

AINet: Association Implantation for Superpixel Segmentation

Yaxiong Wang^{1,2} Yunchao Wei^{3‡} Xueming Qian^{1‡} Li Zhu^{1†} Yi Yang⁴
¹Xi'an Jiaotong University ²Baidu Research
³Beijing Jiaotong University ⁴Zhejiang University

Abstract

Recently, some approaches are proposed to harness deep convolutional networks to facilitate superpixel segmentation. The common practice is to first evenly divide the image into a pre-defined number of grids and then learn to associate each pixel with its surrounding grids. However, simply applying a series of convolution operations with limited receptive fields can only implicitly perceive the relations between the pixel and its surrounding grids. Consequently, existing methods often fail to provide an effective context when inferring the association map. To remedy this issue, we propose a novel **Association Implantation (AI)** module to enable the network to explicitly capture the relations between the pixel and its surrounding grids. The proposed AI module directly implants the grid features to the surrounding of its corresponding central pixel, and conducts convolution on the padded window to adaptively transfer knowledge between them. With such an implantation operation, the network could explicitly harvest the pixel-grid level context, which is more in line with the target of superpixel segmentation comparing to the pixel-wise relation. Furthermore, to pursue better boundary precision, we design a boundary-perceiving loss to help the network discriminate the pixels around boundaries in hidden feature level, which could benefit the subsequent inferring modules to accurately identify more boundary pixels. Extensive experiments on *BSDS500* and *NYUv2* datasets show that our method could achieve state-of-the-art performance. Code and pre-trained model are available at <https://github.com/wangyxxjtu/AINet-ICCV2021>.

1. Introduction

Superpixels are image regions formed by grouping image pixels similar in color and other low-level properties, which could be viewed as an over-segmentation of image. The process of extracting superpixels is known as superpixel segmentation. Comparing to pixels, superpixel pro-

[†]Corresponding author

[‡]Co-corresponding author

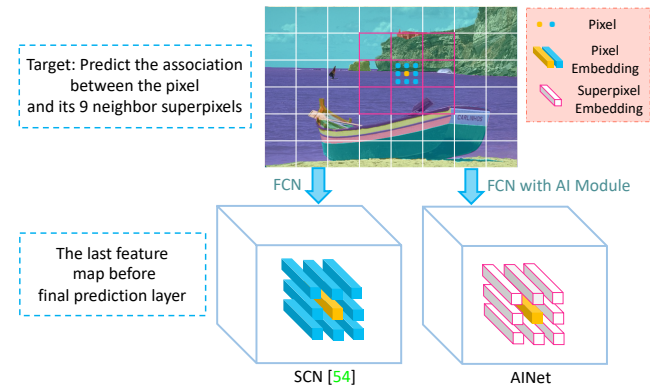


Figure 1: Different from the SCN [54] that implicitly learns the association using the cascaded convolutions, our AINet proposes to implant the corresponding grid features to the surrounding of the pixel to explicitly perceive the relation between each pixel and its neighbor grids.

vides a more effective representation for image data. With such a compact representation, the computational efficiency of vision algorithms could be improved [22, 12, 46]. Consequently, superpixel could benefit many vision tasks like semantic segmentation [13, 53, 60, 58, 59], object detection [10, 40], optical flow estimation [15, 32, 44, 52], and even adversarial attack [8]. In light of the fundamental importance of superpixels in computer vision, superpixel segmentation attracts much attention since it is first introduced by Ren and Malik [38] in 2003.

The common practice for superpixel segmentation is to first split the image into grid cells and then estimate the membership of each pixel to its adjacent grids, by which the grouping could be performed. Meanwhile, the membership estimation plays the key role in superpixel segmentation. Traditional approaches usually utilize the hand-craft features and estimate the relevance of pixel to its neighbor grids based on clustering or graph-based methods [1, 31, 25, 28, 2], however, these methods all suffer from the weakness of the hand-craft features and are difficult to integrate into other trainable deep frameworks. Inspired by the success of deep neural networks in many computer vision problems, researchers recently attempts

to adopt the deep learning technique to superpixel segmentation [18, 54, 47]. As mentioned in abstract, previous deep methods attempt to assign pixels by learning the association of each pixel to its surrounding grids using the fully convolutional networks [41]. The popular solutions like SCN [54], SSN [18] employ the U-net architecture [39] to predict the association, *i.e.*, the 9-way probabilities, for each pixel. Although stacking convolution layers can enlarge the receptive field and help study the *pixel-grid wise* probabilities, introducing low-level features with skip connection in the final layer will pollute the probabilities due to the added *pixel-pixel wise* information, since the ultimate target is to predict the association between the target pixel and *its 9-neighbor grids* instead of its 9-neighbor pixels.

To tackle this weakness, we propose to directly implant the grid features to the surrounding of the corresponding pixel using an association implantation (AI) module. Fig. 1 simply shows the core idea of our AI module, before feeding the last features into the prediction layer, our AI module is performed: for each pixel, we place the corresponding grid features to its neighbors, then a convolution with 3×3 kernel is followed, this convolution is no longer to capture the pixel-pixel relation but *the relation between pixel and its 9 neighbor grids*, providing the consistent context with the target of superpixel segmentation. Our proposed AI module provides a simple and intuitive way to allow the network to harvest the pixel-neighbor grids context in an explicit fashion, which is exactly required by superpixel segmentation. Comparing to existing methods, such a design is more consistent with the target of superpixel segmentation and could give more beneficial support for the subsequent association map inferring.

Besides, a satisfactory superpixel algorithm should actually identify the boundary pixels, however, some designs towards this target still missed among existing works. To pursue better boundary precision, we augment the optimization with a boundary-perceiving loss. To be specific, we first sample a set of small local patches on the pixel embedding map along the boundaries. Then, the features with the same/different labels in each patch are treated as the positive/negative samples, on which a classification procedure is performed to enhance the compactness of the features with the same label while distinguish the different semantic features. Our boundary-perceiving loss encourages the model to pay more attention to discriminate the features around boundaries, consequently, more boundary pixels could be identified.

Quantitative and qualitative results on BSDS500 [3] and NYUv2 [42] datasets demonstrate that the proposed method achieves more outstanding performance against the state-of-the-art superpixel segmentation methods. In summary, we make the following contributions in this work:

- We propose a novel AI module to directly capture the

relation between the pixel and its surrounding grid cells, such a design builds a more consistent architecture with the target of superpixel segmentation.

- A boundary-perceiving loss is designed to discriminate the features with different semantic labels around boundaries, which could help the network accurately identify boundary pixels and improve the boundary precision.

