

# SATS: Self-Attention Transfer for Continual Semantic Segmentation

Yiqiao Qiu<sup>1</sup>, Yixing Shen<sup>1</sup>, Zhuohao Sun<sup>1</sup>, Yanchong Zheng<sup>3</sup>

Xiaobin Chang<sup>2</sup>, Weishi Zheng<sup>1</sup>, Ruixuan Wang<sup>1</sup>,<sup>✉</sup>

1.School of Computer Science and Engineering, Sun Yat-sen University, Guangzhou, Guangdong, China

2.School of Artificial Intelligence, Sun Yat-sen University, Zhuhai, Guangdong, China

3.School of Computer Science, Wuhan University, Wuhan, Hubei, China

## Abstract

Continually learning to segment more and more types of image regions is a desired capability for many intelligent systems. However, such continual semantic segmentation exhibits catastrophic forgetting issues similar to those of continual classification learning. Unlike the existing knowledge distillation strategies for alleviating this problem, transferring a new type of information, namely, the relationships between elements (e.g., pixels) within each image that can capture both within-class and between-class knowledge, is proposed in this study. Such information can be effectively obtained from self-attention maps in a Transformer-style segmentation model. Considering that pixels belonging to the same class in each image typically share similar visual properties, a class-specific region pooling operator is novelly applied to provide reliable relationship information for knowledge transfer. Extensive evaluations on multiple public benchmarks reveal that the proposed self-attention transfer method can effectively alleviate the catastrophic forgetting issue. Furthermore, flexible combinations of the proposed method with widely adopted strategies considerably outperform state-of-the-art solutions.

## 1. Introduction

Continually learning knowledge is a desired capability for intelligent systems in many application scenarios, such as autonomous driving, autonomous stores, and intelligent healthcare. Most studies on continual learning have focused on classification tasks [1], [2], [3], [4] in which the classifier is continually updated to learn to recognize more and more classes over multiple stages of continual learning. In addition to continual classification tasks, continual semantic segmentation of images have been investigated [5], [6], [7], [8] to continually update a segmentation model such that it can learn to segment more and more types of image region. Similar to continual classification, continual semantic segmentation also suffers from the catastrophic forgetting issue [9], that is, the model rapidly forgets the knowledge obtained from previously learned old classes after learning to segment more classes of image regions, particularly over multiple stages of continual learning.

To alleviate the catastrophic forgetting issue, several strategies originally used for continual classification tasks have been directly adopted and empirically proved to be effective for continual semantic segmentation. One strategy is knowledge distillation, which is used to transfer old knowledge by distilling

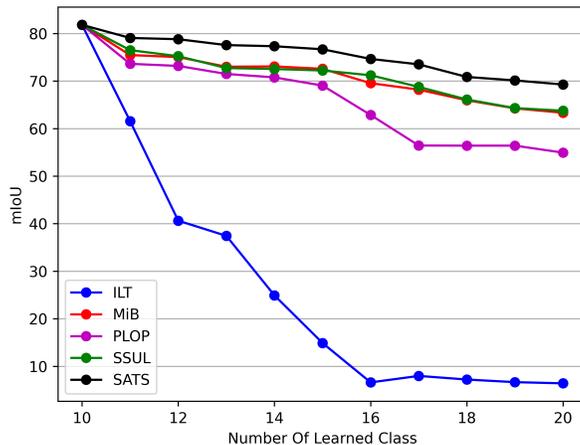


Figure 1: Performance of the proposed SATS (Self-Attention Transfer for continual semantic Segmentation) method and current state-of-the-art methods in continual semantic segmentation. Each model with the same backbone initially learns to segment 10 classes of regions and then continually learns to segment one more class at each new stage of continual learning.

the output of the last layer or multiple layers from the previously updated old model to the current model [7], [10]. Another strategy is the use of stored exemplars for each previously learned class when training the new model, allowing the new model to directly refresh old knowledge based on these limited old data. In addition to directly borrowing ideas from continual classification tasks, researchers have identified a specific background-shifting issue in continual semantic segmentation tasks, that is, background regions in images at one learning stage may contain regions of classes learned at another stage, thus confusing the model when discriminating foreground classes against the background class. Several effective remedies, including pseudo-labeling of the background regions in new training images based on the previously learned old model [7], [11], [12] and excluding salient regions from the background as the unknown (future) class during model learning [12], have been proposed to mitigate this problem.

In this study, we propose a simple and novel knowledge distillation strategy called Self-Attention Transfer for effective continual semantic Segmentation (SATS). Unlike existing knowledge distillation methods in which the model output or visual features are distilled from single or multiple convolutional layers, the proposed SATS *distills the visual relationships* between elements (e.g., pixels) within each image to capture both within-class and between-class knowledge. Such relational knowledge is obtained from self-attention maps in a transformer-style segmentation model, and it has never been used for continual semantic segmentation in previous studies. Furthermore, class-specific region pooling (CRP)

is novelly applied to efficiently elucidate the visual relationship for both within-class and between-class knowledge. To the best of our knowledge, this is the first study which applies self-attention and class-region pooling to continual learning. The proposed SATS can be combined with widely adopted continual learning strategies and has achieved state-of-the-art continual segmentation performance on the VOC [13] and ADE [14] semantic segmentation benchmark datasets. The contributions of this study are as follows:

- A novel knowledge distillation method is proposed for continual semantic segmentation. The proposed distillation of relational knowledge is complementary to existing distillation of visual knowledge.
- CRP is applied to continual semantic segmentation. In particular, CRP is used to extract within-class and between-class knowledge for distillation during continual learning.
- This is the first study in which Transformer is innovatively applied in continual semantic segmentation.
- Extensive evaluations on multiple benchmarks and settings reveal that the proposed SATS method can effectively alleviate the catastrophic forgetting issue, and its flexible combinations with widely adopted strategies outperform state-of-the-art methods.

## **2. Related Work**

### **2.1. Semantic image segmentation**

Semantic image segmentation task aims to automatically divide an image into multiple local regions such that each local region corresponds to an object, part of the object, or part of the background. As for other tasks, such as image classification and object detection, the state-of-the-art segmentation performance is based on the supervised training of a certain deep learning model with a set of labeled training data. Most deep-learning models for semantic segmentation have an encoder–decoder architecture, in which the encoder is responsible for extracting features from the input image and the decoder is used to predict the semantic label of each image pixel. While fully convolutional networks were previously dominant for both the encoder and decoder, as the segmentation models: FCN [15], UNet [16], and the DeepLab series [17], [18], [19], [20]. These CNN-based segmentation models have been outperformed by the transformer-style segmentation models such as SETR [21], SegFormer [22], and MaskFormer [23]. Such transformer-based segmentation models can learn to extract features for each local region by considering visual proximity between every pair of local regions using the self-attention operator at multiple scales in the models.

Besides investigating advanced model backbones for supervised semantic image segmentation, studies have focused on challenging conditions, including semi-supervised, weakly supervised, interactive, and domain-adaptive semantic segmentation. In semi-supervised semantic segmentation, only a small portion of the training dataset is assumed to be labeled. In this case, the pseudo-labeling strategy and consistency regularization are typically applied to unlabeled training data to train a better segmentation model [24], [25], [26]. Unlike semi-supervised segmentation, in weakly supervised segmentation, each training image is weakly supervised, either in the form of points [27], [28], scribbles [29], [30], or bounding box [31], [32], [33], or image-level labels [34], [35], [36]. Depending on the specific form of weak annotations, various strategies, such as the use of a class activation map (CAM) [34], [36], have been proposed to estimate more precise regions for objects of interest in the training images [37], [38], [35], and conventional supervised learning is applied based on updated and precise annotations. As an extension of the weakly supervised segmentation task, interactive image segmentation allows humans to interactively provide feedback to iteratively refine the initial segmentation result [39], [40], [41]. Such interactive segmentation is an efficient strategy for obtaining accurate segmentation results from the model at the inference (i.e., testing) phase and may be practical in real scenarios. Another practical but challenging problem is domain-adaptive semantic segmentation, in which the data domain may be shifted more or less when applying a trained segmentation model to a real-world scenario. In this case, a set of new unlabeled data from the new (target) domains are typically collected to fine-tune the pretrained segmentation model either using the training set from the old (source) domain [42], [43] or not [44], [45]. Pseudolabels are typically estimated for all or part of the unlabeled data, which are then used for model fine-tuning [44], as in the semi-supervised segmentation. Active learning can be applied by estimating a small portion of unlabeled data for humans to annotate [46], and human-annotated data (together with pseudo-labeled data) are then used to fine-tune the model. Even without any human annotations or pseudolabels, the model can be fine-tuned during testing, for example, by modifying the batch normalization parameters in each model layer with statistics for the test data [45]. Techniques for these challenging segmentation tasks have been proposed ( e.g., pseudo-labeling) or could potentially (e.g., active learning) be applied to continual semantic segmentation.

