

# UC Irvine

## UC Irvine Previously Published Works

### Title

Scene-Graph Augmented Data-Driven Risk Assessment of Autonomous Vehicle Decisions

### Permalink

<https://escholarship.org/uc/item/54h8b19h>

### Journal

IEEE Transactions on Intelligent Transportation Systems, 23(7)

### ISSN

1524-9050

### Authors

Yu, Shih-Yuan

Malawade, Arnav Vaibhav

Muthirayan, Deepan

et al.

### Publication Date

2022-07-01

### DOI

10.1109/tits.2021.3074854

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

# Scene-Graph Augmented Data-Driven Risk Assessment of Autonomous Vehicle Decisions

Shih-Yuan Yu, Arnav V. Malawade, Deepan Muthirayan, Pramod P. Khargonekar, Mohammad A. Al Faruque  
 Department of Electrical Engineering & Computer Science, University of California, Irvine CA, USA.  
 {shihyuay,malawada,dmuthira,pramod.khargonekar,alfaruqu}@uci.edu

**Abstract**—Despite impressive advancements in Autonomous Driving Systems (ADS), navigation in complex road conditions remains a challenging problem. There is considerable evidence that evaluating the subjective risk level of various decisions can improve ADS’ safety in both normal and complex driving scenarios. However, existing deep learning-based methods often fail to model the relationships between traffic participants and can suffer when faced with complex real-world scenarios. Besides, these methods lack *transferability* and *explainability*. To address these limitations, we propose a novel data-driven approach that uses *scene-graphs* as intermediate representations. Our approach includes a Multi-Relation Graph Convolution Network, a Long-Short Term Memory Network, and attention layers for modeling the subjective risk of driving maneuvers. To train our model, we formulate this task as a supervised scene classification problem. We consider a typical use case to demonstrate our model’s capabilities: lane changes. We show that our approach achieves a higher classification accuracy than the state-of-the-art approach on both large (96.4% vs. 91.2%) and small (91.8% vs. 71.2%) synthesized datasets, also illustrating that our approach can learn effectively even from smaller datasets. We also show that our model trained on a synthesized dataset achieves an average accuracy of 87.8% when tested on a real-world dataset compared to the 70.3% accuracy achieved by the state-of-the-art model trained on the same synthesized dataset, showing that our approach can more effectively transfer knowledge. Finally, we demonstrate that the use of spatial and temporal attention layers improves our model’s performance by 2.7% and 0.7% respectively and increases its explainability.

**Index Terms**—Autonomous Vehicle, Risk Assessment, Scene Understanding, Graph Convolutional Neural Network.

## I. INTRODUCTION

Autonomous Driving Systems (ADSs) have advanced significantly in recent years. However, navigation is still challenging in complex urban environments since the scenarios are highly variable and complex [1], [2], [3]. The continued reports of autonomous vehicle crash only highlight these challenges [4], [5], [6], [7]. A risk-based approach for autonomous driving has the potential to address this challenge and better assure driving safety. Within this context, the effectiveness of understanding the driving scenes and quantifying the risk of driving decisions becomes particularly crucial for ADSs.

In prior research, the problem of risk assessment for autonomous driving has been tackled by modeling either the *objective risk* or the *subjective risk* [8], [9], [10]. The *objective risk* is defined as the objective probability of an accident occurring and is usually determined by statistical analysis [8]. Some works have focused on minimizing the *objective risk* by

modeling the trajectories of vehicles [11], [12] to guarantee safe driving.

*Subjective risk* refers to the driver’s own perceived risk and is an output of the driver’s cognitive process [9], [10]. Studies suggest that modeling the *subjective risk* can be a better option as drivers typically avoid taking a risky driving action based on their subjective perception [10]. For this reason, our objective in this work is to build a model for subjective risk assessment.

Several papers have leveraged state-of-the-art deep learning architectures for modeling subjective risk [2], [13]. Such methods typically apply Convolutional Neural Networks (CNNs) and Long-Short Term Memory Networks (LSTMs) and have been proven to be effective at capturing features essential for modeling subjective risk in both spatial and temporal domains [13]. However, it is unclear whether these methods can capture critical higher-level information, such as the relationships between the traffic participants in a given scene. Failure in capturing these relationships can result in poor ADS performance in complex scenarios.

Also, training these networks requires large datasets covering a wide range of “corner cases” (especially risky driving scenarios), which are expensive and time-consuming to generate [14]. Many researchers resort to using synthesized datasets containing many examples of these corner cases to address this issue. However, for these to be valuable, a model must be able to transfer the knowledge gained from simulated training data to real-world situations. A standard method for measuring a model’s ability to generalize is *transferability*, where a model’s accuracy on a dataset different from the training dataset is evaluated. If a model can transfer the knowledge gained from a simulated training set to a real-world testing set effectively, it is likely that it will perform better in unseen real-world scenarios.

Even if these existing methods can transfer knowledge well, the predictions of such methods lack *explainability*, which is crucial for establishing trust between ADSs and human drivers [10], [15], [16]. *Explainability* refers to the ability of a model to effectively communicate the factors that influenced its decision-making process for a given input, particularly those that might lead the model to make incorrect decisions [16], [17]. If a model can give attention to the aspects or entities in a traffic scene that make the scenario risky or non-risky, it can improve its decision and its decisions become more explainable [18].

Overall, designing a risk assessment system for ADSs using data-driven approaches presents the following challenges:

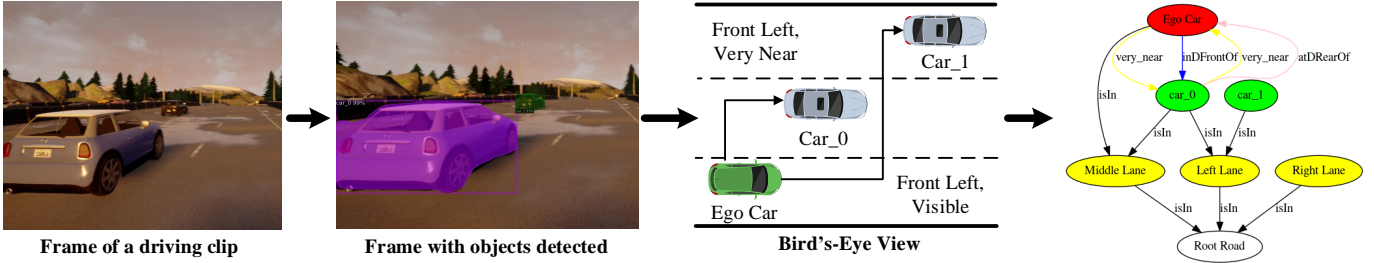


Fig. 1: An illustration of *scene-graph* extraction using the *Real Image Pipeline*. In this process, the first step is to detect a list of objects on each frame of a clip. Then, we project each frame to its bird’s-eye view to better approximate the spatial relations between objects. Finally, we construct a *scene-graph* using the list of detected objects and their attributes.

- Designing a reliable method that can handle a wide range of complex and unpredictable traffic scenarios.
- Building a model that is transferable from the simulation setting to the real-world setting because the real-world datasets for supervised training are limited.
- Building a model that can provide explainable decisions.

In this paper, we present our work on the modeling of *subjective risk* for a specific driving maneuver: lane change. This task by itself is crucial, given that 7.62% of all traffic accidents between light vehicles can be attributed to improper execution of lane change [19]. We propose a *scene-graph* augmented data-driven approach for assessing the subjective risk of driving maneuvers, where the *scene-graphs* serve as intermediate representations (IR) as shown in Figure 1. The key advantage is that using *scene-graphs* as IRs allows us to model the relationships between the participants in a traffic scene, thus potentially improving the model’s understanding of a scene.

