

# INTERPRETABLE SELF-ATTENTION TEMPORAL REASONING FOR DRIVING BEHAVIOR UNDERSTANDING

Yi-Chieh Liu<sup>1,\*</sup>, Yung-An Hsieh<sup>2,\*</sup>, Min-Hung Chen<sup>2</sup>, C.-H. Huck Yang<sup>2</sup>, J. Tegner<sup>3</sup>, Y.-C. James Tsai<sup>1,2</sup>

<sup>1</sup>School of Civil and Environmental Engineering; <sup>2</sup> School of Electrical and Computer Engineering  
Georgia Institute of Technology, Atlanta, GA, USA

<sup>3</sup>Living Systems Laboratory, KAUST, KSA

## ABSTRACT

Performing driving behaviors based on causal reasoning is essential to ensure driving safety. In this work, we investigated how state-of-the-art 3D Convolutional Neural Networks (CNNs) perform on classifying driving behaviors based on causal reasoning. We proposed a perturbation-based visual explanation method to inspect the models' performance visually. By examining the video attention saliency, we found that existing models could not precisely capture the causes (e.g., traffic light) of the specific action (e.g., stopping). Therefore, the Temporal Reasoning Block (TRB) was proposed and introduced to the models. With the TRB models, we achieved the accuracy of **86.3%**, which outperform the state-of-the-art 3D CNNs from previous works. The attention saliency also demonstrated that TRB helped models focus on the causes more precisely. With both numerical and visual evaluations, we concluded that our proposed TRB models were able to provide accurate driving behavior prediction by learning the causal reasoning of the behaviors.

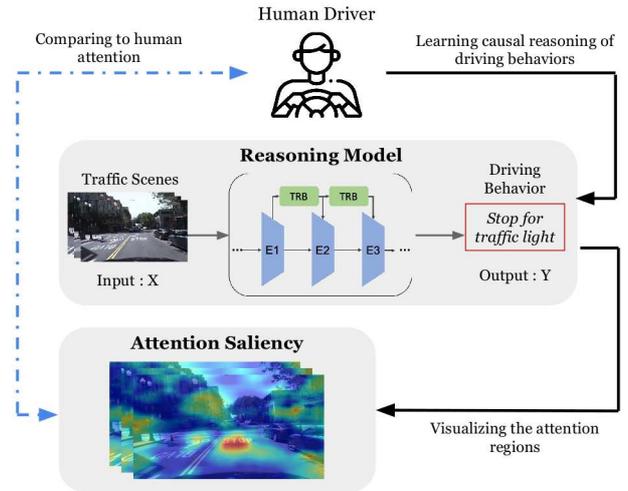
**Index Terms**— Self-driving Vehicles, Driving Behaviors Reasoning, Action Recognition, Self-attention Models, Video Saliency

## 1. INTRODUCTION

Human drivers are required to perform traffic scenes reasoning and act correspondingly to ensure driving safety. Therefore, to design a robust self-driving system, two important factors need to be considered as shown in Fig.1. First, a reasoning model is needed to predict actions based on the reasoning that human drivers perform. Second, a visualization of video attention saliency is required to check and further improve the models on predicting the behaviors based on the correct reasons.

Various driving scene datasets have been introduced [1, 2, 3, 4, 5] to accelerate the pace of self-driving vehicles development. Most of the datasets contain annotations of traffic scenes (e.g. segmentation of pavements [1]) or visual-based control (e.g. steering angles [5]). On the other hand, Ramanish *et al.* [4] proposed the Honda Research Institute Driving Dataset (HDD), which contains the annotations of higher-level driving behavior classes. A 4-layer annotation was defined in the HDD to describe the driving behaviors. To explore the causal reasoning on self-driving vehicles, the Stimulus-driven Action and Cause layers were utilized in this work. A detailed description of the data was shown in Section 4.1.

