

RESEARCH ARTICLE

CSC-Unet: A Novel Convolutional Sparse Coding Strategy Based Neural Network for Semantic Segmentation

HAITONG TANG¹, SHUANG HE¹, MENGDUO YANG², XIA LU³, QIN YU⁴, KAIYUE LIU¹,
HONGJIE YAN⁵, AND NIZHUAN WANG^{1,6,7}

¹School of Geomatics and Marine Information, Jiangsu Ocean University, Lianyungang 222005, China

²School of Information Technology, Suzhou Institute of Trade and Commerce, Suzhou 215009, China

³School of Geography Science and Geomatics Engineering, Suzhou University of Science and Technology, Suzhou 215000, China

⁴School of Computer Engineering, Jiangsu Ocean University, Lianyungang 222005, China

⁵Department of Neurology, Affiliated Lianyungang Hospital of Xuzhou Medical University, Lianyungang 222002, China

⁶School of Biomedical Engineering, ShanghaiTech University, Shanghai 201210, China

⁷Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University, Hong Kong, SAR, China

Corresponding authors: Hongjie Yan (yanhjns@gmail.com) and Nizhuan Wang (wangnizhuan1120@gmail.com)


This work was supported in part by the National Natural Science Foundation of China under Grant 41506106, Grant 61701318, and Grant 82001160; in part by the Natural Science Research Project of Jiangsu Higher Education Institutions under Grant 20KJB520014; in part by the Project of Huaguoshan Mountain Talent Plan—Doctors for Innovation and Entrepreneurship; in part by Jiangsu Province Graduate Research and Practice Innovation Program Project under Grant SY202144X; in part by The First People's Hospital of Lianyungang—Advanced Technology Support Project under XJ1811; and in part by The '14th Five-Year Plan' Youth Medical Talent Project for Science, Education, and Health Engineering in Huaguoshan, Lianyungang City.

ABSTRACT It is still a challenging task to perform the semantic segmentation with high accuracy due to the complexity of real picture scenes. Many semantic segmentation methods based on traditional deep learning insufficiently captured the semantic and appearance information of images, which put limit on their generality and robustness for various application scenes. Thus, in this paper, we proposed a novel strategy that reformulated the popularly used convolution operation to multi-layer convolutional sparse coding block in semantic segmentation method to ease the aforementioned deficiency. To prove the effectiveness of our idea, we chose the widely used U-Net model for the demonstration purpose, and we designed CSC-Unet model series based on U-Net. Through extensive analysis and experiments, we provided credible evidence showing that the multi-layer convolutional sparse coding block enables semantic segmentation model to converge faster, extract finer semantic and appearance information of images, and improve the ability to recover spatial detail information. The best CSC-Unet model significantly outperforms the results of the original U-Net on three public datasets with different scenarios, i.e., 87.14% vs. 84.71% on DeepCrack dataset, 68.91% vs. 67.09% on Nuclei dataset, and 53.68% vs. 48.82% on CamVid dataset, respectively. In addition, the proposed strategy could be possibly used to significantly improve segmentation performance of any semantic segmentation model that involves convolution operations and the corresponding code is available at <https://github.com/NZWANG/CSC-Unet>.

INDEX TERMS U-Net, semantic segmentation, deep learning, convolution operation, convolutional sparse coding (CSC).

I. INTRODUCTION

In reality, the increasing application scenarios require inferring relevant knowledge or semantics from images,

The associate editor coordinating the review of this manuscript and approving it for publication was Mostafa M. Fouda .

as a result, the importance of semantic segmentation for scene understanding is gradually increasing. Semantic segmentation gives us more detailed understanding of images than image classification [1], [2], [3], [4], [5] or object detection [6], [7], [8], [9], [10], [11], [12], [13], [14]. This understanding is crucial in many different domains such

as autonomous driving [15], [16], [17], robotics [18], [19], [20], image search engines [21], [22], [23], etc. Recently, many semantic segmentation methods have emerged. For example, fully convolutional networks (FCN) [24], at an end to end form, has firstly implemented the pixel-wise prediction task based on convolution operation, achieving relatively better results in natural scene image segmentation. SegNet [25] makes the model more efficient than FCN by introducing more skip architectures and max-pooling indexes. PSPNet [26] used dilated convolution and pyramid pooling to improve SegNet. U-Net [27] was proposed in the 2015 ISBI competition, which consists of contracting and symmetrically expanding sub-networks to form a U-shaped architecture. This model was originally designed to solve biomedical image segmentation. Since it requires a small number of training samples to achieve good segmentation results. There are also many variants of the U-Net. For example, Unet++ [28] concatenates U-Net models with different layers that share weights in encoding blocks and features from different layers are also fused through skip connections to achieve better segmentation results. Based on U-Net, FPN [29] and Resnet [5], the U2-Net [30] model is proposed, which achieves surprising results on saliency detection tasks with good real-time performance. Other excellent semantic segmentation models include DeepLab series. For example, Deeplabv1 [31] uses VGG16 [3] with atrous convolutions as the backbone, and adds conditional random fields (CRFs) in post-processing to further improve segmentation performance. Deeplabv2 [32] replaces the backbone in Deeplabv1 with ResNet and proposes atrous spatial pyramid pooling (ASPP) for multi-scale segmentation. Deeplabv3 [33] increases the depth of the backbone without CRFs. It also replaced the convolution with atrous rate of 24 in the ASPP with a 1×1 convolution and added average pooling and batch normalization layers [34]. Deeplabv3+ [35] designs Deeplabv3-based encode-decode models and modifies Xception [36] as the backbone.

All the above semantic segmentation models are based on convolution operations, which have strong feature representation capabilities to extract semantic (global) and appearance (local) information of images. In fact, for the segmentation task of complex image, it is usually limited by the semantic and appearance information extracted from the shallow convolution layers. As a rule, they mostly choose to deepen the network layers so that the semantic segmentation network can better capture the semantic and appearance information of the images to improve the segmentation performance. However, if the network keeps deepening indefinitely, there is a tremendous challenge for both the computational power and the optimizer. Therefore, the main motivation of this study is to address the problem of insufficient feature extraction of convolution operation at the root by optimizing instead of deepening them.

In this paper, we proposed a novel strategy in semantic segmentation model which reformulated convolution operation to multi-layer convolutional sparse coding (ML-CSC)

block. Taking the U-Net as an example, we demonstrated the effectiveness and robustness of ML-CSC block strategy in the designed CSC-Unet model series, and it can also be potentially applied to other excellent convolution-based semantic segmentation networks, such as SegNet, U2-Net, etc. Actually, in the Appendix part, we have also implemented the CSC-Unet++, CSC-Unet3+, and CSC-DeepLabv3+ models corresponding to Unet++, Unet3+ [37], and DeepLabv3+, respectively. Benefit from the advantages of ML-CSC block in information representation compared to convolutional operation, we hypothesize that the CSC-Unet model series has the superiorities of better captured semantic and appearance information of original images, better spatial detail information, and better convergence efficiency without increasing the trainable parameters.

As far as we are aware, it is the first work to explore semantic segmentation based on convolutional sparse coding. We hope that our strategy can provide new insights for designing semantic segmentation models. The main contributions of this paper are summarized as follows:

- 1) We proposed a novel strategy in semantic segmentation networks that used the multi-layer convolutional sparse coding blocks instead of the traditional convolution operations;
- 2) We extended the ML-CSC block to U-Net as CSC-Unet model series, and further demonstrated the advantages and feasibility through extensive experiments;
- 3) We explored the impact of the number of unfoldings in the ML-CSC block on the performance of the semantic segmentation model.

