

Channelized Axial Attention for Semantic Segmentation

Ye Huang, *Member, IEEE*, Wenjing Jia, *Member, IEEE*, Xiangjian He, *Member, IEEE*, Liu Liu, *Member, IEEE*, Yuxin Li, and Dacheng Tao, *Fellow, IEEE*

Abstract—Self-attention and channel attention, modelling the semantic interdependencies in spatial and channel dimensions respectively, have recently been widely used for semantic segmentation. However, computing spatial-attention and channel attention separately and then fusing them directly can cause conflicting feature representations. In this paper, we propose the Channelized Axial Attention (CAA) to seamlessly integrate channel attention and axial attention with reduced computational complexity. After computing axial attention maps, we propose to channelize the intermediate results obtained from the transposed dot-product so that the channel importance of each axial representation is optimized across the whole receptive field. We further develop grouped vectorization, which allows our model to be run with very little memory consumption at a speed comparable to the full vectorization. Comparative experiments conducted on multiple benchmark datasets, including Cityscapes, PASCAL Context and COCO-Stuff, demonstrate that our CAA not only requires much less computation resources compared with other dual attention models such as DANet [1], but also outperforms the state-of-the-art ResNet-101-based segmentation models on all tested datasets.

Index Terms—Semantic Segmentation, Axial Attention, Channelization, Grouped Vectorization.

I. INTRODUCTION

SEMANTIC segmentation is a fundamental task in many computer vision applications, which assigns a class label to each pixel in the image. Most of the existing approaches for semantic segmentation (e.g., [2], [3], [4], [5], [1], [6]) have adopted a pipeline similar to the one that is defined by Fully Convolutional Networks (FCNs) [7] using fully convolutional layers to output pixel-level segmentation results of the input image, and have achieved state-of-the-art performance. After the FCN approach, there have been many approaches dedicated to extracting enhanced pixel representations from the backbone. Earlier approaches, including PSPNet [8] and DeepLab [9], used a Pyramid Pooling Module (PPM) or an Atrous Spatial Pyramid Pooling (ASPP) module to expand the receptive field and capture multiple-range information to enhance the representation capabilities.

The latest research works on segmentation head in the past few years have mainly focused on using the attention mechanisms to improve the performance. During the early days of the attention mechanisms, the Squeeze and Excitation Networks (SENet) [10] introduced a simple and yet efficient channel attention module to explicitly model the interdependencies between channels. Meanwhile, the Non-Local Networks in [11] proposed self-attention to model long-range dependencies in spatial domain, so as to produce more correct

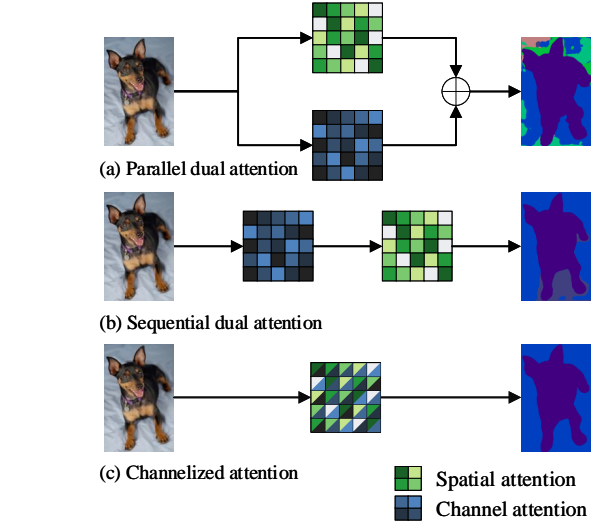


Fig. 1: Different dual attention designs

pixel representations. Thus, for each pixel in the feature maps, spatial self-attention makes its representation more similar to the representations of the pixels that are closer, whereas channel attention finds important channels in the entire feature maps and applies different weights to the extracted features.

To enjoy the advantages of both spatial attention and channel attention, some approaches (e.g., [1]) proposed to directly fuse their results with an element-wise addition (see Fig. 1(a)). Although they have achieved improved performance, the relationship between the contributions of spatial self-attention and channel attention to the final result is unclear.

Moreover, calculating the two attentions separately not only increases the computational complexity, but may also result in conflicting importance of feature representations. For example, some channels may appear to be important in channel attention for a pixel that belongs to a partial region in the feature maps, But spatial attention may have its own perspective, which is calculated by summing up the similarities over the entire feature maps, and weakens the impact of the channel attention.

The existing designs (e.g. [12]) combining channel attention and spatial attention in a cascaded, sequential manner (Fig. 1(b)) have similar issues. Channel attention can **ignore** the partial region representation obtained from the overall perspective, which may be required by spatial attention. Thus, directly fusing the spatial attention results with channel atten-

tion results may yield incorrect importance weights for pixel representations. In the Experiments section of this paper, we develop an approach to visualize the impact of the conflicting feature representation on the final segmentation results.

Attempting to combine the advantages of spatial-attention and channel attention seamlessly and efficiently in a complementary way, we propose a Channelized Axial Attention (CAA), which is based on a redefined axial attention to reduce the computation cost of self-attention. Specifically, when applying the redefined axial attention maps to the input signal [11], we capture the intermediate results of the dot product before they are summed up along the corresponding axes. Capturing these intermediate results allows channel attention to be integrated for each column and each row, instead of computing on the mean or sum of the features in the entire feature maps. More importantly, when applying the attention maps, we propose a novel transposed approach, which allows the channel attention to be conducted in the whole receptive field. Last but not the least, we develop a novel grouped vectorization approach to maximize the computation speed under limited GPU memory.

In summary, our contributions of this paper include:

- We propose a novel Channelized Axial Attention, which integrates spatial self-attention with channel attention seamlessly and efficiently and significantly boosts the segmentation performance with only minor computation overhead of the original axial attention.
- We develop a novel approach to visualize the impact of the conflicting pixel representation of the existing dual attention designs on segmentation.
- To balance the computation speed and GPU memory usage, we propose a novel grouped vectorization approach to compute the channelized attentions, which is particularly advantageous when processing large images with a limited GPU memory.
- Extensive experiments on three challenging benchmark datasets, *i.e.*, PASCAL Context [13], COCO-Stuff [14] and Cityscapes [15], demonstrate the superiority of our approach over the state-of-the-art approaches.

Next, Sect. II briefly summarizes the related work. Then, we illustrate the details of our proposed approach in Sect. III. Sect. IV presents the experiments and ablation studies. The paper concludes in Sect. V.

II. RELATED WORK

Towards using the attention mechanisms to improve the performance of semantic segmentation, many research works have been reported. In this section, we introduce these approaches in the way of their evolution.

A. Capturing Information from Fixed Ranges

The PSPNet [8] proposed a PPM, which used multiple average pooling layers with different sizes together to get average pixel representations in multiple receptive fields, and then upsampled and concatenated them together. Similarly, the ASPP in DeepLab [2], [9] used parallel atrous convolutions with different rates to capture information from multiple

ranges. The core ideas of both models are to utilize the surrounding information of each pixel in multiple ranges to achieve better pixel representations. Both methods have achieved highest scores in some popular public datasets [13], [15]. However, as claimed in [5], fixed receptive fields may lose important information, to which stacking more receptive fields can be a solution, at the cost of dramatically increased computation.