With the ability to consider both spatial and temporal relationships, video recognition models have proven their effectiveness on action recognition tasks [6, 7, 8]. Video recognition models also benefit the development of algorithms for self-driving vehicles [9,



**Fig. 1:** Designing a robust self-driving system includes two important factors. **Reasoning Model** aims at predicting driving behaviors based on human causal reasoning. Our proposed TRB demonstrated its effectiveness in improving the state-of-the-art models. **Attention Saliency** helps examine and optimize the system to align with the mechanism of how humans interact with complex environments.

10]. By providing video inputs, the models can predict control signals, such as steering angles, for self-driving vehicles. In this work, video recognition models were used as the reasoning model of driving behaviors, in which the behaviors based on causal reasoning (e.g. stop for traffic light) were classified. Video recognition models can be categorized into Convolutional Recurrent Neural Networks (CRNNs) [11] and 3D Convolutional Neural Networks (CNNs) [12]. While [4] proposed a CRNN as a reasoning model of driving behaviors, the performance of 3D CNNs on those higher-level behaviors has not been explored. To further improve the models on reasoning driving behaviors, we introduced a **Temporal Reasoning Block (TRB)** to enhance the model understanding of the causes of driving behaviors. The experiment showed that TRB could help the models capture spatial-temporal features and global dependency within videos using a self-attention mechanism.

The visual traffic scenes usually contain complex information, such that different observed objects can serve as the cause of a certain action in different scenarios. For example, both traffic lights and pedestrians occur in two scenarios in which the car is stopping. In one scenario, the stopping is caused by the traffic light, while the other is caused by the pedestrians crossing the road. To inspect

\*Equal contribution.

whether the model is paying attention to the actual cause of action, a visual explanation of CNNs is required. The visual explanation aims at understanding where the models are paying attention to. In other words, it generates the attention saliency of the input data. To inspect the models for self-driving vehicles, we introduced a spatial-temporal explanation method. **Our contributions include:**

- The investigation of state-of-the-art 3D CNNs on the recognition of driving behaviors based on causal reasoning.
- The introduction of the Temporal Reasoning Block (TRB) for improving the state-of-the-art models on classifying reasoning-based driving behaviors.
- The proposition of a perturbation-based visual explanation method for spatial-temporal models, which enables the inspection of self-driving models.

## 2. RELATED WORK

### 2.1. Self-Driving Behavior Recognition

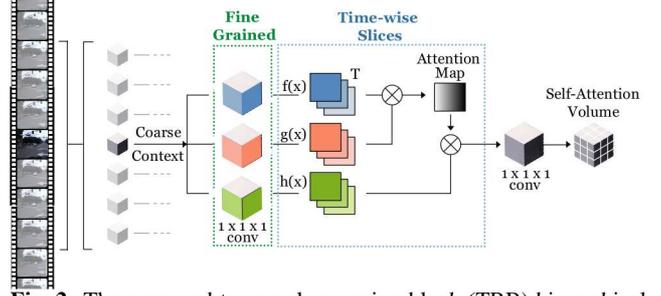
As self-driving technology demonstrated incredible performance in both urban and off-road scenarios [13], the reasoning of self-driving behavior became a needed research problem. ALVINN [14] proposed a shallow neural network to classify actions from images. Researches have also suggested the potential of deep neural networks for tightly coupling the perception and control in simple scenarios [9]. Robust action-perception mapping models were also been explored to recognize the complex visual representation in urban environments [10]. However, these prior efforts formulate the behavior as a goal-oriented task, which is not sufficient to learn how humans drive and interact with traffic scenes. In our work, we explore the cause-aware features by our proposed framework using the HDD [4]. Driving behavior understanding could also be performed by video recognition approaches. One of the approaches is combining CNNs and Recurrent Neural Networks (RNNs) to model the temporal pattern of visual representation [12]. 3D convolutions (C3D) [15] and its improvements [7] also shown great success in space-time features. Also, the 2D residual architecture is also extended into the 3D residual CNNs (3DResNet) [16].

