

---

# SAPA: Similarity-Aware Point Affiliation for Feature Upsampling

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We introduce point affiliation into feature upsampling, a notion that describes the  
2 affiliation of each upsampled point to a semantic cluster formed by local decoder  
3 feature points with semantic similarity. By rethinking point affiliation, we present a  
4 generic formulation for generating upsampling kernels. The kernels encourage not  
5 only semantic smoothness but also boundary sharpness in the upsampled feature  
6 maps. Such properties are particularly useful for some dense prediction tasks such  
7 as semantic segmentation. The key idea of our formulation is to generate similarity-  
8 aware kernels by comparing the similarity between each encoder feature point and  
9 the spatially associated local region of decoder features. In this way, the encoder  
10 feature point can function as a cue to inform the semantic cluster of upsampled  
11 feature points. To embody the formulation, we further instantiate a lightweight  
12 upsampling operator, termed Similarity-Aware Point Affiliation (SAPA), and inves-  
13 tigate its variants. SAPA invites consistent performance improvements on a number  
14 of dense prediction tasks, including semantic segmentation, object detection, depth  
15 estimation, and image matting. Code will be available online.

## 16 1 Introduction

17 We introduce the notion of point affiliation into feature upsampling. Point affiliation defines a relation  
18 between each upsampled point and a *semantic cluster*<sup>1</sup> to which the point should belong. It highlights  
19 the spatial arrangement of upsampled points at the semantic level. Considering an example shown  
20 in Fig. 1 w.r.t.  $\times 2$  upsampling in semantic segmentation, the orange point of low resolution will  
21 correspond to 4 upsampled points of high resolution, in which the red and yellow ones should be  
22 assigned the ‘picture’ cluster and the ‘wall’ cluster, respectively. Designating point affiliation is  
23 difficult and sometimes can be erroneous, however.

24 In  $\times 2$  upsampling, nearest neighbor (NN) interpolation directly copies 4 identical points from the  
25 low-res one, which assigns the same semantic cluster to the 4 points. On regions in need of details,  
26 the 4 points probably do not share the same cluster but are forced to share. Bilinear interpolation  
27 assigns point affiliation with distance priors. Yet, when tackling points of different semantic clusters,  
28 it not only cannot inform clear point affiliation, but also blurs the boundary between different  
29 semantic clusters. Recent dynamic upsampling operators have similar issues. CARAFE [1] judges the  
30 affiliation of an upsampled point with content-aware kernels. Certain semantic clusters will receive  
31 larger weights than the rest and therefore dominate the affiliation of upsampled points. However, the  
32 affiliation near boundaries or on regions full of details can still be ambiguous. As shown in Fig. 2, the  
33 boundaries are unclear in the feature maps upsampled by CARAFE. The reason is that the kernels  
34 are conditioned on decoder features alone; the decoder features carry little useful information about  
35 high-res structure.

---

<sup>1</sup>A semantic cluster is formed by local decoder feature points with the similar semantic meaning.

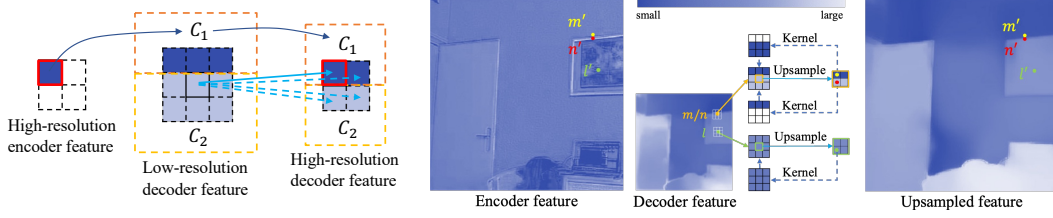


Figure 1: **Left: Similarity between an encoder point and different semantic clusters in the decoder.** **Right: Point affiliation mechanism of ideal upsampling kernels.** Left: If the red-box encoder feature point is classified into the semantic cluster  $C_1$ , then it is more similar to  $C_1$  than  $C_2$ . Right: The upsampling kernels can be a ‘soft’ selector in a local window to assign point affiliation. For an upsampled point, the kernel selects a/some representative points from its most relative semantic cluster. *E.g.*, according to the encoder feature, the red upsampled feature point should belong to the ‘picture’ cluster. Then we expect the kernel can assign large weights on picture-related points and small weights on wall-related points. In this way, after the weighted sum, the upsampled point will be revalued and assigned the ‘picture’ cluster.

36 Inferring structure requires high-res encoder features. For instance, if the orange point in Fig. 1 lies  
 37 on the low-res boundary, it is difficult to judge to which cluster the 4 upsampled points should belong.  
 38 However, the encoder feature in Fig. 1 actually tells that, the yellow point is on the wall, and the red  
 39 one is on the picture, which suggests one may extract useful information from the encoder feature to  
 40 assist the designation of point affiliation. Indeed IndexNet [2] and A2U [3] have attempted to encode  
 41 such information to improve detail delineation in encoder-dependent upsampling kernels; however,  
 42 the encoder feature can easily introduce noise into kernels, engendering discontinuous feature maps  
 43 shown in Fig. 2. Hence, the key problem seems to be how to extract only required information into  
 44 the upsampling kernels from the encoder feature while filtering out the rest.