B. Attention Mechanisms

Spatial Self-Attention. Non-Local networks [11] introduced the self-attention mechanism to examine the pixel relationship in spatial domain. It usually calculates dot-product similarity or cosine similarity to obtain the similarity measurement between every two pixels in feature maps, and recalculate the feature representation of each pixel according to its similarity with others. Spatial self-attention has successfully addressed the feature map coverage issue of multiple fixed-range approaches [2], [8], [5], but it has also introduced huge computation cost for computing the full feature map. This means, for each pixel in the feature maps, its attention similarity concerns all other pixels. Recently, many approaches [16], [17], [18], [19] have been developed to optimize the spatial self-attention. They have not only reduced computation and GPU memory costs but also improved the performance.

Channel Attention. Channel attention [10] examined the relationships between channels, and enhanced the important channels so as to improve the performance. SENets [10] conducted a global average pooling to get mean feature representations, and then went through two fully connected layers, where the first one had reduced channels and the second one recovered the original channels, resulting in channel-wise weights according to the importance of channels. In DANet [1], channel-wise relationships were modelled by a 2D attention matrix, similar to the spatial self-attention mechanism except that it computed the attention with a dimension of $C \times C$ rather than $H \times W \times H \times W$ (C denotes the number of channels, and H and W denote the height and width of the feature maps, respectively).

C. Spatial Attention + Channel Attention

Combining spatial attention and channel attention can provide fully optimized pixel representations in a feature map. However, it is not easy to enjoy both advantages seamlessly. In the DANet [1], the results of the channel attention and spatial attention are directly added together. Supposing that there is a pixel belonging to a semantic class that has a tiny region in the feature maps, spatial-attention can find its similar pixels. However, channel representation of the semantic class with a partial region of the feature maps may not be important in the perspective of entire feature maps, so it may be ignored when conducting channel attention computation. Computing self-attention and channel attention separately (as illustrated in Fig. 1(a)) can cause conflicting results, and thus weaken their performance when both results are summarized together. In the cascaded model (see Fig. 1(b)), the spatial attention module after the channel attention module may pick up the incorrect

pixel representation enhanced by channel attention, as channel attention computes the channel importance according to the entire feature maps.

In our work, we propose a Channelized Axial Attention approach, which first computes the spatial attention row-by-row and column-by-column, and then inserts the channel attention module to integrate both approaches seamlessly, as detailed next.

III. METHODS

A. Formulation of the Spatial Self-Attention

Following [11], [20], a 2D self-attention operation in spatial domain of neural networks can be defined by:

$$\mathbf{y}_{i,j} = \sum_{\forall m,n} f(\mathbf{x}_{i,j}, \mathbf{x}_{m,n}) g(\mathbf{x}_{m,n}). \quad (1)$$

Here, a pairwise function f computes the similarity between the pixel representations $\mathbf{x}_{i,j}$, $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$, at the position (i,j) and the pixel representation $\mathbf{x}_{m,n}$ at all other possible positions (m,n) , producing a spatial attention map over the whole feature maps. The unary function g maps the original representation at position (m,n) to a new domain. In our work, we also take the softmax function as f , *i.e.*,

$$f(\mathbf{x}_{i,j}, \mathbf{x}_{m,n}) = \text{softmax}_{m,n}(\theta(\mathbf{x}_{i,j})^T \theta(\mathbf{x}_{m,n})). \quad (2)$$

Thus, given a feature map output from a backbone network such as ResNet [21], the self-attention operation firstly uses a 1×1 convolution θ to map the feature maps \mathbf{x} to a new domain, and then applies dot-product similarity [11] between every two pixels. Then, using this similarity as the weight, Eq. (1) calculates a weighted global sum over all pixels in the feature maps and outputs a new pixel representation $\mathbf{y}_{i,j}$ at the position (i,j) .

It can be seen from Eq. (2) that, the original similarity map is of $H \times W \times H \times W$ dimensions, and is computed as the dot product over the whole feature maps for each pixel.

Axial Attention, proposed in [19], [22] for NLP and Panoptic Segmentation, has a computation complexity of $O(HW^2 + H^2W)$, smaller than the self attention's $O(H^2W^2)$ because its attention is computed within the same column or row only for each pixel. However, it has not yet had a baseline in semantic segmentation. In this work, in order to take the computation complexity advantage of the axial attention, we redefine the axial attention in [19], [22] and convert it to a specialized semantic segmentation model. In the next section, we first formulate the axial attention and then introduce our proposed channelized axial attention.

B. Formulation of the Spatial Axial Attention

In axial attention, the spatial attention map is calculated along the column axis and row axis, respectively. For the convenience of reference, we call the partial attention map calculated along the Y axis as 'column attention' and 'row attention' for the partial attention map calculated along the X axis. For the j -th column attention, the attention similarity tensor is calculated by the similarity between the current position (i,j) and each of the other positions (m,j) in the j -th

column (instead of all other positions, as in the self-attention), *i.e.*,

$$A_{\text{col}}(\mathbf{x}_{i,j}, \mathbf{x}_{m,j}) = \text{softmax}_m \left(\theta(\mathbf{x}_{i,j})^T \theta(\mathbf{x}_{m,j}) \right), \quad j \in [W].^1 \quad (3)$$

Here, θ represents the learned feature extraction process for the Y axis. Each $A_{\text{col}}(\mathbf{x}_{i,j}, \mathbf{x}_{m,j})$ represents the similarity between $\mathbf{x}_{i,j}$ and $\mathbf{x}_{m,j}$ for $i, m \in [H]$, so each $\mathbf{x}_{i,j}$ corresponds to H column-attention maps $A_{\text{col}}(\mathbf{x}_{i,j}, \mathbf{x}_{m,j})$. Thus, the resultant column attention map A_{col} is a tensor of $W \times H \times H$ dimensions.

Similarly, for the i -th row attention, the similarity attention tensor calculates the similarity between the current position (i,j) and other positions (i,n) in the i -th row, *i.e.*,

$$A_{\text{row}}(\mathbf{x}_{i,j}, \mathbf{x}_{i,n}) = \text{softmax}_n \left(\phi(\mathbf{x}_{i,j})^T \phi(\mathbf{x}_{i,n}) \right), \quad i \in [H], \quad (4)$$

where ϕ represents the learned feature extraction process for the X axis. Similarly, each $\mathbf{x}_{i,j}$ corresponds to W row-attention maps $A_{\text{row}}(\mathbf{x}_{i,j}, \mathbf{x}_{i,n})$. Thus, the resultant row attention map A_{row} is a tensor of $H \times W \times W$ dimensions.