### 2.2. Attention Models

Recently, attention mechanisms have become a reliable method to capture global dependencies [17, 18]. In particular, self-attention [19] represents the importance of different positions in a sequence. The mechanism had been applied to actions recognition tasks in video as a non-local operation [20]. Yet, the potential of self-attention have not been explored on the reasoning tasks of driving behaviors.

### 2.3. Visual Explanation of CNNs

Different methods have been proposed to visually explain where CNNs are paying attention to the input images. Class Class Activation Mapping (CAM) [21] and Grad-CAM [22] generate visual explanation by linear combining activations and class-specific weights. Backpropagation methods were used in DeConvNet [23]. A problem of the above methods is that accessing intermediate layers [21, 22, 24] and/or architectural modification [23] are required. On the other hand, several methods perform visual explanation by perturbing the input images. The dropping of classification score was observed when occluding a small patch on the image in [23], while [25] occluded the image with segmented super-pixels. In [26], the



**Fig. 2:** The proposed temporal reasoning block (TRB) hierarchical structure. The  $\otimes$  denotes matrix multiplication. We applied softmax operation on each row to generate the attention map for each frame and stack them back to spatial-time volume to acquire the temporal-aware self-attention features.

perturbed regions on the image were found by solving an optimization problem. This method can be used on any kind of the model since it is based on explicitly perturbing the images. Therefore, we adapted this method to the visual explanation of self-driving models.

## 3. METHODOLOGY

### 3.1. Temporal Reasoning Block

Most of the non-local based video understanding [20] models optimize the performance with spatial-temporal dependency. However, for the reasoning-oriented dataset, spatial features and temporal integrity should be more weighted to explore the causes of behavior. Inspired by the non-local operation [20], we proposed **temporal reasoning block (TRB)**, a hierarchical self-attention module on different dimensions to reason the driving behavior, as shown in **Fig. 2**. In our proposed method, the TRB consists of two layers which focus on spatial-temporal and spatial domain separately and takes input from previous state-of-art feature extractor. The input  $\mathbf{x} \in \mathbb{R}^{C \times T \times N}$  could be seen as coarse-grained video context. In the first layer of TRB, we used 3D convolution kernel to extract fine-grained spatial-temporal feature  $\mathbf{h}(\mathbf{x}) \in \mathbb{R}^{(C \times T) \times N}$  from  $\mathbf{x}$  to keep the temporal continuity in the second layer, we sliced  $\mathbf{x}$  temporal-wise to  $\mathbf{x}_s \in \mathbb{R}^{C \times N}$  and transformed each  $\mathbf{x}_s$  into two feature spaces  $\mathbf{f}$  and  $\mathbf{g}$  to attain spatial attention, where  $\mathbf{f}(\mathbf{x}) = \mathbf{W}_f \mathbf{x}$ ,  $\mathbf{g}(\mathbf{x}) = \mathbf{W}_g \mathbf{x}$ ,

$$\alpha_{j,i} = \frac{\exp(s_{ij})}{\sum_{i=1}^N \exp(s_{ij})}, s_{ij} = \mathbf{f}(\mathbf{x}_i)^T \mathbf{g}(\mathbf{x}_j) \quad (1)$$

$\alpha_{j,i}$  represents the attention score of the  $i^{th}$  location when interacting with the  $j^{th}$  region.  $C$  is the number of channels and  $N$  is the number of feature locations of the acquired feature map, where  $N = W \times H$ ,  $W$  and  $H$  are the size of the feature map. The output  $\mathbf{o}$  represents the stacked attention map along with time domain, where,

$$\mathbf{o}_j = \sum_{i=1}^N \alpha_{j,i} \mathbf{h}(\mathbf{x}), \mathbf{h}(\mathbf{x}) = \mathbf{W}_h \mathbf{x} \quad (2)$$

$$\mathbf{O}_v = \text{Stack}\{\mathbf{o}_t\}, t = 1 \text{ to } T \quad (3)$$

In above formula,  $\mathbf{O}_v \in \mathbb{R}^{C \times T \times W \times H}$ . For best describing temporal characteristic and memory efficiency, we choose  $\bar{C} = C/8$  in our all experiments. We further scaled the attention volume and added it back to the input  $\mathbf{x}$ . The final output could be represented as,

$$\mathbf{Y}_i = \gamma \mathbf{O}_i + \mathbf{x}_i \quad (4)$$

where  $\gamma$  is a learnable scalar parameter.

