# Learning Multimodal Embeddings for Traffic Accident Prediction and Causal Estimation

### Ziniu Zhang
zhang.zini@northeastern.edu
Northeastern University
Boston, Massachusetts

### Minxuan Duan
duan.mi@northeastern.edu
Northeastern University
Boston, Massachusetts

### Haris N. Koutsopoulos
h.koutsopoulos@northeastern.edu
Northeastern University
Boston, Massachusetts

### Hongyang R. Zhang
ho.zhang@northeastern.edu
Northeastern University
Boston, Massachusetts

## Abstract

We consider analyzing traffic accident patterns using both road network data and satellite images aligned to road graph nodes. Previous work for predicting accident occurrences relies primarily on road network structural features while overlooking physical and environmental information from the road surface and its surroundings. In this work, we construct a large multimodal dataset across six U.S. states, containing nine million traffic accident records from official sources, and one million high-resolution satellite images for each node of the road network. Additionally, every node is annotated with features such as the region's weather statistics and road type (e.g., residential vs. motorway), and each edge is annotated with traffic volume information (i.e., Average Annual Daily Traffic). Utilizing this dataset, we conduct a comprehensive evaluation of multimodal learning methods that integrate both visual and network embeddings. Our findings show that integrating both data modalities improves prediction accuracy, achieving an average AUROC of 90.1%, which is a 3.7% gain over graph neural network models that only utilize graph structures. With the improved embeddings, we conduct a causal analysis based on a matching estimator to estimate the key contributing factors influencing traffic accidents. We find that accident rates rise by 24% under higher precipitation, by 22% on higher-speed roads such as motorways, and by 29% due to seasonal patterns, after adjusting for other confounding factors. Ablation studies confirm that satellite imagery features are essential for achieving accurate prediction. We release the dataset and the experimental code for using this dataset at https://github.com/VirtuosoResearch/MMTraCE.

## CCS Concepts

• **Computing methodologies → Neural networks**; **Learning latent representations**; **Causal reasoning and diagnostics**.

## Keywords

Multimodal Learning, Satellite Imagery, Road Safety, Causal Inference and Analysis

## 1 Introduction

Road safety remains a persistent and pressing challenge in urban cities and districts. According to the World Health Organization, approximately 1.19 million people die in road traffic crashes each year, making it the leading cause of death for children and young adults [24]. In the United States alone, traffic accidents are projected to cause 39,345 deaths in 2024 and $871 billion in annual societal costs [2, 22]. Faced with such a severe reality, it is imperative to shift from reactive analysis to proactive prediction, using data-driven methods to identify high-risk road areas and implement interventions for accident prevention.

Current data-driven models, especially those based on graph neural networks (GNNs), draw on road network topology, traffic flow statistics, and historical accident data [19, 23, 38, 44]. These data sources have improved the modeling of accident patterns by capturing spatial relationships [45]. However, structural features fail to reflect the visual and environmental conditions of the physical road. Factors such as lane width, curvature, surface quality, and nearby land use shape driver behavior and accident likelihood, yet they do not appear in standard network data.

Advances in remote sensing and computer vision provide a path to fill this gap. These tools support a road-level view that links physical scenes to accident patterns. Recent work in transportation and climate modeling has shown clear gains from hybrid models that combine high-dimensional visual inputs with traditional numeric features [21, 32, 34]. Satellite images offer complementary information to network-based representations and make visible the physical road attributes that influence accident risk. Examples are shown in Figure 1 (see also Figure 9 in the Appendix). These attributes are hard to extract from conventional data sources but are closely linked to accident risk [7].

Despite this potential, the use of satellite imagery in state-level traffic accident analysis remains limited. There are two main challenges. The first challenge is the lack of large public multimodal datasets with enough temporal coverage across different regions [10, 33, 37]. This gap comes from the heterogeneity of data sources.

**Figure 1: Example satellite images showing different types of roads. Each image is centered around a road network node and captures both the physical characteristics of the road, such as layout, width, and intersections, and the surrounding context, including vegetation, buildings, and terrain.**

Each state publishes accident and traffic records in its own format. They also use different variable names and hosting structures. Unifying these datasets requires careful preprocessing and normalization. This includes schema alignment and geospatial matching. The second challenge is the need to reduce the dimensionality of high-resolution satellite images. These images must also be integrated with graph-structural features and other low-dimensional inputs. Examples include weather, traffic volume, and speed limits. Early experiments show that simple fusion methods do not work well. Direct concatenation fails to capture nonlinear interactions between different modalities. As a result, models struggle to use both visual and non-visual information in an effective way.

To address these challenges, we begin by constructing a large-scale, multimodal dataset spanning six U.S. states. It comprises over nine million accident records and one million satellite images with rich feature annotations.

- Each road network node is aligned with high-resolution satellite images and supplemented with road-level features, including weather conditions (e.g., temperature, precipitation, wind speed, and atmospheric pressure). Traffic volume indicators of each edge, such as Average Annual Daily Traffic (AADT), are also associated with each node when available.
- Each satellite image has a resolution of $1024 \times 1024$ pixels and covers approximately $200\,\text{m} \times 200\,\text{m}$ of physical space.
- Every road is also associated with historical accident occurrences. Notably, we collect the latest accident data up to the year 2025, with a maximum period covering a temporal span of up to 24 years, allowing models to capture long-term trends and evolving risk patterns. This dense spatial alignment enables models to jointly learn from both visual cues and graph-structural features.

Compared to existing datasets [9, 23, 40], our dataset covers a much wider geographic range. It also provides detailed annotations of

satellite images. These features make it well-suited for large-scale multimodal traffic accident analysis.

Building on this dataset, we develop a multimodal learning framework. The framework combines vision-based encoders with GNNs and supports joint reasoning over visual, structural, and contextual properties of the road environment. We conduct a detailed evaluation of different fusion strategies to measure their impact on traffic accident prediction. By adding satellite imagery to the model, we obtain an average AUROC of 90.1% across six states. The prediction accuracy improves by 3.7% compared to GNNs that use only graph-structural features. Based on this result, we further perform a causal analysis of key environmental and traffic factors. We estimate the average treatment effect on the treated group (ATT) using a matching estimator built on the multimodal embeddings. The ATT scores for seasonal variation, road type, and precipitation are 28.6%, 21.9%, and 24.2%. We further estimate the effects with propensity score-based adjustments (PSM) and the double robust estimator (DR). We run a leave-one-out ablation study to measure the contribution of each feature type. Removing image features leads to a drop of 3.5% AUROC. Excluding weather, traffic volume, and road network features results in decreases of 1.8%, 2.4%, and 3.7%, respectively.

In summary, this paper makes three contributions to the study of traffic accidents on road networks. First, we build a large multimodal dataset with recent traffic accident records from six U.S. states. Each satellite image is aligned with a specific road network node. To our knowledge, this is the most comprehensive dataset for evaluating multimodal learning in the transportation domain. Second, we provide a detailed evaluation of multimodal learning methods that combine visual embeddings with network embeddings. The results show clear gains in accident prediction accuracy. Third, we perform a causal analysis to measure the effect of key factors that influence accident occurrences. We use the learned embeddings to control for confounding variables such as road conditions. We hope that this dataset can support future research on multimodal learning for transportation and road safety.

## 2 Methodology

We now describe our multimodal dataset and the analytical tools along with the dataset. First, we provide an overview of the problem setting. We then introduce the collection of our dataset, including the Satellite images. Next, we describe the multimodal approaches that integrate both network embeddings and visual embeddings from the images. Finally, we design a causal analysis framework on top of the multimodal embeddings.

### 2.1 Preliminaries

We study the problem of accident prediction and analysis on road networks. The road network is modeled as a directed graph $G = (V, E)$, where each node $v \in V$ represents a road intersection or critical point along a road segment, and each directed edge $e \in E$ represents a road connecting two such points. Our goal is to predict the probability or frequency of accidents occurring at each edge, leveraging both spatial structure and temporal dynamics.

Each node is associated with dynamic features that evolve over time, such as weather conditions, including temperature, wind

**Table 1: Statistics of road network graphs for each state in our dataset. We report the total number of edges ($m$), average edge length in meters, road network density ($m / \binom{n}{2}$, where $n$ is the number of nodes), which is the number of edges per unit area, and availability of traffic volume, which indicates the proportion of road segments that are associated with available traffic volume measurements.**

|  | # Edges | Avg Length (m) | Density | Volume (%) |
|---|---|---|---|---|
| Delaware | 116, 196 | 213.0 | $9.7 \times 10^{-5}$ | 3.7 |
| Massachusetts | 706, 402 | 188.8 | $1.7 \times 10^{-5}$ | 0.9 |
| Maryland | 580, 526 | 211.7 | $1.8 \times 10^{-5}$ | 2.1 |
| Nevada | 292, 674 | 280.1 | $4.0 \times 10^{-5}$ | 1.3 |
| Montana | 351, 516 | 859.2 | $3.3 \times 10^{-5}$ | 0.8 |
| Iowa | 707, 072 | 532.4 | $2.2 \times 10^{-5}$ | – |

speed, precipitation, and sea-level pressure. Additionally, each node is linked to a high-resolution satellite image centered at its geographic coordinates. We extract visual features that capture the road structure and the surrounding environment from the image.

Each edge is associated with static and dynamic attributes that characterize the corresponding road segment. These include road length, road type (e.g., residential, motorway), and average annual daily traffic (AADT).

Formally, at each time step $t$, we observe a dynamic graph $G_t = (V, E, X_t)$, where $X_t$ denotes the time-dependent attributes. Given a sequence of previous observations $\{G_{t-T+1}, \ldots, G_t\}$, where $T$ is the observation period, the objective is to predict accident risk at each node for the next time step:

$$\hat{y}_v^{(t+1)} = f(G_{t-T+1}, \ldots, G_t; \mathcal{I}_v),$$

where $\mathcal{I}_v$ denotes the visual embedding derived from the satellite image corresponding to node $v$. The function $f(\cdot)$ is parameterized by a multimodal spatiotemporal model that jointly reasons over the graph structure, node dynamics, edge dynamics, and image content.
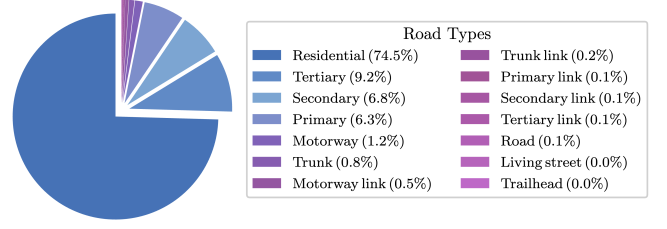
## 2.2 Dataset Collection

The dataset is assembled in several stages. We begin by constructing road network graphs for six states, followed by the collection of high-resolution satellite images that are spatially aligned with the road graph. Finally, we collect historical traffic accident records from official sources and align them with both the road network and satellite images.

**Road network.** For each state, we construct a detailed road network graph based on data obtained from OpenStreetMap (OSM).[1] The road network encompasses five major categories of roadways: city streets, county roads, neighborhood streets, tract roads, and urbanized area roads, ensuring comprehensive spatial coverage across diverse geographic and administrative regions.

Each road segment in the network is represented as a directed edge defined by a start node and an end node, and is annotated with

---

[1]Notice that the original data from OSM divides one edge into several parts. Instead, we have preprocessed the raw data to combine the parts into a single edge. The volume statistic is low due to the large number of residential roads in the data. See Figure 2 for an illustration of the distribution of different types of roads.



Figure 2: The proportion of different road types among six states' road networks. Residential roads account for the vast majority of the total, making up approximately 74.5% of all roads. Other types, such as tertiary, secondary, and primary, contribute much smaller proportions by comparison.

metadata including the road name, whether the segment is one-way, which road type it is, and the physical length of the segment in meters. This representation captures both the structural and functional aspects of the road infrastructure.

Each node in the graph corresponds to a specific geospatial point, uniquely identified by a node_id. It is associated with precise latitude and longitude coordinates. These nodes typically represent road intersections or endpoints, serving as key units for spatial reasoning and image alignment. See Table 1 for an overview of the network statistics.

There are 14 road types in total: living street, motorway, motorway link, primary, primary link, residential, road, secondary, secondary link, tertiary, tertiary link, trailhead, trunk, and trunk link. We report the probability of different road types in Figure 2. This processed road graph serves as the structural backbone for integrating visual and contextual data in our framework.