45 To use encoder features effectively, an important assumption of this paper is that, *an encoder feature*  
 46 *point is most similar to the semantic cluster into which the point will be classified.* Per the left of  
 47 Fig. 1, suppose that the encoder point in the red box is assigned into the cluster  $C_1$  by its semantic  
 48 meaning, then it is similar to  $C_1$ , while not similar to  $C_2$ . As a result, by comparing the similarity  
 49 between the encoder feature point and different semantic clusters in the decoder feature, the affiliation  
 50 of the upsampled point can be informed according to the similarity scores. In particular, we propose  
 51 to generate upsampling kernels with local mutual similarity between encoder and decoder features.  
 52 For every encoder feature point, we compute the similarity score between this point and each decoder  
 53 feature point in the spatially associated local window. For the green point in Fig. 1, since every point  
 54 in the window shares the same semantic cluster, the encoder feature point is as similar as every point  
 55 in the window. In this case we expect an ‘average kernel’ which is the key characteristic to filter  
 56 noise, and the upsampled 4 points would have the same semantic cluster as before. For the yellow  
 57 point in the encoder, since it belongs to the ‘wall’ cluster, it is more similar to the points on the wall  
 58 than those on the picture. In this case we expect a kernel with larger weights on points related to the  
 59 ‘wall’ cluster. This can help to assign the affiliation of the yellow point to be in the ‘wall’ cluster.

60 By modeling the local mutual similarity, we derive a generic form of upsampling kernels and show  
 61 that this form implements our expected upsampling behaviors: encouraging both semantic smoothness  
 62 and boundary sharpness. Following our formulation, we further instantiate a lightweight upsampling  
 63 operator, termed Similarity-Aware Point Affiliation (SAPA), and investigate its variants. We evaluate  
 64 SAPA across a number of mainstream dense prediction tasks, for example: i) *semantic segmentation*:  
 65 we test SAPA on several transformer-based segmentation baselines on the ADE20K dataset [4], such  
 66 as SegFormer [5], MaskFormer [6], and Mask2Former [7], improving the baselines by 1% ~ 2.7%  
 67 mIoU; ii) *object detection*: SAPA improves the performance of Faster RCNN by 0.4% AP on MS  
 68 COCO [8]; iii) *monocular depth estimation*: SAPA reduces the rmse metric of BTS [9] from 0.419  
 69 to 0.408 on NYU Depth V2 [10]; and iv) *image matting*: SAPA outperforms a strong A2U matting  
 70 baseline [3] on the Adobe Composition-1k testing set [11] with a further 3.8% relative error reduction  
 71 in the SAD metric. SAPA also outperforms or at least is on par with other state-of-the-art dynamic  
 72 upsampling operators. Particularly, even without additional parameters, SAPA outperforms the  
 73 previous best upsampling operator CARAFE on semantic segmentation.

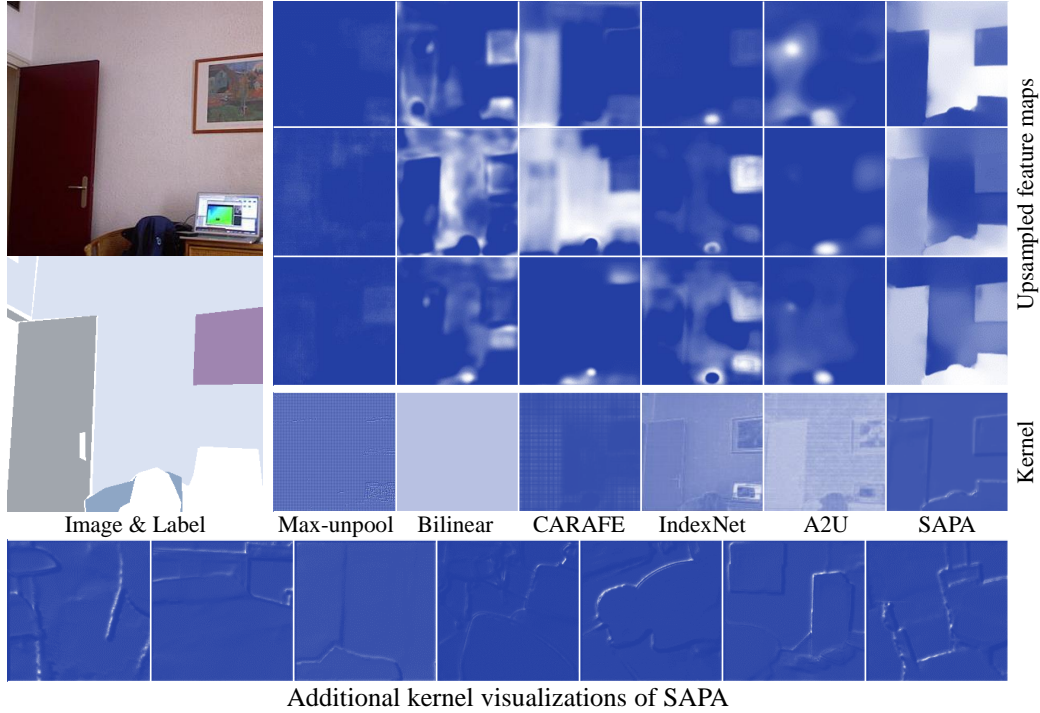


Figure 2: **Top: Upsampled feature maps and upsampling kernel maps of different upsampling operators. Down: Additional upsampling kernel maps of SAPA generated from various samples.** The visualization is produced with SegNet [12] as the baseline on the SUN RGBD [13] dataset. For each upsampling operator, we choose the first three channels from the feature maps of the last upsampling stage for visualization. Only SAPA shows both smooth regions and sharp boundaries. The kernel map of CARAFE is coarse and lacking in details, while IndexNet and A2U generate kernels with much details from the encoder. See the supplementary materials for additional visualizations.