## **2.2. Continual learning and continual semantic segmentation**

Most continual learning studies have focused on classification tasks [2], [1], [3], [4], [47], whereas limited studies have focused on continual semantic segmentation using deep learning models [5], [6], [7], [8], [12]. Generally, there exist two types of continual learning tasks, task-incremental and class-

incremental. Compared to task-incremental learning which often assumes task identification and the corresponding model head is available during inference, class-incremental learning is more challenging because the model needs to be continually updated to predict all learned classes using a single model head. This study focused on the class incremental semantic segmentation problem.

Numerous approaches developed originally for continual classification may be adapted for continual semantic segmentation. A widely used approach of continual classification is knowledge distillation [2], [1], [47], [48]. In this method, the knowledge of previously learned old classes can be demonstrated by the response of the old model to the input data. Therefore, expecting that the new model has a similar response to the same input would gain or maintain a similar old knowledge of the old model. The similarity between the response of the new model and that of the old model can be measured using the cross-entropy loss if the response is the model output, as in the LwF [1] and iCarL methods [2], or generally, by the Euclidean distance or cosine distance when the response is the output of single or multiple intermediate layers, as in the LwM [4] and PODNet methods [49]. Such knowledge distillation has been effectively applied to continual semantic segmentation, for example, by distilling intermediate features [5], model output [6] or spatially pooled outputs at each intermediate layer [7], [10]. In addition to knowledge distillation, the replay strategy has been proven helpful for continual classification [2], [50], [51] and continual semantic segmentation [12]. A small amount of old data is stored for each previously learned class, and the old data are combined with new classes of training data when updating the model at each subsequent learning stage. Although the stored old data are limited compared with new classes of training data at each learning stage, applying them directly to model update and knowledge distillation can substantially alleviate the catastrophic forgetting of old knowledge [52], [53], [54]. Because of the crucial influence of old data, synthetic old data have been generated for continual classification when original data cannot be preserved because of privacy or security concerns [55]; this has been shown to be helpful for continual semantic segmentation [56].

Since the changing of existing model parameters during continual learning is the primary reason of forgetting old knowledge, researchers have attempted to keep model parameters relevant to old knowledge from changing during learning new classes of knowledge [57], [58]. An example is to fix and combine the major parts (feature extractor) of each old model into the new model structure, and thus expecting to retain all previously learned old knowledge [59]. Although such a model-growing method achieved state-of-the-art performance on continual classification, maintaining a balance between model scale and model performance over several stages of continual learning remains challenging [60]. In continual semantic segmentation, while state-of-the-art performance was achieved when keeping the segmentation model fixed except for the last layer [12], the model is too rigid to learn new classes, and its performance

degrades rapidly over more stages. In comparison, our method allows updating of all model parameters for new knowledge learning and provides an (additional) effective method to transfer old knowledge from the old to the new model.

In addition to alleviating the catastrophic forgetting issue, continual semantic segmentation has to solve the background-shifting issue [6], [7], [12], that is, the background regions in images at one learning stage may contain the regions of the classes learned at another stage. At each new learning stage, part of background regions in the new training images can be pseudo-labeled as learned old classes using the old model. Such pseudo-labeling is often helpful for alleviating this issue [7], [10], [12]. Besides pseudo-labeling, special consideration of the background class during knowledge distillation from the old model to the new model can help alleviate the background shift issue as well. Specifically, since the old model can predict each image pixel as one of old classes or the background class whereas the new model can additionally predict each image as one of the newly learned classes, the probability of each image pixel belonging to each new classes or the background class predicted by the new model should be aggregated (i.e., summed) and such aggregated probability is then compared to the probability of the pixel belonging to the background class predicted by the old model during knowledge distillation [6]. In addition, the detection of salient regions from image background and considering them as possible future classes (as a separate ‘unknown’ class at the current stage) may reduce background shifting in the future stages of continual learning [12]. Strategies to alleviate background-shifting issues can be used with those for alleviating catastrophic forgetting in practice. In this study, the proposed method can also be flexibly combined with existing strategies for continual semantic segmentation.

To date, most continual semantic segmentation models have been based on fully convolutional network (FCN) backbone. Besides FCN, the powerful transformer backbone, based on self-attention between elements within each image, has started to exhibit potential ability in solving various computer vision tasks [61], [62], [63] and Transformer-style segmentation models, such as SegFormer [22], have already shown superior performance than FCN models for semantic segmentation [19], [20]. This study is the first to evaluate the performance of the Transformer backbone on continual semantic segmentation and use of the unique self-attention information in the transformer to alleviate the catastrophic forgetting issue in continual semantic segmentation.

### **3. Method**

The objective of class-incremental semantic segmentation is to continually update a segmentation model that can learn to segment more classes of image regions. In each learning stage, a set of training

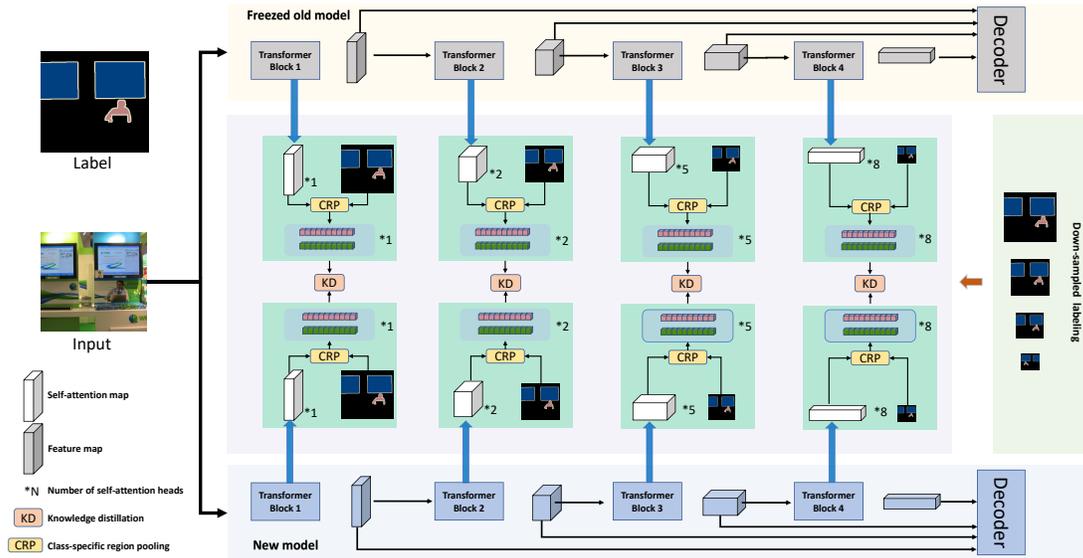


Figure 2: Framework of the proposed self-attention transfer method for continual semantic segmentation. Self-attention maps are obtained from the last multi-head self-attention layer in each transformer encoder block, and CRP (also see Figure 3) is applied to generate a self-attention vector for each foreground class appearing in the input image. The self-attention vectors are then used for within- and between-class knowledge distillation.

images that correspond to a specific number of new classes, sometimes together with a small subset of the stored old data for previously learned old classes are used to update the model. Because few or even no old data of previously learned classes are available, the primary challenge of continual semantic segmentation is alleviating the rapid forgetting of old class knowledge by the updated segmentation model, particularly over multiple stages of continual learning. In this study, by using self-attention information from the recently proposed segmentation backbone, a simple yet effective knowledge distillation strategy is proposed to substantially alleviate the catastrophic forgetting issue (Figure 2).