**Satellite image data.** Building upon the extracted road network, we obtain high-resolution satellite images centered at each node, each of which typically corresponds to a road intersection. The geographic coordinates of these nodes guide the image collection process, ensuring tight spatial alignment between the visual and structural representations. For each state, we collect hundreds of thousands of satellite images, resulting in a large-scale visual dataset. We implement two strategies for image acquisition: (1) directly fetching static images using the Mapbox Static API, and (2) reconstructing large images by stitching together 25 smaller tiles. Each final image has a resolution of $1024 \times 1024$ pixels and covers approximately $200 \, \text{m} \times 200 \, \text{m}$ of physical space, capturing rich visual cues and fine-grained road features. In summary, for each network node in the road network, we align one satellite image to that node, for all six states. Therefore, the number of satellite images is equal to the number of nodes.

**Accident records.** As accident occurrence serves as the primary prediction objective in our framework, we obtain accident records from official sources, specifically the Departments of Transportation (DOT) of each U.S. state, to ensure high reliability and real-world relevance. We collected accident records spanning up to the past two decades, with states providing data ranging from the early 2000s to as recent as 2025. However, these records vary significantly across states in terms of file format, schema, field names, spatial encoding, and level of detail. To enable unified modeling and large-scale

**Table 2: Overview of the collected dataset, including accident counts, period of accident records, and aligned satellite imagery across 6 U.S. states. The start and end years of the accident records are based on the latest released data from the Department of Transportation.**

|               | Start | End  | # Labels  | # Satellite Images |
|---------------|-------|------|-----------|--------------------|
| Delaware      | 2009  | 2024 | 533,112   | 49,023             |
| Massachusetts | 2002  | 2025 | 5,165,834 | 285,942            |
| Maryland      | 2015  | 2023 | 997,532   | 250,565            |
| Nevada        | 2016  | 2025 | 376,252   | 121,392            |
| Montana       | 2016  | 2023 | 140,011   | 145,525            |
| Iowa          | 2014  | 2024 | 594,492   | 253,623            |

analysis, we perform extensive preprocessing and normalization to convert all state-specific datasets into a standardized format. Further, we align their axes with the satellite images.

The overall statistics of our dataset are shown in Table 2, including the period of accidents, total number of accidents, and total number of satellite images. More details about the dataset collection process can be found in Appendix A.

**Remark 2.1.** Compared with the most recent dataset on traffic accidents [23], our dataset now includes over one million satellite images. Additionally, we align the axes of the satellite images with the road networks to perform analysis with both data modalities.

## 2.3 Learning Multimodal Embeddings

We use graph neural networks (GNNs) to learn network embeddings. Consider a graph $G = (V, E)$, annotated with both node and edge features. GNNs learn node representations by iteratively aggregating information from local neighborhoods. Let $x_i^{(k)}$ denote the representation of node $i$ at layer $k$, and $v_{i,j}$ the feature of edge $(i, j)$. The update rule of a message-passing layer is given by:

$$x_i^{(k+1)} = \phi\left(x_i^{(k)}, h\left(\left\{\psi(x_i^{(k)}, x_j^{(k)}, v_{i,j}) \mid j \in N(i)\right\}\right)\right),$$

where $N(i)$ is the set of neighbors of $i$, $h$ is a permutation-invariant aggregator (e.g., sum, mean), and $\phi, \psi$ are neural networks with nonlinearities. This iterative procedure enables GNNs to capture neighborhood connections from the network.

Next, we use a vision model such as Vision Transformer (ViT) to extract semantic information from the images and encode each image into a visual embedding.

Given the GNN and image embeddings, we develop a multimodal learning framework to combine network embeddings and visual features. To incorporate visual semantics into the graph-based modeling process, we extract image features corresponding to each node and combine them with node embeddings generated by a graph neural network. We develop three combination methods to handle different types of features in traffic accident prediction.

- Basic fusion: we build a multilayer perceptron (MLP) to combine both embeddings.
- Gated fusion: we train a scalar gate value to determine the relative importance of the two modalities and then derive a weighted combination of both of them.

---

**Algorithm 1** Mixture-of-Experts Fusion

**Input:** Visual embedding $z$, graph embedding $x^{(\text{GNN})}$, number of experts $K$, and edge embedding $x_{\text{edge}}$
**Require:** Expert networks $\{f_k\}_{k=1}^K$, gating network $f_{\text{gate}}$, prediction head $f_{\text{pred}}$
**Output:** Prediction $\hat{y}$
1: Concatenate visual and graph features: $\tilde{x} \leftarrow \left[x^{(\text{GNN})} \parallel z\right]$
2: **for** $k = 1$ to $K$ **do**
3:     Compute expert output: $e^{(k)} \leftarrow f_k(\tilde{x})$
4: **end for**
5: Compute gating weights: $\lambda \leftarrow \text{softmax}\left(f_{\text{gate}}(\tilde{x})\right)$
6: Fuse expert outputs: $\tilde{e} \leftarrow \sum_{k=1}^K \lambda^{(k)} \cdot e^{(k)}$
7: Predict output: $\hat{y} \leftarrow f_{\text{pred}}(\tilde{e}, x_{\text{edge}})$
8: **return** $\hat{y}$

---

- Mixture of experts: we leverage multiple specialized networks, referred to as experts, that each learn to combine features from different perspectives. We utilize the gated network that computes the probability distribution over experts for each node.

See Algorithm 1 for the complete procedure of the implementation of the mixture of experts. Further details of the fusion methods can be found in Appendix B.1.

## 2.4 Causal Estimation Using Multimodal Embeddings

To provide a more fine-grained and context-aware estimation of causal effects, we perform matching-based causal analysis in the learned multimodal embedding space. Let $x_i \in \mathbb{R}^d$ denote the multimodal embedding of edge $i$, incorporating both visual and structural information. We define a binary treatment indicator $t_i = \mathbb{1}$ for each edge $i$. The corresponding outcome variable $y_i$ represents the accident count or occurrence.

Our goal is to estimate the average treatment effect on the treated group (ATT). Consider the treated set $T = \{i \mid t_i = 1\}$, and a given treated edge $i \in T$, we compute the expected influence by

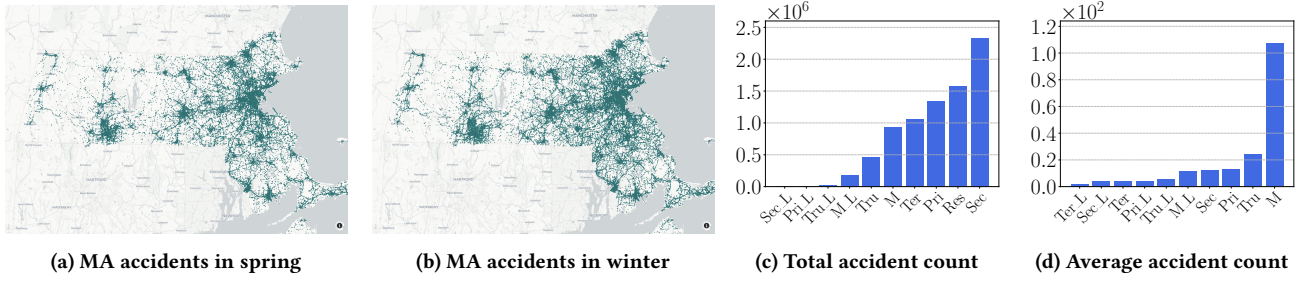$$\tau = \mathbb{E}\left[\hat{y}_{i,1} - \hat{y}_{i,0}\right],$$

where $\hat{y}_{i,1}$ and $\hat{y}_{i,0}$ denote the potential outcomes under treatment and control. Since only one of these outcomes is observed, we approximate the missing counterfactuals through nearest neighbor matching in the embedding space. For each treated sample, we identify its nearest neighbor $j$ from the control set,

$$j = \arg\min_k \|x_i - x_k\|_2,$$

and use its outcome $y_j$ as an estimate of $\hat{y}_{i,0}$. The ATT is then computed by

$$\hat{\tau} = \frac{1}{|T|} \sum_{i \in T} (y_i - y_j).$$

This embedding-based approach allows us to control for confounding factors encoded in the multimodal representation. It lets us compare treatment effects under similar road networks and weather conditions. To strengthen this analysis, we also estimate treatment effects with propensity score matching (PSM), which aligns treated and control edges with similar treatment probabilities. In addition,

(a) MA accidents in spring    (b) MA accidents in winter    (c) Total accident count    (d) Average accident count

**Figure 3: Seasonal comparison of traffic accidents in Massachusetts. 3a, 3b: Accident records in Massachusetts during spring and winter. It is evident that accident points are more densely distributed in winter, indicating a higher frequency of incidents likely due to adverse weather conditions. 3c, 3d: Accident count of motorway (M), motorway link (M_L), primary (Pri), primary link (Pri_L), residential (Res), secondary (Sec), secondary link (Sec_L), tertiary (Ter), tertiary link (Ter_L), trunk (Tru), trunk link (Tru_L), living street, road, and trailhead. Figure 3c gives the top-10 total count on different road types, while Figure 3d provides the top-10 average count on different road types.**

we use a doubly robust (DR) estimator that combines outcome prediction with propensity weighting to provide a more stable estimate. Full details of both procedures are reported in Appendix B.

## 3 Experiments

Given the constructed multimodal dataset and the proposed fusion methods, we evaluate the performance of various graph learning and multimodal strategies for predicting traffic accidents. We aim to assess the contribution of satellite imagery and the impact of different contextual factors on the performance. We begin by detailing the experimental setup, including the baselines and the evaluation metrics. We then report the prediction results across different baselines and fusion strategies. Lastly, we conduct a causal estimation analysis to investigate the influence of key structural and environmental factors on accident occurrences.

### 3.1 Experimental Setup

*3.1.1 Baselines.* Recall that we focus on an edge-level link prediction task. As baselines, we include embedding methods, graph neural networks, and feature fusion approaches.

First, we evaluate multilayer perceptrons (MLPs) using node features such as node degrees, betweenness centrality, road type, weather data, road length, and traffic volume. This setup assesses node features without incorporating the underlying network structure. We test embedding methods, including DeepWalk [25], CLIP, and Vision Transformer, by appending a layer that concatenates the learned node embeddings with the original node features.

Second, we utilize common GNNs architectures including Graph-SAGE, GCN, GIN, and Graph Transformer (Graphormer) with attention [31]. We use a spatial-temporal framework such as DCRNN [19]. We also consider supervised contrastive learning to improve the classification of positive and negative edges.

Then, we consider three approaches to integrate graph-based and vision-based features. The first approach basically employs a three-layer multilayer perceptron that jointly processes the concatenated graph and image representations. The second approach uses a gated fusion network that adaptively learns the contribution weights

of each modality during training. The last one utilizes the MoE architecture to fetch features from different perspectives.

*3.1.2 Implementations.* For each state, we partition the available accident records into training, validation, and test sets based on temporal splits. Specifically, historical data up to a designated cutoff year is used for model training, while accidents occurring after that year are reserved for evaluation. We focus on monthly accident prediction. This avoids future information leakage. We consider both classification and regression tasks. For classification, we report the Area Under the Receiver Operating Characteristic curve (AUROC). For regression, we use the Mean Absolute Error (MAE) to measure the difference between the predicted and observed accident counts on each road segment. We provide additional implementation details, including concrete hyperparameters, formal fusion method definitions, and formal metric definitions in Appendix B.1. The experimental details, including the omitted regression and classification task prediction result, the precision score, and the recall score, are shown in Appendix B.2.

### 3.2 Traffic Accident Prediction Results

We illustrate the experimental results on multiple baselines shown in Section 3.1. Below, we summarize two key insights that we believe hold broader relevance beyond the specific scope of our study.

We evaluate all models on six states and compare variants with and without visual inputs. GCN-based models show clear gains from incorporating satellite imagery. In particular, GCN plus Gated Fusion improves average AUROC by 3.8% over the standard GCN, with a maximum gain of 7.2% across states. GCN plus MoE yields a similar pattern, improving AUROC by an average of 3.9% and up to 7.7% in the best case.

For GIN-based models, visual features also provide consistent benefits. GIN plus Gated Fusion improves over the vanilla GIN by an average of 2.3% AUROC, with a maximum gain of 6.7%. GIN plus MoE improves over the base GIN by an average of 3.6%, with the largest gain of 7.9%. Taken together, these results show that satellite imagery improves performance for all architectures, and that the relative gains are larger for models with lower base expressivity.

