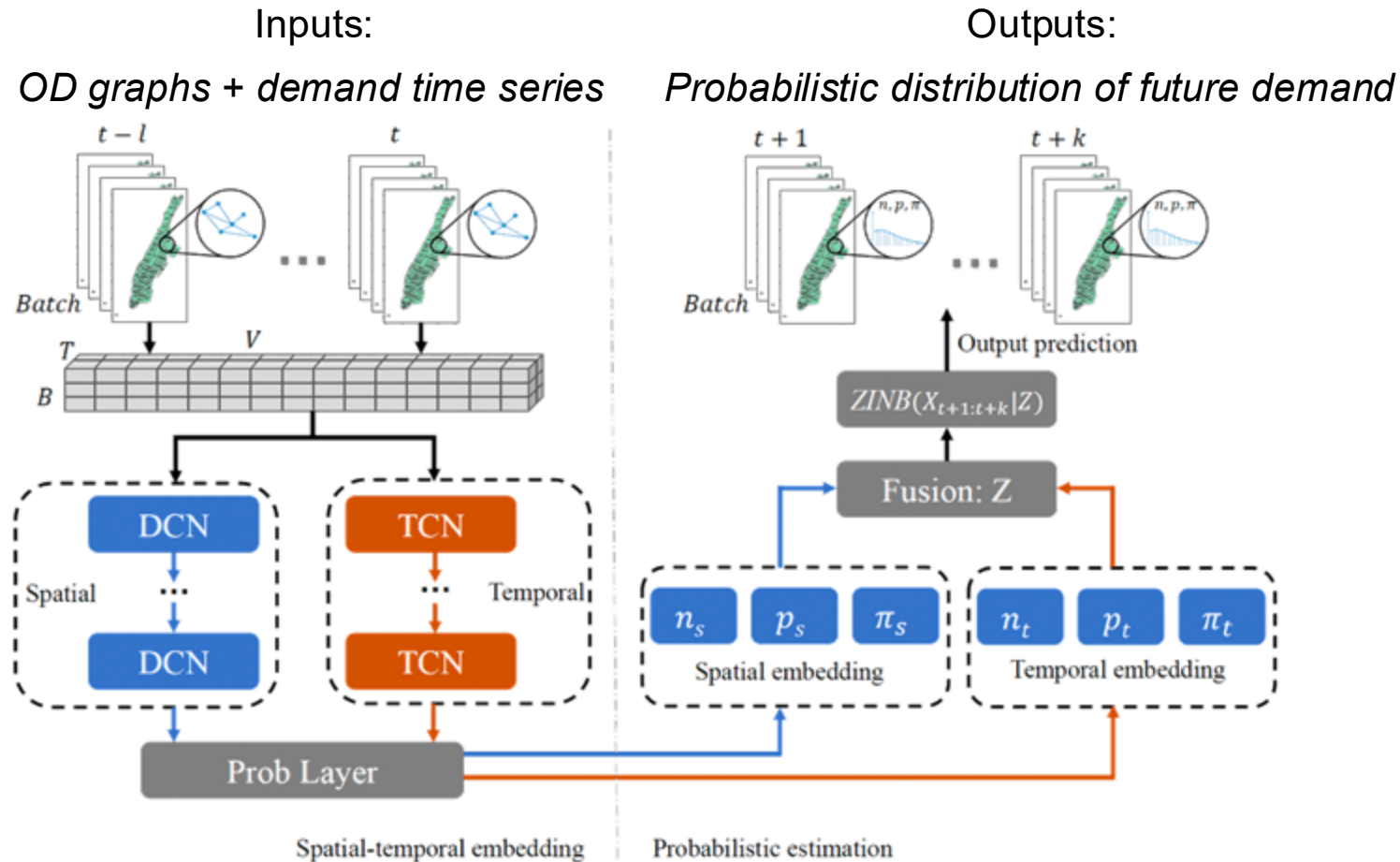
Sparsity is ubiquitous if scaling up spatial/temporal resolutions

Graph Representation of OD Demand



Modeling Sparse OD Demand with Graph-Based Relationships

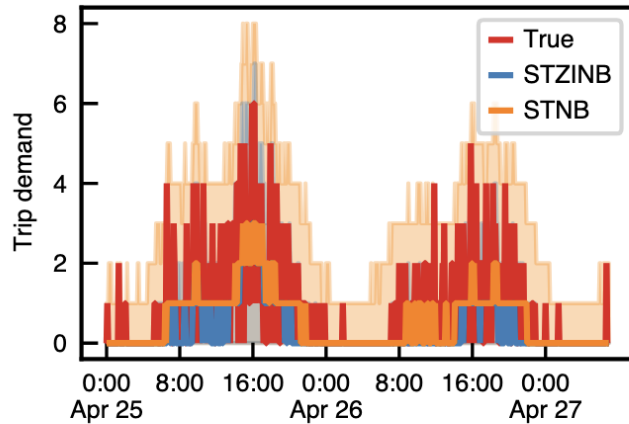
$$X \sim \pi\delta_0 + (1 - \pi)NB(n, p)$$



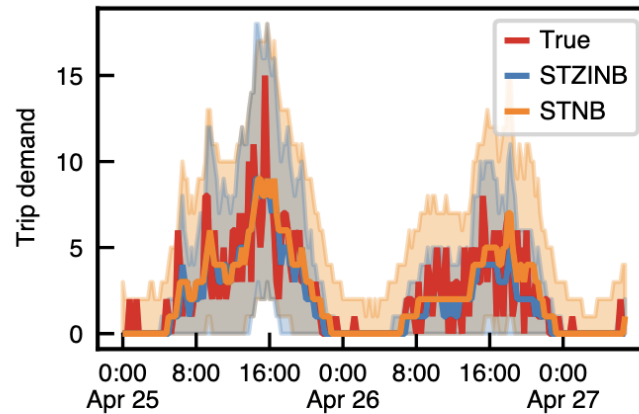
- Transform sparse OD demand into a probabilistic distribution
- Zero-inflated modeling handles excess zeros
- GNN capture spatial relationships among OD pairs

δ_0 : Dirac delta distribution at zero (i.e. point mass at 0)
 NB : Negative binomial distribution

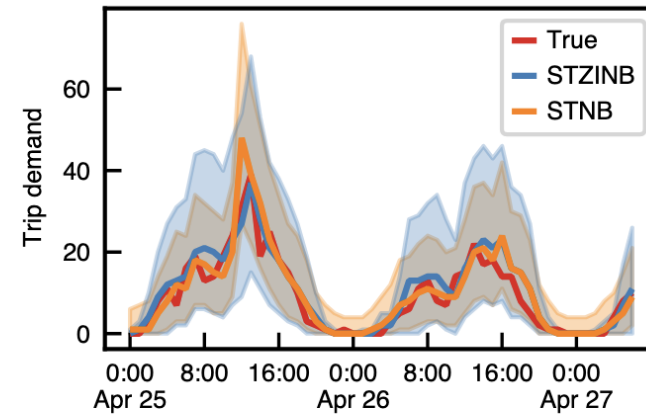
Results: Performance under Sparse Data



(a) SLD_5min case



(b) SLD_15min case



(c) SLD_60min case



Overall **6%** accuracy gains compared to baselines



Handling extreme sparse cases (**90%** data entries being zeroes)



Efficient prediction intervals (\geq **55%** narrower than non-zero-inflated models)

Improper Relationships in GNNs Can Propagate Untrustworthy and Inequitable Results

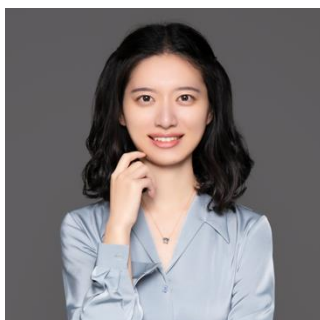
Mitigating Spatial Disparity in Urban Prediction Using Residual-Aware Spatiotemporal Graph Neural Networks: A Chicago Case Study



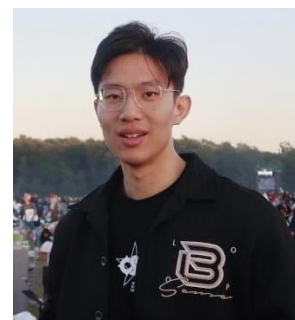
Dingyi Zhuang



Hanyong Xu



Yunhan Zheng



Xiaotong Guo



Shenhao Wang



Jinhua Zhao

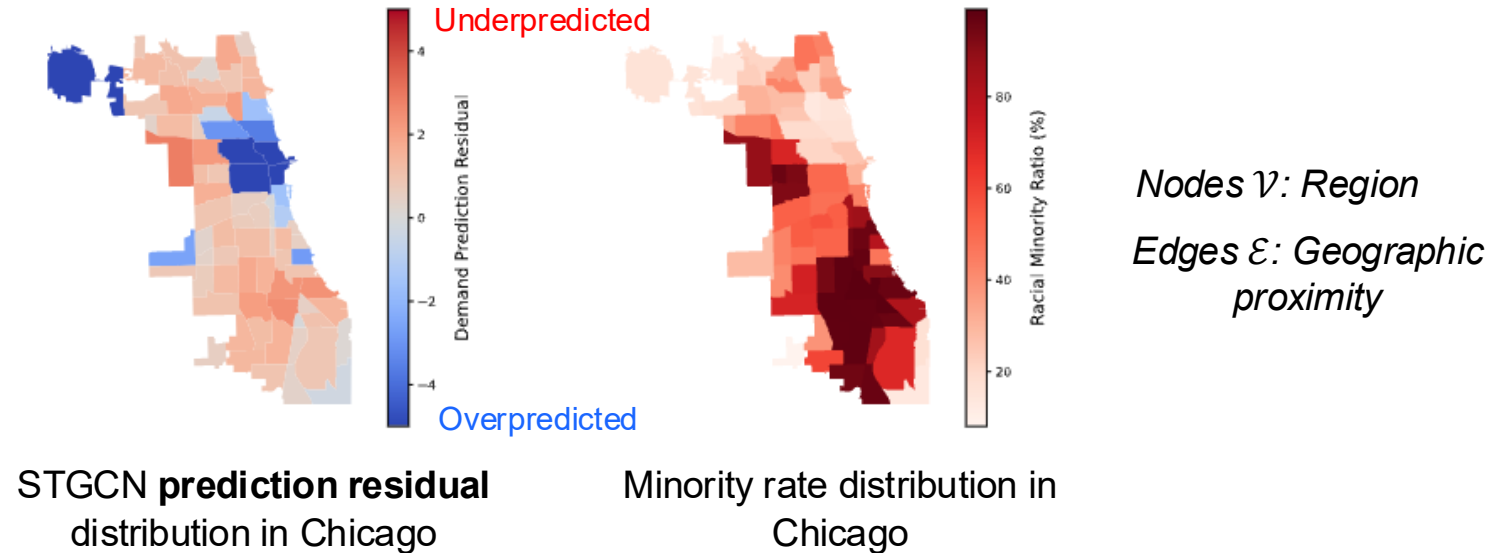
International World Wide Web Conference 2025
🏆 **Best Paper Award at WebST Workshop**

How Spatial Disparity Emerges

Transportation Network Companies (TNCs)