### 3.1. Self-attention transfer

The recently proposed transformer model SegFormer [22] achieved state-of-the-art performance on multiple image semantic segmentation tasks [13], [14]. SegFormer and previous FCN-based segmentation models [16], [18], [19], [20] is the self-attention module at each encoder layer in SegFormer. With self-attention between each element and every other element (corresponding to a single pixel or a small local region) in the input image or feature maps from each SegFormer layer, and the difference between further-

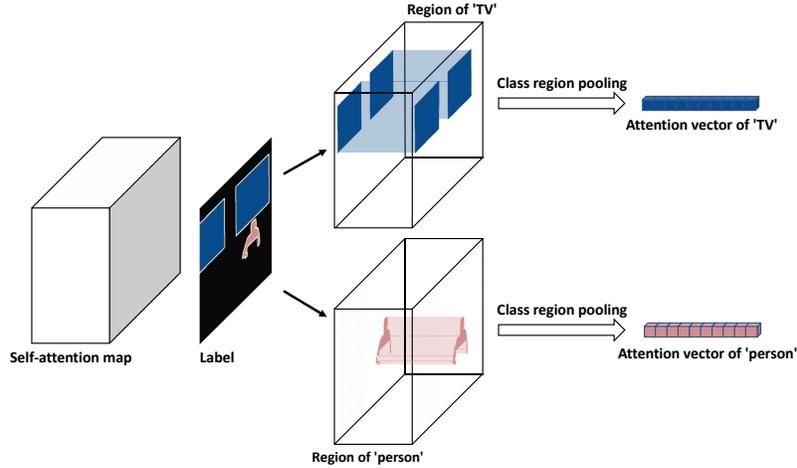


Figure 3: Class-specific region pooling.

apart local regions can be easily learned and globally to help discriminate different classes of pixels or regions. In particular, the elements belonging to the same class will have higher self-attention scores than those from different classes. This result leads to a similar attention-weighted feature vector for the elements that belong to the same classes in the feature map output of the layer. Therefore, self-attention should contain certain essential information from both within-class and between-class information.

Based on these observations, we hypothesize that transferring self-attention information from the old segmentation model to the new model during continual learning may help new model better keep old knowledge. This technique differs from the output feature maps of each layer that represent visual feature information for each input element, self-attention score vectors (from multiple self-attention heads at the layer) for each element contain various relationships between the element and every other element. Thus, self-attention scores within each layer may contain complementary information compared with the feature output of each layer, and distilling such information from the old model to the new model may help the new model remember the knowledge within each old class and across old classes.

Considering that distilling each self-attention score vector at each SegFormer encoder layer could prevent the new model from being flexibly updated to learn new classes of knowledge, we propose applying a class-specific region pooling (CRP) strategy to improve model plasticity during continual learning (Figure 3). Intuitively, given any specific input image, elements belonging to the same class will have similar self-attention vectors at each attention layer. Therefore, pooling self-attention vectors over all elements of the same class results in a single self-attention vector, which is representative of

all the elements belonging to the same class. Such pooled self-attention information may even be more robust to small variations or noise in the individual elements. Furthermore, because the feature map output from the last layer of each encoder block is passed to the decoder, the self-attention information from the last self-attention layer in each encoder block contains more information that is useful for segmentation. Therefore, only self-attention information from the last attention layer in each encoder block is considered for distillation (Figure 2), which avoids layer-wise distillation and improves model plasticity during continual learning.

Formally, we denote  $\mathbf{x}_i$  the  $i$ -th training image,  $\mathbf{A}_{i,j,h}$  the three-dimensional self-attention map for the  $h$ -th attention head at the last attention layer in the  $j$ -th encoder block,  $\mathcal{R}_{i,j,c}$  the set of element locations belonging to the  $c$ -th class in the layer. Subsequently, the pooled self-attention feature vector for each class from each attention head in each encoder block can be obtained by the following equation:

$$\mathbf{f}_{i,j,c,h} = \frac{1}{|\mathcal{R}_{i,j,c}|} \sum_{(u,v) \in \mathcal{R}_{i,j,c}} \mathbf{A}_{i,j,h}(u,v) \quad (1)$$

where  $(u, v)$  represents the location of elements belonging to class  $c$ . For any specific image, only those classes appearing in the image are considered for the pooling. Additionally, the background region in an image may contain various visual information that is difficult to be represented by one type of high-level knowledge. Therefore, the background region is excluded from the pooling procedure. Here,  $\mathbf{f}_{i,c,h}^{t-1}$  and  $\mathbf{f}_{i,c,h}^t$  denote the concatenated pooled self-attention vectors over all the encoder blocks from the SegFormer at the previous  $(t - 1)$ -th learning stage and the new SegFormer at the current  $t$ -th stage, respectively, and  $\mathcal{C}_i$  the set of classes (excluding the background class) appearing in image  $\mathbf{x}_i$ . Next, the distillation of the self-attention information can be achieved by minimizing the loss  $L_a(\boldsymbol{\theta})$  as follows,

$$L_a(\boldsymbol{\theta}) = \frac{1}{N \cdot H} \sum_{i=1}^N \left\{ \frac{1}{|\mathcal{C}_i|} \sum_{c \in \mathcal{C}_i} \sum_{h=1}^H \|\mathbf{f}_{i,c,h}^t - \mathbf{f}_{i,c,h}^{t-1}\|^2 \right\} \quad (2)$$

where  $\boldsymbol{\theta}$  is the new model parameters at the  $t$ -th stage, and  $N$  and  $H$  represent the number of training images available at the  $t$ -th stage and the number of attention heads at the last attention layer in the last encoder block, respectively.

This study is the first to apply a class-region pooling (CRP) strategy to continual semantic segmentation. Unlike existing region pooling [64], [65], [66], which operate on the feature map output of the last convolutional layer (representing the class center in the feature space), our CRP operates on self-attention maps (representing relational knowledge within and between classes) of each block of the transformer-style segmentation model. Another difference is in the role of pooling. The region pooling in related

work [64], [65], [66] is typically part of the segmentation model, and the output of the pooling is used to train a better model under a conventional segmentation setting ( that is, training a segmentation model for all classes in a single learning stage). By contrast, pooling in our method is used to distill previously learned knowledge during continual learning.

Furthermore, knowledge distillation in our SATS method differs from existing knowledge distillation for continual learning. Our SATS distills within-class and between-class relational knowledge from self-attention maps during continual semantic segmentation, while existing distillation strategies distill non-relational visual knowledge from the outputs of one or more model layers. Therefore, our method is complementary to existing knowledge distillation methods and can be used to further distill old knowledge during continual learning. Furthermore, knowledge distillation from our SATS is class-specific region-based, whereas existing knowledge distillation is based on the pooling of whole feature maps (as in ILT [5] and MiB [6]) or pooling along various spatial dimensions (as in PLOP [7]). In other words, our method tries to respectively distill each class of knowledge from each image, while existing methods distill globally averaged visual information in each image. Considering that multiple categories regions appear in each image, globally averaged visual features are less precise in representing multiple types of knowledge in the input rather than the visual features pooled from each class-specific region.

### **3.2. Flexible combination of existing strategies**

The proposed SATS is an independent knowledge distillation strategy and can be easily combined with existing continual learning strategies for semantic segmentation. Studies have revealed that (i) pseudo-labeling of background regions in new classes of images [7], [11], [12], (ii) knowledge distillation from the output of the old segmentation model [6], [7], and (iii) maintaining a small number of exemplar images for each old class [12] keeps model performance in continual learning. Thus, these three strategies were adopted as optional components of continual learning.