It is also worth of pointing out that, in Eqs. (3) and (4), the calculations of column and row attention maps both use the same feature $\mathbf{x}_{i,j}$ extracted from the backbone module as the input, as shown in Fig. 2. This is different from [22], where the row attention map was computed based on the result of the column attention. By using the same feature as the input, the dependency of the final output $\mathbf{y}_{i,j}$ on the feature $\mathbf{x}_{i,j}$ has been enhanced effectively, instead of using skip connections, as in [22].

With the column and row attention maps A_{col} and A_{row} , the final value weighted by the column and row attention maps can be represented as:

$$\mathbf{y}_{i,j} = \sum_{\forall n} \left(A_{\text{row}}(\mathbf{x}_{i,j}, \mathbf{x}_{i,n}) \left(\sum_{\forall m} A_{\text{col}}(\mathbf{x}_{i,j}, \mathbf{x}_{m,j}) g(\mathbf{x}_{m,n}) \right) \right) \quad (5)$$

For the convenience of illustration, we introduce two variables $\alpha_{i,j,m}$ and $\beta_{i,j,n}$ to capture the intermediate, weighted features, respectively, where

$$\alpha_{i,j,m} = A_{\text{col}}(\mathbf{x}_{i,j}, \mathbf{x}_{m,j}) g(\mathbf{x}_{m,j}) \quad (6)$$

and

$$\beta_{i,j,n} = A_{\text{row}}(\mathbf{x}_{i,j}, \mathbf{x}_{i,n}) \sum_{\forall m} \alpha_{i,j,m}. \quad (7)$$

As illustrated later in Sect. III-C, capturing the intermediate attention results brings opportunity to conduct independent channel attentions for each partial attention result.

Thus, Eq. (5) can be simplified as:

$$\mathbf{y}_{i,j} = \sum_{\forall n} \beta_{i,j,n} = \sum_{\forall n} A_{\text{row}}(\mathbf{x}_{i,j}, \mathbf{x}_{i,n}) \left(\sum_{\forall m} \alpha_{i,j,m} \right). \quad (8)$$

The above Eqs. (6), (7) and (8) show that, the computation of the dot product is composed of two steps: 1) The element-wise multiplication for applying the column attention as shown in Eq. (6) and for applying the row attention as shown in Eq. (7) for column and row attentions, respectively; 2) The summarization of the elements along each row and column according to Eq. (8).

¹We use $i \in [n]$ to denote that i is generated from $[n] = \{1, 2, \dots, n\}$.

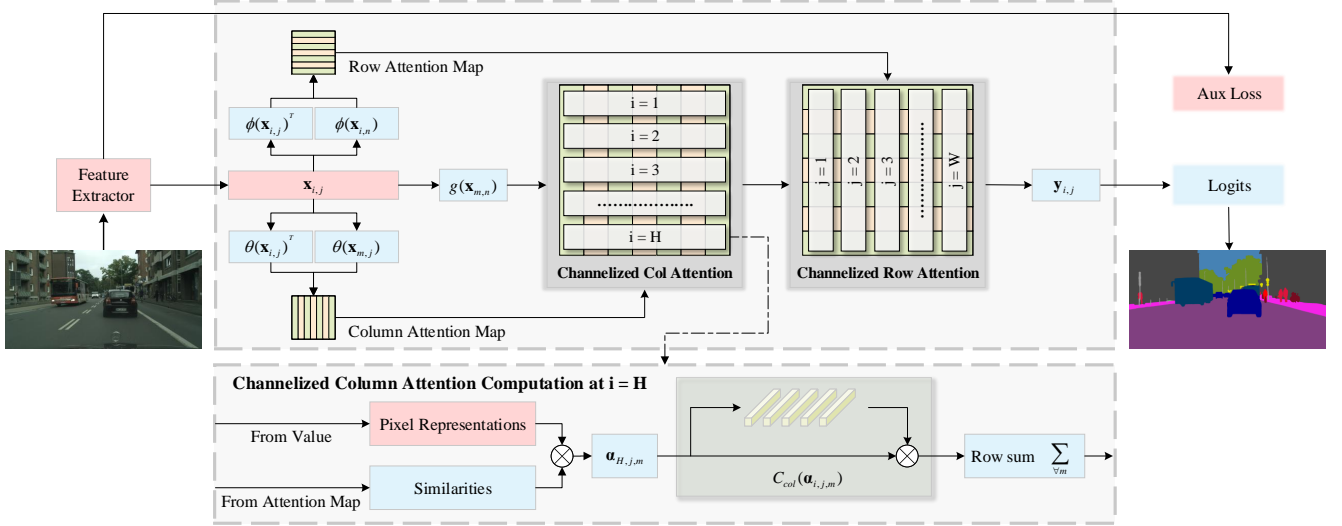


Fig. 2: The detailed architecture of our proposed Channelized Axial Attention model. To obtain $H \times W \times C$ inputs for the channel attention, we apply the resultant column and row attentions in a transposed way. The bottom section illustrates the channelization of the column attention for $i = H$.

C. The Proposed Channelized Axial Attention

In order to address the feature conflicting issue of the dual attention and seamlessly combine the advantages of spatial attention and channel attention, we propose a novel *Channelized Axial Attention*, which takes the intermediate results $\alpha_{i,j,m}$ and $\beta_{i,j,n}$ in Eqs. (6) and (7) as input.

Note that, in Eqs. (6) and (7), we apply the column and row attention maps in a transposed order. That is to say, the column and row attention results are decomposed along the transposed axis (*i.e.*, decomposing $\alpha_{i,j,m}$ along the row direction and $\beta_{i,j,n}$ along the column direction), instead of along the column and row, into multiple 3-dimension column or row attention results for different i or j . This is illustrated in Fig. 2.

This transpositional way of applying the axial attentions not only produces partial column and row attention results with consistent dimensions, but also enables them to capture the dependencies inherent in the other axis so as to conduct channelization in the whole receptive field.

Now, we introduce our channelized attentions C_{col} and C_{row} , corresponding to the column attention and row attention, respectively, as:

$$C_{\text{col}}(\alpha_{i,j,m}) = \text{Sigmod} \left(\text{ReLU} \left(\frac{\sum_{\forall m,j} (\alpha_{i,j,m})}{H \times W} \omega_{c1} \omega_{c2} \right) \right) \alpha_{i,j,m} \quad (9)$$

and

$$C_{\text{row}}(\beta_{i,j,n}) = \text{Sigmod} \left(\text{ReLU} \left(\frac{\sum_{\forall i,n} (\beta_{i,j,n})}{H \times W} \omega_{r1} \omega_{r2} \right) \right) \beta_{i,j,n} \quad (10)$$

where ω_{c1} , ω_{c2} , ω_{r1} and ω_{r2} represent the learned relationships between different channels in $\alpha_{i,j,m}$ and $\beta_{i,j,n}$, which will be discussed later in Sect. III-E.

