

Received November 27, 2020, accepted December 21, 2020, date of publication December 24, 2020, date of current version January 5, 2021.

Digital Object Identifier 10.1109/ACCESS.2020.3047091

The Traffic Scene Understanding and Prediction Based on Image Captioning

WEI LI^{1,2}, ZHAOWEI QU¹, HAIYU SONG^{1,2}, PENGJIE WANG², AND BO XUE²

¹Transportation College, Jilin University, Changchun 130022, China

²School of Computer Science and Engineering, Dalian Minzu University, Dalian 116600, China

Corresponding author: Zhaowei Qu (quzw@jlu.edu.cn)

This work was supported in part by the Liaoning Natural Science Foundation Projects under Grant 2019-ZD-0182, and in part by the Creative Research Talents from Universities of Liaoning Province under Grant LR2016071.

ABSTRACT The traffic scene understanding is the core technology in Intelligent Transportation Systems (ITS) and Advanced Driver Assistance System (ADAS), and it is becoming increasingly important for smart or autonomous vehicles. The recent methods for traffic scene understanding, such as Traffic Sign Recognition (TSR), Pedestrian Detection, and Vehicle Detection, have three major shortcomings. First, most models are customized for recognizing a specific category of traffic target instead of general traffic targets. Second, as for these recognition modules, the task of traffic scene understanding is to recognize objects rather than make driving suggestions or strategies. Third, numerous independent recognition modules disadvantage to fusing multi-modal information to make a comprehensive decision for driving operation in accordance with complicated traffic scenes. In this paper, we first introduce the image captioning model to alleviate the aforementioned shortcomings. Different from existing methods, our primary idea is to accurately identify all categories of traffic objects and understand traffic scenes by making full use of all information, and making the suggestions or strategy for driving operation in natural language by using Long Short Term Memory network (LSTM) rather than keywords. The proposed solution naturally solves the problems of feature fusion, general object recognition, and low-level semantic understanding. We tested the solution on our created traffic scene image dataset for evaluation of image captioning. Extensive experiments including quantitative and qualitative comparisons demonstrate that the proposed solution can identify more objects and produce higher-level semantic information than the state-of-the-arts.

INDEX TERMS Image captioning, traffic scene understanding, intelligent transportation system, advanced driver assistance system, LSTM, driving suggestions, driving strategies.

I. INTRODUCTION

Developing new Intelligent Transportation Systems (ITS) which take into consideration the socio-economical, environmental, and safety factors of the modern society is one of the grand challenges of the next century. Autonomous vehicles are promising for improving the safety and efficiency of transportation and mobility systems by providing vehicle control during normal driving. However, a human driver will still be expected to have control responsibilities for an automated driving car in the case of an emergency. Therefore, the most possible solution to ITS may be mixed autonomous vehicles and human-driven vehicles rather than fully automated without human operation. Therefore, the Advanced Driver Assistance System (ADAS) plays an important role in the Intelligent Transportation System. Nevertheless, decades

of researches for ADAS are concentrated on traffic target detection and classification, rather than complicated traffic scene perception or driving strategic decision.

In summary, the current ADASs have two shortcomings. First, The methods mainly focus on specific problems ranging from traffic sign recognition [1]–[3] and vector detection [4]–[6] to pedestrian detection [7]–[9]. Up to the present, most attention has been paid to the researches on single or multiple target detection. However, there are various targets likely not belonging to the above ones in reality, and the ADAS will omit the new targets to put the driver in danger [58], [59]. Second, the current ADASs could not supply the driver intelligent driving suggestion and strategy of traffic scene ahead yet [55]–[57]. In contrast, the intelligent suggestion and appropriate strategies are capable of curtailing the response time of the driver to avoid the accident. In the area of deep learning, image captioning has attracted much more attention in recent years, which can provide a solution to

The associate editor coordinating the review of this manuscript and approving it for publication was Keli Xiao¹.

the aforementioned problems. Image captioning can generate descriptions of the traffic situations and assist with ADAS under special traffic conditions. Third, the existing methods could not fuse all categories of information to represent comprehensive scenes or abstract semantic concepts.

In this paper, we provide a novel solution for ADAS by utilizing image captioning. To the best of our knowledge, this is the first work that introduces image captioning into ITS(AVs). The main goal of this work is to allow the development of a new generation of ADAS solutions where control could actually be shared between the automated vehicle and the driver. Our contributions are as follows:

- We propose a solution for traffic scene understanding with the image captioning method. Benefiting from image captioning, our method can not only outperform the traditional target detection but also comprehend and discover all the targets threatening to driving.
- Human mistakes due to inattention have been identified as the major crash-contributing factor by large-scale naturalistic studies [10], [11]. Our solution could call the driver's attention when he is distracted. Furthermore, the solution can provide an operation decision or strategy for drivers in urgent and limited time, which can reduce traffic accidents significantly.
- As there is no appropriate dataset about image captioning in the field of traffic, we create a dataset based on the LaRA dataset [51] and perform our experiments on the dataset.