Our architecture consists of three major components: (i) a pipeline to convert the images of a driving clip to a sequence of *scene-graphs*, (ii) a Multi-Relational Graph Convolution Network (MR-GCN) to convert each of the *scene-graphs* to an embedding (a vectorized representation), and (iii) an LSTM for temporally modeling the sequence of embeddings of the respective *scene-graphs*. Our model also contains multiple attention layers: (i) a node attention layer before the embedding of a *scene-graph* is computed, and (ii) an attention layer on top of the LSTM, both of which can further improve its performance and the explainability. For training the model, we propose formulating the problem of subjective risk assessment as a supervised *scene-graph* sequence classification problem.

Our key contributions are as follows:

- We present a novel *scene-graph* augmented data-driven approach for assessing the risk of driving actions in autonomous vehicles.
- We demonstrate that our approach outperforms existing methods at risk assessment across a wide range of scenarios using lane change as a use case.
- We demonstrate that the use of multi-level attention in our proposed approach provides better explainability.
- We demonstrate that our *scene-graph* based approach can better transfer knowledge gained from simulated environments to real-world risk assessment tasks.

The rest of the paper is structured as follows: In Section II, we discuss related works. In Section III we introduce our *scene-graph* augmented approach. In Section IV we discuss our experimental results. Finally, in Section V we present our conclusions.

## II. RELATED WORK

### A. Design Philosophies in ADSs

Two broad approaches for designing ADSs are (i) modular design, (ii) end-to-end design [2]. Most modular approaches comprise a pipeline of separate components from the sensory inputs to the actuator outputs, while end-to-end approaches generate output directly from their sensory inputs [20], [21]. One advantage of a modular design approach is the division of a task into an easier-to-solve set of sub-tasks that have been addressed in other fields such as robotics [22], computer vision [23] and vehicle dynamics [24]. Therefore, prior knowledge from these fields can be leveraged when designing the components corresponding to the sub-tasks. However, one disadvantage of such an approach is the complexity of the whole pipeline [2]. End-to-end approaches can achieve good performance with a smaller network size because they perform feature extraction from sensor inputs implicitly through the network’s hidden layers [21]. However, the authors in [25] point out that the needed level of supervision is too weak for the end-to-end model to learn critical controlling information (e.g., from image to steering angle), so it can fail to handle complicated driving maneuvers.

A third approach was first proposed by DeepDriving [25], called the *direct perception* approach. In their approach, a set of *affordance indicators*, such as the distance to the lane markings and the distance to cars in the current and adjacent lanes, are extracted from an image and serves as an IR for generating the final control output. They prove that the use of this IR is effective for simple driving tasks such as lane following and for generalizing to real-world environments. Authors in [26] use a collection of filtered images, each representing a piece of distinct information, as the IR. They state that the IR used in their approach allows the training to be conducted on real or simulated data, facilitating testing and validation in simulations before testing on a real car. Moreover, they show that it is easier to synthesize perturbations to the driving trajectory at the mid-level representations than at the

level of raw sensors, enabling them to produce non-expert behaviors such as off-road driving and collisions. The authors in [13] use Mask-RCNN [27] to color the vehicles in each input image, producing a form of IR. In contrast to the works mentioned above, our approach uses a *scene-graph* IR that encodes the spatial and semantic relations between all the traffic participants in a frame.

### B. Graph-Based Driving Scene Understanding

Several papers have applied graph-based formulations for driving scene understanding [28], [29], [30]. In [29], the authors propose a 3D-aware egocentric Spatio-temporal interaction framework that uses both an *Ego-Thing* graph and an *Ego-Stuff* graph, which encode how the ego vehicle interacts with both moving and stationary objects in a scene, respectively. In [28], [30], the authors propose a pipeline using a multi-relational graph convolutional network (MR-GCN) for classifying the driving behaviors of traffic participants. The MR-GCN is constructed by combining spatial and temporal information, including relational information between moving objects and landmark objects. Our work is primarily inspired by [28], [30] but differs in the application and network architecture. These papers focus on predicting each object’s behavior in the *scene-graph* while our work focuses on assessing the subjective risk of the entire scenario. Consequently, we propose a network architecture that implements more components such as node-attention, graph pooling layers, and readout operations.

### C. Risk Assessment

Several works have studied subjective risk assessment for autonomous driving systems [9], [10], [13], [31]. In [31], Hidden Markov Models (HMMs) and Language Models are used to detect unsafe lane change events. The approach taken in [13] is the most related to our work as it infers the risk-level of overall driving scenes with a deep Spatio-temporal neural network architecture. By using Mask-RCNN [27] to generate an IR for each image, their approach achieves a 3% performance gain in risk assessment. They show that the architecture with Semantic Mask Transfer (SMT) + CNN + LSTM can perform 25% better than the architecture with Feature Transfer (FT) + Frame-by-Frame (FbF), indicating that capturing the spatial and temporal features just from a monocular camera can be useful in modeling subjective risk. However, this approach only considers the spatial features (the latent vector output of the CNN layers) of a frame instead of the relations between all the traffic participants. Our work uses *scene-graphs* as IRs to capture the high-level relationships between all the traffic participants of a scene.

## III. SCENE-GRAPH AUGMENTED APPROACH FOR RISK ASSESSMENT

In this section, we discuss our proposed approach for *scene-graph* augmented risk assessment. In our work, we make the same assumption used in [13] that the set of driving sequences can be partitioned into two jointly exhaustive and mutually

exclusive subsets: risky and safe. We denote the sequence of images of length  $T$  by  $\mathbf{I} = \{I_1, I_2, I_3, \dots, I_T\}$ . We assume the existence of a spatio-temporal function  $f$  that outputs whether a sequence of driving actions  $x$  is safe or risky via a risk label  $y$ , as given in Equation 1.

$$y = f(\mathbf{I}) = f(\{I_1, I_2, I_3, \dots, I_{T-1}, I_T\}), \quad (1)$$

where

$$y = \begin{cases} (1, 0), & \text{if the driving sequence is safe} \\ (0, 1), & \text{if the driving sequence is risky.} \end{cases} \quad (2)$$

In this section, we propose a suitable model for approximating the function  $f$ . In the model we propose, the first step is the extraction of the *scene-graph*  $G_t$  from each image  $I_t$  of the video clip  $\mathbf{I}$ . This is achieved by a series of processes that we collectively call the *Scene-Graph Extraction Pipeline*, described in Section III-A. In the second step, these *scene-graphs* are passed through graph convolution layers and an attention-based graph pooling layer. The graph-level embeddings of each *scene-graph*,  $\mathbf{h}_{G_t}$ , are then calculated using a graph readout operation. Next, these *scene-graph* embeddings are passed sequentially to LSTM cells to acquire the *spatio-temporal* representation, denoted as  $\mathbf{Z}$ , of each *scene-graph* sequence. Lastly, we use a Multi-Layer Perceptron (MLP) layer with a *Softmax* activation function to acquire the final inference, denoted as  $\hat{y}$ , of the risk for each driving sequence  $\mathbf{I}$ . We describe more details regarding each of these components of our model in Section III-B.

### A. Scene-Graph Extraction Pipeline

Several approaches have been proposed for extracting *scene-graphs* from images by detecting the objects in a scene and then identifying their visual relationships [32], [33]. However, these works have focused on single general images instead of a sequence of images as it arises in autonomous driving, where higher accuracy is demanded. Thus, we adopted a partially rule-based process to extract objects and their attributes from images called the *Real Image Pipeline*. Besides, to evaluate how our approach performs with *scene-graphs* containing ground truth information, we use the *Carla Ground Truth (GT) Pipeline* as a surrogate for the ideal situation where the attributes for each object can be correctly extracted. After the objects in a scene and their attributes have been extracted, the *scene-graphs* are constructed as described in III-A3. We discuss these two approaches in detail below.