**Table 3: Comparing the results of GNNs, vision models, and multimodal fusion strategies. The performance is evaluated using the mean absolute error (MAE) and area under the ROC curve (AUROC) on the test split. A leave-one-out analysis is also provided. To account for variability, each experiment is repeated with three different random seeds, and we report the average results along with their standard deviations.**

| Category | MAE (↓)<br>Average count | Delaware<br>4.59 | Massachusetts<br>7.25 | Maryland<br>1.72 | Nevada<br>1.29 | Montana<br>0.40 | Iowa<br>0.84 |
|---|---|---|---|---|---|---|---|
| GNNs | MLP | 0.5 ± 0.07 | 0.6 ± 0.10 | 0.2 ± 0.08 | 0.3 ± 0.03 | 0.2 ± 0.04 | 0.3 ± 0.09 |
| | DeepWalk | 0.2 ± 0.02 | 0.6 ± 0.06 | 0.3 ± 0.02 | 0.1 ± 0.02 | 0.3 ± 0.01 | 0.2 ± 0.04 |
| | GraphSAGE | 0.1 ± 0.01 | 0.2 ± 0.01 | 0.2 ± 0.00 | 0.1 ± 0.02 | 0.2 ± 0.01 | 0.1 ± 0.02 |
| | Graph Transformer | 0.1 ± 0.03 | 0.2 ± 0.02 | 0.2 ± 0.01 | 0.2 ± 0.01 | 0.2 ± 0.04 | 0.1 ± 0.03 |
| | DCRNN | 0.2 ± 0.02 | 0.3 ± 0.01 | 0.1 ± 0.03 | 0.1 ± 0.03 | 0.2 ± 0.00 | 0.1 ± 0.03 |
| | GCN | 0.1 ± 0.02 | 0.8 ± 0.20 | 0.3 ± 0.01 | 0.2 ± 0.01 | 0.2 ± 0.03 | 0.2 ± 0.01 |
| | SupConGCN | 0.2 ± 0.06 | 0.3 ± 0.03 | 0.3 ± 0.03 | 0.2 ± 0.02 | 0.2 ± 0.03 | 0.2 ± 0.02 |
| | GIN | 0.1 ± 0.02 | 0.5 ± 0.05 | 0.3 ± 0.01 | 0.1 ± 0.02 | 0.1 ± 0.03 | 0.2 ± 0.01 |
| Vision models | CLIP | 0.3 ± 0.11 | 0.5 ± 0.11 | 0.3 ± 0.02 | 0.3 ± 0.03 | 0.1 ± 0.03 | 0.2 ± 0.01 |
| | Vision Transformer | 0.2 ± 0.04 | 0.6 ± 0.03 | 0.2 ± 0.02 | 0.2 ± 0.04 | 0.3 ± 0.01 | 0.2 ± 0.01 |
| Multimodal fusion | GIN + Basic Fusion | 0.1 ± 0.00 | 0.3 ± 0.01 | 0.2 ± 0.01 | 0.2 ± 0.01 | 0.2 ± 0.03 | 0.2 ± 0.01 |
| | GIN + Gated Fusion | 0.1 ± 0.01 | 0.3 ± 0.03 | 0.2 ± 0.01 | 0.1 ± 0.00 | 0.3 ± 0.04 | 0.2 ± 0.01 |
| | GIN + MoE | **0.1 ± 0.00** | **0.3 ± 0.02** | **0.2 ± 0.01** | **0.1 ± 0.00** | **0.2 ± 0.01** | **0.1 ± 0.02** |

| Category | AUROC (↑)<br>Positive rate | Delaware<br>0.32 | Massachusetts<br>0.28 | Maryland<br>0.24 | Nevada<br>0.14 | Montana<br>0.10 | Iowa<br>0.18 |
|---|---|---|---|---|---|---|---|
| GNNs | MLP | 78.5 ± 0.1 | 64.2 ± 0.1 | 60.7 ± 0.9 | 80.8 ± 0.3 | 60.2 ± 1.1 | 66.7 ± 0.4 |
| | DeepWalk | 82.1 ± 0.4 | 82.3 ± 0.6 | 86.3 ± 0.2 | 89.5 ± 0.2 | 78.1 ± 1.5 | 77.3 ± 1.6 |
| | GraphSAGE | 79.5 ± 0.0 | 80.3 ± 0.6 | 78.7 ± 0.0 | 81.8 ± 0.7 | 78.1 ± 1.6 | 77.4 ± 0.2 |
| | Graph Transformer | 85.8 ± 1.5 | 81.0 ± 0.6 | 77.9 ± 3.0 | 88.7 ± 0.9 | 79.2 ± 2.5 | 80.1 ± 1.0 |
| | DCRNN | 91.6 ± 2.5 | 84.3 ± 0.1 | 86.1 ± 0.1 | 85.5 ± 0.7 | 81.3 ± 0.2 | 80.4 ± 1.6 |
| | GCN | 87.3 ± 0.8 | 86.2 ± 0.2 | 85.6 ± 0.2 | 89.4 ± 0.2 | 78.1 ± 1.4 | 77.2 ± 1.6 |
| | SupConGCN | 86.4 ± 1.8 | 79.8 ± 1.3 | 83.8 ± 1.2 | 90.8 ± 0.1 | 83.5 ± 0.6 | 83.2 ± 0.3 |
| | GIN | 91.6 ± 0.7 | 85.7 ± 0.3 | 85.9 ± 0.3 | 93.5 ± 0.2 | 79.6 ± 0.6 | 85.3 ± 0.9 |
| Vision models | CLIP | 86.3 ± 0.5 | 82.6 ± 0.0 | 85.6 ± 0.1 | 91.9 ± 0.0 | 79.6 ± 0.8 | 83.9 ± 0.0 |
| | Vision Transformer | 89.4 ± 0.9 | 82.1 ± 0.5 | 86.4 ± 0.0 | 93.2 ± 0.1 | 80.9 ± 0.1 | 85.9 ± 0.2 |
| Multimodal fusion | GIN + Basic Fusion | 92.3 ± 1.9 | 84.8 ± 0.3 | 88.1 ± 0.2 | 93.5 ± 0.2 | 79.0 ± 3.5 | 86.7 ± 0.4 |
| | GIN + Gated Fusion | 92.8 ± 2.8 | 85.8 ± 0.2 | 88.5 ± 0.3 | 94.1 ± 0.2 | 86.3 ± 1.1 | 87.7 ± 0.2 |
| | GIN + MoE | **96.4 ± 1.0** | **88.1 ± 0.4** | **88.7 ± 0.1** | **94.5 ± 0.2** | **87.5 ± 1.6** | **88.0 ± 0.4** |
| LOO analysis | GCN + Gated Fusion | 88.5 ± 0.2 | 87.7 ± 0.1 | 89.3 ± 0.0 | 93.5 ± 0.2 | 83.7 ± 0.7 | 84.4 ± 0.5 |
| | w/o visual features | 87.3 ± 0.8 | 86.2 ± 0.2 | 87.7 ± 0.3 | 89.4 ± 0.2 | 78.1 ± 1.4 | 77.2 ± 1.6 |
| | w/o weather features | 82.1 ± 1.7 | 86.1 ± 0.2 | 88.9 ± 0.3 | 92.9 ± 0.3 | 82.3 ± 0.2 | 83.8 ± 0.3 |
| | w/o road network features | 77.5 ± 2.4 | 82.3 ± 1.0 | 85.9 ± 0.3 | 93.2 ± 0.3 | 81.9 ± 1.6 | 83.9 ± 0.5 |
| | w/o traffic volume features | 85.4 ± 1.5 | 84.2 ± 0.4 | 86.9 ± 0.4 | 93.1 ± 0.1 | 80.9 ± 2.7 | 82.0 ± 1.2 |
| | w/ speed limit features | 89.5 ± 1.4 | 88.6 ± 0.2 | 89.7 ± 0.3 | 94.2 ± 0.1 | 84.2 ± 0.2 | 84.7 ± 0.2 |

**Ablation studies.** We conduct two ablation studies to understand the behavior of our frameworks. First, we evaluate the contribution of different feature types. Then, we examine the sensitivity to key hyperparameters.

For the feature ablation, we perform a leave-one-out (LOO) analysis on GCN plus Gated Fusion. We consider four feature categories: vision structure, weather observations, road layout, and traffic volumes. In each trial, we remove one category and measure the resulting change in AUROC. The results in Table 3 show that dropping image features leads to a 3.5% decrease. This confirms their strong predictive value. Removing weather features reduces AUROC by 1.8%. Excluding traffic volume produces a 2.4% drop. Eliminating road network features causes the largest decrease of 3.7%. This highlights the central role of road layout. Adding speed limit as an extra feature yields a small but consistent improvement suggesting that it provides complementary information.

For the hyperparameter study, we vary the number of GCN layers and the number of training epochs on the Nevada subset.

Adjusting the depth from 2 to 10 layers changes AUROC by no more than 3%. This shows that the model remains stable across a wide range of depths. Varying the number of training epochs from 15 to 50 shows that performance plateaus around epoch 30, suggesting that extended training offers limited benefit.

**Cross-state transfer.** To evaluate how well the model transfers to new regions, we conduct a cross-state generation experiment. In this setting, we train the GIN plus MoE model on one state and test it on another. The AUROC scores for all train–test pairs are shown in Figure 4. The results reveal strong cross-state consistency. We find that the accident patterns in Maryland and Nevada are more consistent with the shared representations learned by the model. It also indicates that both structural and visual features in these two states align well with those from other regions. Further understanding these transfer patterns is an interesting question for future work [14, 35, 36].

**Runtime.** The fusion structure may introduce extra computation, so we remove one GNN layer in the fusion modules to keep the

|     | DE   | MA   | MD   | NV   | MT   | IA   |
|-----|------|------|------|------|------|------|
| DE  |      | 81.0 | 85.4 | 90.5 | 78.6 | 80.0 |
| MA  | 82.4 |      | 87.8 | 91.8 | 85.2 | 84.5 |
| MD  | 82.3 | 84.8 |      | 91.9 | 85.3 | 84.5 |
| NV  | 82.1 | 83.8 | 87.3 |      | 84.6 | 80.8 |
| MT  | 82.5 | 84.8 | 88.0 | 91.9 |      | 84.5 |
| IA  | 82.8 | 84.3 | 87.9 | 92.0 | 86.9 |      |

**Figure 4: Cross-state AUROC performance of the GIN + MoE model, computed over six states. Each entry shows the score when training on one state (represented by rows) and testing on another state (represented by columns). Darker colors indicate better transferability.**

cost low. We then empirically measure the computational time for evaluating one month of data in each state. The results are reported in Table 4. It shows that the added overhead is small, and the runtime remains almost unchanged across all states.

**Per-class performance.** To understand how the model behaves on different road types, we evaluate different fusion strategies on Delaware. We report the AUROC for major road classes in Table 5. Link roads are not included in this analysis. The results indicate that the model performs better on Residential, Road, and Living street, while relatively lower on Primary, Trunk, and Motorway.

### 3.3 Causal Estimation Results

Building on the predictive performance of our model, we analyze how contextual factors influence accident risk. We first conduct descriptive analyses by directly aggregating accident counts across different conditions, such as season, road type, and traffic volume, to identify coarse-grained trends. To further validate these observations, we apply causal estimation techniques using learned multimodal embeddings to quantify the impact of specific factors.

**Seasonal variation.** We first analyze seasonal variations in accident occurrences by dividing the year into four seasons: winter (December to February), spring (March to May), summer (June to August), and autumn (September to November).

Most states show fewer accidents in spring than in winter. This pattern suggests that winter weather conditions contribute to higher accident risk. Snow and ice are likely major factors. In warmer states such as Nevada, where winter conditions are mild, accident counts remain stable across seasons. This pattern is consistent with the causal estimates in Table 6. The ATT values are higher in colder states like Montana and Iowa, and much lower in Nevada. As an example, Figure 3 shows the seasonal accident pattern in Massachusetts, where accident frequency peaks in winter.

**Road type.** We also examine how different road types relate to accident occurrence. Road classification reflects structural design and traffic regulation levels, both of which can shape accident risk. The categories include living street, motorway, motorway link, primary, primary link, residential, road, secondary, secondary link, tertiary, tertiary link, trailhead, trunk, and trunk link. By comparing

**Table 4: Average prediction time for one month of traffic accident data. We compare the base GIN model with three fusion variants. The results show that adding fusion modules leads to only a small increase in computation.**

| Runtime (s)      | DE  | MA   | MD   | NV   | MT   | IA   |
|------------------|-----|------|------|------|------|------|
| GIN              | 6.2 | 53.1 | 21.5 | 9.5  | 10.8 | 20.5 |
| GIN + Basic      | 6.4 | 53.6 | 22.2 | 9.8  | 12.5 | 20.8 |
| GIN + GatedFusion| 6.5 | 54.4 | 22.4 | 10.0 | 12.6 | 20.9 |
| GIN + MoE        | 6.7 | 54.5 | 23.0 | 10.2 | 12.8 | 21.7 |

**Table 5: AUROC on Delaware for different road types. We compare three fusion strategies of the GIN model. The results show that the model reaches higher accuracy on Residential, Road, and Living street, but performs worse on Primary, Trunk, and Motorway road types.**

| AUROC (%)         | Residential | Tertiary | Secondary | Primary       |
|-------------------|-------------|----------|-----------|---------------|
| GIN+Basic         | 89.8        | 83.3     | 81.1      | 78.0          |
| GIN+GatedFusion   | 90.4        | 82.8     | 81.3      | 78.2          |
| GIN+MoE           | 91.8        | 83.3     | 81.9      | 79.0          |
|                   | Motorway    | Trunk    | Road      | Living street |
| GIN+Basic         | 78.6        | 78.4     | 95.2      | 95.1          |
| GIN+GatedFusion   | 77.4        | 77.1     | 92.7      | 92.2          |
| GIN+MoE           | 80.0        | 78.7     | 94.4      | 96.3          |

accident counts and frequencies across these categories, we aim to understand how differences in road infrastructure contribute to variations in traffic safety.

