# SAPA: Similarity-Aware Point Affiliation for Feature Upsampling

Anonymous Author(s) Affiliation Address email

#### Abstract

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

# 16 **1** Introduction

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

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

<sup>&</sup>lt;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.

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

the upsampling kernels from the encoder feature while filtering out the rest.

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

<sup>59</sup> '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 that this form implements our expected upsampling behaviors: encouraging both semantic smoothness 61 and boundary sharpness. Following our formulation, we further instantiate a lightweight upsampling 62 operator, termed Similarity-Aware Point Affiliation (SAPA), and investigate its variants. We evaluate 63 SAPA across a number of mainstream dense prediction tasks, for example: i) semantic segmentation: 64 we test SAPA on several transformer-based segmentation baselines on the ADE20K dataset [4], such 65 as SegFormer [5], MaskFormer [6], and Mask2Former [7], improving the baselines by  $1\% \sim 2.7\%$ 66 mIoU; ii) *object detection*: SAPA improves the performance of Faster RCNN by 0.4% AP on MS 67 COCO [8]; iii) monocular depth estimation: SAPA reduces the rmse metric of BTS [9] from 0.419 68 to 0.408 on NYU Depth V2 [10]; and iv) *image matting*: SAPA outperforms a strong A2U matting 69 baseline [3] on the Adobe Composition-1k testing set [11] with a further 3.8% relative error reduction 70 in the SAD metric. SAPA also outperforms or at least is on par with other state-of-the-art dynamic 71 upsampling operators. Particularly, even without additional parameters, SAPA outperforms the 72 previous best upsampling operator CARAFE on semantic segmentation. 73



Additional kernel visualizations of SAPA

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 too much detail from encoder to preserve semantics. Please refer to supplementary materials for additional visualizations.

# 74 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 79 affiliation in terms of relative distances in upsampling, which follows fixed rules to designate point 80 affiliation, even if the true affiliation may be different. Max unpooling [12] stores the indices of 81 max-pooled feature points in encoder features and uses the sparse indices to guide upsampling. While 82 it executes location-specific point affiliation which benefits detail recovery, most upsampled points 83 84 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 85 space ones. The notion of point affiliation does not apply to this operator, however. 86

Another stream of upsampling operators implement learning-based upsampling. Among them, 87 transposed convolution or deconvolution [17] is known as an inverse convolutional operator. Based 88 on a novel interpretation of deconvolution, PixelTCL [18] is proposed to alleviate the checkerboard 89 artifacts [19] of standard deconvolution. In addition, bilinear additive upsampling [20] attempts 90 to combine learnable convolutional kernels with hand-crafted upsampling operators to achieve 91 composited upsampling. Recently, DUpsample [21] seeks to reconstruct the upsampled feature map 92 with pre-learned projection matrices, expecting to achieve a data-dependent upsampling behavior. 93 While these operators are learnable, the upsampling kernels are fixed once learned, still resulting in 94 fixed designation of point affiliation. 95

96 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.

<sup>99</sup> 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.

#### **105 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  $\mathcal{X}$ , if we upsample it to the target feature  $\mathcal{X}'$ , which is often twice the size of  $\mathcal{X}$ , 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

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

where  $\psi$  refers to a kernel generation module, which is often implemented by a sub-network used to 112 predict the kernel conditioned on  $\mathcal{N}_{l'}$ , and norm is a normalization function. The feature  $\mathcal{N}_{l'}$  can 113 originate from two sources, to be specific, from the encoder feature  $\mathcal{Y}$  or from the decoder feature  $\mathcal{X}$ . 114 If the encoder feature  $\mathcal{Y}$  is chosen as the source, then a local region of  $\mathcal{Y}$  centered at l' is extracted 115 to be  $\mathcal{N}_{l'}$ . If the decoder feature  $\mathcal{X}$  is chosen, one first needs to compute the projective location 116 of l', *i.e.*,  $l = (\lfloor \frac{i'}{2} \rfloor, \lfloor \frac{j'}{2} \rfloor)$  according to the spatial location correspondence, then a neighborhood region extracted from  $\mathcal{X}_l$  is regarded as  $\mathcal{N}_{l'}$ . The softmax function is often used as the normalization 117 118 function such that relevant points can be softly selected to compute the value of the target point using 119 the weight  $\mathcal{W}_{l'}$ . 120

As illustrated in Fig. 2, the source of  $N_{l'}$  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. Benefited from the encoder feature, the kernel maps generated by IndexNet and A2U have rich details; however, they manifest large similarity to the original encoder feature, which means noise is introduced into the kernel.