2. Related Work

Superpixel Segmentation Superpixel segmentation is a well-defined problem and has a long line of research [43, 33, 23, 48, 6, 20]. Traditional superpixel algorithms can be broadly classified into graph-based and clustering-based approaches. Graph-based methods consider the pixels as nodes and the edges as strength of connectivity between adjacent pixels, respectively. Consequently, the superpixel segmentation could be formulated as a graph-partitioning problem. Wide-used algorithms, Felzenszwalb and Huttenlocher (FH) [9] and the entropy rate superpixels (ERS) [28], belong to this category. On the other hand, clustering-based approaches utilize classic clustering techniques like k -means to compute the connectivity between the anchor pixels and its neighbors. Well-known methods in this category include SLIC [1], LSC [25], Manifold-SLIC [31] and SNIC [2]. Inspired by the success of deep learning techniques, recently, researchers attempt to utilize the deep neural network to learn the membership of each pixel to its surrounding grid cells. Jampani *et al.* [18] develop the first differentiable deep network motivated by the classic SLIC method, and Yang *et al.* [54] further simplify the framework and contribute a more efficient model.

Application of Superpixel The pre-computed superpixel segmentation could be viewed as a type of weak label or prior knowledge to benefit many downstream tasks like image & video segmentation [34, 57, 29, 16, 35, 19, 17, 55, 51, 49, 5], object detection [24, 27, 50, 37] and so on. The superpixels could be integrated into deep learning pipeline to provide guidance so that some important image properties (e.g., boundaries) could be better preserved [11, 45, 56, 4, 26]. For example, KwaK *et al.* [22] utilize the superpixel segmentation to perform a region-wise pooling to make the pooled feature have better semantic compactness. In [36], Cheng *et al.* consider the superpixel as pseudo label and attempt to boost the image segmentation by identifying more semantic boundary. Besides benefiting the image segmentation or feature pooling, superpixel also provides flexible ways to represent the image data. He *et al.* [14] convert 2D images patterns into 1D sequential representation, such a novel representation allows the deep network to explore the long-range context of the image. Liu *et al.* [30] learn the similarity between different superpixels,

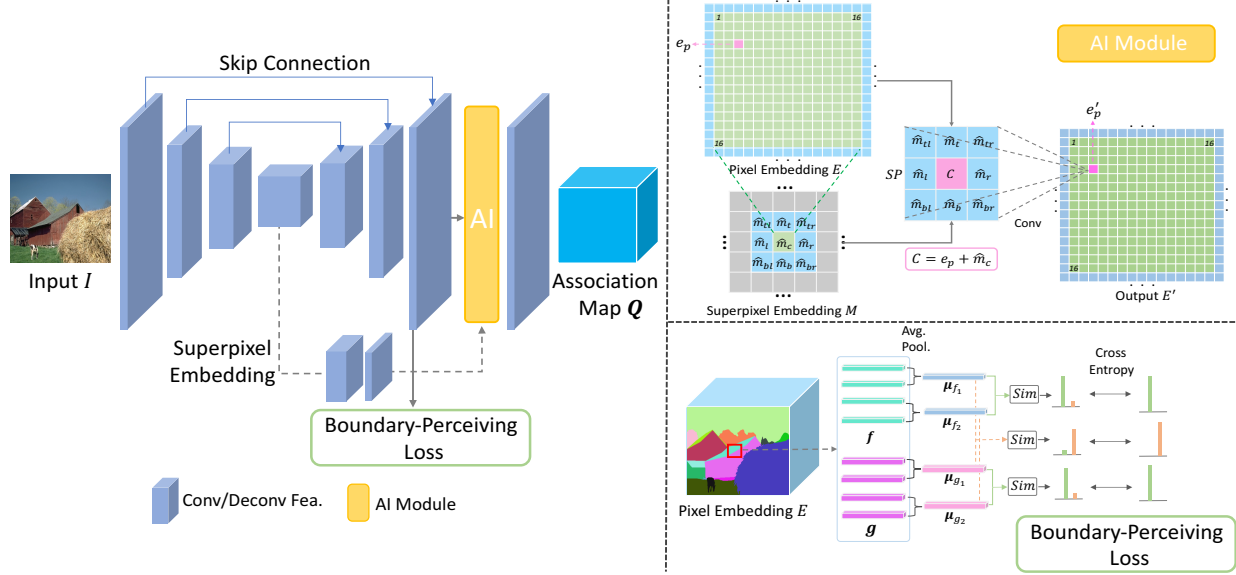


Figure 2: The framework of our AINet. The network takes an image as input and outputs the association map. Meanwhile, the superpixel embedding and pixel embedding are first obtained by the convolutions and then fed into the AI module to obtain the pixel-superpixel context. And the local patch loss is performed on the pixel-wise embeddings to boost the boundary precision. In AI module, the sampling interval is set to 16, and each block indicates a pixel or superpixel embedding.

the developed framework could produce different grained segmentation regions by merging the superpixels according to the learned superpixel similarity.

3. Preliminaries

Before delving into the details of our method, we first introduce the framework of deep-learning based superpixel segmentation, which is also the fundamental theory of this paper. As illustrated in Fig. 1, the image I is partitioned into blocks using a regular size grid, and the grid cell is regarded as the initial superpixel seed. For each pixel p in image I , the superpixel segmentation aims at finding a mapping that assigns each pixel to one of its surrounding grids, *i.e.* 9 neighbors, just as shown in Fig. 1. Mathematically, deep-learning based method feeds the image $I \in \mathcal{R}^{H \times W \times 3}$ to convolution neural network and output an association map $Q \in \mathcal{R}^{H \times W \times 9}$, which indicates the probability of each pixel to its neighbor grids [18, 54]. Since there is no ground-truth for such an output, the supervision for network training is performed in an indirect fashion: the predicted association map Q serves as the intermediate variable to reconstruct the pixel-wise property $l(p)$ like semantic label, position vector, and so on. Consequently, there are two critical steps in training stage.

Step1: Estimate the superpixel property from the surrounding pixels:

$$h(s) = \frac{\sum_{p:s \in N_p} l(p) \cdot q(p, s)}{\sum_{p:s \in N_p} q(p, s)}. \quad (1)$$

Step2: Reconstruct the pixel property according to the superpixel neighbors:

$$l'(p) = \sum_{s \in N_p} h(s) \cdot q(p, s), \quad (2)$$

where the N_p is the set of adjacent superpixels of p , $q(p, s)$ indicates the probability of pixel p assigned to superpixel s . Thus, the training loss is to optimize the distance between the ground-truth property and the reconstructed one:

$$\mathcal{L}(Q) = \sum_p \text{dist}(l(p), l'(p)). \quad (3)$$

Following Yang's practice [54], the properties of pixel in this paper include the semantic label and the position vector, *i.e.*, two-dimension spatial coordinates, which are optimized by the cross-entropy loss and the \mathcal{L}_2 reconstruction loss, respectively.