II. RELATED WORKS

A. TRAFFIC SIGN RECOGNITION

Traffic Sign Recognition (TSR) is the main issue for a driver assistance system as it has a dual role to control the road traffic as well as warning and guiding the driver. The first research on traffic sign recognition can be traced back to 1987 when Akatsuka and Imai [12] attempted to make an early TSR system. The TSR systems are usually comprised of two parts: traffic sign detection and traffic sign classification. The regions of interest are first extracted in detection, and then identified correctly or rejected in classification. The methods that were used for the extraction of regions of interest also include the color-based technique, as in [13]–[16]. In recent years, some researchers have used different architectures and forms of Neural Networks for classification [17]–[20].

B. PEDESTRIAN DETECTION

Pedestrian detection in autonomous driving provides an important guarantee for safe and stable operation of the intelligent driving system in complex and uncertain environments. Traditional algorithms include HOG features [21], SIFT feature [22], Haar feature [23], LBP features [24], SVM, the Deformable Parts Model (DPM) [25], and pedestrian detection based on action features. With the rapid development of deep learning, the recent researches on pedestrian detection have become more and more prominent. The standard Regional Convolutional Neural Network (R-CNN)

detection framework [26] was enhanced with region proposals and CNN classification as in Faster R-CNN [27].

C. VEHICLE DETECTION

Vehicle detection is a very important task for traffic information analysis that can be used in traffic control and management to ensure a safe transportation system. The most recent studies on vehicle detection can be categorized into two groups [28]. The first is based on vehicle appearance, and the second is based on vehicle motion. The appearance-based strategies depend mainly on visual features, including vehicle symmetry, texture, edges, and color [29]–[32]. The concept of motion-based techniques is to extract the moving vehicles based on their motion characteristics that separate them from the background, such as optical flow, frame differences, and background subtraction [33]–[35].

In recent years, there is an increasing concern about vehicle detection. Considerable researches have been devoted to improving the training efficiency by combining the CNN with the circular pairwise classification for vehicle recognition [36].

D. IMAGE CAPTIONING

Being capable of feature representation, deep convolutional neural networks have achieved dramatic progress in object detection. Many famous region proposal-based methods, such as R-CNN [26], Fast R-CNN [37], and Faster R-CNN [27] mentioned above, have great performance in object detection. In recent years, The regression-based single networks include YOLO [38], YOLO2 [39], SSD [40], and these types of detection methods use a single convolutional network to predict bounding boxes and class probabilities simultaneously from an input image. Although many successes have been achieved in traffic target detection, the categories of targets are specific and determined. However, there are many targets that the existing detectors can not recognize in the traffic environment. Furthermore, the previous methods of object detection can not supply the strategies and prediction for driving.

It is a relatively new task to let a computer use a human-like sentence to automatically describe an image. Image captioning will lead to a great number of possible applications, such as producing natural human-robot interactions, early childhood education, information retrieval, and visually impaired assistance, and so on. As a challenging and meaningful research field in artificial intelligence, image captioning is attracting more and more attention and becoming increasingly important. But this task is significantly harder than the well-studied image classification or object recognition tasks, which have been a main focus in the computer vision community [41]. The image captioning models not only must be powerful enough to solve the computer vision challenges of determining which objects are in an image, but they must also be capable of capturing and expressing their relationships in a natural language. For this reason, image captioning is viewed as a difficult problem and it is a very important challenge for

machine learning and deep-learning algorithms. So there has been a recent surge of research interest in attacking the image captioning problems [42]–[45].

Recent works about image captioning mainly focus on language-based methods with an encoder-decoder framework [41], [46]–[50], where a convolutional neural network (CNN) encodes images into visual features, and a Recurrent Neural Network (RNN) decodes features into sentences [41].

Karpathy [48] proposed a system to provide natural language descriptions for image regions. He employed a CNN to extract features from image regions and an RNN to generate a short sentence for each region. Vinyals [41] also proposed a holistic method based neural image caption generator, known as NIC. He integrated a deep CNN as the image encoder for vision feature learning and an RNN for caption generating. Significantly higher BLEU scores were achieved in comparison with those state-of-the-art methods when evaluated using MSCOCO, Flickr8k, and Flickr30k datasets. Xu introduced an attention-based model automatically learning to describe the content of the image. He showed how the model was able to automatically learn to fix its gaze on salient objects while generating the corresponding words in the output caption [50].

III. MODEL

To apply the image captioning method in traffic scene understanding, we should focus on describing the factors influencing driving. Inspired by Vinyals [41], we introduce the image captioning method about the encoder-decoder model called Neural Image Caption (NIC). The encoder and the decoder are constructed by CNN and RNN respectively. The model is an end-to-end system that is fully trainable using stochastic gradient descent (SGD).

