# Image object detection and semantic segmentation based on convolutional neural network

Laigang Zhang[1] · Zhou Sheng[2] · Yibin Li[3] · Qun Sun[1] · Ying Zhao[1] · Deying Feng[1]

## Abstract

In this paper, an unsupervised co-segmentation algorithm is proposed, which can be applied to the image with multiple foreground objects simultaneously and the background changes dramatically. The color edge image in RGB space is extracted for semantic extraction. This method can effectively distinguish foreground and background by recursively modeling the appearance distribution of pixels and regions. The coherence of image foreground and background model is enhanced by using the correlation between different image regions and image interior. Experimental results show that deep convolutional neural network can effectively realize semantic classification of scene images by end-to-end feature learning and achieve accurate semantic segmentation of scene images.

**Keywords** CNN · AdaBoost · Image object detection · Semantic segmentation

## 1 Introduction

Scene understanding is a hot topic in the field of computer vision and artificial intelligence. Its research results have been widely used in many fields such as robot navigation,

✉ Zhou Sheng
  yushin@stu.cqe.edu.cn

  Laigang Zhang
  zhanglaigang@lcu.edu.cn

  Yibin Li
  liyb@edu.cn

  Qun Sun
  sunqun@lcu.edu.cn

  Ying Zhao
  zhaoying@lcu.edu.cn

  Deying Feng
  fengdeying@lcu.edu.cn

[1] School of Mechanical and Automotive Engineering, Liaocheng University, Liaocheng, Shandong, China

[2] School of Management, Wuhan Donghu University, Wuhan, China

[3] Shandong University, Jinan, Shandong, China

network search, security monitoring, medical care and so on. Various branch tasks of scene understanding, such as target detection, image semantic segmentation and so on, have made a breakthrough in recent years, but there are still many shortcomings. For example, it is difficult to obtain reliable and robust features for dynamic target classification in the scene because of the deformation of the target itself and the interference of external factors [1].

Object detection is a wide area problem in the field of computer and machine vision. Complex background also increases the range of challenges and errors as well as problems. Many algorithms used for object detection have difficulty in matching the influence of occlusion and pixel moment. Hence, a highly robust algorithm, ORBTRIAN, is proposed for low-resolution images, and gradient enhancement machine learning algorithm is used to detect ORB. This work has been compared with the technology based on AdaBoost and Surf. Analysis shows that the performance of the early model improved to 3.8%. The feature points extracted from the ORB method are further processed to further reduce the processing. Only those points farthest from its centroid triangle are selected, and only one feature point is selected. The result was about 28, much faster than earlier calculations [2, 3]. Tree-based GB

has been implemented in this algorithm. As more feature points need to be identified and more classes are needed, the execution of the computation requires unreasonable effort and time. Therefore, their algorithm is used to allocate some nearby classes at the same level to reduce the number of tree nodes. Experiments show that the overall performance of the proposed algorithm shows a significant increase in computational time efficiency.

Bernal et al. [4] proposed the sparse key point detector which is used to present object detection, in which interesting points are calculated from the oriented surface. The surface is constructed by a normal vector corresponding to the uniform surface, which is obtained from the range image. The probabilistic dense function (pdf) analysis applied to the normal vector contained in the directional surface allows us to select high highlights to obtain the profile of the scene. In addition, the probabilistic quality function (pmf) analysis applied to the information contained around these points is designed to select the region of interest in which the object to be detected is found. Finally, the Chess distance is applied to the points of interest contained in the region of interest to detect objects [5]. The proposed test involves qualitative and quantitative analysis using Middlebury and DSPLab datasets.

Motion detection is a challenging problem in video sequels with complex background. Different background subtraction techniques are introduced to detect moving objects in video sequels, but these techniques are not sufficient to solve the complex properties of dynamic scenes in actual monitoring tasks. Therefore, Vijayan et al. [6] proposed a simple and effective vector based approach to solve the practical monitoring challenges. In their research, the concept of linear dependence of vectors is used to construct background models corresponding to each pixel. Continuously, linear independence is used to detect moving objects from an input video sequel. The proposed method is based on space–time background subtraction because it uses spatial and temporal information to construct background models. An adaptive updating rate updating background model with regional diffusion is used to generate effective results in dynamic background. The proposed algorithm is used to test the datum dataset obtained from complex scenes, and the results are compared with the recent background subtraction method. The proposed method achieves better performance in terms of qualitative and quantitative results [7].

The major challenge in object detection is to accurately identify the position of objects in the image space, but an algorithm with a set of parameters is usually not enough, and the fusion of multiple algorithms and/or parameters can lead to more robust results. Wei et al. [8] proposed a computational intelligence fusion method based on dynamic analysis of object detection output protocol. In addition, they proposed an online and only training image enhancement strategy. Experiments with and without fusion results are presented. They proved that enhancement and fusion combinations are the best results for higher accuracy and reduced outliers. This method is proved in the background of conical, pedestrian and box detection of advanced driving aid system (ADAS) [9].

Foreign scholars in the field of computer vision research on video text extraction technology started earlier. In 1995, R. Lienhart and F. Stuber of Mannheim University proposed a method of automatic segmentation and recognition of subtitle text in digital video [10].

The rise of the smartphone industry has increased the demand for mobile capture of document images and embedded application. In the field of image processing, it is the first step of many algorithms to detect the key points and calculate the related features. However, the key point detector is mainly designed for real-world images and is usually not good for document images. Royer et al. [11] put forward an idea, which includes using document layout information to guide key point extraction. They compared this approach to the CORE algorithm for obfuscation reduction on three classical descriptors. The results showed that the matching quality and processing time were improved.

