

ReCal-Net: Joint Region-Channel-Wise Calibrated Network for Semantic Segmentation in Cataract Surgery Videos*

Negin Ghamsarian¹, Mario Taschwer¹, Doris Putzgruber-Adamitsch², Stephanie Sarny², Yosuf El-Shabrawi², and Klaus Schöffmann¹[0000-0002-9218-1704]

¹ Department of Information Technology, Alpen-Adria-Universität Klagenfurt
{negin,mt,ks}@itec.aau.at

² Department of Ophthalmology, Klinikum Klagenfurt
{doris.putzgruber-adamitsch,stephanie.sarny,Yosuf.El-Shabrawi}@kabeg.at

Abstract. Semantic segmentation in surgical videos is a prerequisite for a broad range of applications towards improving surgical outcomes and surgical video analysis. However, semantic segmentation in surgical videos involves many challenges. In particular, in cataract surgery, various features of the relevant objects such as blunt edges, color and context variation, reflection, transparency, and motion blur pose a challenge for semantic segmentation. In this paper, we propose a novel convolutional module termed as *ReCal* module, which can calibrate the feature maps by employing region intra-and-inter-dependencies and channel-region cross-dependencies. This calibration strategy can effectively enhance semantic representation by correlating different representations of the same semantic label, considering a multi-angle local view centering around each pixel. Thus the proposed module can deal with distant visual characteristics of unique objects as well as cross-similarities in the visual characteristics of different objects. Moreover, we propose a novel network architecture based on the proposed module termed as *ReCal-Net*. Experimental results confirm the superiority of ReCal-Net compared to rival state-of-the-art approaches for all relevant objects in cataract surgery. Moreover, ablation studies reveal the effectiveness of the ReCal module in boosting semantic segmentation accuracy.

Keywords: Cataract Surgery · Semantic Segmentation · Feature Map Calibration.

1 Introduction

Cataract surgery is the procedure of returning a clear vision to the eye by removing the occluded natural lens, followed by implanting an intraocular lens (IOL). Being one of the most frequently performed surgeries, enhancing the outcomes of cataract surgery and diminishing its potential intra-operative and post-operative risks is of great importance. Accordingly, a large body of research has been focused on computerized surgical workflow analysis in cataract surgery [24,10,8,17,16,9], with a majority of approaches

* This work was funded by the FWF Austrian Science Fund under grant P 31486-N31.

relying on semantic segmentation. Hence, improving semantic segmentation accuracy in cataract surgery videos can play a leading role in the development of a reliable computerized clinical diagnosis or surgical analysis approach [19,18].

Semantic segmentation of the relevant objects in cataract surgery videos is quite challenging due to (i) transparency of the intraocular lens, (ii) color and contextual variation of the pupil and iris, (iii) blunt edges of the iris, and (iv) severe motion blur and reflection distortion of the instruments. In this paper, we propose a novel module for joint Region-channel-wise Calibration, termed as *ReCal* module. The proposed module can simultaneously deal with the various segmentation challenges in cataract surgery videos. In particular, the ReCal module is able to (1) employ multi-angle pyramid features centered around each pixel position to deal with transparency, blunt edges, and motion blur, (2) employ cross region-channel dependencies to handle texture and color variation through interconnecting the distant feature vectors corresponding to the same object. The proposed module can be added on top of every convolutional layer without changing the output feature dimensions. Moreover, the ReCal module does not impose a significant number of trainable parameters on the network and thus can be used after several layers to calibrate their output feature maps. Besides, we propose a novel semantic segmentation network based on the ReCal module termed as *ReCal-Net*. The experimental results show significant improvement in semantic segmentation of the relevant objects with ReCal-Net compared to the best-performing rival approach (85.38% compared to 83.32% overall IoU (intersection over union) for ReCal-Net vs. UNet++).

The rest of this paper is organized as follows. In Section 2, we briefly review state-of-the-art semantic segmentation approaches in the medical domain. In Section 3, we first discuss two convolutional blocks from which the proposed approach is inspired, and then delineate the proposed ReCal-Net and ReCal module. We detail the experimental settings in Section 4 and present the experimental results in Section 5. We finally conclude the paper in Section 6.

2 Related Work

Since many automatic medical diagnosis and image analysis applications entail semantic segmentation, considerable research has been conducted to improve medical image segmentation accuracy. In particular, U-Net [22] achieved outstanding segmentation accuracy being attributed to its skip connections. In recent years, many U-Net-base approaches have been proposed to address the weaknesses of the U-Net baseline [12,14,26,5,27]. UNet++ [27] is proposed to tackle automatic depth optimization via ensembling varying-depth U-Nets. MultiResUNet [14] factorizes large and computationally expensive kernels by fusing sequential convolutional feature maps. CPFNet [5] fuses the output features of parallel dilated convolutions (with different dilation rates) for scale awareness. The SegSE block [21] and scSE block [23], inspired by Squeeze-and-Excitation (SE) block [13], aim to recalibrate the channels via extracting inter-channel dependencies. The scSE block [23] further enhances feature representation via spatial channel pooling. Furthermore, many approaches are proposed to enhance semantic segmentation accuracy for particular medically relevant objects, in-

