# ADD: An Automatic Desensitization Fisheye Dataset for Autonomous Driving

Zizhang Wu[a,*], Xinyuan Chen[a], Hongyang Wei[a,b], Fan Song[a] and Tianhao Xu[a,*]

[a] the Zongmu Technology (Shanghai) Co., Ltd, Building 10, Zhangjiang Artificial Intelligence Island, Lane 55, Chuanhe Road, Pudong New Area, Shanghai, China

[b] the Xinjiang University, College of Software, No. 14 Shengli Road, Urumqi, Xinjiang Uygur Autonomous Region, China

## ABSTRACT

Autonomous driving systems require many images for analyzing the surrounding environment. However, there is fewer data protection for private information among these captured images, such as pedestrian faces or vehicle license plates, which has become a significant issue. In this paper, in response to the call for data security laws and regulations and based on the advantages of large Field of View(FoV) of the fisheye camera, we build the first **A**utopilot **D**esensitization **D**ataset, called **ADD**, and formulate the first deep-learning-based image desensitization framework, to promote the study of image desensitization in autonomous driving scenarios. The compiled dataset consists of 650K images, including different face and vehicle license plate information captured by the surround-view fisheye camera. It covers various autonomous driving scenarios, including diverse facial characteristics and license plate colors. Then, we propose an efficient multitask desensitization network called **DesCenterNet** as a benchmark on the **ADD** dataset, which can perform face and vehicle license plate detection and desensitization tasks. Based on **ADD**, we further provide an evaluation criterion for desensitization performance, and extensive comparison experiments have verified the effectiveness and superiority of our method on image desensitization.

## 1. Introduction

With the massive explosion of big data, the in-depth improvement of algorithms, and the continuous development of hardware technology, artificial intelligence technology represented by deep learning has ushered in a new round of prosperity. Although the current stage is still in the era of quasi-artificial intelligence, and artificial intelligence cannot be granted legal subject qualification, a series of problems caused by artificial intelligence technology still need to be responded to and regulated by law.

The technology and application of artificial intelligence have greatly enhanced the value of data. However, the problems of illegal acquisition, illegal trading, and illegal leakage of data are also particularly prominent and have seriously affected the data security of every member of society. According to the regulation proposed by [2] and the viewpoint of Ling [39], sensitive data generated by intelligent driving vehicles should be masked before updating to other devices. Data use must first undergo desensitization processing, that is, deprivacy processing of data to protect sensitive information, which can not only effectively use data but also ensure the security of data use. Therefore, data desensitization is necessary to filter sensitive information before further analysis.

The rapid development of autonomous driving is accompanied by the generation of many images from the surrounding environment. These images will be served for subsequent tasks, such as pedestrian detection[59, 57, 27, 61, 62, 30], route planning[5, 15, 18, 45], and automatic parking[72, 55, 29, 64, 6, 31]. To achieve this, most of these images will be directly transferred into remote servers or computing clouds for further analysis, typically driven by deep-learning-based image processing and analysis techniques. However, these images frequently include sensitive personal information, such as faces or license plates. This private information is typically used for data transmission, storage, and analysis in current autonomous driving systems, with no privacy protection.

In autonomous driving, pedestrian and vehicle information are vital and fundamental, so ensuring the privacy of pedestrian faces and license plate numbers is especially critical. In this paper, we propose a task for desensitization of face and vehicle license plate information in automatic driving. Data desensitization [66, 14] detects sensitive information and then masks or obscures it before data transmission or analysis. Given the original or background image $I$, a desensitized image $O_{Des}$ can be modelled as a combination of a detection module $I_{Det}$ with Joint desensitization method $M$ and segmentation module $I_{Seg}$,

$$O_{des} = M(I_{seg}, I_{det}) \quad if \quad Iou_{(I_{seg}, I_{det})} > 0.5 \quad (1)$$

The $M$ is a Joint desensitization method where each module $I$ deals with a desensitized region if $IOU_{I_{seg}, I_{det}} > 0.5$, and otherwise belongs to the negative sample. $IOU_{I_{seg}, I_{det}}$ calculates iou with $I_{det}$ by finding the minimum bounding box of the segmented region. $O_{Des}$ can be any image to cover the private information, such as a cartoon icon or grey rectangle square, as shown in Fig. 1.

Current datasets and benchmarks for face and vehicle license plate detection, which frequently include urban, field,

---

[*]Corresponding author.

E-mail address: zizhang.wu@zongmutech.com (Z. Wu), lan.chen@zongmutech.com (X. Chen), weihy@stu.xju.edu.cn (H. Wei), fan.song@zongmutech.com (F. Song), tobias.xu@zongmutech.com (T. Xu).
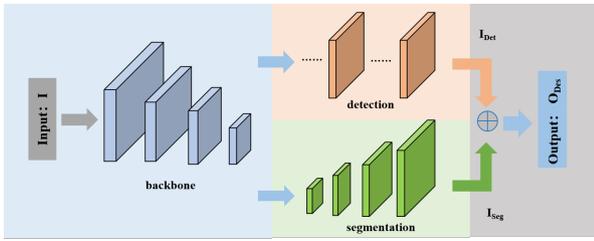
ORCID(s):

**Figure 1:** Desensitization task sketch.

or highway scenes, include Wider Face[68],CCPD[67],CRPD[22], etc. However, existing research lacks a dataset of face and vehicle license plate desensitization in autonomous driving scenarios and does not propose the relevant task of face and vehicle license plate information desensitization. To complement this field, we propose the first **A**utopilot **D**esensitization **D**ataset (**ADD**) and a desensitization framework to promote the study of the desensitization of face and vehicle license plate information. Table 1 summarizes our benchmarks and provides a comparison to existing datasets.

First, we compile a new large-scale fisheye dataset that is provided to facilitate studies dealing with face and vehicle license plate detection and desensitization. In addition, it is subdivided into three sub-datasets: **F**ace **D**esensitization dataset(**FD**), **V**ehicle license plate **D**esensitization dataset(**VD**), and **mix**ed **D**esensitization dataset(**mixD**), each serving different training or evaluation purposes. The face desensitization dataset considers face information of different ages, genders, hairstyles, expressions, and accessories. The vehicle license plate desensitization dataset considers vehicle license plate information of different colors. The hybrid desensitization dataset contains both face and vehicle license plate information, which is more informative. To ensure the diversity of the data, we collected data from different scenes, periods, and attributes. Unlike other face and vehicle license plate recognition datasets[67, 22, 82, 81, 74, 44, 8, 42, 58, 54, 50, 76], ours focuses more on the desensitization of face and vehicle license plate information and provides 650K images of autonomous driving scenarios.

We propose a new practical two-stage multitask desensitization network (MDN) for face and vehicle license plate desensitization tasks to better apply the dataset to real autonomous driving scenarios. MDN can detect the face and vehicle license plate and the desensitization task of the face and vehicle license plate information. We also established a new desensitization evaluation index to verify the desensitization effect. We divide the facial features into three regions and assign different weights to represent the importance of facial features (IOFF), similar to AP [38] and other indicators, to evaluate the desensitization effect.

Our contributions are summarized as follows:

- To the best of our knowledge, the first fisheye dataset for face and vehicle license plate desensitization contains 650K images with adequate face information and vehicle license plate information in the autonomous driving scenario.

- We build the face and vehicle license plate desensitization benchmark on the **ADD** dataset and propose a novel MDN to achieve detection and desensitization tasks.

The rest of the paper is organized as follows. Sec. 2 introduces the relevant methods, while Sec. 3 covers the **ADD** dataset in detail. Then we provide an overview of the proposed desensitization network in Sec. 4, and Sec. 5 describes the comparative analysis. Finally, Sec. 6 concludes this paper.

## 2. Related Work

### 2.1. Image processing-based datasets

Various datasets have been compiled to meet the needs of deep learning models for data. In terms of image processing, there are object detection and segmentation datasets [38, 17, 20], image denoising datasets [1, 77, 73], image deraining datasets [33, 60, 75], image definition datasets [79], etc. The popularity of fisheye cameras has also produced many automatic driving fisheye datasets[71, 51, 35, 69, 48, 52]. These datasets, combined with deep learning models, have greatly improved image processing performance. However, there is no dataset for image desensitization. Therefore, we propose new vehicle license plates and face desensitization datasets for image desensitization. It is very different from common vehicle license plate recognition datasets [67, 22, 82] and face recognition datasets [81, 74, 44, 8, 42, 58, 54, 50, 76]. It is not simply detection and identification. Our desensitization dataset aims to combine the desensitization method with the convolutional neural network to achieve the desensitization of the face and vehicle license plate.