### 3.2. Visual Explanation for Self-Driving Models

In [26], the visual explanation was done by finding the regions to perturb the original image which makes the classifier model,  $f_c$ , to produce a minimal score on the target class. This is done by defining a mask  $m : \Lambda \rightarrow [0, 1]$ , where  $\Lambda = \{1, \dots, H\} \times \{1, \dots, W\}$ . The perturbed image can be defined as:

$$[\Phi \{x_0; m\}](u) = m(u)x_0(u) + (1 - m(u))x_p(u) \quad (5)$$

where  $u \in \Lambda$  is each pixel location and  $x_p$  is generated by performing perturbation on the whole original image  $x_0$ . In other words, the more  $m(u)$  is closer to 0, the more severe the pixel  $u$  will be perturbed. Different kinds of perturbation operations can be used, such as add noise or blurring the image. To find the optimal mask  $m^*$  that minimizes the classification score  $f_c(\Phi\{x_0; m\})$ , a objective function is defined as:

$$\min_{m \in [0,1]^\Lambda} f_c(\Phi\{x_0; m\}) + \lambda_1 \|1 - m\|_1 + \lambda_2 \sum_{u \in \Lambda} \|\nabla m(u)\|_\beta^\beta \quad (6)$$

where the L1 term aims at making the mask as small as possible and the total variation (TV) regularization aims at smoothing the mask. To prevent the mask from over-fitting the network, a lower resolution mask is solved during the optimization process. The mask is up-sampled to the image size when performing the perturbation.

To visually inspect self-driving models, the perturbation-based method needs to be extended to perform the spatial-temporal explanation. A three dimensional mask is now used, such that  $m : (\Lambda, T) \rightarrow [0, 1]$ , where  $T$  is the size of the additional temporal dimension. To consider the relationship between each frame of the input video, a temporal TV regularization term is introduced:

$$\sum_{t \in T} \lambda_t \|\nabla m(:, t)\|_\beta^\beta \quad (7)$$

where  $\lambda_t$  is the regularization coefficient. By adding the third dimension in eq. (6) and the temporal tv regularization, the objective function can now be written as:

$$\min_{m \in [0,1]^{(\Lambda, T)}} f_c(\Phi\{x_0; m\}) + \lambda_1 \|1 - m\|_1 + \sum_{t \in T} \left( \lambda_s \sum_{u \in (\Lambda, t)} \|\nabla m(u, t)\|_\beta^\beta + \lambda_t \|\nabla m(:, t)\|_\beta^\beta \right) \quad (8)$$

where  $\lambda_s$  is the coefficient of the spatial tv regularization. Similarly, a lower resolution mask in both spatial and temporal dimensions is used to prevent the mask from over-fitting the network.

## 4. EXPERIMENTAL ENVIRONMENT

Two experiments were conducted to evaluate and inspect the proposed TRB models on the reasoning of driving behaviors. A baseline CRNN and the 3D CNNs, both with and without the proposed TRB, were trained and evaluated. Then the proposed perturbation-based explanation method was applied to generate the video attention saliency, which was used to inspect where the models were paying attention to.

### 4.1. Dataset Preparation

The annotation layers of Stimulus-driven Action and Cause in the HDD were utilized. The Stimulus-driven Action contains *stop* and *deviate*, while Cause contains the reasons of these two actions. To

focus the experiment on the reasoning for driving behaviors, we simplified the task to consider a single action class. Stopping is a fundamental but critical behavior to ensure driving safety. Therefore, the stop action was selected, with the causes of traffic light (stop4light), pedestrian (stop4ped), stop sign (stop4sign), and congestion (stop4cong).