## 2 Related Work

We review work related to feature upsampling. Feature upsampling is a fundamental procedure in encoder-decoder architectures used to recover the spatial resolution of low-res decoder feature maps and has been extensively used in dense prediction tasks such as semantic segmentation [12, 5, 6, 7] and depth estimation [14, 15, 9].

Standard upsampling operators are hand-crafted. NN and bilinear interpolation measure the semantic affiliation in terms of relative distances in upsampling, which follows fixed rules to designate point affiliation, even if the true affiliation may be different. Max unpooling [12] stores the indices of max-pooled feature points in encoder features and uses the sparse indices to guide upsampling. While it executes location-specific point affiliation which benefits detail recovery, most upsampled points are assigned with null affiliation due to zero filling. Pixel Shuffle [16] is widely-used in image/video super-resolution. Its upsampling only includes memory operation – reshaping depth channels to space ones. The notion of point affiliation does not apply to this operator, however.

Another stream of upsampling operators implement learning-based upsampling. Among them, transposed convolution or deconvolution [17] is known as an inverse convolutional operator. Based on a novel interpretation of deconvolution, PixelTCL [18] is proposed to alleviate the checkerboard artifacts [19] of standard deconvolution. In addition, bilinear additive upsampling [20] attempts to combine learnable convolutional kernels with hand-crafted upsampling operators to achieve composited upsampling. Recently, DUpsample [21] seeks to reconstruct the upsampled feature map with pre-learned projection matrices, expecting to achieve a data-dependent upsampling behavior. While these operators are learnable, the upsampling kernels are fixed once learned, still resulting in fixed designation of point affiliation.

In learning-based upsampling, some recent work introduces the idea of generating content-aware dynamic kernels. Instead of learning the parameters of the kernels, they learn how to predict the kernels. In particular, CARAFE [1] predicts dynamic kernels conditioned on the decoder features. IndexNet [2] and A2U [3], by contrast, generate encoder-dependent kernels. While they significantly outperform previous upsampling operators in various tasks, they still can cause uncertain point affiliation, resulting in either unclear predictions near boundaries or fragile predictions in regions.

Our work is closely related to dynamic kernel-based upsampling. We also seek to predict dynamic kernels; however, we aim to address the uncertain point affiliation in prior arts to achieve simultaneous region smoothness and boundary sharpness.

### 3 Dynamic Upsampling Revisited

We first revisit two key components shared by existing dynamic upsampling operators: kernel generation and feature assembly.

**Kernel Generation** Given the decoder feature  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$ , if we upsample it to the target feature  $\mathcal{X}' \in \mathbb{R}^{2H \times 2W \times C}$ , then for any point at the location  $l' = (i', j')$  in  $\mathcal{X}'$ , we generate a kernel  $\mathcal{W}_{l'}$  based on the neighborhood feature  $\mathcal{N}_{l'}$  that spatially corresponds to  $l'$ . In this way, kernel generation can be defined by

$$\mathcal{W}_{l'} = \text{norm}(\psi(\mathcal{N}_{l'})), \quad (1)$$

where  $\psi$  refers to a kernel generation module, which is often implemented by a sub-network used to predict the kernel conditioned on  $\mathcal{N}_{l'}$ , and  $\text{norm}$  is a normalization function. The feature  $\mathcal{N}_{l'}$  can originate from two sources, to be specific, from the encoder feature  $\mathcal{Y} \in \mathbb{R}^{2H \times 2W \times C}$  or from the decoder feature  $\mathcal{X}$ . If the encoder feature  $\mathcal{Y}$  is chosen as the source, then a local region of  $\mathcal{Y}$  centered at  $l'$  is extracted to be  $\mathcal{N}_{l'}$ . If the decoder feature  $\mathcal{X}$  is chosen, one first needs to compute the projective location of  $l'$ , i.e.,  $l = (i, j) = \lfloor \frac{l'}{2} \rfloor = (\lfloor \frac{i'}{2} \rfloor, \lfloor \frac{j'}{2} \rfloor)$  according to the spatial location correspondence, then a local region of  $\mathcal{X}$  centered at  $l$  is regarded as  $\mathcal{N}_{l'}$ . The  $\text{softmax}$  function is often used as the normalization function such that relevant points can be softly selected to compute the value of the target point using the weight  $\mathcal{W}_{l'}$ .

According to Fig. 2, the source of  $\mathcal{N}_{l'}$  can affect the predicted kernel. The kernels predicted by CARAFE, IndexNet, and A2U show significantly distinct characteristics. With the decoder feature alone, the kernel map is coarse and lacking in details. Benefiting from the encoder feature, the kernel maps generated by IndexNet and A2U have rich details; however, they manifest high similarity to the encoder feature, which means noise is introduced into the kernel.

**Feature Assembly** For each target feature point at  $l'$ , we assemble the corresponding sub-region of decoder feature with the predicted  $K \times K$  kernel  $\mathcal{W}_{l'}$ , whose weight is denoted by  $\mathcal{W}_{l',(u,v)}$ ,  $u, v = -r, \dots, r, r = \lfloor \frac{K}{2} \rfloor$ , to obtain the value of the target point  $\mathcal{X}'_{l'}$  by a weighted sum

$$\mathcal{X}'_{l'} = \sum_{u=-r}^r \sum_{v=-r}^r \mathcal{W}_{l',(u,v)} \cdot \mathcal{X}_{l+(u,v)}. \quad (2)$$

By executing feature assembly on every target feature point, we can obtain the target upsampled feature map. As shown in Fig. 2, the upsampled feature has a close relation to the kernel. A well-predicted kernel can encourage both semantic smoothness and boundary sharpness; a kernel without encoding details or with too many details encoded can introduce noise. We consider an ideal kernel should only response at the position in need of details, while do not response (appearing as an average value over an area) at good semantic regions. More importantly, an ideal kernel should assign weights reasonably so that each point can be designated to a correct semantic cluster.