4. Methodology

An overview of our proposed AINet is shown in Fig. 2. In general, the overall architecture is an encoder-decoder style paradigm, the encoder module compresses the input image and outputs a feature map called superpixel embedding, whose pixels exactly encode the features of grid cells. Subsequently, the superpixel embedding is further fed into the decoder module to produce the association map. Meanwhile, the superpixel embedding and the pixel embedding

in decoding stage are integrated to perform the association implantation, and the boundary-perceiving loss also acts on the pixel embedding. Hereinafter, we elaborate the details of our proposed AI module and boundary-perceiving loss.

4.1. Association Implantation Module

To enable the network to explicitly perceive the relation between each pixel and its surrounding grid cells, this work proposes an association implantation module to perform a direct interaction between the pixel and its neighbor grids. As shown in the top right of Fig. 2, we first obtain the embeddings of superpixels and pixels by the convolution network. Then, for each pixel embedding, the corresponding neighbor superpixel features are picked and implanted to its surrounding. Finally, a convolution with kernel size 3×3 is conducted on the expanded pixel embedding to achieve the knowledge propagation.

Formally, let $e_p \in \mathcal{R}^D$ be the embedding of pixel p from the pixel embedding $E \in \mathcal{R}^{H \times W \times D}$, which is obtained by the deep neural network as shown in Fig. 2. To obtain the embeddings of the grid cells, *i.e.*, superpixel embedding, we compress the input image by $\log_2 S$ times using multiple convolutions and max-pooling operations, where S is the sampling interval for the grid cell. For example, if the sampling interval is 16, then, we downsample the image 4 times. This would result in a feature map $M \in \mathcal{R}^{h \times w \times D'}$ whose pixels exactly encode the features of grid cells, where $h = H/S$, and $w = W/S$. To perform the implantation operation on the pixel embedding, we first adjust the channels of M using two 3×3 convolutions, producing a new map $\hat{M} \in \mathcal{R}^{H \times W \times D}$. Then, for the pixel p , we pick up its 9 adjacent superpixel embeddings from left to right and top down: $\{\hat{m}_{tl}, \hat{m}_t, \hat{m}_{tr}, \hat{m}_l, \hat{m}_c, \hat{m}_r, \hat{m}_{bl}, \hat{m}_b, \hat{m}_{br}\}$ from \hat{M} . To allow the network could explicitly capture the relation between pixel p and its neighbor grids, we directly implant the superpixel embeddings into the surrounding of the pixel p to provide pixel-superpixel context:

$$SP = \begin{bmatrix} \hat{m}_{tl} & \hat{m}_t & \hat{m}_{tr} \\ \hat{m}_l & \hat{m}_c + e_p & \hat{m}_r \\ \hat{m}_{bl} & \hat{m}_b & \hat{m}_{br} \end{bmatrix}. \quad (4)$$

It is worth noting that the pixels in the same initial grid would share the same surrounding superpixels, since they would degrade into one element in superpixel view. We then adopt a 3×3 convolution to adaptively distill information from the expanded window to benefit the subsequent association map inferring:

$$e'_p = \sum_{ij} SP_{ij} \times w_{ij} + b, \quad (5)$$

where w and b are the convolution weight and bias, respectively. We traverse all of the pixel embeddings in E and

apply the operations in Eq. 4-5, thus, we could obtain a new pixel embedding E' whose elements capture the pixel-superpixel level context. In the following, the feature map E' is fed through a convolution layer to predict the association map Q .

As shown in Eq. 4-5, our AI module directly places the neighbor grid embeddings in the surrounding of the pixel to provide the context required by superpixel segmentation, which is an intuitive and reasonable solution. Comparing to the existing methods that use the stacked convolutions to accumulate the pixel-wise relation, the pixel-superpixel context captured by our AI module is more in line with the target of superpixel segmentation.

4.2. Boundary-Perceiving Loss

Our boundary-perceiving loss is proposed to help the network appropriately assign the pixels around boundaries. As shown in the bottom right of Fig. 2, we first sample a series of patches with a certain size (5×5 , for example) around boundaries in the pixel embedding map, and then a classification procedure is conducted to improve the discrimination of the different semantic features.

Formally, let $E \in \mathcal{R}^{H \times W \times D}$ be the pixel-wise embedding map, since the ground-truth label is available during training stage, we could sample a local patch $B \in \mathcal{R}^{K \times K \times D}$ surrounding a boundary pixel from E . For the sake of simplification, the patch B only covers the pixels from two different semantic regions, that is, $B = \{f_1, \dots, f_m, g_1, \dots, g_n\}$, where $f, g \in \mathcal{R}^D$, $m+n = K^2$. Intuitively, we attempt to make the features in the same categories be closer, while the embeddings from different labels should be far away from each other. To this end, we evenly partition the features in the same categories into two groups, f^1, f^2, g^1, g^2 , and employ a classification-based loss to enhance the discrimination of the features:

$$\begin{aligned} \mathcal{L}_B = & -\frac{1}{2}(\log(\text{sim}(\boldsymbol{\mu}_{f1}, \boldsymbol{\mu}_{f2})) + \log(1 - \text{sim}(\boldsymbol{\mu}_{f1}, \boldsymbol{\mu}_{g1}))) \\ & - \frac{1}{2}(\log(\text{sim}(\boldsymbol{\mu}_{g1}, \boldsymbol{\mu}_{g2})) + \log(1 - \text{sim}(\boldsymbol{\mu}_{f2}, \boldsymbol{\mu}_{g2}))), \end{aligned} \quad (6)$$

where the $\boldsymbol{\mu}_{f1}$ is the average representation for f^1 , and the function $\text{sim}(\cdot, \cdot)$ is the similarity measure for two vectors:

$$\boldsymbol{\mu}_{f1} = \frac{1}{|f^1|} \sum_{f \in f^1} f, \quad (7)$$

$$\text{sim}(f, g) = \frac{2}{1 + \exp(\|f - g\|_1)}, \quad (8)$$

Taking all of the sampled patches \mathcal{B} into consideration, our full boundary-perceiving loss is formulated as follow:

$$\mathcal{L}_B = \frac{1}{|\mathcal{B}|} \sum_{B \in \mathcal{B}} \mathcal{L}_B. \quad (9)$$

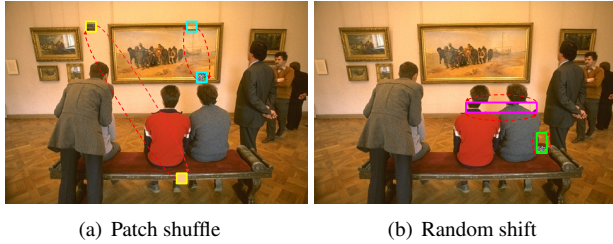


Figure 3: The illustrations for our patch jitter augmentation, patch shuffle and random shift. Color frames indicate the changed regions.