We start by briefly describing the encoder-decoder framework [41]. Given an image and the corresponding caption, the model directly maximizes the following objective.

$$\theta^* = \arg \max_{\theta} \sum_{(I,S)} \log p(S|I; \theta) \quad (1)$$

where θ are the parameters of the model, I is an image, and $S = \{S_1, S_2, \dots, S_t\}$ is the correct caption. S represents any sentence and its length is unbounded. Using the chain rule, the log-likelihood of the joint probability distribution can be decomposed into ordered conditionals:

$$\log p(S|I) = \sum_{t=0}^N p(S_t|I, S_0, \dots, S_{t-1}) \quad (2)$$

where N is the length of this particular example and we drop the dependency on model parameters for convenience. At training time, (S, I) is a training example pair, and the sum of the log probabilities is optimized as described in (2) over the whole training set using stochastic gradient descent (SGD).

To model $p(S_t|I, S_0, \dots, S_{t-1})$, it is natural to apply an RNN, where the variable number of words we condition upon up to

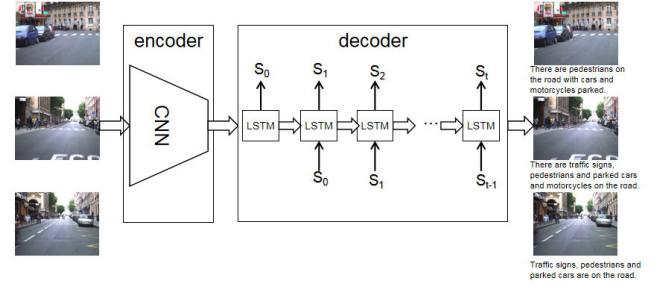


FIGURE 1. The overall system architecture.

$t-1$ is expressed by a fixed-length hidden state or memory h_t . This memory is updated after seeing a new input x_t by using a non-linear function f :

$$h_{t+1} = f(h_t, x_t) \quad (3)$$

Though RNN has been proven successful on tasks such as speech recognition and text generation, it can be difficult to learn long-term dynamics being in part to the vanishing and exploding gradients problem. So in this stage, we utilize a Long-Short Term Memory (LSTM) net, which has shown state-of-the-art performance on sequence tasks such as translation. Formula (3) is replaced by LSTM as following:

$$h_{t+1} = \text{LSTM}(h_t, x_t) \quad (4)$$

Our goal is to describe traffic scenes using image captioning based on the above framework, so we construct the architecture of our model (see Figure 1). Furthermore, as shown in Figure 1, abundant targets are described in the sentence to provide as much as possible information about the traffic environment.

A. THE ENCODER

To encode images, we use a CNN architecture as a feature extractor. In the encoder stage of the image captioning model, we use the VGG-16 network [54] to extract the feature vectors of an image.

We extract features from a lower convolutional layer of the VGG-16 network to get the correspondence between the feature vectors and portions of the 2-D image. The work is different from the previous methods which instead used a fully connected layer [42], [44]. Thus, the decoder can selectively focus on a certain part of an image by selecting a subset of all the feature vectors.

B. THE DECODER

The features extracted by the VGG16 network will be sent to the LSTM network to decode them and generate the corresponding image content description. The advantage of LSTM is modeling and predicting the implied context-dependency of information sequence. Figure 3 shows the basic memory unit structure of LSTM. The core of the LSTM model is a memory cell encoding knowledge at each time step for what inputs have been observed up to this step.

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f) \quad (5)$$

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i) \quad (6)$$

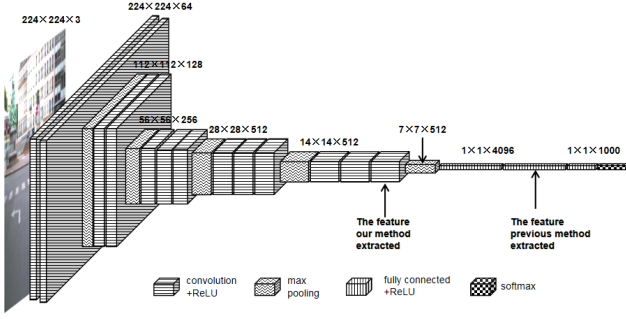


FIGURE 2. The difference between our method and the previous methods in VGG16 network.

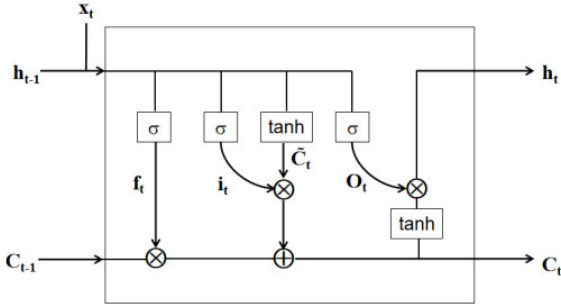


FIGURE 3. The LSTM memory cell.

$$\tilde{C}_t = \tanh(W_c[h_{t-1}, x_t] + b_c) \quad (7)$$

$$C_t = f_t C_{t-1} + i_t \tilde{C}_t \quad (8)$$

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o) \quad (9)$$

$$h_t = o_t \tanh(C_t) \quad (10)$$

where $i_t, f_t, o_t, \tilde{C}_t, C_t, h_t$ are the input, forget, output, temporary memory, memory, and hidden state of the LSTM. The behavior of the cell is controlled by the gates which can either keep a value from the gated layer if the gate is 1 or zero value if the gate is 0. Furthermore, the forget gate is used to control whether to forget the current cell value, the input gate is used to read its input, and the output gate is used to decide whether to output the new cell value. Where σ is the sigmoid activation function. The LSTM can learn exceedingly complex and long-term temporal dynamics from these additional cells. Additional depth can be added to LSTMs by stacking them on top of each other, using the hidden state of the LSTM in layer $t-1$ as the input to the LSTM in layer t .