Text extraction is an important problem in image processing from optical character recognition to autopilot. Most traditional text segmentation algorithms consider separating text from a simple background (usually having a different color from text). In this work, S It's consider separating text from textured background with similar color to text. They look at the problem from the perspective of signal decomposition and consider a more realistic scenario in which the signal components are superimposed on each other (rather than added together). When the signal overlaps, we need to find a binary mask to display each component in order to separate the signal components. Because it is difficult to solve binary mask directly, we extend this problem to approximate continuous problem and solve it by alternating optimization method. We prove that the proposed algorithm achieves better results than other recent works on several challenging images [12–14].

Image-based text extraction is a popular and challenging research field in computer vision recently. Due to background clutter, unstructured scenarios, directions, fuzziness, and so on, urgent aspects such as text recognition and extraction of natural scenes have been studied. For text recognition, contrast enhancement is accomplished by applying LUV channel to input image to obtain a perfect and stable region. Then, the standard segmentation technique MSER is used to select L channel for region segmentation. In addition, the classification of connected components is performed to obtain the segmented image

through the fusion of two feature descriptors, LBP and T-HOG. Firstly, two feature descriptors are classified by linear SVM (s) mathContainer Loading Mathjax. Secondly, the results are combined into text/nontext parts by using weighted and fusion techniques. In text recognition, text regions are identified and marked with novel CNN networks. CNN output is stored in a text file to generate text words. Finally, if necessary, search a text file through a dictionary to obtain an appropriately optimized scenario text word containing hamming distance (error correction) techniques [15].

Convolutional neural network (CNN) is a multilayer artificial neural network, and its network structure is more special and complex than the traditional neural network. Convolutional neural network originated from the concept of receptive field (receptive field) and neural cognitive machine (neocognitron), which is the first application of receptive field in the field of artificial neural network. The neural cognitive machine decomposes a complete visual feature into many subfeatures and then processes the subfeatures on the hierarchical and hierarchical feature plane. This method attempts to model the visual system so that the recognition of features cannot be disturbed by object displacement or slight deformation [16, 17].

At present, the convolutional neural network has become a hot topic in the field of speech analysis and image classification and recognition, and many achievements have been obtained.

Periodic inspection of plant components is important to ensure safe operation. However, current practice is time-consuming, tedious and subjective, involving human technicians examining videos and identifying reactor cracks. Some vision-based crack detection methods have been developed for metal surfaces and they usually perform poorly when used to analyze nuclear inspection videos. Detecting these cracks is a challenging task because of their small size and the presence of noise patterns on the surface of the components. A deep learning framework is proposed based on convolutional neural network (CNN) and Bayesian data fusion scheme called NB-CNN, analysis of a single video frame for crack detection. At the same time, a new data fusion scheme is proposed to aggregate the information extracted from each video frame to enhance the overall performance and robustness of the system [18]. For this reason, a CNN is proposed to detect the fissures in each video frame, and the proposed data fusion scheme maintains the temporal and spatial coherence of the cracks in the video and effectively discards the false positives.

The prediction of visual attention data from any type of media is valuable to content creators and is used to drive coding algorithms effectively. With the current trend in the field of virtual reality (VR), the adaptation of known

technologies to this new media begins to gain momentum. An architecture extension is proposed for any convolutional neural network (CNN) to fine-tune traditional 2D significant prediction to omnidirectional image (ODI) in an end-to-end manner. The results show that each step in the proposed pipeline is aimed at making the generated salient map more accurate than the ground real data [19].

Inspired by the latest advances in deep learning. A new iterative confidence propagation convolutional neural network (BP-CNN) architecture was proposed for channel decoding under correlated noise. The architecture connects a trained CNN to a standard BP decoder. Standard BP decoders are used to estimate coding bits, followed by CNN to eliminate estimation errors of BP decoders and obtain more accurate channel noise estimates. The iteration between BP and CNN will gradually improve the decoding SNR and thus lead to better decoding performance.

## 2 Proposed method

### 2.1 Image segmentation processing

When the benchmark algorithm is initialized, the ground and sky regions are extracted from the scene by the plane geometry annotation algorithm. Then, an image block dictionary is generated by an unsupervised segmentation method based on different parameter settings for the area perpendicular to the ground in the scene [20]:

$$\{B_1, B_2, B_3, \ldots, B_i, \ldots, B_{K-1}, B_K\} \tag{1}$$

The benchmark algorithm not only estimates the visual angle type of each image block, but also trains a superpixel classifier to learn the mapping relationship between image features and object density. In order to realize the rough estimation of image block density. The estimation methods of visual angle class and density class of image block are very similar: (1) The general methods of these two kinds of estimators are given:

$$P(y_j|I) \propto \sum_{i \in B_j} P(B_j|I) P(\overline{y_j}|B_j, I) \tag{2}$$

In Eq. (2), the confidence P (y|I) of the label $y_i$(corresponding to the view class or the density class) of the superpixel i in the image $I$ can be marked by the probability of the image block $Bj$ corresponding to the superpixel i The weighted average is obtained, and the weight $P(B_j|I)$ expresses the homogeneity of the block $P(B_j|I)$. Both $P(\tilde{y}|B_j)$ and $P(B_j|I)$ are estimated by training LogoRegression (LR) classifier based on AdaBoost. The classifier uses the decision tree as a weak classifier and achieves the corresponding probability output through effective selection of features [21].