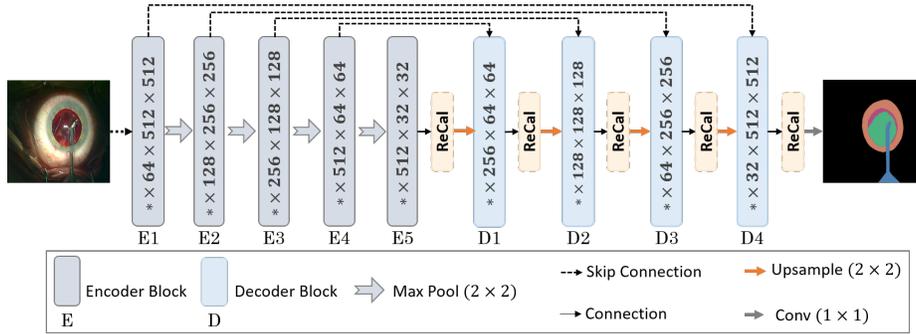


Fig. 1: The overall architecture of ReCal-Net containing five ReCal blocks.

cluding but not limited to liver lesion [2], surgical instruments [20,19,18], pulmonary vessel [3], and lung tumor [15].

3 Methodology

Notations. Everywhere in this paper, we show convolutional layer with the kernel-size of $(m \times n)$, P output channels, and g groups as $*_{(m \times n)}^{P,g}$ (we consider the default dilation rate of 1 for this layer). Besides, we show average-pooling layer with a kernel-size of $(m \times n)$ and a stride of s pixels as $\sum_{(m \times n)}^s$, and global average pooling as \sum^G .

Feature Map Recalibration. The Squeeze-and-Excitation (SE) block [13] was proposed to model inter-channel dependencies through squeezing the spatial features into a channel descriptor, applying fully-connected layers, and rescaling the input feature map via multiplication. This low-complexity operation unit has proved to be effective, especially for semantic segmentation. However, the SE block does not consider pixel-wise features in recalibration. Accordingly, scSE block [23] was proposed to exploit pixel-wise and channel-wise information concurrently. This block can be split into two parallel operations: (1) spatial squeeze and channel excitation, exactly the same as the SE block, and (2) channel squeeze and spatial excitation. The latter operation is conducted by applying a pixel-wise convolution with one output channel to the input feature map, followed by multiplication. The final feature maps of these two parallel computational units are then merged by selecting the maximum feature in each location.

ReCal-Net. Fig. 1 depicts the architecture of ReCal-Net. Overall, the network consists of three types of blocks: (i) encoder blocks that transform low-level features to semantic features while compressing the spatial representation, (ii) decoder blocks that are responsible for improving the semantic features in higher resolutions by employing the symmetric low-level feature maps from the encoder blocks, (iii) and *ReCal* modules that account for calibrating the semantic feature maps. We use the VGG16 network as the encoder network. The i th encoder block ($E_i, i \in \{1, 2, 3, 4\}$) in Fig. 1 correspond

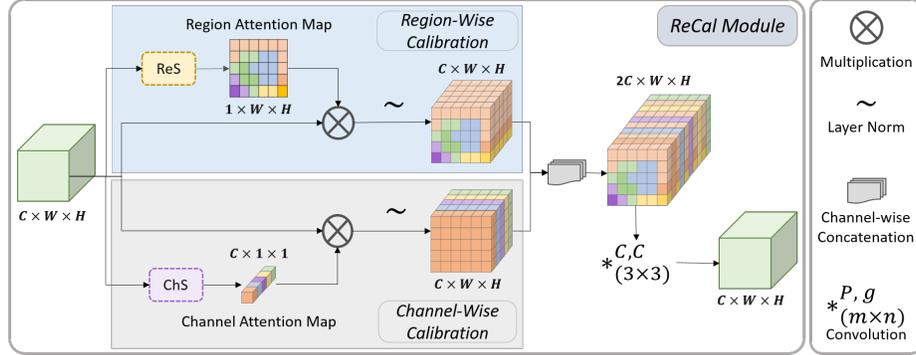


Fig. 2: The detailed architecture of ReCal block containing regional squeeze block (ReS) and channel squeeze block (ChS).

to all layers between the $i-1$ th and i th max-pooling layers in the VGG16 network (max-pooling layers are indicated with gray arrows). The last encoder block (E5) corresponds to the layers between the last max-pooling layer and the average pooling layer. Each decoder block follows the same architecture of decoder blocks in U-Net [22], including two convolutional layers, each of which being followed by batch normalization and ReLU.