In the pseudo-labeling strategy, the old model from the previous learning stage is used to annotate the possible regions of the old classes in the background areas of the new image classes. At each learning stage, only regions of new classes were annotated in the new training images. Therefore, the background regions in the new images may contain regions of previously learned old classes. Annotating and collecting such regions of old classes would increase the amount of training data for old classes and benefit the training of the new segmentation model. If the pseudolabels are not available here, the new model will be confused to segment the old classes as background. Following the pseudo-labeling strategy in the PLOP method [7], only the background regions with confident predictions as old classes were pseudo-labeled.

For the output distillation strategy, special considerations were made for the background class [6]. Because the image regions of new classes are typically learned as part of background regions by old models in previous learning stages, for each training image at the new learning stage, the MiB method [6] adds the prediction probabilities of the background class and all the new classes from the new model as the modified unbiased prediction of the background class, and such modified prediction is compared to the output of the old model for each input image during knowledge distillation.

As in continual classification tasks, storing a small subset of images for each old class can substantially improve the performance of the new segmentation model for old classes. In the state-of-the-art SSUL method [6], a memory buffer with a limited size is provided to store old data, and an equivalent number of images for each old class were stored in memory. Some of the stored old data may need to be discarded to store data for more recently learned classes.

It is worth noting these three strategies are supplementary to each other and can be combined with the proposed method. In existing state-of-the-art methods, two or more strategies are always combined during continual segmentation learning. For example, in MiB and PLOP, knowledge distillation and pseudo-labeling strategies are used, whereas in SSUL, knowledge distillation, pseudo-labeling, and example replay strategies are used. In this study, the pseudo-labeling strategy and knowledge distillation from the model output are combined with the proposed method by default, and the inclusion of the memory buffer to store small old data is optional. Overall, the new segmentation model at each new learning stage is updated by minimizing the combined loss function  $L(\theta)$  over all current new classes of training images at the new learning stage and the stored small old data for each old class as follows:

$$L(\theta) = L_c(\theta) + \lambda_a L_a(\theta) + \lambda_d L_d(\theta) \quad (3)$$

where  $L_c(\theta)$  is the conventional cross-entropy loss, and  $L_a(\theta)$  is the self-attention transfer loss (Equation 2) based on the CRP of the self-attention maps from the old and the new models; here,  $L_d(\theta)$  is the unbiased knowledge distillation loss based on the outputs of both the old and new models [6], and  $\lambda_a$  and  $\lambda_d$  are hyperparameters to balance the loss terms.

## 4. Experiments

### 4.1. Settings

**Dataset:** Following existing work on continual semantic segmentation, we used classical image segmentation benchmark datasets named Pascal VOC 2012 [13] and ADE20K [14] for both quantitative and

qualitative evaluations. The VOC dataset includes 20 foreground classes and one background class, with 10582 images for training, 1449 for validation, and 1449 for testing. The ADE20K dataset includes 150 foreground classes and one background class, 20200 images for training, 2000 for validation, and 2000 for testing.

**Protocols:** Following the “m-n” protocols used in previous studies, we used the first m foreground classes and the background class in the dataset to train an initial segmentation model. The model was then updated over the stages of continual learning by each continual learning method, with each stage learning new classes. We used protocols VOC 15-1, 15-5, 5-3, and 10-1 and ADE 100-10 and 25-25. The segmentation model was initially trained and then updated in six stages (one initial stage plus five continual learning stages), two stages, six stages, and 11 stages for the four protocols on VOC and six stages on ADE. On the ADE20K dataset, the ADE 25-25 and 100-10 protocols were used because the 25-25 protocol is more challenging and practical in real applications. Based on previous studies, we used the *overlapped* setting, that is, in each training and validation image at each stage of learning, only the image regions belonging to classes that are learned at the current stage are considered as foreground regions, and other regions are considered as background, even if some of the regions belong to foreground classes that have been learned or will be learned in the future learning stages. Similarly, for test images, only the foreground classes that have been learned at the current or previous stages are considered foreground regions, and all other regions are considered background.

**Metrics:** Standard mean Intersection-over-Union (mIoU) over classes at each learning stage was used to evaluate the performance of each of these methods. For detailed investigation, the mIoUs were calculated respectively over the first set of classes in the initial learning stage, and the other classes over subsequent continual learning stages, and all the classes for each method with each protocol.

**Implementation details:** SegFormer-B2 [22] was used as the default model, whose backbone encoder MixTransformer-B2 was pretrained on ImageNet [67] following previous studies [5], [6], [7], [8], [12] on continual semantic segmentation, and the SegFormer decoder was randomly initialized. For the baseline methods using Deeplab V3, the model backbone is ResNet-101, according to previous studies [6], [7]. Also following previous studies [6], [7], during model training and updating of SegFormer-B2 or DeepLab V3, the SGD optimizer with a batch size of 24 for VOC, 32 for ADE, and momentum of 0.9, was used with an initial learning rate of 0.01 for the first (initial) learning stage, 0.001 for subsequent learning stages, and then exponentially decreased with a decay rate of 0.9 over epochs. Each model was trained over 30 epochs and was selected from those with the best performance on the validation set. The hyperparameters  $\lambda_a$  and  $\lambda_d$  were empirically set to 20. For the experiments using memory to store old data, the memory size was set to 100 for VOC and 300 for ADE for all methods following the setting from SSUL [12].

TABLE 1: Semantic segmentation performance after finishing the last stage of continual learning. With each protocol, model performance on the set of classes learned at the first learning stage (e.g., class 0-15 with protocol VOC 15-1), on the other classes learned over stages of continual learning (e.g., class 16-20 with VOC 15-1), and on all the classes (All) were reported for each method. Methods with \* indicate the results were directly obtained from the corresponding original work, and all the other results were based on our re-implementations of these methods. Methods with “-M” denote that the memory of limited old data was used. In each column, the numbers in bold represent the highest performance, and the underlined numbers represent the second-highest performance.

Method	Network	VOC 15-1 (6 Tasks)			VOC 15-5 (2 Tasks)			VOC 5-3 (6 Tasks)			VOC 10-1 (11 Tasks)		
		0-15	16-20	All	0-15	16-20	All	0-5	6-20	All	0-10	11-20	All
Joint Training	DeepLab V3	79.77	72.35	77.43	79.77	72.35	77.43	76.91	77.63	77.43	78.41	76.35	77.43
LwF-MC* [2]	DeepLab V3	6.40	8.90	6.90	58.10	35.00	52.30	20.91	36.67	24.66	4.65	5.90	4.95
ILT* [5]	DeepLab V3	8.75	7.99	8.56	67.08	39.23	60.45	22.51	31.66	29.04	7.15	3.67	5.50
MiB* [6]	DeepLab V3	35.1	13.5	29.7	75.5	49.4	69.0	57.1	42.5	46.7	12.2	13.1	12.6
PLOP* [7]	DeepLab V3	66.25	24.86	56.4	75.49	49.66	69.34	17.48	19.16	18.68	44.03	15.51	30.45
SDR* [8]	DeepLab V3+	44.7	21.8	39.2	75.4	52.6	69.9	N/A	N/A	N/A	N/A	N/A	N/A
RECALL* [56]	DeepLab V2	65.7	47.8	62.7	66.6	50.9	64.0	N/A	N/A	N/A	59.5	46.7	54.8
ST-CIL* [11]	DeepLab V3	71.4	40.0	63.6	76.7	54.3	71.1	N/A	N/A	N/A	N/A	N/A	N/A
SSUL* [12]	DeepLab V3	78.06	28.54	66.27	77.42	47.16	70.21	71.17	45.38	52.75	73.78	41.13	58.23
SSUL-M* [12]	DeepLab V3	78.92	43.86	70.58	79.53	52.87	73.19	72.91	49.02	55.85	<u>74.79</u>	48.87	62.45
Joint Training	SegFormer B2	80.84	74.97	79.44	80.84	74.97	79.44	78.36	79.87	79.44	80.46	78.32	79.44
ILT [5]	SegFormer B2	17.44	12.13	16.18	49.07	53.97	50.24	13.20	15.43	14.79	6.67	6.13	6.41
MiB [6]	SegFormer B2	73.21	37.93	64.81	78.78	60.93	74.53	61.12	58.02	58.56	48.7	39.58	44.36
PLOP [7]	SegFormer B2	64.59	37.23	58.08	72.51	48.37	66.76	35.65	32.71	33.54	48.53	33.71	41.47
SSUL [12]	SegFormer B2	79.91	40.56	70.54	79.91	56.83	74.41	74.33	60.79	64.66	74.06	51.85	63.48
ILT-M	SegFormer B2	15.15	11.01	14.16	49.85	53.49	50.72	12.91	15.43	14.71	6.75	6.07	6.42
MiB-M	SegFormer B2	73.89	58.39	70.72	79.91	63.56	75.20	70.11	65.17	66.58	69.73	56.28	63.33
PLOP-M	SegFormer B2	71.11	52.61	66.70	78.53	<u>65.58</u>	75.44	67.29	62.91	64.16	57.94	51.64	54.94
SSUL-M	SegFormer B2	<u>79.84</u>	49.33	72.58	79.84	55.82	74.12	<b>76.04</b>	61.95	65.98	74.23	52.24	<u>63.76</u>
SATS (ours)	SegFormer B2	78.38	<u>62.02</u>	<u>74.48</u>	<u>80.24</u>	61.17	<u>75.70</u>	75.43	<u>64.13</u>	<u>67.36</u>	64.27	<u>58.66</u>	61.60
SATS-M (ours)	SegFormer B2	<b>80.37</b>	<b>64.54</b>	<b>76.61</b>	<b>81.44</b>	<b>70.02</b>	<b>78.72</b>	<u>75.58</u>	<b>69.67</b>	<b>71.36</b>	<b>76.21</b>	<b>61.62</b>	<b>69.27</b>