The rest of this paper is organized as follows. Section II will give a brief introduction of sparse coding, ML-CSC and block of ML-CSC. In Section III, we will present the design details of CSC-Unet model series. The procedure and results of the experiments are presented in Section IV. Section V is the conclusion and future work that we will carry out.

II. REVIEW OF MULTI-LAYER CONVOLUTIONAL SPARSE CODING

In this section, we firstly reviewed the sparse coding, multi-layer convolutional sparse coding and its solution algorithms, then we presented the details of the designed ML-CSC block.

A. SPARSE CODING

Sparse coding represents the signal with few non-zero coefficients as possible and has been used in a wide variety of applications [38], [39], [40], [41], [42], [43]. In which, the image \mathbf{y} is considered as a linear combination of a set of basis vectors \mathbf{d}_i , where most of the coefficients γ_i are zero.

$$\mathbf{y} = \mathbf{D}\boldsymbol{\Gamma} = [\mathbf{d}_1 \quad \mathbf{d}_2 \quad \cdots \quad \mathbf{d}_K] \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_K \end{bmatrix} \quad (1)$$

However, when processing image signal, sparse coding first decomposes the whole image into a set of overlapping

image blocks and then operates these blocks independently, which leads to that the sparse representation of the image is highly redundant and loss of detailed information during image recovery [44].

B. MULTI-LAYER CONVOLUTIONAL SPARSE CODING

Assumes that the input image \mathbf{y} satisfies ML-CSC model, it can be denoted as:

$$\begin{aligned} \mathbf{y} &= \mathbf{D}_1 \mathbf{\Gamma}_1, \\ \mathbf{\Gamma}_1 &= \mathbf{D}_2 \mathbf{\Gamma}_2, \\ &\vdots \\ \mathbf{\Gamma}_{L-1} &= \mathbf{D}_L \mathbf{\Gamma}_L. \end{aligned} \quad (2)$$

where $\{\mathbf{D}_i\}_{i=1}^L$ are special dictionaries, each \mathbf{D}_i is a transpose of a convolutional operator matrix \mathbf{W}_i .

$$\mathbf{D}_i = \mathbf{W}_i^T \quad (3)$$

Convolutional sparse coding (CSC) [45] applies convolutional filters to reconstruct the whole image. Since the image is processed as a whole, it bridges the above gap in sparse coding. ML-CSC is an extension of convolutional CSC which sets that the $\{\mathbf{\Gamma}_i\}_{i=1}^L$ also satisfy CSC model to form multi-layer representation method about original image.

C. THE SOLVER OF MULTI-LAYER CONVOLUTIONAL SPARSE CODING MODEL

The process of solving the $\{\mathbf{\Gamma}_i\}_{i=1}^L$ in ML-CSC model can be formulated as an optimization problem as Equation 4. Where $\hat{\mathbf{\Gamma}}_0 = \mathbf{y}$, and $\|\cdot\|_0$ is sparsity regularization term [46]. $\{\lambda_i\}_{i=1}^L$ are regularization parameters to balance the sparsity and accuracy of $\{\mathbf{\Gamma}_i\}_{i=1}^L$.

$$\begin{aligned} \hat{\mathbf{\Gamma}}_1 &= \arg \min_{\mathbf{\Gamma}_1} \frac{1}{2} \left\| \hat{\mathbf{\Gamma}}_0 - \mathbf{D}_1 \mathbf{\Gamma}_1 \right\|_2^2 + \lambda_1 \|\mathbf{\Gamma}_1\|_0, \\ \hat{\mathbf{\Gamma}}_2 &= \arg \min_{\mathbf{\Gamma}_2} \frac{1}{2} \left\| \hat{\mathbf{\Gamma}}_1 - \mathbf{D}_2 \mathbf{\Gamma}_2 \right\|_2^2 + \lambda_2 \|\mathbf{\Gamma}_2\|_0, \\ &\vdots \\ \hat{\mathbf{\Gamma}}_L &= \arg \min_{\mathbf{\Gamma}_L} \frac{1}{2} \left\| \hat{\mathbf{\Gamma}}_{L-1} - \mathbf{D}_L \mathbf{\Gamma}_L \right\|_2^2 + \lambda_L \|\mathbf{\Gamma}_L\|_0. \end{aligned} \quad (4)$$

Finding $\{\mathbf{\Gamma}_i\}_{i=1}^L$ at once is NP-hard and challenging in computation and concept because of the inclusion of the 0-norm [47]. Candes et al. [48] have shown that the 0-norm can be deflated to 1-norm, turning Equation 4 into a convex optimization problem, as shown in Equation 5.

$$\begin{aligned} \hat{\mathbf{\Gamma}}_1 &= \arg \min_{\mathbf{\Gamma}_1} \frac{1}{2} \left\| \hat{\mathbf{\Gamma}}_0 - \mathbf{D}_1 \mathbf{\Gamma}_1 \right\|_2^2 + \lambda_1 \|\mathbf{\Gamma}_1\|_1, \\ \hat{\mathbf{\Gamma}}_2 &= \arg \min_{\mathbf{\Gamma}_2} \frac{1}{2} \left\| \hat{\mathbf{\Gamma}}_1 - \mathbf{D}_2 \mathbf{\Gamma}_2 \right\|_2^2 + \lambda_2 \|\mathbf{\Gamma}_2\|_1, \\ &\vdots \\ \hat{\mathbf{\Gamma}}_L &= \arg \min_{\mathbf{\Gamma}_L} \frac{1}{2} \left\| \hat{\mathbf{\Gamma}}_{L-1} - \mathbf{D}_L \mathbf{\Gamma}_L \right\|_2^2 + \lambda_L \|\mathbf{\Gamma}_L\|_1. \end{aligned} \quad (5)$$

1) LAYERED THRESHOLDING ALGORITHM

Papayan et al. [49] proposed the layered thresholding algorithm, that use thresholding algorithm [46] to solve the sparse vectors $\{\mathbf{\Gamma}_i\}_{i=1}^L$ step by step in different layers. The solver can be written as follow:

$$\begin{aligned} \hat{\mathbf{\Gamma}}_1 &= h_{\theta_1} \left(\mathbf{D}_1^T \mathbf{y} \right), \\ \hat{\mathbf{\Gamma}}_2 &= h_{\theta_2} \left(\mathbf{D}_2^T \hat{\mathbf{\Gamma}}_1 \right), \\ &\vdots \\ \hat{\mathbf{\Gamma}}_L &= h_{\theta_L} \left(\mathbf{D}_L^T \hat{\mathbf{\Gamma}}_{L-1} \right). \end{aligned} \quad (6)$$

The soft non-negative threshold operator $\{h_{\theta_i}\}_{i=1}^L$ can be viewed as a translation of the activation function rectified linear unit (ReLU) activation function [50] by $\{\theta_i\}_{i=1}^L$ units [49]. Combined with Equation 3, the above equation can be written as follow:

$$\begin{aligned} \hat{\mathbf{\Gamma}}_1 &= \text{ReLU} \left(\mathbf{W}_1 \mathbf{y} + \theta_1 \right), \\ \hat{\mathbf{\Gamma}}_2 &= \text{ReLU} \left(\mathbf{W}_2 \hat{\mathbf{\Gamma}}_1 + \theta_2 \right), \\ &\vdots \\ \hat{\mathbf{\Gamma}}_L &= \text{ReLU} \left(\mathbf{W}_L \hat{\mathbf{\Gamma}}_{L-1} + \theta_L \right). \end{aligned} \quad (7)$$

where and $\{\theta_i\}_{i=1}^L$ are trainable parameters.