Overall, the full losses for our network training comprise three components, *i.e.*, cross-entropy (CE) and \mathcal{L}_2 reconstruction losses for the semantic label and position vector according to the Eq. 3, and our boundary-perceiving loss:

$$\mathcal{L} = \sum_p CE(l'_s(p), l_s(p)) + \lambda \|p - p'\|_2^2 + \alpha \mathcal{L}_B \quad (10)$$

where $l'_s(p)$ is the reconstructed semantic label from the predicted association map Q and the ground-truth label $l_s(p)$ according to Eq. 1-2, and λ, α are two trade-off weights.

5. Experiments

Datasets. We conduct experiments on two public benchmarks, **BSDS500** [3] and **NYUv2** [42] to evaluate the effectiveness of our method. BSDS500 comprises 200 training, 100 validation and 200 test images, and each image is annotated by multiple semantic labels from different experts. To make a fair comparison, we follow previous works [54, 18, 47] and treat each annotation as an individual sample. Consequently, 1,087 training, 546 validation samples and 1,063 testing samples could be obtained. NYUv2 is an indoor scene understanding dataset and contains 1,449 images with object instance labels. To evaluate the superpixel methods, Stutz *et al.* [43] remove the unlabelled regions near the boundary and collect a subset of 400 test images with size 608×448 for superpixel evaluation.

Following Yang’s practice [54], we conduct a standard train and test pipeline on the BSDS500 dataset. On the subject of NYUv2 dataset, we directly apply the model trained on BSDS500 and report the performance on the 400 tests to evaluate the generality of the model.

Augmentation via Patch Jitter. To further improve the performance and enhance the generality of our model, we propose to augment the data by jittering the image patches. Specifically, the proposed patch jitter augmentation comprises two components, *i.e.*, patch shuffle and random shift. Fig. 3 shows the respective examples for these two types of data augmentation. The **patch shuffle** first samples two image patches with shape $S \times S$ and then randomly ex-

change them to extend the image patterns, the corresponding ground-truth patches are also exchanged accordingly to maintain the consistency. To further augment the data, we randomly pick up one of the selected two patches and replace it with a random patch, whose ground-truth is assigned with a new label. While the **random shift** could be conducted along with the horizontal or vertical directions. For horizontal random shift, we first randomly sample a patch with shape $S \times L$, where $L = \text{rand_int}(S, W)$, and a random offset $o = \text{rand_int}(0, S)$. Then, we conduct a cycle translation on the patch by o offset towards left or right. Meanwhile, the random patch trick in patch shuffle could also be adopted. Finally, the augmentation is done by replacing the original patch with the new one. Analogously, the augmentation along vertical direction could be done similarly. The patch jitter augmentation is repeated 2 times during training. SSN [54] is adopted as our baseline method in our experiments.

Implementation Details. In training stage, the image is randomly cropped to 208×208 as input, and the network is trained using the adam optimizer [21] for 4k iterations with batch size 16. The learning rate starts with $8e-5$ and is discounted by 0.5 for every 2K iterations. The sampling interval is fixed as 16, consequently, the encoder part employs 4 convolution&pooling operations to get the superpixel embedding with shape $13 \times 13 \times 256$. The following decoder module produces the pixel embedding with shape $208 \times 208 \times 16$ using 4 convolution&deconvolution operations. Then, the channels of superpixel embedding are first compressed by two convolution layers: $256 \Rightarrow 64 \Rightarrow 16$, then our AI module is performed. The boundary-perceiving loss also acts on the pixel embedding, where the patch size is set to 5, *i.e.*, $K = 5$. In the following, two convolution layers are stacked to predict the association map Q with shape $208 \times 208 \times 9$. In our practice, simultaneously equipping the boundary-perceiving loss and AI Module could not make the performance step further, therefore, we first train the network using the first two items in Eq. 10 for 3K iterations, and use the boundary-perceiving loss to finetune 1K. Following Yang’s practice [54], the weight of position reconstruction loss is set to $0.003/16$, while the weight for our boundary-perceiving loss is fixed to 0.5, *i.e.*, $\lambda = 0.003/16, \alpha = 0.5$. In testing, we employ the same strategy as [54] to produce varying numbers of superpixels.

Several methods are considered for performance comparison, including classic methods, SLIC [1], LSC [25], ERS [28], SEEDS [7], SNIC [2] and deep learning-based methods, SEAL [47], SSN [18], SCN [54]. We simply use the OpenCV implementation for methods SLIC, LSC and SEEDS. For other methods, we use the official implementations with the recommended parameters from the authors.

Evaluation Metrics. We use three popular metrics includ-

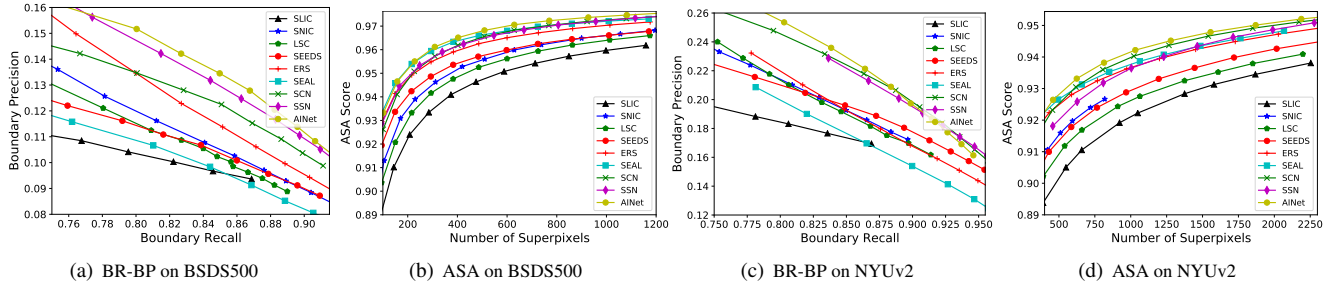


Figure 4: Performance comparison on datasets BSDS500 and NYUv2.