ReCal Module. Despite the effectiveness of SE and scSE blocks in boosting feature representation, both fail to exploit region-wise dependencies. However, employing region-wise inter-dependencies and intra-dependencies can significantly enhance semantic segmentation performance. We propose a joint region-channel-wise calibration (ReCal) module to calibrate the feature maps based on joint region-wise and channel-wise dependencies. Fig. 2 demonstrates the architecture of the proposed ReCal module inspired by [13,23]. This module aims to reinforce a semantic representation considering inter-channel dependencies, inter-region and intra-region dependencies, and channel-region cross-dependencies. The input feature map of ReCal module $\mathcal{F}_{In} \in \mathbb{R}^{C \times W \times H}$ is first fed into two parallel blocks: (1) the Region-wise Squeeze block (ReS), and (2) the Channel-wise Squeeze block (ChS). Afterward, the region-wise and channel-wise calibrated features ($\mathcal{F}_{Re} \in \mathbb{R}^{C \times W \times H}$ and $\mathcal{F}_{Ch} \in \mathbb{R}^{C \times W \times H}$) are obtained by multiplying (\otimes) the input feature map to the region-attention map and channel-attention map, respectively, followed by the layer normalization function. In this stage, each particular channel $\mathcal{F}_{In}(C_j) \in \mathbb{R}^{W \times H}$ in the input feature map of a ReCal module has corresponding region-wise and channel-wise calibrated channels ($\mathcal{F}_{Re}(C_j) \in \mathbb{R}^{W \times H}$ and $\mathcal{F}_{Ch}(C_j) \in \mathbb{R}^{W \times H}$). To enable the utilization of cross-dependencies between the region-wise and channel-wise calibrated features, we concatenate these two feature maps in a depth-wise manner. Indeed, the concatenated feature map (\mathcal{F}_{Concat}) for each $p \in [1, C]$, $x \in [1, W]$, and $y \in [1, H]$ can be formulated as (1).

$$\begin{cases} \mathcal{F}_{Concat}(2p, x, y) = \mathcal{F}_{Re}(p, x, y) \\ \mathcal{F}_{Concat}(2p - 1, x, y) = \mathcal{F}_{Ch}(p, x, y) \end{cases} \quad (1)$$

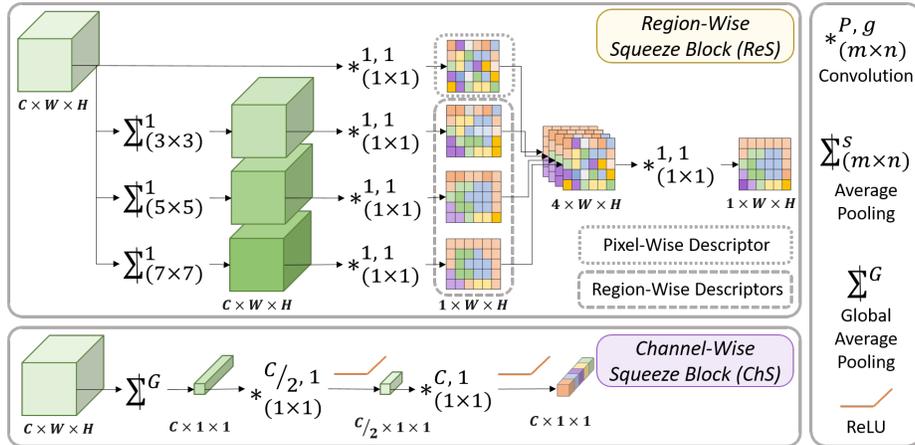


Fig. 3: Demonstration of regional squeeze block (ReS) and channel squeeze block (CS).

The cross-dependency between region-wise and channel-wise calibrated features is computed using a convolutional layer with C groups. More concretely, every two consecutive channels in the concatenated feature map undergo a distinct convolution with a kernel-size of (3×3) . This convolutional layer considers the local contextual features around each pixel (a 3×3 window around each pixel) to determine the contribution of each of region-wise and channel-wise calibrated features in the output features. Using a kernel size greater than one unit allows jointly considering inter-region dependencies.

Region-Wise Squeeze Block. Fig. 3 details the architecture of the ReS block, which is responsible for providing the region attention map. The region attention map is obtained by taking advantage of multi-angle local content based on narrow to wider views around each distinct pixel in the input feature map. We model multi-angle local features using average pooling layers with different kernel sizes and the stride of one pixel. The average pooling layers do not any number of impose trainable parameters on the network and thus ease using the ReS block and ReCal module in multiple locations. Besides, the stride of one pixel in the average pooling layer can stimulate a local view centered around each distinctive pixel. We use three average pooling layers with kernel-sizes of (3×3) , (5×5) , and (7×7) , followed by pixel-wise convolutions with one output channel ($\ast_{(1 \times 1)}^{1,1}$) to obtain the region-wise descriptors. In parallel, the input feature map undergoes another convolutional layer to obtain the pixel-wise descriptor. The local features can indicate if some particular features (could be similar or dissimilar to the centering pixel) exist in its neighborhood, and how large is the neighborhood of each pixel containing particular features. The four attention maps are then concatenated and fed into a convolutional layer ($\ast_{(1 \times 1)}^{1,1}$) that is responsible for determining the contribution of each spatial descriptor in the final region-wise attention map.