1) *Real Image Pipeline*: In this pipeline, object attributes and bounding boxes are extracted directly from images using state-of-the-art image processing techniques. As Figure 1 shows, we first convert each image  $I_t$  into a collection of objects  $O_t$  using Detectron2, a state-of-the-art object detection model based on Faster RCNN [34], [35]. Next, we use OpenCV’s perspective transformation library to generate a top-down perspective of the image, commonly known as a “birds-eye view” projection [36]. This projection allows us to approximate each object’s location relative to the road markings and the ego vehicle. Next, for each detected object in  $O_t$ , we use its estimated location and class type (cars, motorcycles, pedestrians, lanes, etc.) to compute the attributes required in building the *scene-graph*.



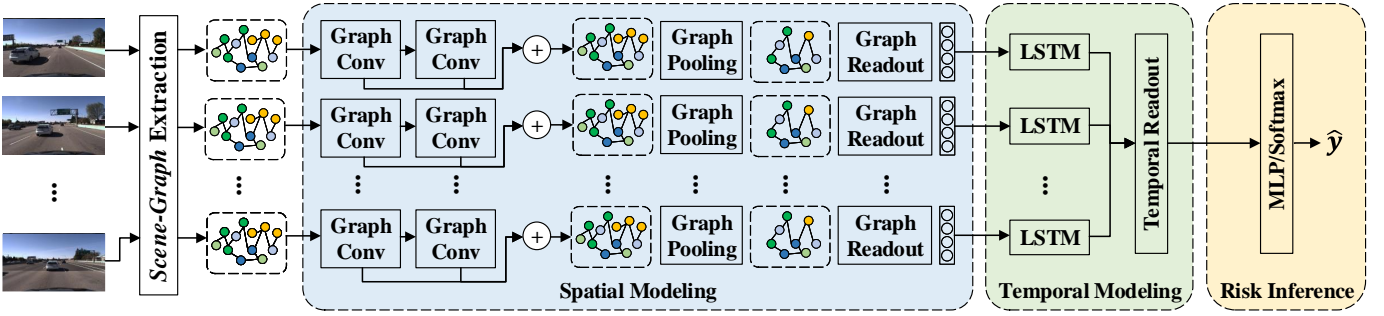


Fig. 2: An illustration of our model’s architecture. It first converts each image  $I_t \in \mathbf{I}$  to a *scene-graph*  $G_t$  via the *Scene-Graph Extraction Pipeline*. Next, it converts each  $G_t$  to its corresponding *scene-graph* embedding  $\mathbf{h}_{G_t}$  with layers in *Spatial Modeling*. Then, it temporally models these *scene-graph* embeddings to acquire the spatio-temporal representation  $\mathbf{Z}$  for a *scene-graph* sequence. Finally, the risk inference  $\hat{y}$  of a clip is calculated from  $\mathbf{Z}$  using an MLP with a *Softmax* activation function.

2) *Carla Ground Truth Pipeline*: Object detection and location estimation with solely monocular camera can be unstable because of factors such as weather and camera position [37], which can impact the correctness of our image-based *scene-graph* construction pipeline and thus our approach’s performance. To evaluate our methodology under the assumption that object attributes can be extracted without error, we build our *scene-graphs* using the ground-truth location and class information for each vehicle in the *Carla GT Pipeline*. We extract this information directly from Carla simulator [14] without using any image processing steps.

3) *Graph Construction*: After collecting the list of objects in each image and their attributes, we begin constructing the corresponding *scene-graphs* as follows. For each image  $I_t$ , we denote the corresponding *scene-graph* by  $G_t = \{O_t, A_t\}$  and model it as a directed multi-graph where multiple types of edges connect nodes. The nodes of a *scene-graph*, denoted as  $O_t$ , represent the objects in a scene such as lanes, roads, traffic signs, vehicles, pedestrians, etc. The edges of a *scene-graph* are represented by the corresponding adjacency matrix  $A_t$ , where each value in  $A_t$  represents the type of the edges. The edges between two nodes represent the different kinds of relations between them (e.g., near, Front\_Left, isIn, etc.).

In assessing the risk of driving behaviors, traffic participants’ relations that we consider to be useful are the distance relations and the directional relations. The assumption made here is that the local proximity and positional information of one object will influence the other’s motion only if they are within a certain distance. Therefore, in this work, we extract only the location information for each object and adopt a simple rule to determine the relations between the objects using their attributes (e.g., relative location to the ego car), as shown in Figure 1. For distance relations, we assume two objects are related by one of the relations  $r \in \{Near\_Collision (4 \text{ ft.}), Super\_Near (7 \text{ ft.}), Very\_Near (10 \text{ ft.}), Near (16 \text{ ft.}), Visible (25 \text{ ft.})\}$  if the objects are physically separated by a distance that is within that relation’s threshold. In the case of the directional relations, we assume two objects are related by the relation  $r \in \{Front\_Left, Left\_Front, Left\_Rear, Rear\_Left, Rear\_Right, Right\_Rear, Right\_Front, Front\_Right\}$  based on their relative positions if they are within the *Near* threshold distance from one another.

In addition to directional and distance relations, we also implement the *isIn* relation that connects vehicles with their respective lanes. For the *Carla GT Pipeline*, we extract the ground-truth lane assignments for each vehicle from the simulator directly. For the *Real Image Pipeline*, we use each vehicle’s horizontal displacement relative to the ego vehicle to assign vehicles to either the *Left Lane*, *Middle Lane*, or *Right Lane* based on a known lane width. Our abstraction only includes these three-lane areas and, as such, we map vehicles in all left lanes to the same *Left Lane* node and all vehicles in right lanes to the *Right Lane* node. If a vehicle overlaps two lanes (i.e., during a lane change), we assign it an *isIn* relation to both lanes. Figure 1 illustrates an example of resultant *scene-graph*.

### B. Scene-graph Sequence Based Risk Assessment Model

The model we propose consists of three major components: spatial model, temporal model, and risk inference. The spatial model outputs the embedding  $h_{G_t}$  for each *scene-graph*  $G_t$ . The temporal model processes the sequence of *scene-graph* embeddings  $\mathbf{h}_I = \{h_{G_1}, h_{G_2}, \dots, h_{G_T}\}$  and produces the spatio-temporal embedding  $\mathbf{Z}$ . The risk inference component outputs each driving clip’s final risk assessment, denoted as  $\hat{Y}$ , by processing the Spatio-temporal embedding  $\mathbf{Z}$ . The overall network architecture is shown in Figure 2. We discuss each of these components in detail below.

1) *Spatial Modeling*: The spatial model we propose uses MR-GCN layers to compute the embedding for a *scene-graph*. The use of MR-GCN allows us to capture multiple types of relations on each *scene-graph*  $G_t = \{O_t, A_t\}$ . In the *Message Propagation* phase, a collection of node embeddings and their adjacency information serve as the inputs to the MR-GCN layer. Specifically, the  $l$ -th MR-GCN layer updates the node embedding, denoted as  $\mathbf{h}_v^{(l)}$ , for each node  $v$  as follows:

$$\mathbf{h}_v^{(l)} = \Phi_0 \cdot \mathbf{h}_v^{(l-1)} + \sum_{r \in A_t} \sum_{u \in N_r(v)} \frac{1}{|N_r(v)|} \Phi_r \cdot \mathbf{h}_u^{(l-1)}, \quad (3)$$

where  $N_r(v)$  denotes the set of neighbor indices of node  $v$  with the relation  $r \in A_t$ .  $\Phi_r$  is a trainable relation-specific transformation for relation  $r$  in MR-GCN layer. Since the information in  $(l - 1)$ -th layer can directly influence the

representation of the node at  $l$ -th layer, MR-GCN uses another trainable transformation  $\Phi_0$  to account for the self-connection of each node using a special relation [38]. Here, we initialize each node embedding  $\mathbf{h}_v^{(0)}$ ,  $\forall v \in O_t$ , by directly converting the node's type information to its corresponding one-hot vector.