The results are shown in Figure 3. Motorways have the highest average accident frequency, and the gap compared to other road types is large. This is reasonable because vehicles on motorways travel at high speeds, which reduces reaction time during unexpected events. Trunk roads show the second highest accident frequency. They also carry fast-moving traffic, which creates similar road accident risks.

Secondary, residential, and primary roads have the highest total number of accidents. This is mainly due to their broad coverage in the road network. Their large presence leads to more cumulative accidents, while less common road types contribute fewer cases. Figure 2 shows the distribution of road types and supports this overall pattern.

In our causal analysis, we treat motorways as the intervention group and obtain an ATT of 21.9%. This result indicates that road type has a measurable effect on accident occurrence.

**Precipitation.** We align accident counts with precipitation levels and observe a clear rising trend in Figure 5a.

When precipitation is below 40 mm, accident totals stay relatively low. Once precipitation exceeds 60 mm, the numbers increase rapidly and reach more than $1.6 \times 10^6$ at 160 mm. This pattern reflects the higher risks during heavy rainfall, likely caused by reduced traction, limited visibility, and longer stopping distances. The causal estimation analysis supports this result, showing an average ATT of 24.2%.

**Table 6: Average treatment effect on the treated (ATT) among all six states. We analyze the effect of seasonal variation, road type, and precipitation. We vary for different years to compute the mean and standard deviations.**

| ATT (%) | DE | MA | MD | NV | MT | IA |
|---|---|---|---|---|---|---|
| Season | $23.2_{\pm0.7}$ | $28.5_{\pm1.1}$ | $29.7_{\pm0.2}$ | $15.7_{\pm1.4}$ | $38.3_{\pm1.6}$ | $35.9_{\pm0.4}$ |
| Road type | $25.8_{\pm1.2}$ | $23.3_{\pm1.1}$ | $19.2_{\pm1.7}$ | $22.4_{\pm0.8}$ | $21.5_{\pm2.7}$ | $19.1_{\pm1.4}$ |
| Precipitation | $18.1_{\pm2.4}$ | $25.1_{\pm1.2}$ | $24.9_{\pm0.4}$ | $28.3_{\pm0.3}$ | $23.3_{\pm1.7}$ | $25.2_{\pm1.5}$ |

**Traffic volume.** We analyze the relationship between traffic volume and accident occurrence by examining the daily traffic volume corresponding to each recorded accident.

Our analysis reveals that as traffic volume increases up to approximately 200 vehicles per day, the number of accidents rises accordingly. However, beyond this threshold, the accident rate begins to decline gradually. This trend is reasonable: higher traffic volumes often lead to congestion, which reduces vehicle speeds and, consequently, lowers the likelihood of severe accidents. To understand the relationship between traffic volume and accident frequency, we visualize the distribution in Figure 5b.

## 4 Related Work

Traffic accident prediction is a dynamic task that has been explored using various foundational models. The critical role of spatial and temporal features in time-series forecasting is well established in prior work [19]. In addition, image-based features have proven valuable in related applications such as road attribute inference [6, 7]. We now discuss several areas of research that are most related to our distributions.

*Generalization of graph neural networks and language models.* GraphGPT [26] shows that LLMs can be aligned with graph structure through instruction tuning. UrbanGPT [20] further shows that instruction tuning can align spatio-temporal signals with LLMs, allowing LLMs to handle numeric time-series data and achieve strong zero-shot forecasting across urban tasks. Also, HiGPT [27] demonstrates that instruction-tuned LLMs can generalize across heterogeneous graphs by encoding node and edge semantics in natural language. Ju et al. [11] analyze the generalization in graph neural networks using PAC-Bayesian bounds. Li et al. [13] use higher-order task affinities to boost multitask learning on graphs. Li et al. [15] use gradient-based estimation to accelerate this task affinity computation. GradEx [16] introduces a first-order approach for scalable model fine-tuning. EnsembleLoRA [17] uses GradEx in multitask learning. GradSel [41] extends to the in-context learning setting and evaluates the generalization ability of language models.

*Spatiotemporal mining for traffic prediction.* The importance of jointly modeling spatial and temporal dependencies over graph-structured data for time series prediction is well recognized. Recent advances have focused on developing powerful deep learning architectures to capture these complex patterns from structured data like traffic flow and road graphs. For instance, DCRNN [19] treats traffic flow as a diffusion process on directed graphs and incorporates this into a sequence-to-sequence recurrent architecture with scheduled sampling. To improve training efficiency, STGCN [39]



(a) Precipitation    (b) Traffic volume

**Figure 5: Accident records with different ranges of precipitation and traffic volume.**

eliminates recurrent units by stacking graph and temporal convolutions into a fully convolutional framework. TEMPO [3] introduces a novel prompt-based generative pre-training framework that integrates time series decomposition with transformer architectures. MG-TAR [28] leverages dangerous driving statistics as near-miss indicators within a multi-view GNN for citywide risk prediction.

The advent of high-resolution remote sensing has opened up new possibilities for analyzing the physical state of road networks at scale [1, 18, 30, 42]. Satellite imagery provides visual cues that are strongly correlated with accident risk, but are difficult to capture in traditional datasets. Researchers have used this data for a range of safety-related tasks. Some approaches use imagery to assess proxy variables for safety, such as monitoring road surface conditions or automatically detecting infrastructure like pedestrian crossings.

*Multimodal road network analysis.* Our approach builds on prior research that combines visual models with graph-based learning for road network analysis. The RoadTagger system [6] introduced an end-to-end architecture that combines a CNN with a GNN to infer road attributes, such as lane counts, from satellite imagery. By propagating visual features along the road graph, the GNN overcomes the limited receptive field of a CNN, enabling robust inference in the presence of occlusions. Subsequently, the work on Inferring High-Resolution Traffic Accident Risk Maps [8] developed a deep learning framework to fuse multiple data sources, including satellite imagery, GPS trajectories, road maps, and historical accidents, to generate fine-grained (5m resolution) risk maps. This work demonstrated how to overcome data sparsity challenges by leveraging a rich, multimodal context. The challenge of integrating heterogeneous data modalities on a graph is also studied in multi-modal graph learning [5]. Our work contributes to this area in the context of road safety modeling.

Inspired by progress in computer vision, several large-scale satellite imagery datasets have recently been introduced [1, 18, 30, 42]. OAM-TCD [30] focuses on high-resolution tree crown delineation with global diversity, while SolarCube [18] integrates satellite and ground data for solar radiation forecasting across continents. Road-Tracer [1] addresses road network extraction via iterative graph construction from aerial images.

**Table 7: A summary of our dataset and several existing datasets. Our dataset combines large-scale traffic accident records with aligned satellite imagery, while prior datasets either lack image data or have limited volume.**

| Dataset | Year | Volume | Spatial | Time Series | Tabular | Satellite Images | Category |
|---|---|---|---|---|---|---|---|
| METR-LA [19] | 2018 | $6M+$ | ✓ | ✓ | ✓ | ✗ | Traffic Forcast |
| RoadTracer [1] | 2018 | 300 | ✗ | ✗ | ✗ | ✓ | Semantic Segmentation |
| SEVIR [29] | 2020 | $10K+$ | ✓ | ✓ | ✗ | ✓ | Weather Forcast |
| ML4Roadsafety [23] | 2023 | $9M+$ | ✓ | ✓ | ✓ | ✗ | Traffic Accident Analysis |
| CrashFormer [12] | 2023 | $6M+$ | ✓ | ✓ | ✓ | ✗ | Traffic Accident Analysis |
| OAM-TCD [30] | 2024 | $5K+$ | ✓ | ✓ | ✗ | ✓ | Semantic Segmentation |
| SolarCube [18] | 2024 | $600K+$ | ✓ | ✓ | ✗ | ✓ | Weather Forcast |
| FT-AED [4] | 2024 | $3M+$ | ✓ | ✓ | ✓ | ✗ | Traffic Detection |
| Tumtraffic-QA [43] | 2025 | $90K+$ | ✓ | ✓ | ✗ | ✓ | Traffic Detection |
| MMTraCE (This work) | 2025 | $10M+$ | ✓ | ✓ | ✓ | ✓ | Traffic Accident Analysis |

## 5 Discussion

The satellite images in our dataset capture road characteristics that change very slowly, such as lane width, curvature, and intersection density. These features remain steady across long periods, so small differences between the imagery date and the accident records do not create clear sources of bias. The dynamic elements of accident risk come from other modalities. Weather observations and time-varying traffic volume (AADT) provide the signals that describe short-term changes in road conditions.

This combination allows us to model relative risk across the network in a stable and interpretable way. Our goal is not to forecast the exact time or location of an accident. The predicted values instead represent the likelihood or hazard level of each road segment. They offer a practical measure of traffic risk that can support driver awareness, road maintenance decisions, and safety planning.

Seen from this angle, our dataset opens up several directions for future work. It provides a strong benchmark for multimodal learning methods that operate on both graph and image data. It also demonstrates that visual cues, structural features, and traffic signals can work together to reveal fine-grained patterns in accident risk. This creates space for developing stronger fusion models, exploring more detailed temporal patterns, and building tools that help agencies use data-driven insights for safety interventions.

**The MMTraCE package.** To support future research, we release our dataset on Huggingface along with an easy-to-use Python package for streamlined access. The full dataset is available via the DATASETS library and includes road network graphs, accident records, weather data, traffic volume statistics, and satellite image embeddings. The embeddings are generated using Vision Transformer and CLIP models. We also include a set of satellite images with the package for all states.

With a single line of code, users can load the complete dataset. To retrieve data for a specific state, one simply specifies the state name, and the package will automatically download, cache, and return the corresponding dataset object. We also provide functionality to extract accident records and features for any given month. Additionally, the package includes a trainer module for training and evaluating baseline models used in our framework, which encompasses basic GNN models, spatial-temporal models, and fusion models. We add some examples of the usage in Appendix B.3.

## 6 Conclusion

We present a large and up-to-date multimodal traffic dataset that includes road networks, weather data, traffic volumes, accident records, and satellite images from six U.S. states. The dataset is designed to support research on traffic accident prediction and to provide a comprehensive view of the factors that shape road safety. It brings together long-term records from multiple public sources and aligns them at the road-segment level, making it suitable for both predictive modeling and causal analysis.

On top of this dataset, we build a GNN-based framework that integrates visual, structural, and temporal features. The model combines information from all modalities and achieves strong predictive performance. The best baseline reaches an average AUROC of 90.1% across states, showing that multimodal representations are effective for accident prediction. Our analysis also includes causal estimation and feature ablation. These studies help clarify how different contextual factors contribute to accident risk and how each feature type affects model performance.