Thus, instead of directly using $\alpha_{i,j,m}$ and $\beta_{i,j,n}$ as in Eq. (8), for each column and row, we obtain the seamlessly mixed attention results for spatial channels, where the interme-

diate results $\alpha_{i,j,m}$ and $\beta_{i,j,n}$ are weighted by the channelized axial attention defined in Eqs. (9) and (10) as:

$$\mathbf{y}_{i,j} = \sum_{\forall n} C_{\text{row}} \left(A_{\text{row}}(\mathbf{x}_{i,j}, \mathbf{x}_{i,n}) \left(\sum_{\forall m} C_{\text{col}}(\alpha_{i,j,m}) \right) \right). \quad (11)$$

The bottom section in Fig. 2 illustrates the channelization of the column attention at $i = H$. Later in Sect. IV-B (TABLE II and Fig. 5), we will show with ablation experiments and visualized feature maps the impact of the channelization on improving the performance of the segmentation.

D. Grouped Vectorization

Computing spatial attention row by row and column by column can save computation but it is still too slow even with parallelization. Vectorization can achieve a very high speed but it has a high requirement on GPU memory for storing the intermediate partial axial attention results α (which has a dimension of $H \times H \times W \times C$) and β (which has a dimension of $W \times H \times W \times C$) in Eqs. (6) and (7). To enjoy the high speed benefit of the vectorized computation with reduced GPU memory usage, in our implementation we propose *grouped* vectorization to dynamically batch rows and columns into multiple groups, and then perform vectorization for each group respectively. Algorithm 1 shows the pseudo code of implementing the grouped vectorization.

E. Going Deeper in Channel Attention

The channel attention in our method firstly uses a fully connected layer with a smaller rate to compress channels, and then uses another fully connected layer with the same rate as the original channels, followed by a sigmoid function to generate the final channel attention weights. To further boost the performance, we explore the design of more powerful channel attention modules in channelization.

Algorithm 1 Our proposed grouped vectorization algorithm

Require: G : Group Number, A : Attention Map $[N, H, H, W]$, X : Feature Map $[N, C, H, W]$

- 1: $padding \leftarrow H \% G$
- 2: $A \leftarrow \text{Transpose } A \text{ into } [H, N, H, W]$
- 3: $H^+ \leftarrow H + padding$
- 4: $A \leftarrow \text{padding zero to } A \text{ into } [H^+, N, H, W]$
- 5: $A \leftarrow \text{Reshape } A \text{ into } [G, H^+ // G, N, H, W]$
- 6: **for** $g \in G$ **do**
- 7: $Y_g \leftarrow \text{Channelization } (X, A_g), Y_g \in [H^+ // G, N, C, W]$
- 8: **end for**
- 9: $Y \leftarrow \text{Concat}(Y_{0,1,\dots,G}), Y \in [G, H^+ // G, N, C, W]$
- 10: $Y \leftarrow \text{Reshape } Y \text{ into } [H^+, N, C, W]$
- 11: $Y \leftarrow \text{Remove padding from } Y \text{ into } [H, N, C, W]$
- 12: $Y \leftarrow \text{Transpose } Y \text{ into } [N, C, H, W]$

return Y

The simplest way of gaining performance is enhancing the representation ability of the neural networks, and it is usually achieved by increasing the depth and width of the networks. Here, we simply add more hidden layers before the last layer. This design allows channel attention to find better relationship between channels and find more important channels for each axial attention’s intermediate results. We also find that it is not effective to increase the width (*i.e.*, adding more hidden units to each layer except for the last layer), so we keep the original settings.

Furthermore, in the spatial domain, each channel of a pixel contains unique information that can lead to unique semantic representation. In our channel attention module, we find that using Leaky ReLU [23], instead of ReLU, is more effective in preventing the loss of information along deeper activations [24]. Apparently, this replacement only works in our channel attention module.

IV. EXPERIMENTS

To demonstrate the performance of our proposed CAA, comprehensive experiments are conducted with results compared with the state-of-the-art results on three benchmark datasets, *i.e.*, PASCAL Context [13], COCO-Stuff [14] and Cityscapes [15].

The same as the other existing works [1], [2], [6], [4], we measure the segmentation accuracy using mIOU (Mean Intersection Over Union). Moreover, to demonstrate the efficiency of our CAA, we also report and compare the FLOPS (Floating Point Operations per Second) of different approaches. Note that, a higher mIOU value means more accurate segmentation, whereas a lower FLOPS value indicates less computation operations. Experimental results show that our CAA outperforms the state-of-the-art performance on all tested datasets.

Moreover, as we have analysed earlier in Sect. II-C, computing spatial attention and channel attention separately and then fusing them together directly can cause conflicting feature representations. This means, channels important for a class in spatial domain may not dominate and therefore can be ignored in the computation of the global channel attention. In our experiments, to illustrate the feature conflicting issue caused by existing dual attention approaches, we design a simple

way to visualize the effects of spatial attention and channel attention on pixel representation.

For the parallel dual attention design such as DANet [1], since it has two auxiliary losses for spatial attention and channel attention respectively, we directly use their logits during inference and generate their segmentation results to compare with the result generated by the main logits. For the sequential dual attention design, we add an extra branch that directly uses the pixel representation obtained from channel attention to perform the segmentation logits. Note that, since the original sequential design does not have independent logits of channel attention, we stop the gradient from the main branch to make sure our newly added branch has no effect on the main branch.

Next, we first present the implementation details. This is followed by a series of ablation experiments conducted on the PASCAL Context dataset showing the effectiveness of each of our proposed ideas. Then, we report the comparative results obtained on PASCAL Context [13], COCO-Stuff [14] and Cityscapes [15] datasets, respectively. For fair comparison, we only compare with the methods that use ResNet-101 and naive $8 \times$ bilinear upsampling.

A. Implementation Details

Backbone: Our network is built on ResNet-101 [21] pre-trained on ImageNet. The original ResNet results in a feature map of $1/32$ of the input size. Following other similar works [9], [22], [6], we apply dilated convolution at the output stride = 16 during training for most of the ablation experiments. We conduct experiments with the output stride = 8 during training to compare with the state of the arts.

Segmentation Head: We use a 3×3 convolution to reduce the number of feature map channels from 2,048 to 512 [1], [6], which is then followed by our proposed Channelized Axial Attention module. Note that, our Axial Attention generates the column attention map and the row attention map from the same feature maps, instead of generating one based on the computation results of the other, as in [22]. Also, after the computation of the attention maps, we do not add the original pixel representations to the resultant feature maps. In the end, we directly upsample our logits to the input size by applying bilinear interpolation.

Training Settings: We employ SGD (Stochastic Gradient Descent) for optimization, where the poly decay learning rate policy $(1 - \frac{iter}{maxiter})^{0.9}$ is applied with an initial learning rate = 0.007. We use synchronized batch normalization during training. Our experiments are conducted on TPUv3 and V100. For data argumentation, we only apply the most basic data argumentation strategies in [9] including random flip, random scale and random crop, same as in the other works.