Typically, the node embedding becomes more refined and global as the number of graph convolutional layers,  $L$ , increases. However, the authors in [39] also suggest that the features generated in earlier iterations might generalize the learning better. Therefore, we consider the node embeddings generated from all the MR-GCN layers. To be more specific, we calculate the embedding of node  $v$  at the final layer, denoted as  $\mathbf{H}_v^L$ , by concatenating the features generated from all the MR-GCN layers, as follows,

$$\mathbf{H}_v^L = \text{CONCAT}(\{\mathbf{h}_v^{(l)}\} | l = 0, 1, \dots, L). \quad (4)$$

We denote the collection of node embeddings of *scene-graph*  $G_t$  after passing through  $L$  layers of MR-GCN as  $\mathbf{X}_t^{\text{prop}}$  ( $L$  can be 1, 2 or 3).

The node embedding  $\mathbf{X}_t^{\text{prop}}$  is further processed with an attention-based graph pooling layer. As stated in [17], such an attention-based pooling layer can improve the explainability of predictions and is typically considered as a part of a unified computational block of a graph neural network (GNN) pipeline. In this layer, nodes are pooled according to the scores predicted from either a trainable simple linear projection [40] or a separate trainable GNN layer [41]. We denote the graph pooling layer that uses the **SCORE** function in [40] as *TopkPool* and the one that uses the **SCORE** function in [41] as *SAGPool*. The calculation of the overall process is presented as follows:

$$\alpha = \text{SCORE}(\mathbf{X}_t^{\text{prop}}, \mathbf{A}_t), \quad (5)$$

$$\mathbf{P} = \text{top}_k(\alpha), \quad (6)$$

where  $\alpha$  stands for the coefficients predicted by the graph pooling layer for nodes in  $G_t$  and  $\mathbf{P}$  represents the indices of the pooled nodes which are selected from the top  $k$  of the nodes ranked according to  $\alpha$ . The number  $k$  of the nodes to be pooled is calculated by a pre-defined pooling ratio,  $pr$ , and using  $k = pr \times |O_t|$ , where we consider only a constant fraction  $pr$  of the embeddings of the nodes of a scene-graph to be relevant (i.e., 0.25, 0.5, 0.75). We denote the node embeddings and edge adjacency information after pooling by  $\mathbf{X}_t^{\text{pool}}$  and  $\mathbf{A}_t^{\text{pool}}$  and are calculated as follows:

$$\mathbf{X}_t^{\text{pool}} = (\mathbf{X}_t^{\text{prop}} \odot \tanh(\alpha))_{\mathbf{P}}, \quad (7)$$

$$\mathbf{A}_t^{\text{pool}} = \mathbf{A}_t^{\text{prop}}_{(\mathbf{P}, \mathbf{P})}. \quad (8)$$

where  $\odot$  represents an element-wise multiplication,  $()_{\mathbf{P}}$  refers to the operation that extracts a subset of nodes based on  $\mathbf{P}$  and  $()_{(\mathbf{P}, \mathbf{P})}$  refers to the formation of the adjacency matrix between the nodes in this subset.

Finally, our model aggregates the node embeddings of the graph pooling layer,  $\mathbf{X}_t^{\text{pool}}$ , using a graph **READOUT** operation, to produce the final graph-level embedding  $\mathbf{h}_{G_t}$  for each *scene-graph*  $G_t$  as given by

$$\mathbf{h}_{G_t} = \text{READOUT}(\mathbf{X}_t^{\text{pool}}), \quad (9)$$

where the **READOUT** operation can be either summation, averaging, or selecting the maximum of each feature dimension, over all the node embeddings, known as *sum-pooling*, *mean-pooling*, or *max-pooling*, respectively. The process until this point is repeated across all images in  $\mathbf{I}$  to produce the sequence of embedding,  $\mathbf{h}_I$ .

2) *Temporal Modeling*: The temporal model we propose uses an LSTM for converting the sequence of scene-graph embeddings  $\mathbf{h}_I$  to the combined spatio-temporal embedding  $\mathbf{Z}$ . For each timestamp  $t$ , the LSTM updates the hidden state  $p_t$  and cell state  $c_t$  as follows,

$$p_t, c_t = \text{LSTM}(\mathbf{h}_{G_t}, c_{t-1}), \quad (10)$$

where  $\mathbf{h}_{G_t}$  is the final *scene-graph* embedding from timestamp  $t$ . After the LSTM processes all the scene-graph embeddings, a temporal readout operation is applied to the resultant output sequence to compute the final Spatio-temporal embedding  $\mathbf{Z}$  given by

$$\mathbf{Z} = \text{TEMPORAL\_READOUT}(p_1, p_2, \dots, p_T) \quad (11)$$

where the **TEMPORAL\_READOUT** operation could be extracting only the last hidden state  $p_T$  (LSTM-last), or be a temporal attention layer (LSTM-attn).

In [15], adding an attention layer  $b$  to the encoder-decoder based LSTM architecture is shown to achieve better performance in Neural Machine Translation (NMT) tasks. For the same reason, we include *LSTM-attn* in our architecture. *LSTM-attn* calculates a context vector  $q$  using the hidden state sequence  $\{p_1, p_2, \dots, p_T\}$  returned from the LSTM encoder layer as given by

$$q = \sum_{t=1}^T \beta_t p_t \quad (12)$$

where the probability  $\beta_t$  reflects the importance of  $p_t$  in generating  $q$ . The probability  $\beta_t$  is computed by a *Softmax* output of an energy function vector  $e$ , whose component  $e_t$  is the energy corresponding to  $p_t$ . Thus, the probability  $\beta_t$  is formally given by

$$\beta_t = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)}, \quad (13)$$

where the energy  $e_t$  associated with  $p_t$  is given by  $e_t = b(s_0, p_t)$ . The temporal attention layer  $b$  scores the importance of the hidden state  $p_t$  to the final output, which in our case is the risk assessment. The variable  $s_0$  in the temporal attention layer  $b$  is computed from the last hidden representation  $p_T$ . The final Spatio-temporal embedding for a video clip,  $\mathbf{Z}$ , is computed by feeding the context vector  $q$  to another LSTM decoder layer.

3) *Risk Inference*: The last piece of our model is the risk inference component that computes the risk assessment prediction  $\hat{Y}$  using the spatio-temporal embedding  $\mathbf{Z}$ . This component is composed of a MLP layer followed by a *Softmax* activation function. Thus, the prediction  $\hat{Y}$  is given by

$$\hat{Y} = \text{Softmax}(\text{MLP}(\mathbf{Z})) \quad (14)$$

The loss for the prediction is calculated as follows,

$$\arg \min \text{CrossEntropyLoss}(Y, \hat{Y}) \quad (15)$$

For training our model, we use a mini-batch gradient descent algorithm that updates its parameters by training on a batch of *scene-graph* sequences. To account for label imbalance, we apply class weighting when calculating loss. Besides, several dropout layers are inserted into the network to reduce overfitting.

#### IV. EXPERIMENTAL RESULTS

In this section, we provide extensive experimental results to illustrate the accuracy of our model and the ability of our model to transfer knowledge (transferability). We do this by providing comparisons between our model and a state-of-the-art SMT+CNN+LSTM based risk assessment model [13]; we refer to this model as the *Baseline*. Besides, we provide results for our model’s best hyper-parameter setting and perform an ablation study to evaluate the contribution of each major component in our model.