Pick-up demand from TNCs

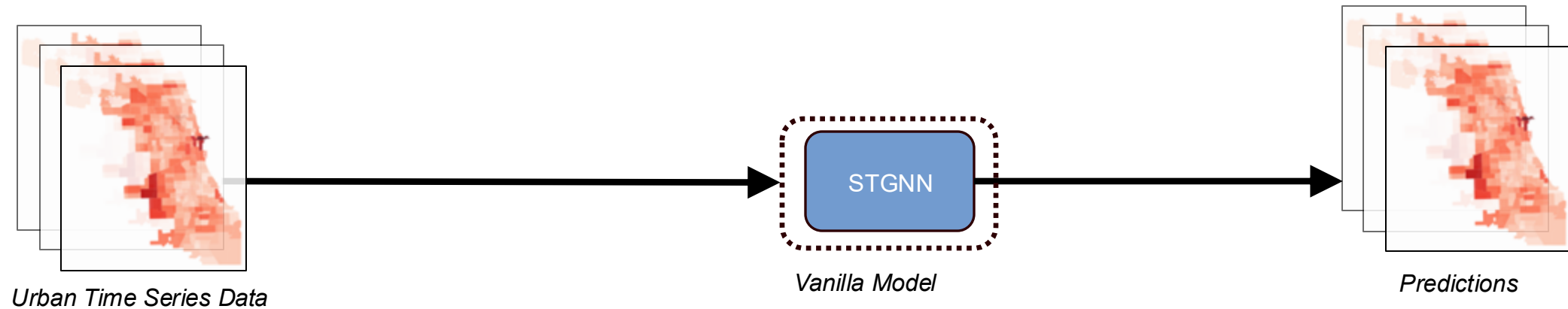


- Adjacency matrix propagates information from neighbors.
- If your neighbors are **underpredicted**, so are you. Leading to **underserved** demand.
- In transportation, this means entire regions can be systematically underpredicted.
- ***Accurate predictions do not necessarily lead to equitable outcomes***

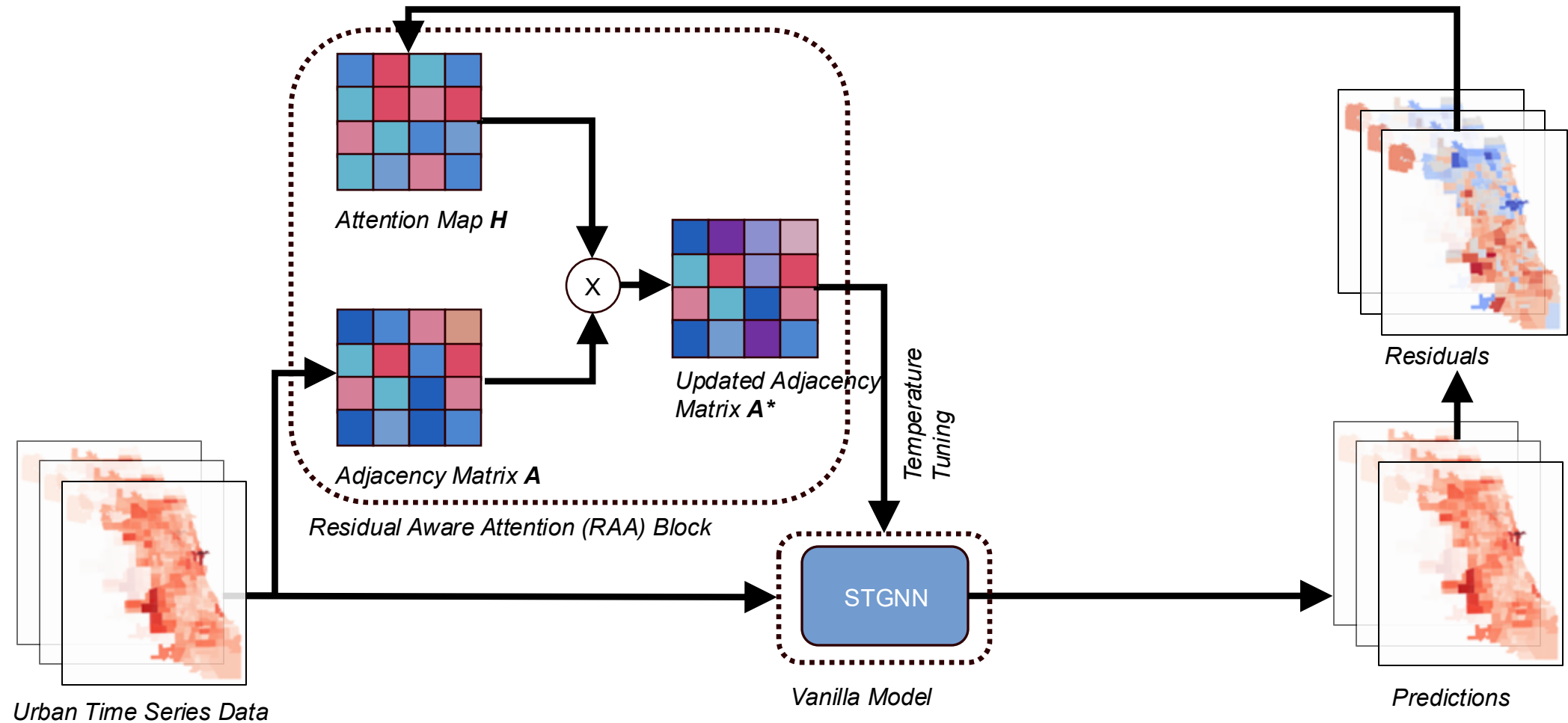
What do you observe?

Residual-Aware Attention: Rethinking Adjacency

Residual-aware block: adjusts adjacency weights using residual signs



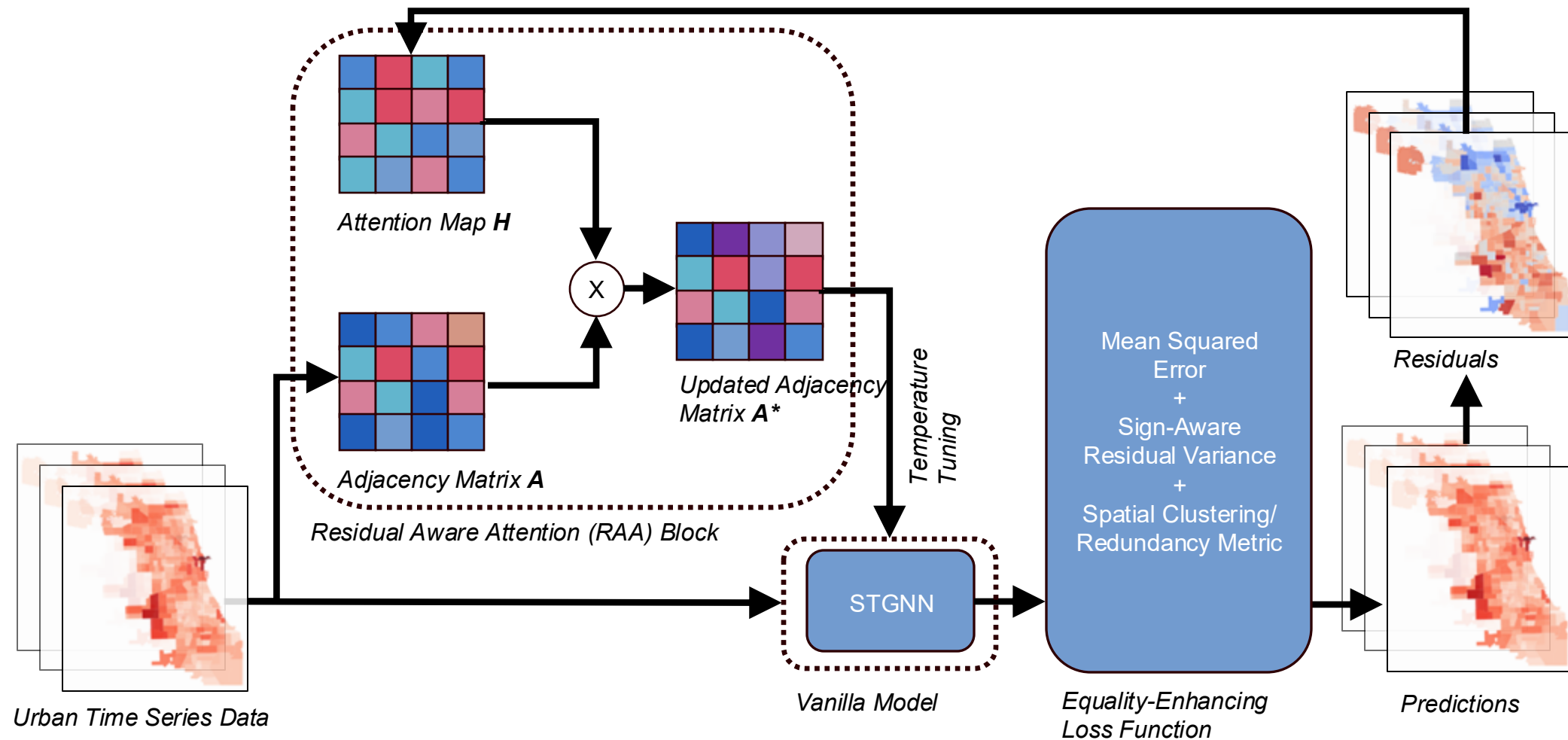
Residual-Aware Attention: Training



$$A_{adapted} = A \odot H$$

$$H = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

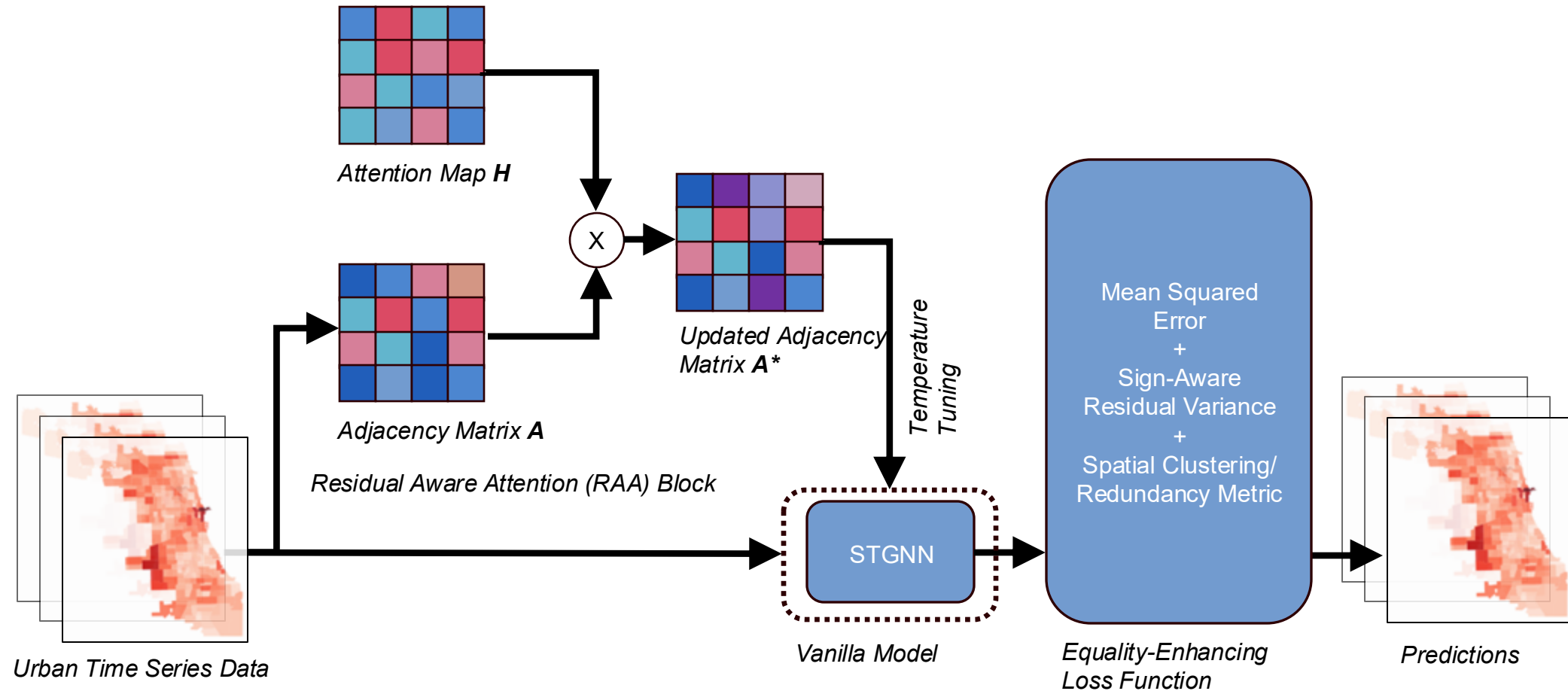
Residual-Aware Attention: Training



Include the spatial disparity D_s and fairness metrics D_f like Moran's I into loss function