### 4 Rethinking Point Affiliation with Local Mutual Similarity

To obtain an expected upsampling kernel mentioned above, we first derive a generic formulation for generating the upsampling kernel by exploiting local mutual similarity, then explain why the formulation encourages semantic smoothness and boundary sharpness, and finally present an upsampling operator, SAPA, that embodies our formulation.

#### 4.1 Local Mutual Similarity

We rethink point affiliation from the view of local mutual similarity between encoder and decoder features. With a detailed analysis, we explain why such similarity can assist point affiliation.

We first define a generic similarity function  $\text{sim}(\mathbf{x}, \mathbf{y}) : \mathbb{R}^C \times \mathbb{R}^C \rightarrow \mathbb{R}$ . It scores the similarity between a vector  $\mathbf{x}$  and a vector  $\mathbf{y}$  of the same dimension  $C$ . We also define a normalization function involving  $n$  real numbers  $x_1, x_2, \dots, x_n$  by  $\text{norm}(x_i : x_1, x_2, \dots, x_n) = \frac{h(x_i)}{\sum_{j=1}^n h(x_j)}$ , where  $h(x) : \mathbb{R} \rightarrow \mathbb{R}$  is an arbitrary function, ignoring zero division. Given  $\text{sim}(\mathbf{x}, \mathbf{y})$  and  $h(x)$ , we can define a generic formulation for generating the upsampling kernel

$$w = \frac{h(\text{sim}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{x}' \in N_{l'}} h(\text{sim}(\mathbf{x}', \mathbf{y}))}, \quad (3)$$

where  $w$  is the kernel value specific to  $\mathbf{x}$  and  $\mathbf{y}$ . To analyze the upsampling behavior of the kernel, we further define the following notations.

Let  $\mathbf{y}_{l'} \in \mathbb{R}^C$  denote the encoder feature vector at position  $l'$  and  $\mathbf{x}_l \in \mathbb{R}^C$  be the decoder feature vector at position  $l$ , where  $C$  is the number of channels. Our operation will be done within a local window of size  $K \times K$ , between each encoder vector  $\mathbf{y}_{l'}$  and all its spatially associated decoder feature vectors,  $\mathbf{x}_{l+(u,v)}$ 's, where  $u, v = -r, \dots, r, r = \lfloor \frac{K}{2} \rfloor$ .

To simplify our analysis, we also assume local smoothness. That is, points with the same semantic cluster will have a similar value, which means a local region will share the same value on every channel of the feature map. As shown in Fig. 1, we define  $\mathbf{x}_a = (R_1^a, R_2^a, \dots, R_C^a)^T$  and  $\mathbf{x}_p = (R_1^p, R_2^p, \dots, R_C^p)^T$  as the semantic clusters related to 'wall' and 'picture', respectively, where  $R_c^a$  and  $R_c^p$  are constants, and  $c = 1, 2, \dots, C$ . For ease of analysis, we define two types of windows distinguished by their contents. When all the points inside a window belong to the same semantic cluster, it is called a smooth window; while different semantic clusters appear in a window, it is defined as a detail window.

Next, we explain why the kernel can filter out noise, why it encourages semantic smoothness in a smooth window, and why it can help to recover details when dealing with boundaries/textures in a detail window.

**Upsampling in a Smooth Window** Without loss of generality, we consider an encoder point at the position  $l'$  in Fig. 1. Its corresponding window is a smooth window of the semantic cluster 'picture', thus  $\mathbf{x}_{l+(u,v)} = \mathbf{x}_p, u, v = -r, \dots, r, r = \lfloor \frac{K}{2} \rfloor$ . Then the upsampling kernel weight w.r.t. the upsampled point  $l'$  at the position  $(u, v)$  takes the form

$$\mathcal{W}_{l', (u,v)} = \text{norm}(\text{sim}(\mathbf{x}_{l+(u,v)}, \mathbf{y}_{l'})) = \frac{h(\text{sim}(\mathbf{x}_{l+(u,v)}, \mathbf{y}_{l'}))}{\sum_{s=-r}^r \sum_{t=-r}^r h(\text{sim}(\mathbf{x}_{l+(s,t)}, \mathbf{y}_{l'}))} = \frac{h(\text{sim}(\mathbf{x}_p, \mathbf{y}_{l'}))}{K^2 h(\text{sim}(\mathbf{x}_p, \mathbf{y}_{l'}))} = \frac{1}{K^2}, \quad (4)$$

which has nothing to do with  $l, u$ , and  $v$ . Eq. (4) reveals a key characteristic of local mutual similarity in a smooth window: the kernel weight is a constant regardless of  $\mathbf{y}$ . Therefore, the kernel fundamentally can filter out noise from encoder features with an 'average' kernel.

Note that, in the derivation above the necessary conditions include: i)  $\mathbf{x}$  is from a local window in the decoder feature map; ii) a normalization function in the form of  $\frac{h(x_i)}{\sum_j h(x_j)}$ .