Channel-Wise Squeeze Block. For ChS Block, we follow a similar scheme as in [13]. At first, we apply global average pooling (\sum^G) on the input convolutional feature map. Afterward, we form a bottleneck via a pixel-wise convolution with C/r output channels ($*_{(1 \times 1)}^{C/r, 1}$) followed by ReLU non-linearity. The scaling parameter r can curbs the computational complexity. Besides, it can act as a smoothing factor that can yield a better-generalized model by preventing the network from learning outliers. In experiments, we set $r = 2$ as it is proved to have the best performance [23]. Finally, another pixel-wise convolution with C output channels ($*_{(1 \times 1)}^{C, 1}$) followed by ReLU non-linearity is used to provide the channel attention map.

Module Complexity. Suppose we have an intermediate layer in the network with convolutional response map $\mathcal{X} \in \mathbb{R}^{C \times H \times W}$. Adding a ReCal module on top of this layer with its scaling parameter being equal to 2, amounts to “ $C^2 + 22C + 4$ ” additional trainable weights. More specifically, each convolutional layer $*_{(m \times n)}^{P, g}$ applied to C input channels amounts to $((m \times n) \times C \times P)/g$ trainable weights. Accordingly, we need “ $4C + 4$ ” weights for the ReS block, “ C^2 ” weights for the ChS block, and “ $18C$ ” weights for the last convolution operation of the ReCal module. In our proposed architecture, adding five ReCal modules on convolutional feature maps with 512, 256, 128, 64, and 32 channels sums up to 371K additional weights, and only 21K more trainable parameters compared to the SE block [13] and scSE block [23].

4 Experimental Settings

Datasets. We use four datasets in this study. The iris dataset is created by annotating the cornea and pupil from 14 cataract surgery videos using “supervisely” platform. The iris annotations are then obtained by subtracting the convex-hull of the pupil segment from the cornea segment. This dataset contains 124 frames from 12 videos for training and 23 frames from two videos for testing³. For lens and pupil segmentation, we employ the two public datasets of the LensID framework [6], containing the annotation of the intraocular lens and pupil. The lens dataset consists of lens annotation in 401 frames sampled from 27 videos. From these annotations, 292 frames from 21 videos are used for training, and 109 frames from the remaining six videos are used for testing. The pupil segmentation dataset contains 189 frames from 16 videos. The training set consists of 141 frames from 13 videos, and the testing set contains 48 frames from three remaining videos. For instrument segmentation, we use the instrument annotations of the CaDIS dataset [11]. We use 3190 frames from 18 videos for training and 459 frames from three other videos for testing.

Rival Approaches. Table 1 lists the specifications of the rival state-of-the-art approaches used in our evaluations. In “Upsampling” column, “Trans Conv” stands for *Transposed Convolution*. To enable direct comparison between the ReCal module and scSE block, we have formed scSE-Net by replacing the ReCal modules in ReCal-Net with scSE modules. Indeed, the baseline of both approaches are the same, and the only

³ The dataset will be released with the acceptance of this paper.

Table 1: Specifications of the proposed and rival segmentation approaches.

Model	Backbone	Params	Upsampling	Reference	Year
UNet++ (DS)	VGG16	24.24 M	Bilinear	[27]	2020
MultiResUNet	✗	9.34 M	Trans Conv	[14]	2020
BARNet	ResNet34	24.90 M	Bilinear	[19]	2020
PAANet	ResNet34	22.43 M	Trans Conv & Bilinear	[18]	2020
CPFNet	ResNet34	34.66 M	Bilinear	[5]	2020
dU-Net	✗	31.98 M	Trans Conv	[26]	2020
CE-Net	ResNet34	29.90 M	Trans Conv	[12]	2019
scSE-Net	VGG16	22.90 M	Bilinear	[23]	2019
U-Net	✗	17.26 M	Bilinear	[22]	2015
ReCal-Net	VGG16	22.92 M	Bilinear	Proposed	

difference is the use of scSE blocks in scSE-Net at the position of ReCal modules in ReCal-Net.

Data Augmentation Methods. We use the Albumentations [1] library for image and mask augmentation during training. Considering the inherent features of the relevant objects and problems of the recorded videos [7], we apply motion blur, median blur, brightness and contrast change, shifting, scaling, and rotation for augmentation. We use the same augmentation pipeline for the proposed and rival approaches.

Neural Network Settings. We initialize the parameters of backbones for the proposed and rival approaches (in case of having a backbone) with ImageNet [4] training weights. We set the input size of all networks to $(3 \times 512 \times 512)$.

Training Settings. During training with all networks, a threshold of 0.1 is applied for gradient clipping. This strategy can prevent the gradient from exploding and result in a more appropriate behavior during learning in the case of irregularities in the loss landscape. Considering the different depths and connections of the proposed and rival approaches, all networks are trained with two different initial learning rates ($lr \in \{0.005, 0.002\}$) for 30 epochs with SGD optimizer. The learning scheduler decreases the learning rate every other epoch with a factor of 0.8. We list the results with the highest IoU for each network.