### 2.2. Detection and segmentation-based tasks

With the rapid development of object detection [46, 65, 12, 19] and semantic segmentation technology [41, 80, 3, 78], research on face detection [34, 40, 32], car plate detection [53, 56], and instance segmentation [11, 23, 26, 70] has increased, and increasingly more research on pedestrian detection [43, 28] and vehicle detection [7] based on depth learning has effectively solved the density problems and occlusion in the natural environment. Both faces and license plates are vital signs of pedestrians and vehicles and have unique attributes representing their parents, which have a unique role in identifying specific pedestrians and vehicles. Therefore, it is necessary to use pedestrian and vehicle detection and segmentation methods to enhance the extraction and recognition of face and license plate information. The desensitization task is similar to instance segmentation tasks, such as Mask R-CNN [26] and FCIS [70]. The difference is that it proposes a new desensitization evaluation method for the critical features of the face and license plate, and the effect is particularly prominent in the instance segmentation task.

**Table 1**
Comparison of current State-of-the-Art Benchmarks and Datasets.

**detection and segmentation**

| Datasets | Setting | Sensor type | Ground Label | #categories | avg. #labels/category | #images | Resolution |
|---|---|---|---|---|---|---|---|
| COCO[38] | indoor/outdoor | camera | pedestrian/car | 80 | 18K | 220K | 800*600 |
| VOC[17] | indoor/outdoor | camera | pedestrian/car | 20 | 22K | 44K | 500*375 |
| KITTI[20] | outdoor | camera | pedestrian/car | 2 | 80K | 160K | 1242*375 |
| **ADD** | indoor/outdoor | fisheye camera | pedestrian/car/vehicle license plate/face | 4 | 162K | 650K | 1920*1280 |

**vehicle license plate detection and recognition**

| Datasets | Setting | Sensor type | Ground Label | #categories | #images | Resolution |
|---|---|---|---|---|---|---|
| CCPD[67] | indoor/outdoor | camera | vehicle license plate | 1 | 250K | 720*1160 |
| CPRD[22] | indoor/outdoor | camera | vehicle license plate | 1 | 30K | 640*640 |
| ReID[82] | outdoor | camera | vehicle license plate | 1 | 76K | 200*40 |

**face detection and recognition**

| Datasets | Setting | Sensor type | Ground Label | #categories | #images |
|---|---|---|---|---|---|
| WebFace260M[81] | real | camera | face | 1 | 260M |
| MS-Celeb-1M[24] | real | camera | face | 1 | 10M |
| MF2[42] | real | camera | face | 1 | 4.7M |
| LFW[50] | real | camera | face | 1 | 13K |

**autonomous driving datasets**

| Datasets | Setting | image type | Ground Label | numFisheyeCameras | Resolution |
|---|---|---|---|---|---|
| WoodScape[71] | real | fisheye image | pedestrian/car | 4 | 1280*966 |
| SynWoodScape[51] | synthetic | fisheye image | pedestrian/car | 4 | 1280*966 |
| KITTI 360[35] | real | fisheye image | pedestrian/car | 2 | 1400*1400 |
| FisheyeCityScapes[69] | real | fisheye image | pedestrian/car | 1 | 600*600 |
| THEODORE[48] | synthetic | fisheye image | pedestrian/car | 1 | 1024*1024 |
| OmniScape[52] | synthetic | fisheye image | pedestrian/car | 2 | 1024*1024 |
| **ADD** | real | fisheye image | pedestrian/car/vehicle license plate/face | 4 | 1920*1280 |

## 2.3. Traditional desensitization-based tasks

The explosive growth of data promotes the protection of sensitive information in various forms of data. The traditional recognition method of sensitive information based on rules and regular expressions [4, 10] requires much expert knowledge, has poor mobility, and the recognition pattern is relatively rigid. Unstructured data (text, image, etc.) desensitization technology based on deep learning and machine learning came into being [66, 14]. However, tasks applied to image desensitization are very few [25, 9]. For an image to be desensitized, mage areas having content related to the content of a partial image area that needs to be desensitized are automatically recognized based on the partial image area. Alternatively, an image area to be desensitized is automatically recognized based on a desensitization rule corresponding to a service scenario with which the image to be desensitized is associated. Alternatively, a complete set of image areas to be desensitized is acquired by automatic extension based on a selected partial image area. This paper combines object detection with semantic segmentation and realizes the desensitization task of face and license plate information. Compared with the previous desensitization methods, the combined desensitization method proposed in this paper is more novel and significant.

## 3. Autopilot Desensitization Dataset

In this section, we introduce our **ADD** dataset in detail, including the data collection and preprocessing, annotation protocols, informative statistics, dataset characteristics, and dataset application.

### 3.1. Data Collection and Preprocessing

As shown in Table 2 (a), we collect a total of three cities, three scenarios (parking, street, and highway), and two periods (morning and evening) and capture 200 videos. These videos last from 1 hour to 6 hours with an average of 2 hours, where the density of faces or vehicle license plates is at least 3 (The density of the dataset is shown in Table 2 (a)). We perform frame extraction processing on the video (with a frame rate of 2 FPS after frame extraction), store the video as frame images, and select images containing facial and license plate instances. For high-quality images, we restrict the visible range of the fisheye camera and further remove distorted and blurred face and vehicle license plate images. It is worth noting that we do not distort the images collected by the fisheye camera and filter the images we need through human eye judgment. We also enhance the image data, such as copy-paste[21] that can increase the density of faces/license plates in an image. Taking facial copy-paste as an example, we select a basic background image that may contain 0 pedestrian targets or at least one pedestrian target. By pasting some diverse features of pedestrians (as shown in Table 2 (b)) into
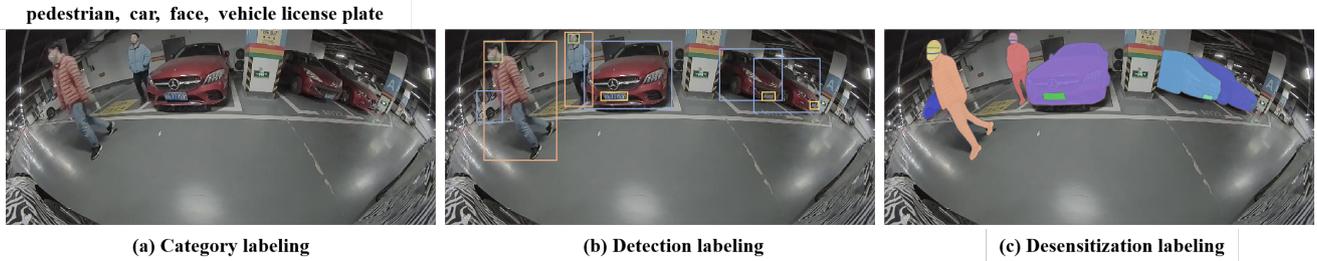
**pedestrian, car, face, vehicle license plate**



(a) Category labeling          (b) Detection labeling          (c) Desensitization labeling

**Figure 2:** Our image annotation pipeline is split into three primary worker tasks: (a) Labelling the categories present in the image, (b) Locating and marking all instances of the labelled categories, and (c) Segmenting each object instance.

the background image, we enrich its pedestrian and facial diversity. The same applies to the license plate copy-paste protocol. Then, the collected data will be divided according to the density of people and vehicles to ensure the diversity of **ADD**. Finally, for the convenience of other studies, we divide the data into face, license plate sub-datasets and a mixed sub-dataset with face and license plate.

### 3.2. Data Annotation

We next describe how we annotated our image collection. Due to our desire to label over 250K images, the design of a cost-efficient yet high-quality annotation pipeline was critical, as Fig.2 shows. Our marking task is distributed to a total of 5 marking engineers in the form of crowdsourcing, and each engineer is responsible for marking a batch of data. The marking personnel are all professional and use the marking tools to mark effectively and, finally, return to the marking engineer for acceptance. General acceptance is to check whether our requirements for labelled pipes are met. The labelling method of this paper is similar to the COCO dataset[38] detection and segmentation labelling method. We performed data cleaning in the collection phase, reducing the time and cost of labelling personnel. The first task in annotating our dataset is determining which object categories are present in each image, Fig. 2(a). In the next stage, we label the location and category of the object, as shown in Fig. 2(b). Finally, we add masks to each target, but in particular, we divide the face mask into three areas: the area above the eyes, the area from the eyes to the nose, and the area below the nose. The regions help us effectively evaluate face desensitization, Fig. 2(c). The basic principle of all labelling is that only visible targets can be labelled, but if the visible target area is less than half of the target itself, it will not be marked. Additionally, we do not cover all pedestrians' continuous moving processes for redundant annotations but select some highly visible targets as tracking objects and generate pedestrian trajectory annotation information and video sequence information (the annotated temporal sequences length of 40). Finally, we collect the first large-scale fisheye desensitization dataset for autopilot application.

