

# BPKD: Boundary Privileged Knowledge Distillation For Semantic Segmentation

Liyang Liu<sup>1</sup> Zihan Wang<sup>2</sup> Minh Hieu Phan<sup>1</sup> Bowen Zhang<sup>1</sup> Yifan Liu<sup>1\*</sup>

<sup>1</sup> The University of Adelaide, Australia

<sup>2</sup> The University of Queensland, Australia

{akide.liu, vuminhhie.phan, b.zhang, yifan.liu04}@adelaide.edu.au, zihan.wang@uq.edu.au

## Abstract

Current approaches for knowledge distillation in semantic segmentation tend to adopt a holistic approach that treats all spatial locations equally. However, for dense prediction tasks, it is crucial to consider the knowledge representation for different spatial locations in a different manner. Furthermore, edge regions between adjacent categories are highly uncertain due to context information leakage, which is particularly pronounced for compact networks. To address this challenge, this paper proposes a novel approach called boundary-privileged knowledge distillation (BPKD). BPKD distills the knowledge of the teacher model’s body and edges separately from the compact student model. Specifically, we employ two distinct loss functions: 1) Edge Loss, which aims to distinguish between ambiguous classes at the pixel level in edge regions, and 2) Body Loss, which utilizes shape constraints and selectively attends to the inner-semantic regions. Our experiments demonstrate that the proposed BPKD method provides extensive refinements and aggregation for edge and body regions. Additionally, the method achieves state-of-the-art distillation performance for semantic segmentation on three popular benchmark datasets, highlighting its effectiveness and generalization ability. BPKD shows consistent improvements over various lightweight semantic segmentation structures. The code is available at <https://github.com/AkideLiu/BPKD>.

## 1. Introduction

Semantic segmentation is a complex computer vision task that involves assigning unique categories to each pixel of an input frame. In recent years, deep learning models with large numbers of parameters have achieved remarkable performance in semantic segmentation [14, 58, 32]. However, such models are impractical for resource-constrained

\*Corresponding author.

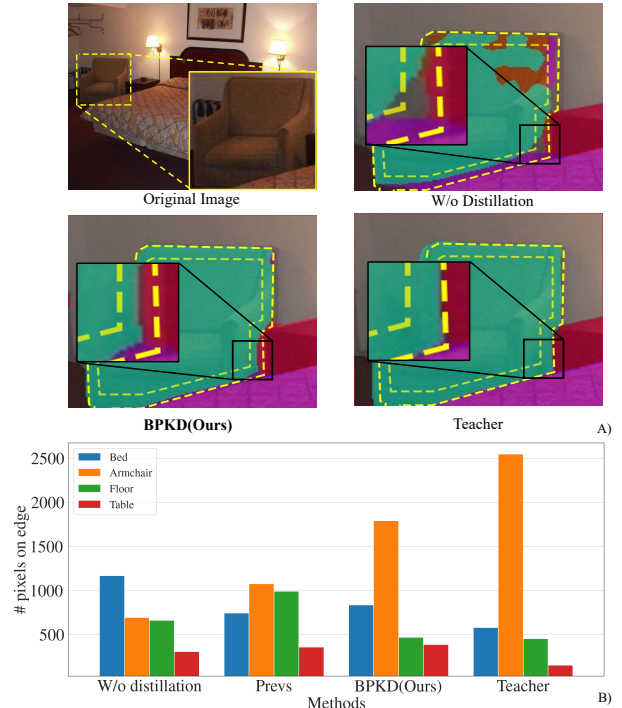


Figure 1: Illustration of context information leaking: A) The image shows the prediction images for the student, teacher, and BPKD. The calculated edge region is enclosed within a yellow dotted box. B) The image depicts the distribution of pixels in the edge area. The ground truth belong to the Armchair. It is evident that the compact model needs to improve in distinguishing adjacent categories.

devices like mobile devices and robotics due to their high computational complexity [45, 54, 49]. To address this issue, lightweight base models, such as MobileNet [18], ShuffleNet [34], and EfficientNet [42], have been used for real-time semantic segmentation.

Designing compression and acceleration techniques for compact networks is challenging but crucial. Knowledge distillation approaches, such as those introduced in

[16, 55, 23, 15], train a smaller student network to mimic the more complex teacher network by minimizing the soft probabilities distance, typically measured by Kullback–Leibler divergence, between the student and teacher. In [53, 26, 57], authors have attempted to distill hidden knowledge by utilizing network and data relations, with a focus on classification tasks, achieving impressive results.

Pioneering knowledge distillation methods for semantic segmentation [51, 41, 47, 29] focus more on capturing the correlational information among pixels, channels, and images. Liu et al. [29] suggest that hidden knowledge in semantic segmentation is constructed through structured representation. Structured knowledge is more suitable for pairwise similarity reduction and holistic distillation. Wang et al. [47] propose encoding the knowledge based on semantic masks. In [41], authors refine distillation by emphasizing the alignment of the most salient region of each channel between the teacher and student.

In semantic segmentation, accurately distinguishing adjacent categories in sensitive regions, particularly edge slices, is a challenging task. Previous studies have primarily focused on transferring knowledge representations of the entire image [29, 41, 22, 47, 51, 52, 1], neglecting the importance of distinct knowledge representations at different spatial locations.

As illustrated in Figure 1, edge pixels in semantic segmentation tasks extract contextual information for both spatially adjacent categories [50], leading to the problem of ‘context information leaking’. This issue results in higher uncertainty levels between spatially adjacent categories for edge pixels, particularly in low-capacity student networks. Therefore, developing methods to address this issue is crucial for semantic segmentation knowledge distillation.

In this study, we propose a novel approach called Boundary Privileged Knowledge Distillation (BPKD). To tackle the issue of context information leakage in existing methods, we bifurcate the knowledge distillation process into two subsections: the edge distillation and the body distillation sections. This is demonstrated in Figure 2. Our proposed BPKD approach explicitly enhances the quality of edge regions and object boundaries by decoupling knowledge distillation and using teacher soft labels. The edge distillation loss involves spatial probability alignment and aggregation of context information to refine the boundaries. Furthermore, boundaries provide prior knowledge of the shape of an object’s inner regions, and the body region can exploit this knowledge to eliminate high-uncertainty boundary samples and smooth the learning curves. Consequently, we observed that the object center received greater attention due to the implicit shape constraints.

Through empirical analysis, we have demonstrated that our proposed approach effectively guides the student network to learn from the teacher network’s knowledge, re-

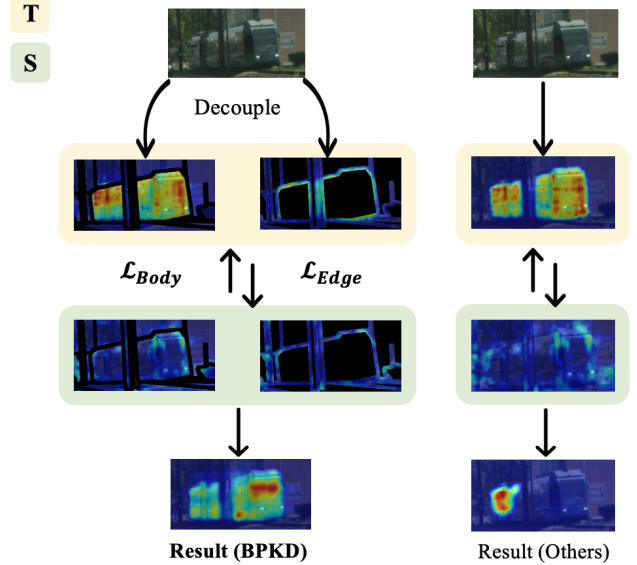


Figure 2: BPKD (left) decouples the edge and body information from the input image and creates two parallel distillation pipelines compared to pioneering works, such as CWD [41], SKDS [29], IFVD [47], CIRKD [51]. The compositional loss forces the student to learn each part separately. The result shows performance gain explicitly.

sulting in improved segmentation performance. We evaluate our method on popular architectures using three segmentation benchmark datasets: Cityscapes [9], ADE20K [61], and Pascal Context [11]. Experimental results indicate that BPKD outperforms other state-of-the-art distillation approaches. In particular, we reduce the performance gap between the student and teacher networks from 7.42% to 2.80%, surpassing the previous method CWD by 1.71%, and achieving competitiveness with specialized real-time segmentation techniques [12].

Our main contributions are summarized as follows:

- We identify and address the problem of ‘context information leaking’, which leads to uncertainty in out-of-distribution data, by proposing solutions to mitigate this issue. Notably, our work is the first to recognize and investigate this problem.
- We introduce a novel knowledge distillation approach that separately distills body and edge knowledge. Our Edge loss provides explicit enhancement to the edge slices, delivering strong shape constraints to the body regions. This reduces the uncertainty of object slices in knowledge distillation and enhances attention to the inner object region.
- Our proposed method achieves state-of-the-art distillation performance for semantic segmentation on three popular benchmark datasets. Additionally, we observe a significant improvement in the prediction quality in the edge region compared to other methods.