B. Experiments on PASCAL Context Dataset

PASCAL Context [25] dataset has 59 classes with 4,998 images for training and 5,105 images for testing. We train the network model on PASCAL Context Training set with the batch size = 16 with 70k iterations. During training, we set the output stride = 16 and use an output stride = 8 for

TABLE I: Comparison results with different segmentation heads in the PASCAL Context dataset [25]

Methods	mIOU%	FLOPS
ResNet-101 [21]	-	59.85G
FCN [7]	48.12	+0G
ASPP [9]	50.47	+16.7G
Non-Local [11]	50.42	+11.18G
Redefined Axial Attention	50.27 (± 0.2)	+8.85G

inference. Later in TABLE VI, we present our CAA results with an output stride = 16 and 8, where it can be seen clearly a 1.4% increase can be observed with the output stride = 8.

Next, we first present a series of ablation experiments conducted on the PASCAL Context dataset to show the effectiveness of our proposed channelized axial attention. Then, quantitative and qualitative comparisons with the state of the arts are presented.

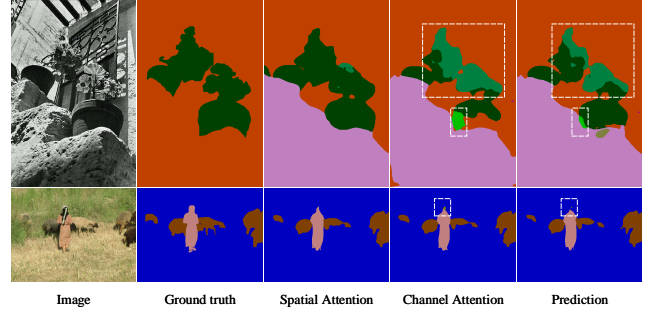
1) *Axial Attention for Semantic Segmentation*: To verify the effectiveness of Axial Attention for semantic segmentation (see Sect. III-B), we compare the mIOU and FLOPS achieved with our channelized Axial Attention with other segmentation heads implemented by us, as shown in TABLE I. Note that our redefined Axial Attention used for semantic segmentation is different from [22], as mentioned in Sect. III-A, and in this table we only compare with the methods that are independent with backbone [21]. Also, all results in this table are obtained with an output stride = 16.

From TABLE I, we can easily see that our redefined Axial Attention improves mIOU a lot compared to the Dilation-FCN (50.27 vs 48.12), which has a naive segmentation head. The mIOU obtained with our redefined axial attention is also comparable with other approaches, such as ASPP [2], [9] and Non-Local [11]. However, the redefined axial attention has much lower FLOPS than the original self-attention [11] (an increase of 8.85G vs 11.18G over the baseline), which demonstrates that the redefined axial attention for semantic segmentation can achieve comparable performance with the original self-attention at much lower computation cost.

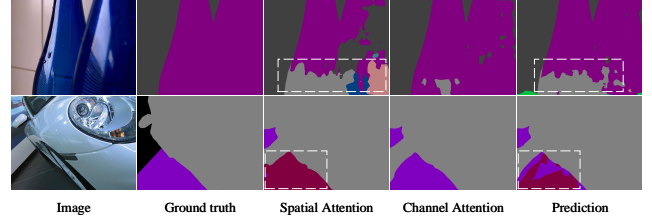
2) *Examples of Conflicting Feature Representations*: To visualize the impact of the feature conflicting issue of the existing dual attention designs (see Sect. II-C), Fig. 3 shows two groups of examples of the segmentation results obtained with the conflicting features in the parallel dual attention design (see Figs. 3a and 3b) and the sequential dual attention design (see Fig. 3c).

As it can be observed from Figs. 3a and 3b, the parallel design of dual attentions directly sums up the pixel representations obtained from spatial attention and channel attention. With this approach, the advantages of the pixel representations obtained from one can be weakened by the other.

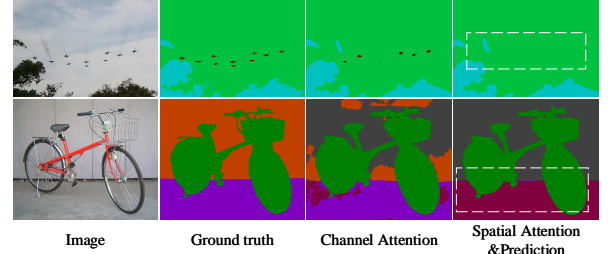
The sequential way of combining the dual attentions avoids taking their average but still has its issue. As shown in Fig. 3, the pixel representation obtained from the spatial attention abandons the correct pixel representation obtained from the channel attention, and worsens the prediction result.



(a) The bad channel attention representation negatively influences the good spatial attention representation. See the highlighted areas.



(b) The bad spatial attention representation negatively influences the good channel attention representation. See the highlighted areas.



(c) The spatial attention representation abandons the correct channel attention representation.

Fig. 3: Examples of conflicting feature representation in the parallel (a and b) and sequential dual attention (c) designs.

TABLE II: Result comparison without using channelization (Row 1) and using channelization with different layer counts and channel numbers.

Layer Counts				# of Channels			mIOU%	FLOPS
1	3	5	7	64	128	256		
-	-	-	-	-	-	-	50.27(± 0.2)	68.7G
✓					✓		50.75(± 0.2)	+0.00024G
	✓				✓		50.85(± 0.2)	+0.00027G
		✓			✓		51.06(± 0.2)	+0.00030G
			✓		✓		50.40(± 0.3)	+0.00043G
				✓			50.12(± 0.2)	+0.00015G
						✓	50.35(± 0.4)	+0.00098G

3) *Effectiveness of the Proposed Channelization*: We then use our proposed channelized dot product to replace the naive dot product in Axial Attention (see Sect. III-C). We report the impact of adding channelized dot product and with different depth and width in TABLE II, where ‘-’ for the baseline result indicates no channelization is performed.

As it can be seen from this table, our proposed channel-

TABLE III: Result comparison between axial attention, channelized axial attention and axial attention + SE [10]

Axial Attention	+ Channelization	+ SE
50.27(± 0.2)	51.06(± 0.2)	50.37(± 0.2)

TABLE IV: Ablation study of applying our Channelized Attention on self-attention with ResNet-101 [21]. **Eval OS:** Output strides [9] during evaluation.

Attention Base	Eval OS	Channelized	mIOU%
Axial Attention	16		50.27
	16	✓	51.06
Self Attention	16		50.42
	16	✓	51.09

ization improves the mIOU performance over the baseline regardless of the layer counts and the number of channels used. In particular, a best performance is achieved when the Layer Counts = 5 and the number of Channels = 128.

We also compare our model with the sequential design of Axial Attention + SE, as shown in TABLE III. We repeated the experiments many times but found the sequential design only brings slightly contributions on the performance, indicating that our purposed channelization method can combine the advantages of both spatial attention and channel attention effectively. Also note that, Fig. 1(b) shows a failure result of the sequential design.

4) *Channelized Self-Attention:* In this section, we conduct additional experiments on the PASCAL Context testing set by applying channelization to the original self-attention. We report its single-scale performance in TABLE IV with ResNet-101 [21].

We can see from the table that our proposed channelized method can further improve the performance of self-attention slightly by 0.67%. It also shows the current channelized design is more effective for our Axial Attention (0.79% vs 0.67%).