2) MULTI-LAYER ITERATIVE SOFT THRESHOLDING ALGORITHM

An attractive approximate solver of Equation 5 is the multi-layer iterative soft thresholding algorithm (ML-ISTA) [51], which uses an iterative soft thresholding algorithm [52] at each layer.

$$\begin{aligned} \hat{\mathbf{\Gamma}}_1 &= \hat{\mathbf{\Gamma}}_1^{k+1} = \text{ReLU} \left(\tilde{\mathbf{\Gamma}}_1 - \mu_1 \mathbf{D}_1^T \left(\mathbf{D}_1 \tilde{\mathbf{\Gamma}}_1 - \hat{\mathbf{\Gamma}}_0 \right) + \theta_1 \right), \\ \hat{\mathbf{\Gamma}}_2 &= \hat{\mathbf{\Gamma}}_2^{k+1} = \text{ReLU} \left(\tilde{\mathbf{\Gamma}}_2 - \mu_2 \mathbf{D}_2^T \left(\mathbf{D}_2 \tilde{\mathbf{\Gamma}}_2 - \hat{\mathbf{\Gamma}}_1 \right) + \theta_2 \right), \\ &\vdots \\ \hat{\mathbf{\Gamma}}_L &= \hat{\mathbf{\Gamma}}_L^{k+1} = \text{ReLU} \left(\tilde{\mathbf{\Gamma}}_L - \mu_L \mathbf{D}_L^T \left(\mathbf{D}_L \tilde{\mathbf{\Gamma}}_L - \hat{\mathbf{\Gamma}}_{L-1} \right) + \theta_L \right). \end{aligned} \quad (8)$$

It is well known that the layered thresholding algorithm is the simplest and dumbest pursuit algorithm for ML-CSC [49], [53]. This is because it only uses the thresholding algorithm independently at each layer and does not take into account the connection between layers. In contrast, ML-ISTA solves the above problem by constructing $\tilde{\mathbf{\Gamma}}_i = \mathbf{D}_i \mathbf{D}_{i+1} \cdots \mathbf{D}_L \hat{\mathbf{\Gamma}}_L^k$ in each layer which takes full account of the connections between the different layers.

D. MULTI-LAYER CONVOLUTIONAL SPARSE CODING BLOCK

In the ML-CSC model, we summarize the solving process as shown in Figure 1, and name it as ML-CSC block. $\{\mathbf{W}_i\}_{i=1}^L$ denotes the convolution operation, and $\{\mathbf{W}_i^T\}_{i=1}^L$ denotes

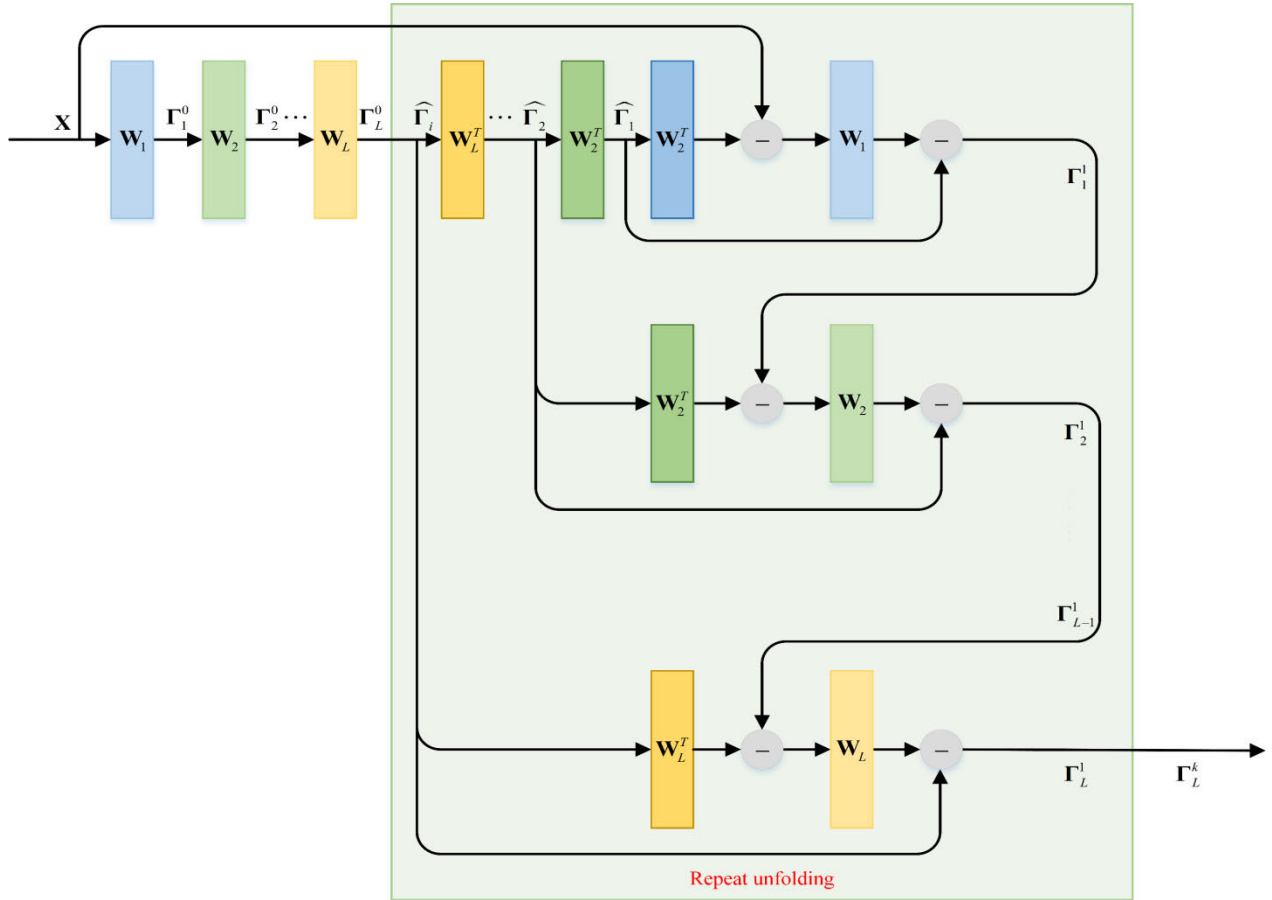


FIGURE 1. The ML-CSC block. $\{W_i\}_{i=1}^L$ denote the convolution operation, $\{W_i^T\}_{i=1}^L$ denote the deconvolution operation, L denotes the number of layers, and k denotes unfolding number.

the deconvolution operation, where $\{\mu_i\}_{i=1}^L$ and $\{\theta_i\}_{i=1}^L$ are trainable parameters. Solving $\{\Gamma_i\}_{i=1}^L$ can be seen as the process of extracting features from multi-layer convolution operation as follow:

$$\Gamma_L = \text{ReLU}(W_L \cdots \text{ReLU}(W_2 \text{ReLU}(W_1 y))) \quad (9)$$

When the number of unfoldings is set to 0, the layered thresholding algorithm is performed in the ML-CSC block. When the unfolding number is greater than 0, the ML-ISTA algorithm is executed. From the sparse point of view, due to ML-ISTA algorithm is superior to the layered thresholding algorithm, the ML-CSC block will extract more accurate feature compared with multi-layer convolution operation, which is beneficial to the forward propagation of the neural network, and also can better capture the semantic and appearance information of the image to improve the segmentation performance.