## 2. Related Work

**Semantic segmentation.** Many early works incorporated probabilistic graphical models to integrate more semantic context, such as Conditional Random Fields (CRFs) and Markov Random Fields (MRFs) [2, 40, 60]. However, recent approaches are predominantly based on Fully Convolutional Networks (FCNs) [33, 28, 56], which have achieved state-of-the-art performance in many segmentation benchmarks. Some of the most accurate methods for semantic segmentation are based on large models with complex architectures. For example, PSPNet [59] proposes a pyramid pooling module (PPM) to enlarge the receptive field in a multi-scale context. Meanwhile, the DeepLab series [2, 3, 4, 5] uses atrous spatial pyramid pooling (ASPP) to capture multi-scale contextual information. HRNet [45] proposes a parallel backbone to maintain high-resolution representations throughout the network.

However, the high computational cost of the aforementioned models makes them unsuitable for real-time inference on edge devices, prompting a surge of interest in lightweight models that can perform efficient segmentation. Several popular approaches have emerged, including ENet [37], which employs early downsampling, small decoder size, and filter factorization, and SqueezeNet [21], which utilizes fire modules and parallel convolutions for efficient segmentation. ESPNet [35] adopts efficient spatial pyramid and filter factorization techniques to reduce computational costs. Alongside segmentation network design, lightweight feature extraction networks such as MobileNet [19], MobileNetV2 [39], and MobileNetV3 [18] have also demonstrated effectiveness in efficient semantic segmentation.

**Knowledge distillation.** Knowledge distillation (KD) is a technique that compresses knowledge from one or more teacher models into a smaller student model [16, 13]. Current KD methods are mostly focused on fundamental vision tasks and can be categorized into response-based, features-based, and relation-based knowledge. The response-based method, proposed by Hinton et al. [16], is the most popular, where Kullback-Leibler divergence is minimized to transfer informative dark knowledge between the teacher and the student network. FitNet [38] is a representative features-based method that matches the feature activations of the teacher and the student network, while additive methods match features indirectly [55, 23, 15]. Relation-based knowledge distillation [53, 26, 57] further explores the relationships between different layers or data samples. However, these KD methods are mainly designed to improve image classification performance and are not tailored to pixel-level segmentation problems.

Over the past few years, knowledge distillation (KD) methods have been increasingly employed in semantic segmentation for model compression. One such method, structural knowledge distillation [29, 30], defines semantic seg-

mentation as a structured prediction problem and transfers knowledge through pairwise similarities and holistic adversarial enhancement. Channel-wise knowledge distillation [41] obtains a channel-wise probability map and focuses on the most salient regions of different channels between the teacher and the student. Intra-class feature variation distillation [47] combines pixel-level and class-wise intra-class feature variation to improve segmentation accuracy. Cross-Image Relational knowledge distillation [51] reweighs the importance of global semantic relations among pixels across various images and delivers improved global structured features from the teacher network. Masked Generative Distillation [52] improves student representation by utilizing the teacher to guide the student’s feature recovery. Recent work [1] also shows that Pearson correlation coefficient is an alternative to KL divergence for distillation. Empirical studies have shown that these KD methods effectively enhance the performance of semantic segmentation.

In recent years, several techniques for knowledge distillation in semantic segmentation have been proposed, each focusing on different aspects of the problem. However, these methods have often neglected the significance of regional feature differences and their relative importance. To address this limitation, our approach involves splitting the knowledge distillation pipelines and improving performance.

## 3. Methods

In this section, we first provide an overview of the workflow of the Boundary Privileged Knowledge Distillation (BPKD) framework (Section 3.1), followed by a detailed description of two key implementation aspects of our approach. Specifically, Section 3.2 outlines the process of edge knowledge distillation, which involves pre-mask filtering and post-mask filtering. Section 3.3 introduces the distillation loss for body enhancements.

### 3.1. BPKD Framework

Existing feature distillation techniques [41, 51] transfer the whole-view representations from the teacher while overlooking the effects of noisy edge features on the distillation process. In this framework, we carefully consider the sensitivity of edge representations and introduce the novel boundary-privileged knowledge distillation (BPKD) that transfers knowledge in the body and the edge regions separately. Distilling edge regions individually enhances the quality of object boundaries explicitly. Furthermore, the edge distillation loss provides prior shape knowledge for the object’s inner regions. For instance, given a vehicle’s boundary constraint, the model can easily determine pixel categories for its inner region. The body distillation loss has two key benefits from the prior boundary knowledge:

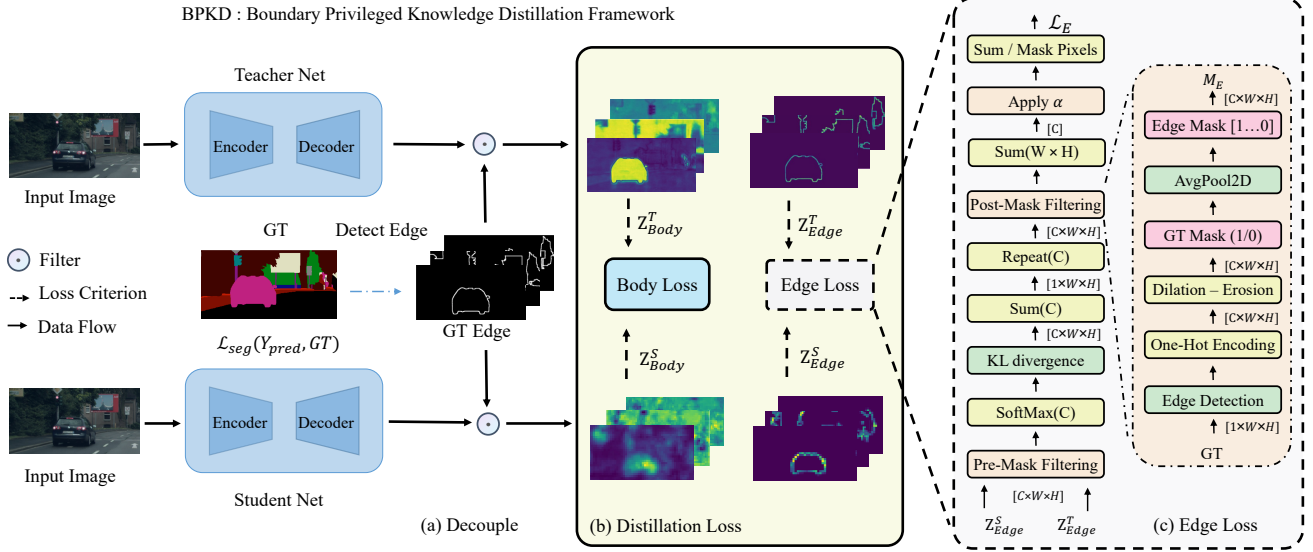


Figure 3: Illustration of our proposed **Boundary Privileged Knowledge Distillation** framework and architecture. (a) the decoupling process involves the Edge Detection on the ground truth to generate  $GT_{Edge}$  Mask, followed by applying the mask filter to obtain the Teacher and Student logits masks. This step ensures that the information from the boundary region is isolated and appropriately conveyed to the Student. (B) shows that distillation comprises two terms: body loss and edge loss. The body loss term captures the categorized similarities, whereas the edge loss term concentrates on the boundary regions’ transfer. (C) shows edge loss calculation is performed in two stages: pre-mask filtering and post-mask filtering. The pre-mask filtering step shapes the probability distribution to contain only edge information. Subsequently, the post-mask filtering step aggregates context information between adjacent categories to produce the final edge loss.

(1) reducing learning difficulty by mimicking the teacher’s logit probability distribution since high-uncertainty boundaries are removed, and (2) leveraging higher attention on the object center through implicit shape constraints.

Our approach uses an edge detection technique to generate edge masks  $M_E$  for each class by processing the ground truth and the segmentation logit map. Let  $\mathcal{Z} \in \mathbb{R}^{H \times W \times C}$  denote a network’s logit map, where  $C$  corresponds to the number of channels and  $H \times W$  represents the spatial resolution. The edge masks  $M_E$  are applied to separate the logit map  $\mathcal{Z}$  into two components: the body component  $\mathcal{Z}_B$  and the edge component  $\mathcal{Z}_E$ , which adhere to an additive rule, denoted by  $\mathcal{Z} = \mathcal{Z}_B + \mathcal{Z}_E$ . Our BPKD framework separately transfers the edge and body knowledge encoded in these two components to the student. As the edge slices have fewer knowledge representations, we introduce categorized awareness to balance the importance of different spatial perspectives. Together, these techniques play a crucial role in improving the overall performance of the model.

In this study, we propose a novel approach for decomposing the distillation loss into two distinct components, namely the body loss  $\ell_B$  and the edge loss  $\ell_E$ , as expressed by Equation 2. We include the body loss weight  $\lambda_b$  and the edge loss weight  $\lambda_e$  to control the contributions of each loss term. This decoupling strategy allows us to examine the sensitivity of edge learning in the knowledge distillation

process, which has been overlooked in the literature. The loss objective is defined as:

$$\ell = \lambda_b \cdot \ell_B(\mathcal{Z}_B^S, \mathcal{Z}_B^T) + \lambda_e \cdot \ell_E(\mathcal{Z}_E^S, \mathcal{Z}_E^T). \quad (1)$$

### 3.2. Edge knowledge representation.

Our framework aims to minimize the discrepancy between the teacher’s and the student’s features on the edge areas. To achieve this, we extract an edge knowledge representation using soft edge masks. These masks are applied to the logit map to produce masked feature representations for both the teacher’s and the student’s edge regions. The edge map  $M_E$  is created through two stages: Pre-Mask Filtering (PRM), which captures edge discrepancy for all classes, and Post-Mask Filtering (POM), which extracts edge discrepancy for each individual class. This approach allows the model to extract a more accurate and precise representation of the edge knowledge, leading to improved performance in detailed classification.