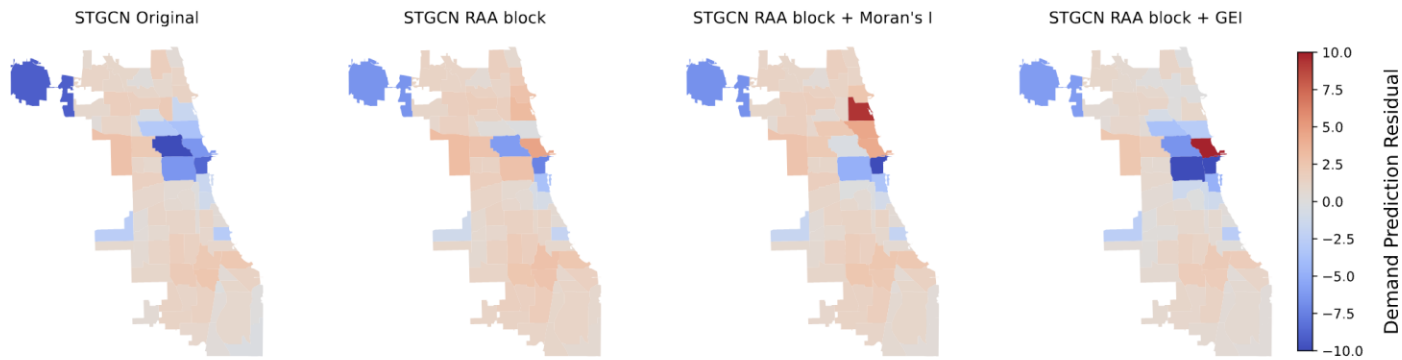
$$\mathcal{L}_{joint} = \mathcal{L}_{prediction} + \lambda_s D_s + \lambda_d D_f$$

Residual-Aware Attention: Inference

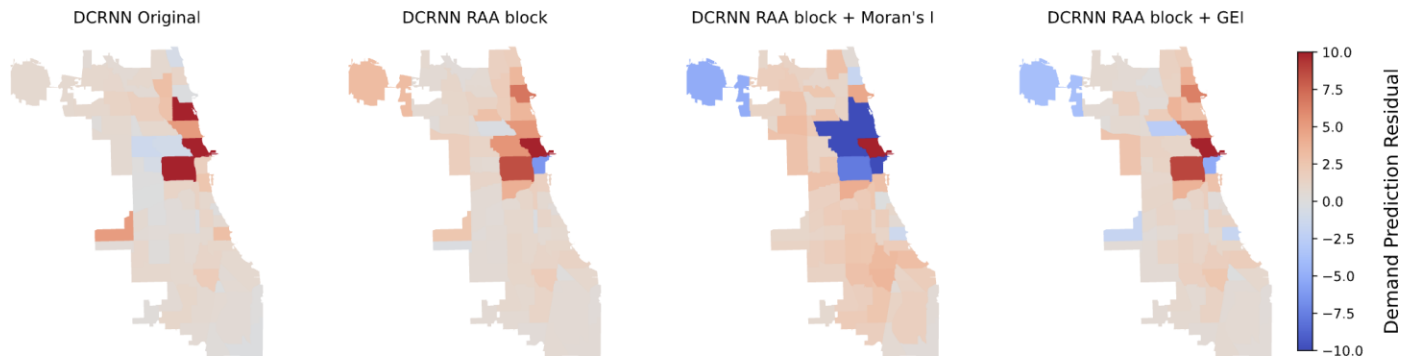


The attention map is fixed during model inference

Results: More Equitable with Minimal Trade-off



(a) STGCN Residual Distribution



(b) DCRNN Residual Distribution

Small accuracy loss

MAE

7%

SMAPE

12%

Fairness metrics improvements

GEI

18%

Moran's I

80%

SDI

47%

MAE: Mean Absolute Error

SMAPE: Symmetric Mean Absolute Percentage Error

GEI: Generalized Entropy Index, measures spatial disparity

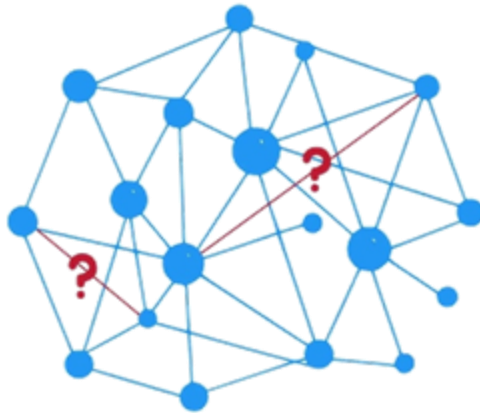
Moran's I: measures spatial autocorrelation

SDI: Scaled Disparity Index, measures demographic disparity

From Learning to Trust

Can We Really Rely on These Models?

Graphs as a lens



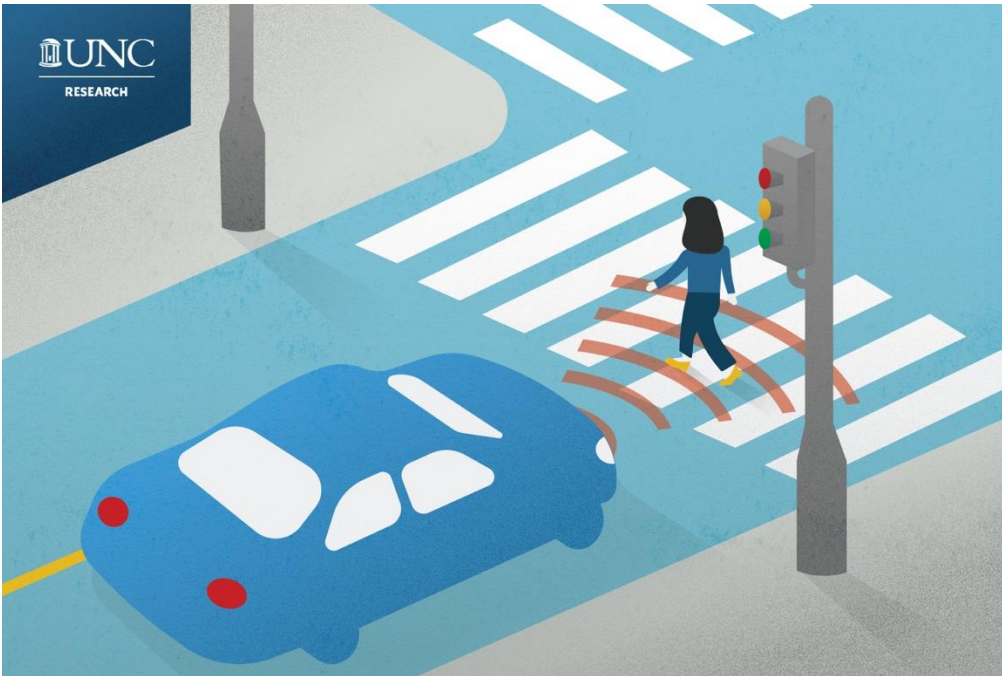
Trust?



- Graphs are powerful representations.
- They learn hidden patterns.
- They could make accurate predictions.
- But... can we trust their predictions?

What is Calibration?

Calibration = aligning model confidence with reality



50% chance a pedestrian is on the crosswalk across 100 similar simulations/training



80% of the labels show the pedestrian is on the crosswalk, you are **underconfident!**

Calibration bridges the gaps

Calibrating Graph Link Predictions for Trustworthy Topology in Autonomous Driving

GNN-Based Topology Refinement for Map Generation



Dingyi Zhuang



Xiaoqi Wang



David Paz



Wenbin He

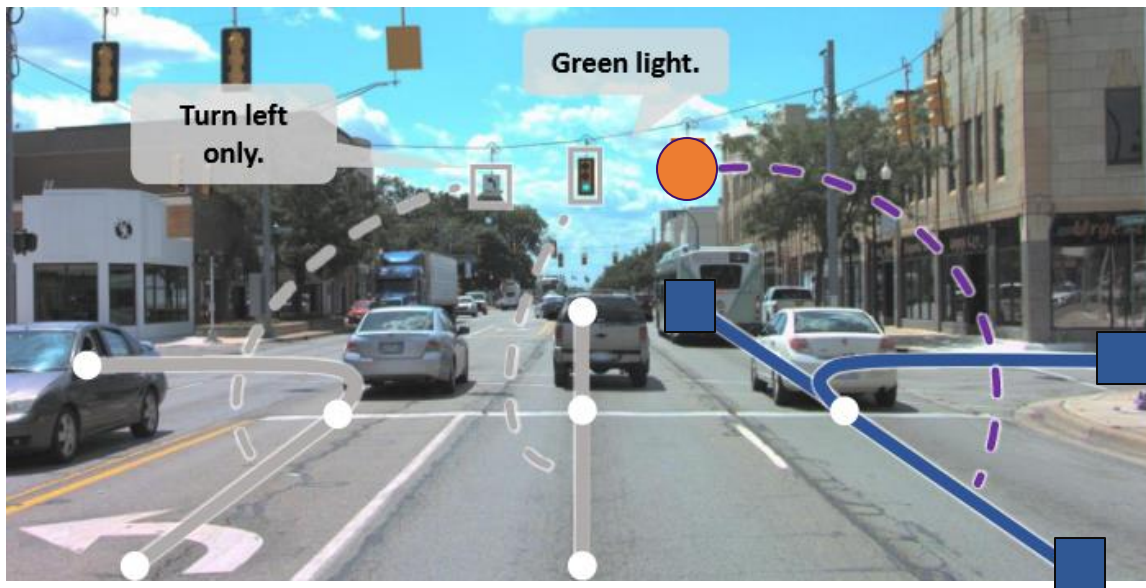


Liu Ren

*Internship work at Bosch Center for Artificial Intelligence during Summer 2025
In submission to ICLR 2026*

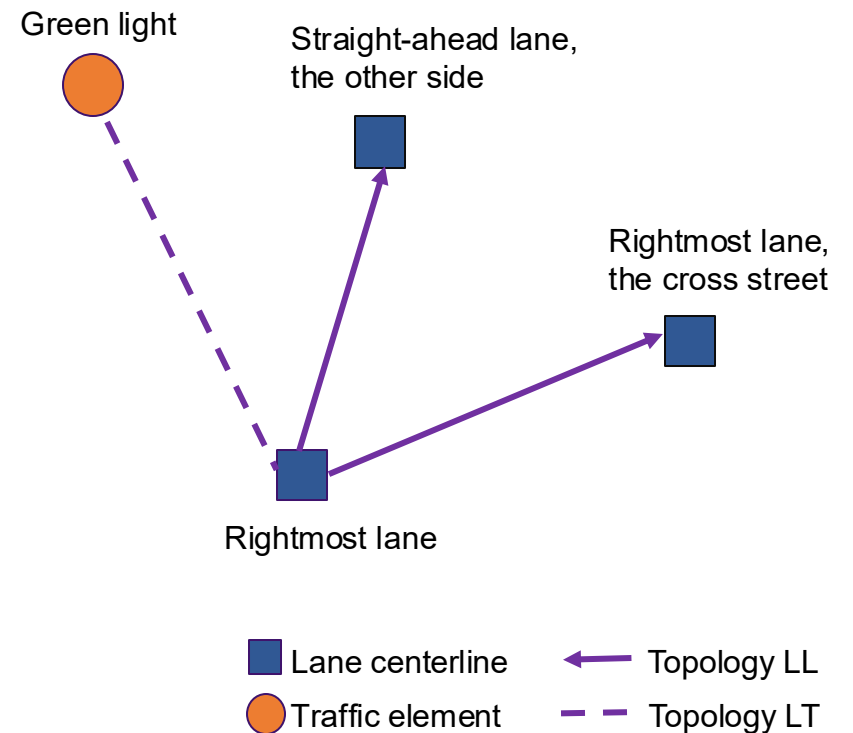
Driving Scene Graph

An autonomous vehicle navigating towards an intersection, this is what front-view camera sees.