The unrolled LSTMs memories of the model are shown in Figure 4 and all LSTMs share the same parameters. The output h_{t-1} of the LSTM at time $t-1$ is fed to the LSTM at time t . All recurrent connections are transformed to feed-forward connections in the unrolled version as following:

$$x_{-1} = CNN(I) \quad (11)$$

$$x_t = W_e S_t, \quad t \in \{0 \dots N-1\} \quad (12)$$

$$p_{t+1} = LSTM(x_t), \quad t \in \{0 \dots, N-1\} \quad (13)$$

where W_e is word embedding and the image I is only input once at x_{-1} to inform the LSTM about the image contents.

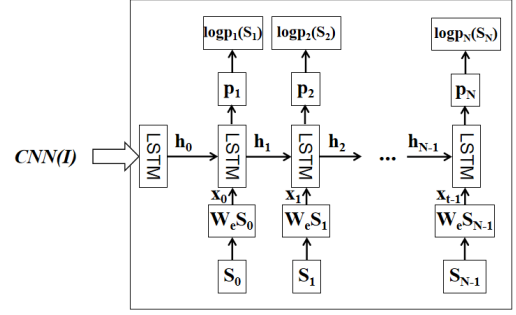


FIGURE 4. The decoder framework.

The loss function is the sum of the negative log-likelihood of the correct word at each step as follows:

$$L(I, S) = - \sum_{t=1}^N \log p_t(S_t) \quad (14)$$

The loss is minimized with respect to all the parameters of the LSTM, the top layer of the image embedder CNN and word embedding W_e (see Figure 4).

IV. EXPERIMENTS

A. DATA PREPARATION

As the first work that introduces image captioning into traffic, to evaluate this work, we created an image captioning dataset for traffic scene understanding based on the LaRA dataset [51] and annotated each image with an English sentence. This dataset includes abundant traffic scenes with caption as a clue for the strategy or the prediction (see Table 1). We annotated 11178 images in the LaRA dataset manually and selected one-tenth of them randomly as the test set.









B. TRAINING DETAIL

The experiments are implemented in Anaconda with PyTorch open-source machine learning framework. We choose the widely used VGG-16 network architecture as a feature extractor. The input images were set to 224×224 according to the VGG-16 network, the size of the feature map extracted from the VGG-16 network is 14×14 , and the dimension of the feature map is 512. The feature map was linearized into a vector which was the first input information into the following LSTM network.

In the training process, we applied the pre-trained VGG-16 network to extract the feature map of an image in the encoder. In the decoder, the LSTM knew the caption (ground truth) of the image, so that each LSTM cell had the source word and the destination word (the input and the output of an LSTM cell). The parameters of our image captioning model were trained with stochastic gradient descent using adaptive learning rate algorithms.

However, in the testing process, the difference with the training process was that the decoder didn't know the real caption of the image, and LSTM cells had no source word and destination word. The start of the caption was generated from the first LSTM cell based on the feature vector extracted from the VGG-16 network. The first word (S_1 in Figure 4) was

TABLE 1. The samples of our created dataset.

Image	Caption	Image	Caption
	There is car ahead and some motors by the road.		There are some cars on both sides of the road, the traffic light is red, pedestrians on the road and a bus passing by.
	There are some motors and cars on both sides of the road.		There are traffic signs, pedestrians and parked cars and motorcycles on the road.
	There is a white car in the middle of the road and some cars on both sides of the road.		There is a car in ahead, some pedestrians are crossing the road and some cars in the distance, the light is red.
	There are sidewalks, cars, bicycles, billboards, traffic signs, red lights on both sides.		There are traffic signs, pedestrians, red traffic lights and cars on the road.

generated base on the start of the caption and the second word (S_2 in Figure 4) was generated base on the first word. When the destination word is 'end' or the length of the generating caption exceeds the predefined threshold, the testing process can generate a caption of the input image.

In most cases, we must handle overfitting when training models. Nonetheless, we explored several techniques to deal with overfitting. The most obvious way to not overfit is to initialize the weights of the CNN component of our system to a pre-trained model (e.g., on ImageNet). We did this in our experiments and it did help quite a lot in terms of generalization. Another set of weights that could be sensibly initialized are W_e , the word embeddings. As no significant gains were observed after trying some initialization method such as mentioned in paper [52], we uninitialized the weights for simplicity.

