

# CCSPNet-Joint: Efficient Joint Training Method for Traffic Sign Detection Under Extreme Conditions

1<sup>st</sup> Haoqin Hong

*Institute of Innovation & Entrepreneurship, Hanhong College  
Southwest University  
Chongqing, China  
honghaoqin@email.swu.edu.cn*

2<sup>nd</sup> Yue Zhou (✉)

*School of Artificial Intelligence  
Southwest University  
Chongqing, China  
zhouyuenju@163.com*

3<sup>rd</sup> Xiangyu Shu

*College of Computer and Information Science  
Southwest University  
Chongqing, China  
shuxy6263@163.com*

4<sup>th</sup> Xiaofang Hu

*School of Artificial Intelligence  
Southwest University  
Chongqing, China  
huxf@swu.edu.cn*

**Abstract**—Traffic sign detection is an important research direction in intelligent driving. Unfortunately, existing methods often overlook training methods specifically designed for extreme conditions such as fog, rain, and motion blur. Moreover, the end-to-end training strategy for image denoising and object detection models fails to utilize inter-model information effectively. To address these issues, we propose CCSPNet, an efficient feature extraction module based on Contextual Transformer and CNN, capable of effectively utilizing the static and dynamic features of images, achieving faster inference speed and providing stronger feature enhancement capabilities. Furthermore, we establish the correlation between object detection and image denoising tasks and propose a joint training model, CCSPNet-Joint, to improve data efficiency and generalization. Finally, to validate our approach, we create the CCTSDB-AUG dataset for traffic sign detection in extreme scenarios. Extensive experiments have shown that CCSPNet achieves state-of-the-art performance in traffic sign detection under extreme conditions. Compared to end-to-end methods, CCSPNet-Joint achieves a 5.32% improvement in precision and an 18.09% improvement in mAP@.5.

**Index Terms**—Traffic sign detection, Joint training method

## I. INTRODUCTION

Traffic sign detection (TSD) plays a significant role in the field of intelligent driving by providing vital road information to intelligent driving systems, enabling accurate recognition for subsequent decision-making processes[7, 32]. Traffic sign detection algorithms utilize computer vision techniques to rapidly and accurately identify and extract information from images or video data pertaining to traffic signs[24]. The application of such algorithms aids intelligent vehicles in the real-time acquisition of road sign information, enhancing driving safety and overall driving efficiency.

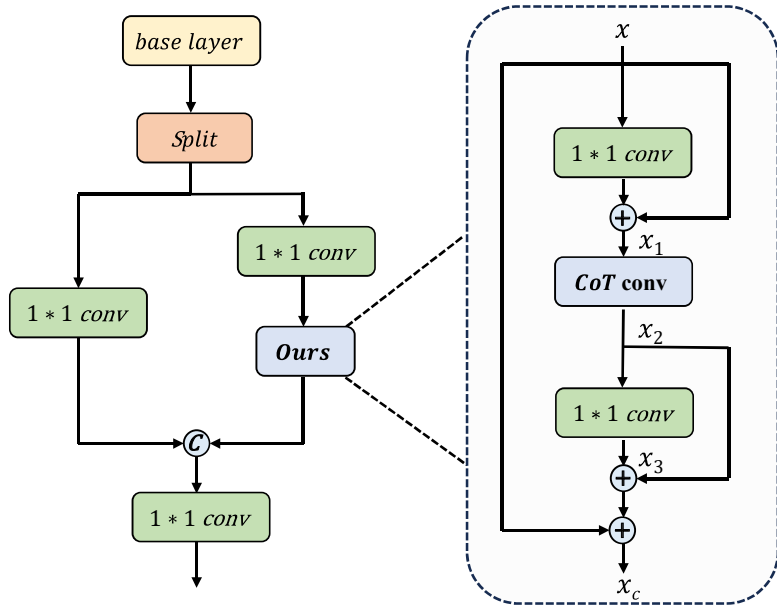
Early traffic sign detection algorithms primarily relied on traditional computer vision techniques, such as image edge