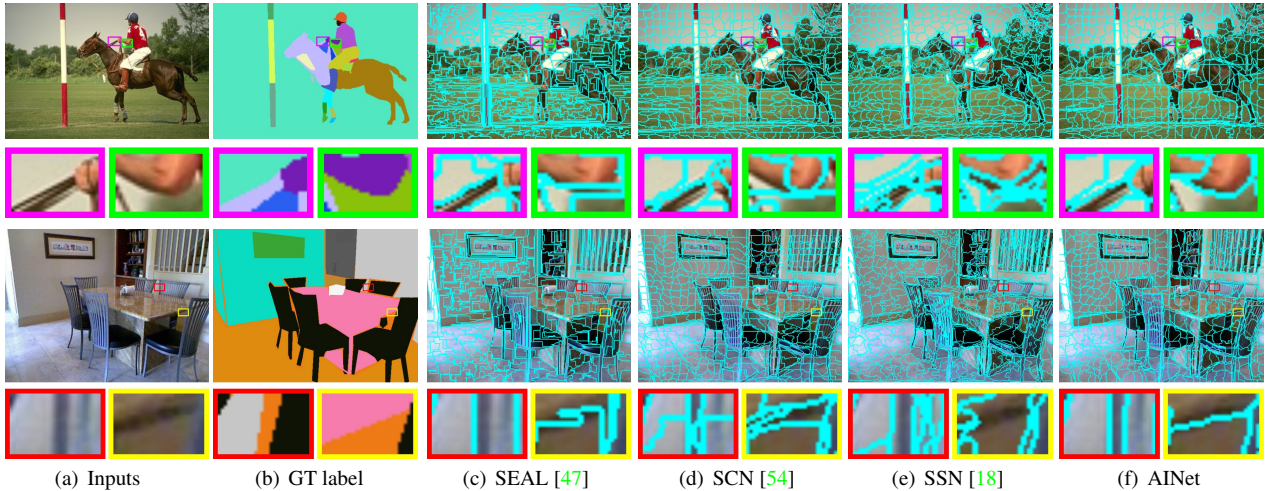


Figure 5: Qualitative results of four SOTA superpixel methods, SEAL, SCN, SSN, and our AINet. The top row exhibits the results from BSDS500 dataset, while the bottom row shows the superpixels on NYUv2 dataset.

ing achievable segmentation accuracy (ASA), boundary recall (BR) and boundary precision (BP) to evaluate the performance of superpixel. ASA score studies the upper bound on the achievable segmentation accuracy using superpixel as pre-processing step, while BR and BP focus on accessing how well the superpixel model could identify the semantic boundaries. The higher value of these metrics indicates better superpixel segmentation performance.

5.1. Comparison with the state-of-the-arts

Fig. 4 reports the quantitative comparison results on BSDS500 and NYUv2 test sets. As indicated in Fig. 4, our AINet attains the best ASA score and BR-BP on both datasets. With the help of deep convolution networks, the methods, SEAL, SCN, SSN, and AINet could achieve superior or comparable performance against the traditional superpixel algorithms, and our AINet is the best model among them. From Fig. 4 (a)-(b), the AINet could surpass the traditional methods by a large margin on BSDS500 dataset. By harvesting the pixel-sueppixel level context and highlighting the boundaries, AINet could also outperform the deep methods SEAL, SCN and SSN. Fig. 4 (c)-(d) shows the performance when adapting to the NYUv2 test set, we

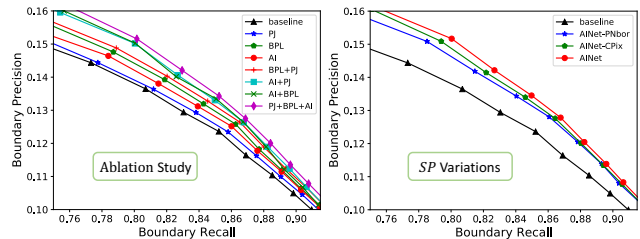
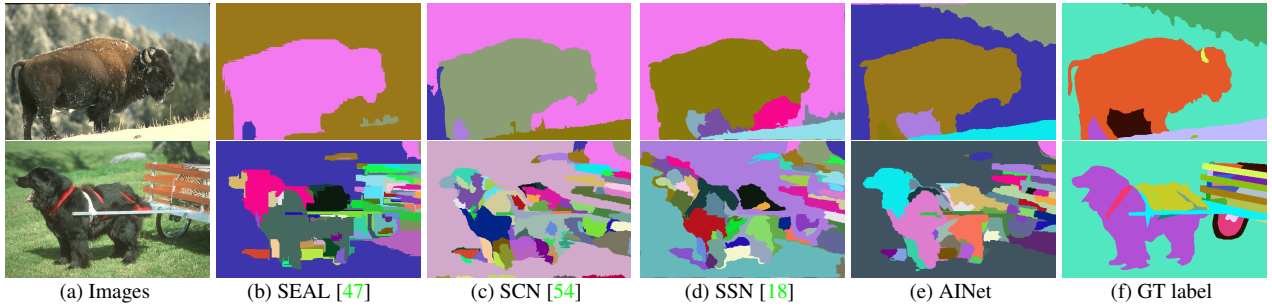


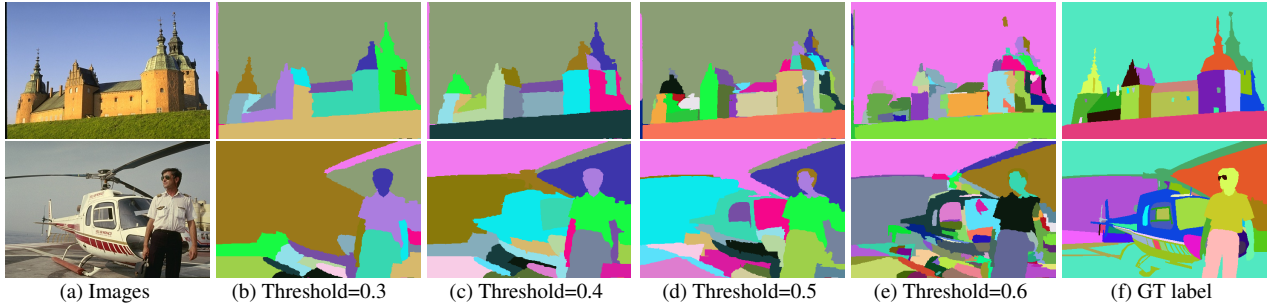
Figure 6: Ablation study on BSDS500. The left figure shows the contributions of each component in our system, while the right one discusses two variations of SP (Eq. 4).

can observe that the AINet also shows better generality. Although the BR-BP is comparable with the SCN and SSN, our ASA score is more outstanding than all of the competitive methods.