III. METHOD

A. U-NET MODEL

The architecture of U-Net model [27] is displayed in Figure 2(a). For convenience, we use 3×3 convolution layer with padding to keep the same size before and after

convolution operation, thus the input size of model is equal to the output size. The up-sampling is performed by 3×3 transposed convolution operation. Batch normalization [34] is after convolution operation and before ReLU activation function.

B. CSC-UNET MODEL SERIES

In the encode and decode side of the U-Net model, both of which can be seen as a composition of blocks containing two layers of convolution operation, as shown in Table 1. To fairly demonstrate our strategy, we set the number of layers in the ML-CSC block to 2 as well. According to the characteristics of encoding and decoding structure in U-Net, we designed the CSC-Unet-model series including CSC-Unet-Encode, CSC-Unet-Decode, and CSC-Unet-All model. The details of the CSC-Unet model series were also shown in Table 1.

1) CSC-UNET-ENCODE MODEL

The encoding side of U-Net is used to extract the semantic and appearance information of the input image. In order to explore this kind of ability of ML-CSC block, we replaced the convolution operations of the encoding side of the

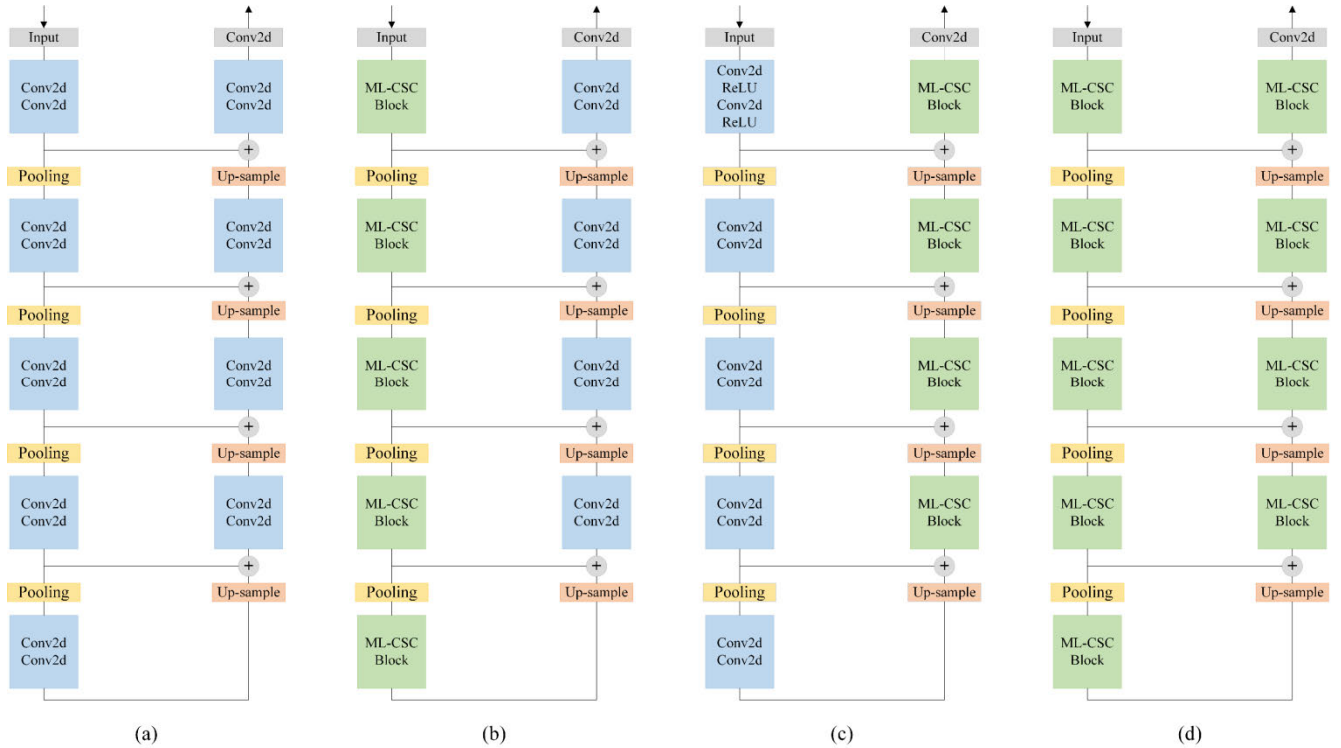


FIGURE 2. Structures of (a) U-Net, (b) our CSC-Unet-Encode, (c) our CSC-Unet-Decode, and (d) our CSC-Unet-All, respectively.

TABLE 1. The details of models.

Model	U-Net	CSC-Unet-Encode	CSC-Unet-Decode	CSC-Unet-All
Encoder	$[2 \times \text{Conv2d}] \times 5$	$[\text{ML-CSC block}] \times 5$	$[2 \times \text{Conv2d}] \times 5$	$[\text{ML-CSC block}] \times 5$
Decoder	$[2 \times \text{Conv2d}] \times 4$	$[2 \times \text{Conv2d}] \times 4$	$[\text{ML-CSC block}] \times 4$	$[\text{ML-CSC block}] \times 4$

U-Net with the ML-CSC blocks to form CSC-Unet-Encode model, and the corresponding architecture was shown in Figure 2(b).

2) CSC-UNET-DECODE MODEL

In the expanding sub-network of U-Net model, it firstly uses skip connection to combine appearance information from the shallow layers and semantic information from the deep layers. Then, the decoding side of U-Net precisely locates the segmentation boundary and gradually recovers the spatial detail information of the image. To explore the ability of ML-CSC block to recover the spatial detail information of image, we introduced the ML-CSC block into the decoding side of U-Net to form CSC-Unet-Encode model, and the corresponding architecture was shown in Figure 2(c).

3) CSC-UNET-ALL MODEL

To explore the impact of ML-CSC block on the overall segmentation performance of U-Net model, we added this block to both the encoding and decoding side of the U-Net to form CSC-Unet-All model, and the corresponding architecture was shown in Figure 2(d).

IV. EXPERIMENT AND ANALYSIS

In this study, the computing platform are Ubuntu 18.04.5 LTS 64-bit OS, 32G RAM, and Nvidia GeForce GTX 1080 Ti GPU with 11 GB memory. The deep learning framework is based on PyTorch [54] and the program language is Python.

A. DATASETS

For perform a fair evaluation, we select different scenarios datasets for testing to obtain evaluation metrics of comparison models. CamVid dataset [55] is one of the first datasets used for autonomous driving. It is assembled from 5 video sequences taken by the on-board camera from the driver’s perspective. DeepCrack dataset [56] is a public benchmark dataset containing cracks at multiple scales and scenarios to evaluate crack detection systems. Nuclei dataset¹ is the dataset in the 2018 Kaggle Data Science Bowl which is acquired under a variety of conditions and variations in the cell type, magnification, and imaging modality. The details of the datasets were shown in Table 2.