### 3.3. Dataset Description

To allow different sub-datasets to be used for the training of different desensitization tasks, we specifically divide

them into three sub-datasets to meet this premise: (1) **F**ace **D**esensitization (**FD**); (2) **V**ehicle license plate **D**esensitization (**VD**); and (3) **mix**ed **D**esensitization (**mixD**), as shown in Table 2(b). A total of more than 650K images comprise three sub-datasets: **FD**, **VD**, and **mixD**, with amounts of 200K, 200K, and 250K, respectively. The image resolution is 1920x1280, and the average object density of each image is 5-6, 3-4, and 8-10 respectively. To complete the face and license plate desensitization task, we decided to use the **mixD** dataset for further desensitization research. To further describe the basic properties of different targets in the **mixD** dataset, we have made detailed statistics in Fig. 3, and all the data meet the requirements of the actual scene as much as possible.

### 3.4. Dataset Characteristics

Our **ADD** dataset exhibits differences from existing datasets in image representation style, scenarios, quantity, and diversity. Below, we elaborate on the four main characteristics of our **ADD** dataset.

**Fisheye image representation.** The **ADD** dataset consists of fisheye images, different from the common pinhole images of public datasets. Fisheye images provide a larger field-of-view (FoV), which is more suitable for close-range and low-lying face and vehicle license plate detection and recognition.

**Specifically for autopilot scenarios.** The **ADD** dataset focuses on the face and vehicle license plate desensitization of high-speed or low-speed parking scenarios in the autonomous driving field, which is different from the natural scene of public datasets. The environmental conditions in autonomous driving scenarios, such as dark and opaque conditions, greatly increase the difficulty of desensitization of face and vehicle license plates. Regarding research on the application of desensitization in autonomous driving scenarios, the **ADD** dataset can promote the study of face desensitization and vehicle license plate desensitization in the real world.

**Great quantity.** Our **ADD** dataset obtains more than 650k data from more than 200-hour autopilot scene video clips. We constantly collect various scenarios, including facial and vehicle license plates, and finally reach more than hundreds of thousands of face and vehicle license plate data.

**Table 2**
Properties and statistics of **ADD**

**(a) Statistics**

| Dataset | Sub-dataset | Image resolution | #images | | | density(/image) |
|---|---|---|---|---|---|---|
| | | | train | val | test | |
| Autopilot Desensitization | Face Desensitization | 1920×1280 | 100k | 60k | 40k | 5∼6 |
| | Vehicle license plate Desensitization | | 100k | 60k | 40k | 3∼4 |
| | Face and vehicle license plate Desensitization(mixD) | | 125k | 75k | 50k | 8∼10 |

**(b) Properties**

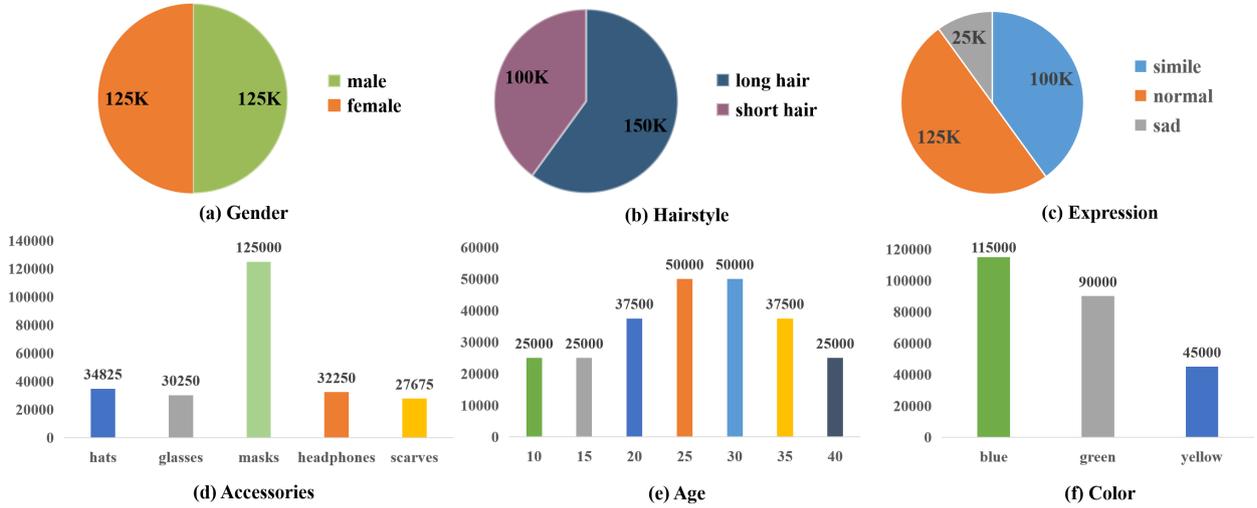| Scenario | Time period | Category | Attribute | Explanation |
|---|---|---|---|---|
| Parking | Daytime | Face | Age | Faces of different ages |
| | | | Gender | Faces of different genders |
| Street | | | Hairstyle | Faces with different hairstyles |
| | Night | | Expression | Faces with different expressions |
| Highway | | | Accessories | Such as hats, glasses, masks, headphones, scarves, etc. |
| | | Vehicle license plate | Color | vehicle license plates in different colors |



**Figure 3:** Statistics of various attributes in the **mixD** subdataset. Attributes of Faces include (a) Gender. (b) Hairstyle. (c) Facial expressions. (d) Different accessories. (e) Age. (f) Different colors are the main attributes of the vehicle license plate.

**High quality and diversity.** Our **ADD** dataset covers three cities, three scenarios from different periods, and different face and vehicle license plate cases. In addition, we selected many people's facial attributes, including age and gender. We also selected many vehicle license plate attributes to include different colors. Finally, we carefully select high-quality, high-resolution images to ensure the dataset's advantages.

### 3.5. Evaluation Metrics
#### 3.5.1. Detection Evaluation

In this paper, we first use AP, AP50, and AP75 to evaluate the detection performance of the model. AP (Average Precision) is an essential measure of the accuracy of the object detection algorithm, by evaluating the object's position in the image and predicting the object class.

$$Recall = \frac{TP}{TP + FN}, \ Precision = \frac{TP}{TP + FP} \quad (2)$$

$$AP = \frac{\sum P_c}{N_{c\_total}}, c \in [0, C) \quad (3)$$

where TP represents **T**rue **P**ositives (samples correctly classified as positive), FN is **F**alse **N**egatives (samples incorrectly classified as negative), and FP is false positives (samples incorrectly classified as positive). $C$ is the number of categories in the dataset and $P_c$ indicates the precision of category $c$, and $N_{c\_total}$ represents the total number of categories in the test set of category $c$. AP is obtained through Precision and Recall. AP50 is the average precision obtained when the detector threshold is more significant than 50%. AP75 is the average precision obtained when the detector threshold exceeds 75%.

**Table 3**
Three regional weights of average importance of facial features (IoFF)

| region | above the eyes | from eyes to nose | below the nose |
|--------|----------------|-------------------|----------------|
| weight | 25% | 50% | 25% |

### 3.5.2. Desensitization Evaluation

To achieve the goal of this task, we use evaluation metrics similar to the Intersection over Union (IoU). The difference is that IoU evaluates every region of the objects with the same weight, and we assign different weights to different regions according to the importance of the facial features and plate letters.

For the face desensitization task, to wipe the critical facial feature against face recovery technologies, we evaluate the importance of different facial features by collecting 1000 faces and dividing them into three different regions: above the eyes, from eyes to nose, and below the nose. We mask each of them with three masking regions and test the similarity between the original face and the masked face using the network proposed by Schroff[49], and the confidence of the similarity is treated as the importance of the facial features, with a confidence ratio close to 1:2:1. Therefore, we assigned weights of the same ratio (as shown in Table 3).

Table 3 shows the three regional weights of the Importance of Facial Features ($IoFF$). It is worth noting that the $IoFF$ positively correlates to the distance to the eyes, and the most critical region to identify a person is the eyes, which takes up to 50% importance to identify a person. To calculate the final score of face desensitization, we sum all the scores of different masked regions to perform the same calculation of mean Average Precision (mAP) as illustrated in equation 5[16].

$$IoFF = \begin{cases} \sum_{i=1}^{N} W_{f_i} * F_{IoU_i} & if \quad c = face \\ IoU & if \quad c = plate \end{cases} \quad (4)$$

$$mIOFF = \begin{cases} \dfrac{\sum_{i=1}^{N} IoFF(i)}{N} & if \quad c = face \\ IoFF & if \quad c = plate \end{cases} \quad (5)$$

where $W_{f_i}$ and $F_{IoU}$ are the $i_{th}$ weight and IOU of the corresponding facial features, respectively, c is the category of the prediction result, and N is the number of face regions.

For the vehicle license plate desensitization task, since the importance of each letter is the same, we assign the same weight to each of them and perform the same calculation as mAP.

## 4. Methodology

In this section, we introduce a new end-to-end face and vehicle license plate desensitization framework, which can be taken as a baseline model for future desensitization research.