Fig. 5 shows the qualitative results of four state-of-the-art methods on dataset BSDS500 and NYUv2, comparing to the competing methods, the boundaries of our results are more accurate and clearer, which intuitively shows the superiority of our method.



I: The generated proposals from DEL [30] using different superpixels.



II: The generated proposals using different thresholds (1 as upper bound).

Figure 7: Qualitative proposals from DEL [30] using different superpixels (I), and the results of DEL [30] with our superpixel using different thresholds (II), where threshold=0.3 mean merging the adjacent superpixels if their similarity is above 0.3.

5.2. Ablation Study

To validate the respective contributions of our proposed modules including the data augmentation trick, AI module, and the boundary-perceiving loss, we conduct ablation study on BSDS500 dataset to thoroughly study their effectiveness. The left figure in Fig. 6 reports the performances of all methods, where the BPL means the boundary-perceiving loss, and BPL+PJ stands for the baseline simultaneously equipped with the boundary-perceiving loss and the patch jitter augmentation. From Fig. 6, we can observe that individually applying the three modules on the baseline method could all boost the performance, and the boundary-perceiving loss could contribute the most performance gains. The combination of the patch jitter augmentation and the BPL or AI Module could make the performance step further, and the AI module equipped with the data augmentation achieves better performance. When simultaneously employing the three modules, we could harvest the best BR-BP.

Besides, we also give a discussion for two alternative choices of SP (Eq. 4): a greedy version of SP that further adds the neighbor pixels to the corresponding surrounding superpixels like the central position, for example, \hat{m}_t is replaced by $\hat{m}_t + e_t$; And a simplified version that ignoring the central superpixel, *i.e.*, $\hat{m}_c + e_p$ changes to e_p . The models with the above two versions of SP are marked as AINet-

PNbor and AINet-CPix, respectively. The right figure of Fig. 6 shows the results, we can observe that AINet-PNbor and AINet-CPix could both surpass the baseline but perform a litter worse than AINet. By summing the neighbor pixels, the AINet-PNbor could integrate the pixel-wise relation, on the other hand, the sum operation would also reduce the force of superpixel embedding, which would conspire against capturing the pixel-superpixel context. For AINet-CPix, the excluded \hat{m}_c is also one of the neighbor superpixels, directly abandoning \hat{m}_c would fail to explicitly perceive the relation between pixel e_p and central superpixel \hat{m}_c . Consequently, the above two variations of SP are both not effective to capture the super context.

5.3. Inference Efficiency

Besides the performance, the inference speed is also a concerned aspect. Therefore, we conduct experiments on BSDS500 dataset to investigate the inference efficiency of four deep learning-based methods. To make a fair comparison, we only count the time of network inference and post-processing steps (if available). All methods run on the same workstation with NVIDIA 1080Ti GPU and Intel E5 CPU.

The time costs of four deep learning-based methods, SEAL, SCN, SSN and our AINet are reported in Fig. 8. The method SCN achieves the best inference efficiency due to its simple architecture, while our AINet introduces more

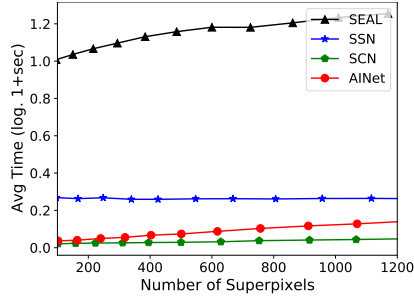


Figure 8: The average time costs of four deep learning based methods w.r.t number of superpixels. The runtime is added with 1 and scaled by logarithmic to show positive values and a clear tendency.

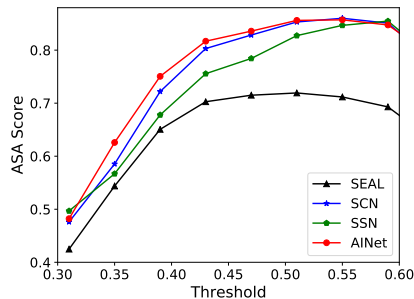


Figure 9: The ASA scores of four state-of-the-art methods on object proposal generation.

layers and operations, consequently, the inference is slightly slower than the SCN. The superpixel segmentation of SEAL and SSN is much complex comparing to the SCN and our AINet, SEAL needs first output the learned deep features and then feed them to a traditional algorithm to conduct superpixel segmentation, and SSN further performs the K -means iteration after obtaining the pixel affinity. As a result, SEAL and SSN both cost much more time in inference stage. Although the SCN is faster, the performance of AINet is much better than SCN. Comparing to these competing methods, our AINet achieves a good trade-off between the performance and the inference efficiency.

5.4. Application on Object Proposal Generation

Image annotation is one of the important application scenarios for superpixels, since it could identify the semantic boundaries and provide the outlines of many semantic regions. To generate the object proposals, Liu *et al.* [30] propose a model named DEL, they first estimate the similarities between superpixels and merge them according to a certain threshold, by which the proposed method could flexibly control the grain size of object proposal. In this subsection, we feed the superpixels from four state-of-the-art methods, SEAL, SCN, SSN and our AINet to the framework of [30] to further investigate the superiority of our AINet. To evaluate the performance, we use the ASA score to measure how

well the produced object proposals cover the ground-truth labels:

$$ASA(O, G) = \frac{1}{N} \sum_{O_k} \max_{G_k} \{|O_k \cap G_k|\}, \quad (11)$$

where N is the number of generated object proposal O , and G is the ground-truth semantic label.

The performance of all methods is reported in Fig. 9, from which we can observe that the average performance of our AINet is more outstanding. Fig. 7 I shows three results of DEL [30] with the superpixels from four methods, different thresholds are used to produce varied size proposals: the adjacent superpixels would be merged if their similarity is above the threshold, which means that higher value would produce finer object proposals. As shown in Fig. 7 I, our AINet could generate more satisfactory object proposals comparing to the competing methods, which validates the effectiveness of our proposed method. Fig. 7 II exhibits the results using the superpixels of our AINet with different thresholds, varying sizes of generated object proposals could be generated by adjusting the threshold.

6. Conclusion

We have presented an association implantation network for superpixel segmentation task. A novel association implantation module is proposed to provide the consistent pixel-superpixel level context for superpixel segmentation task. To pursue better boundary precision, a boundary-perceiving loss is designed to improve the discrimination of pixels around boundaries in hidden feature level, and a data augmentation named patch jitter is developed to further improve the performance. Experiments on two popular benchmarks show that the proposed method could achieve state-of-the-art performance with good generalizability. What's more, the produced superpixels by our method could also perform well when applied to the object proposal generation. In the future, we will continue to study the effectiveness of the proposed AINet on the stereo matching task.