The benchmark algorithm provides the global constraints needed for scene combination and semantic reasoning through the representation of geometric and physical attributes such as static balance support force volume constraint and depth sorting. In this method, the image block is gradually added to the scene by recursion, and the image block is analyzed and processed accordingly, so as to realize the consistent understanding of the scene step-by-step. The advantage is that the image segmentation with high confidence can guide the understanding of other image segmentation after adding the scene, and the understanding of some scenes satisfies all of the above geometric and physical constraints in any recursive step. Which blocks in the image block dictionary are selected and added to the scene and in which order they are added are unknown. For this reason, the benchmark algorithm adopts a brand-new search strategy for the "building block world" $W_t$ composed of image blocks, that is, selecting a subset of $k = 4$ image blocks according to the following five local criteria in each recursion. And added to the scene: (a) surface geometric layout confidence; (b) density estimation; Physical stability of (c); (d) external support relationship; (e) relative depth relationship. Each image block subset is added to obtain the corresponding score through the cost function 5 in the minimization formula. And finally, it is determined that the image block subset with the least cost is added to $W_t$ and constitutes the updated $W_{t+1}$. The process of generating new subsets and estimating the cost is repeated in the $W_{t+1}$ formation until all image partitions are traversed.

For candidate block $B_j$, the benchmark algorithm determines whether it is added to $W$ by minimizing it:

$$E(B_i) = F_{\text{Geometry}}(Vi)$$
$$+ \sum_{R \in \text{ground,sky}} F_{\text{Contacts}}(Vi, R) + F_{\text{Intra\_Stability}}(B_i, V_i, d)$$
$$+ \sum_{j \in \text{block}} F_{\text{Inter\_Stability}}(V_i, S_{ij}, B_j) + F_{\text{Depth}}(V_i, S_{ij}, D)$$

$$(3)$$

1.  In order to overcome the difficulty of solving the problem caused by the too large assumption space of the scene partitioning method, the benchmark algorithm realizes the approximate minimization of the above cost function by the following steps on the basis of recursion: The estimated geometric properties, including the two-dimensional position $f_i$ of its perspective type h and horizon, are inferred by minimizing the energy terms $F_{\text{Geometry}}$ and $F_{\text{Contacts}}$, which can be expressed by Eq. (4):

$$F_{\text{Geometry}} + F_{\text{Contracts}} = \log P(Vi, fi | g, B_i, C_i^G, C_i^S)$$
$$\propto \log P(g | V_i, f_i, B_i) + \log P(C_i^G | V_i, f_i)$$
$$+ \log P(C_i^S | V_i, f_i)$$

$$(4)$$

In Eq. (4), the first term $P(g | V_i, f_i, B_i)$ representing $F_{\text{Geometry}}$ defines the similarity between the plane of the planar geometric notations corresponding to all in $B_i$ and the view type $V_i$ of the predicted $B_i$. $F_{\text{Contacts}}$ represents the geometric attribute consistency between the ground area and the intersection $C_i^G$ perpendicular to the ground area and the intersection point $C_i^S$ of the sky area and the ground area, which is composed of the latter two parts of formula (4). Among them, $P(C_i^G | V_i, f_i)$ and $P(C_i^S | V_i, f_i)$, respectively, represent the similarity between $V_i$, $C_i^G$ and $C_i^S$. The straight lines $l_G$ and $l_S$ are, respectively, fitted by $C_i^G$ and $C_i^S$, the corresponding slope is consistent with the estimated plane geometric label. Sexual estimates are obtained [22].

2.  Using (1), the estimated geometric properties, the stability of $B_i$ is calculated by minimizing the energy terms $F_{\text{Intra\_Stability}}$ and $F_{\text{Inter\_Stability}}$. Among them, $F_{\text{Intra\_Stability}}$ and $F_{\text{Inter\_Stability}}$ measure the intrinsic physical stability of $B_i$ and the relative physical stability between other blocks. When calculating $F_{\text{Intra\_Stability}}$, based on the estimation of the density of $B_i$, the internal stability is judged by rotating $B_i$ along its center of gravity by a small angle $\Delta\theta$ and calculating the change in its potential energy, as shown in Eq. (5):

$$\Delta P_i = \sum_{c \in \{\text{light,medium,heavy}\}} \sum_{s \in B_i} p(d_s = c) \cdot m_c \cdot \Delta h_s \qquad (5)$$

In Eq. (5), $p(d_s = c)$ is the probability that superpixel s in $B_i$ belongs to density class c, $\Delta h_s$ is the height change of s after rotation and $m_c$ is a predefined density class constant, which is obtained by training the density of the concentrated object and different materials. Frequency statistics are obtained. Such a smaller potential energy $B_i$ having a lower top density and a higher bottom density will be selected with a higher priority, while a $B_i$ with a higher potential energy will be discarded.

When calculating $F_{\text{Inter\_Stability}}$, based on the density estimate d and the inter-block support relationship $S_{ij}$, respectively, calculate the moment of each superpixel in $B_i$ and the support force around the intersection line with the block supporting it, and then pass a potential energy analysis method similar to $F_{\text{Intra\_Stability}}$ is calculated to determine the stability of the block [23].

The $S_{ij}$ and the sum of $V_i$ are derived by (1) and (2) to minimize the cost function. The relative depth relationship

of each image segment in three-dimensional space is defined. Different depth image blocks correspond to different energy values. The depth sorting D between blocks needs to satisfy the consistency of the projection region in the image, and the global depth sorting can be obtained by calculating the depth constraints between different block pairs.