Feature Assembly In this step, for each target feature point at l', we assemble the corresponding sub-region of decode feature with the predicted kernel  $W_{l'}$  and assign a value to this target point. If the predicted kernel has a kernel size of  $K \times K$ , then the sub-region is also of size  $K \times K$  centered at l, where  $l = (i, j) = (\lfloor \frac{i'}{2} \rfloor, \lfloor \frac{j'}{2} \rfloor)$ . We compute the weighted sum between the normalized kernel and the sub-region of decoder feature to obtain the value of  $X'_{l'}$  by

$$\mathcal{X}_{l'}' = \sum_{n=-r}^{r} \sum_{m=-r}^{r} \mathcal{W}_{l'(n,m)} \cdot \mathcal{X}_{(i+n,j+m)}, \ r = \lfloor \frac{K}{2} \rfloor.$$
<sup>(2)</sup>

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 wellpredicted kernel can encourge 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.

# 138 4 Rethinking Point Affiliation with Local Mutual Similarity

To obtain an expected upsampling kernel mentioned above, we first derive a generic formulation for generating upsampling kernels by exploiting local mutual similarity, then explain why the formulation

encourages semantic smoothness and boundary sharpness, and finally present an upsampling operator, 141 SAPA, that embodies our formulation. 142

#### 4.1 Local Mutual Similarity 143

We rethink point affiliation from the view of local mutual similarity between encoder and decoder 144 features. With a detailed analysis, we explain why such similarity can assist point affiliation. 145

We first define a generic similarity function  $sim(x, y) : \mathbb{R}^C \times \mathbb{R}^C \to \mathbb{R}$ . It scores the similarity 146 between a vector x and a vector y of the same dimension C. We also define a normalization 147 function involving n real numbers  $x_1, x_2, ..., x_n$  by  $\operatorname{norm}(x_i : x_1, x_2, ..., x_n) = \frac{h(x_i)}{\sum_{j=1}^n h(x_j)}$ , where 148  $h(x) : \mathbb{R} \to \mathbb{R}$  is an arbitrary function, ignoring zero division. Given sim(x, y) and h(x), we can define a generic formulation for generating the upsampling kernel 149 150

$$w = \frac{h\left(\operatorname{sim}(\boldsymbol{x}, \boldsymbol{y})\right)}{\sum\limits_{\boldsymbol{x}' \in N_{l'}} h\left(\operatorname{sim}(\boldsymbol{x}', \boldsymbol{y})\right)},$$
(3)

where w is the kernel value specific to x and y. To analyze the upsampling behavior of the kernel, 151 we further define the following notations. 152

Let  $\mathcal{Y} \in \mathbb{R}^{2H \times 2W \times C}$  and  $\mathcal{X} \in \mathbb{R}^{H \times W \times C}$  be the encoder and decoder feature maps, respectively. Let  $\boldsymbol{y}_{i,j} = \mathcal{Y}_{(i,j)} \in \mathbb{R}^C, (i,j) \in \{0, 1, ..., 2H - 1\} \times \{0, 1, ..., 2W - 1\}$  denote the encoder feature vector and  $\boldsymbol{x}_{k,l} = \mathcal{X}_{(k,l)} \in \mathbb{R}^C, (k,l) \in \{0, 1, ..., H - 1\} \times \{0, 1, ..., W - 1\}$  be the decoder feature vector, where *C* is the number of channels. 153 154 155 156

Let K be the upsampling kernel size. Our operation will be done within a local window of size 157  $K \times K$ , between each encoder vector  $y_{i,j}$  and all its spatially associated decoder feature vectors, 158

159 
$$\boldsymbol{x}_{\lfloor \frac{i}{2} \rfloor + m, \lfloor \frac{j}{2} \rfloor + n}$$
's, where  $m, n = -\lfloor \frac{K}{2} \rfloor, ..., \lfloor \frac{K}{2} \rfloor$ .