All experiments were performed on PyTorch 1.8 with CUDA 10.2 using two NVIDIA V100 GPUs.

## 4.2. Quantitative evaluation

Multiple continual learning settings (15-1, 15-5, 5-3, 10-1 on Pascal VOC2012, and 100-10, 25-25 on ADE20K) were adopted to compare our SATS method with state-of-the-art methods for the two datasets, VOC2012 and ADE20K, respectively. For the VOC2012 dataset, among the prior methods for continual semantic segmentation with the DeepLab series (mainly DeepLab V3), SSUL performed the best regardless of memory usage. However, the performance of SSUL on continually learned classes (that is, except for the classes learned in the first learning stage) was outperformed by other methods that used large amounts of auxiliary unlabeled data, such as RECALL [56], with the setting VOC 15-1 (classes 16-20, 47.8% vs. 43.86%), and ST-CIL with VOC 15-5 (class 16-20, 54.3% vs. 52.87%). This phenomenon is mainly because SSUL fixes the majority of the models and only updates the segmentation head during continual learning. However, this limits its ability to learn new classes. Compared with the

reported results (Table 1, rows 2-10) of the strong baseline from their original work, where the DeepLab backbone was used, the proposed SATS without using the memory of old data (Table 1, second last row, “SATS”) outperformed all existing methods without using memory by a significant margin and even outperformed the state-of-the-art method (SSUL-M\*), which used stored small old data during continual learning, except for the VOC 10-1 setting.

Using the same SegFormer backbone for a fair comparison, our SATS method without memory (“SATS”) still outperformed the re-implemented representative methods ILT, MiB, PLOP, and SSUL. (Table 1, rows 12-15), except for the VOC 10-1 setting, in which the proposed method was slightly worse than SSUL for the average classification of all classes (61.60% vs. 63.48%) but significantly better than SSUL for the learned new classes (classes 11–20, 58.66% vs. 51.85%) over the ten stages of continual learning. Again, this phenomenon could be caused by SSUL’s fixing of the model parameters for the first set of 11 classes (ten foreground classes plus one background class) and only adds channels at the last layer for new classes. As a result, SSUL often exhibits superior performance on the initially learned old classes, but the model is too stable to learn new classes, particularly in more stages of learning. Compared with the freezing strategy in SSUL, our proposed method allows flexible updates of all model parameters and enables the updated segmentation model to easily learn new knowledge over continual learning.

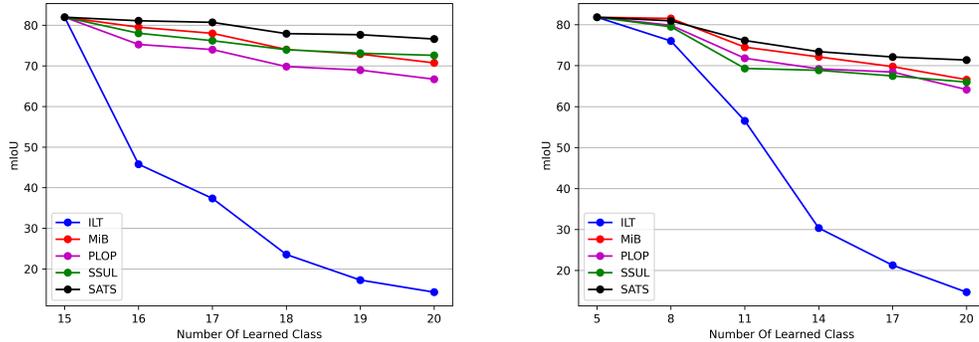


Figure 4: Semantic segmentation performance at each stage of continual learning. Our SATS method always outperforms all existing methods after each stage of learning, with the protocol VOC 15-1 (left) and 5-3 (right). All methods used memory of the same size to store limited old data during continual learning.

This result is further confirmed by another series of experiments in which all the methods used stored small old data during continual learning (Table 1, rows 16-19 and the last row). With the help of stored old data, both our method (“SATS-M”) and baselines MiB-M and PLOP-M substantially boosted performance. However, the improvement in the SSUL method (“SSUL-M”) is limited, particularly for the settings VOC

TABLE 2: Performance comparison on ADE20k dataset.

Method	ADE 100-10 (6 Tasks)			ADE 25-25 (6 Tasks)		
	0-100	101-150	All	0-25	26-150	All
MiB	40.30	17.27	32.67	54.04	23.6	28.84
PLOP	39.16	15.08	31.19	57.44	12.76	20.41
SSUL	<u>42.51</u>	16.03	33.74	<u>59.43</u>	13.87	21.88
MiB-M	41.12	18.86	33.75	54.59	24.25	29.47
PLOP-M	41.24	14.94	32.36	57.25	14.42	21.80
SSUL-M	<b>42.79</b>	15.84	33.86	<b>60.12</b>	16.89	24.33
<b>SATS (ours)</b>	41.42	<u>19.09</u>	<u>34.18</u>	57.12	<u>26.23</u>	<u>31.56</u>
<b>SATS-M (ours)</b>	41.55	<b>23.13</b>	<b>35.45</b>	57.42	<b>27.14</b>	<b>32.36</b>

15-5, 5-3, and 10-1. Because the SSUL fixes model parameters for old classes, storing small amounts of old data would play a limited role in keeping old knowledge from forgetting.

In comparison, our proposed method allows the model to be updated flexibly for new class learning and simultaneously makes superior use of stored old data to alleviate the forgetting of old knowledge. Consequently, our method (‘SATS-M’) achieves state-of-the-art performance in all four VOC settings for continual semantic segmentation. This result is further confirmed in Figures 1 and 4, which reveal that the method outperformed the representative methods at each stage of continual learning and the performance gap between existing methods and gradually increases with more stages of learning.

Similar results were obtained for the ADE20K dataset (see Table 2). For both settings, 100-10 and 25-25, notice that our method (last two rows) outperformed strong baselines on the new classes (columns 2 and 5) and all classes (columns 3 and 6), regardless of whether memory buffer is used (rows 4-6 vs. row 8) (rows 1-3 vs. row 7). These results support the generalizability of our method to different continual semantic segmentation tasks.

### 4.3. Visualization analysis

Besides the quantitative evaluation, qualitative evaluation was performed with the VOC 15-1 protocol. A set of representative and challenging test input images were selected. In these images, the foreground classes in these images were learned in the first (initial) stage of learning. The images were segmented after each segmentation model completed all stages of continual learning. Therefore, if old knowledge is catastrophically forgotten over learning stages, the model would not effectively segment these images containing earlier learned old classes.