Video clips were prepared as the input of the models. Each clip contains 20 continuous frames with 4 frame-per-second, and each frame image has a size of 360 by 360. Each clip also has a target label from the 4 classes. The numbers of video clips for each class were shown in **Table 1**.

Data splits	stop4light	stop4ped	stop4sign	stop4cong
Train	100	45	170	170
Validation	10	6	20	20
Test	13	10	30	30

**Table 1:** The number of video clips for each class.

### 4.2. Model Implementation

**Baseline CRNN:** The *CNN conv* model in [4] was set up as the baseline in the experiment. The InceptionResnet-V2 model was served as the CNN encoder. The encoded features were flattened and inputted to a Long-Short Term Memory (LSTM) network. The output of the model was generated from the last LSTM hidden state by connecting a fully connected (FC) layer with the size of four.

**3D models:** All 3D CNNs were implemented base on the *Conv3d* layer in Pytorch. The C3D [15], I3D [7], and 3DResnet [16] models were implemented and fine-tuned on the prepared HDD. The Resnet34 [27] architecture was used in the 3DResnet model. We observed that adding TRB in shallow layers of the models did not improve the performance. Therefore, we added four TRBs in the dur layers of each model.

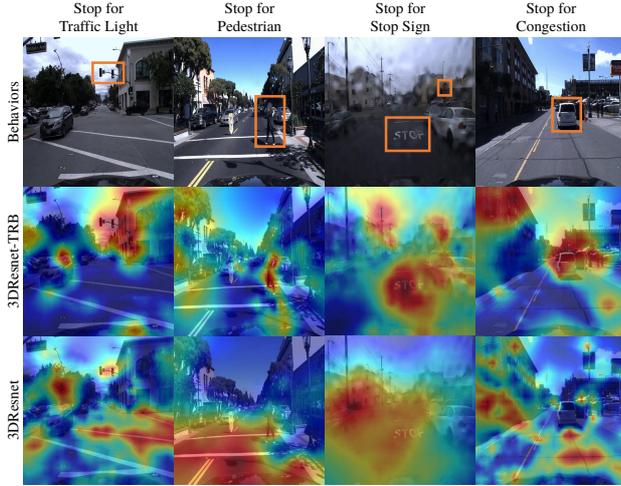
**Training:** The models were trained on a single GeForce RTX 2080Ti GPU. Cross entropy loss was used as the objective function for the recognition task. Stochastic gradient descent (SGD) was used as the optimizer, with a learning rate of 0.0001, a momentum of 0.9, and a weight decay of 0.0001. The models were trained with 150 epochs and the batch size set to 4. Early stopping was applied to select the weights of the model that were used in the evaluation stage. To deal with the imbalanced training samples for different classes, oversampling of data was used during the training process.

**Evaluation:** Different from [4], which performed per-frame evaluation on testing videos and calculated the average precision(AP), we measured the model performances based on the video recognition tasks. Video clips, where each contains a single target label, were used for evaluation. The accuracy of classifying the label for each video clip in the test set was computed for each model.

### 4.3. Visual Explanation

The visual explanation was generated using Adam optimizer to minimize eq. (8). The following parameters were used in the experiment: learning rate = 0.01, iteration = 500,  $\lambda_1 = 0.001$ ,  $\lambda_s = 0.2$ ,  $\lambda_t = 0.1$ , and  $\beta = 3$ . The image was perturbed by combining both Gaussian and median blur. The perturbation mask  $m$  was set to the size of  $(28 \times 28 \times 10)$ . The mask was up-sampled to the size of the video clip and served as an attention saliency  $M$ , which larger intensity indicated more attention.