Loss Function. To provide a fair comparison, we use the same loss function for all networks. The loss function is set to a weighted sum of binary cross-entropy (BCE) and the logarithm of soft Dice coefficient as follows.

$$\mathcal{L} = (\lambda) \times BCE(\mathcal{X}_{true}(i, j), \mathcal{X}_{pred}(i, j)) - (1 - \lambda) \times \left(\log \frac{2 \sum \mathcal{X}_{true} \odot \mathcal{X}_{pred} + \sigma}{\sum \mathcal{X}_{true} + \sum \mathcal{X}_{pred} + \sigma} \right) \quad (2)$$

Soft Dice refers to the dice coefficient computed directly based on predicted probabilities rather than the predicted binary masks after thresholding. In (2), \mathcal{X}_{true} refers to the ground truth mask, \mathcal{X}_{pred} refers to the predicted mask, \odot refers to Hadamard product

Table 2: Quantitative comparisons among the semantic segmentation results of ReCal-Net and rival approaches based on IoU(%).

Network	Lens	Pupil	Iris	Instruments	Overall
U-Net	61.89 \pm 20.93	83.51 \pm 20.24	65.89 \pm 16.93	60.78 \pm 26.04	68.01 \pm 21.03
CE-Net	78.51 \pm 11.56	92.07 \pm 4.24	71.74 \pm 6.19	69.44 \pm 17.94	77.94 \pm 9.98
dU-Net	60.39 \pm 29.36	68.03 \pm 35.95	70.21 \pm 12.97	61.24 \pm 27.64	64.96 \pm 26.48
scSE-Net	86.04 \pm 11.36	96.13 \pm 2.10	78.58 \pm 9.61	71.03 \pm 23.25	82.94 \pm 11.58
CPFNet	80.65 \pm 12.16	93.76 \pm 2.87	77.93 \pm 5.42	69.46 \pm 17.88	80.45 \pm 9.58
BARNet	80.23 \pm 14.57	93.64 \pm 4.11	75.80 \pm 8.68	69.76 \pm 21.29	79.86 \pm 12.16
PAANet	80.30 \pm 11.73	94.35 \pm 3.88	75.73 \pm 11.67	68.01 \pm 22.29	79.59 \pm 12.39
MultiResUNet	61.42 \pm 19.91	76.46 \pm 29.43	49.99 \pm 28.73	61.01 \pm 26.94	62.22 \pm 26.25
UNet++/DS	84.53 \pm 13.42	96.18 \pm 2.62	74.01 \pm 13.13	65.99 \pm 25.66	79.42 \pm 14.75
UNet++	85.74 \pm 11.16	96.50 \pm 1.51	81.98 \pm 6.96	69.07 \pm 23.89	83.32 \pm 10.88
ReCal-Net	87.94 \pm 10.72	96.58 \pm 1.30	85.13 \pm 3.98	71.89 \pm 19.93	85.38 \pm 8.98

Table 3: Quantitative comparisons among the semantic segmentation results of ReCal-Net and rival approaches based on Dice(%).

Network	Lens	Pupil	Iris	Instruments	Overall
U-Net	73.86 \pm 20.39	89.36 \pm 15.07	78.12 \pm 13.01	71.50 \pm 25.77	78.21 \pm 18.56
CE-Net	87.32 \pm 9.98	95.81 \pm 2.39	83.39 \pm 4.25	80.30 \pm 15.97	86.70 \pm 8.15
dU-Net	69.99 \pm 29.40	73.72 \pm 34.24	81.76 \pm 9.73	71.30 \pm 27.62	74.19 \pm 25.24
scSE-Net	91.95 \pm 9.14	98.01 \pm 1.10	87.66 \pm 6.35	80.18 \pm 21.49	89.45 \pm 9.52
CPFNet	88.61 \pm 10.20	96.76 \pm 1.53	87.48 \pm 3.60	80.33 \pm 15.85	88.29 \pm 7.79
BARNet	88.16 \pm 10.87	96.66 \pm 2.30	85.95 \pm 5.73	79.72 \pm 19.95	87.62 \pm 9.71
PAANet	88.46 \pm 9.59	97.05 \pm 2.16	85.62 \pm 8.50	78.15 \pm 21.51	87.32 \pm 10.44
MultiResUNet	73.88 \pm 18.26	82.45 \pm 25.49	61.78 \pm 25.96	71.35 \pm 26.88	72.36 \pm 24.14
UNet++/DS	90.80 \pm 11.41	98.03 \pm 1.41	84.38 \pm 9.06	75.64 \pm 25.38	87.21 \pm 11.81
UNet++	91.80 \pm 11.16	98.26 \pm 0.79	89.93 \pm 4.51	78.54 \pm 22.76	89.63 \pm 9.80
ReCal-Net	93.09 \pm 8.56	98.26 \pm 0.68	91.91 \pm 2.47	81.62 \pm 17.75	91.22 \pm 7.36

(element-wise multiplication), and σ refers to the smoothing factor. In experiments, we set $\lambda = 0.8$ and $\sigma = 1$.