### 4.1. Overall Framework

Fig. 5 shows our framework design where we desensitize faces, and vehicle license plates based on the multitask CenterNet [13] model. The backbone network uses the DLA34 network [13]. Different from the CenterNet model, we propose a new desensitization module (in Section 4.2) to achieve face and vehicle license plate desensitization tasks.

First, we resize the input image to 640×480 and fully extract the object features through the DLA34 encoder, whose parameters are shared by the following three tasks. These tasks are comprised of face & plate detection, face & plate segmentation, and pedestrian & vehicle detection, aiming to promote desensitization performance. In addition, to further address the problem of miss & negative desensitization, we also add a post-processing module that will be discussed in detail in the following section.

### 4.2. Desensitization Method

Our desensitization model is a two-stage model. In the first stage, our DesCenterNet serves for the detection of vehicle license plates or faces. The second stage is a small network FCN [41] for segmentation of the desensitization region and then mosaicing them. After that, we can obtain the final desensitized output. The total pipeline of desensitization of the face is shown in Fig. 5.

### 4.2.1. DesCenterNet

In the first stage, DesCenterNet consists of the backbone(DLA34), DLA-UP, Head, Des-UP, and Head. The module DLA-UP and Head server predict cars and pedestrians, the same as CenterNet [13]. Based on them, we add special DES-UP and Head for face and car-plate object detection. As shown in Fig. 6, our heads consist of regression of the heatmap of the center, size of the object, and offset. The desensitization task requires only the classification of the desensitized region and background. This task can be viewed as a two-classification task, and the output channel of the heatmap is set as 2. Furthermore, with multi heads, our DesCenterNet can predict general objects, vehicle license plates, and face regions.

### 4.2.2. Segmentation Module

The second stage is a segmentation model for the pixel classification of the desensitization region(face and vehicle license plate) and background. Its input is a feature map extracted through the DLA34 network, and the network weights are shared. Learning the feature information of facial and license plate regions is of great help in extracting semantic information from sensitive regions. We choose a simple and effective FCN model to find the desensitization region and mosaic it. As shown in Table 6, the FCN achieves the best trade-off of speed and precision. This conforms to the requirements of real-time autonomous driving.

**Figure 4:** Mask R-CNN (top) vs. DesCenterNet (bottom, Hourglass-104). Mask R-CNN has a false detection problem at the edge of the image.
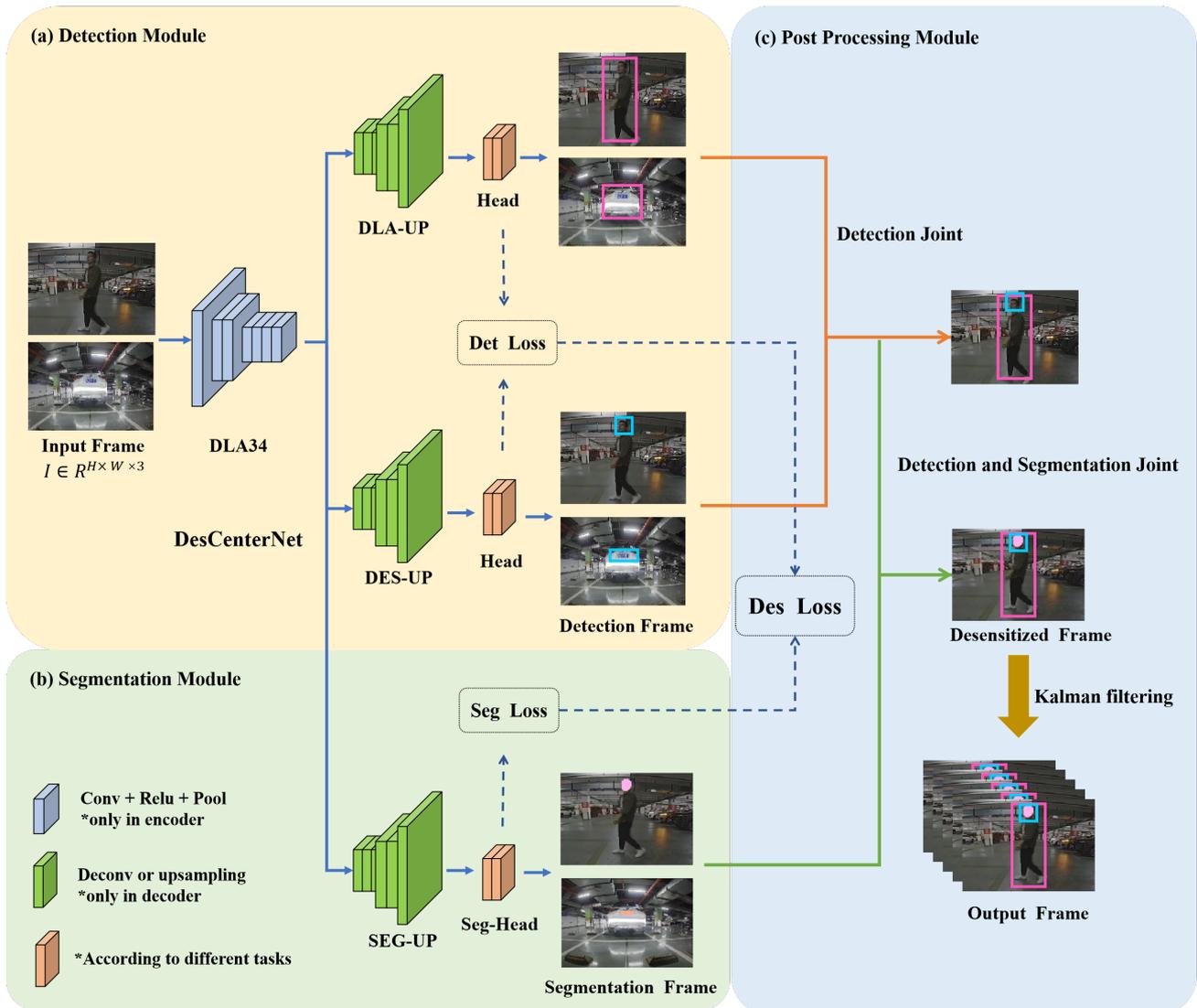


**Figure 5:** Overview of our desensitization model. This is a two-stage model including (a) a detection module, (b) a segmentation module, and (c) a post-processing module. In the first stage, DesCenterNet consists of the backbone(DLA34), DLA-UP, and Head, which are the same as CenterNet [13]. We add special DES-UP and Head for face and car-plate object detection. With multi heads, our DesCenterNet can predict the general object, vehicle license plate, and face region. On the other hand, for the location of the desensitization region, we will mark them by a small segmentation network, including SEG_UP and Seg_Head, and output the frame after desensitization. Finally, we provided a post-processing module that combined the joint desensitization method with Kalman filtering to improve face and vehicle license plate desensitization accuracy.
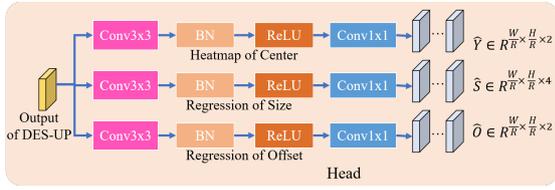
**Figure 6:** Overview of our desensitization head. Inspired by CenterNet, we propose new heads for the detection of desensitization regions. Our heads consist of regression of a heatmap of the center, size of the object, and offset.

### 4.2.3. Post-processing Module

This work aims to conceal the sensitive area as much as possible while minimizing error detection. Meanwhile, 100% desensitized using only the detection and segmentation results can be very hard. As a result, we propose a novel post-processing module to ensure performance. The module consists of a local & global detection sub-module, detection & segmentation sub-module, and Kalman filtering sub-module that utilizes multitasks results from the network and outputs the final result jointly.

Regarding the local & global detection sub-module, we aim to promote the performance by using both local & global results (face & person, plate & vehicle). For instance, we first use both face and pedestrian detection results to overcome the face miss detection, Furthermore, to overcome the mistake detection, we output those face predictions with low confidence by checking if the face boxes are within the person box. The joint desensitization rules are shown in Fig. 7. It is worth noting that we treat the pedestrian detection box and the face detection box as a set of detection pairs and calculate their iou sequences. The red pedestrian detection box and the yellow face detection box in Fig. 7 are considered unqualified detection pairs if their iou is less than 0.5, while the green pedestrian detection box and face detection box are considered qualified detection pairs and serve as the final output.

In addition, the detection & segmentation sub-module further improves the performance. For the detection & segmentation sub module, there are two specific methods. One is to obtain the minimum bounding box from the segmented area and perform iou calculation with the detection area, which can filter out some negative examples. Another approach that We output the final result jointly by considering the confidence of the same area on both tasks. We accept it when both are high and drop it when the results are relatively low. For example, when the detection confidence is greater than 0.7 and the segmentation iou is greater than 0.7. In this paper, we use the first method. Following this, a Kalman filtering algorithm is performed to smooth the result.