We trained all sets of weights using stochastic gradient descent with a learning rate of 0.005. All weights were randomly initialized except for the CNN weight, which we trained it by ImageNet. After pretraining with ImageNet, we fine-tuned the model with our created dataset. We used 512 dimensions for the embeddings and the size of the LSTM memory. We set the mini-batch size as 10.

C. EXPERIMENTAL RESULTS

The most commonly used metric so far in the image description literature has been the BLEU (bilingual evaluation understudy) scores [53], which is a form of precision of word n-grams between generated and reference sentences. The BLEU evaluation metrics, considered as standard evaluation metrics, have been widely used to evaluate image captioning methods in recent years. Following most previous works [41]–[43], we also used BLEU scores as evaluation metrics.

Since there is no open traffic dataset on image captioning, to verify the performance of our method, we performed the

TABLE 2. Performance of our method on flickr30K dataset in terms of BLEU scores.

methods	BLEU-1	BLEU-2	BLEU-3	BLEU-4
NIC[41]	66.3	42.3	27.7	18.3
Att-NIC[43]	66.9	43.9	29.6	19.9
ours	64.7	46.4	31.2	22.4

experiment on our created dataset. To further fairly evaluate our method, we compare ours with the baseline and state-of-art methods on the public dataset Flickr30K. As illustrated in Table 2, our method can achieve comparable performance compared with the two other methods.

We first trained the model on the Flickr30k dataset and fine-tuned the model on our dataset. The performance of our method could be enhanced if our dataset is optimized when it is annotated more suited for traffic scenes. We compared our method with the traffic target detection in ADAS on the same images. As shown in Table 3, in contrast to the existing works that only recognize a specific category of traffic targets, our proposed method not only recognizes general traffic targets affecting traffic factors but also generates a clue for driving suggestion or strategy. Apparently, the information generated by our method is much richer than that of traditional methods, such as TSR, Pedestrian Detection, Vehicle Detection.



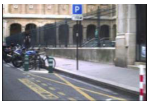



Furthermore, our model can identify more abundant semantic information from the traffic scene (see Table 4). For example, the textual caption generated by our method contains the relative position of traffic targets and the risk factors on the road. Thus, the caption can provide more information for the driver. The results show that the method is more valuable with the abundant semantic information of the traffic scene.

From the results, obviously, our image captioning method can describe general traffic targets and more abundant semantic information at the same time. Since the image caption information is a clue for strategy or suggestion when

TABLE 3. Qualitative comparison with the existing methods of traffic scene understanding.

Raw Image	Pedestrian Detection	Vehicle Detection	Traffic Sign Recognition	Image Caption Detection
			Parking No Stopping Beware Pedestrian	There are Parking sign and pedestrians on the road with cars and motorcycles parked.
			Parking	There are traffic signs, pedestrians, parked cars, and motorcycles on the road.
			No Stopping Beware Pedestrian	There are pedestrians, traffic signs, and parked cars on the road.
			No Passing No Stopping	There are cars, traffic signs, sidewalks, and pedestrians on the road.

TABLE 4. The examples of image captioning generated by Our method.

Test image	Caption	Test image	Caption
	There are some cars on both sides of the road.		There are some cars, some motors and some pedestrians crossing the road.
	There are some motors under the parking sign.		There are two cars, a motor ahead, some pedestrians are crossing the road and the traffic light is red.
	There are some cars on the right sides of the road.		There are some motors, some pedestrians and a zebra crossings ahead.

operating a car, our method can be directly utilized in ADAS. Especially for some emergency conditions, the direct operation strategy or suggestion generated by image captioning is significant for the driver to avoid traffic accidents.

V. CONCLUSION

The ITS is an active topic in the fields of traffic and artificial intelligence. However, most works focus on only recognizing certain targets involved in driving rather than providing driving suggestions. Image captioning originates from Machine Translating [60] and Natural Language Processing(NPL) [43], and mainly applied in human-robot interactions [61], early childhood education [62], information retrieval [63], and visually impaired assistance [64]. In this paper, we attempt to introduce the image captioning method, the new emerging technology in artificial intelligence, to ITS so as to generate high-level semantic information for driving suggestions or strategies. To train the model for ITS,

we created a traffic image dataset for image captioning by associating sentences describing traffic scenes and providing driving suggestions with each image. The experiments on the traffic image dataset demonstrate its effectiveness, which we compare our method with traditional traffic target detection and recognition. The experiment results suggest that this method is a promising solution to ITS.

Currently, there is little research on describing general traffic targets and generating driving suggestions simultaneously. Furthermore, no previous study has been found introducing image captioning to ITS. As shown in our experiments, image captioning gives a new perspective on ITS and has great potential for ITS. However, directly applying the existing model of image captioning to ITS is not an optimal solution. It would seem, therefore, that further research focusing on specific model for the traffic domain is required in the future. In addition, the large-scale image dataset for image captioning needs collecting in ITS.