Because of the original fixed segmentation of the algorithm, the superpixel converges to synthesize the wrong or incomplete blocks. For example, because of the occlusion of an object, two regions belonging to the same object are separated and cannot be recognized as the same object. To this end, based on the previously estimated depth order D and view type $V_i$, if two or more nonadjacent image blocks have co-contiguous image blocks and are estimated to be in front of them at the current view angle, then the merged strategy is used to merge the above images into blocks and the image block is judged as an occlusion. If the block attributes obtained by reasoning are quite different from those obtained by minimization, then they are divided into two or more image blocks with similar attributes [24].

## 2.2 Image text semantic extraction

The text in the video has a strong contrast with the background, and there is a very high frequency area at the intersection of the text and the background. Domain, text extraction can be done by analyzing the edge image.

The color of two adjacent pixels may be completely different in the case of similar gray scale. In this case, the gray edge operator cannot detect these edges. In this paper, the color edge images in RGB space are extracted for analysis. The Sobel edge operator is used to detect the edge of the RG and B color channels. The operator consists of gradient templates in four directions. The four templates are convolution with the original image, and the output results are described as follows: $(S_H, S_V, S_{LD}, S_{RD})_h$, where $h = 1, 2, 3$ represents three color channels. In this way, the subedge images in four directions are obtained, namely: $(S_H, S_V, S_{LD}, S_{RD})_h = \text{MAX}(|S_H|, |S_V|, |S_{LD}|, |S_{RD}|)_h$

Define the edge image as:

$$\text{Edgemap}(x,y) = \begin{cases} S(x,y), & \text{if } S(x,y) > T_e \quad \text{if } S(x,y) > T_e \\ 0 & \text{else} \end{cases}$$

(6)

where $T_e$ is the threshold; in this article, $T_e = 25$. The location of the text is realized by analyzing the edge points, so the binary result of the edge image is related to whether the correct text boundary can be detected. If the threshold value is too large, some text edges may be missed, which will lead to the missing detection. Too small threshold will introduce more nontext edges, resulting in false detection.

Therefore, it is better to choose the threshold dynamically according to the specific situation. For this reason, an adaptive threshold selection algorithm based on the background complexity threshold is proposed; that is, when the background is simple, a lower threshold is set (Fig. 1).

When the background is complex, a higher threshold should be adopted, which ensures that the nontext pixels can be suppressed and the text pixels can be kept as much as possible. The edge map is scanned using the window shown in Figs. 2 and 4, with a horizontal and vertical step value of N. Based on the horizontal projection analysis of large window coverage region, the background complexity of small window coverage area is calculated.

Assume that $PH_m(k)$ is the number of edge points of the kth line at the mth move, $k = 1…3N$, $m = 1…M$. If $PH_m(k) < \delta$, the kth line is considered to be a blank line, and the maximum number of consecutive blank lines from $PH_m(1)$ to $PH_m(N)$ is counted as $MAXBL_m$. Define the background complexity of the BComplexity $(m)$ to represent the coverage area of a small window on the mth move:

$$BComplexity(m) = \begin{cases} 0 & \text{if} \quad \sum_{k=1}^{N} PH_m(k) < \varepsilon \\ 1 & \text{if} \quad MAXBL_m \geq T_1 \\ 2 & \text{if} \quad MAXBL_m < T_1 \end{cases}$$

(7)

Set the threshold based on the complexity of the window background

$$T = \begin{cases} 255 & \text{if} \quad BComplexity = 0 \\ T_{Low}^*, & \text{if} \quad BComplexity = 1 \\ T_{High}^*, & \text{if} \quad BComplexity = 2 \end{cases}$$

(8)

where $T_{Low}^*$ and $T_{High}^*$ are obtained by using the Otsu algorithm three times in a large window. First, the Otsu iteration is performed on the gray scale [0–255] to obtain $T_{Middle}^*$ and then iteratively obtain $[0, T_{Middle}^*]$ and $[T_{Middle}^*, 255]$ in $T_{Low}^*$ and $T_{High}^*$, respectively, where $\sigma_B^2$ represents the variance between classes.

## 2.3 Model of convolutional neural network

As a multilayer network model, convolutional neural network is also a kind of artificial neural network. Because convolutional neural network has made a great breakthrough in speech and image recognition and has been successfully trained, it has become a major research hot spot in the current society [25]. The network structure of the convolutional neural network is as follows (Fig. 2):

In general, layer C is a convolution layer, and the input of each neuron is connected to the local receptive field in the previous layer, and the local features are extracted; the s layer is a downsampling layer, and each feature computing layer of the network consists of multiple feature mapping layers. Each feature is mapped into a plane on