**Upsampling in a Detail Window** Again we consider two encoder points at the position  $m'$  and  $n'$  in Fig. 1. Ideally  $\mathbf{y}_{m'}$  and  $\mathbf{y}_{n'}$  should be classified into the semantic cluster of 'wall' and 'picture', respectively. Taking  $\mathbf{y}_{m'}$  as an example, following our assumption, it is more similar to points of the 'wall' cluster rather than the 'picture' cluster. From Eq. (4), we can tell that  $\text{sim}(\mathbf{x}_{m+a}, \mathbf{y}_{m'})$  is larger than  $\text{sim}(\mathbf{x}_{m+p}, \mathbf{y}_{m'})$ , where  $a$  and  $p$  are the offsets in the window such that  $m+a$  and  $m+p$  are within the 'wall' and the 'picture' cluster, respectively. Therefore, after computing similarity scores and normalization, one can acquire a kernel with significantly larger weights on points with the semantic cluster of 'wall' than that of 'picture', i.e.,  $\mathcal{W}_{m',a} \gg \mathcal{W}_{m',p}$ . After applying the kernel to the corresponding window, the upsampled point at  $m'$  will be revalued and assigned to the semantic cluster of 'wall'. Similarly, the upsampled point at  $n'$  will be assigned into 'picture'.

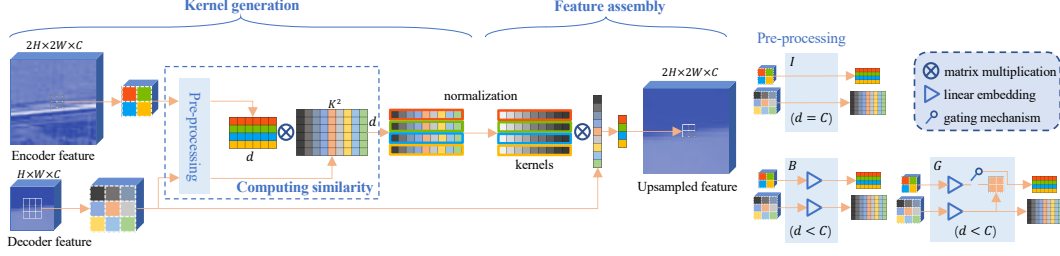


Figure 3: **Feature upsampling of SAPA.** We compute mutual similarity scores with a similarity function between each encoder feature point and the spatially associated local decoder feature points (features are normalized with LayerNorm before similarity computation). The scores are then transformed into upsampling kernel weights after kernel normalization. The kernel is then used to assemble the upsampled feature points. We illustrate three variants: SAPA-I with inner-product similarity, SAPA-B with (low-rank) bilinear similarity, and SAPA-G with gated (low-rank) bilinear similarity. They differ in the pre-processing step. See detailed definition in Section 4.2.

185 Note that, in Eq. (4) we have no constraint on  $\mathbf{y}$ . But here in a detail window,  $\mathbf{y}$  as an encoder  
 186 feature vector can play a vital role for designating correct point affiliation. Next in our concrete  
 187 implementation, we discuss how to appropriately use  $\mathbf{y}$  in the similarity function.

## 188 4.2 SAPA: Similarity-Aware Point Affiliation

189 Here we embody Eq. (3) by investigating different similarity and normalization functions.

190 **Normalization Function** Though we do not constrain  $h(x)$  in theory, in reality it must be carefully  
 191 chosen. For example, to avoid zero division, the scope for the choice of  $h(x)$  is narrowed. Following  
 192 existing practices [1, 2, 3], we choose  $h(x) = e^x$  by default, which is equivalent to softmax  
 193 normalization. We also test some other  $h(x)$ 's, such as  $h(x) = \text{ReLU}(x)$ ,  $h(x) = \text{sigmoid}(x)$ , and  
 194  $h(x) = \text{softplus}(x)$ . Their performance comparisons will be given in ablation studies.

195 **Similarity Function** We study three types of similarity functions:

- 196 • Inner-product similarity:  $\text{sim}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$ ,
- 197 • (Low-rank) bilinear similarity [22]:  $\text{sim}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T P_x^T P_y \mathbf{y}$ ,
- 198 • Gated (low-rank) bilinear similarity:  $\text{sim}(\mathbf{x}, \mathbf{y}) = g \mathbf{x}^T P_x^T P_{xy} \mathbf{y} + (1 - g) \mathbf{x}^T P_x^T P_{xx} \mathbf{x}$ ,

199 where  $P_x \in \mathbb{R}^{d \times C}$ ,  $P_y \in \mathbb{R}^{d \times C}$ ,  $P_{xy} \in \mathbb{R}^{d \times C}$ , and  $P_{xx} \in \mathbb{R}^{d \times C}$  are the linear projection matrices,  
 200  $d$  is the embedding dimension, and  $g \in (0, 1)$  is a gate unit learned by linear projection. We model the  
 201 projection matrices to be low-rank, *i.e.*,  $d < C$ , to reduce the number of parameters and computational  
 202 complexity. The gate-modulated bilinear similarity is designed to further filter out the encoder noise.  
 203 The gate is generated by learning a linear projection matrix used to project the decoder feature  $\mathcal{X}$   
 204 to a single-channel mask, and then it is normalized to  $(0, 1)$  by sigmoid function. Finally, we have  
 205  $\mathcal{Y} = G\mathcal{Y} + (1 - G)\mathcal{X}$  (nearest neighbor interpolation is used for matching the resolution), where  
 206  $G$  is the matrix form of the gate unit. The use of the gate implies we reduce noise in some spatial  
 207 regions by replacing encoder features  $\mathcal{Y}$  with decoder features  $\mathcal{X}$ . Its vector form explains that in the  
 208 area of noise it tends to switch to the self-similarity mode. We will prove the effectiveness of the  
 209 gating mechanism by comparing the gate-modulated bilinear similarity with a baseline without the  
 210 gating mechanism.

211 In practice, encoder and decoder features may have different data distributions, which is not suitable  
 212 to compare similarity directly. Therefore we apply LayerNorm for both encoder and decoder  
 213 features before similarity computation. Indeed we observe that the loss does not decrease without  
 214 normalization.