Acknowledgements This work was supported by the National Key Research and Development Project under Grant 2019YFB2102500 and China Scholarship Council (CSC), No: 201806280398, in part by the NSFC under Grant 161772407, Grant 61701391, Grant 61732008, and Grant 61902309; in part by the Australian Research Council (ARC) Discovery Early Career Researcher Award (DECRA) under Grant DE190101315; in part by the ARC under Grant DP200100938; in part by the Emergence Mechanism and Calculation Method of Group Intelligence Based on Internet under Grant 2018AAA0101100; and Pazhou Lab,Guangzhou.

References

- [1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurélien Lucchi, Pascal Fua, and Sabine Süsstrunk. SLIC superpixels compared to state-of-the-art superpixel methods. *TPAMI*, 34(11):2274–2282, 2012.
- [2] Radhakrishna Achanta and Sabine Süsstrunk. Superpixels and polygons using simple non-iterative clustering. In *CVPR*, pages 4895–4904, June 2017.
- [3] Pablo Arbelaez, Michael Maire, Charless C. Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *TPAMI*, 33(5):898–916, 2011.
- [4] András Bódis-Szomorú, Hayko Riemenschneider, and Luc Van Gool. Superpixel meshes for fast edge-preserving surface reconstruction. In *CVPR*, pages 2011–2020, June 2015.
- [5] Zixuan Chen, Huajun Zhou, Jianhuang Lai, Lingxiao Yang, and Xiaohua Xie. Contour-aware loss: Boundary-aware learning for salient object segmentation. *TIP*, 30:431–443, 2021.
- [6] Dorin Comaniciu and Peter Meer. Mean shift: A robust approach toward feature space analysis. *TPAMI*, 24(5):603–619, 2002.
- [7] Michael Van den Bergh, Xavier Boix, Gemma Roig, and Luc Van Gool. SEEDS: superpixels extracted via energy-driven sampling. *IJCV*, 111(3):298–314, 2015.
- [8] Xiaoyi Dong, Jiangfan Han, Dongdong Chen, Jiayang Liu, Huanyu Bian, Zehua Ma, Hongsheng Li, Xiaogang Wang, Weiming Zhang, and Nenghai Yu. Robust superpixel-guided attentional adversarial attack. In *CVPR*, pages 12892–12901, June 2020.
- [9] Pedro F. Felzenszwalb and Daniel P. Huttenlocher. Efficient graph-based image segmentation. *IJCV*, 59(2):167–181, 2004.
- [10] Raghudeep Gade, Varun Jampani, Martin Kiefel, Daniel Kappler, and Peter V. Gehler. Superpixel convolutional networks using bilateral inceptions. In *ECCV*, pages 597–613, Oct. 2016.
- [11] Utkarsh Gaur and B. S. Manjunath. Superpixel embedding network. *TIP*, 29:3199–3212, 2020.
- [12] Shengfeng He, Rynson W. H. Lau, Wenxi Liu, Zhe Huang, and Qingxiong Yang. Supercnn: A superpixelwise convolutional neural network for salient object detection. *IJCV*, 115(3):330–344, 2015.
- [13] Shengfeng He, Rynson W. H. Lau, Wenxi Liu, Zhe Huang, and Qingxiong Yang. Supercnn: A superpixelwise convolutional neural network for salient object detection. *IJCV*, 115(3):330–344, 2015.
- [14] Shengfeng He, Rynson W. H. Lau, Wenxi Liu, Zhe Huang, and Qingxiong Yang. Supercnn: A superpixelwise convolutional neural network for salient object detection. *IJCV*, 115(3):330–344, 2015.
- [15] Yinlin Hu, Rui Song, Yunsong Li, Peng Rao, and Yangli Wang. Highly accurate optical flow estimation on superpixel tree. *Image Vis. Comput.*, 52:167–177, 2016.
- [16] Zilong Huang, Xinggang Wang, Yunchao Wei, Lichao Huang, Humphrey Shi, Wenyu Liu, and Thomas S. Huang. Ccnet: Criss-cross attention for semantic segmentation. *TPAMI*, pages 1–1, 2020.
- [17] Zilong Huang, Yunchao Wei, Xinggang Wang, Humphrey Shi, Wenyu Liu, and Thomas S. Huang. Alignseg: Feature-aligned segmentation networks. *TPAMI*, pages 1–1, 2021.
- [18] Varun Jampani, Deqing Sun, Ming-Yu Liu, Ming-Hsuan Yang, and Jan Kautz. Superpixel sampling networks. In *ECCV*, pages 363–380, Sep. 2018.
- [19] Guoliang Kang, Yunchao Wei, Yi Yang, Yueting Zhuang, and Alexander G. Hauptmann. Pixel-level cycle association: A new perspective for domain adaptive semantic segmentation. In *NeurIPS*, 2020.
- [20] Xuejing Kang, Lei Zhu, and Anlong Ming. Dynamic random walk for superpixel segmentation. *TIP*, 29:3871–3884, 2020.
- [21] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, May 2015.
- [22] Suha Kwak, Seunghoon Hong, and Bohyung Han. Weakly supervised semantic segmentation using superpixel pooling network. In *AAAI*, pages 4111–4117, Feb. 2017.
- [23] Alex Levinstein, Adrian Stere, Kiriakos N. Kutulakos, David J. Fleet, Sven J. Dickinson, and Kaleem Siddiqi. Turbopixels: Fast superpixels using geometric flows. *TPAMI*, 31(12):2290–2297, 2009.
- [24] Xin Li, Shenqi Lai, and Xueming Qian. Dbcfac: Towards pure convolutional neural network face detection. *TCSVT*, pages 1–1, 2021.
- [25] Zhengqin Li and Jiansheng Chen. Superpixel segmentation using linear spectral clustering. In *CVPR*, pages 1356–1363, 2015.
- [26] Zhenguo Li, Xiao-Ming Wu, and Shih-Fu Chang. Segmentation using superpixels: A bipartite graph partitioning approach. In *CVPR*, pages 789–796, June 2012.
- [27] Chengxu Liu, Yuanzhi Liang, Yao Xue, Xueming Qian, and Jianlong Fu. Food and ingredient joint learning for fine-grained recognition. *TCSVT*, 31(6):2480–2493, 2021.
- [28] Ming-Yu Liu, Oncel Tuzel, Srikumar Ramalingam, and Rama Chellappa. Entropy rate superpixel segmentation. In *CVPR*, pages 2097–2104, June 2011.
- [29] Si Liu, Tianrui Hui, Shaofei Huang, Yunchao Wei, Bo Li, and Guanbin Li. Cross-modal progressive comprehension for referring segmentation. *TPAMI*, pages 1–1, 2021.
- [30] Yun Liu, Peng-Tao Jiang, Vahan Petrosyan, Shi-Jie Li, Jiawang Bian, Le Zhang, and Ming-Ming Cheng. DEL: deep embedding learning for efficient image segmentation. In *IJ-CAI*, pages 864–870, July 2018.
- [31] Yong-Jin Liu, Cheng-Chi Yu, Minjing Yu, and Ying He. Manifold SLIC: A fast method to compute content-sensitive superpixels. In *CVPR*, pages 651–659, 2016.
- [32] Jiangbo Lu, Hongsheng Yang, Dongbo Min, and Minh N. Do. Patch match filter: Efficient edge-aware filtering meets randomized search for fast correspondence field estimation. In *CVPR*, pages 1854–1861, June 2013.
- [33] Vaia Machairas, Matthieu Faessel, David Cárdenas-Peña, Théodore Chabardès, Thomas Walter, and Etienne Decencière. Waterpixels. *TIP*, 24(11):3707–3716, 2015.
- [34] Jiayu Miao, Yunchao Wei, Yu Wu, Chen Liang, Guangrui Li, and Yi Yang. Vspw: A large-scale dataset for video scene parsing in the wild. In *CVPR*, 2021.