During the edge detection process, we utilize an adjustable Trimap algorithm [44] to detect edge presentations  $\mathcal{Z}_E$  from the ground truth. Specifically, we calculate the binary edge mask  $GT_{edge} = \text{dilation}(GT) - \text{erosion}(GT)$  as the difference between the dilation and erosion operations applied to the ground truth  $GT$ . The resulting binary mask  $GT_{edge} \in \mathbb{R}^{C \times H \times W}$  is then subjected to average pooling



to obtain  $M_E \in \mathbb{R}^{C \times H' \times W'}$ , which has the same shape as the logits prediction  $Z_{\text{pred}}$ . The size of  $M_E$  is determined by the output stride  $S$  of the segmentation network, where  $W' = W/S$ .

- **Pre-mask filtering (PRM).** To obtain the logits map, we apply  $M_E$  to both student and teacher logits:  $Z_E = Z_{\text{pred}} \cdot M_E$ . Specifically, we apply the edge mask for each channel  $c$  so that we can concentrate the logits on the overlapping edge regions. For example, if a frame displays a dog and a cat standing nearby, only the logits activation of the dog and the cat class will be considered, while all other activations will be suppressed. This operation forces the student to focus more on the correlations between adjacent ambiguous classes. A spatial-level KL divergence loss is applied to the filtered logits  $Z_E^T$  and  $Z_E^S$ :

$$\varphi(Z_E^T, Z_E^S) = \sum_{c=1}^C \phi(Z_{E,i}^T) \cdot \log \frac{\phi(Z_{E,i}^T)}{\phi(Z_{E,i}^S)}, \quad (2)$$

where  $\phi$  is the softmax operation for each pixel.  $\varphi(Z_E^T, Z_E^S)$  represents the edge-masked KL distances for all spatial locations.

- **Post-mask filtering (POM).** We further separate the edge loss for each class and perform normalization based on the edge area by Post-mask Filtering (POM). Let  $Z_{E,i,c}^T$  and  $Z_{E,i,c}^S$  denote the logits for the  $c$ -th class at pixel  $i$  in the teacher and student models, respectively. Let  $M_{E,c}$  be the soft mask obtained by average-pooling the ground truth binary edge mask for the  $c$ -th class, and  $n_c$  denotes the number of non-zero pixels in  $M_{E,c}$  for this class. Our POM term can be formulated as follows:

$$\ell_E = \sum_{c=1}^C \frac{\alpha_c}{n_c} \sum_{i=1}^{W \cdot H} \varphi(Z_{E,i,c}^T, Z_{E,i,c}^S) \cdot M_{E,c}, \quad (3)$$

By re-weighting the loss based on the edge area of each class, we prioritize the center of the edge, where the most important information is often located. This approach ensures that the student model focuses on learning the correct edge positions and shapes for each class.

The Soft Edge Masks  $M_E$  play a critical role in the Edge Loss, and our approach to generating them involves two specialized designs: 1) converting the binary mask  $GT_E$  into a weighted discrete space, and 2) generating masks per channel instead of a unified mask. Instead of directly applying binary masks, we use average pooling to generate softer masks. This helps minimize unconfident bias and noise, providing a more accurate representation of the edge knowledge. By carefully considering overlapping masks, we aim

to exclude confident regions and focus on areas with uncertainty, thereby refining the knowledge distributions.

In summary, the proposed Pre-Mask Filtering (PRM) and Post-Mask Filtering (POM) stages in the edge region refine the knowledge distillation process by identifying edge discrepancies for each class and re-weighting the loss based on the edge area. This method ensures that the student model learns the correct edge positions and shapes for each class, while providing shape prior knowledge for the body knowledge representation.

### 3.3. Body knowledge representation.

This section focuses on the body knowledge distillation. Previous works mainly consider whole-view distillation, which includes noisy edge representations in the body knowledge. To overcome this challenge, we utilize the reversed edge binary mask to extract body masks. By removing the edge region, we exploit implicit shape constraints and reduce uncertainty, allowing the body loss to focus on assigning the large inner regions of objects to their corresponding categories. We propose a region alignment approach that synthesizes channel-level activations to obtain semantically rich sections.

Since we defined  $Z = Z_B + Z_E$ , the body logits are obtained by  $Z_B = Z \times (1 - M_E)$ . A pixel-wise loss for the body region would introduce unexpected noise due to hard constraints. Therefore, we employ a loose constraint of channel-wise distillation [41] for the body part. The body enhancement loss (BEL) is defined as:

$$\ell_B = \frac{\mathcal{T}^2}{C} \sum_{c=1}^C \sum_{i=1}^{W \cdot H} \phi(Z_{B,c,i}^T) \cdot \log \left[ \frac{\phi(Z_{B,c,i}^T)}{\phi(Z_{B,c,i}^S)} \right], \quad (4)$$

## 4. Experiments

### 4.1. Experimental Setup

**Dataset.** We conducted the experiments on three benchmark datasets for semantic segmentation: Cityscapes [9], Pascal Context 2010 [11], and ADE20K [61].

**ADE20K [61].** ADE20K contains 20k/2k/3k images for train/val/test with 150 semantic classes. It was constructed as the benchmark for scene parsing and instance segmentation.

**Cityscapes [9].** Cityscapes is an urban scene parsing dataset that contains 2975/500/1525 finely annotated images used for train/val/test. The performance is evaluated on 19 classes.

**Pascal Context [11].** The Pascal Context dataset provides dense annotations, which include 4998/5105/9637 train/val/test images. We use 59 object categories for training and testing. Our results are reported on the validation set.

Table 1: Performance comparison of different distillation methods with state-of-the-art techniques. We test these methods on various segmentation networks for both student and teacher models, using datasets including Cityscapes[9], ADE20K [61], and Pascal Context [11]. The FLOPs are obtained on  $512 \times 512$  resolutions. Our BPKD outperforms all previous methods in a large margins across multiple datasets and network architectures.

Methods	FLOPs(G)	Parameters(M)	ADE20K 80k 512*512		Cityscapes 80k 1024*512		Pascal Context 59 80k 480*480	
			mIoU(%)	mAcc(%)	mIoU	mAcc(%)	mIoU	mAcc(%)
T: PSPNet-R101[59]	256.89	68.07	44.39	54.75	79.74	86.56	52.47	63.15
S:PSPnet-R18[59]	54.53	12.82	33.30	42.58	74.23	81.45	43.79	54.46
SKDS [29]	54.53	12.82	34.49	44.28	76.13	82.58	45.08	55.56
IFVD [17]	54.53	12.82	34.54	44.26	75.35	82.86	45.97	56.6
CIRKD [51]	54.53	12.82	35.07	45.38	76.03	82.56	45.62	56.15
CWD [41]	54.53	12.82	37.02	46.33	76.26	83.04	45.99	55.56
<b>BPKD(Ours)</b>	54.53	12.82	<b>38.51</b>	<b>47.70</b>	<b>77.57</b>	<b>84.47</b>	<b>46.82</b>	<b>56.29</b>
T: HRNetV2P-W48 [45]	95.64	65.95	42.02	53.52	80.65	87.39	51.12	61.39
S:HRNetV2P-W18S [45]	10.49	3.97	31.38	41.39	75.31	83.71	40.62	51.43
SKDS [29]	10.49	3.97	32.57	43.22	77.27	84.77	41.54	52.18
IFVD [17]	10.49	3.97	32.66	43.23	77.18	84.74	41.55	52.24
CIRKD [51]	10.49	3.97	33.06	44.30	77.36	84.97	42.02	52.88
CWD [41]	10.49	3.97	34.00	42.76	77.87	84.98	42.89	53.37
<b>BPKD(Ours)</b>	10.49	3.97	<b>35.31</b>	<b>46.11</b>	<b>78.58</b>	<b>85.78</b>	<b>43.96</b>	<b>54.51</b>
T:DeeplabV3P-R101 [4]	255.67	62.68	45.47	56.41	80.98	88.7	53.20	64.04
S:DeeplabV3P+MV2 [39]	69.60	15.35	31.56	45.14	75.29	83.11	41.01	52.92
SKDS [29]	69.60	15.35	32.49	46.47	76.05	84.14	42.07	55.06
IFVD [17]	69.60	15.35	32.11	46.07	76.97	84.85	41.73	54.34
CIRKD [51]	69.60	15.35	32.24	46.09	77.71	85.33	42.25	55.12
CWD [41]	69.60	15.35	35.12	49.76	77.97	86.68	43.74	56.37
<b>BPKD(Ours)</b>	69.60	15.35	<b>35.49</b>	<b>53.84</b>	<b>78.59</b>	<b>86.45</b>	<b>46.23</b>	<b>58.12</b>
T:ISANet-R101 [20]	228.21	56.80	43.80	54.39	80.61	88.29	53.41	64.04
S: ISANet-R18 [20]	54.33	12.46	31.15	41.21	73.62	80.36	44.05	54.67
SKDS [29]	54.33	12.46	32.16	41.80	74.99	82.61	45.69	56.27
IFVD [17]	54.33	12.46	32.78	42.61	75.35	82.86	46.75	56.4
CIRKD [51]	54.33	12.46	32.82	42.71	75.41	82.92	45.83	56.11
CWD [41]	54.33	12.46	37.56	45.79	75.43	82.64	46.76	56.48
<b>BPKD(Ours)</b>	54.33	12.46	<b>38.73</b>	<b>47.92</b>	<b>75.72</b>	<b>83.65</b>	<b>47.25</b>	<b>56.81</b>