##### A. Dataset Preparation

We prepared two types of datasets for our experiments (i) real-world driving datasets and (ii) synthesized datasets. We generated the real-world dataset by extracting lane change clips from the Honda Driving Dataset (HDD) [42]. To create the synthesized datasets, we developed a tool to generate lane changing clips using the Carla<sup>1</sup> and Carla Scenario Runner<sup>2</sup>.

Carla is an open-source driving simulator [14] that allows users to control a vehicle in either manual mode or autopilot mode. The Carla Scenario Runner contains a set of atomic controllers that enable users to control a car in a driving scene and perform complex driving maneuvers. We modified the user script in Carla so that it can (i) select one autonomous car randomly and switch its mode to manual mode, and then (ii) utilize Scenario Runner’s function to force the vehicle to change lanes.

The data generating tool allows us to fabricate lane changing clips directly instead of extracting them from long driving clips. We generated a wide range of simulated lane changes using the various presets in Carla that allowed us to specify the number of cars, pedestrians, weather and lighting conditions, etc. Also, through the APIs provided by the Traffic Manager (TM) of the Carla simulator, we were able to customize the driving characteristics of every autonomous vehicle, such as the intended speed considering the current speed limit, the chance of ignoring the traffic lights, or the chance of neglecting collisions with other vehicles. Overall, this allowed us to simulate a wide range of very realistic urban driving environments and generate synthesized datasets suitable for training and testing a model.

To label the lane change clips in both the real-world and synthesized datasets, we performed an annotation process similar to the one used in [13]. First, in this process, human annotators were asked to assign a risk score to each clip that ranges from -2 and 2, where 2 implies a highly risky lane change and -2 means the safest lane change. Then, for each

lane change clip, all the risk labels of annotators were averaged and converted to the binary label  $y$  as follows: if the average is  $\leq 0$ , then assign  $y = 0$  (safe); else assign  $y = 1$  (risky).

We generated two synthesized datasets: a *271-syn* dataset and a *1043-syn* dataset, containing 271 and 1,043 lane changing clips, respectively. In addition, we sub-sampled the *271-syn* and *1043-syn* datasets further to create two balanced datasets that have a 1:1 distribution of risky to safe lane changes: *96-syn* and *306-syn*. We call the real-driving dataset as *571-honda* as it contains 571 lane changing video clips.

We randomly split each dataset into a training set and a testing set by the ratio 7:3 such that the split is stratified, i.e., the proportion of risky to safe lane change clips in each of the splits is the same. The models are evaluated on a dataset by training on the training set and evaluating their performance on the testing set. The final score of a model on a dataset is computed by averaging over the testing set scores for ten different stratified train-test splits of the dataset.

##### B. Training and Model Specification

Our models were implemented using *PyTorch* and *PyTorch-Geometric* [43], [44]. We used the ADAM optimizer for the training algorithm. We considered three learning rates:  $\{0.0005, 0.0001, 0.00005\}$ , and a weight decaying rate of  $5 \times 10^{-4}$  per epoch. We used a batch size of 16 sequences for each training epoch. In our experiments, we trained each model for 200 epochs. Regarding the setting of hyper-parameters, we considered the options described in Section III-B. To ensure a fair comparison between our model and the baseline model, we reported the model’s performance with the lowest validation loss throughout the training for scoring.

All the experiments were conducted on a server with one NVIDIA TITAN-XP graphics card and one NVIDIA GeForce GTX 1080 graphics card. For implementing the baseline model [13], we used the source code available at their open-source repository<sup>3</sup>. Our model implementation is also available at [https://github.com/louisccc/av\\_av](https://github.com/louisccc/av_av).

##### C. Experiments on Risk Assessment

We evaluate each model’s performance by measuring its classification accuracy and the Area Under the Curve (AUC) of the Receiver Operating Characteristic (ROC) for each dataset. The classification accuracy is the ratio of the number of correct predictions on the test set of a dataset to the total number of samples in the testing set. AUC, sometimes referred to as a balanced accuracy measure [45], measures the probability that a binary classifier ranks a positive sample more highly than a random negative sample. This is a more balanced measure for measuring accuracy, especially with imbalanced datasets (i.e. *271-syn*, *1043-syn*, *571-honda*).

From our experimentation, we found that the best option for the hyper-parameters of our model is a mini-batch size of 16 sequences, a learning rate of 0.00005, two MR-GCN layers with 100 hidden units, a SAGPool pooling layer with ratio

<sup>1</sup><https://github.com/carla-simulator/carla>

<sup>2</sup>[https://github.com/carla-simulator/scenario\\_runner](https://github.com/carla-simulator/scenario_runner)

<sup>3</sup><https://github.com/Ekim-Yurtsever/DeepTL-Lane-Change-Classification>

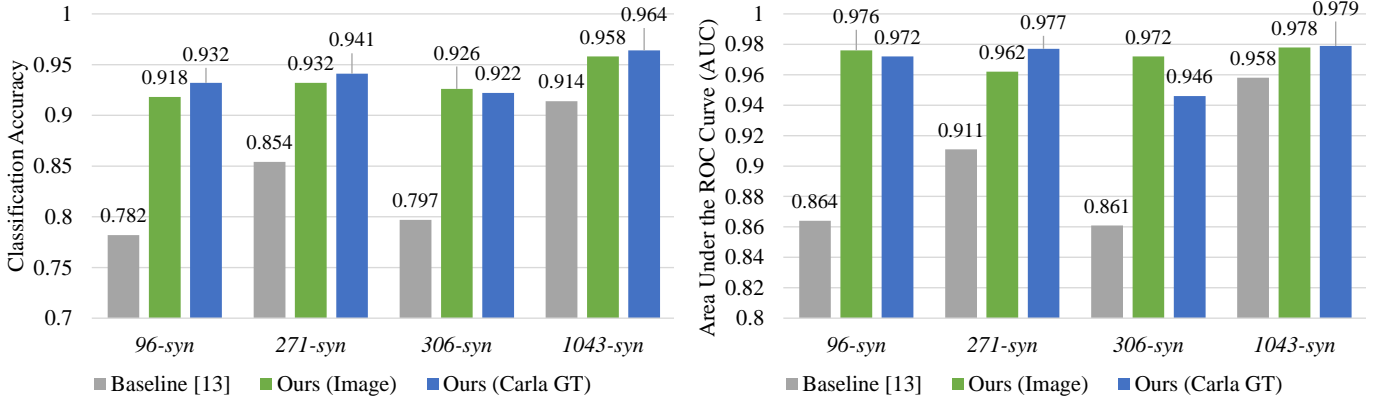


Fig. 3: Accuracy and AUC comparison between our approaches (Real Image and Carla GT) and [13] on different datasets. In these experiments, we trained the model using the hyper-parameter settings reported in Section IV-C.

0.5, *sum-pooling* for graph readout operation and *LSTM-attn* for temporal modeling.

Figure 3 shows the comparison between our model’s performance and the baseline model [13] for all the synthetic datasets. The results show that our approach consistently outperforms [13] across all the datasets in terms of both classification accuracy and AUC. Particularly, on the *1043-syn* dataset, our Image-based and GT pipelines outperform [13] in classification accuracy by 4.4% and 5% respectively (i.e., an accuracy of 95.8% and 96.4% compared to 91.4% for the baseline).

We found that the performance difference between our approach and the baseline increased when the training datasets were smaller. Figure 3 shows that the difference in the accuracy between our approach using the GT pipeline and the baseline [13] is 5% for the *1043-syn* dataset and 8.7% for the *271-syn* dataset. This indicates that our approach can learn an accurate model even from a smaller dataset. We postulate this is a direct result of its use of a scene-graph based IR.