To quantitatively examine how much the models pay attention to a certain object, we defined an attention score,  $S_{Atten}$ , for each



**Fig. 3:** Attention saliency of each driving behavior class. The first row showed the frames with the causes labeled by bounding boxes. The second and third rows showed the attention saliency of 3DResnet-TRB and 3DResnet, respectively.

object  $o$  in a single frame  $I_t$  as:

$$S_{Atten} = \sum_{u \in I_t} s_u, \text{ where } s_u = \begin{cases} M(u, t) & d(u, c_o) \leq r_o \\ \frac{M(u, t)}{d(u, c_o)} & d(u, c_o) > r_o \end{cases} \quad (9)$$

In eq. (9),  $M(u, t)$  is the value of the attention saliency at a pixel location  $u$  in time  $t$ .  $d(u, c_o)$  is the distance between  $u$  and the center of the object  $c_o$ , while  $r_o$  is the radius of the circle that encloses the object. To compensate the different sizes of the objects, we normalized the attention score when comparing different objects.

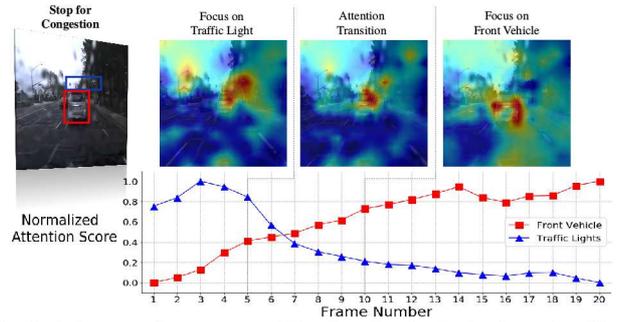
## 5. EXPERIMENTAL RESULTS

### 5.1. Driving Behavior Recognition

The evaluation results of each model on classifying the driving behaviors were presented in **Table 2**. Among the 3D CNNs, both I3D and 3DResnet showed better performance than the baseline CRNN. This indicated that the state-of-the-art 3D CNNs were able to more precisely classify the reasoning-based driving behaviors. With the highest accuracy among the models without TRB, 3DResnet also demonstrated the effectiveness of the residual architecture on this task. As shown in the table, the proposed TRB largely improved the performance of all models. This indicated that the self-attention mechanism in TRB effectively helped the models to capture the global dependency within the videos. The result also showed that the TRB can be flexibly applied to different models of driving behavior recognition to provide improvement.

### 5.2. Attention Saliency of Driving Behaviors

Examples of generated attention saliency using the proposed visual explanation method were shown in **Fig. 3**. The frame images were shown in the first row, with the causes of the action (e.g. stop signs) being labeled by bounding boxes. The attention saliency was presented as heatmaps in the second and third rows. As shown in the figure, adding the proposed TRB to 3DResnet made the model capture the causes more precisely. For each of the class, 3DResnet-TRB



**Fig. 4:** The attention scores of the front vehicle (red) and traffic lights (blue) in a stop for congestion behavior occurred in a rainy scenario. As the vehicle approached the front vehicle, the attention transition from traffic lights to the front vehicle was observed.

Model	Accuracy	Model	Accuracy
CRNN	73.49%	CRNN-TRB	<b>78.31%</b>
C3D	60.71%	C3D-TRB	<b>69.88%</b>
I3D	77.11%	I3D-TRB	<b>83.13%</b>
3DResnet	83.56%	3DResnet-TRB	<b>86.30%</b>

**Table 2:** The recognition results of driving behaviors based on causal reasoning.

paid attention to most of the regions around the bounding boxes. This demonstrated the effectiveness of TRB on the reasoning of driving behaviors.

**Fig. 4** presented a quantitative evaluation of attention saliency generated from 3DResnet-TRB. A stop for congestion behavior in the rainy scenario was showed. The attention score of the front vehicle and the average score of multiple traffic lights (enclosed by the blue bounding box) were computed. A reasonable attention transition can be observed from both the attention scores and attention saliency. When the vehicle approached the front vehicle, the model gradually shifted its attention from the traffic lights to the front vehicle, which was the cause of the stopping behavior. This example indicated that the proposed TRB helped the models pay attention to reasonable regions when predicting driving behaviors.