Figure 5 reveals that all baselines except the ILT-M method can find the foreground regions. However, all of these incorrectly segmented parts of the background regions (particularly those background regions

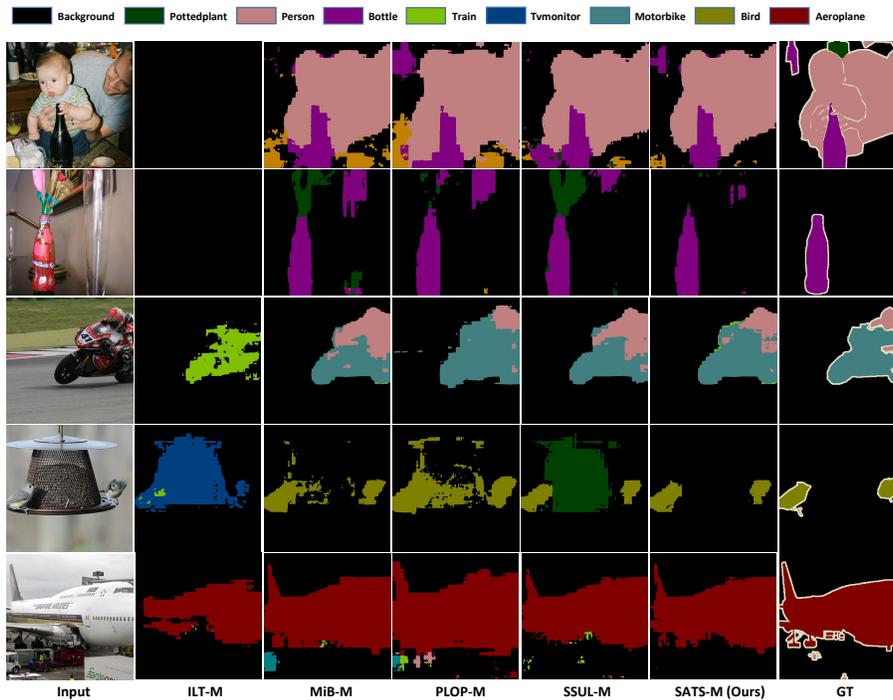


Figure 5: Demonstration of segmentation results based on our SATS method and four representative baselines.

close to the boundary of foreground objects) as foreground classes which even did not appear in the images. Our proposed method (second last column) can accurately discriminate the foreground regions against the background regions, regardless of whether the foreground object is large (first and last rows) or small (fourth row). Another observation is that the proposed method can find more accurate boundaries of foreground objects particularly when multiple foreground objects appear in the images (first, third, and fourth rows, respectively). Superior discrimination between background and foreground classes and between different foreground classes of regions using our SATS method may partly originate from the self-attention transfer, which can capture between-class relationships for knowledge distillation, whereas all existing methods do not directly consider such relationships when distilling old knowledge.

It’s worth noting that this study focuses on alleviating the catastrophic forgetting of old knowledge in continual semantic segmentation rather than on the more accurate segmentation of challenging objects (e.g., very small objects). The missing “person” from the segmentation result (Figure 5, last row, sixth

TABLE 3: Ablation study of the proposed SATS. Left: effect of SATS components on continual learning with three VOC settings. Right: the necessity of using multiscale self-attention transfer. PL: pseudo-labeling.

PL	$L_a$	$L_d$	VOC 15-1			VOC 5-3			VOC 10-1			Blocks for SATS	VOC 15-1		
			0-15	16-20	all	0-5	6-20	all	0-10	11-20	all		0-15	16-20	all
✗	✓	✓	74.04	66.18	72.17	76.56	65.89	68.94	70.92	58.65	65.08	Last one block	78.54	48.96	71.50
✓	✗	✓	76.88	58.84	72.58	75.43	62.51	66.2	74.77	52.86	64.33	Last two blocks	79.23	53.40	73.10
✓	✓	✗	78.14	49.14	71.24	68.05	64.9	65.78	72.03	60.20	66.40	Last three blocks	79.47	61.83	75.28
✓	✓	✓	<b>80.37</b>	<b>64.54</b>	<b>76.61</b>	<b>75.58</b>	<b>69.67</b>	<b>71.36</b>	<b>76.21</b>	<b>61.62</b>	<b>69.27</b>	All blocks	<b>80.37</b>	<b>64.54</b>	<b>76.61</b>

column) is probably because the adopted SegFormer backbone may not be sufficiently powerful for segmentation of very small objects in images. The other methods with SegFormer (Figure 5, columns 2-5) do not accurately segment the very small region of “person”. If a more powerful segmentation model can segment small objects, the proposed CRP would better extract the visual knowledge of small objects (from the corresponding class-specific region) and then transfer old knowledge to the new model during continual learning. Therefore, a more accurate segmentation of small objects can be considered as one segmentation challenge, which is independent of the catastrophic forgetting issue in continual learning.

#### 4.4. Ablation and sensitivity study

The effect of each component of the proposed SATS method is confirmed through an ablation study. In particular, when removing pseudo-labeling (Table 3, left, first row), the self-attention transfer loss  $\mathcal{L}_a$  (left, second row), unbiased knowledge distillation loss  $\mathcal{L}_d$  (left, third row), and model performance all decreased substantially (by 4%–5% in mIoU) for all three VOC settings (15-1, 5-3, 10-1) compared with using all (the last row). This phenomenon suggests that self-attention transfer may be used as the default component in future studies to improve continual learning performance, which is similar to the widely used knowledge distillation based on the model output and the proposed pseudo-labeling strategy.

In addition, the necessity of using multiple scales of self-attention maps has been confirmed. As displayed in Table 3 (right), when using self-attention information from only the last one, two, or three SegFormer blocks, the model performance degraded substantially from a mIoU of 76.6% to 71.5%, 73.1%, and 75.28%. This result is reasonable because SegFormer collects visual features from all four scales of encoder blocks for accurate segmentation, indicating that the self-attention information from the four encoder blocks may capture within- and between-class relationships at different visual scales, and such multiscale information would represent more relational knowledge.

Furthermore, a comparison with distillation using the encoder’s feature maps supports that the distillation of self-attention is more effective (Table 4, row 1 vs. row 3), and the results of distilling both is

TABLE 4: Ablation study of self-attention distillation, class-specific region pooling (CRP), and the model backbone. "Feature": encoder's feature maps for distillation. GP: global pooling of self-attention over the whole image region. NP: no pooling.

Setting		VOC 15-1 (6 Tasks)			VOC 10-1 (11 Tasks)		
		0-15	16-20	All	0-10	11-20	all
Distillation	Feature	78.93	59.77	74.37	74.73	59.16	67.32
	Self-Attention + Feature	<b>80.88</b>	61.78	76.34	<b>76.52</b>	59.89	68.60
	Self-Attention (Ours)	80.37	<b>64.54</b>	<b>76.61</b>	76.21	<b>61.62</b>	<b>69.27</b>
Pooling	GP	80.32	61.67	75.88	74.38	60.81	67.92
	NP	80.07	61.54	75.66	73.51	58.92	66.57
	CRP (Ours)	<b>80.37</b>	<b>64.54</b>	<b>76.61</b>	<b>76.21</b>	<b>61.62</b>	<b>69.27</b>
Backbone	CRP with DeepLab V3	78.00	45.65	70.3	66.61	44.43	56.05
	CRP with Transformer (Ours)	<b>80.37</b>	<b>64.54</b>	<b>76.61</b>	<b>76.21</b>	<b>61.62</b>	<b>69.27</b>

slightly worse than our results (Table 4, row 2 vs. row 3). Furthermore, the ablation of CRP (Table 4, rows 4-5 vs. row 6) shows that simply pooling all pixels for self-attention distillation (GP in Table 4) or trivially distilling all self-attention vectors without pooling (NP in Table 4) downgraded model performance. These results support the necessity of distillation of self-attention information and the performance improvement from the region pooling operators.