5 Experimental Results

Table 2 and Table 3 compare the segmentation performance of ReCal-Net and ten rival state-of-the-art approaches based on the average and standard deviation of IoU and Dice coefficient, respectively ⁴. Overall, ReCal-Net, UNet++, scSE-Net, and CPFNet have shown the top four segmentation results. Moreover, the experimental results reveal that ReCal-Net has achieved the highest average IoU and Dice coefficient for all relevant objects compared to state-of-the-art approaches. Considering the IoU report, ReCal-Net has gained considerable enhancement in segmentation performance compared to the second-best approach in lens segmentation (87.94% vs. 86.04% for scSE-Net) and iris segmentation (85.13% vs. 81.98% for UNet++). Having only 21k more trainable parameters than scSE-Net (0.08% additive trainable parameters), ReCal-Net has achieved

⁴ The ‘‘Overall’’ column in Table 2 and Table 3 is the average of the other four values.

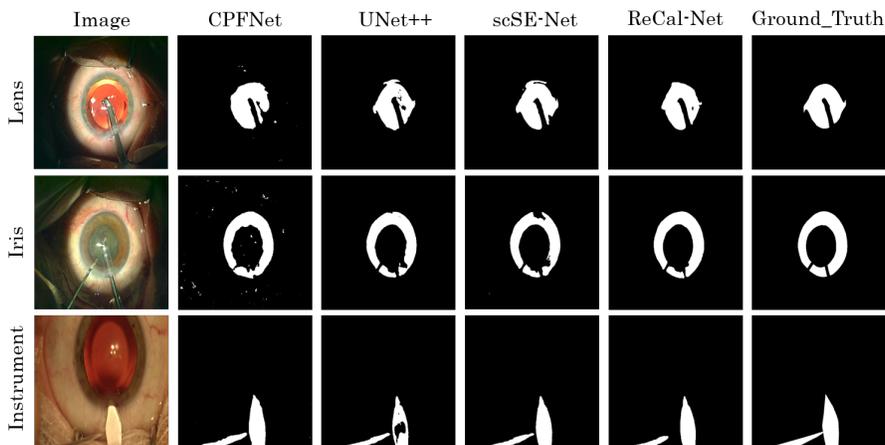


Fig. 4: Qualitative comparisons among the top four segmentation approaches.

8.3% relative improvement in iris segmentation, 2.9% relative improvement in instrument segmentation, and 2.2% relative improvement in lens segmentation in comparison with scSE-Net. Regarding the Dice coefficient, ReCal-Net and UNet++ show very similar performance in pupil segmentation. However, with 1.32M fewer parameters than UNet++ as the second-best approach, ReCal-Net shows 1.7% relative improvement in overall Dice coefficient (91.22% vs. 89.63%). Surprisingly, replacing the scSE blocks with the ReCal modules results in 4.25% higher Dice coefficient for iris segmentation and 1.44% higher Dice coefficient for instrument segmentation.

Fig. 4 provides qualitative comparisons among the top four segmentation approaches for lens, iris, and instrument segmentation. Comparing the visual segmentation results of ReCal-Net and scSE-Net further highlights the effectiveness of region-wise and cross channel-region calibration in boosting semantic segmentation performance.

Table 4 reports the ablation study by comparing the segmentation performance of the baseline approach with ReCal-Net considering two different learning rates. The baseline approach refers to the network obtained after removing all ReCal modules of ReCal-Net in Fig. 1. These results approve of the ReCal module’s effectiveness regardless of the learning rate.

To further investigate the impact of the ReCal modules on segmentation performance, we have visualized two intermediate filter response maps for iris segmentation in Fig. 5. The E5 output corresponds to the filter response map of the last encoder block, and the D1 output corresponds to the filter response map of the first decoder block (see Fig. 1). A comparison between the filter response maps of the baseline and ReCal-Net in the same locations indicated the positive impact of the ReCal modules on the network’s semantic discrimination capability. Indeed, employing the correlations between the pixel-wise, region-wise, and channel-wise descriptors can reinforce the network’s semantic interpretation.

Table 4: Impact of adding ReCal modules on the segmentation accuracy based on IoU(%).

Learning Rate	Network	Lens	Iris	Instrument
0.002	Baseline	84.83 \pm 11.62	81.49 \pm 6.82	70.04 \pm 23.94
	ReCal-Net	85.77 \pm 12.33	83.29 \pm 5.82	71.89 \pm 19.93
0.005	Baseline	86.13 \pm 11.63	81.00 \pm 8.06	67.16 \pm 24.67
	ReCal-Net	87.94 \pm 10.72	85.13 \pm 3.98	70.43 \pm 21.17

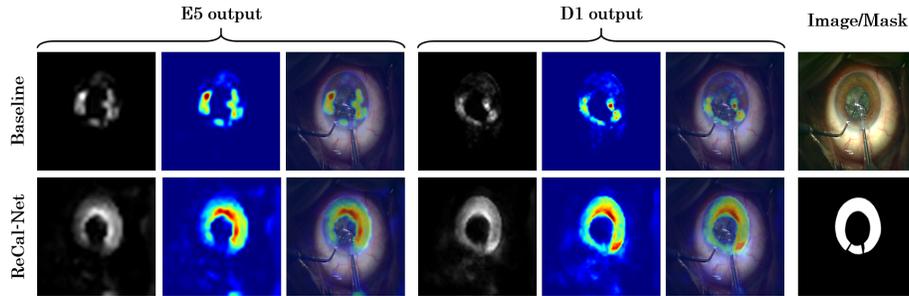


Fig. 5: Visualizations of the intermediate outputs in the baseline approach and ReCal-Net based on class activation maps [25]. For each output, the figures from left to right represent the gray-scale activation maps, heatmaps, and heatmaps on images.