**Implementation details.** Our implementation is based on the open-source toolbox MMSegmentation [7, 8] with PyTorch 1.11.0. We employ standard data augmentation, including random flipping, cropping, and scaling in the range of [0.5, 2]. All experiments are optimized by SGD with a momentum of 0.9 and a batch size of 16. We use a crop of  $512 \times 512$ ,  $512 \times 1024$ , and  $480 \times 480$  for ADE20k, Cityscapes, and Pascal Context, respectively. We use an initial learning rate of 0.01 for ADE20K and Cityscapes. In addition, we use an initial learning rate of 0.004 for Pascal Context. The number of total training iterations is 80K. Following previous methods [5, 59], we use the poly learning rate policy and report the single-scale testing result. We set up a fair comparison by assigning identical parameters for each method with the same dataset. Mean Intersection-over-Union (mIoU), Trimap mIoU, and pixel mean accu-

racy (mAcc) are employed as the main evaluation metrics. GFLOPs and No. Parameters are also reported for various student networks that we tested. All reported computational costs are measured using fvcare<sup>1</sup>.

## 4.2. Comparison with State-of-the-Art Methods

For a fair comparison, we re-implemented previous knowledge distillation methods [29, 47, 41, 51], as well as our BPKD, and applied them to various compact networks: PSPNet with ResNet18 backbone [59], HRNet-W18 [45], Deeplab-V3+ with MobileNetV2 backbone [39], and ISANet with ResNet18 [20]. All the distillation methods were trained with the same configurations. We conducted all experiments using 4 NVIDIA A100 GPUs.

<sup>1</sup><https://github.com/facebookresearch/fvcare>

Method						
Teacher: PSP-ResNet101			79.74%			
Student: PSP-ResNet18			Standard: 68.99%    Trimap: 55.34%			
Channel Wise Distillation [41]			Standard: 74.29%    Trimap: 57.34%			
Pixel Wise Distillation [29]			Standard: 69.33%    Trimap: 53.82%			
		Body(C)	Body(P)	Edge(P)	Edge(C)	Ours
mIoU(%)	Standard	74.17 (▲5.18)	72.70 (▲3.71)	71.63 (▲2.64)	66.83 (▼2.16)	<b>75.94 (▲6.95)</b>
	Trimap	56.20 (▲0.86)	54.12 (▼1.22)	61.37 (▲6.03)	51.76 (▼3.58)	<b>62.91 (▲7.57)</b>

Table 2: The effectiveness of the decoupling whole-view knowledge representation. The results show knowledge representation for different spatial locations should be considered, separately. C and P denote channel-wise and pixel-wise knowledge distillation, respectively.

**Performance.** Table 1 shows the performance on the ADE20K validation set. The proposed BPKD achieved the state-of-the-art (SOTA) performance across various student networks. Furthermore, the distillation improved the students’ mIoU by 15.64%, 12.58%, 12.45%, and 24.33%, respectively. In particular, our BPKD consistently outperformed the current SOTA method, CWD, by up to 3.87% across all network structures. These results demonstrate that our methods guide the students to increase their capacity without changing the network architecture for efficiency. The table also summarizes our results on the Pascal Context validation set, where the performance shows an average increase of 2% compared to the SOTA method across different network structures.

### 4.3. Ablation Study

In this section, we comprehensively evaluate our BPKD under different settings. All ablation experiments are conducted on the Cityscapes dataset with T: PSPNet-R101 and S: PSPNet-R18. We set a reduced training setting to reduce the computational cost, including a crop size reduced to  $512 \times 512$  and a training schedule of 40k iterations. More experimental results are shown in the supplementary material. **Effectiveness of the decoupled knowledge:** To verify the effectiveness of the proposed knowledge distillation approach, we evaluate the segmentation performance in the edge region as presented in Table 2.

In this study, we employed the popular Trimap mIoU metric [5], which is a segmentation edge-enhancing metric [24, 25], to evaluate our edge quality. We employed both channel-wise and pixel-wise normalization techniques to enable the student network to learn comprehensive knowledge representations from the teacher network. We observed that channel-wise normalization for distillation leads to a reduction in the standard mIoU by 2.16% and Trimap mIoU by 3.58%. This is because the edge region contains fewer pixels, making it more sensitive. Thus, a channel-wise strategy may result in misdirection for the student network when enumerating cropped boundary knowledge representation.

In contrast, pixel-wise distillation has significantly improved Trimap mIoU by 6.03%, indicating that it preserves fine edge details better than the channel-wise approach. Our Edge loss has played a crucial role in its specialized design by involving spatial probability alignment and context information aggregation, leading to explicit refinement of boundary quality.

We have also observed that the Body loss function has provided a performance boost of 5.18% and 3.71% when calculating the loss in a channel-wise and spatial-wise manner, respectively. This suggests that categorical distillation is more suitable for the body region. However, no significant improvement in Trimap performance was observed for our body loss, as it ignores boundary information.

In summary, our combination of techniques has yielded the best mIoU results of 75.94 (+6.95)% on standard evaluation and 62.91 (+7.59)% on Trimap measurement. This indicates that BPKD has significantly improved the quality of semantic boundaries by enabling pixel-level alignment and context information aggregation. Leveraging the prior knowledge constraints provided by the Edge loss terms, the body distillation process has enhanced the larger body region with better concentration, producing the best results.

Method	mIoU(%)	IMP.(%)	mAcc(%)
Teacher	79.74	-	86.56
ResNet18	68.99	-	75.19
+ PRM	70.37	▲1.98	76.95
+ POM	71.63	▲2.64	78.47
+ BEL	74.17	▲5.18	80.47
+ PRM + BEL	74.81	▲5.92	81.52
+ POM + BEL	74.62	▲5.63	80.98
+ PRM + POM + BEL	<b>75.94</b>	<b>▲6.95</b>	<b>82.62</b>

Table 3: The different locations apply the mask in the proposed method. PRM = Pre-Mask filtering, BEL = Body Loss, POM = Post-Mask filtering. IMP. refers to the improvement achieved by the student network.

**Compare edge filter locations.** From Table 3, the numer-

Edge Width	mIoU	IMP.	mAcc
3	73.89	▲4.90	81.24
5	74.26	▲5.27	82.00
7	<b>75.94</b>	<b>▲6.95</b>	<b>82.62</b>
10	74.70	▲5.71	82.13
15	74.36	▲5.37	82.02

Table 4: Performance comparison for various different Edge Width with PSPNet-R18 on the Cityscapes validation set, for a fair comparison, we rerun the same setting 3 times and measured mean mIoU for evaluation.

ical results demonstrated the effectiveness of our proposed method. We further analysed the impact of applying the edge filter for different locations. Applying Per Mask filter, the performance slightly improved by 1.98% compared to the student without distillation. In contrast, Post Mask filtering improves the performance by 2.64%, because POM extracts edge discrepancy specifically for each type of class. The Body Enhancement Module takes care of non-edge information during our distillation setting, and the performance is raised by 5.18%. Afterwards, we explore the performance by applying PRM or POM to BEL. The combination of three terms archives best mIoU that 75.94% on the Cityscapes validation set.

**The Impact of different edge widths.** Edge width is an important hyper-parameter for Edge Detection Module, in which a larger edge unit produces a wider edge. Where a wider edge incorporates more pixels into the Edge Delegation Loss terms, at the same time, the larger amount of pixels will be from the body region. We analyze the effects of edge width size on the distillation performance in Table 4. We find that setting edge width as 7 units produces the optimal performance with 6.95% improvement.

**CAM Visualization Analysis.** Figure 5 illustrates the explicit refinement of semantic boundaries on multiple objects. The shape constraints are evident, such as the strong attention given to the pillow outlines. Despite BPKD distillation, the student network cannot perfectly segment the horse in the second row, but it has highly attended to the bone silhouette. This shows that BPKD has tried its best to distill knowledge to the student network, but due to its limited size and capacity, the student can only learn the surface-level capacity of the teacher. The body sections have been smoothly affiliated to a single category, reducing high-uncertainty edge and incorporating shape prior knowledge from the edge loss pressure. More qualitative segmentation results in supplement visually demonstrate our BPKD’s effectiveness for both tiny and large objects with explicit boundaries enhancement.