Which lane to drive?

Which traffic signal to follow?

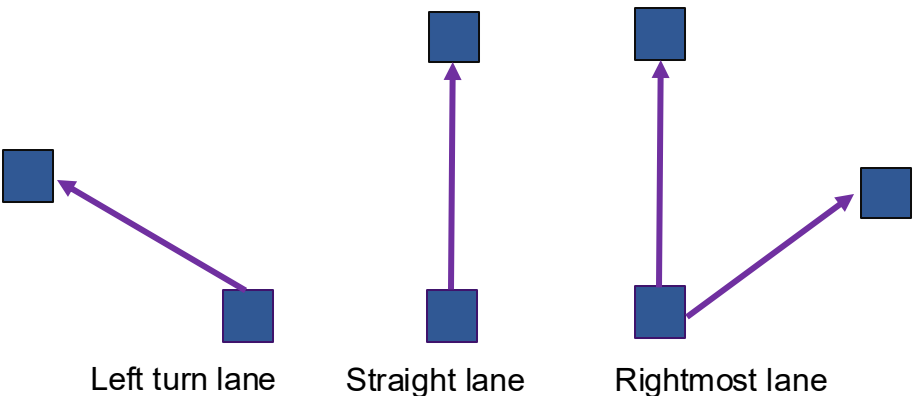


Driving Scene Topology Graph

How to ensure topology is reliable?

Lane topology graph

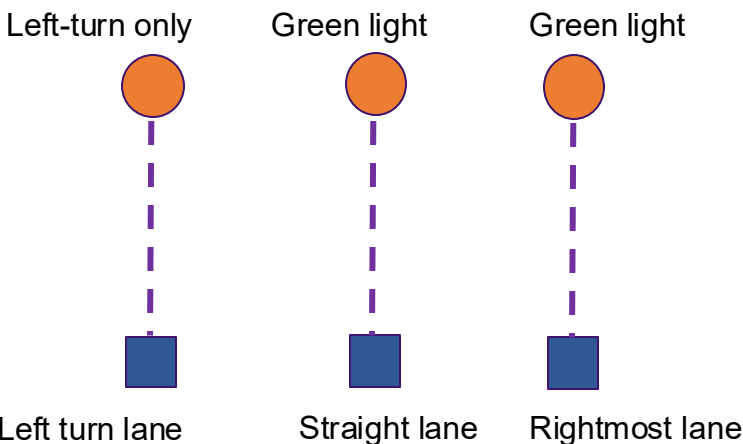
Consists of centerlines as well as their connectivity



Full lane topology graph of the example

Lane-traffic element topology graph

Lane-to-traffic element assignment
(e.g., lights, signs, markers)



Full lane-traffic element topology graph of the example

■ Lane centerline

● Traffic element

← Topology LL

- - Topology LT

HD Maps – Current Solutions

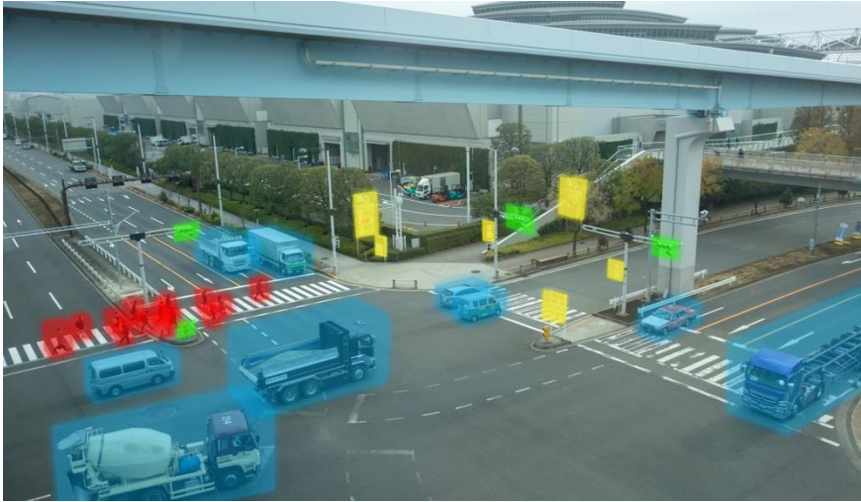


Fig: Vehicle navigation with High Definition (HD) mapping



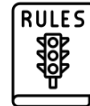
Precise geometry of each lane



Capture how lanes connect (e.g., left turn)



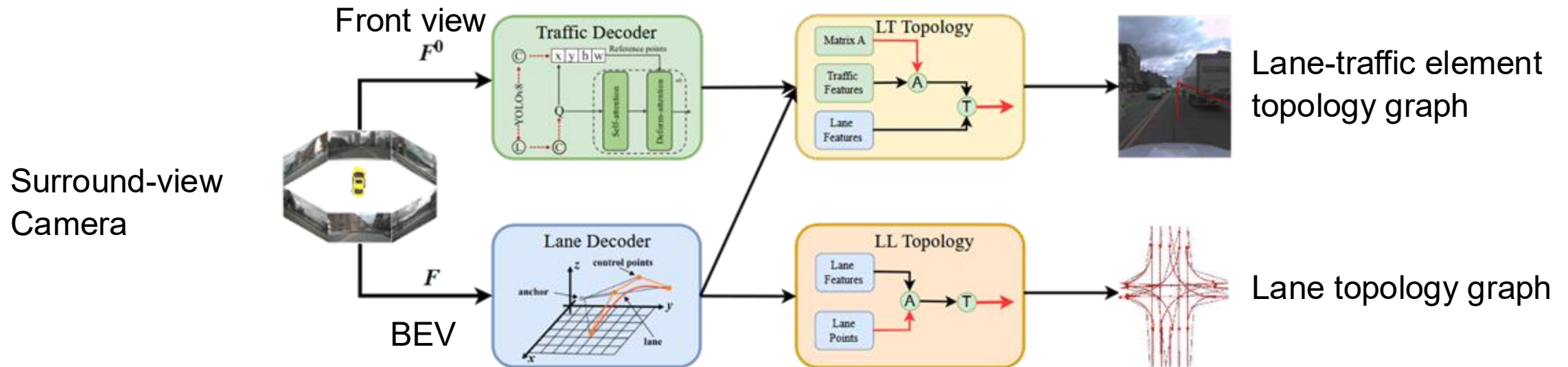
How traffic elements controls lanes



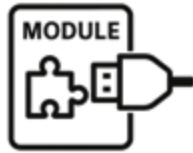
Driving rules: speed limits, priority rules at intersections, and so on

Topology Reasoning – Concept

Real-time understanding of objects and predicting topology graphs in driving scenes

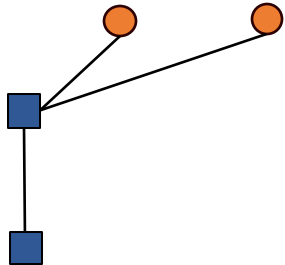


Graph Self-Supervised Learning (GSSL) based Refinement/Calibration

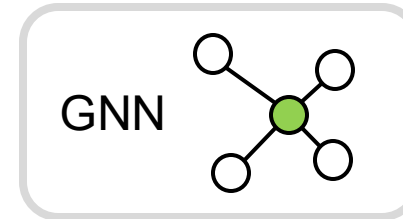
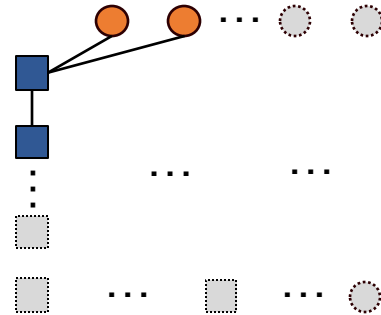


A **universal refinement module, plug-and-play** with existing topology reasoning models to calibrate topology graph

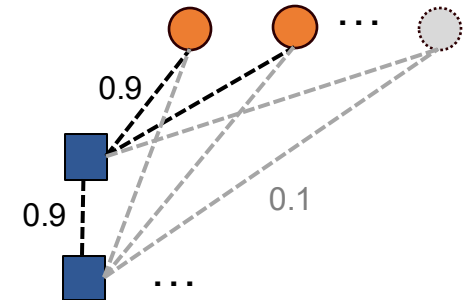
Ground-truth topology
from annotation



Graph augmentation



- Link prediction
- Differentiate **true** v.s. **fake** connections
- **Calibrate** confidence of the relationship



SSL: The Engine Behind Modern AI

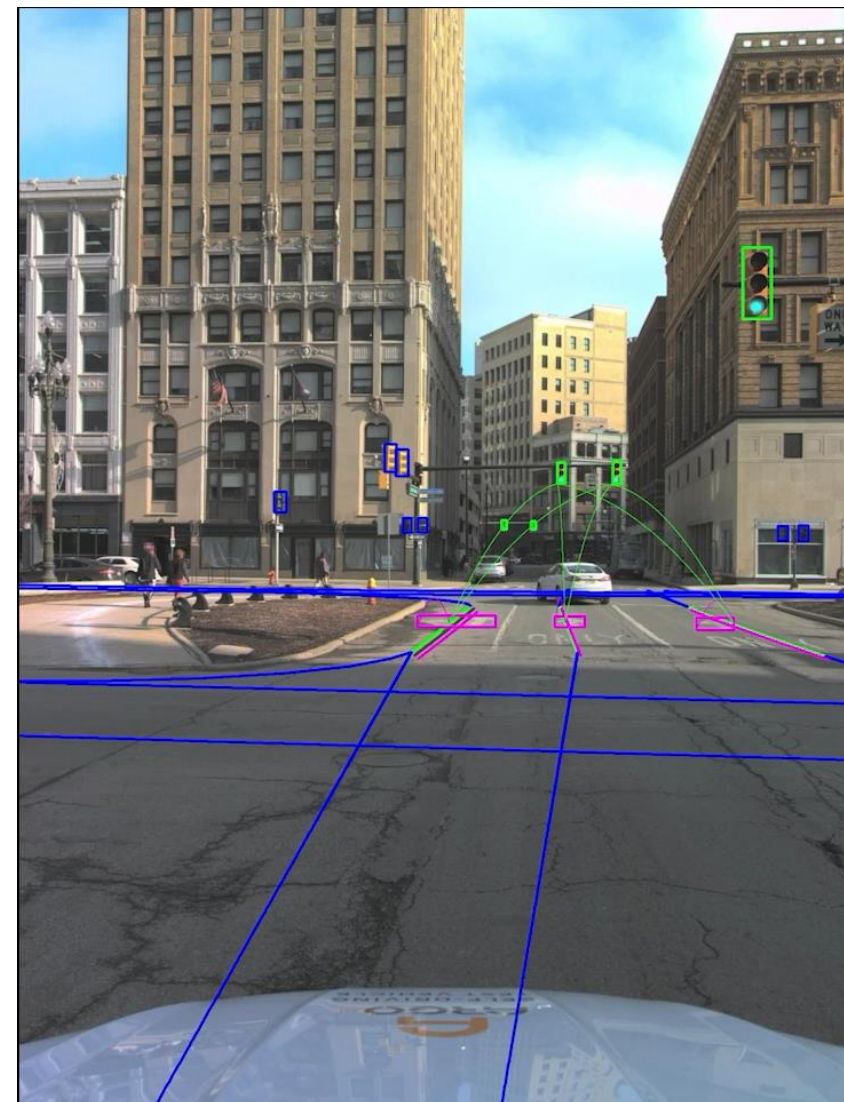
- Learn from raw data without human labels
- Pretext tasks: predict missing words (text), masked pixels (vision), missing edges (graphs)
- Rich representations
- Foundation of most frontier AI systems today



OpenLaneV2 Dataset and Evaluations

- Combine **Argoverse 2** and **nuScence** datasets.
- ~2,000 annotated road scenes with 72K frames @ 2Hz
- Annotate **lanes polylines**, **traffic elements bounding boxes**, and the driving scene topology.