We will further improve the performance of image captioning in the following research work so that the strategy and prediction of driving could be more accurate. We will investigate how to collect and annotate a larger traffic dataset suitable for future traffic image captioning research.

REFERENCES

- [1] P. S. K. Pandey and R. Kulkarni, "Traffic sign detection for advanced driver assistance system," in *Proc. Int. Conf. Adv. Commun. Comput. Technol. (ICACCT)*, Feb. 2018, pp. 182–185.
- [2] J. Li and Z. Wang, "Real-time traffic sign recognition based on efficient CNNs in the wild," *IEEE Trans. Intell. Transp. Syst.*, vol. 20, no. 3, pp. 975–984, Mar. 2019.
- [3] N. Ben Romdhane, H. Mliki, and M. Hammami, "An improved traffic signs recognition and tracking method for driver assistance system," in *Proc. IEEE/ACIS 15th Int. Conf. Comput. Inf. Sci. (ICIS)*, Jun. 2016, pp. 1–6.
- [4] A. Gomaa, M. M. Abdelwahab, M. Abo-Zahhad, T. Minematsu, and R.-I. Taniguchi, "Robust vehicle detection and counting algorithm employing a convolution neural network and optical flow," *Sensors*, vol. 19, no. 20, p. 4588, Oct. 2019.
- [5] A. Appathurai, R. Sundarasekar, C. Raja, E. J. Alex, C. A. Palagan, and A. Nithya, "An efficient optimal neural network-based moving vehicle detection in traffic video surveillance system," *Circuits, Syst., Signal Process.*, vol. 39, no. 2, pp. 734–756, Feb. 2020.
- [6] V. Murugan and V. R. Vijaykumar, "Automatic moving vehicle detection and classification based on artificial neural fuzzy inference system," *Wireless Pers. Commun.*, vol. 100, no. 3, pp. 745–766, Jun. 2018.
- [7] L. Chen, N. Ma, P. Wang, J. Li, P. Wang, G. Pang, and X. Shi, "Survey of pedestrian action recognition techniques for autonomous driving," *Tsinghua Sci. Technol.*, vol. 25, no. 4, pp. 458–470, Aug. 2020.
- [8] A. Brunetti, D. Buongiorno, G. F. Trotta, and V. Bevilacqua, "Computer vision and deep learning techniques for pedestrian detection and tracking: A survey," *Neurocomputing*, vol. 300, pp. 17–33, Jul. 2018.
- [9] A. Rasouli and J. K. Tsotsos, "Autonomous vehicles that interact with pedestrians: A survey of theory and practice," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 3, pp. 900–918, Mar. 2020.
- [10] S. G. Klauer, T. A. Dingus, V. L. Neale, J. D. Sudweeks, and D. J. Ramsey, "The impact of driver inattention on near-crash/crash risk: An analysis using the 100-car naturalistic driving study data," NHTSA, Richmond, VA, USA, Tech. Rep. DOT/HS/810/584, 2006.
- [11] T. Victor, M. Dozza, J. Bärgerman, C. N. Boda, J. Engström, C. Flannagan, J. D. Lee, and G. Markkula, "Analysis of naturalistic driving study data: Safer glances, driver inattention, and crash risk," *Transp. Res. Board*, Washington, DC, USA, Tech. Rep. S2-S08A-RW-1, 2014.
- [12] M.-Y. Fu and Y.-S. Huang, "A survey of traffic sign recognition," in *Proc. Int. Conf. Wavelet Anal. Pattern Recognit.*, Jul. 2010, pp. 119–124.
- [13] S. Maldonado-Bascon, S. Lafuente-Arroyo, P. Gil-Jimenez, H. Gomez-Moreno, and F. Lopez-Ferreras, "Road-sign detection and recognition based on support vector machines," *IEEE Trans. Intell. Transp. Syst.*, vol. 8, no. 2, pp. 264–278, Jun. 2007.
- [14] R. Timofte, K. Zimmermann, and L. Van Gool, "Multi-view traffic sign detection, recognition, and 3D localisation," *Mach. Vis. Appl.*, vol. 25, no. 3, pp. 633–647, Apr. 2014.
- [15] J. Miura, T. Kanda, S. Nakatani, and Y. Shirai, "An active vision system for on-line traffic sign recognition," *IEICE Trans. Inf. Syst.*, vol. 85, no. 11, pp. 1784–1792, 2002.
- [16] A. de la Escalera, J. M. Armingol, and M. Mata, "Traffic sign recognition and analysis for intelligent vehicles," *Image Vis. Comput.*, vol. 21, no. 3, pp. 247–258, 2003.
- [17] P. Dhar, M. Z. Abedin, T. Biswas, and A. Datta, "Traffic sign detection—A new approach and recognition using convolution neural network," in *Proc. IEEE Region 10 Humanitarian Technol. Conf. (R-HTC)*, Dec. 2017, pp. 416–419.
- [18] S. Jung, U. Lee, J. Jung, and D. H. Shim, "Real-time traffic sign recognition system with deep convolutional neural network," in *Proc. 13th Int. Conf. Ubiquitous Robots Ambient Intell. (URAI)*, Aug. 2016, pp. 31–34.
- [19] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Intell. Signal Process.*, vol. 86, no. 11, pp. 306–351, 2001.
- [20] P. Sermanet and Y. LeCun, "Traffic sign recognition with multi-scale convolutional networks," in *Proc. Int. Joint Conf. Neural Netw.*, Jul. 2011, pp. 2809–2813.
- [21] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2005, pp. 886–893.
- [22] Y. Tsuduki and H. Fujiyoshi, "A method for visualizing pedestrian traffic flow using sift feature point tracking," in *Proc. 3rd Pacific Rim Symp. Adv. Image Video Technol. (PSIVT)*, 2009, pp. 25–36.
- [23] S. Zhang, C. Bauckhage, and A. B. Cremers, "Informed haar-like features improve pedestrian detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2014, pp. 947–954.
- [24] P. Geismann and A. Knoll, "Speeding up hog and LBP features for pedestrian detection by multiresolution techniques," in *Proc. 6th Int. Conf. Adv. Vis. Comput. (ISVC)*, 2010, pp. 243–252.
- [25] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645, Sep. 2010.
- [26] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Region-based convolutional networks for accurate object detection and segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 1, pp. 142–158, Jan. 2016.
- [27] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, Jun. 2017.
- [28] Z. Yang and L. S. C. Pun-Cheng, "Vehicle detection in intelligent transportation systems and its applications under varying environments: A review," *Image Vis. Comput.*, vol. 69, pp. 143–154, Jan. 2018.
- [29] L. Unzueta, M. Nieto, A. Cortes, J. Barandiaran, O. Otaegui, and P. Sanchez, "Adaptive multicue background subtraction for robust vehicle counting and classification," *IEEE Trans. Intell. Transp. Syst.*, vol. 13, no. 2, pp. 527–540, Jun. 2012.
- [30] Y. Jia and C. Zhang, "Front-view vehicle detection by Markov chain Monte Carlo method," *Pattern Recognit.*, vol. 42, no. 3, pp. 313–321, Mar. 2009.
- [31] L.-W. Tsai, J.-W. Hsieh, and K.-C. Fan, "Vehicle detection using normalized color and edge map," *IEEE Trans. Image Process.*, vol. 16, no. 3, pp. 850–864, Mar. 2007.
- [32] J. Chen, W. Xu, W. Peng, W. Bu, B. Xing, and G. Liu, "Road object detection using a disparity-based fusion model," *IEEE Access*, vol. 6, pp. 19654–19663, 2018.
- [33] C. Stauffer and W. E. L. Grimson, "Adaptive background mixture models for real-time tracking," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 1999, pp. 246–252.
- [34] S. Kamkar and R. Safabakhsh, "Vehicle detection, counting and classification in various conditions," *IET Intell. Transp. Syst.*, vol. 10, no. 6, pp. 406–413, Aug. 2016.
- [35] L. Maddalena and A. Petrosino, "Background subtraction for moving object detection in RGBD data: A survey," *J. Imag.*, vol. 4, no. 5, p. 71, May 2018.
- [36] H. Tayara, K. Gil Soo, and K. T. Chong, "Vehicle detection and counting in high-resolution aerial images using convolutional regression neural network," *IEEE Access*, vol. 6, pp. 2220–2230, 2018.
- [37] R. Girshick, "Fast R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 1440–1448.
- [38] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 779–788.
- [39] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 6517–6525.
- [40] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 21–37.
- [41] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: Lessons learned from the 2015 MSCOCO image captioning challenge," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 4, pp. 652–663, Apr. 2017.
- [42] J. Lu, C. Xiong, D. Parikh, and R. Socher, "Knowing when to look: Adaptive attention via a visual sentinel for image captioning," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 3242–3250.
- [43] X. Liu, H. Li, J. Shao, D. Chen, and X. Wang, "Show, tell and discriminate: Image captioning by self-retrieval with partially labeled data," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*, vol. 2018, pp. 353–369.