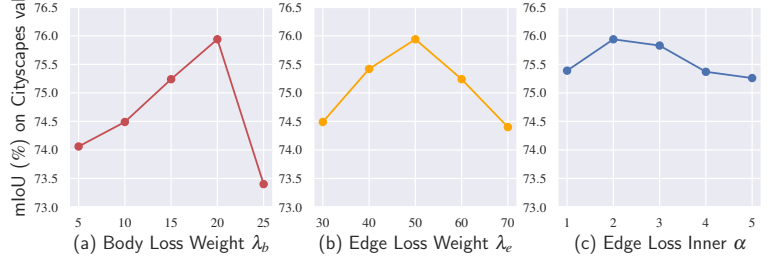


Figure 4: Impact of the (a) Body Loss weight  $\lambda_b$  and (b) Edge Loss weight  $\lambda_e$  and (c) Edge Loss Inner weight  $\alpha$  on Cityscapes val. We found optimal combination by board range study that  $\lambda_b = 20, \lambda_e = 50, \alpha = 2$ .

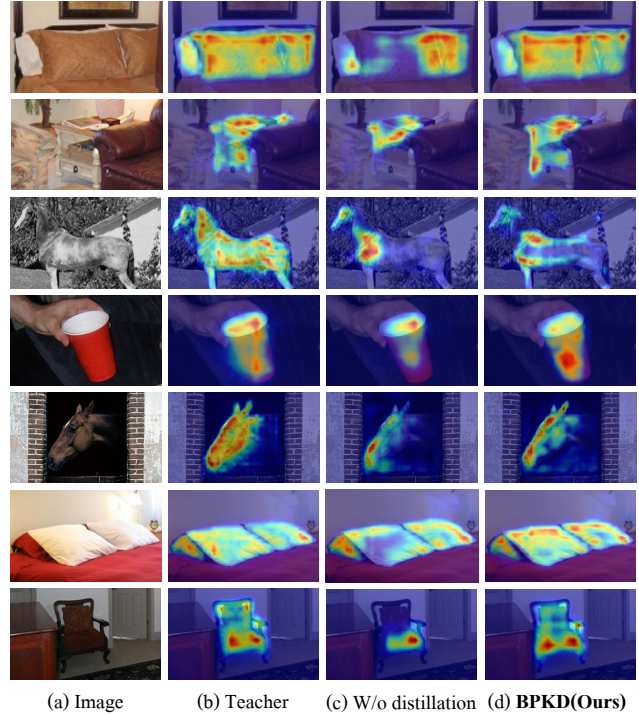


Figure 5: Comparison of CAM visualizations between (b) the teacher model, (c) the student model without distillation, and (d) the BPKD model. Activation maps were extracted from the last block of the corresponding ResNet backbones using HiResCAM [10]. The results indicate that BPKD shows superior refinement of boundaries and higher attention to semantic bodies. For better visualization, zoom in is recommended. The figure is best viewed in color.

## 5. Conclusion

This work presents a novel boundary-privileged knowledge distillation (BPKD) method for semantic segmentation, which transfers the cumbersome body and edge knowledge of the teacher model to the compact student model separately. Extensive experiments demonstrate the impor-



tance of considering the knowledge representation in the body and edge regions differently. The edge region requires focus on distinguishing between uncertain classes for each pixel, while the body region needs to prioritize localizing and connecting object structures. Experimental results consistently show that the proposed distillation method outperforms state-of-the-art methods on various public benchmark datasets. It significantly improves the overall mIoU and achieves remarkable performance in the edge region.

## 6. Acknowledgments

Y. Liu acknowledges the support of start-up funding from The University of Adelaide for their participation in this work. We express our gratitude to The University of Adelaide High-Performance Computing Services for providing the GPU Compute Resources, and to Mr. Wang Hui and Dr. Fabien Voisin for their valuable technical support. We would like to thank Mr. Hanwen Wang for providing the visualization collection used in the supplementary.

### A. Performance On Reduced Schedules

In order to demonstrate the efficacy of leveraging soft-labels from teachers to accelerate convergence speed, we conducted a comparison experiment with reduced schedules, 40k iteration training, and a  $512 \times 512$  resolution. The results of this experiment, as well as comparisons with other state-of-the-art algorithms, are reported in Table 6. To ensure a fair comparison, our proposed knowledge distillation framework was applied to different teachers. Our students were able to learn knowledge from the teacher network, resulting in significant performance gains and achieving state-of-the-art results on multiple datasets. As demonstrated in Table 6, on ADE20K our method outperformed strong baseline channel-wise distillation by 2.34%, 1.46%, 1.66%, 1.57%, and 1.59%, respectively. Furthermore, our methods were able to enhance the lightweight student network without increasing computational capacity, improving performance by 6.45%, 4.13%, 6.08%, and 8.49% compared to raw students. On Cityscape our method outperformed strong baseline channel-wise distillation by 2.22%, 1.35%, 4.40%, 1.52%, respectively. Furthermore, our methods were able to enhance the lightweight student network without increasing computational capacity, improving performance by 10.07%, 3.78%, 8.64%, and 1.75% compared to raw students. On Pascal Context our method outperformed strong baseline channel-wise distillation by 2.34%, 1.46%, 1.66%, 1.57%, and 1.59%, respectively. Furthermore, our methods were able to enhance the lightweight student network without increasing computational capacity, improving performance by 10.07%, 3.78%, 6.08%, and 8.49% compared to raw students. These numerical results suggest that our method is not dependent on a specific model struc-

ture, and that it produces significant performance gains with the pure student network, even without ImageNet Pre-train shows in table 5. To further demonstrate the effectiveness of our method, qualitative segmentation results are visualized in Figure 9.

Table 5: Performance Comparison with state-of-the-art distillation methods on ADE20K dataset, the student backbone is not pre-trained on ImageNet.

Methods	mIoU	mAcc(%)
T:PSPNet-R101 [59]	44.39	54.74
S:PSPnet-R18 [59]	17.11	22.99
SKDS [29]	20.79	27.74
IFVD [17]	20.75	27.6
CIRKD [51]	22.90	30.68
CWD [41]	24.79	31.44
<b>BPKD(Ours)</b>	<b>27.46</b>	<b>36.10</b>

### B. Implementation Details of Edge Loss

This section presents the implementation details of the Edge Loss, aimed at facilitating reproducibility. Our implementation consists of well-annotated components that are incorporated into our distillation system. To begin with, the Pre Mask Filter operation is applied to the logits obtained from both teachers and students. Subsequently, a dimensional rearrangement is performed to optimize the process. The KL divergence serves as the core component to estimate the distance between the student and teacher probability distributions, which establishes an embryonic reference for calculating the loss. The Post Mask Filter takes input from the unreduced criterion and aggregates the distance on the channel dimension. It is followed by an additional spatial expansion that repeats the spatial information based on the given channels. Finally, the inner weights vector is applied to corresponding categories to enhance the learning capacity for hard edge samples. The EDM loss terms are then finalized by overall weights adaption and average reduction. Furthermore, the Edge Detection Module is another crucial sub-component that employs multiple-level edge masks by providing input and ground truth labels. We utilize dilation and erosion to retrieve the edges, with the hyperparameter, *kernel\_size*, controlling the edge width. To address the computational pressure arising from a large kernel, the *compute\_iters* is introduced for GPU memory optimization. Additionally, the edge detection module employs average pooling as a downsampling policy, considering progressively decreasing importance from internal boundaries to outlines.

Table 6: Performance comparison of different distillation methods with state-of-the-art techniques in **Reduced schedules**. We set a reduced training setting to reduce the computational cost, including crop size reduced to  $512 \times 512$ , and training schedule to 40k iterations. We test these methods on various segmentation networks for both student and teacher models, using datasets including Cityscapes[9], ADE20K [61], and Pascal Context [11]. The FLOPs are obtained on  $512 \times 512$  resolutions. Our BPKD outperforms all previous methods in a large margins across multiple datasets and network architectures.