One additional ablation study was performed by replacing the transformer backbone with the DeepLab V3 backbone while keeping the other components (including CRP on feature map to perform knowledge distillation) unchanged during continual semantic segmentation. Table 4 (row 7 vs. row 8) shows that the ablated version exhibits inferior performance to the original version, supporting the fact that the relational knowledge from self-attention maps in the transformer backbone is crucial during continual semantic segmentation.

Finally, the sensitivity of the hyperparameters  $\lambda_a$  and  $\lambda_d$  in the loss function is investigated. As Figure 6 (left) shows that when  $\lambda_a$  varies from 5 to 40, the model performance fluctuates within a small range. Similarly, varying  $\lambda_d$  within the range [10, 30] results in stable model performance (Figure 6, right). This sensitivity study reveals that our method is robust to the setting of each hyperparameter within a large range of values.

#### 4.5. Limitations

Although extensive evaluations confirmed the effectiveness of the proposed SATS method on continual semantic segmentation, this study has several limitations. First, all experiments were conducted according

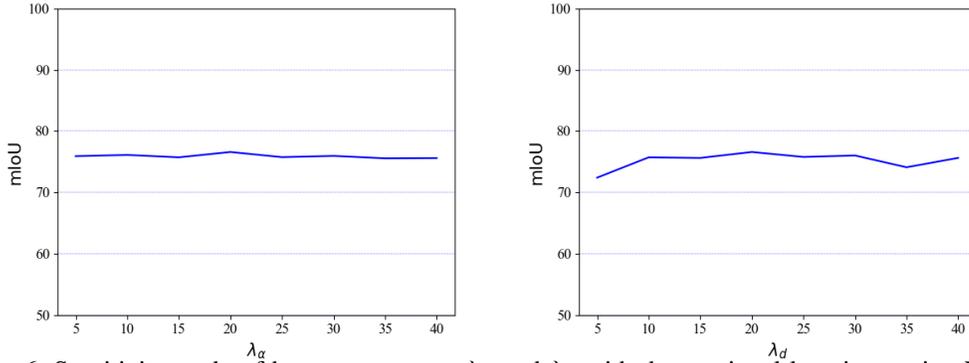


Figure 6: Sensitivity study of hyper-parameters  $\lambda_a$  and  $\lambda_d$  with the continual learning setting VOC 15-1. When varying one hyper-parameter, the other one is fixed to 20.

to previous studies, where the maximum number of total learning stages is low. It remains unclear how effective our method and other existing strategies are in more (e.g. 20 or 50) stages of continual learning. Second, the setting of the hyperparameter  $\lambda_a$  for the self-attention transfer loss  $L_a$  is affected by multiple factors. For example, without pseudo-labeling or storing a small amount of old data, the number of image regions belonging to old classes decreases, causing loss term  $L_a$  much smaller (compared with the cross-entropy loss term,  $L_c$ ). In this case,  $\lambda_a$  needs to be set to a larger constant for the self-attention transfer loss  $L_a$  to play its role during continual learning. A more adaptive setting of the hyperparameter  $\lambda_a$  should be further explored. Furthermore, the proposed SATS method can be applied to other continual tasks including continual classification learning. However, such an extension has not been investigated in this study.

## 5. Conclusion

In this study, a novel method, SATS, is proposed to effectively transfer within-class and between-class relational knowledge information from the old model to the new model during continual semantic segmentation. The relationship information can be obtained using self-attention maps from vision transformer models such as SegFormer, which provides complementary knowledge not satisfactorily captured by the output of the conventional convolution layers or fully connected layers in deep learning models. The proposed class-specific region pooling over self-attention maps can provide more efficient representations of both the within- and between-class knowledge in each image. Such an efficient representation allows the model to be flexibly updated. Therefore, a superior trade-off exists between model stability and flexibility during continual learning. The proposed SATS method is not only complementary to some

knowledge distillation strategies but also to other types of continual semantic segmentation strategies including the pseudo-labeling strategy, and an example replay strategy. We believe that the proposed SATS can become a useful plug-and-play component that can be flexibly embedded in existing or future continual semantic segmentation frameworks. One by-product of this study is the observation of the superior performance of the transformer model over conventional convolutional models on continual semantic segmentation regardless of strategies used for continual learning. The proposed self-attention transfer strategy and the transformer-style models would also work for other continual learning tasks, such as continual classification and detection learning. However, this study is limited in exploring the effectiveness of the proposed SATS method under more continual learning stages or continual learning without storing any old data, which will be investigated in future work.

### Acknowledgement

This work is supported in part by the National Natural Science Foundation of China (grant No. 62071502) and the Guangdong Key Research and Development Program (grant No. 2020B1111190001).

### References

- [1] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2935–2947, 2018.
- [2] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "iCaRL: Incremental classifier and representation learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2001–2010.
- [3] S. Hou, X. Pan, C. C. Loy, Z. Wang, and D. Lin, "Learning a unified classifier incrementally via rebalancing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 831–839.
- [4] P. Dhar, R. V. Singh, K.-C. Peng, Z. Wu, and R. Chellappa, "Learning without memorizing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5138–5146.
- [5] U. Michieli and P. Zanuttigh, "Incremental learning techniques for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019.
- [6] F. Cermelli, M. Mancini, S. R. Bulò, E. Ricci, and B. Caputo, "Modeling the background for incremental learning in semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9233–9242.
- [7] A. Douillard, Y. Chen, A. Dapogny, and M. Cord, "PLOP: Learning without forgetting for continual semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4040–4050.
- [8] U. Michieli and P. Zanuttigh, "Continual semantic segmentation via repulsion-attraction of sparse and disentangled latent representations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 1114–1124.
- [9] M. McCloskey and N. J. Cohen, "Catastrophic interference in connectionist networks: The sequential learning problem," in *Psychology of Learning and Motivation*, vol. 24, 1989, pp. 109–165.
- [10] A. Douillard, Y. Chen, A. Dapogny, and M. Cord, "Tackling catastrophic forgetting and background shift in continual semantic segmentation," *arXiv preprint arXiv:2106.15287*.
- [11] L. Yu, X. Liu, and J. van de Weijer, "Self-training for class-incremental semantic segmentation," *IEEE Transactions on Neural Networks and Learning Systems*, vol. PP, 2022.
- [12] S. Cha, Y. Yoo, T. Moon *et al.*, "SSUL: Semantic segmentation with unknown label for exemplar-based class-incremental learning," in *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [13] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [14] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba, "Scene parsing through ade20k dataset," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 633–641.