## Acknowledgments

## References

[1] Favyen Bastani, Songtao He, Sofiane Abbar, Mohammad Alizadeh, Hari Balakrishnan, Sanjay Chawla, Sam Madden, and David DeWitt. 2018. Roadtracer: Automatic extraction of road networks from aerial images. In *Computer Vision and Pattern Recognition (CVPR)*. 4720–4728. 8, 9

[2] Lawrence J Blincoe, Ted R Miller, Eduard Zaloshnja, and Bruce Lawrence. 2015. *The economic and societal impact of motor vehicle crashes, 2010 (Revised)*. Technical Report. United States. Department of Transportation. National Highway Traffic Safety Administration. 1

[3] Defu Cao, Furong Jia, Sercan O Arik, Tomas Pfister, Yixiang Zheng, Wen Ye, and Yan Liu. 2024. TEMPO: Prompt-based Generative Pre-trained Transformer for Time Series Forecasting. In *International Conference on Learning Representations (ICLR)*. 8

[4] Austin Coursey, Junyi Ji, Marcos Quinones Grueiro, William Barbour, Yuhang Zhang, Tyler Derr, Gautam Biswas, and Daniel Work. 2024. FT-AED: Benchmark dataset for early freeway traffic anomalous event detection. *Advances in Neural Information Processing Systems (NeurIPS)* 37 (2024), 15526–15549. 9

[5] Yasha Ektefaie, George Dasoulas, Ayush Noori, Maha Farhat, and Marinka Zitnik. 2023. Multimodal learning with graphs. *Nature Machine Intelligence* 5, 4 (2023), 340–350. 8

[6] Songtao He, Favyen Bastani, Satvat Jagwani, Edward Park, Sofiane Abbar, Mohammad Alizadeh, Hari Balakrishnan, Sanjay Chawla, Samuel Madden, and Mohammad Amin Sadeghi. 2020. Roadtagger: Robust road attribute inference with graph neural networks. In *Association for the Advancement of Artificial Intelligence (AAAI)*. 10965–10972. 8

[7] Songtao He, Mohammad Amin Sadeghi, Sanjay Chawla, Mohammad Alizadeh, Hari Balakrishnan, and Samuel Madden. 2021. Inferring high-resolution traffic accident risk maps based on satellite imagery and GPS trajectories. In *International Conference on Computer Vision (ICCV)*. 11957–11965. 1, 8

[8] Songtao He, Mohammad Amin Sadeghi, Sanjay Chawla, Mohammad Alizadeh, Hari Balakrishnan, and Samuel Madden. 2021. Inferring high-resolution traffic accident risk maps based on satellite imagery and GPS trajectories. In *International Conference on Computer Vision (ICCV)*. 11977–11985. 8

[9] Baixiang Huang, Bryan Hooi, and Kai Shu. 2023. Tap: A comprehensive data repository for traffic accident prediction in road networks. In *Proceedings of the 31st ACM International Conference on Advances in Geographic Information Systems*. 1–4. 2

[10] Erik Jenelius and Haris N Koutsopoulos. 2013. Travel time estimation for urban road networks using low frequency probe vehicle data. *Transportation Research Part B: Methodological* 53 (2013), 64–81. 1

[11] Haotian Ju, Dongyue Li, Aneesh Sharma, and Hongyang R Zhang. 2023. Generalization in graph neural networks: Improved pac-bayesian bounds on graph diffusion. In *International conference on artificial intelligence and statistics*. PMLR, 6314–6341. 8

[12] Amin Karimi Monsefi, Pouya Shiri, Ahmad Mohammadshirazi, Nastaran Karimi Monsefi, Ron Davies, Sobhan Moosavi, and Rajiv Ramnath. 2023. Crashformer: A multimodal architecture to predict the risk of crash. In *ACM SIGSPATIAL International Workshop on Advances in Urban-AI*. 42–51. 9

[13] Dongyue Li, Haotian Ju, Aneesh Sharma, and Hongyang R Zhang. 2023. Boosting multitask learning on graphs through higher-order task affinities. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 1213–1222. 8

[14] Dongyue Li, Huy Nguyen, and Hongyang R. Zhang. 2024. Identification of Negative Transfers in Multitask Learning Using Surrogate Models. *Transactions on Machine Learning Research* (2024). 6

[15] Dongyue Li, Aneesh Sharma, and Hongyang R Zhang. 2024. Scalable Multitask Learning Using Gradient-based Estimation of Task Affinity. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 1542–1553. 8

[16] Dongyue Li, Ziniu Zhang, Lu Wang, and Hongyang R Zhang. 2024. Scalable Fine-tuning from Multiple Data Sources: A First-Order Approximation Approach. *Findings of the Association for Computational Linguistics: EMNLP 2024* (2024). 8

[17] Dongyue Li, Ziniu Zhang, Lu Wang, and Hongyang R Zhang. 2025. Efficient Ensemble for Fine-tuning Language Models on Multiple Datasets. In *Association for Computational Linguistics (ACL)*. 8

[18] Ruohan Li, Yiqun Xie, Xiaowei Jia, Dongdong Wang, Yanhua Li, Yingxue Zhang, Zhihao Wang, and Zhili Li. 2024. SolarCube: An Integrative Benchmark Dataset Harnessing Satellite and In-situ Observations for Large-scale Solar Energy Forecasting. *Advances in Neural Information Processing Systems (NeurIPS)* 37 (2024), 3499–3513. 8, 9

[19] Yaguang Li, Rose Yu, Cyrus Shahabi, and Yan Liu. 2018. Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting. In *International Conference on Learning Representations (ICLR)*. 1, 5, 8, 9

[20] Zhonghang Li, Lianghao Xia, Jiabin Tang, Yong Xu, Lei Shi, Long Xia, Dawei Yin, and Chao Huang. 2024. Urbangpt: Spatio-temporal large language models. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 5351–5362. 8

[21] Chuizheng Meng, Sam Griesemer, Defu Cao, Sungyong Seo, and Yan Liu. 2025. When physics meets machine learning: A survey of physics-informed machine learning. *Machine Learning for Computational Science and Engineering* 1, 1 (2025), 20. 1

[22] National Highway Traffic Safety Administration. 2025. *Early Estimate of Motor Vehicle Traffic Fatalities for 2024*. Technical Report. U.S. Department of Transportation. 1

[23] Abhinav Nippani, Dongyue Li, Haotian Ju, Haris Koutsopoulos, and Hongyang Zhang. 2023. Graph neural networks for road safety modeling: Datasets and evaluations for accident analysis. *Advances in neural information processing systems (NeurIPS)* 36 (2023), 52009–52032. 1, 2, 4, 9

[24] World Health Organization. 2019. *Global status report on road safety 2018*. World Health Organization, Geneva, Switzerland. 1

[25] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations. In *International Conference on Knowledge Discovery and Data Mining (KDD)*. ACM, New York City, USA, 701–710. 5

[26] Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Lixin Su, Suqi Cheng, Dawei Yin, and Chao Huang. 2024. Graphgpt: Graph instruction tuning for large language models. In *ACM SIGIR Conference on Research and Development in Information Retrieval*. 491–500. 8

[27] Jiabin Tang, Yuhao Yang, Wei Wei, Lei Shi, Long Xia, Dawei Yin, and Chao Huang. 2024. Higpt: Heterogeneous graph language model. In *ACM SIGKDD conference on knowledge discovery and data mining (KDD)*. 2842–2853. 8

[28] Patara Trirat, Susik Yoon, and Jae-Gil Lee. 2023. MG-TAR: Multi-view graph convolutional networks for traffic accident risk prediction. *IEEE Transactions on Intelligent Transportation Systems* 24, 4 (2023), 3779–3794. 8

[29] Mark Veillette, Siddharth Samsi, and Chris Mattioli. 2020. Sevir: A storm event imagery dataset for deep learning applications in radar and satellite meteorology. *Advances in Neural Information Processing Systems (NeurIPS)* 33 (2020), 22009–22019. 9

[30] Josh Veitch-Michaelis, Andrew Cottam, Daniella Schweizer, Eben Broadbent, David Dao, Ce Zhang, Angelica Almeyda Zambrano, and Simeon Max. 2024. OAM-TCD: A globally diverse dataset of high-resolution tree cover maps. *Advances in Neural Information Processing Systems (NeurIPS)* 37 (2024), 49749–49767. 8, 9

[31] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. *stat* 1050, 20 (2017), 10–48550. 5

[32] Qingyi Wang, Shenhao Wang, Yunhan Zheng, Hongzhou Lin, Xiaohu Zhang, Jinhua Zhao, and Joan Walker. 2024. Deep hybrid model with satellite imagery: How to combine demand modeling and computer vision for travel behavior analysis? *Transportation Research Part B: Methodological* 181 (2024), 102914. 1

[33] Shenhao Wang, Baichuan Mo, Yunhan Zheng, Stephane Hess, and Jinhua Zhao. 2024. Comparing hundreds of machine learning and discrete choice models for travel demand modeling: An empirical benchmark. *Transportation Research Part B: Methodological* 190 (2024), 103061. 1

[34] Qiming Wu, Zichen Chen, Will Corcoran, Misha Sra, and Ambuj Singh. 2025. GraphEval36K: Benchmarking Coding and Reasoning Capabilities of Large Language Models on Graph Datasets. In *Findings of the Association for Computational Linguistics: NAACL 2025*. 8095–8117. 1

[35] Sen Wu, Hongyang R Zhang, and Christopher Ré. 2020. Understanding and Improving Information Transfer in Multi-Task Learning. In *International Conference on Learning Representations*. 6

[36] Fan Yang, Hongyang R Zhang, Sen Wu, Christopher Ré, and Weijie J Su. 2025. Precise high-dimensional asymptotics for quantifying heterogeneous transfers. *Journal of Machine Learning Research* 26, 113 (2025), 1–88. 6

[37] Qi Yang and Haris N Koutsopoulos. 1996. A microscopic traffic simulator for evaluation of dynamic traffic management systems. *Transportation research part C: emerging technologies* 4, 3 (1996), 113–129. 1

[38] Huaxiu Yao, Xianfeng Tang, Hua Wei, Guanjie Zheng, and Zhenhui Li. 2019. Revisiting spatial-temporal similarity: A deep learning framework for traffic prediction. In *Association for the Advancement of Artificial Intelligence (AAAI)*. 5668–5675. 1

[39] Bing Yu, Haoteng Yin, and Zhanxing Zhu. 2018. Spatio-Temporal Graph Convolutional Networks: A Deep Learning Framework for Traffic Forecasting. In *International Joint Conference on Artificial Intelligence (IJCAI)*. 3634–3640. 8

[40] Zhuoning Yuan, Xun Zhou, and Tianbao Yang. 2018. Hetero-convlstm: A deep learning approach to traffic accident prediction on heterogeneous spatio-temporal data. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*. 984–992. 2

[41] Ziniu Zhang, Zhenshuo Zhang, Dongyue Li, Lu Wang, Jennifer Dy, and Hongyang R Zhang. 2025. Linear-Time Demonstration Selection for In-Context Learning via Gradient Estimation. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 16470–16488. 8

[42] Hangyu Zhou, Chia Hsiang Kao, Cheng Perng Phoo, Utkarsh Mall, Bharath Hariharan, and Kavita Bala. 2024. AllClear: A Comprehensive Dataset and Benchmark for Cloud Removal in Satellite Imagery. *Advances in Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track* 37 (2024), 53571–53597. 8

[43] Xingcheng Zhou, Konstantinos Larintzakis, Hao Guo, Walter Zimmer, Mingyu Liu, Hu Cao, Jiajie Zhang, Venkatnarayanan Lakshminarasimhan, Leah Strand, and Alois C Knoll. 2025. Tumtraffic-videoqa: A benchmark for unified spatio-temporal video understanding in traffic scenes. In *International Conference on Machine Learning (ICML)*. PMLR. 9

[44] Dingyi Zhuang, Yuheng Bu, Guang Wang, Shenhao Wang, and Jinhua Zhao. 2024. Sauc: Sparsity-aware uncertainty calibration for spatiotemporal prediction with graph neural networks. In *ACM International Conference on Advances in Geographic Information Systems (SIGSPATIAL)*. 160–172. 1

[45] Dingyi Zhuang, Shenhao Wang, Haris Koutsopoulos, and Jinhua Zhao. 2022. Uncertainty quantification of sparse travel demand prediction with spatial-temporal graph neural networks. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*. 4639–4647. 1

# A  Data Collection Procedure

In this section, we provide more details about our data collection procedure, including how to construct the road network, how to align the location of accident records, how to normalize traffic volumes from different formats, and how to obtain weather features. Then, we summarize all the features we collected.

## A.1  Road Networks

We construct the road network using the OSMnx Street Network Dataverse dataset. For each state, we load street networks at multiple administrative levels, including cities, counties, neighborhoods, census tracts, and urbanized areas. These individual networks are then concatenated to form a comprehensive statewide road graph.

From the raw data, we extract a list of nodes and edges. For each node, we retain the node ID along with its geographic coordinates (latitude and longitude). For each edge, we record the start node ID, end node ID, one-way or not, the corresponding road type, and the length of the road.

## A.2  Traffic Accident Records

We collect the traffic accident data provided by each state's DOT, which varies in different formats.

In our analysis, we focus specifically on extracting the latitude and longitude information from these records. When geographic coordinates are not directly available in the original dataset, we use the textual address descriptions, such as street names or intersections, and query the corresponding latitude and longitude using the Google Maps API. This step ensures that all accident locations are consistently represented in a geospatial format, which is critical for downstream spatial analysis and visualization.