Methods	FLOPs(G)	Parameters(M)	ADE20K 40k 512*512		Cityscapes 40k 512*512		Pascal Context 59 40k 512*512	
			mIoU(%)	mAcc(%)	mIoU	mAcc(%)	mIoU	mAcc(%)
T: PSPNet-R101	256.89	68.07	44.39	54.75	79.74	86.56	52.47	63.15
S:PSPnet-R18	54.53	12.82	29.42	38.48	68.99	75.19	43.07	53.79
SKDS	54.53	12.82	31.80	42.25	69.33	75.37	43.93	54.01
IFVD	54.53	12.82	32.15	42.53	71.08	77.46	44.75	54.99
CIRKD	54.53	12.82	32.25	43.02	72.23	78.79	44.83	55.3
CWD	54.53	12.82	33.53	41.71	74.29	80.95	45.92	55.50
<b>BPKD(Ours)</b>	54.53	12.82	<b>35.87</b>	<b>45.42</b>	<b>75.94</b>	<b>82.62</b>	<b>47.16</b>	<b>57.61</b>
T: HRNetV2P-W48	95.64	65.95	42.02	53.52	80.65	87.39	51.12	61.39
S:HRNetV2P-W18S	10.49	3.97	28.69	37.86	73.77	82.89	40.82	51.70
SKDS	10.49	3.97	30.49	40.19	74.75	83.23	42.91	53.63
IFVD	10.49	3.97	30.57	40.42	75.33	83.83	43.12	54.03
CIRKD	10.49	3.97	31.34	41.45	74.63	83.72	43.45	54.10
CWD	10.49	3.97	31.36	40.24	75.54	84.08	45.50	56.01
<b>BPKD(Ours)</b>	10.49	3.97	<b>32.82</b>	<b>43.49</b>	<b>76.56</b>	<b>85.34</b>	<b>46.12</b>	<b>57.63</b>
T: DeeplabV3P-R101	255.67	62.68	45.47	56.41	80.98	88.7	53.2	64.04
S: DeeplabV3P+MV2	69.6	15.35	22.38	31.71	70.49	80.11	37.16	49.1
SKDS	69.6	15.35	24.65	35.07	70.81	79.31	39.18	51.13
IFVD	69.6	15.35	24.53	35.13	71.82	80.88	38.8	50.79
CIRKD	69.6	15.35	25.21	36.17	72.39	81.84	39.99	52.66
CWD	69.6	15.35	26.89	35.79	73.35	82.41	42.52	53.24
<b>BPKD(Ours)</b>	69.6	15.35	<b>28.46</b>	<b>41.45</b>	<b>76.58</b>	<b>84.14</b>	<b>44.32</b>	<b>56.04</b>
T: ISANet-R101	228.21	56.8	43.8	54.39	80.61	88.29	52.94	63.52
S: ISANet-R18	54.33	12.46	27.68	36.92	71.45	78.65	41.08	50.62
SKDS	54.33	12.46	28.70	38.51	70.65	77.53	42.87	52.89
IFVD	54.33	12.46	29.66	38.80	70.30	77.79	43.19	53.46
CIRKD	54.33	12.46	29.79	40.48	72.00	79.32	43.49	53.89
CWD	54.33	12.46	34.58	43.04	71.61	80.02	44.63	55.01
<b>BPKD(Ours)</b>	54.33	12.46	<b>36.17</b>	<b>45.26</b>	<b>72.72</b>	<b>81.50</b>	<b>45.50</b>	<b>56.55</b>

Class	road	sidewalk	building	wall	fence	pole	light	sign	vege.	terrain
Baseline	80.88	62.44	68.72	24.82	35.91	59.73	58.39	63.58	69.64	42.58
CWD	81.84	65.36	70.14	29.22	37.98	60.72	60.43	65.12	70.91	45.28
<b>BPKD(Ours)</b>	86.29	70.73	76.45	35.13	42.12	63.96	66.25	71.09	77.65	49.87
Class	sky	person	rider	car	truck	bus	train	moto.	bicycle	average
Baseline	77.72	65.94	48.45	74.55	33.63	50.58	34.58	39.83	59.56	55.34
CWD	79.2	67.58	49.11	75.76	36.53	51.37	38.29	43.82	60.75	57.34
<b>BPKD(Ours)</b>	84.87	73.56	54.4	82.28	44.49	57.27	41.27	50.14	67.37	62.91

Table 7: Illustration of our proposed **Boundary Privilege Knowledge Distillation** schemes in terms of class **Trimap IoU** metrics with PSPnet + Resnet18 network architecture over the Cityscapes validation set.

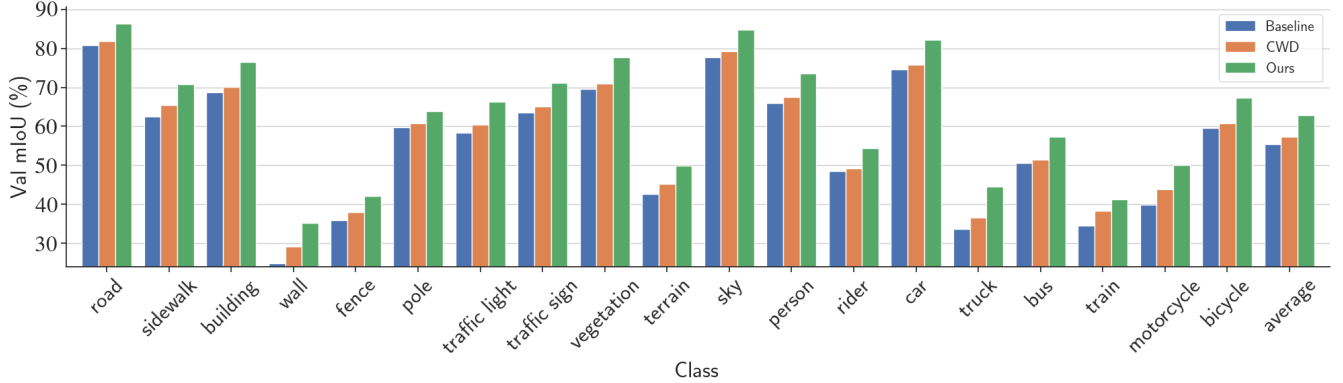


Figure 6: Illustration of our proposed **Boundary Privilege Knowledge Distillation** schemes in terms of class **Trimap IoU** metrics with PSPnet + Resnet18 network architecture over the Cityscapes validation set. It can be seen from the figure that our method has different degrees of improvement for all categories, meanwhile, we have a significant improvement for categories that are difficult to distinguish by boundaries.

### C. Categorical trimap Performance

This section presents an evaluation of the categorical Edge IoU using the PSPnet encoder and Resnet18 backbone architecture on the Cityscape validation set. Figure 6 and Table 7 illustrate the results. Our proposed methods demonstrated significant improvement in most classes based on edge factors compared to the raw student model and strong baseline channel-wise distillation method. We also observed explicit improvement for categories that are difficult to distinguish by boundaries. For instance, we achieved a 10.31% improvement for the wall category and a 6.69% improvement for the bus category.

### D. Instance segmentation

Methods	$AP$	$AP_{50}$	$AP_{75}$	$AP_S$	$AP_M$	$AP_L$
T:SOLOv2-X101	41.7	63.2	45.1	18.0	45.0	61.6
S:SOLOv2-R18*	26.7	44.1	27.5	6.50	27.2	45.8
CWD	28.4	47.1	29.6	9.90	30.3	44.2
<b>BPKD(Ours)</b>	<b>33.2</b>	<b>53.6</b>	<b>35.0</b>	<b>13.4</b>	<b>36.0</b>	<b>50.6</b>

Table 8: The instance segmentation results presented in this report were obtained on the COCO[27] validation set using single-model results. All distillation methods and the student network baseline were trained using a 1x schedule with multiple scale training disabled. The table below demonstrates that our distillation method can easily adapt to instance segmentation tasks and outperforms previous methods in a small-scale training setting.

We conducted experiments using SOLOv2 on the COCO dataset to demonstrate the general adaptability of our method. Specifically, we selected SOLOV2 [46] X-101(DCN) as the teacher and Light SOLOV2 [46] R-18

as the student. Table 8 presents the results, which show that our method improved the raw student by 6.5%, 9.5%, 7.5%, 6.9%, 8.8%, and 4.8% on the corresponding metrics. The results demonstrate that our distillation method can easily adapt to instance segmentation tasks and outperforms previous methods in the small-scale training setting. While instance segmentation and semantic segmentation tasks have distinct differences, they share similar properties in that they predict masks for target senses and given pixel-level annotations. During the experiments, we applied knowledge distillation methods multiple times on pyramid classification logits and masked representations, followed by calculating the average across all levels of sub-terms to obtain the distillation loss. From the numerical results, AP metrics explicitly increased for small and medium objects. However, there was no performance improvement for large objects, indicating that our method has room for improvement in instance segmentation tasks. As a future prospect, we aim to adjust the current method and design a specialized loss term for instance-level knowledge distillation.

### E. Effectiveness on Vision Transformers based methods

Table 9: mIoU and Trimap of Swin Transformers [31] and DeiT [43] with ViT Adapter (DeiT-Ada) [6] on ADE20K with UPerNet [48] decoder. Distil. forward speed (DFS.), training time (TT.). Costs estimated on DeiT-Adapter.

	DFS.(S)↑	TT.(H)↓	VMem(G)↓	Swin↑	Trimap↑	DeiT-Ada.↑	Trimap↑
Teacher-Base	9.52	11.26	8.32	50.13	40.1	48.80	39.72
Student-Tiny	12.8	8.44	3.87	43.57	32.78	41.10	32.15
SKDS	8.72	11.36	4.45	43.58	33.04	41.90	32.25
IFVD	6.06	16.45	8.97	43.75	32.90	41.16	32.25
CIRKD	7.70	16.35	10.7	43.32	32.68	41.64	32.23
CWD	8.76	11.15	4.45	44.99	33.73	44.25	33.49
<b>BPKD</b>	<b>7.84</b>	<b>13.49</b>	<b>5.49</b>	<b>46.13</b>	<b>38.11</b>	<b>45.25</b>	<b>37.05</b>

The results presented in Table 9 provide valuable insights into the performance of various knowledge distillation techniques applied to Swin Transformers and DeiT-Ada models on the ADE20K dataset. In terms of trimap scores, which measure boundary localization quality, we find that both Swin Transformers and DeiT-Ada models exhibit similar performance trends across different knowledge distillation techniques. Among the techniques, BPKD achieves the highest trimap scores for both model architectures, indicating its effectiveness in capturing fine-grained details and preserving sharp object boundaries. When comparing the distillation forward speed (DFS) of the different techniques, we notice that IFVD has the fastest DFS while CIRKD has the slowest. This suggests that IFVD may be more efficient in terms of inference time, although its mIoU and trimap scores are not the highest among the techniques. On the other hand, CIRKD consumes more time during inference but does not necessarily offer better segmentation or boundary localization performance. Regarding training time (TT), Student-Tiny has the shortest training time among all techniques, making it a potentially attractive option for resource-constrained settings. However, it should be noted that its mIoU and trimap scores are also lower compared to other techniques such as BPKD, which has longer training times but yields better performance. Finally, we can observe that the GPU memory consumption (VMem) varies significantly among the techniques. Student-Tiny has the lowest VMem, which may be appealing in situations with limited GPU memory availability. However, as mentioned earlier, its performance is not on par with more resource-intensive techniques such as BPKD.