- [15] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [16] O. Ronneberger, P. Fischer, and T. Brox, “U-Net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2015, pp. 234–241.
- [17] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, “Semantic image segmentation with deep convolutional nets and fully connected crfs,” *arXiv preprint arXiv:1412.7062*.
- [18] —, “Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 4, pp. 834–848, 2018.
- [19] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, “Rethinking atrous convolution for semantic image segmentation,” *arXiv preprint arXiv:1706.05587*.
- [20] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, “Encoder-decoder with atrous separable convolution for semantic image segmentation,” in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 801–818.
- [21] S. Zheng, J. Lu, H. Zhao, X. Zhu, Z. Luo, Y. Wang, Y. Fu, J. Feng, T. Xiang, P. H. Torr, and L. Zhang, “Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 6881–6890.
- [22] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, “Segformer: Simple and efficient design for semantic segmentation with transformers,” in *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [23] B. Cheng, A. G. Schwing, and A. Kirillov, “Per-pixel classification is not all you need for semantic segmentation,” in *Advances in Neural Information Processing Systems*, 2021.
- [24] Y. Wang, H. Wang, Y. Shen, J. Fei, W. Li, G. Jin, L. Wu, R. Zhao, and X. Le, “Semi-supervised semantic segmentation using unreliable pseudo-labels,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4238–4247.
- [25] L. Yang, W. Zhuo, L. Qi, Y. Shi, and Y. Gao, “St++: Make self-training work better for semi-supervised semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 4258–4267.
- [26] J. Peng, G. Estrada, M. Pedersoli, and C. Desrosiers, “Deep co-training for semi-supervised image segmentation,” *Pattern Recognition*, vol. 107, p. 107269, 2020.
- [27] A. Bearman, O. Russakovsky, V. Ferrari, and L. Fei-Fei, “What’s the point: Semantic segmentation with point supervision,” *arXiv preprint, arXiv:1506.02106*, 2015.
- [28] R. A. McEver and B. S. Manjunath, “Pcams: Weakly supervised semantic segmentation using point supervision,” *arXiv preprint arXiv:2007.05615*, 2020.
- [29] Z. Pan, P. Jiang, Y. Wang, C. Tu, and A. G. Cohn, “Scribble-supervised semantic segmentation by uncertainty reduction on neural representation and self-supervision on neural eigenspace,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 7396–7405.
- [30] D. Lin, J. Dai, J. Jia, K. He, and J. Sun, “Scribblesup: Scribble-supervised convolutional networks for semantic segmentation,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 3159–3167.
- [31] J. Lee, J. Yi, C. Shin, and S. Yoon, “Bbam: Bounding box attribution map for weakly supervised semantic and instance segmentation,” *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2643–2651, 2021.
- [32] J. Dai, K. He, and J. Sun, “Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation,” in *IEEE International Conference on Computer Vision*, 2015, pp. 1635–1643.
- [33] T. Zheng, Q. Wang, Y. Shen, X. Ma, and X. Lin, “High-resolution rectified gradient-based visual explanations for weakly supervised segmentation,” *Pattern Recognition*, vol. 129, p. 108724, 2022.
- [34] B. Zhang, J. Xiao, Y. Wei, K. Huang, S. Luo, and Y. Zhao, “End-to-end weakly supervised semantic segmentation with reliable region mining,” *Pattern Recognition*, vol. 128, p. 108663, 2022.
- [35] P. Jiang, Y. Yang, Q. Hou, and Y. Wei, “L2G: A simple local-to-global knowledge transfer framework for weakly supervised semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022, pp. 16865–16875.
- [36] P.-T. Jiang, L.-H. Han, Q. Hou, M.-M. Cheng, and Y. Wei, “Online attention accumulation for weakly supervised semantic segmentation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 10, pp. 7062–7077, 2022.
- [37] S. Yi, H. Ma, X. Wang, T. Hu, X. Li, and Y. Wang, “Weakly-supervised semantic segmentation with superpixel guided local and global consistency,” *Pattern Recognition*, vol. 124, p. 108504, 2022.

- [38] S. Kho, P. Lee, W. Lee, M. Ki, and H. Byun, "Exploiting shape cues for weakly supervised semantic segmentation," *Pattern Recognition*, vol. 132, p. 108953, 2022.
- [39] Z. Lin, Z. Zhang, L.-Z. Chen, M.-M. Cheng, and S.-P. Lu, "Interactive image segmentation with first click attention," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 336–13 345.
- [40] X. Chen, Z. Zhao, Y. Zhang, M. Duan, D. Qi, and H. Zhao, "Focalclick: Towards practical interactive image segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, June 2022, pp. 1300–1309.
- [41] X. Luo, G. Wang, T. Song, J. Zhang, M. Aertsen, J. Deprest, S. Ourselin, T. Vercauteren, and S. Zhang, "Mideepseg: Minimally interactive segmentation of unseen objects from medical images using deep learning," *Medical Image Analysis*, vol. 72, p. 102102, 2021.
- [42] M. Chen, H. Xue, and D. Cai, "Domain adaptation for semantic segmentation with maximum squares loss," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 2090–2099.
- [43] Z. Zheng and Y. Yang, "Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation," *International Journal of Computer Vision*, vol. 129, no. 4, p. 1106–1120, 2021.
- [44] F. You, J. Li, L. Zhu, Z. Chen, and Z. Huang, "Domain adaptive semantic segmentation without source data," in *Proceedings of the 29th ACM International Conference on Multimedia*, ser. MM '21, 2021, p. 3293–3302.
- [45] F. You, J. Li, and Z. Zhao, "Test-time batch statistics calibration for covariate shift," *arXiv preprint arXiv:2110.04065*, 2021.
- [46] F. You, J. Li, Z. Chen, and L. Zhu, "Pixel exclusion: Uncertainty-aware boundary discovery for active cross-domain semantic segmentation," in *Proceedings of the 30th ACM International Conference on Multimedia*, ser. MM '22, 2022, p. 1866–1874.
- [47] F. M. Castro, M. J. Marin-Jimenez, N. Guil, C. Schmid, and K. Alahari, "End-to-end incremental learning," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 233–248.
- [48] S. Hou, X. Pan, C. C. Loy, Z. Wang, and D. Lin, "Lifelong learning via progressive distillation and retrospection," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 437–452.
- [49] A. Douillard, M. Cord, C. Ollion, T. Robert, and E. Valle, "PODNet: Pooled outputs distillation for small-tasks incremental learning," in *Proceedings of the European Conference on Computer Vision*, 2020, pp. 86–102.
- [50] E. Belouadah and A. Popescu, "Il2m: Class incremental learning with dual memory," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 583–592.
- [51] C. Zhuang, S. Huang, G. Cheng, and J. Ning, "Multi-criteria selection of rehearsal samples for continual learning," *Pattern Recognition*, vol. 132, p. 108907, 2022.
- [52] S. Mittal, S. Galesso, and T. Brox, "Essentials for class incremental learning," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2021, pp. 3508–3517.
- [53] Y. Liu, B. Schiele, and Q. Sun, "Adaptive aggregation networks for class-incremental learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 2544–2553.
- [54] Z. Li, C. Zhong, S. Liu, R. Wang, and W.-S. Zheng, "Preserving earlier knowledge in continual learning with the help of all previous feature extractors," *arXiv preprint arXiv:2104.13614*, 2021.
- [55] J. Smith, Y.-C. Hsu, J. Balloch, Y. Shen, H. Jin, and Z. Kira, "Always be dreaming: A new approach for data-free class-incremental learning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 9374–9384.
- [56] A. Maracani, U. Michieli, M. Toldo, and P. Zanuttigh, "RECALL: Replay-based continual learning in semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 7026–7035.
- [57] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska, D. Hassabis, C. Clopath, D. Kumaran, and R. Hadsell, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.
- [58] J. Zhang, J. Zhang, S. Ghosh, D. Li, S. Tasci, L. Heck, H. Zhang, and C.-C. J. Kuo, "Class-incremental learning via deep model consolidation," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*, 2020, pp. 1131–1140.
- [59] S. Yan, J. Xie, and X. He, "DER: Dynamically expandable representation for class incremental learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3014–3023.
- [60] Y. Yang, B. Chen, and H. Liu, "Bayesian compression for dynamically expandable networks," *Pattern Recognition*, vol. 122, p. 108260, 2022.
- [61] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, "Pyramid vision transformer: A versatile backbone for dense prediction without convolutions," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 548–558.

- [62] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE International Conference on Computer Vision*, 2021, pp. 10012–10022.
- [63] X. Yu, J. Wang, Y. Zhao, and Y. Gao, "Mix-vit: Mixing attentive vision transformer for ultra-fine-grained visual categorization," *Pattern Recognition*, vol. 135, p. 109131, 2023.
- [64] F. Zhang, Y. Chen, Z. Li, Z. Hong, J. Liu, F. Ma, J. Han, and E. Ding, "Acfnet: Attentional class feature network for semantic segmentation," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, p. 6798–6807.
- [65] X. He, J. Liu, J. Fu, X. Zhu, J. Wang, and H. Lu, "Consistent-separable feature representation for semantic segmentation," in *AAAI Conference on Artificial Intelligence*, 2021, pp. 1531–1539.
- [66] H. Ma, X. Lin, Z. Wu, and Y. Yu, "Coarse-to-fine domain adaptive semantic segmentation with photometric alignment and category-center regularization," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 4050–4059.
- [67] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.