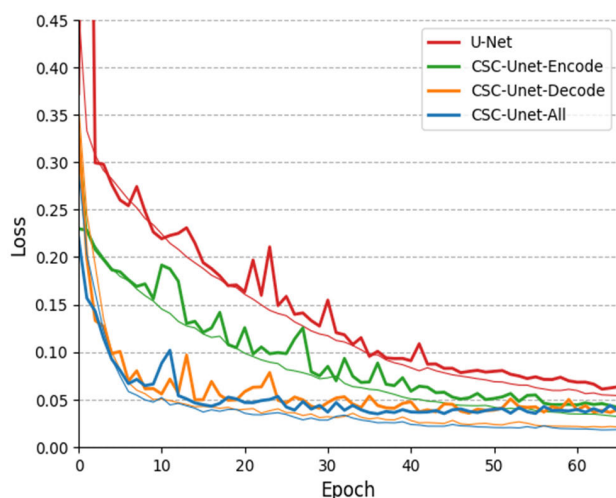
¹<https://www.kaggle.com/c/data-science-bowl-2018/overview>

TABLE 2. The details of datasets used in the experiment.

Dataset	Classes	Samples (training)	Samples (validation)	Samples (test)	Samples (total)
CamVid	11	367	101	233	701
DeepCrack	2	322	107	108	537
Nuclei	3	402	134	134	670

TABLE 3. Result on DeepCrack, Nuclei and CamVid test set of U-Net and CSC-Unet-Encode models with different unfolding number.

Method	DeepCrack Mean IoU (%)	Nuclei Mean IoU (%)	CamVid Mean IoU (%)
U-Net (CSC-Unet-Encode-0)	84.71	67.09	48.82
CSC-Unet-Encode-1	86.41	67.26	52.31
CSC-Unet-Encode-2	86.90	68.44	53.29
CSC-Unet-Encode-3	86.20	67.71	52.43

**FIGURE 3.** Training and validation on DeepCrack. Fine curves indicate the loss of training, and thick curves indicate the loss of validation. Unfolding number is uniformly set to 2 and the trainable parameters are same for all models.

B. THE SETTING OF TRAINING PARAMETERS

The number of all epochs were empirically set to 200 in this experiment. To improve the generalization ability of the model, before each epoch, we randomly disrupted the training data to make it more consistent with the sample distribution under natural conditions. The batch size was set to 4. The loss function was negative log-likelihood, and the input parameters were activated by the log-SoftMax function. The model used Adam [57] algorithm as the optimizer, each 50 epochs, and the learning rate dropped by half. In the CamVid dataset. The initial learning rate was set to 10^{-4} , for the DeepCrack and Nuclei, the initial learning rate was set to 10^{-5} , respectively.

C. EXPERIMENT AND ANALYSIS

1) THE SPEED OF MODEL CONVERGENCE

We investigated the effect of ML-CSC blocks on the convergence speed of segmentation models. On the DeepCrack dataset, we compared the training and validation loss of CSC-Unet model series in the training phase, and the results were shown in Figure 3. We found that ML-CSC blocks

can accelerate the convergence of the semantic segmentation model and reduce the loss value. Adding ML-CSC blocks at the decoding side of the U-Net model converged faster than adding them at the encoding side, and the model converged fastest when ML-CSC blocks were added at both sides of U-Net.

2) THE EXTRACTION OF SEMANTIC AND APPEARANCE INFORMATION

We have assessed the influence of ML-CSC block on CSC-Unet model series compared with U-Net to extract semantic and appearance information, and the results were shown in Table 3, where the number after the model indicated the unfolding number of ML-CSC block. When the number of unfoldings was zero, the CSC-Unet-Encode was equivalent to the U-Net model. The results showed that CSC-Unet-Encode models outperformed U-Net model on all three datasets when the number of unfoldings was greater than 0. This indicated that the ML-CSC block can indeed improve the ability of the semantic segmentation model to capture the semantic and appearance information of the image. Furthermore, we found that it was not always true that the larger number of unfoldings of the ML-CSC block implied better performance.

Semantic segmentation model first extracts feature information at the encoding side, and then based on feature information, the model gradually recovers the spatial detail information of the image at the decoding side. As the number of unfoldings increases, the feature information conveys in the model becomes sparser, which is not beneficial for the recovery process at the decoding side. Therefore, we should find a balance point between the extraction of feature information and the recovery process. According to Table 3, we inferred that the balance point was reached when the unfolding number was two among most datasets. For the convenience of performance demonstration, in our all experiments, the maximum number of unfoldings was set to three.

3) THE ABILITY TO RECOVER SPATIAL DETAIL INFORMATION

Next, we explored the ability of ML-CSC block to recover spatial detail information of image at the decoding side and the results were shown in Table 4. CSC-Unet-Decode models with unfolding number greater than 0 on different datasets

TABLE 4. Result on DeepCrack, Nuclei and CamVid test set of U-Net and CSC-Unet-Decode models with different unfolding number.

Method	DeepCrack Mean IoU (%)	Nuclei Mean IoU (%)	CamVid Mean IoU (%)
U-Net (CSC-Unet-Decode-0)	84.71	67.09	48.42
CSC-Unet-Decode-1	86.29	68.20	51.49
CSC-Unet-Decode-2	85.52	67.24	52.33
CSC-Unet-Decode-3	85.21	67.25	53.18

TABLE 5. Result on DeepCrack and Nuclei test set of U-Net and CSC-Unet-All models with different unfolding number.

Method	DeepCrack		Nuclei	
	Pixel Acc (%)	Mean IoU (%)	Pixel Acc (%)	Mean IoU (%)
U-Net (CSC-Unet-All-0-0)	98.53	84.71	96.64	67.09
CSC-Unet-All-2-1 (best)	98.74	87.14	96.81	68.91
CSC-Unet-All-1-1	98.62	86.61	96.67	67.30
CSC-Unet-All-2-2	98.72	87.04	96.73	68.31

TABLE 6. Results on CamVid test set of CSC-Unet-All models (1) U-Net(CSC-Unet-All-0-0), (2) CSC-Unet-All-2-1, (3) CSC-Unet-All-1-1 and (4) CSC-Unet-All-2-2.

Model	Sky	Building	Pole	Road	Sidewalk	Tree	Sign	Fence	Car	Pedestrian	Building	Class avg. (%)	Mean IOU (%)
(1)	96.6	75.9	27.7	96.4	75.1	81.2	50.1	16.3	87.0	43.6	10.9	60.1	48.82
(2)	96.1	84.6	33.4	97.6	85.2	81.2	45.6	18.4	86.0	54.8	12.5	63.2	53.56
(3)	95.9	80.5	33.6	97.8	85.8	82.8	41.8	15.2	85.5	60.6	14.5	63.1	52.76
(4)	95.5	87.5	36.6	97.6	83.3	82.9	38.4	16.5	86.0	59.6	12.2	63.3	53.68

TABLE 7. The computational cost of U-Net and CSC-Unet-model series with different unfolding number.

Model	Training Memory Usage (GB)		
	CamVid	DeepCrack	Nuclei
Unet	5.8	4.9	4.9
CSC-Unet-Encode-1	7.0	5.8	5.8
CSC-Unet-Encode-2	8.2	6.7	6.6
CSC-Unet-Encode-3	9.3	7.6	7.6
CSC-Unet-Decode-1	7.3	6.1	6.2
CSC-Unet-Decode-2	8.8	7.1	7.2
CSC-Unet-Decode-3	10.2	8.1	8.3
CSC-Unet-All-1-1	8.5	6.9	7.1
CSC-Unet-All-2-1	9.6	7.8	8.0
CSC-Unet-All-2-2	11.1	8.9	9.0

were better than U-Net, which implied that ML-CSC block improved the recovery of spatial detail information compared to the convolution operation. The best unfolding was 1 in DeepCrack and Nuclei, and in CamVid was 3, respectively. We speculated that phenomenon was probably related to the complexity of the image, where the categories of DeepCrack and Nuclei were relatively few and the case of one unfolding was enough, but CamVid was relatively more complex and required a higher number of unfoldings.