In conclusion, the results presented in Table 9 highlight the trade-offs between model performance and computational resources when applying different knowledge distillation techniques to Swin Transformers and DeiT-Ada models on the ADE20K dataset. BPKD appears to be the best performing technique in terms of mIoU and trimap scores, although at the cost of higher training time and GPU memory consumption. Other techniques, such as IFVD and Student-Tiny, offer faster inference and shorter training times but may not provide the same level of segmentation and boundary localization quality.

## References

- [1] Weihao Cao, Yifan Zhang, Jianfei Gao, Anda Cheng, Ke Cheng, and Jian Cheng. Pkd: General distillation framework for object detectors via pearson correlation coefficient. *arXiv preprint arXiv:2207.02039*, 2022. 2, 3
- [2] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. *arXiv preprint arXiv:1412.7062*, 2014. 3
- [3] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017. 3
- [4] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*, 2017. 3, 6
- [5] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 801–818, 2018. 3, 6, 7
- [6] Zhe Chen, Yuchen Duan, Wenhao Wang, Junjun He, Tong Lu, Jifeng Dai, and Yu Qiao. Vision transformer adapter for dense predictions. *arXiv preprint arXiv:2205.08534*, 2022. 11
- [7] MMSegmentation Contributors. MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark. <https://github.com/open-mmlab/mms Segmentation>, 2020. 6
- [8] MMRazor Contributors. Openmmlab model compression toolbox and benchmark. <https://github.com/open-mmlab/mmrazor>, 2021. 6
- [9] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 2, 5, 6, 10, 16
- [10] Rachel Lea Draelos and Lawrence Carin. Hirescam: Faithful location representation in visual attention for explainable 3d medical image classification. *arXiv preprint arXiv:2011.08891*, 2020. 8
- [11] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 2, 5, 6, 10
- [12] Mingyuan Fan, Shenqi Lai, Junshi Huang, Xiaoming Wei, Zhenhua Chai, Junfeng Luo, and Xiaolin Wei. Rethinking bisenet for real-time semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9716–9725, 2021. 2
- [13] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021. 3



- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 1
- [15] Byeongho Heo, Minsik Lee, Sangdoo Yun, and Jin Young Choi. Knowledge transfer via distillation of activation boundaries formed by hidden neurons. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3779–3787, 2019. 2, 3
- [16] Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. Distilling the knowledge in a neural network. 2015. 2, 3
- [17] Yuenan Hou, Zheng Ma, Chunxiao Liu, Tak-Wai Hui, and Chen Change Loy. Inter-region affinity distillation for road marking segmentation. pages 12486–12495, 2020. 6, 9
- [18] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1314–1324, 2019. 1, 3
- [19] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 3
- [20] Lang Huang, Yuhui Yuan, Jianyuan Guo, Chao Zhang, Xilin Chen, and Jingdong Wang. Interlaced sparse self-attention for semantic segmentation. *arXiv preprint arXiv:1907.12273*, 2019. 6
- [21] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and 0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016. 3
- [22] Deyi Ji, Haoran Wang, Mingyuan Tao, Jianqiang Huang, Xian-Sheng Hua, and Hongtao Lu. Structural and statistical texture knowledge distillation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16876–16885, 2022. 2
- [23] Jangho Kim, SeongUk Park, and Nojun Kwak. Paraphrasing complex network: Network compression via factor transfer. *Advances in neural information processing systems*, 31, 2018. 2, 3
- [24] Pushmeet Kohli, Philip HS Torr, et al. Robust higher order potentials for enforcing label consistency. *International Journal of Computer Vision*, 82(3):302–324, 2009. 7
- [25] Philipp Krähenbühl and Vladlen Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. *Advances in neural information processing systems*, 24, 2011. 7
- [26] Seung Hyun Lee, Dae Ha Kim, and Byung Cheol Song. Self-supervised knowledge distillation using singular value decomposition. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 335–350, 2018. 2, 3
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 11
- [28] Akide Liu and Zihan Wang. Cv 3315 is all you need: Semantic segmentation competition. *arXiv preprint arXiv:2206.12571*, 2022. 3
- [29] Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang. Structured knowledge distillation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2604–2613, 2019. 2, 3, 6, 7, 9
- [30] Yifan Liu, Changyong Shu, Jingdong Wang, and Chunhua Shen. Structured knowledge distillation for dense prediction. *IEEE transactions on pattern analysis and machine intelligence*, 2020. 3
- [31] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. 11
- [32] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11976–11986, 2022. 1
- [33] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015. 3
- [34] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. Shufflenet v2: Practical guidelines for efficient cnn architecture design. In *Proceedings of the European conference on computer vision (ECCV)*, pages 116–131, 2018. 1
- [35] Sachin Mehta, Mohammad Rastegari, Anat Caspi, Linda Shapiro, and Hannaneh Hajishirzi. Espnet: Efficient spatial pyramid of dilated convolutions for semantic segmentation. In *Proceedings of the european conference on computer vision (ECCV)*, pages 552–568, 2018. 3
- [36] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 17
- [37] Adam Paszke, Abhishek Chaurasia, Sangpil Kim, and Eugenio Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. *arXiv preprint arXiv:1606.02147*, 2016. 3
- [38] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014. 3
- [39] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 3, 6
- [40] Alexander G Schwing and Raquel Urtasun. Fully connected deep structured networks. *arXiv preprint arXiv:1503.02351*, 2015. 3

- [41] Changyong Shu, Yifan Liu, Jianfei Gao, Zheng Yan, and Chunhua Shen. Channel-wise knowledge distillation for dense prediction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5311–5320, 2021. [2](#), [3](#), [5](#), [6](#), [7](#), [9](#), [15](#), [16](#), [17](#)
- [42] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. [1](#)
- [43] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pages 10347–10357. PMLR, 2021. [11](#)
- [44] Jue Wang, Michael F Cohen, et al. Image and video matting: a survey. *Foundations and Trends® in Computer Graphics and Vision*, 3(2):97–175, 2008. [4](#)
- [45] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 43(10):3349–3364, 2020. [1](#), [3](#), [6](#)
- [46] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *Advances in Neural information processing systems*, 33:17721–17732, 2020. [11](#)
- [47] Yukang Wang, Wei Zhou, Tao Jiang, Xiang Bai, and Yongchao Xu. Intra-class feature variation distillation for semantic segmentation. In *European Conference on Computer Vision*, pages 346–362. Springer, 2020. [2](#), [3](#), [6](#)
- [48] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing for scene understanding. In *Proceedings of the European conference on computer vision (ECCV)*, pages 418–434, 2018. [11](#)
- [49] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. [1](#)
- [50] Danna Xue, Fei Yang, Pei Wang, Luis Herranz, Jinqiu Sun, Yu Zhu, and Yanning Zhang. Slimseg: Slimmable semantic segmentation with boundary supervision. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 6539–6548, 2022. [2](#)
- [51] Chuanguang Yang, Helong Zhou, Zhulin An, Xue Jiang, Yongjun Xu, and Qian Zhang. Cross-image relational knowledge distillation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12319–12328, 2022. [2](#), [3](#), [6](#), [9](#)
- [52] Zhendong Yang, Zhe Li, Mingqi Shao, Dachuan Shi, Zehuan Yuan, and Chun Yuan. Masked generative distillation. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XI*, pages 53–69. Springer, 2022. [2](#), [3](#)
- [53] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4133–4141, 2017. [2](#), [3](#)
- [54] Yuhui Yuan, Xiaokang Chen, Xilin Chen, and Jingdong Wang. Segmentation transformer: Object-contextual representations for semantic segmentation. *arXiv preprint arXiv:1909.11065*, 2019. [1](#)
- [55] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016. [2](#), [3](#)
- [56] Bowen Zhang, Zhi Tian, Quan Tang, Xiangxiang Chu, Xiaolin Wei, Chunhua Shen, and Yifan Liu. Segvit: Semantic segmentation with plain vision transformers. *arXiv preprint arXiv:2210.05844*, 2022. [3](#)
- [57] Chenrui Zhang and Yuxin Peng. Better and faster: knowledge transfer from multiple self-supervised learning tasks via graph distillation for video classification. *arXiv preprint arXiv:1804.10069*, 2018. [2](#), [3](#)
- [58] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Haibin Lin, Zhi Zhang, Yue Sun, Tong He, Jonas Mueller, R Manmatha, et al. Resnest: Split-attention networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2736–2746, 2022. [1](#)
- [59] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017. [3](#), [6](#), [9](#), [15](#), [16](#), [17](#)
- [60] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip HS Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1529–1537, 2015. [3](#)
- [61] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 633–641, 2017. [2](#), [5](#), [6](#), [10](#), [15](#)