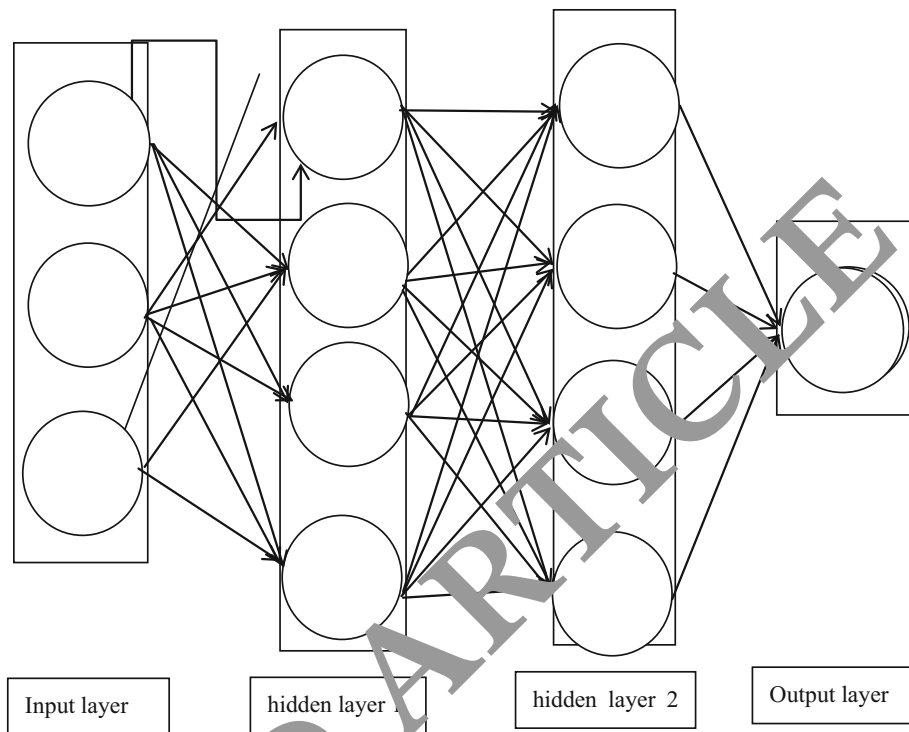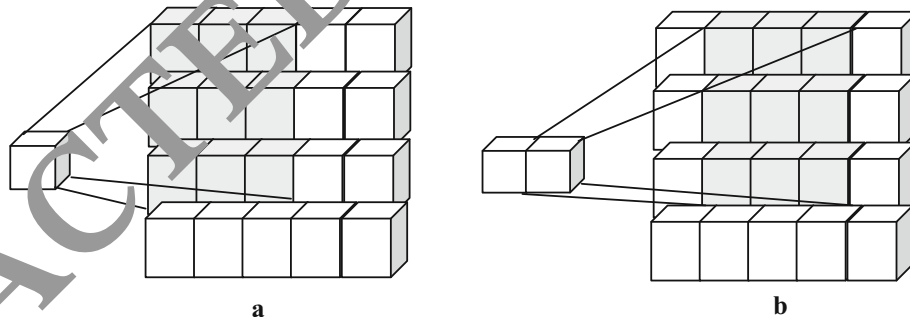**Fig. 1** Image processing process of convolutional neural network



Input layer   hidden layer 1   hidden layer 2   Output layer

**Fig. 2** Structure diagram of convolutional neural network(CNN)



a                    b

which all neurons share weights. In the convolution layer, the former layer of neurons and the convolutional nucleus are convoluted, and then the result is calculated by the activation function, and each output feature graph and the previous layer characteristic graph are calculated and outputted [6]:

$$a_j^l = f\left(\sum_{i \in M_j^l} a_i^{l-1} * k_{ij}^l + b_j^l\right) \tag{9}$$

where $F$ is the activation function, $l$ is the number of network layers, $M_j^l$ is the input feature map set, $a_j^l$ is the $j$th feature map of the $l$th convolutional layer, $k$ is the convolution kernel and $b$ is the offset. $F$ is the activation function, $l$ is the number of network layers, $M$ is the input feature map set, $a$ is the $j$th feature map of the $L$th

convolutional layer, $K$ is the convolution kernel and b is the offset.

At the downsampling level, the output feature map becomes smaller, but does not change the number after the input feature is operated on:

$$a_j^l = f\left(\beta_j^l \text{down}(x_j^{l-1}) + b_j^l\right) \tag{10}$$

where down($\cdot$) is the downsampling operation, which is to reduce the dimension of a small area of the data. The training of convolutional neural network is divided into two stages: forward propagation and backward propagation. In the forward propagation phase, a sample $(x, y_p)$ is taken from the sample set, X is used as the input of the convolutional neural network and the input is transformed layer by layer to the output layer:

$$o_p = F_n(\ldots(F_2(F_1(xW_1)W_2)\ldots)W_n) \tag{11}$$

In the backward propagation phase, i.e., the error propagation phase, the actual output $o_p$ and $y_p$ label errors are

$$E_p = \frac{1}{2} \sum_j (y_{pj} - o_{pj})^2 \tag{12}$$

The weights of convolutional neural networks are adjusted according to the method of minimum error. The weight sharing makes the convolutional neural network more similar to the biological network and further reduces the complexity of the network model. The main purpose of convolutional neural network is to identify a multilayer perceptron for two-dimensional images. The network model has good translation rotation invariance scaling invariance and other deformation invariance. Then, the convolution layer, activation function, pool layer and full connection layer in the convolutional neural network are introduced in detail [27].

## 3 Experiments

### 3.1 Street scene recognition

1. Data collection: the Cityscapes dataset focuses on the urban road environment in real-world situations, including 30 categories of roads, vehicles, pedestrians, traffic signs, etc. The network training uses 2975 images with fine labeling in the dataset.
2. Data processing: (1) because the semantic category of the dataset is too many, and some of the data are too small, such as bridges, tunnels, etc., this paper integrates some categories and divides them into eight categories, that is, roads, pedestrians, motor vehicles, buildings, skies, vegetation, poles and background. (2) Because the original image is too large (2048–1024) and the hardware condition of the experiment is limited, it is necessary to cut the image. The method of clipping is to intercept the image of 960G960 randomly in the original image, then reduce the image to 384kn384, and then input it into the network for training.