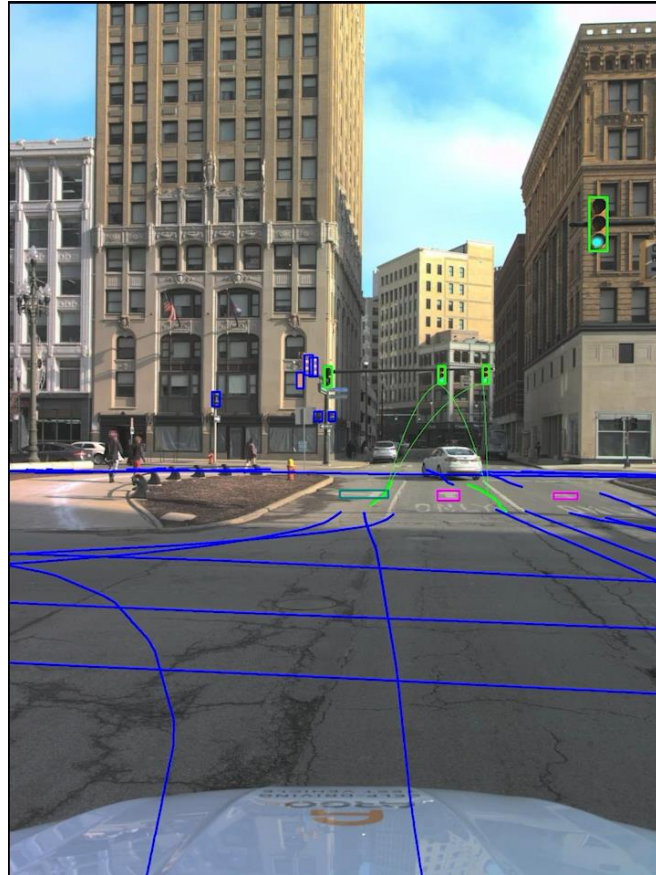
Top_{ll} & Top_{lt} mAP of predicted topology



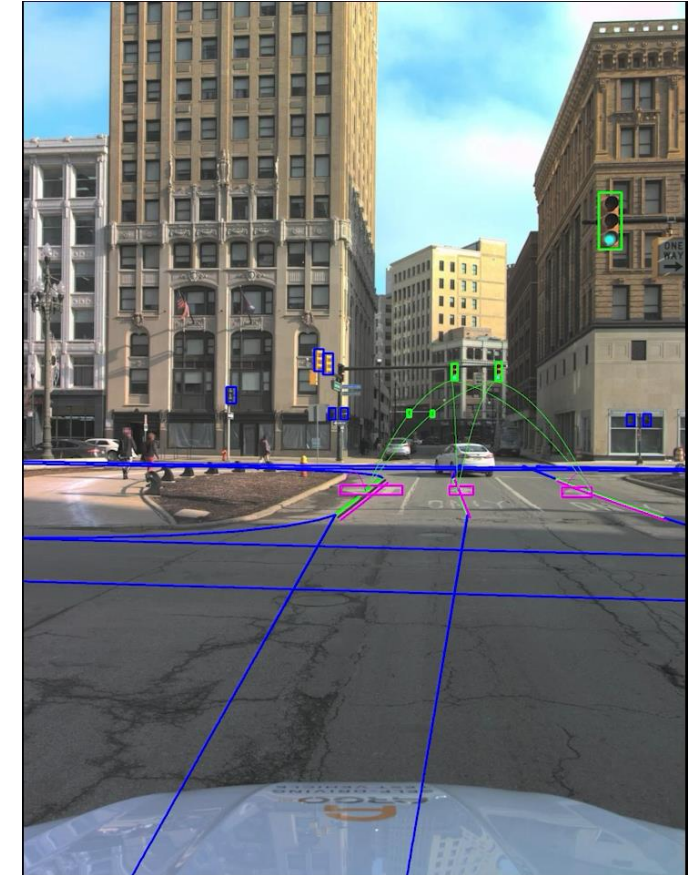
Uncalibrated Results

- Lane-traffic elements topology not well detected
- Lanes-lane topology inconsistent

TopoNet



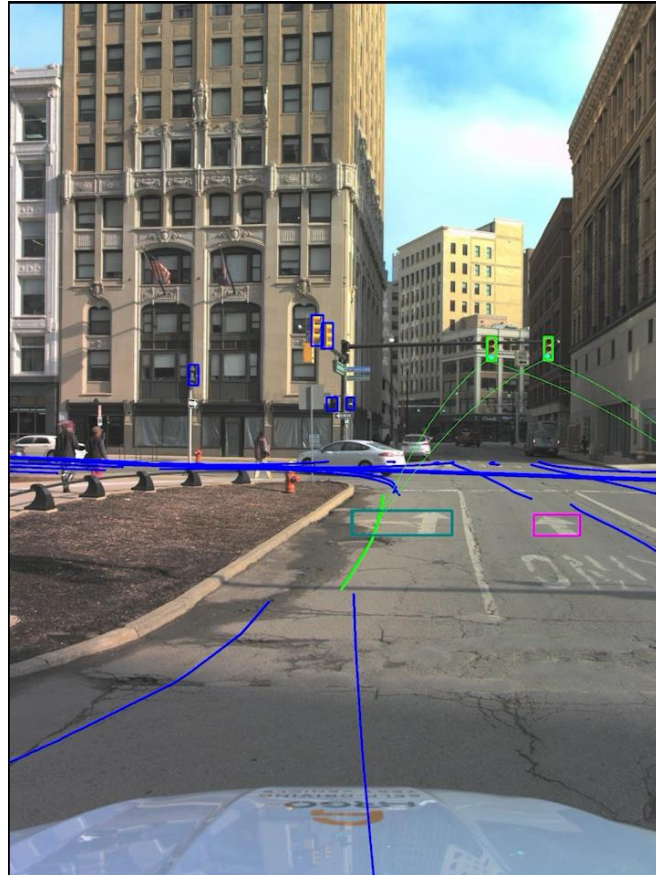
Ground-truth



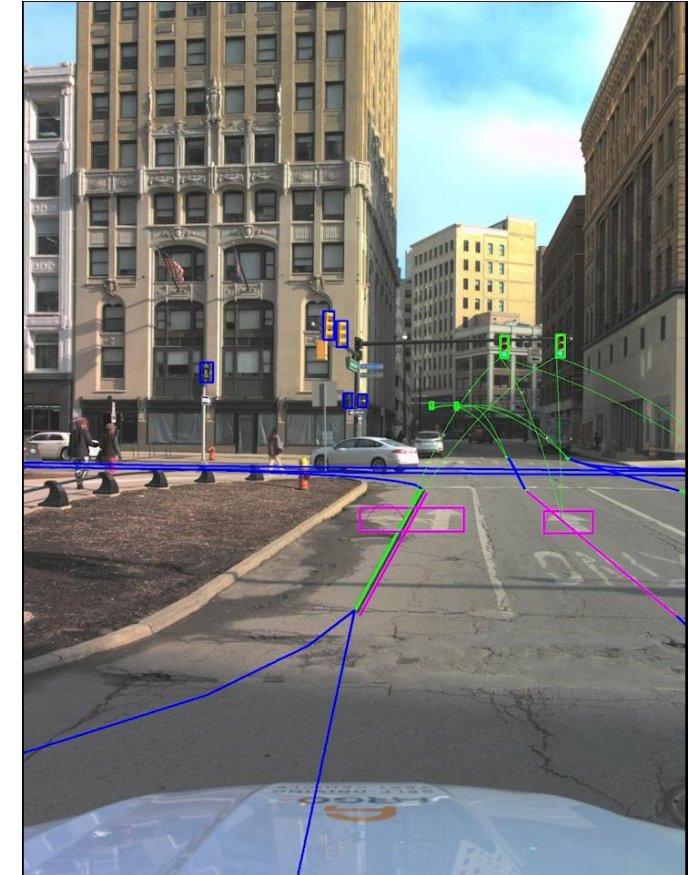
Uncalibrated Results

- Lane-traffic elements topology not well detected
- Lanes-lane topology inconsistent

TopoNet



Ground-truth



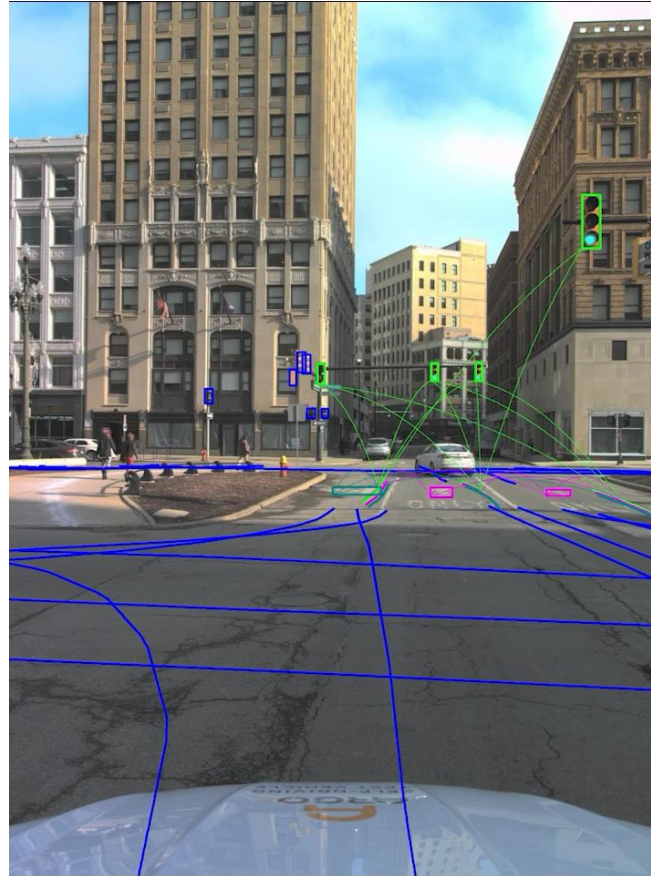
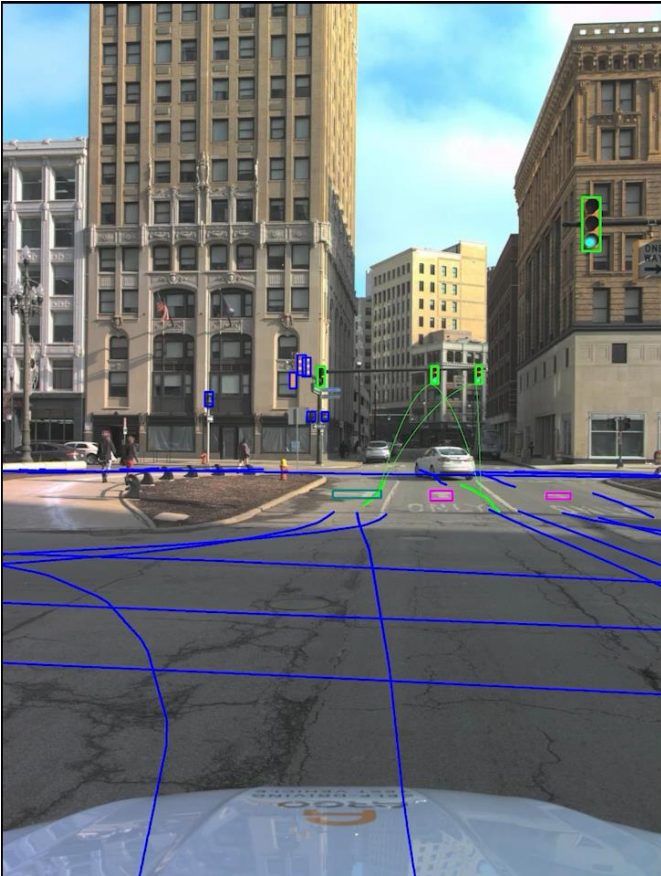
Calibrated Results