6 Conclusion

This paper presents a novel convolutional module, termed as ReCal module, that can adaptively calibrate feature maps considering pixel-wise, region-wise, and channel-wise descriptors. The ReCal module can effectively correlate intra-region information and cross-channel-region information to deal with severe contextual variations in the same semantic labels and contextual similarities between different semantic labels. The proposed region-channel recalibration module is a very light-weight computational unit that can be applied to any feature map $\mathcal{X} \in \mathbb{R}^{C \times H \times W}$ and output a recalibrated feature map $\mathcal{Y} \in \mathbb{R}^{C \times H \times W}$. Moreover, we have proposed a novel network architecture based on the ReCal module for semantic segmentation in cataract surgery videos, termed as ReCal-Net. The experimental evaluations confirm the effectiveness of the proposed ReCal module and ReCal-Net in dealing with various segmentation challenges in cataract surgery. The proposed ReCal module and ReCal-Net can be adopted for various medical image segmentation and general semantic segmentation problems.

References

1. Buslaev, A., Iglovikov, V.I., Khvedchenya, E., Parinov, A., Druzhinin, M., Kalinin, A.A.: Albumentations: Fast and flexible image augmentations. *Information* **11**(2), 125 (Feb 2020). <https://doi.org/10.3390/info11020125>, <http://dx.doi.org/10.3390/info11020125>

2. Chen, X., Zhang, R., Yan, P.: Feature fusion encoder decoder network for automatic liver lesion segmentation. In: 2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019). pp. 430–433 (2019). <https://doi.org/10.1109/ISBI.2019.8759555>
3. Cui, H., Liu, X., Huang, N.: Pulmonary vessel segmentation based on orthogonal fused u-net++ of chest ct images. In: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (eds.) Medical Image Computing and Computer Assisted Intervention – MICCAI 2019. pp. 293–300. Springer International Publishing, Cham (2019)
4. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
5. Feng, S., Zhao, H., Shi, F., Cheng, X., Wang, M., Ma, Y., Xiang, D., Zhu, W., Chen, X.: Cpfnet: Context pyramid fusion network for medical image segmentation. *IEEE Transactions on Medical Imaging* **39**(10), 3008–3018 (2020). <https://doi.org/10.1109/TMI.2020.2983721>
6. Ghamsarian, N., Taschwer, M., Putzgruber-Adamitsch, D., Sarny, S., El-Shabrawi, Y., Schoeffmann, K.: Lensid: A cnn-rnn-based framework towards lens irregularity detection. In: 24th International Conference on Medical Image Computing & Computer Assisted Interventions (MICCAI 2021). p. to appear (2021)
7. Ghamsarian, N., Taschwer, M., Schoeffmann, K.: Deblurring cataract surgery videos using a multi-scale deconvolutional neural network. In: 2020 IEEE 17th International Symposium on Biomedical Imaging (ISBI). pp. 872–876 (2020)
8. Ghamsarian, N., Amirpourazarian, H., Timmerer, C., Taschwer, M., Schöffmann, K.: Relevance-based compression of cataract surgery videos using convolutional neural networks. In: Proceedings of the 28th ACM International Conference on Multimedia. p. 3577–3585. MM '20, Association for Computing Machinery, New York, NY, USA (2020). <https://doi.org/10.1145/3394171.3413658>, <https://doi.org/10.1145/3394171.3413658>
9. Ghamsarian, N., Amirpourazarian, H., Timmerer, C., Taschwer, M., Schöffmann, K.: Relevance-based compression of cataract surgery videos using convolutional neural networks. In: Proceedings of the 28th ACM International Conference on Multimedia. p. 3577–3585. MM '20, Association for Computing Machinery, New York, NY, USA (2020). <https://doi.org/10.1145/3394171.3413658>, <https://doi.org/10.1145/3394171.3413658>
10. Ghamsarian, N., Taschwer, M., Putzgruber-Adamitsch, D., Sarny, S., Schoeffmann, K.: Relevance detection in cataract surgery videos by spatio-temporal action localization. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 10720–10727 (2021). <https://doi.org/10.1109/ICPR48806.2021.9412525>
11. Grammatikopoulou, M., Flouty, E., Kadkhodamohammadi, A., Quelled, G., Chow, A., Nehme, J., Luengo, I., Stoyanov, D.: Cadis: Cataract dataset for image segmentation (2020)
12. Gu, Z., Cheng, J., Fu, H., Zhou, K., Hao, H., Zhao, Y., Zhang, T., Gao, S., Liu, J.: Ce-net: Context encoder network for 2d medical image segmentation. *IEEE Transactions on Medical Imaging* **38**(10), 2281–2292 (2019). <https://doi.org/10.1109/TMI.2019.2903562>
13. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7132–7141 (2018). <https://doi.org/10.1109/CVPR.2018.00745>
14. Ibtehaz, N., Rahman, M.S.: Multiresunet : Rethinking the u-net architecture for multimodal biomedical image segmentation. *Neural Networks* **121**, 74–87 (2020). <https://doi.org/https://doi.org/10.1016/j.neunet.2019.08.025>, <https://www.sciencedirect.com/science/article/pii/S0893608019302503>
15. Jiang, J., Hu, Y.C., Liu, C.J., Halpenny, D., Hellmann, M.D., Deasy, J.O., Mageras, G., Veeraghavan, H.: Multiple resolution residually connected feature streams for automatic lung