### 4.3. Desensitization Loss

Inspired by the loss function of CenterNet, we proposed our loss function of desensitization, which consists of loss of keypoint heatmap, L1 loss of size, and offset. Our model will also output the keypoint heatmap $\hat{Y} \in [0,1]^{\frac{W}{R} \times \frac{H}{R} \times C}$,



(a) Detection Joint     (b) Detection and Segmentation Joint

**Figure 7:** Filtration operation of combined desensitization. Green represents an excellent detection box, red represents a poor detection box, yellow represents a face detection box, and pink represents a face segmentation area. (a) Detection joint desensitization rule between pedestrian detection frame and face detection frame. (b) It is a joint desensitization rule between the detected joint desensitization and face segmentation.

where the $R$ is the output stride, and $C$ is the number of keypoint classes. In the desensitization task, the $C$ is the 3, person's face, car plate, and background. According to the literature [37], the down sampling of the output prediction factor R is 4. We select the focal loss [37] in Equation 6 to compute the distance between the predicted heatmap $\hat{Y}$ and the ground truth heatmap $Y \in [0,1]^{\frac{W}{R} \times \frac{H}{R} \times C}$, on which we transfer each target keypoint using a Gaussian kernel. The $\hat{Y}_{x,y,c} = 1$ indicates the predicted keypoint, while the $\hat{Y}_{x,y,c} = 0$ is the background. Thus, the training object can be viewed as a penalty-reduced pixel logistic regression with focal loss [37]:

$$L_k = \frac{-1}{N} \sum_{xyc} \begin{cases} \left(1 - \hat{Y}_{xyc}\right)^\alpha \log\left(\hat{Y}_{xyc}\right) & \text{if } Y_{xyc} = 1 \\ \left(1 - Y_{xyc}\right)^\beta \left(\hat{Y}_{xyc}\right)^\alpha & \\ \log\left(1 - \hat{Y}_{xyc}\right) & \text{otherwise} \end{cases} \quad (6)$$

where $\alpha$ and $\beta$ are hyperparameters of the focal loss, and $N$ is the number of key points in the input image $I$. Following CenterNet [13], we use $\alpha = 2$ and $\beta = 4$.

The output stride will cause the loss of the decimal part during the down-sampling. To recover this discretization error, we import the local offset loss [37]:

$$L_{off} = \frac{1}{N} \sum_p |\hat{O}_{\tilde{p}} - \left(\frac{p}{R} - \tilde{p}\right)| \quad (7)$$

, where $\hat{O}_{\tilde{p}}$ is predicted offset. Moreover, $p \in \mathcal{R}^2$ of class $c$ is the ground truth point, and its location after down-sampling in output is $\tilde{p} = \lfloor \frac{p}{R} \rfloor$. The difference between them is the target offset.

After the prediction of the center point, we need the regression of the size of objects by the L1 loss:

$$L_{size} = \frac{1}{N} \sum_{k=1}^{N} |\hat{S}_{p_k} - s_k|. \quad (8)$$

where $\hat{S}_{p_k} \in \mathcal{R}^{\frac{W}{R} \times \frac{H}{R} \times 2}$ is the predicted size and $s_k = (x_2^k - x_1^k, y_2^k - y_1^k)$ is the ground truth size of object $k$ and is created by target bounding box $(x_1^k, y_1^k, x_2^k - y_2^k)$.

To predict the region of desensitization, we import the loss of segmentation in equation 9. The ground truth mask is $T_{xyc} \in [0, 1]^{\frac{W}{R} \times \frac{H}{R} \times 3}$, and its predicted semantic map is $\hat{T}_{xyc}$. Following the keypoint map loss, we also use focal loss [37] to balance between hard and easy examples. We set $\alpha = 2$ and $\beta = 4$.

$$
L_{seg} = \frac{-1}{N} \sum_{xyc} \begin{cases} \left(1 - \hat{T}_{xyc}\right)^\alpha \log\left(\hat{T}_{xyc}\right) & \text{if } T_{xyc} = 1 \\ \left(1 - T_{xyc}\right)^\beta \left(\hat{T}_{xyc}\right)^\alpha \\ \quad \log\left(1 - \hat{T}_{xyc}\right) & \text{otherwise} \end{cases}
\tag{9}
$$

After the computation of all losses of multi heads, our model requires the overall loss function:

$$
L = L_k + \lambda_{off} L_{off} + \lambda_{size} L_{size} + \lambda_{seg} L_{seg}.
\tag{10}
$$

where we set $\lambda_{off} = 1$, $\lambda_{size} = 0.1$ and $\lambda_{seg} = 0.5$ in all experiments and desensitization taska. Our network will predict the keypoint heatmap $\hat{Y}$, offset $\hat{O}$, size $\hat{S}$ and pixel-level class $\hat{T}$. They all share a common backbone network, such as DLA34 [13].

# 5. Experiments

In this section, we validate DesCenterNet's strength on our **ADD** dataset. We also perform extensive ablation and contrast experiments to demonstrate the effectiveness of the desensitization framework proposed in this paper.

## 5.1. Experimental Setup

The experimental procedure is set up under the Ubuntu 18.04.06 system, with python 3.8.5 and Pytorch 1.8. We train the **mixD** Dataset with 120 epochs. The learning rate is 0.12, momentum is 0.9, the learning decay rate is 0.001, the optimizer is SGD [47], and the batch size is set to 32. It is worth noting that this experiment uses the mixD dataset, which will be open-sourced after publication with another two sub-datasets.

## 5.2. Main Results

We compare DesCenterNet to the state-of-the-art methods in instance segmentation in Table 4. All instantiations of our model outperform baseline variants of previous state-of-the-art models, including Mask R-CNN[26] and FCIS[70]. There are many improvements on Mask R-CNN, and we also hope there can be more improvements beyond the scope of work. DesCenterNet outputs are visualized in Fig. 8. DesCenterNet achieves good results even under challenging conditions. In Fig. 4, we compare DesCenterNet and Mask R-CNN. Mask R-CNN is prone to false detection in the image edge. DesCenterNet shows no such situation. It is worth noting that all experiments maintain consistent training and

evaluation protocols. When training Mask R-CNN, replace the segmentation head with Seg-UP and Seg-Head to obtain segmentation results of the same resolution size.

## 5.3. Ablation Experiments
### 5.3.1. Comparison with different backbones

We use CenterNet[13] as the baseline and propose DesCenterNet to incorporate an MDN to achieve face and vehicle license plate desensitization tasks. Table 5 shows the results of our validation set with different backbones. The FPS has been tested on our embedding machine, and the maximum computing capacity is approximately 1 Tera Operations Per Second (TOPS), much smaller than a standard Nvidia GPUS (32 TOPS). Hourglass-104 achieves the best accuracy at a relatively good speed, with a 61% mIOFF in 14 FPS. Using ResNet-101, we achieve the minimum real-time requirement for autonomous driving and outstanding performance with 56.2% mIOFF.

### 5.3.2. Comparison with different segmentation models

Here, we study the results from different segmentation modules, including FCN [41], SegNet [3], PSPNet [78], RefineNet [36]. The metric of segmentation is the same as the PASCAL VOC 2012. This result is shown in Table 6, and RefineNet achieved the highest performance with $IOFF_{50}$ of up to 82.4%. However, its speed is only 10 FPS, which is the lowest. In addition, FCN can infer 14 FPS under the same inputs, which is the fastest compared to others. Because the desensitization task requires the inference model's high speed, we choose the simplest FCN as the segmentation model after trade-off performance and speed.

### 5.3.3. The importance of joint desensitization

To cover up the sensitive parts of the face and license plates as much as possible, we proposed a joint desensitization method in Section 4.2 and refined the combined desensitization method into three modules. One is a joint detection module (DJ) for pedestrian and vehicle detection boxes, faces, and license plates. The other is the joint of the result obtained by DJ and the segmentation module (DSJ). The last module is the Kalman filter joint module (KFJ). To verify the effectiveness of different joint methods, we conducted ablation experiments with different joint desensitization methods, as shown in Table 7. In particular, we named only the DJ method w/ DJ, and the method combined with DJ and KFJ was named w/ DJ&KFJ. Others were the same.

We use the method without the combined desensitization module as the baseline for ablation experiments. From Table 7, we can see that the joint desensitization module's method has achieved good performance improvement. With only the DJ methods, mIOFF and $IOFF_{75}$ increased by 0.1% and 0.5%, respectively, but the improvement effect was small. Our analysis shows that the return accuracy of the detection box is significant, and the overlap rate of the detection box is significant, but the desensitization accuracy is limited. In the method of using DSJ, we can see that mIOFF, $IOFF_{50}$, and $IOFF_{75}$ have increased by 5.0%, 0.7%, and 1.0% compared with baseline, respectively. The facial and

**Figure 8:** Visualized results of our desensitization method on **ADD's mixD** test set. Green represents the pedestrian detection frame, orange represents the car detection frame, red represents the license plate detection frame, and yellow represents face detection frame. Other masks, such as pedestrian, car, face, and vehicle license plate, are masks after segmentation.