4) THE OVERALL IMPROVEMENT OF MODEL PERFORMANCE
We first introduced the nomenclature of CSC-Unet-All-a-b, where a denoted the unfolding number at the encoding

side and b denoted the number of unfoldings at the decoding side. For example, CSC-Unet-All-0-0 was equivalent to U-Net model, CSC-Unet-All-a-0 was represented as CSC-Unet-Encode-a, and CSC-Unet-All-0-b can be represented as CSC-Unet-Decode-b. Through Table 3 and 4, we found that on the DeepCrack dataset the CSC-Unet-Encode-2 captured more semantic and appearance information, and CSC-Unet-Decode-1 maximized the ability of the model to recover spatially detailed information. Thus, we argued that CSC-Unet-All-2-1 could maximize the segmentation performance of the U-Net model. Similarly, on Nuclei and CamVid, CSC-Unet-All-2-1 and CSC-Unet-All-2-3 should achieve the best segmentation performance. However, due to the GPU memory size limitation (11 GB), if we use the CSC-Unet-All-2-3 model on the CamVid dataset, the batch size needs to be halved, or we can use GPU parallelism to maintain the batch size. This is something that we do not expect to see. We want to compare the performance of the models under the same conditions; thus, we used CSC-Unet-All-2-2 instead of CSC-Unet-All-2-3. We also set CSC-Unet-All-0-0 (U-Net) and CSC-Unet-All-1-1 for performance comparison. The results on the three datasets were shown in Table 5 and 6. Compared to U-Net we have improved by 2.43% (from 84.71% to 87.14% in DeepCrack), 1.82% (from 67.09% to 68.91% in Nuclei), and 4.86% (from 48.82% to 53.68% in CamVid) on three datasets in terms of Mean Intersection over Union (MIoU), respectively. To better illustrate the

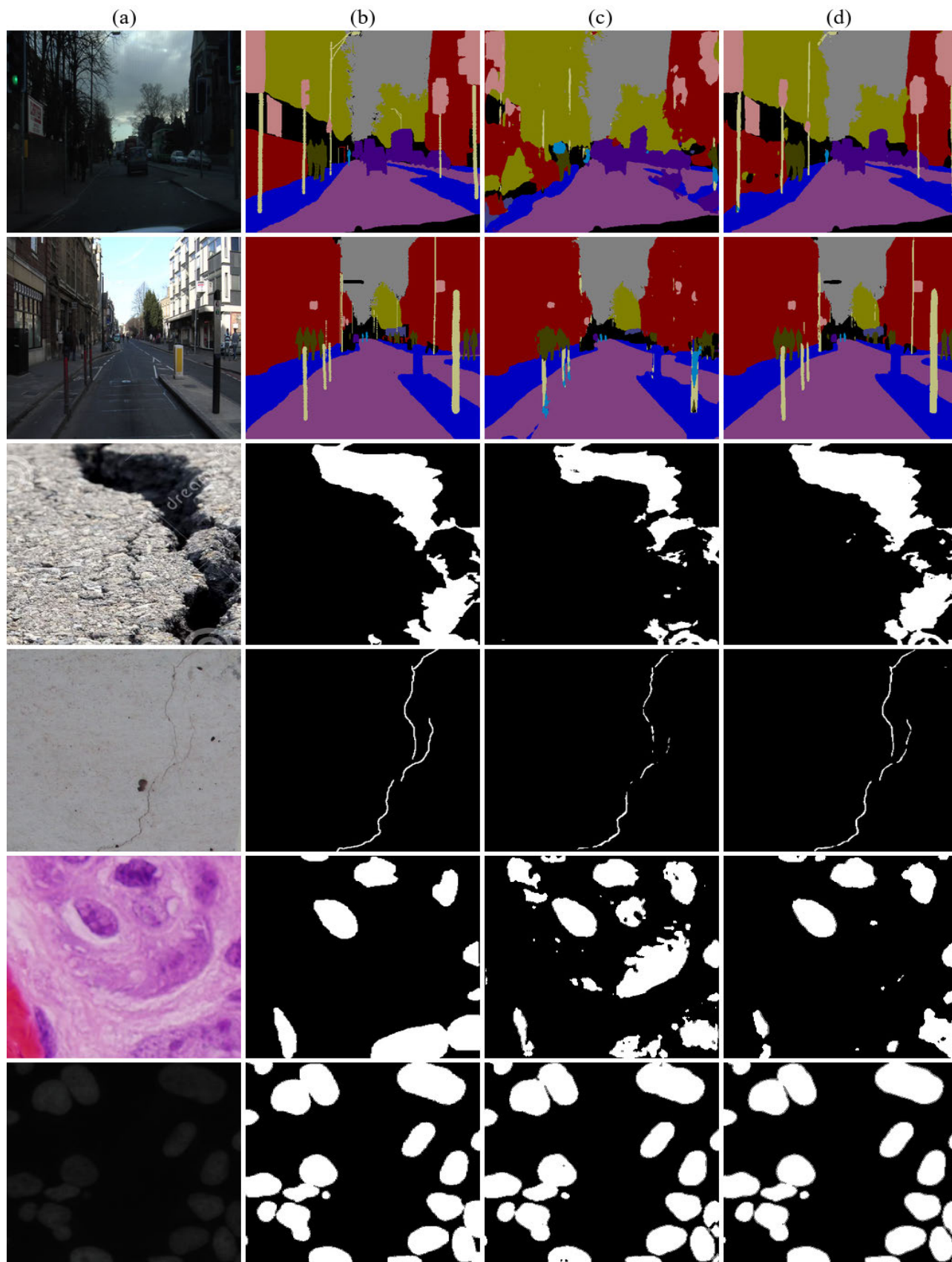


FIGURE 4. Examples of semantic segmentation results on CamVid, DeepCrack, and Nuclei test set. (a) Input images, (b) Ground truths, (c) Results of U-Net, and (d) Results of CSC-Unet-All (best).

TABLE 8. The model performance and computational cost of Unet++, Unet+ + +, DeepLabv3+ and their corresponding improved versions.

Model	Unfoldings	Mean IoU (%)			Training Memory Usage (GB)		
		CamVid	DeepCrack	Nuclei	CamVid	DeepCrack	Nuclei
Unet++	-	58.04	98.75	88.50	5.2	4.1	7.5
CSC-Unet++	1	65.44	98.98	89.91	9.0	6.9	13.3
CSC-Unet++	2	65.87	98.96	89.74	12.7	9.7	18.0
Unet3+	-	60.09	98.92	89.46	14.6	11.3	16.1
CSC-Unet3+	1	66.89	99.06	90.06	15.9	12.2	18.1
CSC-Unet3+	2	66.17	99.07	90.10	16.9	13.0	19.6
DeepLabv3+	-	58.98	98.81	86.92	8.5	6.5	6.4
CSC-DeepLabv3+	1	63.28	98.89	87.37	13.6	10.3	10.0
CSC-DeepLabv3+	2	64.69	99.12	88.01	17.6	13.1	12.0

results, the visualizations on the three datasets were shown in Figure 4.

V. CONCLUSION AND DISCUSSION

In this paper, we proposed a novel strategy, which used ML-CSC block instead of convolutional operation to improve the performance of semantic segmentation model. This strategy could possibly apply to any semantic segmentation model that involved convolutional operation. We used the U-Net model as an example to validate this strategy and designed CSC-Unet model series. We found that using ML-CSC blocks instead of convolution operations could accelerate the convergence of the semantic segmentation model, improve the ability of the model to capture the semantic and appearance information of the image, and improve the ability to recover spatial detail information. We concluded that ML-CSC block was a better operation compared to convolutional operation in semantic segmentation.