215 As shown in Fig. 3, our implementation is similar to the previous dynamic upsampling operators,  
 216 which first generates the upsampling kernels and then assembles the decoder feature conditioned on  
 217 the kernels. We highlight the kernel generation module. By setting a kernel size of  $K$ , for each encoder

Table 1: Computational complexity and parameters of CARAFE and SAPA. **I**: inner-product similarity; **B**: bilinear similarity; **G**: gated bilinear similarity.

Module	Operation	FLOPs ( $\times HW$ )	Params
CARAFE	Kernel generation	$Cd + 36K^2d$	$Cd + 36K^2d$
	Feature assembly	$4K^2C$	0
	<b>Total</b>	$Cd + 36K^2d + 4K^2C$	$Cd + 36K^2d$
IndexNet	Kernel generation	$128C^2 + 512C$	$32C^2 + 128C$
	Feature assembly	$4C$	0
	<b>Total</b>	$128C^2 + 516C$	$32C^2 + 128C$
A2U	Kernel generation	$2K^2C + 4K^2 + C$	$4K^2C + 2C$
	Feature assembly	$4K^2C$	0
	<b>Total</b>	$6K^2C + 4K^2 + C$	$4K^2C + 2C$
SAPA	$\mathcal{Y}$ embedding	$4Cd$	$Cd$
	$\mathcal{X}$ embedding	$Cd$	$Cd$
	Gated addition	$C + 8d$	$C$
	Kernel generation	$4K^2d$	0
	Feature assembly	$4K^2C$	0
	<b>I Total</b>	$8K^2C$	0
	<b>B Total</b>	$5Cd + 4K^2d + 4K^2C$	$2Cd$
	<b>G Total</b>	$5Cd + 4K^2d + 4K^2C + C + 5d$	$2Cd + C$

feature point, we compute the similarity scores between this point and each of  $K \times K$  neighbors in the decoder. Then, the softmax normalization is applied to generate the upsampling kernels. SAPA is lightweight and can even work without additional parameters (with inner-product similarity). To intuitively understand its lightweight property, we compare the computational complexity and number of parameters of different dynamic upsampling operators in Table 1. For example, when  $C = 256$  and  $d = 64$ , the FLOPs are  $H * W * 199K$ ,  $H * W * 17M$ ,  $H * W * 28K$ , and  $H * W * 228K$  for CARAFE, IndexNet, A2U, and SAPA-B, respectively.

We visualize the feature maps of upsampling kernels and upsampled features in Fig. 2. Our upsampling kernels show more favorable responses than other upsampling operators, with weights highlighted on boundaries and noise suppressed in regions, which visually supports our proposition and is a concrete embodiment of Eq. (4).

## 5 Experiments

We first focus our experiments on semantic segmentation to justify the effectiveness of SAPA. We then showcase its universality across three additional dense prediction tasks, including object detection, depth estimation, and image matting. All our experiments are run on a server with 8 NVIDIA GeForce RTX 3090 GPUs.

### 5.1 Data Sets, Metrics, Baselines, and Protocols

For semantic segmentation, we conduct experiments on the ADE20K dataset [4] and report the mIoU metric. Three strong transformer-based models are adopted as the baselines, including SegFormer-B1 [5], MaskFormer-Swin-Base [6] and Mask2Former-Swin-Base [7]. All training settings and implementation details are kept the same as the original papers. We only modify the upsampling stages with specific upsampling operators.

For object detection, we use the MS COCO [8] dataset, which involves 80 object categories. We use AP as the evaluation metric. Faster RCNN [23] with ResNet-50 [24] is adopted as the baseline. We use mmdetection [25] and follow its training configurations.

For depth estimation, we use NYU Depth V2 dataset [10] and its default train/test split. We choose BTS [9] with ResNet-50 as the baseline and follow its training configurations. The inlier measure  $\delta_1 < 1.25$  and RMSE are reported as the evaluation metrics. We replace all upsampling stages but the last one for SAPA, due to no available high-res feature map for the last stage.

For image matting, we train the model on the Adobe Image Matting dataset [11] and report four metrics on the Composition-1k test set, including SAD, MSE, Gradient, and Connectivity [26]. A2U

Table 2: Semantic segmentation results on ADE20K. **I**: inner-product similarity; **B**: bilinear similarity; **G**: gated bilinear similarity. Best performance is in boldface and second best is underlined.

	SegFormer B1 [5]			MaskFormer [6]			Mask2Former [7]		
	mIoU $\uparrow$	FLOPs	Params	mIoU $\uparrow$	FLOPs	Params	mIoU $\uparrow$	FLOPs	Params
Nearest	–	–	–	52.70	195	102	–	–	–
Bilinear	41.68	15.91	13.74	–	–	–	53.90	223	107
CARAFE [1]	42.82	+1.83	+0.44	<u>53.53</u>	+1.67	+0.22	53.94	+6.02	+0.07
IndexNet [2]	41.50	+30.66	+12.60	52.92	+17.60	+6.30	54.71	+13.40	+2.10
A2U [3]	41.45	+0.41	+0.12	52.73	+0.24	+0.06	54.40	+0.72	+0.02
SAPA-I	43.05	+1.50	+0	53.25	+0.86	+0	<u>55.05</u>	+2.62	+0
SAPA-B	<u>43.20</u>	+3.32	+0.20	53.15	+1.91	+0.10	54.98	+5.83	+0.03
SAPA-G	<b>44.39</b>	+3.34	+0.20	<b>53.78</b>	+1.92	+0.10	<b>55.22</b>	+5.86	+0.03

Table 3: Evaluation of object detection on MS COCO, monocular depth estimation on NYU Depth V2, and image matting on Adobe Composition-1k. Best results are in boldface and second best are underlined. Full metrics can be found in the supplementary material.