- [35] Jiayu Miao, Yunchao Wei, and Yi Yang. Memory aggregation networks for efficient interactive video object segmentation. In *CVPR*, pages 10363–10372. IEEE, 2020.
- [36] Cheng Ouyang, Carlo Biffi, Chen Chen, Turkey Kart, Huaqi Qiu, and Daniel Rueckert. Self-supervision with superpixels: Training few-shot medical image segmentation without annotation. In *ECCV*, pages 762–780, Aug. 2020.
- [37] Rong Quan, Junwei Han, Dingwen Zhang, Feiping Nie, Xueming Qian, and Xuelong Li. Unsupervised salient object detection via inferring from imperfect saliency models. *TMM*, 20(5):1101–1112, 2018.
- [38] Xiaofeng Ren and Jitendra Malik. Learning a classification model for segmentation. In *ICCV*, pages 10–17, Oct. 2003.
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241, Oct. 2015.
- [40] Abhishek Sharma, Oncel Tuzel, and Ming-Yu Liu. Recursive context propagation network for semantic scene labeling. In *NeurIPS*, pages 2447–2455, Dec. 2014.
- [41] Evan Shelhamer, Jonathan Long, and Trevor Darrell. Fully convolutional networks for semantic segmentation. *TPAMI*, 39(4):640–651, 2017.
- [42] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from RGBD images. In *ECCV*, pages 746–760, Oct. 2012.
- [43] David Stutz, Alexander Hermans, and Bastian Leibe. Superpixels: An evaluation of the state-of-the-art. *Comput. Vis. Image Underst.*, 166:1–27, 2018.
- [44] Deqing Sun, Ce Liu, and Hanspeter Pfister. Local layering for joint motion estimation and occlusion detection. In *CVPR*, pages 1098–1105, June 2014.
- [45] Wen Sun, Qingmin Liao, Jing-Hao Xue, and Fei Zhou. SP-SIM: A superpixel-based similarity index for full-reference image quality assessment. *TIP*, 27(9):4232–4244, 2018.
- [46] Teppei Suzuki, Shuichi Akizuki, Naoki Kato, and Yoshimitsu Aoki. Superpixel convolution for segmentation. In *ICIP*, pages 3249–3253, Oct 2018.
- [47] Wei-Chih Tu, Ming-Yu Liu, Varun Jampani, Deqing Sun, Shao-Yi Chien, Ming-Hsuan Yang, and Jan Kautz. Learning superpixels with segmentation-aware affinity loss. In *CVPR*, pages 568–576, June 2018.
- [48] Olga Veksler, Yuri Boykov, and Paria Mehrani. Superpixels and supervoxels in an energy optimization framework. In *ECCV*, pages 211–224, Sep. 2010.
- [49] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. Deep high-resolution representation learning for visual recognition. *TPAMI*, pages 1–1, 2020.
- [50] Xuan Wang, Shenqi Lai, Zhenhua Chai, Xingjun Zhang, and Xueming Qian. Spgnet: Serial and parallel group network. *TMM*, pages 1–1, 2021.
- [51] Yunchao Wei, Xiaodan Liang, Yunpeng Chen, Xiaohui Shen, Ming-Ming Cheng, Jiashi Feng, Yao Zhao, and Shuicheng Yan. STC: A simple to complex framework for weakly-supervised semantic segmentation. *TPAMI*, 39(11):2314–2320.
- [52] Koichiro Yamaguchi, David A. McAllester, and Raquel Urtasun. Robust monocular epipolar flow estimation. In *CVPR*, pages 1862–1869, June 2013.
- [53] Chuan Yang, Lihe Zhang, Huchuan Lu, Xiang Ruan, and Ming-Hsuan Yang. Saliency detection via graph-based manifold ranking. In *CVPR*, pages 3166–3173, June 2013.
- [54] Fengting Yang, Qian Sun, Hailin Jin, and Zihan Zhou. Superpixel segmentation with fully convolutional networks. In *CVPR*, pages 13961–13970, June 2020.
- [55] Zongxin Yang, Yunchao Wei, and Yi Yang. Collaborative video object segmentation by multi-scale foreground-background integration. *TPAMI*, pages 1–1, 2021.
- [56] Donghun Yeo, Jeany Son, Bohyung Han, and Joon Hee Han. Superpixel-based tracking-by-segmentation using markov chains. In *CVPR*, pages 511–520, June 2017.
- [57] Xiaolin Zhang, Yunchao Wei, Zhao Li, Chenggang Yan, and Yi Yang. Rich embedding features for one-shot semantic segmentation. *TNNLS*, pages 1–10, 2021.
- [58] Zhedong Zheng and Yi Yang. Unsupervised scene adaptation with memory regularization in vivo. In *IJCAI*, 2020.
- [59] Zhedong Zheng and Yi Yang. Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. *IJCV*, 2021. doi:[10.1007/s11263-020-01395-y](https://doi.org/10.1007/s11263-020-01395-y).
- [60] Wangjiang Zhu, Shuang Liang, Yichen Wei, and Jian Sun. Saliency optimization from robust background detection. In *CVPR*, pages 2814–2821, June 2014.