

Scribble-Supervised Semantic Segmentation by Random Walk on Neural Representation and Self-Supervision on Neural Eigenspace

Zhiyi Pan¹, Peng Jiang^{1*}, Changhe Tu^{1,2}

¹Shandong University

²AICFVE, Beijing Film Academy

panzhiyi1996@gmail.com, sdujump@gmail.com, changhe.tu@gmail.com

Abstract

Scribble-supervised semantic segmentation has gained much attention recently for its promising performance without high-quality annotations. Many approaches have been proposed. Typically, they handle this problem to either introduce a well-labeled dataset from another related task, turn to iterative refinement and post-processing with the graphical model, or manipulate the scribble label. This work aims to achieve semantic segmentation supervised by scribble label directly without auxiliary information and other intermediate manipulation. Specifically, we impose diffusion on neural representation by random walk and consistency on neural eigenspace by self-supervision, which forces the neural network to produce dense and consistent predictions over the whole dataset. The random walk embedded in the network will compute a probabilistic transition matrix, with which the neural representation diffused to be uniform. Moreover, given the probabilistic transition matrix, we apply the self-supervision on its eigenspace for consistency in the image's main parts. In addition to comparing the common scribble dataset, we also conduct experiments on the modified datasets that randomly shrink and even drop the scribbles on image objects. The results demonstrate the superiority of the proposed method and are even comparable to some full-label supervised ones. The code and datasets are available at <https://github.com/panzhiyi/RW-SS>.

Introduction

In recent years, the use of neural networks, especially convolutional neural networks, has dramatically improved semantic classification, detection, and segmentation (He et al. 2016; Long, Shelhamer, and Darrell 2015; Wang et al. 2018). As one of the most fine-grained ways to understand the scene, typically, semantic segmentation demands large-scale data with high-quality annotations to feed the network. However, the pixel-level annotating process for semantic segmentation is costly and tedious, limiting its flexibility and usability on some tasks that require rapid deployment (Lin et al. 2016). As a consequence, the scribble annotations, which are more easily available, become popular.

Correspondingly, some datasets and approaches are developed. Lin et al. 2016 proposed a scribble-annotated dataset based on the pascal VOC (Everingham et al. 2015)

and adopted the classic graphical model as post-processing to obtain the final dense predictions. Vernaza and Chandraker 2017 and Wang et al. 2019a optimize the performance of semantic segmentation by introducing an auxiliary task, edge detection. However, this edge detection task needs to be trained on another well-labeled dataset, so these methods have not relieved the heavy labor of annotation yet. To avoid the post-processing and dependence on another well-labeled dataset, Tang et al. 2018a and Tang et al. 2018b design graphical model based regularized loss to make the predictions between similar pixels consistent. But these two works only measured similarities in color and texture, did not consider semantic similarity for regularization. Moreover, most of these methods require every existing object in the image labeled, which is too strict for dataset preparation. In this work, we intend to propose a more flexible approach to get rid of the above issues holistically.

The proposed approach tackles scribble-supervised semantic segmentation by leveraging the uniform and consistency of neural representation (features). A representative result is shown in Fig. 1(a). The assumption is that, given the sparse scribble labels, if we guide the neural representation to be uniform within each image's objects and consistent between related images, then the neural network's best solution to the cross-entropy loss is the dense semantic predictions. To realize this kind of learning trend, we impose the diffusion on neural representation by random walk and consistency on neural eigenspace by self-supervision during training.

The random walk on neural representation has been studied in some works and could produce uniform and dense semantic predictions within each object (Jiang et al. 2018). The key to this kind of random walk is forming a probabilistic transition matrix that measures the similarity between neural representations. Then, with the transition matrix, the neural representation is diffused to be uniform. Fig. 1(b) shows several neural representations before and after the random walk. In addition to the uniform neural representation within each object, it is also essential to have consistent neural representation over related images to produce dense and confident semantic predictions. Namely, when the image is transformed, the neural representation should be the same as the original one with the corresponding transform. This kind of consistency is usually referred to as self-

*Corresponding Author

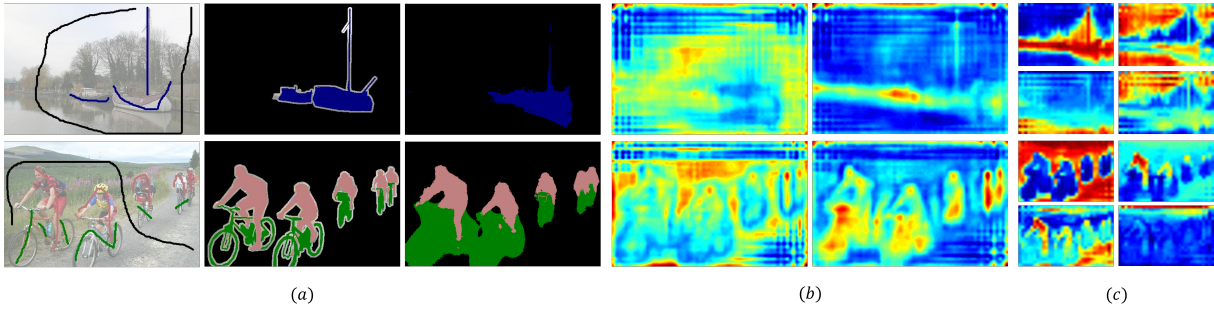


Figure 1: Results and Intermediate Visualizations. (a) From left to right: scribble-annotation, ground-truth, and our prediction. (b) From left to right: Neural representation before and after the random walk. (c) Leading eigenvectors of the transition matrix.