## 6. CONCLUSIONS

In this work, we proposed the **Temporal Reasoning Block (TRB)** to improve the performance of video recognition models on reasoning driving behaviors. The results showed that TRB can largely improve the state-of-the-art models, both CRNN and 3D CNNs, on classifying driving behaviors. By adding TRB to the 3DResnet, we achieved the highest accuracy of 86.3%. A perturbation-based visual explanation method was also proposed to generate video attention saliency. The examination of the attention saliency demonstrated that 3DResnet-TRB was able to focus on reasonable objects when classifying driving behaviors. With both numerical and visual evaluations, we concluded that our proposed TRB models were able to provide accurate driving behavior prediction by learning the causal reasoning of the behaviors.

## 7. REFERENCES

- [1] Andreas Geiger, Philip Lenz, and Raquel Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012, pp. 3354–3361.
- [2] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 3213–3223.
- [3] Vashisht Madhavan and Trevor Darrell, *The BDD-Nexar Collective: A Large-Scale, Crowdsourced, Dataset of Driving Scenes*, Ph.D. thesis, Master's thesis, EECS Department, University of California, Berkeley, 2017.
- [4] Vasili Ramanishka, Yi-Ting Chen, Teruhisa Misu, and Kate Saenko, "Toward driving scene understanding: A dataset for learning driver behavior and causal reasoning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7699–7707.
- [5] Udacity, "Udacity self-driving car dataset," 2017.
- [6] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah, "Ucf101: A dataset of 101 human actions classes from videos in the wild," *arXiv preprint arXiv:1212.0402*, 2012.
- [7] Joao Carreira and Andrew Zisserman, "Quo vadis, action recognition? a new model and the kinetics dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6299–6308.
- [8] Chih-Yao Ma, Min-Hung Chen, Zsolt Kira, and Ghassan Al-Regib, "Ts-lstm and temporal-inception: Exploiting spatiotemporal dynamics for activity recognition," *Signal Processing: Image Communication*, vol. 71, pp. 76–87, 2019.
- [9] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al., "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.
- [10] Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell, "End-to-end learning of driving models from large-scale video datasets," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2174–2182.
- [11] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2625–2634.
- [12] Joe Yue-Hei Ng, Matthew Hausknecht, Sudheendra Vijayanarasimhan, Oriol Vinyals, Rajat Monga, and George Toderici, "Beyond short snippets: Deep networks for video classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 4694–4702.
- [13] Martin Buehler, Karl Iagnemma, and Sanjiv Singh, *The DARPA Urban Challenge: Autonomous Vehicles in City Traffic*, Springer Publishing Company, Incorporated, 1st edition, 2009.
- [14] Dean A. Pomerleau, "Advances in neural information processing systems 1," chapter ALVINN: An Autonomous Land Vehicle in a Neural Network, pp. 305–313. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1989.
- [15] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4489–4497.
- [16] Kensho Hara, Hirokatsu Kataoka, and Yutaka Satoh, "Can spatiotemporal 3d cnns retrace the history of 2d cnns and imagenet?," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 6546–6555.
- [17] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [18] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *International conference on machine learning*, 2015, pp. 2048–2057.
- [19] Ankur P Parikh, Oscar Täckström, Dipanjan Das, and Jakob Uszkoreit, "A decomposable attention model for natural language inference," *arXiv preprint arXiv:1606.01933*, 2016.
- [20] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He, "Non-local neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7794–7803.
- [21] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2921–2929.
- [22] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 618–626.
- [23] Matthew D Zeiler and Rob Fergus, "Visualizing and understanding convolutional networks," in *European conference on computer vision*. Springer, 2014, pp. 818–833.
- [24] Chao-Han Huck Yang, Yi-Chieh Liu, Pin-Yu Chen, Xiaoli Ma, and Yi-Chang James Tsai, "When causal intervention meets adversarial examples and image masking for deep neural networks," in *2019 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2019, pp. 3811–3815.
- [25] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin, "Why should i trust you?: Explaining the predictions of any classifier," in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, 2016, pp. 1135–1144.
- [26] Ruth C Fong and Andrea Vedaldi, "Interpretable explanations of black boxes by meaningful perturbation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3429–3437.
- [27] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.