The current CSC-Unet models have achieved significant improvement in segmentation performance compared with the original U-Net model. However, they still face great challenges. Therefore, in the follow-up study, we will consider the following aspects.

1) To fairly demonstrate our strategy, the added ML-CSC block to the U-Net model was a two-layer convolutional sparse coding model. We speculate that more layers will be beneficial to improve the performance of semantic segmentation. Thus, in the next study, we will design a new type of semantic segmentation model based on multi-layer global convolutional sparse coding block.

2) We found that there is a balance between the extraction of feature information and the recovery process. Thus, how to find the best balance between all the unfolding number in separate encoder and decoder part will be one focus of our future research. In addition, the presence of the unfolding number in the solution algorithm slightly increases the model's computation cost, as shown in Table 7. So, the solution algorithm without unfolding number will be designed in our subsequent study.

3) The ML-CSC block can be extended to not only the field of semantic segmentation, but also possibly to other fields such as object detection [12], [13], [14], generative

adversarial networks (GAN) [58], natural language processing (NLP) [59], etc.

APPENDIX

The results of improving Unet++, Unet3+, and deeplabv3+ are shown in TABLE 8, with the same training strategy and input model sizes as CSC-Unet, and the code of CSC-Unet++, CSC-Unet+ + +, and CSC-DeepLabv3+ is also available at <https://github.com/NZWANG/CSC-Unet>.

REFERENCES

- [1] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 1–9, doi: [10.1109/CVPR.2015.7298594](https://doi.org/10.1109/CVPR.2015.7298594).
- [2] H. Tang, J. Shi, X. Lu, Z. Yin, L. Huang, D. Jia, and N. Wang, "Comparison of convolutional sparse coding network and convolutional neural network for pavement crack classification: A validation study," *J. Phys., Conf. Ser.*, vol. 1682, no. 1, Nov. 2020, Art. no. 012016, doi: [10.1088/1742-6596/1682/1/012016](https://doi.org/10.1088/1742-6596/1682/1/012016).
- [3] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, *arXiv:1409.1556*.
- [4] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Commun. ACM*, vol. 60, no. 6, pp. 84–90, May 2017.
- [5] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 770–778, doi: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- [6] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2014, pp. 580–587.
- [7] R. Girshick, "Fast R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 1440–1448.
- [8] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," 2015, *arXiv:1506.01497*.
- [9] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot MultiBox detector," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 21–37.
- [10] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 779–788, doi: [10.1109/CVPR.2016.91](https://doi.org/10.1109/CVPR.2016.91).
- [11] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 6517–6525, doi: [10.1109/CVPR.2017.690](https://doi.org/10.1109/CVPR.2017.690).
- [12] J. Redmon and A. Farhadi, "YOLOv3: An incremental improvement," 2018, *arXiv:1804.02767*.
- [13] A. Bochkovskiy, C.-Y. Wang, and H.-Y. Mark Liao, "YOLOv4: Optimal speed and accuracy of object detection," 2020, *arXiv:2004.10934*.
- [14] K. Liu, H. Tang, S. He, Q. Yu, Y. Xiong, and N. Wang, "Performance validation of YOLO variants for object detection," in *Proc. Int. Conf. Bioinf. Intell. Comput.*, Jan. 2021, pp. 239–243.

- [15] Y. Li, S. Yang, Y. Zheng, and H. Lu, "Improved point-voxel region convolutional neural network: 3D object detectors for autonomous driving," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 7, pp. 9311–9317, Jul. 2022.
- [16] H. Cui, V. Radosavljevic, F.-C. Chou, T.-H. Lin, T. Nguyen, T.-K. Huang, J. Schneider, and N. Djuric, "Multimodal trajectory predictions for autonomous driving using deep convolutional networks," in *Proc. Int. Conf. Robot. Autom. (ICRA)*, May 2019, pp. 2090–2096.
- [17] Z. Sheng, Y. Xu, S. Xue, and D. Li, "Graph-based spatial-temporal convolutional network for vehicle trajectory prediction in autonomous driving," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 10, pp. 17654–17665, Oct. 2022.
- [18] R. R. Devaram, G. Beraldo, R. De Benedictis, M. Mongiovi, and A. Cesta, "LEMON: A lightweight facial emotion recognition system for assistive robotics based on dilated residual convolutional neural networks," *Sensors*, vol. 22, no. 9, p. 3366, Apr. 2022.
- [19] M. Engelcke, D. Rao, D. Z. Wang, C. H. Tong, and I. Posner, "Vote3Deep: Fast object detection in 3D point clouds using efficient convolutional neural networks," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, May 2017, pp. 1355–1361.
- [20] C. McCool, T. Perez, and B. Upcroft, "Mixtures of lightweight deep convolutional neural networks: Applied to agricultural robotics," *IEEE Robot. Autom. Lett.*, vol. 2, no. 3, pp. 1344–1351, Jul. 2017.
- [21] X. Chen and A. Gupta, "Webly supervised learning of convolutional networks," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 1431–1439.
- [22] H. Dozono, K. Toyozumi, K. Yoshioka, and G. Niina, "The visualization system of image search base on convolutional spherical self organizing map implemented using WebGL," in *Proc. 6th Int. Congr. Inf. Commun. Technol.*, 2022, pp. 493–502.
- [23] R. Mishra and S. P. Tripathi, "Deep learning based search engine for biomedical images using convolutional neural networks," *Multimedia Tools Appl.*, vol. 80, no. 10, pp. 15057–15065, Apr. 2021.
- [24] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 3431–3440.
- [25] V. Badrinarayanan, A. Kendall, and R. Cipolla, "SegNet: A deep convolutional encoder-decoder architecture for image segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 12, pp. 2481–2495, Dec. 2017, doi: [10.1109/TPAMI.2016.2644615](https://doi.org/10.1109/TPAMI.2016.2644615).
- [26] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2017, pp. 2881–2890.
- [27] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI (Lecture Notes in Computer Science, Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 9351, 2015, pp. 234–241, doi: [10.1007/978-3-319-24574-4_28](https://doi.org/10.1007/978-3-319-24574-4_28).
- [28] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "UNet++: Redesigning skip connections to exploit multiscale features in image segmentation," *IEEE Trans. Med. Imag.*, vol. 39, no. 6, pp. 1856–1867, Jun. 2020.
- [29] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 936–944.
- [30] X. Qin, Z. Zhang, C. Huang, M. Dehghan, O. R. Zaiane, and M. Jagersand, "U2-Net: Going deeper with nested U-structure for salient object detection," *Pattern Recognit.*, vol. 106, Oct. 2020, Art. no. 107404.
- [31] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Semantic image segmentation with deep convolutional nets and fully connected CRFs," 2014, *arXiv:1412.7062*.
- [32] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "DeepLab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, Apr. 2018.
- [33] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," 2017, *arXiv:1706.05587*.
- [34] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 448–456.
- [35] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," 2018, *arXiv:1802.02611*.
- [36] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 1800–1807.
- [37] H. Huang, L. Lin, R. Tong, H. Hu, Q. Zhang, Y. Iwamoto, X. Han, Y.-W. Chen, and J. Wu, "UNet 3+: A full-scale connected UNet for medical image segmentation," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2020, pp. 1055–1059.
- [38] W. Dong, L. Zhang, G. Shi, and X. Li, "Nonlocally centralized sparse representation for image restoration," *IEEE Trans. Image Process.*, vol. 22, no. 4, pp. 1620–1630, Apr. 2013.
- [39] J. Wang, C. Lu, M. Wang, P. Li, S. Yan, and X. Hu, "Robust face recognition via adaptive sparse representation," *IEEE Trans. Cybern.*, vol. 44, no. 12, pp. 2368–2378, Dec. 2014.
- [40] J. Yang, J. Wright, T. S. Huang, and Y. Ma, "Image super-resolution via sparse representation," *IEEE Trans. Image Process.*, vol. 19, no. 11, pp. 2861–2873, Nov. 2010.
- [41] N. Wang, W. Zeng, and L. Chen, "SACICA: A sparse approximation coefficient-based ICA model for functional magnetic resonance imaging data analysis," *J. Neurosci. Methods*, vol. 216, no. 1, pp. 49–61, May 2013.
- [42] N. Wang, W. Zeng, Y. Shi, T. Ren, Y. Jing, J. Yin, and J. Yang, "WASICA: An effective wavelet-shrinkage based ICA model for brain fMRI data analysis," *J. Neurosci. Methods*, vol. 246, pp. 75–96, May 2015.
- [43] N. Wang, W. Zeng, and D. Chen, "A novel sparse dictionary learning separation (SDLS) model with adaptive dictionary mutual incoherence constraint for fMRI data analysis," *IEEE Trans. Biomed. Eng.*, vol. 63, no. 11, pp. 2376–2389, Nov. 2016.
- [44] S. Gu, W. Zuo, Q. Xie, D. Meng, X. Feng, and L. Zhang, "Convolutional sparse coding for image super-resolution," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Dec. 2015, pp. 1823–1831.
- [45] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus, "Deconvolutional networks," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, Jun. 2010, pp. 2528–2535.
- [46] M. Elad, *Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing*. New York, NY, USA: Springer, 2010.
- [47] B. K. Natarajan, "Sparse approximate solutions to linear systems," *SIAM J. Comput.*, vol. 24, no. 2, pp. 227–234, Apr. 1995.
- [48] E. J. Candes, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, Feb. 2006.
- [49] V. Pappas, Y. Romano, and M. Elad, "Convolutional neural networks analyzed via convolutional sparse coding," *J. Mach. Learn. Res.*, vol. 18, pp. 1–52, Jul. 2017.
- [50] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proc. 14th Int. Conf. Artif. Intell. Statist.*, 2011, pp. 315–323.
- [51] J. Sulam, A. Aberdam, A. Beck, and M. Elad, "On multi-layer basis pursuit, efficient algorithms and convolutional neural networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 8, pp. 1968–1980, Aug. 2020.
- [52] K. Bredies and D. A. Lorenz, "Linear convergence of iterative soft-thresholding," *J. Fourier Anal. Appl.*, vol. 14, nos. 5–6, pp. 813–837, Dec. 2008.
- [53] Z. Zhang, Y. Xu, J. Yang, X. Li, and D. Zhang, "A survey of sparse representation: Algorithms and applications," *IEEE Access*, vol. 3, pp. 490–530, 2015.
- [54] A. Paszke et al., "PyTorch: An imperative style, high-performance deep learning library," 2019, *arXiv:1912.01703*.
- [55] G. J. Brostow, J. Shotton, J. Fauqueur, and R. Cipolla, "Segmentation and recognition using structure from motion point clouds," in *Computer Vision—ECCV (Lecture Notes in Computer Science, Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 5302, 2008, pp. 44–57, doi: [10.1007/978-3-540-88682-2_5](https://doi.org/10.1007/978-3-540-88682-2_5).
- [56] Y. Liu, J. Yao, X. Lu, R. Xie, and L. Li, "DeepCrack: A deep hierarchical feature learning architecture for crack segmentation," *Neurocomputing*, vol. 338, pp. 139–153, Apr. 2019, doi: [10.1016/j.neucom.2019.01.036](https://doi.org/10.1016/j.neucom.2019.01.036).
- [57] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*.