- [44] P. Anderson, B. Fernando, M. Johnson, and S. Gould, "Spice: Semantic propositional image caption evaluation," in *Proc. 14th Eur. Conf. Comput. Vis.*, 2016, pp. 382–398.
- [45] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang, *Bottom-Up and Top-Down Attention for Image Captioning and VQA*. Amsterdam, The Netherlands: Springer, 2017.
- [46] X. Chen and C. L. Zitnick, "Mind's eye: A recurrent visual representation for image caption generation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 2422–2431.
- [47] J. Gu, G. Wang, J. Cai, and T. Chen, "An empirical study of language CNN for image captioning," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2017, pp. 1231–1240.
- [48] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 3128–3137.
- [49] Y. Wang, Z. Lin, X. Shen, S. Cohen, and G. W. Cottrell, "Skeleton key: Image captioning by skeleton-attribute decomposition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 7378–7387.
- [50] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," *Comput. Sci.*, vol. 2015, pp. 2048–2057, Feb. 2015.
- [51] G. Siogkas, E. Skodras, and E. Dermatas, "Traffic lights detection in adverse conditions using color, symmetry and spatiotemporal information," in *Proc. VISAPP*, 2012, pp. 620–627.
- [52] T. Mikolov, K. Chen, G. S. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," in *Proc. ICLR (Workshop Poster)*, 2013, pp. 1–12.
- [53] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "BLEU: A method for automatic evaluation of machine translation," in *Proc. 40th Annu. Meeting Assoc. Comput. Linguistics (ACL)*, 2001, pp. 311–318.
- [54] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2015, pp. 1–14.
- [55] K. Abrinia, M. Ayati, H. N. Shirvan, A. Abazari, A. H. Tabrizi, M. Shahbazzadeh, Z. Maroufi, M. Akrami, E. Safaei, A. O. Fasakhodi, and A. Panahi, "Advanced driving assistance system using distributed computation on single-board computers," in *Proc. 7th Int. Conf. Robot. Mechatronics (ICRoM)*, Nov. 2019, pp. 392–397.
- [56] T. Hiraoka and Y. Motomura, "Advanced driver assistance system by combination of avoidance guidance based on haptic shared control and automatic collision avoidance control," *Trans. Soc. Automot. Eng. Jpn.*, vol. 50, no. 3, pp. 897–903, 2019.
- [57] M. Pikus and J. Was, "The application of virtual logic models to simulate real environment for testing advanced driving-assistance systems," in *Proc. 24th Int. Conf. Methods Models Autom. Robot. (MMAR)*, Aug. 2019, pp. 544–547.
- [58] P. Morignot, J. P. Rastelli, and F. Nashashibi, "Arbitration for balancing control between the driver and ADAS systems in an automated vehicle: Survey and approach," in *Proc. IEEE Intell. Vehicles Symp.*, Jun. 2014, pp. 575–580.
- [59] G. Andrienko, N. Andrienko, W. Chen, R. Maciejewski, and Y. Zhao, "Visual analytics of mobility and transportation: State of the art and further research directions," *IEEE Trans. Intell. Transp. Syst.*, vol. 18, no. 8, pp. 2232–2249, Aug. 2017.
- [60] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Proc. Adv. Neural Inf. Process. Syst. (NIPS)*, vol. 27, 2014, pp. 3104–3112.
- [61] X. Yang, M.-L. Shyu, H.-Q. Yu, S.-M. Sun, N.-S. Yin, and W. Chen, "Integrating image and textual information in human-robot interactions for children with autism spectrum disorder," *IEEE Trans. Multimedia*, vol. 21, no. 3, pp. 746–759, Mar. 2019.
- [62] K. Khurana and S. Mundada, "Image caption generation a survey," *Int. J. Comput. Sci. Eng.*, vol. 6, no. 3, pp. 256–262, 2018.
- [63] S. Bai and S. An, "A survey on automatic image caption generation," *Neurocomputing*, vol. 311, pp. 291–304, Oct. 2018.
- [64] B. Makav and V. Kilic, "A new image captioning approach for visually impaired people," in *Proc. 11th Int. Conf. Electr. Electron. Eng. (ELECO)*, Nov. 2019, pp. 945–949.