5) *Impact of the Testing Strategies:* We report and compare the performance and computation cost of our proposed model against the baseline and the DANet with different testing strategies. This is shown in TABLE V. Same as the settings in other works [8], [1], we add multi-scale, left-right flip and aux loss [8], [1] during inference. Note that, in this table, we report the mean mIOU figures with a dynamic range to show the stability of our algorithm. As it shows in this table, We found our proposed CAA can be further boosted with OS = 8 since the channel attention can learn and optimize three times more pixels.

6) *Comparison with the State of the Arts:* Finally, we compare our proposed approach with the state-of-the-art approaches. The results on the PASCAL Context dataset is shown in TABLE VI. Like other similar works, we apply multi-scale and left-right flip during inference. For fair comparison, we only compare with the methods that use ResNet-101 and naive decoder (directly upsampling logits). Also note that, in this and

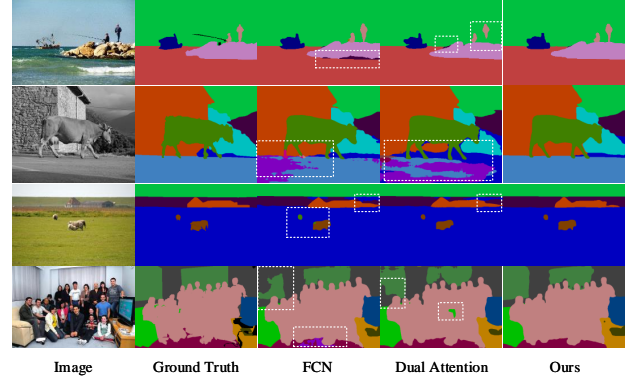


Fig. 4: Examples of the segmentation results obtained on the PASCAL Context dataset [25] with our proposed CAA approach in comparison to the results obtained with FCN [7], DANet [1] and the ground truth. All results are inferred with an output stride of 8.

the following tables, we report the best results of our approach obtained in experiments.

As shown in this table, our proposed CAA method achieves the highest score in the methods trained with an output stride = 16 with ResNet-101 and naive decoder, and even outperforms some methods trained with an output stride = 8. Moreover, after we train our model with an output stride = 8, the performance of our model has been further improved and outperforms all of the state-of-the-art models, including the ones recently published in CVPR2019 and CVPR2020.

In Fig. 4, we provide the visualizations of the prediction results obtained with our CAA model in comparison with the state-of-the-art approaches. As shown in the figure, our model is able to segment objects very well without requiring any post-processing.

To further demonstrate the effectiveness of our proposed channelization, in Fig. 5 we visualize the feature maps obtained after applying the column attention and row attention maps and the difference between the corresponding feature maps with and without applying the channel attentions.

7) *Alternative Backbones:* In previous sections, we have reported our CAA's performance using ResNet-101 [21] as backbone, which is widely used in semantic segmentation [9], [1], [17], [6], [5], [4], [26], [30], [8], [31], [18]. In this section, we conduct additional experiments on Pascal Context by attaching our CAA module with some other backbones. We report our results obtained with single scale without flipping in TABLE VII.

8) *Result with EfficientNet:* As mentioned in Sect. IV-B6, our CAA outperforms the SOTA methods [30], [6] with the same settings (ResNet-101 w/o decoder). Furthermore, TABLE VII shows the universality of our proposed CAA with different backbones. In this section, we report our CAA's performance with EfficientNet-B7 [33] in TABLE VIII. Note that, this is not a fair comparison, since the listed methods

TABLE V: Comparison results with different testing strategies. **Train OS**: Output stride in training. **Eval OS**: Output stride in inference. **MS**: Apply multi-scale during inference. **Aux loss**: Add auxiliary loss during training. “+” refers to the FLOPs over the baseline FLOPs of ResNet-101.

Methods	Train OS 16	8	Eval OS 16	8	Strategies MS flip	Aux Loss	mIOU%	FLOPs
ResNet-101 [21]	-	-	✓	✓		-	-	59.85G 190.70G
DANet [1]		✓		✓	✓	✓	- 52.60	+101.25G -
Our CAA	✓		✓		✓		51.06(±0.2) 53.09(±0.3)	+8.85G -
Our CAA + Aux loss	✓		✓		✓	✓	51.80(±0.2) 53.52(±0.2)	+8.85G -
		✓		✓	✓	✓	53.48(±0.3)	+34.33G
		✓		✓	✓	✓	54.65(±0.4)	-

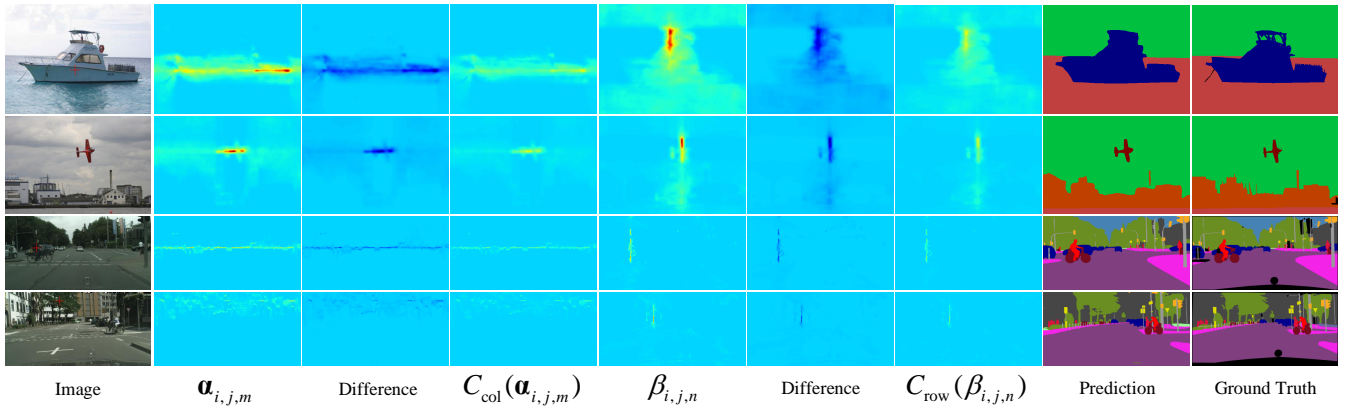


Fig. 5: Visualization of the feature maps ($\alpha_{i,j,m}$, $C_{\text{col}}(\alpha_{i,j,m})$, $\beta_{i,j,n}$ and $C_{\text{row}}(\beta_{i,j,n})$) on PASCAL Context [25] (top two rows) and Cityscapes [15] (bottom two rows). For each input image, we list the feature maps obtained after applying the column attention map and row attention map, the difference between the corresponding feature maps with and without applying the channel attentions, as well as our prediction and the ground truth segmentation, respectively. For more details, please refer to Sect. 2.

TABLE VI: Result comparison with the state-of-the-art approaches on the PASCAL Context testing set for multi-scale prediction. For fair comparison, we only compare with the methods that use ResNet-101 and naive decoder.