	Faster RCNN [23]			BTS [9]		A2U Matting [3]		
	AP $\uparrow$	Params	RMSE $\downarrow$	$\delta_1 < 1.25 \uparrow$	Params	SAD $\downarrow$	Grad $\downarrow$	Params
Nearest	37.4	41.53	0.419	0.865	49.53	37.51	19.07	8.05
CARAFE [1]	<b>38.6</b>	+0.22	0.418	0.864	+0.41	41.01	21.39	+0.26
IndexNet [2]	37.6	+6.30	0.416	0.866	+44.20	34.28	15.94	+12.26
A2U [3]	37.3	+0.12	0.429	0.860	+0.15	32.15	16.39	+0.04
SAPA-I	37.7	+0	n/a	n/a	n/a	34.25	18.93	+0
SAPA-B	<u>37.8</u>	+0.10	<u>0.410</u>	<u>0.871</u>	+0.31	<u>31.19</u>	<b>15.48</b>	+0.07
SAPA-G	<u>37.8</u>	+0.10	<b>0.408</b>	<b>0.872</b>	+0.49	<b>30.98</b>	<u>15.59</u>	+0.07

matting [3] is adopted as the baseline. We use the code provided by the authors and follow the same training settings as in the original paper.

## 5.2 Main Results

We compare SAPA and its variants against different upsampling operators on the three strong segmentation baselines. Results are shown in Table 2, from which we can see that SAPA consistently outperforms other upsampling operators. Note that SAPA can work well even without parameters and achieves the best performance with only few additional #Params and #Flops.

Results on other three dense prediction tasks are shown in Table 3. SAPA outperforms other upsampling operators on depth estimation and image matting, but falls behind CARAFE on object detection. One plausible explanation is that the demand of details in object detection is low (Section 6 presents a further in-depth analysis). In detail-sensitive tasks like image matting, SAPA significantly outperforms CARAFE. Qualitative comparisons are provided in supplementary material.

## 5.3 Ablation Study

Here we conduct ablation studies to compare the choices of similarity function and normalization function, the effect of different kernel sizes, and the number of embedding dimension. For the default setting, we use the bilinear similarity function, apply the normalization function  $h(x) = e^x$ , set the kernel size  $K = 5$  and the embedding dimension  $d = 64$ . Quantitative results are shown in Table 4.

**Similarity Function** We investigate three types of similarity function aforementioned and an additional ‘plain addition’ baseline. It ablates the gating mechanism from the gated bilinear similarity. Among them, gated bilinear similarity generates the best performance, which highlights the complementarity of semantic regions and boundary details in kernel generation.



Table 4: Ablation studies. **I**: inner-product similarity; **B**: bilinear similarity; **P**: plain addition; **G**: gated bilinear similarity.

SegFormer-B1		Sim func				Embedding dim			
Setting		<b>I</b>	<b>B</b>	<b>P</b>	<b>G</b>	16	32	64	128
mIoU		43.1	43.2	43.2	44.4	43.0	43.4	43.2	43.4
SegFormer-B1		Norm func				Kernel size			
Setting		None	$e^x$	relu	sigmoid	softplus	3	5	7
mIoU		41.45	44.4	42.2	43.5	43.3	43.1	43.2	42.5

**Normalization Function** We also investigate different normalization functions. ‘None’ indicates no normalization is used. Among validated functions,  $h(x) = e^x$  works the best, which means normalization matters. The other three have to play with the *epsilon* trick to prevent zero division.

**Kernel Size** The kernel size controls the region that each point in the upsampled feature can attend to. Results show that, compared with a large kernel, a small local window is sufficient to distinguish the semantic clusters.

**Embedding Dimension** We further study the influence of embedding dimension in the range of 16, 32, 64, and 128. Interestingly, results suggest that SAPA is not sensitive to the embedding dimension. This also verifies that SAPA extracts existing knowledge rather than learn unknown knowledge.

## 6 Discussion

**Understanding SAPA from a Backward Perspective** We have explained SAPA in the forward pass, here we further discuss how it may work during training. In fact, originally the model does not know to which the semantic cluster an encoder point should belong. The working mechanism of SAPA assigns each encoder feature point a possibility to choose a cluster. During training, the ground truths produce gradients, thus changes the assignment possibility. In SAPA, for every encoder point, the semantic clusters in the corresponding local window in decoder features serve as *implicit labels* and cooperate with the ground truths. The correct cluster is the positive label, while the wrong one is negative. In the preliminary stage of training, if an encoder feature point is more similar to a wrong cluster, it will be punished by gradients and engender large losses, and vice versa. Therefore, the encoder feature points can gradually learn its affiliation. We find this process is fast by visualizing the epoch-to-epoch feature maps.

**Limitations compared with CARAFE** CARAFE is a purely semantic-driven upsampling operator able to mend semantic information with a single-input flow. Such a mechanism in CARAFE brings advantages in smoothing semantic regions. *E.g.*, we observe it mends holes in a continuous region. However, due to its single-input flow, it cannot compensate the details lost in downsampling. Our SAPA, by contrast, mainly aims to compensate details such as textures and boundaries. SAPA characters in two aspects: semantic preservation and detail delineation. However, as Eq. (4) suggests, we do not add any semantic mending mechanism in SAPA. This explains why SAPA is worse than CARAFE on object detection, because detection has less demand for details but more for regional integrity. In short, CARAFE highlights *semantic mending*, while SAPA highlights *semantic preserving* and *detail delineation*.

