

DSA Research Experiences for Undergraduates

Research Project

Section1: Faculty Information

Full Name	Yuxuan Liang	Tel	(020) 8833 5733
Thrust/Hub	System Hub / Thrust of Intelligent Transportation Info Hub / Thrust of Data Science and Analytics	Office	W2 L5 511
Email	yuxuanliang@hkust-gz.edu.cn		

Section2: Research Project Proposal

Project Title	Exploring Vision-Language Models for Spatio-Temporal Data Prediction
Project Description (max 800 words)	<p>1. Background</p> <p>Spatio-temporal prediction, such as traffic flow forecasting and weather modeling, involves capturing complex interactions between temporal dynamics and spatial structures. While existing methods have made significant progress, we propose a novel direction that leverages the power of Vision-Language Models (VLMs) to enhance spatio-temporal forecasting. Our approach transforms raw spatio-temporal data into image-like frames, enabling the use of pre-trained VLMs (e.g., CLIP, Flamingo) to integrate visual patterns (e.g., traffic heatmaps) and semantic context (e.g., event descriptions) for more accurate and interpretable predictions. By bridging the gap between spatio-temporal data and multimodal learning, this project aims to unlock new possibilities for urban dynamics modeling and beyond.</p> <p>2. Objectives</p> <p>This project aims to:</p> <ul style="list-style-type: none"> - Develop a spatio-temporal-to-visual mapping framework to convert raw spatio-temporal data (e.g., traffic grids) into image-like frames. - Design a multimodal fusion architecture integrating pre-trained VLMs for joint visual-textual reasoning. - Validate the framework on public datasets (e.g., PeMS traffic data) to achieve state-of-the-art prediction accuracy with interpretable outputs. <p>3. Methodology</p> <p>Step 1: Spatio-Temporal Visual Encoding</p> <ul style="list-style-type: none"> - Convert spatio-temporal tensors (time × height × width × channels) into multi-channel image sequences. (Example: Map traffic flow grids to RGB images where channels encode speed, density, and direction.)

	<ul style="list-style-type: none"> - Use video pretrained models (e.g., TimeSformer) to extract spatio-temporal features. <p>Step 2: Textual Context Injection</p> <ul style="list-style-type: none"> - Generate text prompts from metadata (e.g., "Weekday morning rush hour with light rain"). - Align visual features and text embeddings via CLIP-style contrastive learning. <p>Step 3: Multimodal Fusion</p> <ul style="list-style-type: none"> - Combine visual features, text embeddings, and raw temporal signals using gated cross-attention. - Train a lightweight predictor (e.g., MLP or 1D-CNN) to generate future frames. <p>4. Innovation</p> <ul style="list-style-type: none"> - First integration of VLMs for spatio-temporal prediction, leveraging both visual priors and semantic context. - Interpretable outputs: Visualize attention maps to explain how text prompts (e.g., "accident on Highway 101") influence predictions. <p>5. Expected Outcomes</p> <ul style="list-style-type: none"> - A lightweight library for spatio-temporal visual-textual data conversion. - A benchmark comparison against ST-GCN, ConvLSTM, and ST-Transformer on traffic datasets. - A research paper (targeting CIKM or KDD workshops) and open-source code.
Proposed Research Duration	<p>Start Date: _1_ / _June_ / 2025__</p> <p>End Date: __1__ / _December_ / __2025_</p>
Student/Researcher Duties	<p>Phase 1: Data Preprocessing & Visualization</p> <ul style="list-style-type: none"> - Convert raw spatio-temporal data (e.g., CSV files) into image frames using Python/OpenCV. - Annotate text prompts based on metadata (e.g., weather, holidays). <p>Phase 2: Model Implementation</p> <ul style="list-style-type: none"> - Fine-tune pre-trained VLMs (e.g., CLIP-ViT) on spatio-temporal image sequences. - Implement fusion modules (e.g., cross-attention) using PyTorch. <p>Phase 3: Evaluation & Visualization</p> <ul style="list-style-type: none"> - Compare prediction accuracy (MAE/RMSE) against baselines. - Generate Grad-CAM visualizations to interpret model decisions.

	<p>Phase 4: Paper writing & Deployment</p> <ul style="list-style-type: none"> - Write a paper following the standard format of conference. - Develop and deploy a visual interactive system.
<p>Technical Skills Required</p>	<input checked="" type="checkbox"/> Python <input checked="" type="checkbox"/> Machine Learning <input checked="" type="checkbox"/> Big Data <input type="checkbox"/> R <input checked="" type="checkbox"/> Deep Learning <input type="checkbox"/> SQL <input type="checkbox"/> C/C++ <input type="checkbox"/> Other: _____
<p>Preferred Student/Researcher Background</p>	<p>Academic: Undergraduate students in Computer Science, Data Science, or GIS-related fields.</p> <p>Experience:</p> <ul style="list-style-type: none"> - Coursework or projects in machine learning (especially CV/NLP). - Hands-on experience with PyTorch. <p>Bonus:</p> <ul style="list-style-type: none"> - Prior exposure to spatio-temporal prediction (e.g., traffic forecasting). - Familiarity with linux system for large-scale training.
<p>Maximum Number of Students/Researchers</p>	<input type="checkbox"/> 1 <input checked="" type="checkbox"/> 2

Section3: Pre-Application Research Exposure Meeting

Faculty members are encouraged to schedule a Research Exposure Meeting to introduce students to their projects.

Preferred Date	Each Friday
Preferred Time	14:00-15:00
Meeting Mode	<input checked="" type="checkbox"/> In-Person <input type="checkbox"/> Online
Venue (if in-person)	W2-511
Meeting Link (if online)	