Methods	mIOU%	Ref
FCN [7]	50.8	CVPR2015
ENet [26]	51.7	CVPR2018
Deeplab [9]	52.7	ECCV2018
ANet [18]	52.8	ICCV2019
EMANet [6]	53.1	ICCV2019
SVCNet [27]	53.2	CVPR2019
SPYGR [28]	52.8	CVPR2020
CPN [29]	53.9	CVPR2020
CFNet [30]	54.0	CVPR2019
DANet [1]	52.6	CVPR2019
Our CAA (OS = 16)	53.7	-
Our CAA (OS = 8)	55.0	-

were not trained under the same settings, or using the same backbone. The results show that our method can still

improve the performance even with a strong CNN backbone Efficientnet-B7, and outperform the latest Transformer [34] based hybrid models such as SETR [35] and DPT [36].

C. Results on the COCO-Stuff 10K Dataset

1) *Comparison with the State of the Arts*: Following the other works [6], [4], [1], we demonstrate that our model can handle complex images with a large number of classes. We further evaluate our model on the COCO-Stuff 10K dataset [14], which contains 9,000 training images and 1,000 testing images, as shown in TABLE IX. As it can be seen from the table, our proposed CAA outperforms all other state-of-the-art approaches by a large margin of 1.3%.

We also report results obtained with our CAA with Efficientnet-b7 [33] in TABLE X.

2) *Visualization of the Segmentation Results*: Fig. 6 show some examples of the segmentation results obtained on the COCO-Stuff 10K dataset [14] with our proposed CAA in comparison to the results of FCN [7], DANet [1] and the

TABLE VII: Ablation study of applying our Channelized Axial Attention to other backbones. All results are obtained in single scale without flipping. **Axial Attention**: Using our Axial Attention after backbone. **Channelized**: Applying our Channelized approach. **Eval OS**: Output strides [9] during evaluation.

Backbone	Eval OS	Axial Attention	Channelized	mIOU%
ResNet-50 [21]	16			46.92
	16	✓		49.73
	16	✓	✓	50.23
ResNet-101 [21]	16			48.12
	16	✓		50.27
	16	✓	✓	51.06
Xception65 [32], [9]	16			49.40
	16	✓		52.42
	16	✓	✓	52.65
EfficientNetB7 [33]	16			56.80
	16	✓		57.24
	16	✓	✓	57.93
	8	✓	✓	58.40

TABLE VIII: Result comparison with the state-of-the-art approaches on the PASCAL Context testing set for multi-scale prediction. Note that, this is not a fair comparison, since all listed methods were not trained under same settings, or using same backbone.

Methods	mIOU%
SETR-MLA [35]	55.83
HRNetV2 + OCR [27]	56.2
ResNeSt-269 [37] + DeepLab V3+ [9]	58.9
HRNetV2 + OCR + RMI [4]	59.6
DPT-Hybrid [36]	60.46
Our CAA (EfficientNet-B7, w/o decoder)	60.12
Our CAA (EfficientNet-B7 + simple decoder [9])	60.50

ground truth. All results are inferred with an output stride of 8. As it can be seen, our CAA can segment common objects such as building, human, or sea very well.

D. Results on the Cityscapes Dataset

The Cityscapes dataset [15] has 19 classes. Its *fine* set contains high quality pixel-level annotations of 5,000 images, where there are 2,975, 500 and 1,525 images in the Training, Validation, and Test sets, respectively. Like other works [28], [1], We only use *fine* set with a crop size 769×769 during training, and our training iteration is set to 90k. We report our results on *test* set in TABLE XI and also visualize our feature maps and results in Fig. 5 (the bottom two rows).

E. Effectiveness of Our Grouped Vectorization

In Sect. III-D, we developed the grouped vectorization to split tensors into multiple groups so as to reduce the GPU memory usage when performing channel attention in Eqs. (9) and (10). The more groups used in group vectorization, the proportionally less GPU memory is needed for the computation, yet with longer inference time. In this section, we conduct experiments to show the variation of the inference

TABLE IX: Comparison results with the state-of-the-art approaches on the COCO-Stuff 10K testing set for multi-scale prediction. For fair comparison, we only compare with the methods that use ResNet-101 and naive decoder.

Methods	mIOU%	Ref
DSSPN [38]	38.9	CVPR2018
SVCNet [27]	39.6	CVPR2019
EMANet [6]	39.9	ICCV2019
SPYGR [28]	39.9	CVPR2020
OCR [4]	39.5	ECCV2020
DANet [1]	39.7	CVPR2019
Our CAA	41.2	-

TABLE X: Result comparison with the state-of-the-art approaches on the COCO-Stuff-10K testing set for multi-scale prediction. Note that, this is not a fair comparison, since all listed methods were not trained under same settings, or using same backbone.

Methods	mIOU%
HRNetV2 + OCR [27]	40.5
DRAN	41.2
HRNetV2 + OCR + RMI [4]	45.2
Our CAA (EfficientNet-B7)	45.4

time (seconds/image) when different numbers of groups are used in group vectorization.

Fig. 7 shows the results where three different input resolutions are tested. As shown in this graph, when splitting the vectorization into smaller numbers of groups, *e.g.*, 2 or 4, our grouped vectorization can achieve comparable inference speed with one half or one quarter of the original spacial complexity.

V. CONCLUSION

In this paper, aiming to combine the advantages of the popular spatial-attention and channel attention, we have proposed a novel and effective Channelized Axial Attention approach for semantic segmentation. After computing column and row attentions, we proposed to capture the intermediate results and perform the corresponding channel attention on each of them. Our proposed approach of applying the column and row attentions transpositionally has allowed the channelization to be conducted in the whole respective field. Experiments on the three popular benchmark datasets have demonstrated the superiority and effectiveness of our proposed channelized axial attention in terms of both segmentation performance and computational complexity.

TABLE VII shows that both our Axial Attention and Channelization approaches have improved the mIOU of the baseline in multiple well-know backbones. We also find that our Channelization approach is more effective with ResNet and EfficientNet, whereas the improvement on Xception65 is relatively small.

REFERENCES

- [1] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, and H. Lu, “Dual attention network for scene segmentation,” in *Conference on Computer Vision and Pattern Recognition*, 2019.



Fig. 6: Examples of the segmentation results obtained on the COCO-Stuff 10K dataset [14] with our proposed CAA approach in comparison to the results obtained with FCN [7], DANet [1] and the ground truth. All results are inferred with an output stride of 8.

TABLE XI: Comparison results with other state-of-the-art approaches on the Cityscapes Test set for multi-scale prediction. For fair comparison, we only compare with the methods that use ResNet-101 and naive decoder.