## 7 Conclusion

In this paper, we introduce similarity-aware point affiliation, *i.e.*, SAPA. It not only indicates a lightweight but effective upsampling operator suitable for tasks like semantic segmentation, but also expresses a high-level concept that characterizes feature upsampling. SAPA can serve as an universal substitution for conventional upsampling operators. Experiments show the effectiveness of SAPA and also indicate its limitation: it is more suitable for tasks that favor detail delineation.

## References

- [1] Jiaqi Wang, Kai Chen, Rui Xu, Ziwei Liu, Chen Change Loy, and Dahua Lin. CARAFE: Context-aware reassembly of features. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, pages 3007–3016, 2019.
- [2] Hao Lu, Yutong Dai, Chunhua Shen, and Songcen Xu. Indices matter: Learning to index for deep image matting. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*, pages 3266–3275, 2019.
- [3] Yutong Dai, Hao Lu, and Chunhua Shen. Learning affinity-aware upsampling for deep image matting. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, pages 6841–6850, 2021.
- [4] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, 2017.
- [5] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [6] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, 2021.
- [7] Bowen Cheng, Ishan Misra, Alexander G. Schwing, Alexander Kirillov, and Rohit Girdhar. Masked-attention mask transformer for universal image segmentation. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, 2022.
- [8] 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 *Proceedings of European Conference on Computer Vision (ECCV)*, pages 740–755, 2014.
- [9] Jin Han Lee, Myung-Kyu Han, Dong Wook Ko, and Il Hong Suh. From big to small: Multi-scale local planar guidance for monocular depth estimation. *arXiv Computer Research Repository*, 2019.
- [10] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from RGBD images. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 746–760, 2012.
- [11] Ning Xu, Brian Price, Scott Cohen, and Thomas Huang. Deep image matting. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, pages 2970–2979, 2017.
- [12] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. SegNet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(12):2481–2495, 2017.
- [13] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. SUN RGB-D: A RGB-D scene understanding benchmark suite. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, pages 567–576, 2015.
- [14] Diana Wofk, Fangchang Ma, Tien-Ju Yang, Sertac Karaman, and Vivienne Sze. FastDepth: Fast monocular depth estimation on embedded systems. In *Proceedings of IEEE International Conference on Robotics and Automation (ICRA)*, 2019.
- [15] Ruibo Li, Ke Xian, Chunhua Shen, Zhiguo Cao, Hao Lu, and Lingxiao Hang. Deep attention-based classification network for robust depth prediction. In *Proceedings of Asian Conference on Computer Vision (ACCV)*, pages 663–678, 2018.
- [16] Wenzhe Shi, Jose Caballero, Ferenc Huszár, Johannes Totz, Andrew P Aitken, Rob Bishop, Daniel Rueckert, and Zehan Wang. Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, pages 1874–1883, 2016.

- [17] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(4):640–651, 2015.
- [18] H. Gao, H. Yuan, Z. Wang, and S. Ji. Pixel transposed convolutional networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2019.
- [19] Augustus Odena, Vincent Dumoulin, and Chris Olah. Deconvolution and checkerboard artifacts. *Distill*, 1(10), 2016.
- [20] Zbigniew Wojna, Vittorio Ferrari, Sergio Guadarrama, Nathan Silberman, Liang-Chieh Chen, Alireza Fathi, and Jasper Uijlings. The devil is in the decoder: Classification, regression and gans. *International Journal of Computer Vision*, 127(11):1694–1706, 2019.
- [21] Z. Tian, T. He, C. Shen, and Y. Yan. Decoders matter for semantic segmentation: Data-dependent decoding enables flexible feature aggregation. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, 2020.
- [22] Hamed Pirsiavash, Deva Ramanan, and Charless Fowlkes. Bilinear classifiers for visual recognition. *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 22, 2009.
- [23] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 28, 2015.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, pages 770–778, 2016.
- [25] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. *arXiv Computer Research Repository*, 2019.
- [26] Christoph Rhemann, Carsten Rother, Jue Wang, Margrit Gelautz, Pushmeet Kohli, and Pamela Rott. A perceptually motivated online benchmark for image matting. In *Proceedings of IEEE Conference on Computer Vision Pattern Recognition (CVPR)*, pages 1826–1833, 2009.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
  - (b) Did you describe the limitations of your work? [\[Yes\]](#) Please see Section 6.
  - (c) Did you discuss any potential negative societal impacts of your work? [\[N/A\]](#)
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
  - (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#)
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[No\]](#) Our method generates stable results.

- 403 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
404 of GPUs, internal cluster, or cloud provider)? [Yes]
- 405 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 406 (a) If your work uses existing assets, did you cite the creators? [Yes] For each of the  
407 datasets used, we cite the creators, as shown in Section 5.1.
- 408 (b) Did you mention the license of the assets? [Yes] The data used in our work is open  
409 source and can be used for academic research.
- 410 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
411 Our code for this work will appear in the supplementary material.
- 412 (d) Did you discuss whether and how consent was obtained from people whose data you're  
413 using/curating? [N/A] The data used in our work is open source and can be used for  
414 academic research.
- 415 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
416 information or offensive content? [N/A] The data do not contain personally identifiable  
417 information or offensive content.
- 418 5. If you used crowdsourcing or conducted research with human subjects...
- 419 (a) Did you include the full text of instructions given to participants and screenshots, if  
420 applicable? [N/A]
- 421 (b) Did you describe any potential participant risks, with links to Institutional Review  
422 Board (IRB) approvals, if applicable? [N/A]
- 423 (c) Did you include the estimated hourly wage paid to participants and the total amount  
424 spent on participant compensation? [N/A]