- [58] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," 2014, *arXiv:1406.2661*.
- [59] R. Sequiera, G. Baruah, Z. Tu, S. Mohammed, J. Rao, H. Zhang, and J. Lin, "Exploring the effectiveness of convolutional neural networks for answer selection in end-to-end question answering," 2017, *arXiv:1707.07804*.



HAITONG TANG received the master's degree in surveying and mapping engineering from Jiangsu Ocean University, in 2022. He is currently a Deep Learning Algorithm Engineer with Jiangsu Dongying Intelligent Engineering Research Institute affiliated with Southeast University. His research interests include deep learning, computer vision, and convolutional sparse coding.



SHUANG HE received the B.S. and M.S. degrees from Jiangsu Ocean University, Lianyungang, China, in 2018 and 2022, respectively. He is currently pursuing the Ph.D. degree with the International Institute for Earth System Science, Nanjing University. His research interests include deep learning, remote sensing image processing, and hyperspectral vegetation quantitative analysis.



MENGDUO YANG received the B.S. degree in computer science and technology and the Ph.D. degree in software engineering from Soochow University, Suzhou, China, in 2011 and 2016, respectively.

Her research interests include medical image analysis and few-shot learning.



XIA LU received the B.S. degree in exploration engineering from Jilin University, Jilin, China, in 1997, and the M.S. degree in cartography and geographic information engineering and the Ph.D. degree in photogrammetry and remote sensing from China University of Geoscience (Beijing), Beijing, China, in 2004 and 2008, respectively.

She is currently a Professor with the School of Geography Science and Geomatics Engineering, Suzhou University of Science and Technology. Her research interests include three-dimensional monitoring and evaluation of marine environments and quantitative remote sensing of coastal wetland ecology.



QIN YU received the master's degree in pattern recognition and intelligent systems from Jiangsu Ocean University, in 2022. He is currently pursuing the Ph.D. degree in biomedical engineering with Shenzhen University. His research interests include neuroscience, image analysis, and emotions.



KAIYUE LIU received the B.S. degree from Shanxi University of Engineering and Technology, in 2018, and the M.S. degree from Jiangsu Ocean University, in 2023. He is currently pursuing the Ph.D. degree with Shenzhen University. His research interests include remote sensing, machine learning, object detection, and ROS systems.



HONGJIE YAN received the Ph.D. degree in clinical medicine from Shimane University, Japan, in 2016. She is currently an Attending Physician with the Department of Neurology, Affiliated Lianyungang Hospital of Xuzhou Medical University, Lianyungang, China.

Her research interests include AI-based diagnosis of neurological disorders, neuroimaging and aging, and neurodegenerative diseases.



NIZHUAN WANG received the B.S. degree in computer science and technology from Heilongjiang University, Harbin, China, in 2010, and the M.S. degree in computer application technology and the Ph.D. degree in communications, information engineering, and control from Shanghai Maritime University, Shanghai, China, in 2012 and January 2016, respectively.

Currently, he is a Research Assistant Professor with the Department of Chinese and Bilingual Studies, The Hong Kong Polytechnic University. Before that, he was with Shenzhen University (as an Assistant Professor), Jiangsu Ocean University (as a Full Professor), and ShanghaiTech University (as an Associate Researcher), from 2016 to 2024. He has published more than 60 scientific articles in many journals, such as *Human Brain Mapping*, *Journal of Neuroscience Methods*, *Magnetic Resonance Imaging*, and *IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS*. His research interests include intelligent computation, neural engineering, medical imaging, and hyperspectral imaging.

...