Next, we align each accident record to a specific road segment in the network. Let the geocoded accident location be denoted as $c \in \mathbb{R}^2$. For each edge $e \in E$, associated with endpoints $e_a$ and $e_b$, we define $D(\cdot, \cdot)$ as the Euclidean distance function. To assign the accident to the most plausible road segment, we adopt the following heuristic:

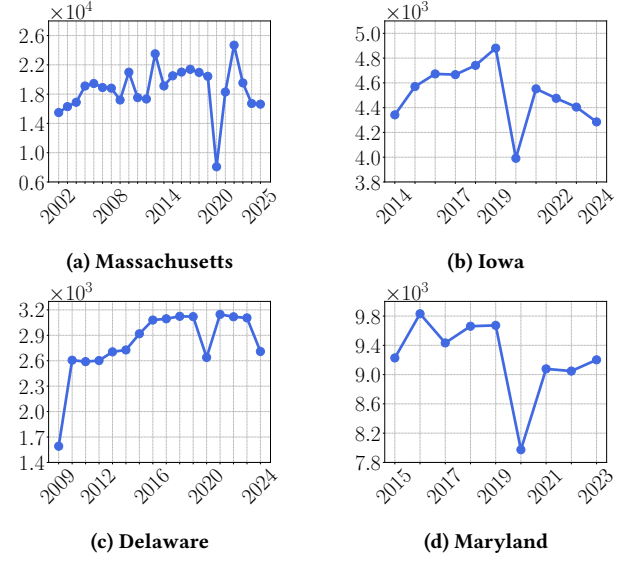$$e_{\text{acc}} = \arg\max_{e \in E} \left( D(e_a, e_b) - (D(e_a, c) + D(e_b, c)) \right).$$

This formulation prioritizes edges for which the accident point lies closest to the segment span between $e_a$ and $e_b$.

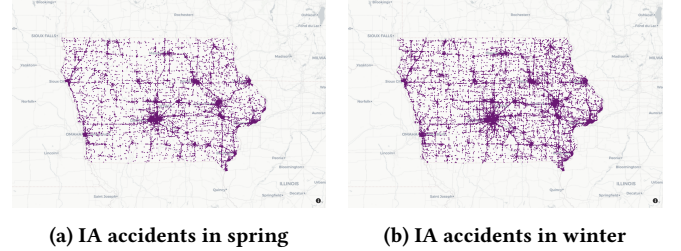We illustrate the traffic accident records in Figure 6 and Figure 7 below.

## A.3  Traffic Volume

We obtain traffic volume data from official state-level transportation departments and use it as a dynamic edge-level feature in our multimodal framework.

Some states provide clean CSV exports containing longitude and latitude fields, while others (such as Massachusetts and Montana) publish data in more complex formats, including KML and GeoJSON, with inconsistent column names and nested spatial structures. In these cases, we developed custom parsers for each state and manually mapped the relevant fields to extract usable geospatial information. Additionally, the temporal resolution differs by state, and in some cases, the data spans dozens of measurement points per segment, which requires iterative merging and normalization.



(a) Massachusetts

(b) Iowa

(c) Delaware

(d) Maryland

**Figure 6: The average number of accidents per month for each year in Massachusetts, Iowa, Delaware, and Maryland. The sharp drop in 2020 is due to the impact of COVID-19.**



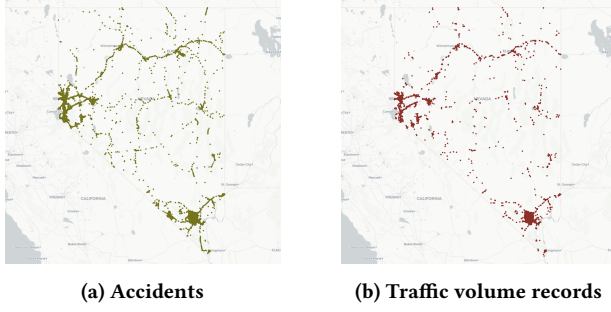(a) IA accidents in spring

(b) IA accidents in winter

**Figure 7: Seasonal comparison of traffic accidents in Iowa during spring and winter, respectively, showing that accidents in winter occur more than in spring.**

Often, the same road segment is represented by multiple sensor points along its path, with no explicit intersection or highway labeling, so we heuristically aggregated these measurements to associate one traffic volume value with each edge.

We report a distribution map of accidents and AADT records in Figure 8. It indicates that the traffic volume monitor has covered nearly all potential road segments with accidents.

## A.4  Weather Statistics

The weather data is extracted using the Meteostat API, focusing on monthly weather statistics. For each node in the road network, the nearest weather station is identified based on geographic coordinates (latitude and longitude). We then retrieve monthly weather data for each location over a specified time range. The time range aligns with the time range of the traffic accident record. After fetching the data, it merges this with the corresponding node ID and coordinates, and stores the results in CSV files. Finally, all files are

(a) Accidents      (b) Traffic volume records

**Figure 8: The distribution of accidents and traffic volume monitor coverage in Nevada. (a) marks the historical accident records across the state, and (b) shows the road segments with recorded traffic volumes. The traffic volume records cover nearly all road segments where accidents have occurred.**

concatenated to create a comprehensive dataset matching historical weather conditions to each road network node.

We collect six types of features: the average daily air temperature in °C, the average daily minimum air temperature in °C, the average daily maximum air temperature in °C, the monthly precipitation total in mm, the average wind speed in km/h, and the average sea-level air pressure in hPa.

In summary, we report all the features we collect from our road networks in this section. For node features, we collect 9 types of features: latitude, longitude, and satellite image for static features, average surface temperature, max surface temperature, min surface temperature, total precipitation, average wind speed, and sea level air pressure for dynamic features. For edge features, we collect 3 types of features: location, length, and road type for static features, and traffic volume (annual average daily traffic) for dynamic features.

## B Experimental Details

### B.1 Implementation Details

We report the hyperparameters in our experiments and give the formal definition of the fusion methods and evaluation metrics in the main result.

*Basic fusion method.* The basic fusion method uses multilayer perceptrons to combine different features. Specifically, we first obtain the structure-aware representation $x_i^{(\text{GNN})}$ for each node $i$ by passing its attributes and graph connectivity through a multi-layer GNN. Independently, we extract a visual embedding $z_i = f_{\text{vision}}(I_i)$ for the image $I_i$ associated with node $i$ using a pretrained vision backbone such as CLIP or Vision Transformer (ViT), where $f_{\text{vision}}$ denotes the vision encoder.

To produce the final representation for downstream prediction, we concatenate the GNN embedding and the visual embedding:

$$\tilde{x}_i = [x_i^{(\text{GNN})} \parallel z_i], \tag{1}$$

and feed the result into a multilayer perceptron (MLP) $f_{\text{mlp}}$, trained jointly with the GNN:

$$\hat{y}_i = f_{\text{mlp}}(\tilde{x}_i). \tag{2}$$

The fusion MLP consists of four fully connected layers with non-linear activation functions (e.g., ReLU) and is optimized end-to-end alongside the GNN during training. This multimodal setup allows the model to leverage both spatial-structural and high-level visual cues, potentially improving predictive performance on downstream tasks such as accident risk prediction or traffic pattern classification.

*Gated fusion network.* While MLP-based concatenation provides a straightforward way to combine graph and visual features, it treats both modalities equally for all nodes. To enable more flexible and context-aware fusion, we adopt a gated mechanism that adaptively weighs the contribution of each modality.

Since the vision and GNN embeddings may differ in dimensionality, the visual feature $z$ is first projected to match the GNN embedding space via a learnable linear transformation: $z' = W_{\text{proj}}z$, where $W_{\text{proj}}$ is a trainable matrix.

A scalar gate value $\lambda \in [0, 1]$ is then computed to determine the relative importance of the two modalities. The gate is derived from their concatenation:

$$\lambda = \sigma\left(f_{\text{gate}}\left([x^{(\text{GNN})} \parallel z]\right)\right),$$

where $f_{\text{gate}}$ is a small MLP and $\sigma$ denotes the sigmoid function.

The final fused representation is a convex combination of the two:

$$\tilde{x} = \lambda \cdot x^{(\text{GNN})} + (1 - \lambda) \cdot z'.$$

This fused embedding $\tilde{x}$ is then fed into an MLP for prediction.

*Mixture of Experts.* Finally, we implement the mixture of experts (MoE) approach, which leverages multiple specialized networks, referred to as experts, that each learn to combine features from different perspectives. Given the visual feature $z$ and the graph embedding $x^{(\text{GNN})}$, we first construct a shared representation by concatenating the two modalities: $\tilde{x} = [x^{(\text{GNN})} \parallel z']$. This shared feature vector is then fed into a set of $K$ expert networks, each parameterized by its own learnable weights. Each expert $f_k$ produces an output representation:

$$e^{(k)} = f_k(\tilde{x}), \quad k = 1, \ldots, K.$$

To dynamically select and combine the outputs of these experts, we utilize the gated network that computes a probability distribution over experts for each node. This gating mechanism is realized as a small multilayer perceptron followed by a softmax activation:

$$\lambda = \text{softmax}(f_{\text{gate}}(\tilde{x})),$$

where $\lambda \in \mathbb{R}^K$ and $\sum_{k=1}^{K} \lambda^{(k)} = 1$. The final fused representation is then computed as a weighted sum of the expert outputs:

$$\tilde{e} = \sum_{k=1}^{K} \lambda^{(k)} \cdot e^{(k)}.$$

The fused expert output $\tilde{e}$ is subsequently passed through a prediction head along with the edge feature $x_{\text{edge}}$ to generate the model's output:

$$\hat{y} = f_{\text{pred}}(\tilde{e}, x_{\text{edge}}).$$

This MoE architecture enables the model to dynamically select among multiple fusion pathways, capturing complex relationships between different modalities. With a learned gating mechanism,

the MoE framework can adaptively choose the most relevant experts for each node, resulting in richer, more flexible multimodal representations.

*Experiment hyperparameters.* For all baselines, we construct edge representations by concatenating the embeddings of the two connected nodes with the associated edge features. Node embeddings are fixed to a dimension of 128. We use two-layer MLP and GNN architectures, both with a hidden size of 256. The GNN, MLP, and fusion modules are fully trainable. All models are optimized using the Adam algorithm with a learning rate of 0.001 and trained for 30 epochs.

We split all the accident records into training, validation, and test sets for each state based on different periods of years. The summarization is shown in Table 8. Since our current analysis is at the monthly level, we note that our framework can also be utilized for conducting analysis at the yearly level.

*Evaluation metrics.* Formally, for the regression task, let $y_i \in \mathbb{R}$ denote the ground-truth number of accidents for road segment $i$, and let $\hat{y}_i$ be the corresponding model prediction. The Mean Absolute Error (MAE) is computed as: $\mathcal{L}_{\text{MAE}} = \frac{1}{N} \sum_{i=1}^{N} |\hat{y}_i - y_i|$, where $N$ is the total number of road segments in the evaluation set.

For the classification task, let $y_i \in \{0, 1\}$ be the binary indicator of whether an accident occurred on road segment $i$, and let $p_i \in [0, 1]$ denote the predicted probability score. The AUROC measures the probability that a randomly chosen positive example is assigned a higher score than a randomly chosen negative one, and is formally defined as:

$$\mathcal{L}_{\text{AUROC}} = \frac{1}{|S^+||S^-|} \sum_{i \in S^+} \sum_{j \in S^-} f(p_i > p_j),$$

where $S^+$ and $S^-$ denote the sets of indices corresponding to positive and negative samples, respectively, and $f(\cdot)$ is the indicator function.

*Omitted causal analysis.* Beyond direct matching, we also estimate causal effects with propensity score-based adjustments (PSM). We model the treatment assignment probability with

$$e(x_i) = \Pr(t_i = 1 \mid x_i),$$

using logistic regression on the multimodal embeddings. Matching in this score space aligns treated and control samples with similar treatment likelihoods and yields the PSM estimate

$$\hat{\tau}_{\text{PSM}} = \frac{1}{|T|} \sum_{i \in T} (y_i - y_{j(i)}),$$

where $j(i)$ is the matched control with the closest propensity score.