Methods	mIOU%	Ref
PSPNet [8]	78.4	CVPR2017
CFNet [30]	79.6	CVPR2019
ANNet [18]	81.3	ICCV2019
CCNet [17]	81.4	ICCV2019
CPN [29]	81.3	CVPR2020
SPYGR [28]	81.6	CVPR2020
OCR [4]	81.8	ECCV2020
DANet [1]	81.5	CVPR2019
Our CAA	82.6	-

- [2] 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 Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [3] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, “Rethinking atrous convolution for semantic image segmentation,” 2017.
- [4] Y. Yuan, X. Chen, and J. Wang, “Object-contextual representations for semantic segmentation,” in *European Conference on Computer Vision*, 2020.
- [5] M. Yang, K. Yu, C. Zhang, Z. Li, and K. Yang, “Denseaspp for semantic segmentation in street scenes,” in *Conference on Computer Vision and Pattern Recognition*, 2018.
- [6] X. Li, Z. Zhong, J. Wu, Y. Yang, Z. Lin, and H. Liu, “Expectation-maximization attention networks for semantic segmentation,” in *International Conference on Computer Vision*, 2019.
- [7] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.

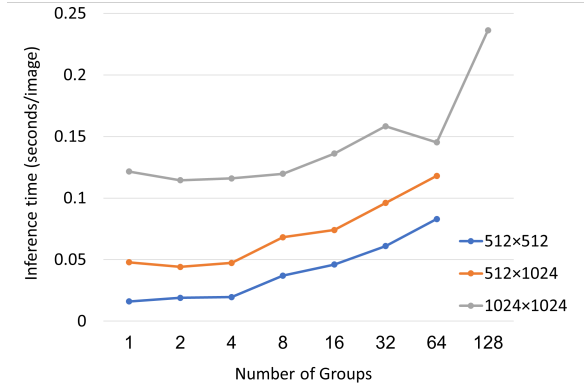


Fig. 7: Inference time (seconds/image) when applying different numbers of groups in grouped vectorization.

- [8] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” in *Conference on Computer Vision and Pattern Recognition*, 2017.
- [9] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, “Encoder-decoder with atrous separable convolution for semantic image segmentation,” in *European Conference on Computer Vision*, 2018.
- [10] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Conference on Computer Vision and Pattern Recognition*, 2018.
- [11] X. Wang, R. Girshick, A. Gupta, and K. He, “Non-local neural networks,” in *Conference on Computer Vision and Pattern Recognition*, 2018.
- [12] S. Woo, J. Park, J.-Y. Lee, and S. Kweon, “Convolutional block attention module,” in *European Conference on Computer Vision*, 2018.
- [13] M. Everingham, L. V. Gool, C. K.I. Williams, J. Winn, and A. Zisserman, “The pascal visual object classes (voc) challenge,” *International Journal of Computer Vision*, 2009.
- [14] H. Caesar, J. Uijlings, and V. Ferrari, “Coco-stuff: Thing and stuff classes in context,” in *Conference on Computer Vision and Pattern Recognition*, 2018.
- [15] C. Marius, O. Mohamed, R. Sebastian, R. Timo, E. Markus, B. Rodrigo, F. Uwe, S. Roth, and S. Bernt, “The cityscapes dataset for semantic urban scene understanding,” in *Conference on Computer Vision and Pattern Recognition*, 2016.
- [16] Y. Chen, Y. Kalantidis, J. Li, S. Yan, and J. Feng, “A2-nets: Double attention networks,” in *Conference on Neural Information Processing Systems*, 2018.
- [17] Z. Huang, X. Wang, Y. Wei, L. Huang, H. Shi, W. Liu, and T. S. Huang, “Cnet: Criss-cross attention for semantic segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [18] Z. Zhu, M. Xu, S. Bai, T. Huang, and X. Bai, “Asymmetric non-local neural networks for semantic segmentation,” in *International Conference on Computer Vision*, 2019.
- [19] J. Ho, N. Kalchbrenner, D. Weissenborn, and T. Salimans, “Axial attention in multidimensional transformers,” 2019.
- [20] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens, “Stand-alone self-attention in vision models,” in *Conference on Neural Information Processing Systems*, 2019.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Conference on Computer Vision and Pattern Recognition*, 2016.
- [22] H. Wang, Y. Zhu, B. Green, H. Adam, A. Yuille, and L.-C. Chen, “Axial-deeplab: Stand-alone axial-attention for panoptic segmentation,” in *European Conference on Computer Vision*, 2020.
- [23] A. L. Mass, A. Y. Hannun, and A. Y. Ng, “Rectifier nonlinearities improve neural network acoustic models,” in *International Conference on Machine Learning*, 2013.
- [24] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, “Mobilenetv2: Inverted residuals and linear bottlenecks,” in *Conference on Computer Vision and Pattern Recognition*, 2018.
- [25] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille, “The role of context for object detection and semantic segmentation in the wild,” in *Conference on Computer Vision and Pattern Recognition*, 2014.

- [26] H. Zhang, K. Dana, J. Shi, Z. Zhang, X. Wang, A. Tyagi, and A. Agrawal, "Context encoding for semantic segmentation," in *Conference on Computer Vision and Pattern Recognition*, 2018.
- [27] H. Ding, X. Jiang, B. Shuai, A. Q. Liu, and G. Wang, "Semantic correlation promoted shape-variant context for segmentation," in *Conference on Computer Vision and Pattern Recognition*, 2019.
- [28] X. Li, Y. Yang, Q. Zhao, T. Shen, Z. Lin, and H. Liu, "Spatial pyramid based graph reasoning for semantic segmentation," in *Conference on Computer Vision and Pattern Recognition*, 2020.
- [29] C. Yu, J. Wang, C. Gao, G. Yu, C. Shen, and N. Sang, "Context prior for scene segmentation," in *Conference on Computer Vision and Pattern Recognition*, 2020.
- [30] H. Zhang, H. Zhan, C. Wang, and J. Xie, "Semantic correlation promoted shape-variant context for segmentation," in *Conference on Computer Vision and Pattern Recognition*, 2019.
- [31] Z. Zhong, Z. Q. Lin, R. Bidart, X. Hu, I. B. Daya, and Z. Li, "Squeeze-and-attention networks for semantic segmentation," in *Conference on Computer Vision and Pattern Recognition*, 2020.
- [32] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Conference on Computer Vision and Pattern Recognition*, 2017.
- [33] T. Mingxing and L. Quoc, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*, 2019.
- [34] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2021.
- [35] Z. Sixiao, L. Jiachen, Z. Hengshuang, Z. Xiatian, L. Zekun, W. Yabiao, F. Yanwei, F. Jianfeng, X. Tao, T. P. H.S., and Z. Li, "Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers," in *Conference on Computer Vision and Pattern Recognition*, 2021.
- [36] R. Ranftl, A. Bochkovskiy, and V. Koltun, "Vision transformers for dense prediction," *ArXiv preprint*, 2021.
- [37] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, Z. Zhang, H. Lin, Y. Sun, T. He, J. Muller, R. Manmatha, M. Li, and A. Smola, "Resnest: Split-attention networks," *arXiv preprint arXiv:2004.08955*, 2020.
- [38] X. Liang, H. Zhou, and E. Xing, "Dynamic-structured semantic propagation network," in *Conference on Computer Vision and Pattern Recognition*, 2018.