**Table 4**

Our method DesCenterNet $vs$ Mask R-CNN with bounding boxes AP and IOFF.

| Method | backbone | $AP$ | $AP_{50}^{bb}$ | $AP_{75}^{bb}$ | $mIOFF$ | $IOFF_{50}$ | $IOFF_{75}$ |
|---|---|---|---|---|---|---|---|
| FCIS[70] +OHEM | ResNet-101-C5-Dilated | 39.1 | 59.8 | 39.8 | 42.2 | 56.2 | 45.7 |
| FCIS+++[70] +OHEM | ResNet-101-C5-Dilated | 42.2 | 62.7 | 44.9 | 48.1 | 62.1 | 51.2 |
| Mask R-CNN [26] | ResNet-101-C4 | 42.1 | 64.8 | 45.9 | 49.2 | 61.2 | 52.7 |
| Mask R-CNN [26] | ResNet-101-FPN | 43.2 | 65.6 | 46.1 | 52.1 | 65.1 | 55.1 |
| Mask R-CNN [26] | ResNeXt-101-FPN | 45.0 | 66.9 | 47.9 | 54.1 | 68.8 | 59.2 |
| **DesCenterNet** | ResNet-101 | 46.1 | 66.9 | 49.1 | 56.2 | 70.4 | 62.7 |
| **DesCenterNet** | Hourglass-104 | **48.7** | **70.1** | **49.2** | **61.0** | **70.8** | **65.7** |

**Table 5**

Speed/accuracy trade-off for different backbones on the validation set.

| Backbone | $mIOFF$ | $IOFF_{50}$ | $IOFF_{75}$ | $FPS$ |
|---|---|---|---|---|
| Hourglass-104 | 61 | 70.8 | 65.7 | 14 |
| DLA-34 | 55 | 72.9 | 60.3 | 52 |
| ResNet-101 | 56.2 | 70.4 | 62.7 | 60 |
| ResNet-18 | 51.2 | 60.4 | 55.7 | 153 |

**Table 6**

The comparison of different segmentation models on our desensitization dataset.

| Dataset | Method | $IOFF_{50}$ | $FPS$ |
|---|---|---|---|
| mixD | FCN-8s [41] | 70.8 | 14 |
| | SegNet [3] | 69.82 | 12 |
| | DeepLab [80] | 70.6 | 7 |
| | PSPNet [78] | 67.4 | 6 |
| | RefineNet [36] | 82.4 | 10 |

license plate characteristics obtained by the segmentation module can effectively filter the noise in the predictive information, improving the recognition rate and desensitization

accuracy of sensitive information. We used DSJ and KFJ in combination and found that $IOFF_{75}$ increased by 3.7%,
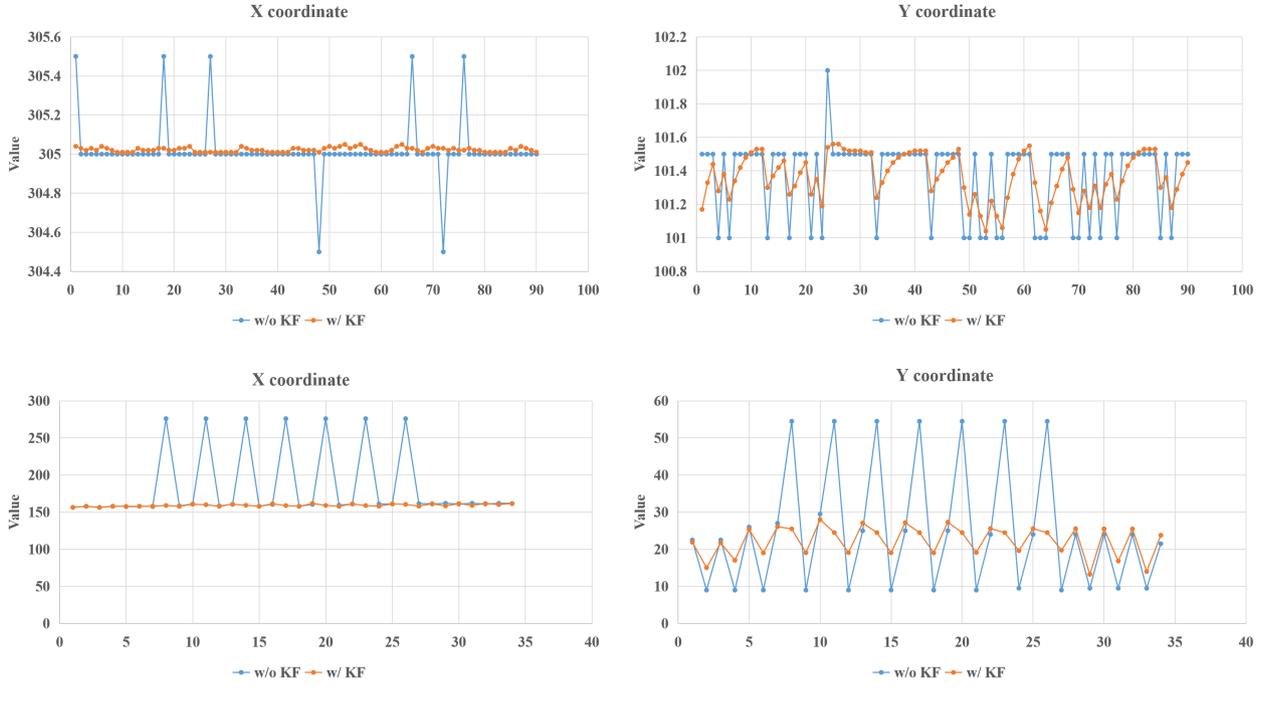
**Figure 9:** The center point coordinate renderings of with Kalman Filter(w/ KF) and without Kalman Filter(w/o KF), the first row is the rendering of the center point of the vehicle license plate, and the second row is the rendering of the center point of the face.

**Table 7**
Different combined desensitization methods of ablation experiments.

| Dataset | DJ | DSJ | KFJ | mIOFF | $IOFF_{50}$ | $IOFF_{75}$ |
|---------|----|-----|-----|-------|-------------|-------------|
|         | -  | -   | -   | 60.9  | 70.8        | 65.3        |
|         | ✓  | -   | -   | 61.0  | 70.8        | 65.8        |
| mixD    | -  | ✓   | -   | 65.9  | 71.5        | 66.3        |
|         | -  | -   | ✓   | 65.3  | 71.0        | 66.0        |
|         | ✓  | -   | ✓   | 66.0  | 71.4        | 66.4        |
|         | -  | ✓   | ✓   | **69.2** | **74.0** | **70.1**    |

reaching the highest score. DSJ can effectively detect sensitive information, and KFJ can optimize the missed detection rate, significantly improving desensitization performance.

### 5.3.4. The importance of Kalman filter

**Indicator improvement.** As shown in Table 7, in w/ KFJ method, mIOFF, $IOFF_{50}$, and $IOFF_{75}$ increased by 4.4%, 0.2%, and 0.7% compared with Baseline, respectively. It can be seen that KFJ can effectively optimize the missed inspection effect.

However, compared to w/ DSJ method, mIOFF was reduced by 0.6%, and mIOFF was increased by 4.3%compared to the w/ DJ method. DJ can detect face and license plate information but also cover useless background information. DSJ can enhance the detection and desensitization effect of face and license plates from pixel-level and characteristic levels. KFJ is a relatively lightweight tracking method. Finally, we combined KFJ, DJ, and DSJ respectively and

found that w/ DJ&KFJ and w/ DSJ&KFJ increased by 0.7% and 3.9% compared with the mIOFF of w/ KFJ, respectively, to the highest. DSJ can effectively detect sensitive information, and KFJ can optimize the missed inspection rate. The two complement each other, which significantly improves desensitization performance.

**Stability and smoothness.** We highlight the critical performance of Kalman filtering in the post-processing in Fig. 9. Among the desensitization results obtained, the post-treatment of Kalman filtering makes the deviation of the center point of the face and the center of the license plate smaller, and the desensitization performance is more stable.

### 5.4. Different Kalman filter previous frame selection

To promote the smoothness of the desensitization task, we apply Kalman filtering proposed by [63], which aims to predict the target's location using previous frames when the network miss-detects the objects. In this section, we study how many previous frames should be included to predict the target because they determine the target motion, while too many frames cause false detection cases and insufficient frames caused missed detection cases.