Top_{ll} 100% ↗
Top_{lt} 20% ↗

TopoNet

After calibration

Ground-truth



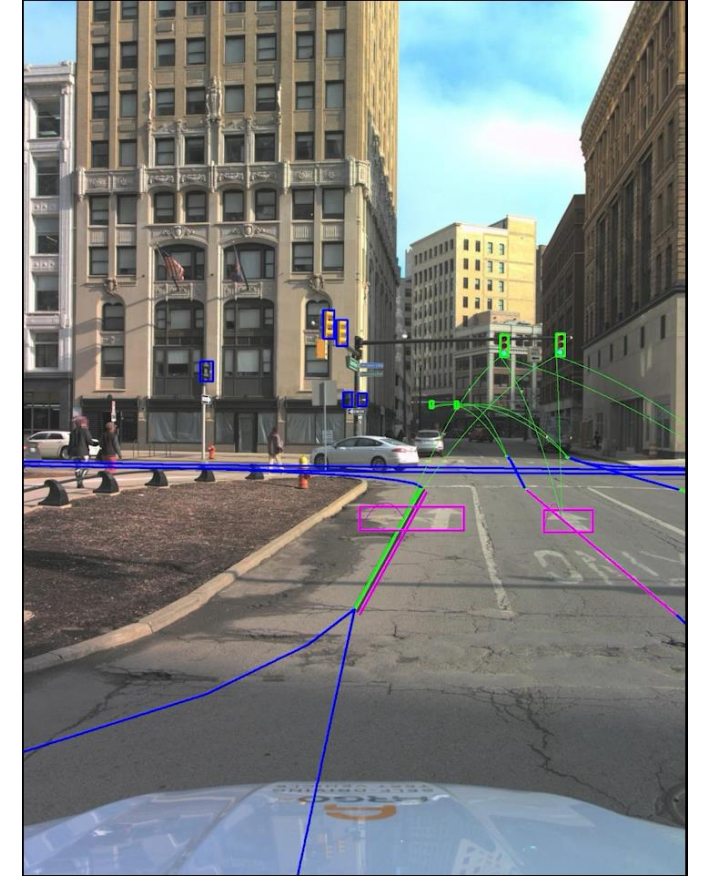
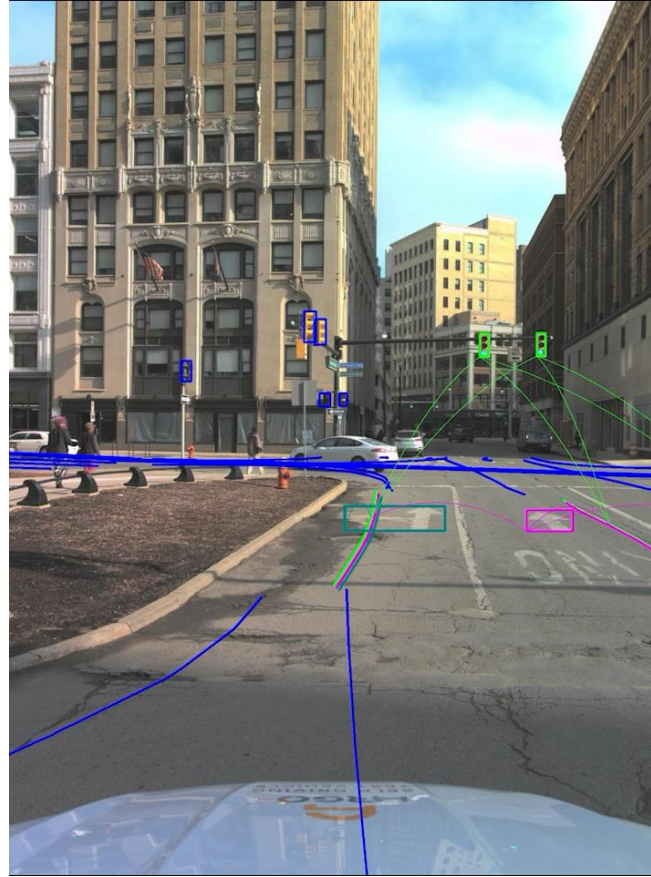
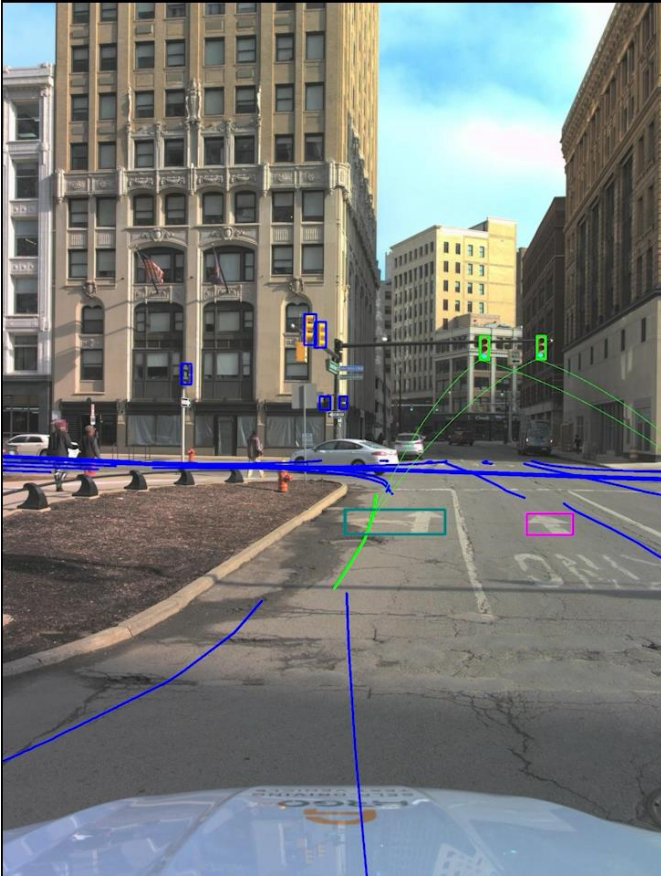
Calibrated Results

Top_{ll} 100% ↗
Top_{lt} 20% ↗

TopoNet

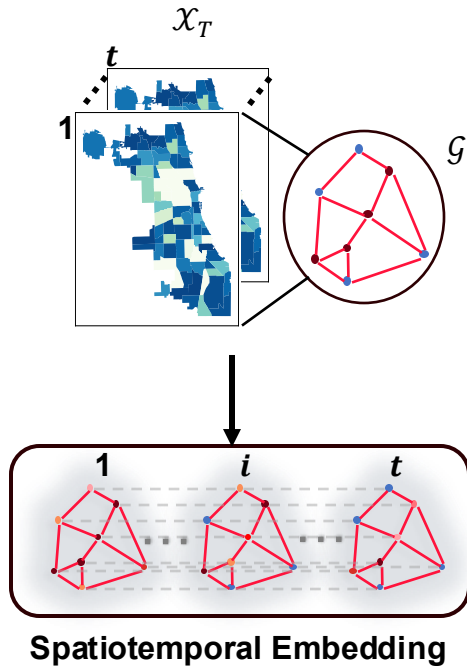
After calibration

Ground-truth



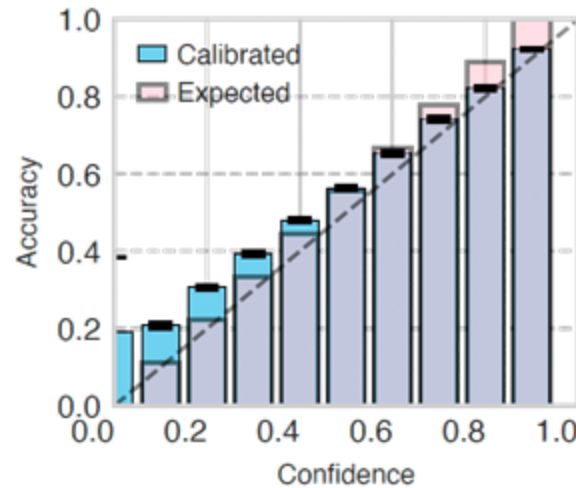
Learning + Calibrating = Trustworthy Patterns

Learning



Learning patterns in
spatiotemporal data

Calibrating



Calibrating model
confidences and uncertainty

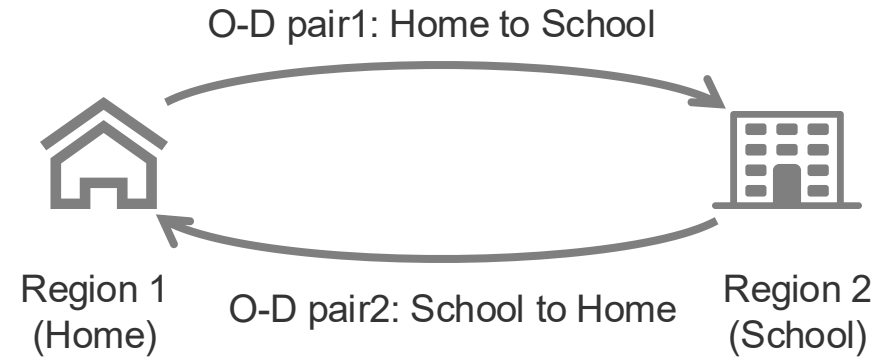
*Trustworthy Patterns:
Predictions we can rely on*

Patterns ≠ Intelligence

Patterns Tell Us What; Reasoning Explains Why

Patterns

(what tends to happen)