To increase robustness, we further apply a doubly robust estimator (DR) that combines outcome regression and propensity weighting. Let $\hat{m}_1(x_i)$ and $\hat{m}_0(x_i)$ be predicted outcomes under treatment and control. The DR estimator is written in a compact form as

$$\hat{\tau}_{\text{DR}} = \frac{1}{|T|} \sum_{i \in T} \left( \hat{m}_1(x_i) - \hat{m}_0(x_i) \right) + \frac{1}{|T|} \sum_{i \in T} \omega_i,$$

where the correction term is

$$\omega_i = \frac{t_i \left( y_i - \hat{m}_1(x_i) \right)}{\hat{e}(x_i)} - \frac{(1 - t_i) \left( y_i - \hat{m}_0(x_i) \right)}{1 - \hat{e}(x_i)}.$$

**Table 8: The train, validation, and test data splitting of Delaware, Massachusetts, Maryland, Nevada, Montana, and Iowa in our framework. We report the period and associated accident records.**

| | Train | | Valid | | Test | |
|---|---|---|---|---|---|---|
| | period | records | period | records | period | records |
| DE | 2009 – 2013 | 145127 | 2014 – 2018 | 179335 | 2019 – 2024 | 208650 |
| MA | 2002 – 2014 | 2887528 | 2015 – 2020 | 1348597 | 2021 – 2025 | 970504 |
| MD | 2015 – 2017 | 341902 | 2018 – 2019 | 232002 | 2020 – 2024 | 423628 |
| NV | 2016 – 2018 | 154834 | 2019 – 2020 | 92359 | 2021 – 2025 | 129059 |
| MT | 2016 – 2018 | 60717 | 2019 – 2020 | 15400 | 2021 – 2023 | 63894 |
| IA | 2014 – 2017 | 219013 | 2018 – 2020 | 161186 | 2021 – 2023 | 161186 |

This embedding-based approach allows the causal analysis to control for confounding factors encoded in the multimodal representations. The PSM and DR adjustments provide more stable estimates when treatment imbalance exists, and support comparisons of treatment effects under similar structural, visual, and weather conditions.

## B.2 Omitted Experiment Results

In addition to AUROC, we also evaluate the performance of different baselines using the precision and recall scores as additional metrics.

*Results using more metrics.* Table 11 reports the test precision and recall scores across six U.S. states using a range of baselines, including node embedding methods, vision-based models, graph neural networks, and their multimodal extensions via fusion techniques. Across all models, GIN + MoE in Delaware (67.04%), while the highest recall is observed in the same model for Montana (99.76%), suggesting that mixture-of-experts can significantly boost detection capability under certain regional conditions.

In general, GIN-based fusion models show clear improvements over their unimodal counterparts. On average across six states, GIN + MoE improves precision by 2.41% and recall by 0.03% compared to the basic GIN model. These gains highlight the effectiveness of mixture-of-experts multimodal fusion in enhancing detection precision while maintaining strong recall performance across diverse traffic environments.

Nonetheless, we observe substantial variance in precision across regions, with some models showing high recall but relatively low precision, especially in Montana and Iowa. This indicates that while these models successfully capture most positive cases, they also generate a significant number of false positives. These findings underscore the effectiveness of multimodal fusion strategies while also highlighting challenges in balancing precision and recall, especially in regions with sparse accident data or heterogeneous road conditions.

*Causal analysis results of PSM and DR..* To enhance the causal analysis, we add Propensity Score Matching (PSM) and Doubly Robust (DR) estimators for more reliable treatment effect estimation, using one nearest neighbor and evaluating effects by road type. The results show a clear increase in accident risk for motorway segments, with effects of around 20%. To assess robustness, we also vary the number of nearest neighbors in the matching step.

**Table 9: Average treatment effect on the treated (ATT) across six U.S. states. The top block reports results from PSM and DR under one-nearest-neighbor matching ($k = 1$). The bottom block varies the number of neighbors for DR to assess robustness. We evaluate treatment effects related to road type. We calculate the mean and standard deviation across multiple years.**

|          | DE         | MA         | MD         | NV         | MT         | IA         |
|----------|------------|------------|------------|------------|------------|------------|
| PSM      | 18.1 ± 0.2 | 19.0 ± 0.4 | 16.5 ± 0.1 | 17.2 ± 0.2 | 14.8 ± 0.2 | 12.2 ± 0.1 |
| DR ($k = 1$) | 37.7 ± 4.2 | 21.0 ± 3.1 | 21.7 ± 3.8 | 21.8 ± 3.4 | 19.7 ± 3.1 | 16.3 ± 2.5 |
| DR ($k = 3$) | 37.7 ± 3.4 | 20.9 ± 3.2 | 21.8 ± 3.8 | 21.8 ± 3.4 | 19.7 ± 3.2 | 16.2 ± 2.4 |
| DR ($k = 5$) | 37.7 ± 3.4 | 20.9 ± 3.2 | 21.8 ± 3.8 | 21.8 ± 3.4 | 19.7 ± 3.2 | 16.2 ± 2.4 |

The estimates stay consistent across settings, showing that the DR estimator is stable.

## B.3 Examples of Using MMTraCE Package

Here we report some examples of how to implement MMTraCE package, including data loader and trainer.

```
>>> from mmtrace import Trainer, Evaluator, MMDataset
# Create the dataset
>>> dataset = MMDataset(state_name = "MA")
>>> data = dataset.load_monthly_data(year = 2024, month = 12)
>>> acci, acci_cnt = data["accidents"], data["accident_cnt"]
# Load different types of features
>>> node_attr, edge_attr = data["x"], data["edge_attr"]
>>> img_attr = data["img_embeddings"]
# Get an evaluator for accident prediction
>>> evaluator = Evaluator(type = "classification")
# Initialize a trainer with the model, dataset, and evaluator
>>> trainer = Trainer(model, dataset, evaluator, ...)
# Conduct training and evaluation inside the trainer
>>> log = trainer.train()
```

**Listing 1: MMTraCE data loader and trainer**

From the code in Listing 1, we can access the data from a specific state, year, and month. We can also fetch the static and dynamic features. Node features, edge features, and edge embeddings are all available. Then, we can initialize a trainer for implementation. The model can be applied to any baseline model.

## C Privacy Details and Leakage Audit

All data sources used in this study are publicly available and released under open data licenses by the corresponding agencies. State Departments of Transportation provide traffic accident and traffic volume records through official open data portals that support academic and non-commercial use. The OSMnx Street Network Dataverse is distributed under the Open Database License (ODbL), which allows reuse under standard attribution and share-alike terms. Mapbox imagery is obtained through the official Static Tiles API under the standard developer terms of service. Weather data is collected from the official Meteostat API, which offers open access to historical and near-real-time meteorological observations for research. Table 10 lists all sources, license statements, and links.

All datasets used in this work contain only anonymized records released by public agencies. No record includes personal identifiers, and no field can be linked back to individuals. The spatial resolution of the satellite and grid-level imagery does not allow identification of people, vehicles, or private property details. Our processing pipeline uses only node-level and segment-level information derived from public road networks, which ensures that the final dataset remains fully anonymous and consistent with public data release standards.

In addition to these safeguards, we conduct a full leakage audit to verify that no hidden pathways for re-identification remain. We check that accident counts are aggregated at the segment or monthly level, that imagery does not contain fine-grained visual cues tied to individuals, and that no timestamps, device traces, or location histories enter the dataset during pre-processing. We also confirm that structural information from OSMnx cannot be cross-matched with any external source to recover personal movement patterns. These checks show that the final dataset cannot be used to infer identities or sensitive attributes. In our experiments, we divide the entire time series into three non-overlapping segments for training, validation, and testing, as shown in Table 8. The splits are strictly chronological, ensuring that no future information leaks into the training or validation sets.
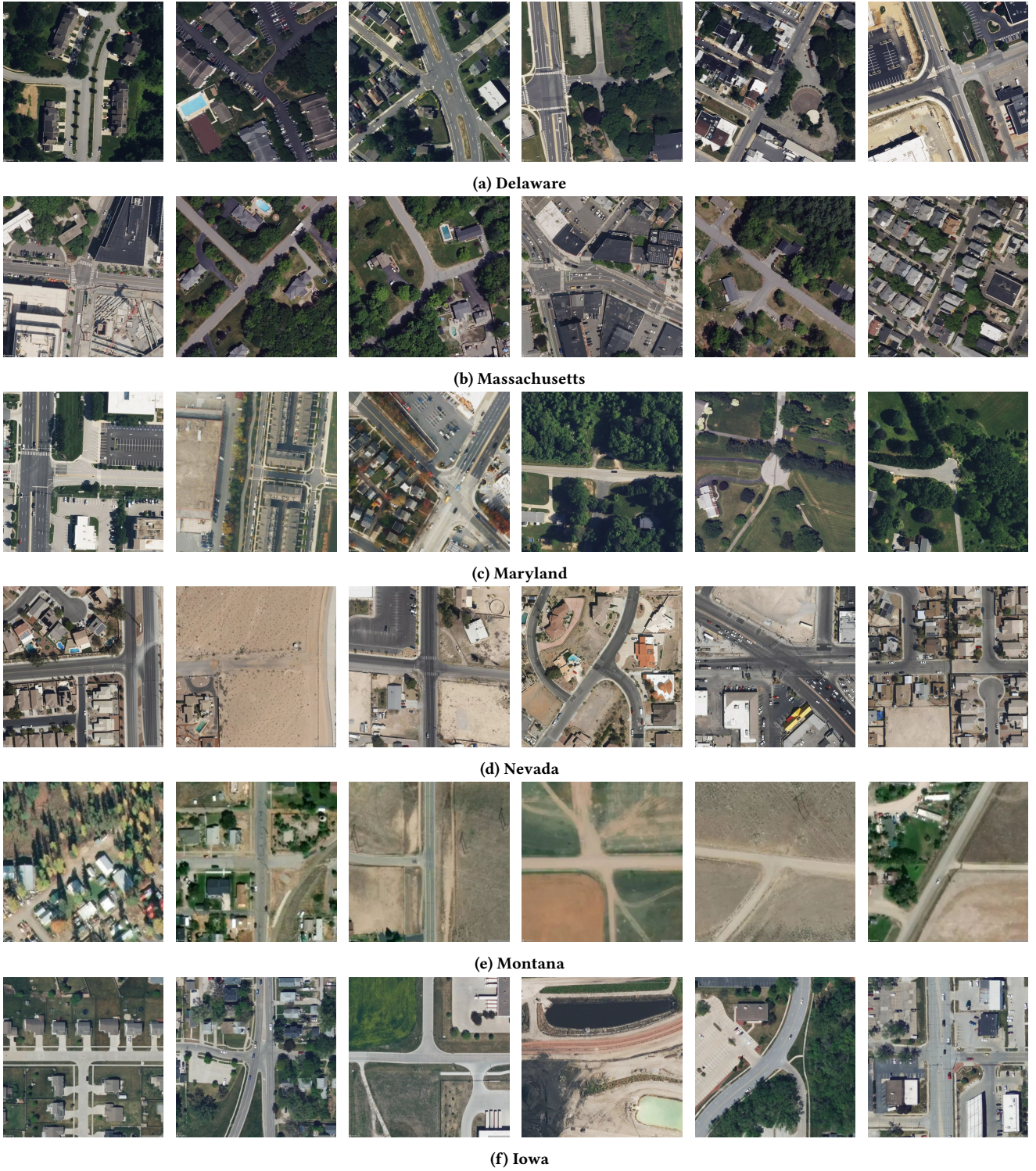
**Table 10: Summary of the data sources used in our dataset construction, including links to official repositories for traffic accident records, traffic volume statistics, and APIs to road network data, weather observations, and satellite imagery. They are official records and serve as the foundation for building our multimodal traffic safety dataset.**

| | |
|---|---|
| **Traffic Accident Records** | |
| Delaware DOT | https://data.delaware.gov/Transportation/Public-Crash-Data/827n-m6xc/about_data |
| Massachusetts DOT | https://massdot-impact-crashes-vhb.opendata.arcgis.com/search |
| Maryland DOT | https://mdsp.maryland.gov/Pages/Dashboards/CrashDataDownload.aspx |
| Nevada DOT | https://geohub-ndot.hub.arcgis.com/datasets/NDOT::crashdata-opendata/ |
| Montana DOT | https://www.mdt.mt.gov/publications/datastats/crashdata.aspx |
| | https://www.mdt.mt.gov/publications/docs/datastats/crashdata/PublicCrashData2019-2023.xlsx |
| Iowa DOT | https://icat.iowadot.gov/ |
| **Satellite Image** | |
| MapBox API | https://console.mapbox.com/ |
| **Road Networks** | |
| OSMnx Street Network Dataverse | https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/CUWWYJ |
| **Traffic Volumes** | |
| Delaware | https://de-firstmap-delaware.hub.arcgis.com/datasets/delaware::delaware-traffic-counts-2-0/ |
| Massachusetts | https://mhd.public.ms2soft.com/tcds/tsearch.asp?loc=Mhd&mod= |
| Maryland | https://data-maryland.opendata.arcgis.com/datasets/maryland::mdot-sha-annual-average-daily-traffic-aadt-locations/ |
| Nevada | https://geohub-ndot.hub.arcgis.com/datasets |
| Montana | https://mdt.public.ms2soft.com/tcds/tsearch.asp?loc=Mdt&mod= |
| Iowa | https://experience.arcgis.com/experience |
| **Weather Observations** | |
| Meteostat API | https://meteostat.net/en/ |
| **Supplementary Websites** | |
| Google Map API | https://maps.google.com/ |