We also found that our approach performs better than the baseline on balanced datasets. Among the datasets used for evaluation of the models, the datasets *271-syn* and *306-syn* contain roughly the same number of clips but differ in the distribution of safe to risky lane changes (2.30:1 for *271-syn* vs. 1:1 for *306-syn*). We found that the performance difference between our image-based approach and the baseline on these datasets is 12.9% on the *306-syn* dataset compared to 7.8% on the *271-syn* dataset, indicating that our approach can discriminate between the two classes better than the baseline.

We also evaluated the contribution of each functional component in our proposed model by conducting an ablation study. The results of the study are shown in Table I. From Table I we find that the simplest of the models, with no MR-GCN layer (replaced with an MLP layer) and a simple average of the embeddings in  $\mathbf{h}_I$  for the temporal model (denoted as *mean* in Table I), achieves a classification accuracy of 75%. We find that replacing *mean* with an LSTM layer for temporal modeling yields a 10.5% increase in performance. We also find that including a single MR-GCN with 64 hidden units and *sum-pooling* to the simplest model results in 14.8% performance gain over the simplest model. The performance

	Spatial Modeling	Temporal Modeling	Avr. Acc.	Avr. AUC
Ablation Study	No MR-GCN	<i>mean</i>	0.762	0.823
	No MR-GCN	<i>LSTM-last</i>	0.867	0.929
	1 MR-GCN	<i>mean</i>	0.910	0.960
	1 MR-GCN	<i>LSTM-last</i>	0.943	0.977
Temporal Attention	No MR-GCN	<i>LSTM-last</i>	0.867	0.929
	No MR-GCN	<i>LSTM-attn</i>	0.868	0.928
	1 MR-GCN	<i>LSTM-last</i>	0.943	0.977
	1 MR-GCN	<i>LSTM-attn</i>	0.950	0.977
Spatial Attention	1 MR-GCN	<i>mean</i>	0.910	0.960
	1 MR-GCN, <i>TopkPool</i>	<i>mean</i>	0.886	0.930
	1 MR-GCN, <i>SAGPool</i>	<i>mean</i>	0.937	0.968

TABLE I: The results of the Carla GT approach on *1043-syn* dataset with various spatial and temporal modeling settings. In these experiments, we used MR-GCN layers with 64 hidden units and *sum-pooling* as the graph readout operation.

gain achieved just by including the MR-GCN layer clearly suggests the effectiveness of modeling the relations between the objects. Finally, we find that the model with one MR-GCN with 64 hidden units and *sum-pooling* plus the LSTM layer for temporal modeling yields the maximum gain of 18.1% over the simplest model. These results clearly demonstrate the importance of each component in the model we propose.

#### D. Evaluation of Attention Mechanisms on Risk Assessment

In this section, we evaluate the various attention components of our proposed model. To evaluate the benefit of attention over the spatial domain, we tested our model with three different graph attention methods: no attention, *SAGPool*, and *TopkPool*. To evaluate the impact of attention on the temporal domain, we tested our model with the following temporal models: *mean*, *LSTM-last*, and *LSTM-attn*. The detailed results that elucidate the effectiveness of these different attention mechanisms are presented in Table I.

For evaluating the benefits of graph attention, we start with an attention-free model: one MR-GCN layer with *sum-pooling* + *mean*. In comparison, the model that uses *SAGPool* for attention on the graph shows a 2.7% performance gain over the attention-free model. This indicates that the use of attention

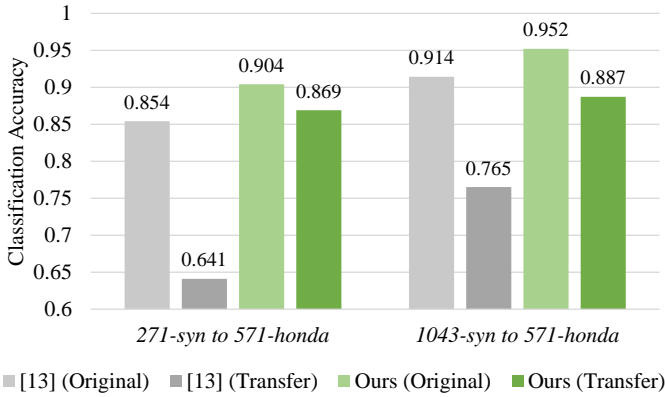


Fig. 4: The results of comparing transferability between our Real Image model and [13]. In this experiment, we trained our model using our best hyper-parameters on both *271-syn* dataset and *1043-syn* dataset. Then we tested the accuracy of our approach on both original dataset and *571-honda* dataset. We followed the same procedure to train and test [13].

over both nodes and relations allows *SAGPool* to better filter out irrelevant nodes from each *scene-graph*. We found that the model using *TopkPool* as the graph-attention layer became relatively unstable, resulting in a 2.4% performance drop compared to the attention-free model. This can be because *TopkPool* ignores the relations of a node when calculating  $\alpha$ . Another reason for this instability could be the random initialization of weights in *TopkPool*, which can exponentially affect the overall performance as stated in [17].

For evaluating the impact of attention on the temporal model, we evaluated the effects of adding a temporal attention layer to the following two models: (i) with no MR-GCN layers and no temporal attention and (ii) with one MR-GCN layer and no temporal attention. Compared to the model with no MR-GCN layer and no temporal attention, the performance of the model with no MR-GCN and *LSTM-attn* was found to be 0.1% higher. We also found that adding *LSTM-attn* to the model with one MR-GCN layer increases its performance by 0.7% over the same model with no temporal attention. These results demonstrate that the inclusion of temporal attention does improve the performance, though only marginally. The reason why we only see a marginal improvement can be that the temporal attention layer is less relevant with the dataset that our model was trained on. When preparing these datasets, we manually removed the frames that are irrelevant to a lane change, exactly the set of frames that temporal attention would have given less attention to, thus minimizing its effect.

The primary benefit of using the graph-attention and the temporal attention is that it improves the explainability of the model’s risk assessment. We demonstrate this capability using the visualization of both graph attention and temporal attention provided in Figure 5. Figure 5 shows the trend of the attention scores  $\beta_1, \beta_2, \dots, \beta_T$  for a risky lane changing clip. Intuitively, the frame with a higher attention score  $\alpha_t$  contributes more to the context vector  $c$  (shown in Equation 12), thus playing a more critical role in calculating  $h_{G_t}$  and contributing to the final risk assessment decision. In this risky lane changing

example, the temporal attention scores progressively increase between frames 19 and 32 during the lane change; and the highest frame attention weights appear in frames 33 and 34, which are the frames immediately before the collision occurs. Figure 5 also shows the projection scores for the node attention layer, where a higher score for a node indicates that it contributes more to the final decision of risk assessment. As shown in this example, as the ego car approaches the yellow vehicle, the node attention weights for the ego car and the yellow vehicle are increased proportionally to the scene’s overall risk. In the first few frames, the risk of collision is low; thus the node attention weights are low; however, in the last few frames, a collision between these two vehicles is imminent; thus the attention weights for the two cars are much higher than for any other nodes in the graph. This example clearly demonstrates our model’s capability to pinpoint the critical factors in a *scene-graph* that contributed to its risk assessment decision. This capability can be valuable for debugging edge cases at design time, thus reducing the chances of ADS making unexpected, erroneous decisions in real-world scenarios and improving human trust in the system.

#### E. Transferability from Virtual To Real Driving

In this section, we demonstrate our approach’s capability to effectively transfer the knowledge learned from a simulated dataset to a real-world dataset. To demonstrate this capability, we use the model weights, and parameters learned from training on the *271-syn* dataset or the *1043-syn* dataset directly for testing on the real-world driving dataset: *571-honda*. We also compare the transferability of our model with that of the baseline method [13]. The results are shown in Figure 4.