OD pairs with nearby origins or destinations tend to have similar demand

Reasoning

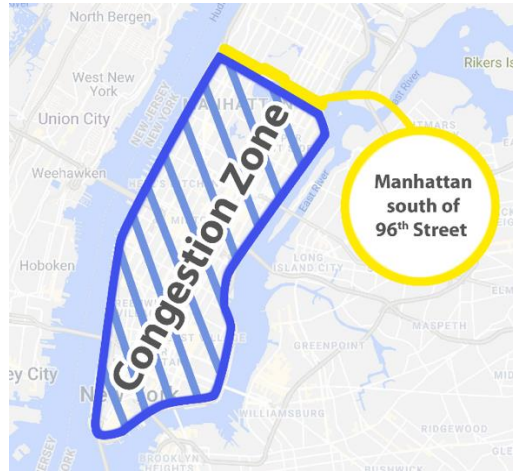
(Why it happens, and when the rule applies)



Nearby origins link to the same job centers, and nearby destinations share accessibility

Why Reasoning Matters

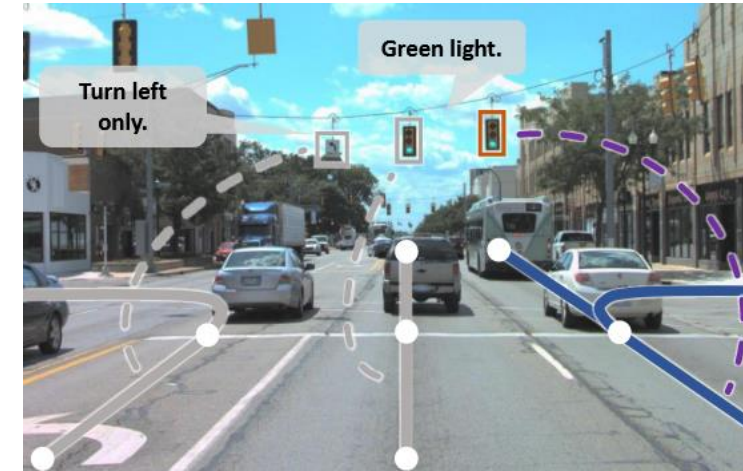
Policy-critical



Urban Planning

requires causal and social reasoning about accessibility and equity.

Safety-critical



Autonomous Driving

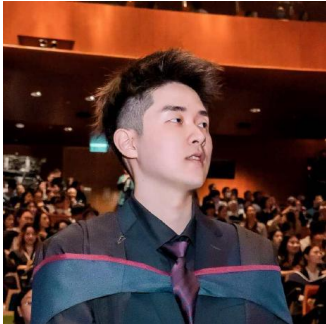
requires reasoning about lanes, pedestrians, traffic rules, and social norm – not just patterns



General AI systems for transportation and urban planning require spatial reasoning to effectively navigate environments and support real-world interactions

Spatial Reasoning: Fine-tuning VLMs

Sparkle: Mastering basic Spatial Capabilities in Vision Language Models
Elicits Generalization to Composite Spatial Reasoning



Yihong Tang*



Ao Qu*



Zhaokai Wang*



Dingyi Zhuang*


Zhaofeng Wu, Wei Ma, Shenhao Wang, Yunhan Zheng, Zhan Zhao, Jinhua Zhao



Accepted at EMNLP Findings 2025
Best Paper Award, MKLM Workshop @ ICJAI 2025

Spatial Reasoning Gap in VLMs


State-of-the-art Vision Language Models (VLMs) fail to solve the pathfinding problem, a simple 2D spatial reasoning task

Image


Question

Given the grid graph with nodes (N1, N2, ... N16), find the shortest path from the green node to the red node.

Answer:
["N5", "N9", "N13"]

 **ChatGPT 4o (Commercial)**

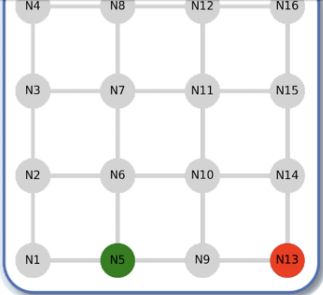
The shortest path from N5 to N13 is:
["N5", "N6", "N10", "N9", "N13"] .

 **InternVL2-Pro (Open-sourced)**

The shortest path can be represented as:
["N5", "N6", "N10", "N11", "N15", "N14", "N13"] .

ChatGPT 4o >

ChatGPT



Given the grid graph with nodes (N1, N2, ..., N16), find the shortest path from the green node to the red node? Return your answer in English and only include the node sequence.

N5 → N6 → N10 → N14 → N13

*ChatGPT 4o still makes mistakes...
Screenshot on Aug 26, 2025*

Key Research Questions



How well do existing models perform on 2D spatial reasoning tasks?

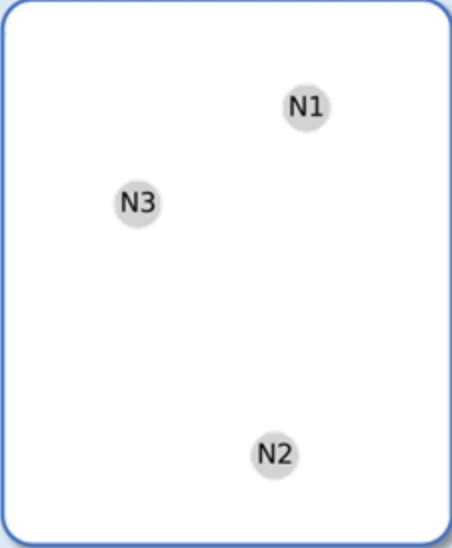


What are the fundamental capabilities that underpin spatial reasoning?



Can mastering basic capabilities lead to better performance on more complex, composite tasks?

Disentangling Spatial Reasoning

Visual Representation	Direction	Localization	Distance
	<p>Q: Determine the direction from the N2 object to the N3 object. A: top left</p> <p>Q: From the N1 object to the N2 object, which direction should you move? A: down</p> <p>Q: What is the direction from the N1 object to the N3 object? A: down left</p>	<p>Q: Which relative location is the N2 located at? A: down right</p> <p>Q: Which relative position is the N1 located at? A: top right</p> <p>Q: Identify the location of the N3. A: top left</p> <p>Q: In a 10x10 image, what is the coordinate of the N3 object? A: (3.59, 7.34)</p>	<p>Q: Which distance is the shortest? N1 and N2, N1 and N3, N2 and N3 A: N1 and N3</p> <p>Q: Compare the distances: N1 and N3 and N2 and N3. Which one is longer? A: N2 and N3</p> <p>Q: In a 10x10 image, what is the distance between the N1 and N3 objects? A: 4.07</p>



Direction: Understanding the relative orientation between objects

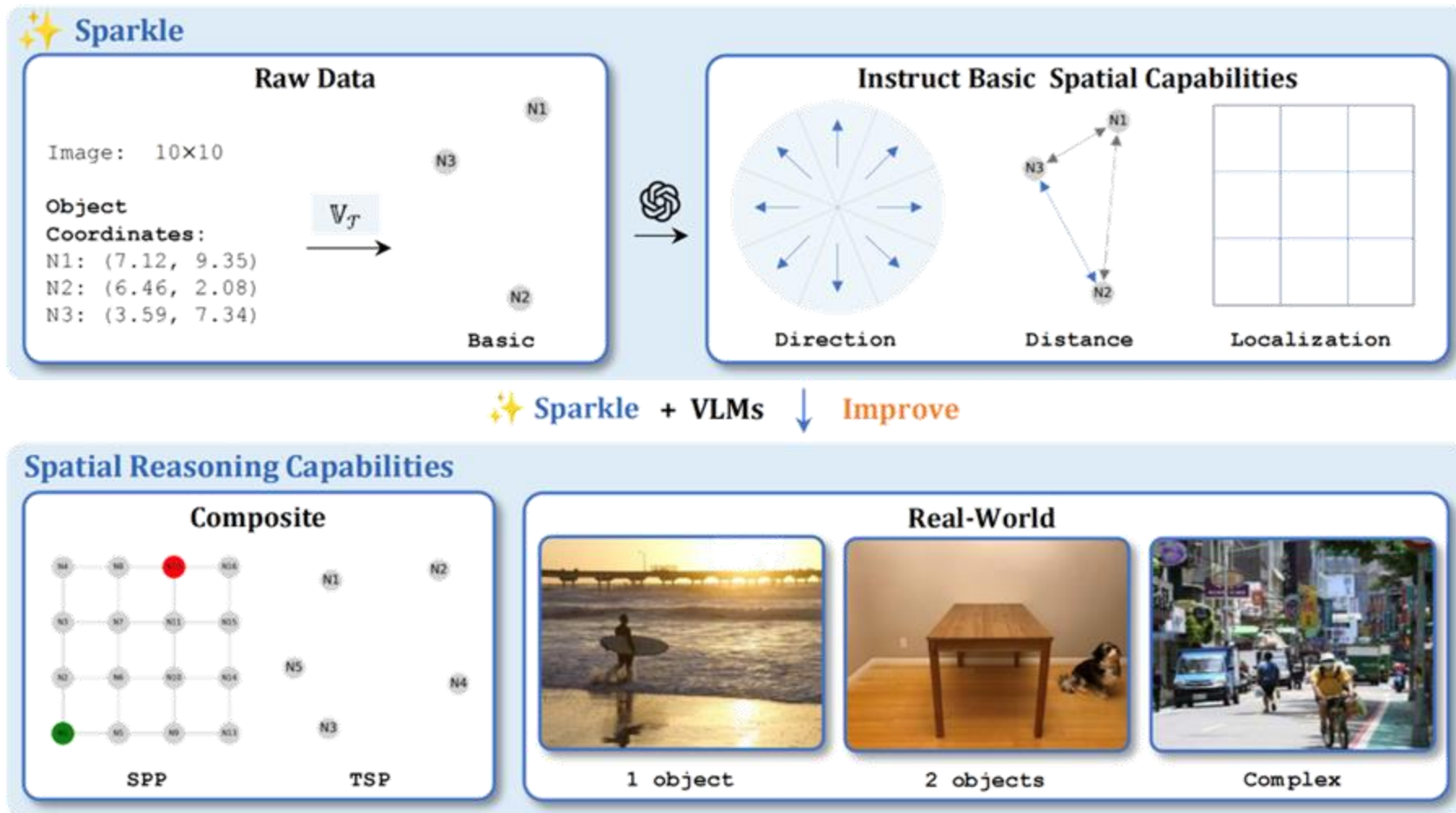


Localization: Determining an object's precise position in space


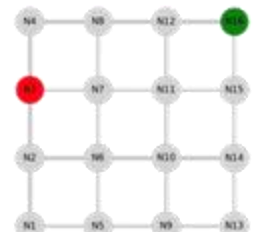


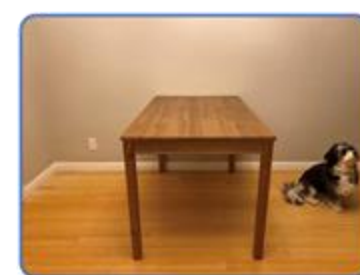


Distance: Measuring the spatial displacement between objects

The Sparkle Framework



A Multi-level Approach

Basic Spatial Relationships Understanding	Shortest Path Problem	Traveling Salesman Problem	General Spatial VQA Tasks (1 Object)	General Spatial VQA Tasks (2 Objects)
 <p>Q: Which distance is the shortest? Options: A. N1 to N4, B. N1 to N3, C. N4 to N3 A: A</p> <p>Q: Determine the direction from N1 to N2. Options: A. top left, B. top right, C. down left, D. down right A: B</p> <p>Q: What is the position of the N4 object? Options: A. top left, B. top, C. top right, D. left, E. center, F. right, G. down left, H. down, I. down right A: I</p>	 <p>Q: The image shows a grid graph where each node is labeled (N1, N2, ... N16) and connected to neighboring nodes.</p> <p>Based on the image, find the shortest path from the start node (green) to the end node (red) without loops or backtracking.</p> <p>Example Output: [N16, N12, N8, N4, N3]</p>	 <p>Q: Given an image with exactly 5 objects, analyze their spatial relationships and find the shortest path that:</p> <ol style="list-style-type: none"> starts at the N1 object visits each object exactly once <p>Example Output: [N1, N3, N4, N5, N2]</p>	 <p>Q: Pick the correct option that matches the image. Options: A. A photo of a fire hydrant on the right, B. A photo of a fire hydrant on the left A: A</p>	 <p>Q: Pick the correct option that matches the image. Options: A. A dog under a table, B. A dog on a table, C. A dog to the left of a table, D. A dog to the right of a table A: D</p>