### 3.2 Text semantic extraction

Text detection and location: it mainly includes three steps: text detection, tracking and location, and the processing process is carried out in time domain and space domain respectively [28]. The processing in the time domain is mainly to detect the frame of each text in the video data

stream, and the processing in the spatial domain is to locate the position of the text region appearing in the video frame.

Text segmentation: On the basis of text detection and localization, the text region is transformed into a binary image. This step mainly includes the enhancement of the text region and the binarization of the text region. The former is to filter the noncharacter field of the text area and enhance the text area to reduce the interference of the background image. The latter is to transform the text area into binary image which can be recognized by the traditional OCR software according to the color shape and texture of the text in the region and to deal with the conglutination and fracture of the characters at the same time.

Text recognition: OCR software is used to complete the final conversion from digital image to character encoding. At present, the OCR technology used in character recognition stage is quite mature, and the commercial OCR software can achieve good recognition results. This part of the content is not included in the research scope of this paper.

### 3.3 Convolutional neural network

The core features of the convolutional neural network are the combination of local receptive field, weight sharing mechanism, time and space subsampling and the purpose of reducing dimension during feature extraction. The weight sharing mechanism not only reduces the complexity of the model, but also effectively controls the number of weights and implements the constraint of the number of parameters by using spatial relations to improve the performance of the forward BP algorithm.

As a multilayer neural network, each layer of convolutional neural network consists of several two-dimensional planes, and each plane contains multiple independent neurons. It mainly includes data input layer, input, feature extraction layer, feature mapping layer, output layer, in which the image is first input through input layer. Then, we do feature extraction in C1 layer, and after filter processing, we do weighted plus bias operation and then convolution to form feature map in C1 layer. Then, we continue to bias the weighted value after summing the pixel values in each group of feature maps. The activation function selected in the whole process is a sigmoid function with a value of $[-1]$ to map to the corresponding feature map in S2 layer, and then the obtained feature map is filtered to C3 layer and the above steps are continued. Until the pixels of these images are rasterized, they are input into the artificial neural network model as vectors and finally output. In the above convolutional neural network diagram, the feature extraction layer is similar to the C layer, and its main function is to realize the local feature extraction of the input data. The successful extraction of this local feature

means that there is an exact location relationship among the features. When making feature mapping, each computing layer is made up of multiple feature maps, similar to S layer, and satisfies the unique corresponding determinacy between feature and mapping plane, and the corresponding neuron weights on the same mapping plane are shared. Because the sigmoid function has small kernel and the displacement invariance of its feature mapping, it is introduced as the activation function of the model in the feature mapping structure of S layer. It is worth noting that in feature extraction every neuron in the same plane shares the same weight. That is, the weight sharing mechanism of convolutional neural network implements the constraint of network parameters and reduces the complexity of the model training process. In addition, the feature extraction layer and the feature computing layer of CNN alternately carry out feature extraction; that is, each C layer is closely connected with a S layer, so the feature extraction function is better realized. So the special structure such as convolutional neural network has a high distortion tolerance to the input layer of the sample when dealing with the problem in the field of recognition [29].

In the era dominated by artificial intelligence, the study of natural language processing (NPL) has attracted great attention from all walks of life at home and abroad and has become one of the most important research direction.

## 4 Discussion

### 4.1 Recognition rate of different objects

From the test results, it can be seen that the network has better segmentation and recognition ability for vehicles, roads, buildings, vegetation and sky, but the segmentation performance is poor for the category of poles with less target pixels (Table 1).

**Table 1** Recognition rates of different objects

| Class | Accuracy rate (%) |
| --- | --- |
| Road | 94.57 |
| Person | 87.65 |
| Vehicle | 83.59 |
| Building | 87.17 |
| Pole | 69.43 |
| Vegetation | 75.19 |
| Sky | 63.48 |
| Background | 79.81 |

On the whole, there is a high recognition accuracy for street scene. The classification, position and scale of objects in the scene can be analyzed by the scene recognition image, and the information of vehicles in front of the vehicle and vegetation on both sides of the road can be obtained, and the road can be calculated. The proportion of pixels taken by objects such as the sky and figures inferred that the road ahead was wider and that there were pedestrians in front of the left.

Text semantic information extraction. The semantic information extracted from the text mainly includes two parts. One is to extract named entities from the text as keywords to annotate the video content. The second is to analyze the emotional tendency of the Chinese text of the video to obtain the emotional semantics of the video, including the semantic tendency recognition of words, the text classification based on emotion, the point of view extraction, the subjective analysis, and so on. This part can be regarded as an independent research direction, which is an important research content.

### 4.2 Recognition effect in different states

The results of the experiment are shown in Fig. 3. The results of the segmentation of the dataset in the daytime and night environments are shown. The three images on the left are the environment of night, and the three pictures on the right are the environment of the day. The classical CNN model is compared with the Kinect method (Fig. 4).

The result of Kinect segmentation in the graph has noise, and the segmentation result does not carry out antagonistic learning. CNN adds the segmentation results of confrontation learning, and the whole structure information of the segmentation results is well preserved, compared with the network results with no confrontation in the above two lines. In the segmentation of small categories, especially the blue character segmentation, CNN model has a good effect.