As shown in Table 8, we conducted ablation experiments with different previous frames and found that the best performance was achieved when the frame number was 4. We think that increasing the number of previous frames increases the model's robustness and the false detection rate. When the number of frames exceeds 4, the possible invalid frames will

**Table 8**
Different Karman Filter previous frames selection.

| Dataset | Previous Frames | $mIOFF$ | $IOFF_{50}$ | $IOFF_{75}$ |
|---------|-----------------|---------|-------------|-------------|
|         | 2               | 67.1    | 71.9        | 67.1        |
|         | 3               | 68.9    | 73.5        | 68.3        |
| mixD    | 4               | **69.2** | **74.0**   | **70.1**    |
|         | 5               | 66.9    | 72.5        | 67.1        |
|         | 6               | 65.1    | 70.9        | 65.2        |
|         | 7               | 63.2    | 70.0        | 63.9        |

affect the desensitization effect, decreasing the IOFF index.

## 6. Discussion

### 6.1. Using desensitization data to train other tasks

It is worth noting that training using desensitization data also has a positive effect on the performance of pedestrian and vehicle-related tasks since those targets are partially masked (face and license plate) and according to the copy-paste method proposed by [21] that these masked objects will promote the robustness of the model. Furthermore, to achieve more effective desensitization tasks, we deploy a multitask training procedure to ensure that the shared backbone can still extract the correct features of sensitive areas.

### 6.2. Performance of desensitization

According to the regulation [2], the sensitive information is erased as long as masks cover more than 50% of the area of the sensitive objects. This gives a fault tolerance to the algorithm. The performance of our proposed desensitization method mainly depends on the effectiveness of detection and segmentation, as the subsequent joint desensitization methods refer to their results, which has a high demand for detection and segmentation networks.

### 6.3. Compare with instance segmentation

The desensitization method is similar to the instance segmentation method, but the difference lies in the desensitization region division and weight setting of the face in the label (as shown in Table 3). In addition, we also propose new evaluation indicators for desensitization tasks and the joint desensitization method mentioned in the post-processing module. In fact, our method will obtain the desensitized area and then cover it with cartoon icons or gray rectangles. We believe that desensitization representations can be freely chosen, and there will be more details to be defined and refined for future desensitization tasks.

### 6.4. Compare with other datasets

Autonomous driving data desensitization is a relatively new topic, which requires that autonomous driving vehicles can effectively process private data. ADD is the first fisheye dataset for the study of autonomous driving data desensitization. It has a strong purpose and is different from other open data sets. Of course, other datasets can also achieve

data desensitization tasks theoretically by additional annotations and methods, but this can not deny our original intention and contribution to the research of autonomous driving data desensitization. We also hope that more researchers will pay attention to the research direction of autonomous driving data desensitization.

## 7. Conclusion

This paper presents the first autopilot desensitization dataset, dubbed **ADD**, which includes face and license plate information. By providing face and vehicle license plates with distinct characteristics, **ADD** intends to aid the industry in the development of a safer advanced autonomous driving system. In addition, we proposed a new framework for desensitization, and extensive experiments demonstrated the superiority and generality of our method. In the future, we hope that **ADD** will stimulate more correlation with desensitization research and promote the protection of more sensitive information for industry applications.

## References

[1] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 1692–1700, 2018.

[2] China Automobile Industry Association. Technical requirements and methods of video and image desensitization for vehicle transmission. In *T/CAAMTB 77-2022*, 2022.

[3] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 39(12):2481–2495, 2017.

[4] David E Bakken, R Rarameswaran, Douglas M Blough, Andy A Franz, and Ty J Palmer. Data obfuscation: Anonymity and desensitization of usable data sets. *IEEE Security & Privacy*, 2(6):34–41, 2004.

[5] Hannah Bast, Daniel Delling, Andrew Goldberg, et al. Route planning in transportation networks. pages 19–80, 2016.

[6] Mahdi Biparva, David Fernández-Llorca, Ruben Izquierdo Gonzalo, et al. Video action recognition for lane-change classification and prediction of surrounding vehicles. *IEEE Transactions on Intelligent Vehicles*, 2022.

[7] Abdelmalek Bouguettaya, Hafed Zarzour, Ahmed Kechida, et al. Vehicle detection from uav imagery with deep learning: a review. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

[8] Qiong Cao, Li Shen, Weidi Xie, et al. Vggface2: A dataset for recognising faces across pose and age. In *IEEE international conference on automatic face & gesture recognition (FG)*, pages 67–74, 2018.

[9] Wilfred A Cassell. Desensitization therapy for body image anxiety. *Canadian Psychiatric Association Journal*, 22(5):239–242, 1977.

[10] Malu Castellanos, Bin Zhang, Ivo Jimenez, Perla Ruiz, Miguel Durazo, Umeshwar Dayal, and Lily Jow. Data desensitization of customer data for use in optimizer performance experiments. In *2010 IEEE 26th International Conference on Data Engineering (ICDE 2010)*, pages 1081–1092. IEEE, 2010.

[11] Bowen Cheng, Omkar Parkhi, and Alexander Kirillov. Pointly-supervised instance segmentation. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 2617–2626, 2022.

[12] Xudong Dong, Shuai Yan, and Chaoqun Duan. A lightweight vehicles detection network model based on yolov5. *Engineering Applications of Artificial Intelligence*, 113:104914, 2022.

[13] Kaiwen Duan, Song Bai, Lingxi Xie, et al. Centernet: Keypoint

triplets for object detection. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 6569–6578, 2019.

[14] Anday Duru, İsmail Rakıp Karaş, Fatih Karayürek, and Aydın Gülses. A deep learning approach for classification of dentinal tubule occlusions. *Applied Artificial Intelligence*, 36(1):2094446, 2022.

[15] Jochen Eisner, Stefan Funke, and Sabine Storandt. Optimal route planning for electric vehicles in large networks. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2022.

[16] Mark Everingham, SM Eslami, Luc Van Gool, et al. The pascal visual object classes challenge: A retrospective. *International journal of computer vision*, 111(1):98–136, 2015.

[17] Mark Everingham, Luc Van Gool, Christopher KI Williams, et al. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.

[18] Mary E Frame, Michaela Schwing, Samuel Johnston, et al. Route planning decisions: evaluating reliance on spatial heuristics under risk. *Spatial Cognition & Computation*, pages 1–26, 2022.

[19] Xianping Fu, Jinjia Peng, Guangqi Jiang, and Huibing Wang. Learning latent features with local channel drop network for vehicle re-identification. *Engineering Applications of Artificial Intelligence*, 107:104540, 2022.

[20] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *Conference on Computer Vision and Pattern Recognition(CVPR)*, pages 3354–3361, 2012.

[21] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, et al. Simple copy-paste is a strong data augmentation method for instance segmentation. In *Conference on Computer Vision and Pattern Recognition(CVPR)*, pages 2918–2928, 2021.

[22] Yanxiang Gong, Linjie Deng, Shuai Tao, et al. Unified chinese license plate detection and recognition with high efficiency. *Journal of Visual Communication and Image Representation*, pages 103541–103543, 2022.

[23] Wenchao Gu, Shuang Bai, and Lingxing Kong. A review on 2d instance segmentation based on deep neural networks. *Image and Vision Computing*, page 104401, 2022.

[24] Yandong Guo, Lei Zhang, Yuxiao Hu, et al. Ms-celeb-1m: A dataset and benchmark for large-scale face recognition. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 87–102, 2016.

[25] Alan Hassard. Eye movement desensitization of body image. *Behavioural and Cognitive Psychotherapy*, 21(2):157–160, 1993.

[26] Kaiming He, Georgia Gkioxari, Piotr Dollár, et al. Mask r-cnn. In *International Conference on Computer Vision(ICCV)*, pages 2961–2969, 2017.

[27] Yifan Jiao, Hantao Yao, and Changsheng Xu. San: selective alignment network for cross-domain pedestrian detection. *IEEE transactions on image processing*, 30:2155–2167, 2021.

[28] Yifan Jiao, Hantao Yao, and Changsheng Xu. San: selective alignment network for cross-domain pedestrian detection. *IEEE transactions on image processing*, 30:2155–2167, 2021.

[29] Bai Li, Zhuyan Yin, Yakun Ouyang, et al. Online trajectory replanning for sudden environmental changes during automated parking: A parallel stitching method. *IEEE Transactions on Intelligent Vehicles*, 2022.

[30] Guofa Li, Zefeng Ji, Xingda Qu, et al. Cross-domain object detection for autonomous driving: A stepwise domain adaptative yolo approach. *IEEE Transactions on Intelligent Vehicles*, 2022.

[31] Liang Li, Ming Yang, Hao Li, et al. Robust localization for intelligent vehicles based on compressed road scene map in urban environments. *IEEE Transactions on Intelligent Vehicles*, 2022.

[32] Menghan Li, Bin Huang, and Guohui Tian. A comprehensive survey on 3d face recognition methods. *Engineering Applications of Artificial Intelligence*, 110:104669, 2022.