As expected, the performance of both our approach and the baseline degrades when tested on *571-honda* dataset. However, as Figure 4 shows, the accuracy of our approach only drops by 6.7% and 3.5% when the model is trained on *271-syn* and *1043-syn*, respectively, while the baseline’s performance drops drastically by a much higher 21.3% and 14.9%, respectively. The results categorically show that our proposed model can transfer knowledge more effectively than the baseline.

## V. CONCLUSION

Subjective risk assessment is a challenging, safety-critical problem that requires a good semantic understanding of many possible road scenarios. Our results show that our scene-graph augmented approach outperforms state-of-the-art techniques at risk assessment tasks in terms of accuracy (95.8% vs. 91.4%) and AUC (0.978 vs. 0.958). We also show that our approach can learn with much less training data than these techniques, as our approach achieves 91.8% accuracy on the *96-syn* dataset compared to 78.2% accuracy achieved by [13]. Additionally, our results show that our approach can better transfer knowledge gained from simulated datasets to real-world datasets (5.0% avg. acc. drop for our approach vs. 18.1% avg. acc. drop for [13]). We also show that the spatial and temporal attention components used in our approach improve both its performance and its explainability.



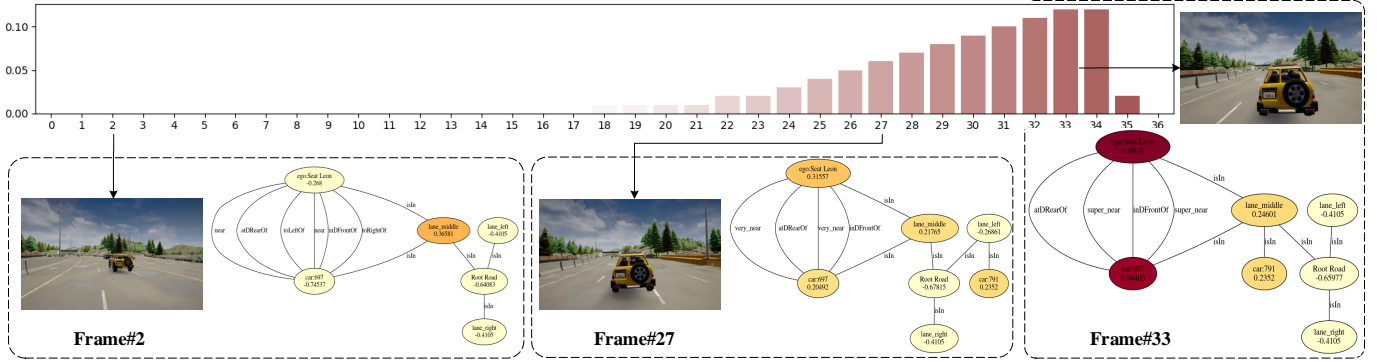


Fig. 5: The visualization of attention weights in both spatial ( $\alpha$ ) and temporal ( $\beta$ ) domains using a risky lane changing clip as an example. We used a gradient color from light yellow to red for visualizing each node’s projection score that indicates its importance in calculating a *scene-graph* embedding. We also used a gradient colored (white to red) bar chart to visualize the temporal attention coefficients  $\beta_1, \beta_2, \dots, \beta_{36}$  used for calculating the context vector  $c$ .

### ACKNOWLEDGMENT

The authors would like to show gratitude to the colleagues of the Embedded & Cyber-Physical Systems (AICPS) Lab, from University of California, Irvine, particularly to Emily Sing Yen Lam and Aung Myat Thu for the contribution during the course of this research. This work was partially supported by NSF under award CMMI-1739503 and the award ECCS-1839429. Any opinions, findings, conclusions, or recommendations expressed in this paper are those of the authors and do not necessarily reflect the views of the funding agency.

### REFERENCES

- [1] R. Mudge, D. Montgomery, E. Groshen *et al.*, “Americas workforce and the self-driving future realizing productivity gains and spurring economic growth,” *no. june*, 2018.
- [2] E. Yurtsever, J. Lambert, A. Carballo *et al.*, “A survey of autonomous driving: Common practices and emerging technologies,” *arXiv preprint arXiv:1906.05113*, 2019.
- [3] B. Paden, M. Čáp, S. Z. Yong *et al.*, “A survey of motion planning and control techniques for self-driving urban vehicles,” *IEEE Transactions on intelligent vehicles*, vol. 1, no. 1, pp. 33–55, 2016.
- [4] D. Lavrinc, “This is how bad self-driving cars suck in rain,” 2014.
- [5] A. Davies, “Google’s self-driving car caused its first crash,” *Wired*, 2016.
- [6] M. McFarland., “Whos responsible when an autonomous car crashes?” *CNN*, 2016. [Online]. Available: <https://money.cnn.com/2016/07/07/technology/tesla-liability-risk/index.html>
- [7] T. B. Lee., “Autopilot was active when a tesla crashed into a truck, killing driver,” *ARS Technica*, 2019. [Online]. Available: <https://arstechnica.com/cars/2019/05/feds-autopilot-was-active-during-deadly-march-tesla-crash/>
- [8] G. Grayson, G. Maycock, J. Groeger *et al.*, “Risk, hazard perception and perceived control,” *TRL Report TRL560*. TRL Ltd., Crowthorne, UK, 2003.
- [9] R. Fuller, “Towards a general theory of driver behaviour,” *Accident analysis & prevention*, vol. 37, no. 3, pp. 461–472, 2005.
- [10] N. Bao, D. Yang, A. Carballo *et al.*, “Personalized safety-focused control by minimizing subjective risk,” in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 3853–3858.
- [11] C. Laugier, I. E. Paromtchik, M. Perrollaz *et al.*, “Probabilistic analysis of dynamic scenes and collision risks assessment to improve driving safety,” *IEEE Intelligent Transportation Systems Magazine*, vol. 3, no. 4, pp. 4–19, 2011.
- [12] H. Woo, Y. Ji, Y. Tamura *et al.*, “Advanced adaptive cruise control based on collision risk assessment,” in *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 939–944.
- [13] E. Yurtsever, Y. Liu, J. Lambert *et al.*, “Risky action recognition in lane change video clips using deep spatiotemporal networks with segmentation mask transfer,” in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*. IEEE, 2019, pp. 3100–3107.
- [14] A. Dosovitskiy, G. Ros, F. Codevilla *et al.*, “Carla: An open urban driving simulator,” *arXiv preprint arXiv:1711.03938*, 2017.
- [15] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” *arXiv preprint arXiv:1409.0473*, 2014.
- [16] A. Adadi and M. Berrada, “Peeking inside the black-box: A survey on explainable artificial intelligence (xai),” *IEEE Access*, vol. 6, pp. 52 138–52 160, 2018.
- [17] B. Knyazev, G. W. Taylor, and M. Amer, “Understanding attention and generalization in graph neural networks,” in *Advances in Neural Information Processing Systems*, 2019, pp. 4202–4212.
- [18] A. Vaswani, N. Shazeer, N. Parmar *et al.*, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [19] W. G. Najm, J. D. Smith, M. Yanagisawa *et al.*, “Pre-crash scenario typology for crash avoidance research,” United States. National Highway Traffic Safety Administration, Tech. Rep., 2007.
- [20] D. A. Pomerleau, “Alvinn: An autonomous land vehicle in a neural network,” in *Advances in neural information processing systems*, 1989, pp. 305–313.
- [21] M. Bojarski, D. Del Testa, D. Dworakowski *et al.*, “End to end learning for self-driving cars,” *arXiv preprint arXiv:1604.07316*, 2016.
- [22] J.-P. Laumond *et al.*, *Robot motion planning and control*. Springer, 1998, vol. 229.
- [23] R. Jain, R. Kasturi, and B. G. Schunck, *Machine vision*. McGraw-Hill New York, 1995, vol. 5.
- [24] R. Rajamani, *Vehicle dynamics and control*. Springer Science & Business Media, 2011.
- [25] C. Chen, A. Seff, A. Kornhauser *et al.*, “Deepdriving: Learning affordance for direct perception in autonomous driving,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2722–2730.
- [26] M. Bansal, A. Krizhevsky, and A. Ogale, “Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst,” *arXiv preprint arXiv:1812.03079*, 2018.
- [27] K. He, G. Gkioxari, P. Dollár *et al.*, “Mask r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2961–2969.
- [28] S. Mylavarapu, M. Sandhu, P. Vijayan *et al.*, “Towards accurate vehicle behaviour classification with multi-relational graph convolutional networks,” *arXiv preprint arXiv:2002.00786*, 2020.
- [29] C. Li, Y. Meng, S. H. Chan *et al.*, “Learning 3d-aware egocentric spatial-temporal interaction via graph convolutional networks,” *arXiv preprint arXiv:1909.09272*, 2019.
- [30] S. Mylavarapu, M. Sandhu, P. Vijayan *et al.*, “Understanding dynamic scenes using graph convolution networks,” *arXiv preprint arXiv:2005.04437*, 2020.
- [31] E. Yurtsever, S. Yamazaki, C. Miyajima *et al.*, “Integrating driving behavior and traffic context through signal symbolization for data reduction and risky lane change detection,” *IEEE Transactions on Intelligent Vehicles*, vol. 3, no. 3, pp. 242–253, 2018.