**Table 11: We report the mean absolute error (MAE), precision, and recall score on the test split using node embedding, graph neural network embeddings, supervised contrastive learning, and feature fusion methods. To account for variability, each experiment is repeated with three different random seeds, and we report the averaged results along with standard deviations.**

| MAE / Average count | Delaware 4.59 | Massachusetts 7.25 | Maryland 1.72 | Nevada 1.29 | Montana 0.40 | Iowa 0.84 |
|---|---|---|---|---|---|---|
| GCN + Basic Fusion | 0.1 ± 0.00 | 0.4 ± 0.02 | 0.2 ± 0.01 | 0.1 ± 0.01 | 0.1 ± 0.02 | 0.3 ± 0.14 |
| GCN + Gated Fusion | 0.1 ± 0.01 | 0.5 ± 0.02 | 0.2 ± 0.00 | 0.1 ± 0.01 | 0.1 ± 0.01 | 0.1 ± 0.03 |
| GCN + MoE | 0.1 ± 0.01 | 0.3 ± 0.01 | 0.1 ± 0.02 | 0.1 ± 0.01 | 0.1 ± 0.01 | 0.1 ± 0.02 |
| **AUROC / Positive rate** | **0.32** | **0.28** | **0.24** | **0.14** | **0.10** | **0.18** |
| GCN + Basic Fusion | 88.5 ± 0.3 | 87.0 ± 0.1 | 89.2 ± 0.1 | 92.9 ± 0.1 | 80.0 ± 0.9 | 82.3 ± 0.6 |
| GCN + Gated Fusion | 88.5 ± 0.2 | 87.7 ± 0.1 | 89.3 ± 0.0 | 93.5 ± 0.1 | 83.7 ± 0.7 | 84.4 ± 0.5 |
| GCN + MoE | 88.7 ± 0.7 | 87.9 ± 0.3 | 88.3 ± 0.0 | 93.7 ± 0.1 | 83.0 ± 0.8 | 84.9 ± 0.4 |
| **Precision** | | | | | | |
| MLP | 6.33 ± 1.1 | 1.99 ± 0.8 | 2.49 ± 0.5 | 5.33 ± 0.6 | 0.75 ± 0.1 | 1.19 ± 0.1 |
| DeepWalk | 16.07 ± 3.9 | 4.19 ± 0.4 | 5.69 ± 0.8 | 6.17 ± 0.5 | 1.73 ± 0.1 | 3.04 ± 0.3 |
| GraphSAGE | 24.11 ± 3.4 | 10.06 ± 2.5 | 11.44 ± 0.2 | 5.24 ± 0.6 | 3.23 ± 0.7 | 5.94 ± 1.5 |
| SupConGCN | 17.21 ± 3.5 | 5.23 ± 0.3 | 4.84 ± 0.8 | 4.89 ± 0.7 | 2.28 ± 0.6 | 2.38 ± 0.4 |
| Graph Transformer | 26.13 ± 5.0 | 8.96 ± 0.9 | 11.25 ± 0.8 | 7.18 ± 1.1 | 0.64 ± 0.2 | 6.21 ± 1.1 |
| DCRNN | 38.08 ± 4.2 | 4.06 ± 0.1 | 5.08 ± 0.2 | 2.11 ± 0.5 | 2.63 ± 0.4 | 6.34 ± 1.2 |
| ViT | 16.84 ± 2.9 | 4.24 ± 0.8 | 5.21 ± 0.2 | 6.00 ± 0.5 | 1.71 ± 0.0 | 2.69 ± 0.3 |
| CLIP | 11.49 ± 2.5 | 3.56 ± 0.2 | 4.37 ± 0.1 | 5.02 ± 0.3 | 1.27 ± 0.1 | 2.55 ± 0.5 |
| GCN | 11.53 ± 0.6 | 6.27 ± 0.3 | 5.62 ± 0.2 | 6.09 ± 0.4 | 1.90 ± 0.3 | 8.66 ± 2.3 |
| GCN + Basic Fusion | 36.69 ± 4.5 | 4.25 ± 0.1 | 4.67 ± 0.3 | 5.67 ± 0.1 | 3.64 ± 0.3 | 2.32 ± 0.3 |
| GCN + Gated Fusion | 46.85 ± 4.3 | 4.17 ± 0.0 | 4.50 ± 0.4 | 5.41 ± 0.1 | 2.95 ± 0.3 | 3.83 ± 2.1 |
| GCN + MoE | 52.43 ± 3.9 | 5.96 ± 0.1 | 4.82 ± 0.5 | 6.71 ± 0.8 | 0.61 ± 0.2 | 3.07 ± 0.4 |
| GIN | 45.94 ± 5.2 | 5.05 ± 0.9 | 5.13 ± 0.2 | 8.19 ± 0.9 | 1.39 ± 0.0 | 2.62 ± 0.3 |
| GIN + Basic Fusion | 49.94 ± 3.3 | 6.97 ± 2.1 | 6.18 ± 0.1 | 5.98 ± 0.5 | 2.32 ± 0.4 | 2.64 ± 0.2 |
| GIN + Gated Fusion | 55.38 ± 3.5 | 5.68 ± 0.3 | 5.88 ± 0.2 | 5.32 ± 0.1 | 1.97 ± 0.5 | 3.15 ± 0.3 |
| GIN + MoE | 67.04 ± 4.1 | 5.16 ± 0.2 | 6.39 ± 1.0 | 6.41 ± 0.9 | 1.41 ± 0.2 | 2.67 ± 0.2 |
| **Recall** | | | | | | |
| MLP | 70.35 ± 6.0 | 53.46 ± 7.7 | 87.66 ± 9.7 | 90.76 ± 4.5 | 79.59 ± 4.3 | 69.57 ± 2.2 |
| DeepWalk | 88.76 ± 1.9 | 81.24 ± 3.4 | 87.55 ± 5.4 | 90.82 ± 2.7 | 67.36 ± 2.4 | 80.99 ± 1.5 |
| GraphSAGE | 44.27 ± 6.1 | 67.49 ± 1.2 | 69.89 ± 7.5 | 80.17 ± 1.5 | 79.47 ± 4.2 | 67.34 ± 4.7 |
| SupConGCN | 76.23 ± 7.0 | 60.85 ± 3.1 | 80.68 ± 3.5 | 87.23 ± 0.8 | 77.98 ± 1.3 | 76.60 ± 1.9 |
| Graph Transformer | 60.39 ± 2.9 | 63.08 ± 3.8 | 69.76 ± 3.7 | 77.35 ± 5.0 | 79.68 ± 5.1 | 77.75 ± 6.5 |
| DCRNN | 79.80 ± 4.3 | 79.02 ± 2.7 | 81.26 ± 0.9 | 87.96 ± 2.3 | 78.34 ± 3.2 | 76.55 ± 2.8 |
| ViT | 92.23 ± 5.4 | 78.99 ± 1.1 | 82.06 ± 1.7 | 90.94 ± 3.1 | 66.41 ± 1.7 | 80.28 ± 2.5 |
| CLIP | 74.40 ± 3.7 | 75.84 ± 3.3 | 82.81 ± 0.7 | 88.62 ± 0.9 | 86.69 ± 1.9 | 76.56 ± 1.5 |
| GCN | 79.47 ± 5.4 | 75.41 ± 1.5 | 78.15 ± 1.2 | 74.80 ± 2.0 | 64.31 ± 2.1 | 58.55 ± 0.5 |
| GCN + Basic Fusion | 85.87 ± 1.3 | 85.24 ± 0.8 | 88.50 ± 0.3 | 87.38 ± 1.0 | 55.01 ± 4.9 | 93.15 ± 4.0 |
| GCN + Gated Fusion | 74.28 ± 3.9 | 84.79 ± 1.4 | 89.58 ± 1.0 | 90.04 ± 0.8 | 98.26 ± 1.0 | 73.12 ± 3.6 |
| GCN + MoE | 79.10 ± 4.5 | 66.88 ± 4.8 | 83.20 ± 1.5 | 87.73 ± 1.7 | 99.76 ± 0.2 | 91.85 ± 1.7 |
| GIN | 93.88 ± 0.3 | 76.69 ± 2.5 | 80.10 ± 0.8 | 84.85 ± 2.5 | 78.86 ± 3.2 | 77.01 ± 0.9 |
| GIN + Basic Fusion | 74.95 ± 3.6 | 68.03 ± 2.5 | 76.99 ± 3.0 | 89.79 ± 0.2 | 56.92 ± 0.4 | 82.50 ± 1.6 |
| GIN + Gated Fusion | 80.00 ± 3.3 | 75.59 ± 4.8 | 80.56 ± 1.4 | 89.75 ± 0.1 | 97.87 ± 1.4 | 81.01 ± 3.3 |
| GIN + MoE | 85.10 ± 4.6 | 72.08 ± 3.3 | 80.57 ± 1.1 | 90.56 ± 0.5 | 91.15 ± 0.7 | 83.55 ± 1.9 |
| **$F_1$-score** | | | | | | |
| MLP | 10.86 ± 2.2 | 2.86 ± 0.2 | 3.94 ± 0.6 | 9.02 ± 0.8 | 1.50 ± 0.1 | 2.31 ± 0.3 |
| DeepWalk | 25.86 ± 5.2 | 7.05 ± 1.0 | 10.46 ± 1.3 | 11.46 ± 0.9 | 3.35 ± 0.1 | 5.72 ± 0.6 |
| ViT | 30.93 ± 2.5 | 6.69 ± 0.5 | 9.67 ± 0.3 | 10.00 ± 0.1 | 3.32 ± 0.1 | 5.14 ± 0.5 |
| CLIP | 19.89 ± 2.3 | 6.83 ± 0.1 | 8.26 ± 0.4 | 9.45 ± 0.4 | 2.49 ± 0.1 | 4.87 ± 0.9 |
| GraphSAGE | 31.57 ± 2.2 | 8.17 ± 0.7 | 8.17 ± 0.7 | 9.74 ± 1.0 | 4.34 ± 1.3 | 5.89 ± 0.6 |
| SupCon | 26.49 ± 9.6 | 8.58 ± 0.2 | 8.90 ± 1.4 | 9.13 ± 1.2 | 4.28 ± 1.0 | 4.56 ± 0.8 |
| GCN | 40.50 ± 1.0 | 7.43 ± 0.5 | 10.43 ± 0.3 | 11.04 ± 0.6 | 3.33 ± 0.4 | 5.21 ± 0.5 |
| GCN + Basic Fusion | 45.42 ± 4.4 | 9.58 ± 0.2 | 10.46 ± 1.3 | 11.41 ± 0.8 | 4.96 ± 1.1 | 5.33 ± 0.5 |
| GCN + Gated Fusion | 50.27 ± 5.2 | 9.92 ± 0.0 | 11.07 ± 0.8 | 11.28 ± 0.1 | 5.03 ± 0.2 | 6.12 ± 0.4 |
| GCN + MoE | 56.78 ± 5.7 | 10.51 ± 0.1 | 11.99 ± 0.7 | 12.38 ± 1.4 | 5.21 ± 0.3 | 6.68 ± 0.6 |
| GIN | 64.01 ± 4.3 | 9.77 ± 2.1 | 9.53 ± 0.5 | 14.78 ± 1.4 | 2.72 ± 0.0 | 5.04 ± 0.5 |
| GIN + Basic Fusion | 70.04 ± 7.5 | 11.11 ± 1.8 | 11.27 ± 0.1 | 15.15 ± 0.8 | 4.39 ± 0.7 | 5.17 ± 0.3 |
| GIN + Gated Fusion | 73.78 ± 4.7 | 10.23 ± 0.3 | 10.87 ± 0.3 | 15.96 ± 0.4 | 4.10 ± 0.9 | 6.00 ± 0.6 |
| GIN + MoE | 76.27 ± 3.8 | 10.73 ± 0.1 | 11.67 ± 1.6 | 15.91 ± 1.7 | 4.74 ± 0.4 | 6.14 ± 0.3 |

**(a) Delaware**

**(b) Massachusetts**

**(c) Maryland**

**(d) Nevada**

**(e) Montana**

**(f) Iowa**

**Figure 9: We showcase representative satellite images sampled from each of the six states included in our dataset. Each state is associated with six figures. These images exhibit diverse geographical and urban characteristics, ranging from dense urban intersections and suburban roadways to rural highways and mountainous terrains. Such variation reflects the heterogeneity of real-world driving environments across regions and provides rich visual cues that are critical for learning robust, transferable vision-based road features. These satellite views serve as the primary modality in our framework for capturing road layout, surrounding context, and surface-level conditions.**