detection [21], image filtering algorithms [10], and morphological image processing [6]. With the advancement of deep learning technology[15], deep learning-based traffic sign detection algorithms have gained widespread adoption. Applying deep learning methods to autonomous driving traffic sign detection improves detection accuracy, speed, and enhances the model's generalization capability. Convolutional neural networks (CNNs) demonstrate strong advantages in traffic sign detection. Compared to other algorithms, CNNs extract more comprehensive features, effectively improving recognition accuracy and speed. CNNs process images through convolution and pooling operations, extracting different features, and automatically learning optimal features, thereby enhancing the accuracy of traffic sign detection. Additionally, CNNs exhibit adaptability, allowing for adjustments in network depth and structure to accommodate different traffic sign detection tasks[14]. However, CNN-based methods alone may not effectively capture global information in complex traffic scenarios. Therefore, it is crucial to design an effective model that can efficiently extract local feature information while improving the capture of global information in traffic scenes.

The backbone network based on Transformer [27] has better capabilities in capturing global information compared to the backbone network based on CNNs. Transformer has been widely used in various visual tasks, including image classification [5, 30], object detection [2], and semantic segmentation[19, 18], among others, for feature extraction. The advantage of Transformer lies in its ability to handle long-range dependencies and parallel computation, effectively capturing crucial information in input sequences. It enables end-to-end training and inference for different tasks. However, existing traffic sign detection models have not fully utilized the advantages of Transformer. In practical traffic scenarios, extracting global features allows for comprehensive consideration of overall shape, color, and other information related to traffic signs. These features possess certain scale and rotation

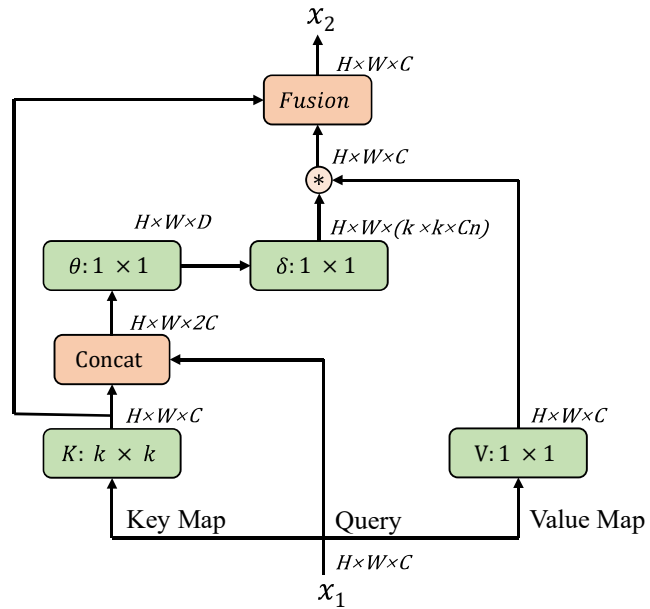
Our code is at <https://github.com/HaoqinHong/CCSPNet-Joint>.

The authors 2<sup>nd</sup> and 4<sup>th</sup> are also with Chongqing Key Laboratory of Brain-inspired Computing and Intelligent Chips.



(a)

(b)



(c)

Traffic sign dataset



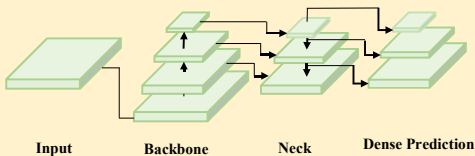
Traffic sign dataset  
for extreme conditions

## Test

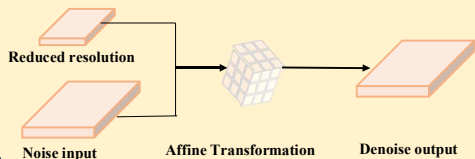
Traffic sign test dataset  
for extreme conditions



### Object Detection



### Image Denoising



Traffic sign image  
with the predicted  
bounding box



## Train

### Back Propagation

$$Loss_1 = \lambda_{coord} L_{coord} + \lambda_{conf} L_{conf} + \lambda_{cls} L_{cls}$$

$$Loss_2 = L_2 = \frac{1}{D} \sum_{i=1}^D \|I_i - J_i\|^2$$

Loss