- [32] D. Xu, Y. Zhu, C. B. Choy *et al.*, "Scene graph generation by iterative message passing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5410–5419.
- [33] J. Yang, J. Lu, S. Lee *et al.*, "Graph r-cnn for scene graph generation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 670–685.
- [34] Y. Wu, A. Kirillov, F. Massa *et al.*, "Detectron2," 2019.
- [35] S. Ren, K. He, R. Girshick *et al.*, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [36] G. Bradski, "The OpenCV Library," *Dr. Dobb's Journal of Software Tools*, 2000.
- [37] H. Caesar, V. Bankiti, A. H. Lang *et al.*, "nuscenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 621–11 631.
- [38] M. Schlichtkrull, T. N. Kipf, P. Bloem *et al.*, "Modeling relational data with graph convolutional networks," in *European Semantic Web Conference*. Springer, 2018, pp. 593–607.
- [39] K. Xu, W. Hu, J. Leskovec *et al.*, "How powerful are graph neural networks?" *arXiv preprint arXiv:1810.00826*, 2018.
- [40] H. Gao and S. Ji, "Graph u-nets," *arXiv preprint arXiv:1905.05178*, 2019.
- [41] J. Lee, I. Lee, and J. Kang, "Self-attention graph pooling," *arXiv preprint arXiv:1904.08082*, 2019.
- [42] V. Ramanishka, Y.-T. Chen, T. Misu *et al.*, "Toward driving scene understanding: A dataset for learning driver behavior and causal reasoning," in *Conference on Computer Vision and Pattern Recognition*, 2018.
- [43] A. Paszke, S. Gross, S. Chintala *et al.*, "Automatic differentiation in pytorch," in *NIPS-W*, 2017.
- [44] M. Fey and J. E. Lenssen, "Fast graph representation learning with PyTorch Geometric," in *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [45] M. Sokolova and G. Lalpalmé, "A systematic analysis of performance measures for classification tasks," *Information processing & management*, vol. 45, no. 4, pp. 427–437, 2009.



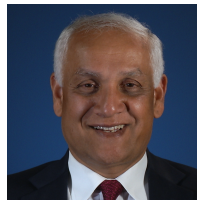
**Shih-Yuan Yu** received the B.S. and M.S. degrees in Computer Science and Information Engineering from the National Taiwan University (NTU) in 2014. He worked at MediaTek for 4 years. Currently he is a Ph.D. student in the University of California, Irvine. Now his research interests are about design automation of embedded systems using data-driven system modeling approaches. It covers incorporating machine learning methods to identify potential security issues in systems.



**Arnav Vaibhav Malawade** received a B.S. in Computer Science and Engineering from the University of California Irvine (UCI) in 2018. He is currently an M.S./Ph.D. Student studying Computer Engineering at UCI under the supervision of Professor Mohammad Al Faruque. His research interests include the design and security of cyber-physical systems in connected/autonomous vehicles, manufacturing, IoT, and healthcare.



**Deepan Muthirayan** is currently a Post-doctoral Researcher in the department of Electrical Engineering and Computer Science at University of California at Irvine. He obtained his Phd from the University of California at Berkeley (2016) and B.Tech/M.tech degree from the Indian Institute of Technology Madras (2010). His doctoral thesis work focussed on market mechanisms for integrating demand flexibility in energy systems. Before his term at UC Irvine he was a post-doctoral associate at Cornell University where his work focussed on online scheduling algorithms for managing demand flexibility. His current research interests include control theory, machine learning, topics at the intersection of learning and control, online learning, online algorithms, game theory, and their application to smart systems.



**Pramod Khargonekar** received B. Tech. Degree in electrical engineering in 1977 from the Indian Institute of Technology, Bombay, India, and M.S. degree in mathematics in 1980 and Ph.D. degree in electrical engineering in 1981 from the University of Florida, respectively. He was Chairman of the Department of Electrical Engineering and Computer Science from 1997 to 2001 and also held the position of Claude E. Shannon Professor of Engineering Science at The University of Michigan. From 2001 to 2009, he was Dean of the College of Engineering and Eckis Professor of Electrical and Computer Engineering at the University of Florida till 2016. After serving briefly as Deputy Director of Technology at ARPA-E in 2012-13, he was appointed by the National Science Foundation (NSF) to serve as Assistant Director for the Directorate of Engineering (ENG) in March 2013, a position he held till June 2016. Currently, he is Vice Chancellor for Research and Distinguished Professor of Electrical Engineering and Computer Science at the University of California, Irvine. His research and teaching interests are centered on theory and applications of systems and control. He has received numerous honors and awards including IEEE Control Systems Award, IEEE Baker Prize, IEEE CSS Axelby Award, NSF Presidential Young Investigator Award, AACC Eckman Award, and is a Fellow of IEEE, IFAC, and AAAS.



**Mohammad Abdullah Al Faruque** (M06, SM15) received his B.Sc. degree in Computer Science and Engineering (CSE) from Bangladesh University of Engineering and Technology (BUET) in 2002, and M.Sc. and Ph.D. degrees in Computer Science from Aachen Technical University and Karlsruhe Institute of Technology, Germany in 2004 and 2009, respectively. He is currently with the University of California Irvine (UCI) as an Associate Professor and Directing the Embedded and Cyber-Physical Systems Lab. He served as an Emulex Career Development Chair from October 2012 till July 2015. Before, he was with Siemens Corporate Research and Technology in Princeton, NJ as a Research Scientist. His current research is focused on the system-level design of embedded and Cyber-Physical-Systems (CPS) with special interest in low-power design, CPS security, data-driven CPS design, etc. He is an ACM senior member. He is the author of 2 published books. Besides many other awards, he is the recipient of the School of Engineering Mid-Career Faculty Award for Research 2019, the IEEE Technical Committee on Cyber-Physical Systems Early-Career Award 2018, and the IEEE CEDA Ernest S. Kuh Early Career Award 2016. He is also the recipient of the UCI Academic Senate Distinguished Early-Career Faculty Award for Research 2017 and the School of Engineering Early-Career Faculty Award for Research 2017. Besides 120+ IEEE/ACM publications in the premier journals and conferences, he holds 8 US patents.