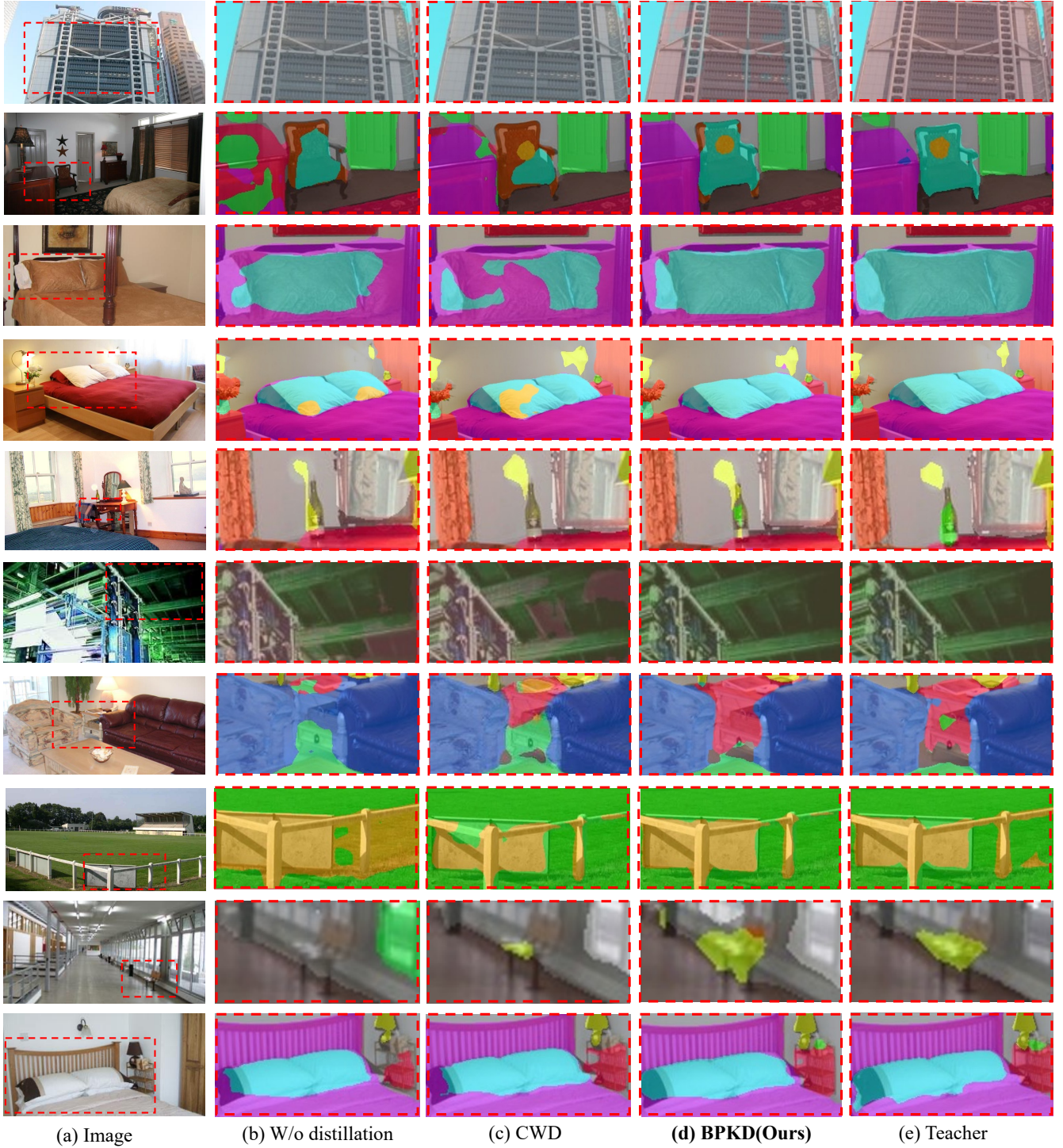


Figure 7: Qualitative results on the ADE20K [61] validation set produced by PSPNet [59] and ResNet18 network architecture: (a) Initial images, (b) w/o distillation scheme, (c) state-of-the-art method channel-wise distillation [41], (d) **BPKD** our method, (e) teacher. This figure shows that our methods segment the small complex objects with explicit boundaries.





Figure 8: Qualitative results on the Cityscapes [9] validation set produced by PSPNet [59] and ResNet18 network architecture: (a) Initial images, (b) w/o distillation, (c) state-of-the-art method channel-wise distillation [41], (d) **BPKD** our method, (e) Ground truth. This figure shows that our methods segment the small complex objects with explicit boundaries. Zoom in for a better view.



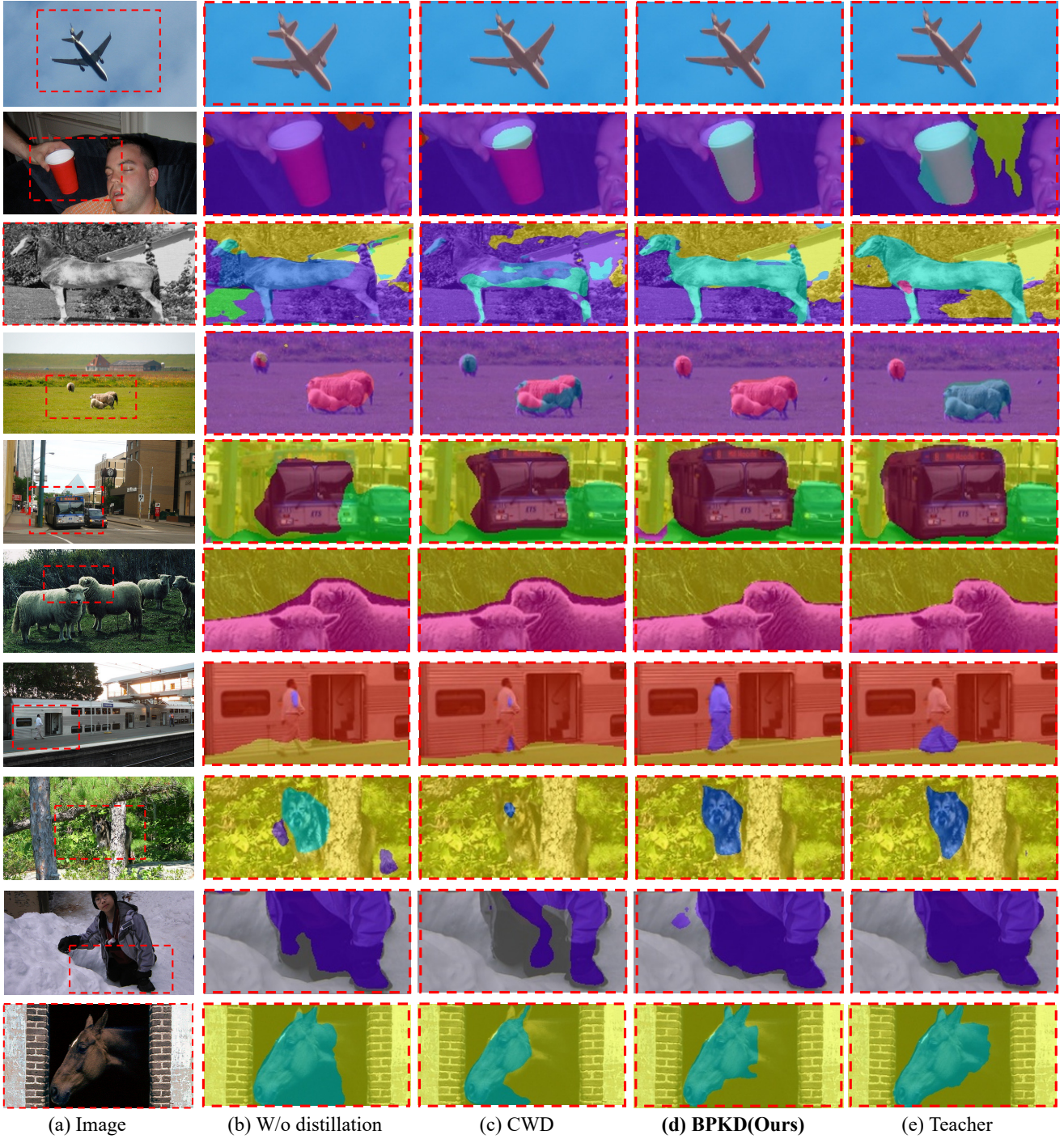


Figure 9: Qualitative results on the Pascal-Context [36] validation set produced by PSPNet [59] and ResNet18 network architecture: (a) Initial images, (b) w/o distillation scheme, (c) state-of-the-art method channel-wise distillation [41], (d) **BPKD** our method, (e) teacher. This figure shows that our methods segment the small complex objects with explicit boundaries.