- tumor segmentation from ct images. *IEEE Transactions on Medical Imaging* **38**(1), 134–144 (2019). <https://doi.org/10.1109/TMI.2018.2857800>
16. Jin, Y., Li, H., Dou, Q., Chen, H., Qin, J., Fu, C.W., Heng, P.A.: Multi-task recurrent convolutional network with correlation loss for surgical video analysis. *Medical Image Analysis* **59**, 101572 (2020). <https://doi.org/https://doi.org/10.1016/j.media.2019.101572>, <https://www.sciencedirect.com/science/article/pii/S1361841519301124>
 17. Marafioti, A., Hayoz, M., Gallardo, M., Neila, P.M., Wolf, S., Zinkernagel, M., Sznitman, R.: Catanet: Predicting remaining cataract surgery duration (2021)
 18. Ni, Z.L., Bian, G.B., Wang, G.A., Zhou, X.H., Hou, Z.G., Chen, H.B., Xie, X.L.: Pyramid attention aggregation network for semantic segmentation of surgical instruments. *Proceedings of the AAAI Conference on Artificial Intelligence* **34**(07), 11782–11790 (Apr 2020). <https://doi.org/10.1609/aaai.v34i07.6850>, <https://ojs.aaai.org/index.php/AAAI/article/view/6850>
 19. Ni, Z.L., Bian, G.B., Wang, G.A., Zhou, X.H., Hou, Z.G., Xie, X.L., Li, Z., Wang, Y.H.: Barnet: Bilinear attention network with adaptive receptive fields for surgical instrument segmentation. In: Bessiere, C. (ed.) *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*. pp. 832–838. International Joint Conferences on Artificial Intelligence Organization (7 2020). <https://doi.org/10.24963/ijcai.2020/116>, <https://doi.org/10.24963/ijcai.2020/116>, main track
 20. Ni, Z.L., Bian, G.B., Zhou, X.H., Hou, Z.G., Xie, X.L., Wang, C., Zhou, Y.J., Li, R.Q., Li, Z.: Raunet: Residual attention u-net for semantic segmentation of cataract surgical instruments. In: Gedeon, T., Wong, K.W., Lee, M. (eds.) *Neural Information Processing*. pp. 139–149. Springer International Publishing, Cham (2019)
 21. Pereira, S., Pinto, A., Amorim, J., Ribeiro, A., Alves, V., Silva, C.A.: Adaptive feature recombination and recalibration for semantic segmentation with fully convolutional networks. *IEEE Transactions on Medical Imaging* **38**(12), 2914–2925 (2019). <https://doi.org/10.1109/TMI.2019.2918096>
 22. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (eds.) *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*. pp. 234–241. Springer International Publishing, Cham (2015)
 23. Roy, A.G., Navab, N., Wachinger, C.: Recalibrating fully convolutional networks with spatial and channel “squeeze and excitation” blocks. *IEEE Transactions on Medical Imaging* **38**(2), 540–549 (2019). <https://doi.org/10.1109/TMI.2018.2867261>
 24. Twinanda, A.P., Shehata, S., Mutter, D., Marescaux, J., de Mathelin, M., Padoy, N.: EndoNet: A Deep Architecture for Recognition Tasks on Laparoscopic Videos. *IEEE Transactions on Medical Imaging* **36**(1), 86–97 (Jan 2017). <https://doi.org/10.1109/TMI.2016.2593957>
 25. Wang, H., Wang, Z., Du, M., Yang, F., Zhang, Z., Ding, S., Mardziel, P., Hu, X.: Score-cam: Score-weighted visual explanations for convolutional neural networks. In: *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. pp. 111–119 (2020). <https://doi.org/10.1109/CVPRW50498.2020.00020>
 26. Zhang, M., Li, X., Xu, M., Li, Q.: Automated semantic segmentation of red blood cells for sickle cell disease. *IEEE Journal of Biomedical and Health Informatics* **24**(11), 3095–3102 (2020). <https://doi.org/10.1109/JBHI.2020.3000484>
 27. Zhou, Z., Siddiquee, M.M.R., Tajbakhsh, N., Liang, J.: Unet++: Redesigning skip connections to exploit multiscale features in image segmentation. *IEEE Transactions on Medical Imaging* **39**(6), 1856–1867 (2020). <https://doi.org/10.1109/TMI.2019.2959609>