w

To simplify our analysis, we also assume local smoothness. That is, points with the same semantic 160 cluster will have a similar value, which means a local region will share the same value on every 161 channel of the feature map. As shown in Fig. 1, we define  $\boldsymbol{x}_w = (R_1^w, R_2^w, ..., R_C^w)^T$  and  $\boldsymbol{x}_p = (R_1^p, R_2^p, ..., R_C^p)^T$  as the semantic clusters related to 'wall' and 'picture', respectively, where  $R_c^w$  and  $R_c^p$  are constants, and c = 1, 2, ..., C. For ease of analysis, we define two types of windows 162 163 164 distinguished by their contents. When all the points inside a window belong to the same semantic 165 cluster, it is called a smooth window; while different semantic clusters appear in a window, it is 166 defined as a detail window. 167

Next, we explain why the kernel can filter noise and can encourage semantic smoothness in a smooth 168 window, and why it can help to recover details when dealing with boundaries/textures in a detail 169 window. 170

**Upsampling in a Smooth Window** Without loss of generality, we consider the point  $y_{i_1,j_1}$  in Fig. 1 in the encoder feature. Its corresponding window is a smooth window of with the semantic cluster 'picture', thus  $x_{\lfloor \frac{i_1}{2} \rfloor + m, \lfloor \frac{j_1}{2} \rfloor + n} = x_p, m, n = -\lfloor \frac{K}{2} \rfloor, ..., \lfloor \frac{K}{2} \rfloor$ . Then the upsampling kernel weight for the upsampled point  $(i_1, j_1)$  at the position (m, n) takes the form 171 172 173

174

$$= \operatorname{norm}(\operatorname{sin}(\boldsymbol{x}_{\lfloor \frac{i_{1}}{2} \rfloor + m, \lfloor \frac{j_{1}}{2} \rfloor + n}, \boldsymbol{y}_{i_{1}, j_{1}}))$$

$$= \frac{h(\operatorname{sim}(\boldsymbol{x}_{\lfloor \frac{i_{1}}{2} \rfloor + m, \lfloor \frac{j_{1}}{2} \rfloor + n}, \boldsymbol{y}_{i_{1}, j_{1}}))}{\sum_{s=-\lfloor \frac{K}{2} \rfloor} \sum_{t=-\lfloor \frac{K}{2} \rfloor} h(\operatorname{sim}(\boldsymbol{x}_{\lfloor \frac{i_{1}}{2} \rfloor + s, \lfloor \frac{j_{1}}{2} \rfloor + t}, \boldsymbol{y}_{i_{1}, j_{1}}))},$$

$$= \frac{h(\operatorname{sim}(\boldsymbol{x}_{p}, \boldsymbol{y}_{i_{1}, j_{1}}))}{K^{2}h(\operatorname{sim}(\boldsymbol{x}_{p}, \boldsymbol{y}_{i_{1}, j_{1}}))}$$

$$= \frac{1}{K^{2}}$$
(4)

which has nothing to do with  $i_1, j_1, m$ , and n. Eq. (4) reveals a key characteristic of local mutual 175 similarity in a smooth window: the kernel weight is a constant regardless of y. Therefore, the kernel 176

fundamentally can filter out noise from encoder features with an 'average' kernel. 177



Figure 3: **Feature upsampling of SAPA.** For every encoder feature point, mutual similarity scores are computed by a similarity function with the spatially associated local region in the decoder features. The scores are then transformed into upsampling kernel weights after normalization. In this way each of the 4 upsampling kernels can select some decoder points from the  $3 \times 3$  decoder window, conditioned on mutual similarity, and assigns the value to the corresponding upsampled point.

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

**Upsampling in a Detail Window** Again we consider two points  $y_{i_{2w},j_{2w}}$  and  $y_{i_{2p},j_{2p}}$  in Fig. 1 in the encoder feature. Ideally  $y_{i_{2w},j_{2w}}$  and  $y_{i_{2p},j_{2p}}$  should be classified into the semantic cluster of 'wall' and 'picture', respectively. Taking  $y_{i_{2w},j_{2w}}$  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 180 181 182 183 tell that  $\sin(\boldsymbol{x}_{\lfloor\frac{i_2}{2}\rfloor+m_w,\lfloor\frac{j_2}{2}\rfloor+m_w}, \boldsymbol{y}_{i_{2w},j_{2w}})$  is larger than  $\sin(\boldsymbol{x}_{\lfloor\frac{i_2}{2}\rfloor+m_p,\lfloor\frac{j_2}{2}\rfloor+m_p}, \boldsymbol{y}_{i_{2w},j_{2w}})$ , where  $\boldsymbol{x}_{\lfloor\frac{i_2}{2}\rfloor+m_w,\lfloor\frac{j_2}{2}\rfloor+m_w}$  and  $\boldsymbol{x}_{\lfloor\frac{i_2}{2}\rfloor+m_p,\lfloor\frac{j_2}{2}\rfloor+m_p}$  belong to the 'wall' cluster and the 'picture' cluster, respectively. Therefore, after computing similarity scores and normalization, one can acquire a kernel 184 185 186 with significantly larger weights on points with the semantic cluster of 'wall' than that of 'picture', 187 *i.e.*,  $w_{i_{2w},j_{2w},m_w,n_w} >> w_{i_{2w},j_{2w},m_p,n_p}$ . After applying the kernel to the corresponding window, the upsampled point at  $(i_{2w},j_{2w})$  will be revalued and assigned to the semantic cluster of 'wall'. 188 189 Similarly, the upsampled point at  $(i_{2p}, j_{2p})$  will be assigned into 'picture'. 190