[33] Siyuan Li, Iago Breno Araujo, Wenqi Ren, et al. Single image deraining: A comprehensive benchmark analysis. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 3838–3847, 2019.

[34] Xin Li, Shenqi Lai, and Xueming Qian. Dbcface: Towards pure convolutional neural network face detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 32(4):1792–1804, 2021.

[35] Yiyi Liao, Jun Xie, and Andreas Geiger. Kitti-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[36] Guosheng Lin, Anton Milan, Chunhua Shen, et al. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 1925–1934, 2017.

[37] Tsung-Yi Lin, Priya Goyal, Ross Girshick, et al. Focal loss for dense object detection. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 2980–2988, 2017.

[38] Tsung-Yi Lin, Michael Maire, Serge Belongie, et al. Microsoft coco: Common objects in context. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 740–755, 2014.

[39] Justin Ling. Autonomous vehicles join the list of us national threats. In *WIRED*, 2022.

[40] Wei Liu, Irtiza Hasan, and Shengcai Liao. Center and scale prediction: Anchor-free approach for pedestrian and face detection. *Pattern Recognition*, page 109071, 2022.

[41] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 3431–3440, 2015.

[42] Aaron Nech and Ira Kemelmacher-Shlizerman. Level playing field for million scale face recognition. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 7044–7053, 2017.

[43] Sweta Panigrahi and USN Raju. Pedestrian detection based on handcrafted features and multi-layer feature fused-resnet model. *International Journal on Artificial Intelligence Tools(IJAIT*, 30(05):2150028, 2021.

[44] Omkar M Parkhi, Andrea Vedaldi, and Andrew Zisserman. Deep face recognition. 2015.

[45] Irina V Pustokhina, Denis A Pustokhin, E Laxmi Lydia, et al. Energy efficient neuro-fuzzy cluster based topology construction with metaheuristic route planning algorithm for unmanned aerial vehicles. *Computer Networks*, 196:108214–108215, 2021.

[46] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Conference on Neural Information Processing Systems*, 28, 2015.

[47] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint arXiv:1609.04747*, 2016.

[48] Tobias Scheck, Roman Seidel, and Gangolf Hirtz. Learning from theodore: A synthetic omnidirectional top-view indoor dataset for deep transfer learning. In *Winter Conference on Applications of Computer Vision(WACV)*, pages 943–952, 2020.

[49] F. Schroff, D. Kalenichenko, and J. Philbin. Facenet: A unified embedding for face recognition and clustering. *IEEE*, 2015.

[50] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 815–823, 2015.

[51] Ahmed Rida Sekkat, Yohan Dupuis, Varun Ravi Kumar, et al. Synwoodscape: Synthetic surround-view fisheye camera dataset for autonomous driving. *arXiv preprint arXiv:2203.05056*, 2022.

[52] Ahmed Rida Sekkat, Yohan Dupuis, Pascal Vasseur, et al. The omniscape dataset. In *International Conference on Robotics and Automation (ICRA)*, pages 1603–1608, 2020.

[53] Mojtaba Shahidi Zandi and Roozbeh Rajabi. Deep learning based framework for iranian license plate detection and recognition. *Multimedia Tools and Applications*, 81(11):15841–15858, 2022.

[54] Yaniv Taigman, Ming Yang, Marc'Aurelio Ranzato, et al. Web-scale training for face identification. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 2746–2754, 2015.

[55] Muhammad Hans Tobi. Design of automatic parking access system based on internet of things (iot). *Brilliance: Research of Artificial Intelligence*, 2(2):62–65, 2022.

[56] Ruben Jose Tom, Awanit Kumar, Syed Basha Shaik, et al. Car license plate detection and recognition using modified u-net deep learning model. In *International Conference on Smart Structures and Systems (ICSSS)*, pages 01–06, 2022.

[57] Napat Wanchaitanawong, Masayuki Tanaka, Takashi Shibata, et al. Multi-modal pedestrian detection with large misalignment based on modal-wise regression and multi-modal iou. In *International Conference on Machine Vision and Applications (ICMVA)*, pages 1–6, 2021.

[58] Fei Wang, Liren Chen, Cheng Li, et al. The devil of face recognition is in the noise. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 765–780, 2018.

[59] Minjun Wang, Houjin Chen, Yanfeng Li, et al. Multi-scale pedestrian detection based on self-attention and adaptively spatial feature fusion. *IET Intelligent Transport Systems*, 15(6):837–849, 2021.

[60] Tianyu Wang, Xin Yang, Ke Xu, et al. Spatial attentive single-image deraining with a high quality real rain dataset. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 12270–12279, 2019.

[61] Xinlong Wang, Tete Xiao, Yuning Jiang, et al. Repulsion loss: Detecting pedestrians in a crowd. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 7774–7783, 2018.

[62] Hongyang Wei, Qianqian Zhang, Jingjing Han, et al. Sarnet: Spatial attention residual network for pedestrian and vehicle detection in large scenes. *Applied Intelligence*, pages 1–16, 2022.

[63] Gregory F Welch. Kalman filter. *Computer Vision: A Reference Guide*, pages 1–3, 2020.

[64] Jingda Wu, Zhiyu Huang, and Chen Lv. Uncertainty-aware model-based reinforcement learning: Methodology and application in autonomous driving. *IEEE Transactions on Intelligent Vehicles*, 2022.

[65] Yue Wu, Yinpeng Chen, Lu Yuan, Zicheng Liu, Lijuan Wang, Hongzhi Li, and Yun Fu. Rethinking classification and localization for object detection. In *Computer Vision and Pattern Recognition*, pages 10186–10195. IEEE, 2020.

[66] Nan Xiang, Xiongtao Zhang, Yajie Dou, Xiangqian Xu, Kewei Yang, and Yuejin Tan. High-end equipment data desensitization method based on improved stackelberg gan. *Expert Systems with Applications*, 180:114989, 2021.

[67] Zhenbo Xu, Wei Yang, Ajin Meng, et al. Towards end-to-end license plate detection and recognition: A large dataset and baseline. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 255–271, 2018.

[68] Shuo Yang, Ping Luo, Chen-Change Loy, et al. Wider face: A face detection benchmark. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 5525–5533, 2016.

[69] Yaozu Ye, Kailun Yang, Kaite Xiang, et al. Universal semantic segmentation for fisheye urban driving images. *Conference on Systems, Man, and Cybernetics (SMC)*, 2020.

[70] Li Yi, Qi Haozhi, Dai Jifeng, et al. Fully convolutional instance-aware semantic segmentation. In *Conference on Computer Vision and Pattern Recognition(CVPR)*, pages 2359–2367, 2017.

[71] Senthil Yogamani, Ciarán Hughes, Jonathan Horgan, et al. Woodscape: A multi-task, multi-camera fisheye dataset for autonomous driving. In *International Conference on Computer Vision(ICCV)*, pages 9308–9318, 2019.

[72] Leiyan Yu, Xianyu Wang, Zeyu Hou, et al. Path planning optimization for driverless vehicle in parallel parking integrating radial basis function neural network. *Applied Sciences*, 11(17):8178–8179, 2021.

[73] Huanjing Yue, Cong Cao, Lei Liao, et al. Supervised raw video denoising with a benchmark dataset on dynamic scenes. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 2301–2310, 2020.

[74] Kaihao Zhang, Dongxu Li, Wenhan Luo, et al. Edface-celeb-1 m: Benchmarking face hallucination with a million-scale dataset. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2022.

[75] Kaihao Zhang, Dongxu Li, Wenhan Luo, et al. Enhanced spatio-temporal interaction learning for video deraining: A faster and better framework. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2022.

[76] Yaobin Zhang, Weihong Deng, Mei Wang, et al. Global-local gcn: Large-scale label noise cleansing for face recognition. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 7731–7740, 2020.

[77] Yide Zhang, Yinhao Zhu, Evan Nichols, et al. A poisson-gaussian denoising dataset with real fluorescence microscopy images. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 11710–11718, 2019.

[78] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, et al. Pyramid scene parsing network. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 2881–2890, 2017.

[79] Zhuoran Zheng, Wenqi Ren, Xiaochun Cao, et al. Ultra-high-definition image dehazing via multi-guided bilateral learning. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 16180–16189, 2021.

[80] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, et al. *Unet++: A nested u-net architecture for medical image segmentation*, pages 3–11. 2018.

[81] Zheng Zhu, Guan Huang, Jiankang Deng, et al. Webface260m: A benchmark unveiling the power of million-scale deep face recognition. In *Proceedings of Computer Vision and Pattern Recognition (CVPR)*, pages 10492–10502, 2021.

[82] Jakub Špaňhel, Jakub Sochor, Roman Juránek, et al. Holistic recognition of low quality license plates by cnn using track annotated data. In *IEEE International Conference on Advanced Video and Signal Based Surveillance*, pages 1–6, 2017.