Compared to the baseline VLM model

Basic tasks

20% - 165%



Composite tasks

20% - 283%



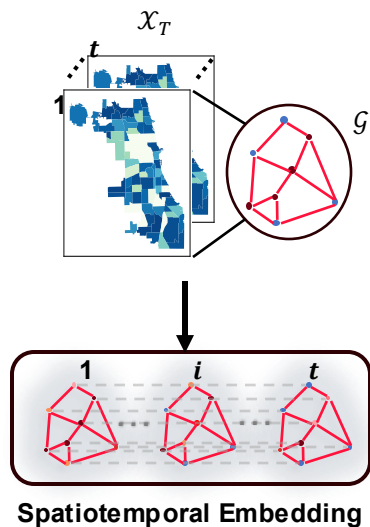
General tasks

5% - 25%

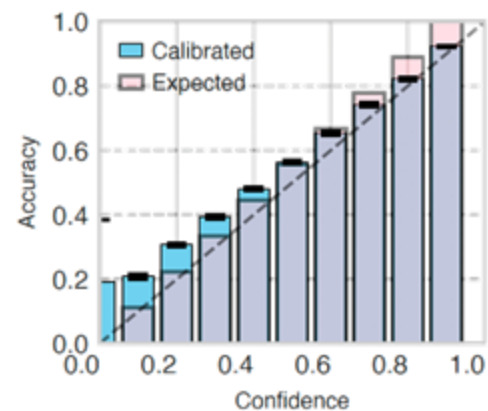


Wrap-up

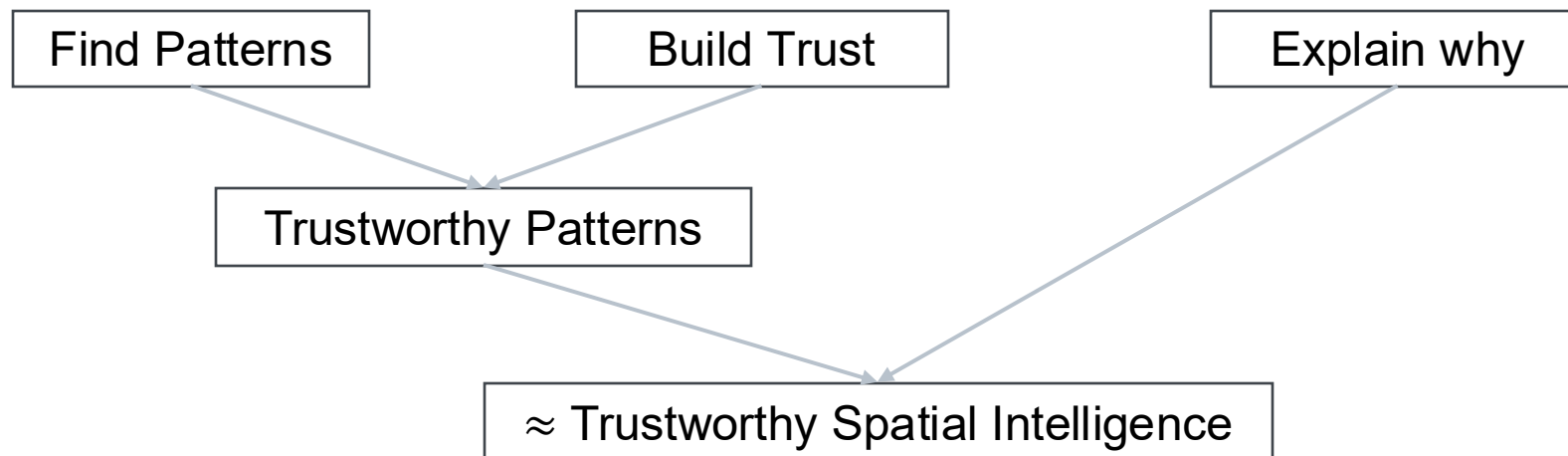
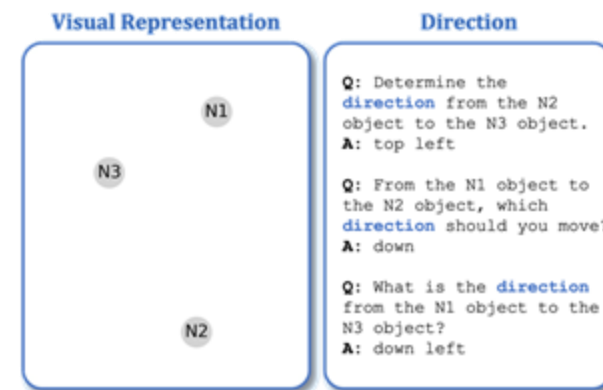
Learning



Calibrating



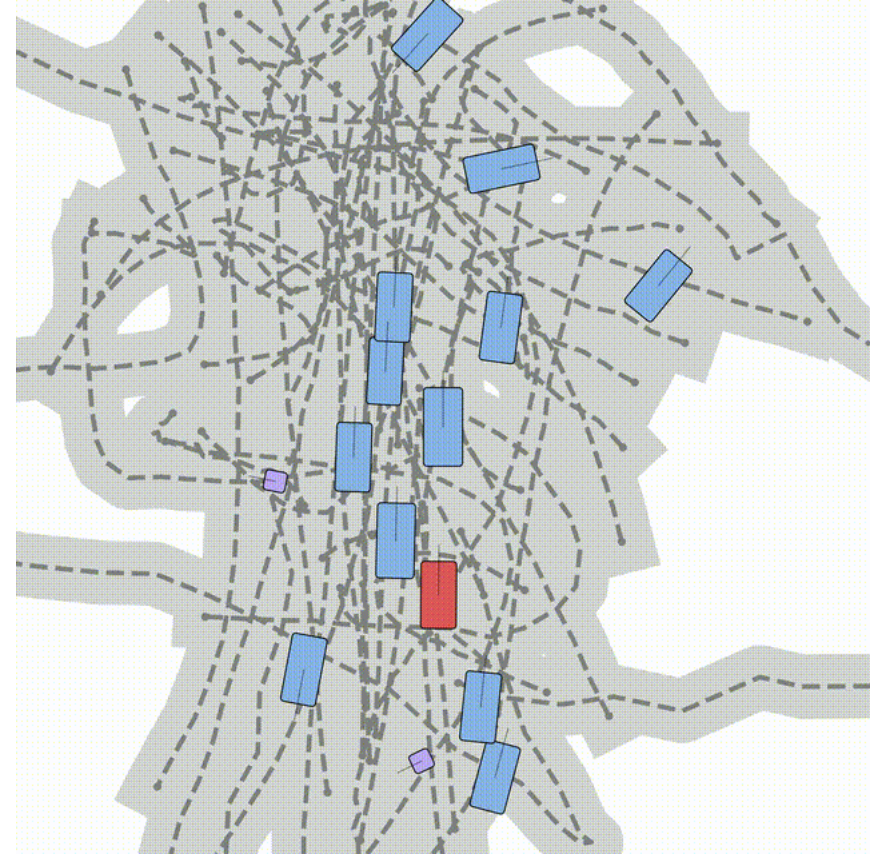
Reasoning



A Bigger Question

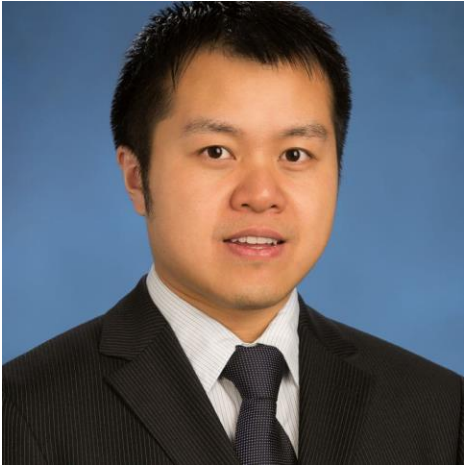
Can we build models for transportation and urban systems that don't just recognize patterns, or even just reason about them — but that actually internalize them, so they can simulate, plan, and ask 'what-if' questions?

Scenario Dreamer



What is “World Model”

David Ha



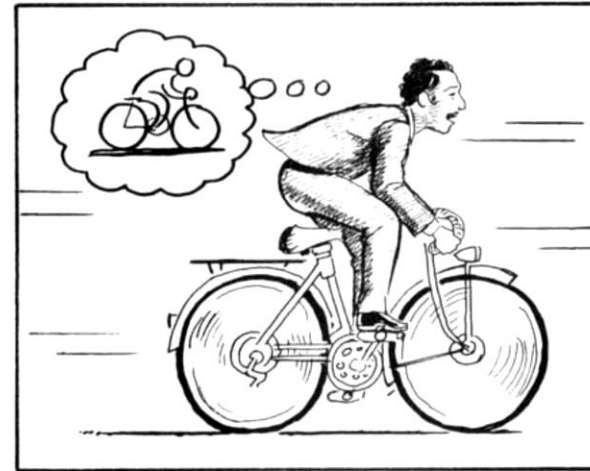
Jürgen Schmidhuber



Yann LeCun



Origin in **cognitive science & robotics**: internal model of the world for prediction & planning.



World models as the foundations of autonomous intelligence; self-supervised learning to internalize physics & dynamics

How Does My Work Connect

Indirectly, but address core ingredients



Contribution

Learning

Calibrating

Reasoning



Still Missing

Richer dynamics:
integrated environment simulator

Calibration under decision-making:
Online, real-time planning with feedbacks

Higer-level reasoning:
Counterfactual simulation, embodied reasoning in urban system

The Road Ahead: World Models of Cities



Transportation data
synthesis & augmentation



Auto-labeling



Multi-modality



Dynamic simulation



Embodied Reasoning

Thank you!

Questions?

Publication Reference

1. Wu, Y., Zhuang, D., Labbe, A., & Sun, L. (2021, May). *Inductive graph neural networks for spatiotemporal kriging*. In *Proceedings of the AAAI conference on artificial intelligence* (Vol. 35, No. 5, pp. 4478-4485).
2. Cao, X., Zhuang, D., Zhao, J., & Wang, S. (2025). *Virtual Nodes Improve Long-term Traffic Prediction*. *arXiv preprint arXiv:2501.10048*.
3. Zhuang, D., Hao, S., Lee, D. H., & Jin, J. G. (2020). *From compound word to metropolitan station: Semantic similarity analysis using smart card data*. *Transportation Research Part C: Emerging Technologies*, 114, 322-337.
4. Zhuang, D., Bu, Y., Wang, G., Wang, S., & Zhao, J. (2024, October). *Sauc: Sparsity-aware uncertainty calibration for spatiotemporal prediction with graph neural networks*. In *Proceedings of the 32nd ACM International Conference on Advances in Geographic Information Systems* (pp. 160-172).
5. Tang, Y., Wang, Z., Qu, A., Yan, Y., Wu, Z., Zhuang, D., ... & Ma, W. (2024). *ItiNera: Integrating spatial optimization with large language models for open-domain urban itinerary planning*. *arXiv preprint arXiv:2402.07204*.