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

#### 194 4.2 SAPA: Similarity-Aware Point Affiliation

<sup>195</sup> Here we embody Eq. (3) by choosing appropriate similarity and normalization functions.

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

Similarity Function We focus on the choices of similarity functions and study three types of functions:

- Inner-product similarity:  $sim(x, y) = x^T y$ ,
- (Low-rank) bilinear similarity [22]:  $sim(x, y) = x^T P_x^T P_y y$ ,
- Gated (low-rank) bilinear similarity:  $sim(\boldsymbol{x}, \boldsymbol{y}) = g\boldsymbol{x}^T P_{\boldsymbol{x}}^T P_{\boldsymbol{x}\boldsymbol{y}} \boldsymbol{y} + (1 g)\boldsymbol{x}^T P_{\boldsymbol{x}}^T P_{\boldsymbol{x}\boldsymbol{x}} \boldsymbol{x}$ ,

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, d is the embedding dimension, and  $g \in (0, 1)$  is a gate unit learned by linear projection. The gatemodulated bilinear similarity is designed to further filter out the encoder noise. The gate is generated by learning a linear projection matrix used to project the decoder feature  $\mathcal{X}$  to a single-channel mask, and then it is normalized to (0, 1) by sigmoid function. Then, we have  $\mathcal{Y} = G\mathcal{Y} + (1 - G)\mathcal{X}$ (nearest neighbor interpolation is used for matching the resolution), where G is the matrix form of the

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

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

gate unit. The use of the gate implies we reduce noise in some spatial regions by replacing encoder features  $\mathcal{Y}$  with decoder features  $\mathcal{X}$ . Its vector-form explains that in the area of noise it tends to switch to self-similarity mode. We will prove the effectiveness of the gating mechanism by comparing the gate-modulated bilinear similarity with a baseline without the gating mechanism.

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 217 on the kernels. We highlight the kernel generation module. By setting a kernel size of K, for each 218 encoder feature point, we compute the similarity scores between this point and each  $K \times K$  neighbor 219 220 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). 221 To intuitively understand its lightweight property, we compare the computational complexity and 222 number of parameters of different dynamic upsampling operators in Table 1. For example, when 223 C = 256 and d = 64, the FLOPs are H \* W \* 199K, H \* W \* 17M, H \* W \* 28K, and H \* W \* 228K224 for CARAFE, IndexNet, A2U, and SAPA-B, respectively. 225

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).

# 230 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.

#### 235 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.

	Seg	Former B1	[5]	M	askFormer	[6]	Mask2Former [7]			
	mIoU↑	FLOPs	Params	mIoU↑	FLOPs	Params	mIoU↑	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	53.53	+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	55.05	+2.62	+0	
SAPA-B	43.20	+3.32	+0.20	53.15	+1.91	+0.10	54.98	+5.83	+0.03	
SAPA-G	44.39	+3.34	+0.20	53.78	+1.92	+0.10	55.22	+5.86	+0.03	

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.

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↑	Params	RMSE↓	$\delta_1 < 1.25 \uparrow$	Params	SAD↓	Grad↓	Params
Nearest	37.4	41.53	0.419	0.865	49.53	37.51	19.07	8.05
CARAFE [1]	38.6	+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	35.14	19.16	+0
SAPA-B	37.8	+0.10	0.410	0.871	+0.31	32.47	16.20	+0.07
SAPA-G	37.8	+0.10	0.408	0.872	+0.49	30.98	15.59	+0.07

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 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.