WEI LI received the B.S. and M.S. degrees in computer software and theory from Jilin University, China, in 2002 and 2005, respectively, where she is currently pursuing the Ph.D. degree with the Transportation College. She is also a Lecturer with the School of Computer Science and Engineering, Dalian Minzu University, China. Her current research interests include image understanding, computer vision, and machine learning.



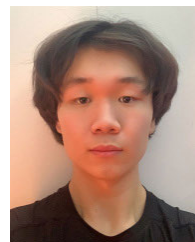
ZHAOWEI QU received the B.S. and M.S. degrees in physics and the Ph.D. degree in traffic control and information engineering from Jilin University, China, in 1984, 1997, and 2009, respectively. He is currently a Professor with the Transportation College, Jilin University. His current research interests include traffic control and video detection of traffic.



HAIYU SONG received the B.S., M.S., and Ph.D. degrees in computer software and theory from Jilin University, in 1996, 2003, and 2012, respectively. He is currently an Associate Professor with the School of Computer Science and Engineering, Dalian Minzu University, China. His current research interests include image understanding, computer vision, and machine learning.



PENGJIE WANG received the B.S. and M.S. degrees in computer software and theory from Jilin University, in 2001 and 2004, respectively, and the Ph.D. degree in computer application from Zhejiang University, in 2011. He is currently a Professor with the School of Computer Science and Engineering, Dalian Minzu University, China. His current research interests include computer vision, virtual reality, and machine learning.



BO XUE is currently pursuing the B.S. degree with the School of Computer Science and Engineering, Dalian Minzu University, China. His current research interests include computer vision, image processing, and deep learning.

• • •