## 5 Conclusions

Compared with the traditional image semantic segmentation method, the method based on the deep convolutional neural network is simple and the segmentation effect is better than the traditional image semantic segmentation method.

Model fusion helps to achieve high accuracy of small objects while still achieving high global accuracy. Image semantic segmentation is a key technology in the field of image processing and computer vision. It is an important part of computer cognitive image content. The quality of semantic segmentation plays a crucial role in subsequent
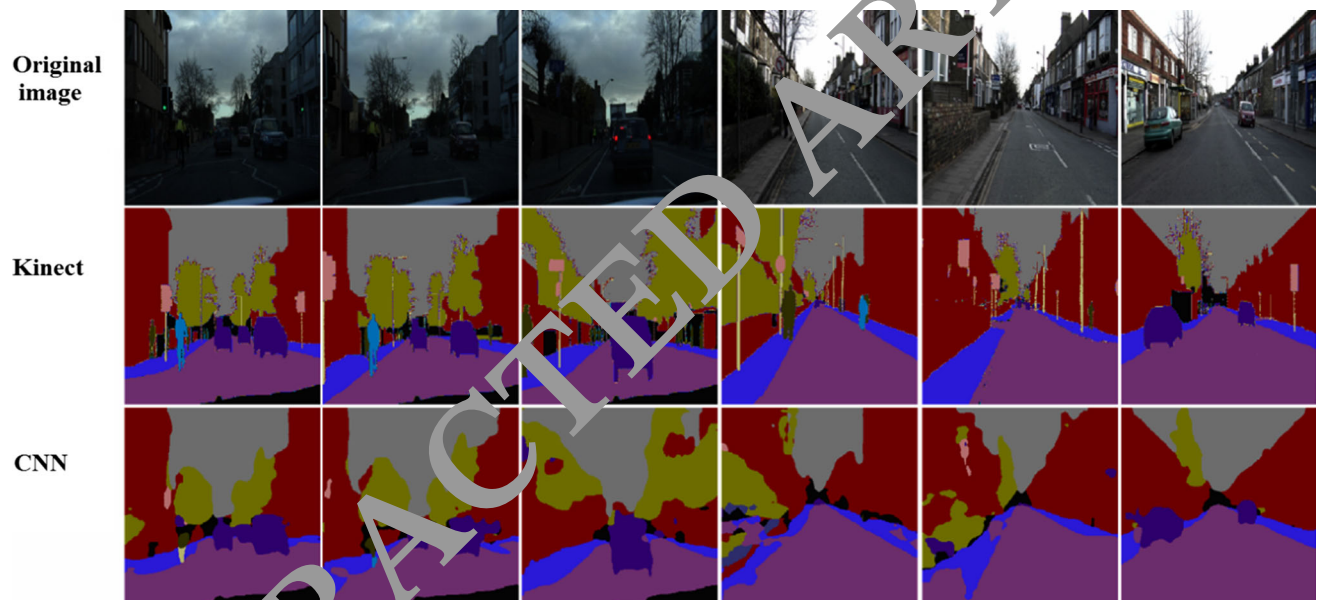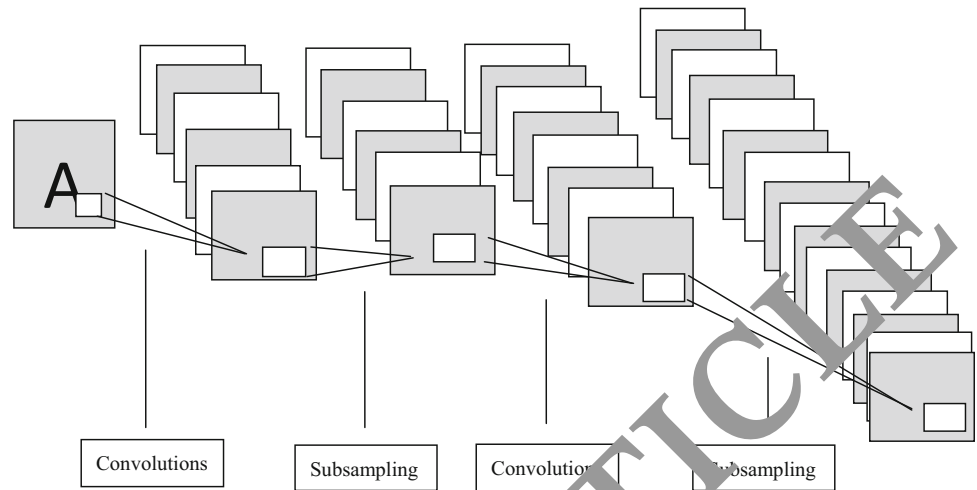
**Fig. 3** Text detection



**Fig. 4** Kinect and CNN recognition effects in different states

tasks such as image understanding, scene analysis and target tracking. Therefore, it is of great practical significance to study an effective image semantic segmentation algorithm. With the continuous development of deep learning, the high accuracy brought by neural networks has been widely studied and applied in many scenes such as image recognition and semantic segmentation. Compared with the traditional semantic segmentation method based on region feature extraction, the image features acquired by the deep convolutional neural network method have stronger representation ability, so the algorithm has better effect. The basic idea of semantic segmentation based on deep convolutional neural network is to extract the semantic features of each pixel in the image by using

neural network, then classify and identify the pixels according to these features, so as to obtain the segmentation image containing semantic information. Therefore, the core of this method is how to improve the recognition accuracy of pixels on the network.