#### 252 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.

#### 262 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.

Pacalina	Sim funa	Norm funa	Kornal siza	Embadding dim	
gated bilinear	similarity.				
Table 4: Ablat	tion studies. I: inner-p	roduct similarity; B: l	oilinear similarity;	P: plain addition; G:	

Baseline	seline Sim func		Norm func		Kernel size			Embedding dim							
SegFormer B1	Ι	В	Р	G	$e^x$	relu	sigmoid	softplus	3	5	7	16	32	64	128
mIoU	43.1	43.2	43.2	44.4	44.4	42.2	43.5	43.3	43.1	43.2	42.5	43.0	43.4	43.2	43.4

267 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.

**Normalization Function** We also investigate different normalization functions. Among validated functions,  $h(x) = e^x$  works the best, which means normalization matters. However, we think normalization mainly affects the optimization of the network, the other three have to play with the *epsilon* trick to prevent zero division, which may affect performance.

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
 between regions and edges.

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.

# 281 6 Discussion

Understanding SAPA from a Backward Perspective We have explained SAPA in the forward 282 pass, here we further discuss how it may work during training. In fact, originally the model does 283 not know to which the semantic cluster an encoder point should belong. The working mechanism 284 of SAPA assigns each encoder feature point a possibility to choose a cluster. During training, the 285 ground truths produce gradients, thus changes the assignment possibility. In SAPA, for every encoder 286 point, the semantic clusters in the corresponding local window in decoder features serve as implicit 287 labels and cooperate with the ground truths. The correct cluster is the positive label, while the wrong 288 one is negative. In the preliminary stage of training, if an encoder feature point is more similar to a 289 wrong cluster, it will be punished by gradients and engender large losses, and vice versa. Therefore, 290 the encoder feature points can gradually learn its affiliation. We find this process is fast by visualizing 291 292 the epoch-to-epoch feature maps.

Limitations compared with CARAFE CARAFE is a purely semantic-driven upsampling operator 293 able to mend semantic information with a single-input flow. Such a mechanism in CARAFE brings 294 advantages in smoothing semantic regions. E.g., we observe it mends holes in a continuous region. 295 However, due to its single-input flow, it cannot compensate the details lost in downsampling. Our 296 SAPA, by contrast, mainly aims to compensate details such as textures and boundaries. SAPA 297 characters in two aspects: semantic preservation and detail delineation. However, as Eq. (4) suggests, 298 we do not add any semantic mending mechanism in SAPA. This explains why SAPA is worse 299 than CARAFE on object detection, because detection has less demand for details but more for 300 regional integrity. In short, CARAFE highlights semantic mending, while SAPA highlights semantic 301 preserving and detail delineation. 302

# 303 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.

#### **309** 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.

- I17] 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.
- 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.

# 387 Checklist

- 388 1. For all authors...
- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's 389 contributions and scope? [Yes] 390 (b) Did you describe the limitations of your work? [Yes] Please see Section 6. 391 (c) Did you discuss any potential negative societal impacts of your work? [N/A] 392 (d) Have you read the ethics review guidelines and ensured that your paper conforms to 393 them? [Yes] 394 2. If you are including theoretical results... 395 (a) Did you state the full set of assumptions of all theoretical results? [N/A] 396 (b) Did you include complete proofs of all theoretical results? [N/A] 397 398 3. If you ran experiments... (a) Did you include the code, data, and instructions needed to reproduce the main experi-399 mental results (either in the supplemental material or as a URL)? [Yes] 400 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they 401 were chosen)? [Yes] 402 (c) Did you report error bars (e.g., with respect to the random seed after running experi-403 ments multiple times)? [No] Our method generates stable results. 404

405 406	<ul><li>(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]</li></ul>
407	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
408 409	(a) If your work uses existing assets, did you cite the creators? [Yes] For each of the datasets used, we cite the creators, as shown in Section 5.1.
410 411	(b) Did you mention the license of the assets? [Yes] The data used in our work is open source and can be used for academic research.
412 413	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes] Our code for this work will appear in the supplementary material.
414 415 416	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] The data used in our work is open source and can be used for academic research.
417 418 419	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] The data do not contain personally identifiable information or offensive content.
420	5. If you used crowdsourcing or conducted research with human subjects
421 422	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
423 424	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
425 426	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]