## Compliance with ethical standards

**Conflict of interest** There are no conflicts of interests in this work.

# References

1. Kumar M, Mao YH, Wang YH, Qiu TR, Yang C, Zhang WP (2017) Fuzzy theoretic approach to signals and systems: static systems. Inf Sci 418-419:668–702

2. Zhang WP, Yang JZ, Fang YL, Chen HY, Mao YH, Kumar M (2017) Analytical fuzzy approach to biological data analysis. Saudi J Biol Sci 24(3):563–573

3. Wang S, Wang M, Zhao X et al (2018) Two-stage object detection based on deep pruning for remote sensing image. In: International conference on knowledge science, engineering and management. Springer, Cham, pp 137–147

4. Suh HP, Kim Y, Suh Y et al (2018) Multidetector computed tomography (CT) analysis of 168 cases in diabetic patients with total superficial femoral artery occlusion: is it safe to use an anterolateral thigh flap without CT angiography in diabetic patients. J Reconstr Microsurg 34(01):065–070

5. Yuan C, Xia Z, Jiang L (2019) Fingerprint liveness detection using an improved CNN with image scale equalization. IEEE Access 7(99):26953–26966

6. Vijayan M, Ramasundaram M (2018) Moving object detection using vector image model. Optik 168:963–973

7. Wei P, Ball JE, Anderson DT (2018) Fusion of an ensemble of augmented image detectors for robust object detection. Sensors 18(3):894

8. Royer E, Bouchara F (2018) Guiding text image keypoints extraction through layout analysis. In: IAPR international conference on document analysis and recognition. IEEE, pp 9–14

9. Minaee S, Wang Y (2018) Text extraction from texture images using masked signal decomposition. In: IEEE global conference on signal and information processing. IEEE, pp 1210–1214

10. Vasilopoulos N, Wasfi Y, Kavallieratou E (2018) Automatic text extraction from arabic newspapers. In: International conference image analysis and recognition. Springer, Cham, pp 505–510

11. Bai D, Wang C, Bo Z (2018) CNN feature boosted SeqSLAM for real-time loop closure detection. Chin J Electron 27(3):488–499

12. Liang F, Shen C, Wu F (2018) An iterative BP-CNN architecture for channel decoding. IEEE J Sel Top Signal Process 12(1):144–159

13. Huan Du, Liu Zhi, Song Hangke, Li Lin, Zheng Xu (2016) Improving RGBD saliency detection using progressive region classification and saliency fusion. IEEE Access 4:8987–8994

14. Dai D, Sakaridis C, Hecker S (2019) Curriculum model adaptation with synthetic and real data for semantic foggy scene understanding. Int J Comput Vision 11:1–23

15. Zhang C, Pan X, Li H (2018) A hybrid MLP-CNN classifier for very fine resolution remotely sensed image classification. ISPRS J Photogram Rem Sens 140:133–144

16. Khatami A, Babaie M, Tizhoosh HR (2010) A sequential search-space shrinking using CNN transfer learning and a radon projection pool for medical image retrieval. Expert Syst Appl 100:224–233

17. Tama BA, Rhee K-H (2019) An in-depth experimental study of anomaly detection using gradient boosted machine. Neural Comput Appl 31(4):955–965

18. Yokota F, Otake Y, Takao M (2018) Automated muscle segmentation from CT images of the hip and thigh using a hierarchical multi-atlas method. Int J Comput Assist Radiol Surg 13(9):1–10

19. Keun CH, Seop KH, Ho JS (2012) Sonography of affected and unaffected shoulders in hemiplegic patients: analysis of the relationship between sonographic imaging data and clinical variables. Ann Rehabil Med 36(6):828–835

20. Dutta A, Zhong Y, Depraetere B (2014) Model-based and model-free learning strategies for wet clutch control. Mechatronics 24(8):1008–1020

21. Kim S, Ji Y, Lee KB (2018) An effective sign language learning with object detection based ROI segmentation. In: IEEE international conference on robotic computing. IEEE Computer Society, pp 330–333

22. Wang K, Liang L, Yan X (2018) Cost-effective object detection: active sample mining with switchable selection criteria. IEEE Trans Neural Netw Learn Syst 99:1–17

23. Khanzadeh T, Hasani MF, Talebi M (2018) Investigation of BAX and BCL2 expression and apoptosis in a resveratrol- and prednisolone-treated human T-ALL cell line, CCRF-CEM. Blood Res 53(1):53–60

24. Fei Y, Wang KCP, Zhang A (2019) Pixel-level cracking detection on 3D asphalt pavement images through deep-learning-based CrackNet-V. IEEE Trans Intell Transp Syst 99:1–12

25. Kim Y, Kang BN, Kim D (2018) Detector with focus: normalizing gradient in image pyramid. In: IEEE international conference on image processing. IEEE, pp 420–424

26. Liang Y, Zhang Y, Chen Z (2019) Re-nucleation and etching of graphene during the cooling stage of chemical vapor deposition. J Electron Mater 48(3):1740–1745

27. Glasner D, Galun M, Alpert S (2012) Viewpoint-aware object detection and continuous pose estimation. Image Vis Comput 30(12):923–933

28. Kang MS, Lim YC (2018) High performance and fast object detection in road environments. In: Seventh international conference on image processing theory, TOOLS and applications. IEEE, pp 1–6

29. Furuta R, Tsubaki I, Yamasaki T (2018) Fast volume seam carving with multipass dynamic programming. IEEE Trans Circuits Syst Video Technol 28(5):1087–1101

RETRACTED ARTICLE