supervision (Laine and Aila 2016; Tarvainen and Valpola 2017; Mittal, Tatarchenko, and Brox 2019) and measured on neural representation. In this way, though with the various and sparse scribble labels, the network will still tend to generate consistent object perception.

However, for semantic segmentation, the consistency over the whole image is not necessary. When some parts of the image are distorted and changed heavily after transform, it is hard for the network to generate consistent neural representation anymore and may confuse the network in some scenarios. In this work, we propose to set the self-supervision on the main parts of images by imposing the consistent loss on the eigenspace of the transition matrix. The idea is inspired by spectral methods (Von Luxburg 2007), which observed that the eigenvectors of the Laplacian matrix have the capability to distinguish the main parts in the images, and some methods use this property for clustering (Ng, Jordan, and Weiss 2002; Nadler and Galun 2007) and saliency detection (Jiang et al. 2019; Yang et al. 2013). Since the eigenspace of the transition matrix has a close relation to the one of the Laplacian matrix, our self-supervision on the transition matrix’s eigenspace will also focus on the main image parts. Several leading eigenvectors are presented in Fig. 1(c).

The computation of eigenspace is time-consuming and unstable, especially during the dynamic optimizing of the neural network. Though some people have developed approximation methods (Dang et al. 2018; Wang et al. 2019b; Sun and Xu 2019), it is better to avoid the explicit eigenspace decomposition. Thus, in our implementation, we only apply a soft consistent loss on the eigenspace. For eigenvalue consistency, according to the fact that the matrix’s trace is equal to the sum of its eigenvalues, we measure the matrix trace consistency instead. Given the consistency on eigenvalue, we compute the Kullback-Leibler Divergence between probabilistic transition matrices to further prompt eigenvectors’ consistency. We also developed convenient ways to compute consistent loss regarding the complicated relationship between probabilistic transition matrices after image transform and modification.

The proposed method demonstrates consistent superiority to others on the common scribble-annotated dataset and is even comparable to some fully supervised ones. Moreover, we further conducted experiments when the scribble-

annotations gradually shrank and dropped. The proposed method could still work reasonably, even the scribble shrank to the point or dropped significantly. Besides, careful ablation study and mechanism study is made to verify the effectiveness of every module. Finally, the code and dataset are open-sourced.

Related Work

Scribble-Supervised Semantic Segmentation

The scribble-supervised semantic segmentation aims to produce dense predictions given only sparse scribble-annotations. Existing deep learning based works usually could be divided into two groups: 1) Two-stage approaches (Lin et al. 2016; Vernaza and Chandraker 2017), which firstly obtain the full mask pseudo-labels by manipulating the scribble annotations, then train the semantic segmentation network as usual with pseudo-labels. 2) Single-stage approaches (Tang et al. 2018a,b), which directly train the network using scribble-annotations by the specific design of loss function and network structure. While two-stage approaches can be formulated as regular semantic segmentation, single-stage approaches are usually defined to minimize L :

$$L = \sum_{p \in \Omega_{\mathcal{L}}} c(s_p, y_p) + \lambda \sum_{p, q \in \Omega} u(s_p, s_q), \quad (1)$$

where Ω is the pixel set, $\Omega_{\mathcal{L}}$ is the pixel set with scribble-annotations, s_i represents the prediction of pixel i , and y_i is the corresponding ground truth. The first term measures the error with respect to the scribble annotations and usually is in the form of cross-entropy. The second term is a pair-wise regularisation to help generate uniform prediction. The two terms are harmonized by a weight parameter λ .

For scribble-supervised semantic segmentation, the graphical model has been prevalently adopted in either two-stage approaches for generating pseudo-label or one-stage approaches for loss design. Lin et al. 2016 iteratively conduct label refinement and network optimization through a graphical model. Vernaza and Chandraker 2017 generate high-quality pseudo-labels for full-label supervised semantic segmentation by optimizing graphical model with edge detector learned from another well-labeled dataset. These two works require iterative optimization or auxiliary dataset.

Instead, Tang et al. 2018a,b add the soft graphical model regularization into loss function and avoid explicitly graphical model optimization. Besides, some works only well on the dataset with every existing object labeled by at least one scribble. In general, most methods have not provided a flexible and efficient solution to scribble supervised semantic segmentation yet.

Random Walk

Uniform neural representation is crucial for semantic segmentation to produce dense predictions no matter scribble supervised or full-label supervised. Typically, embedding a random walk operation in the network would provide help. In this way, Bertasius et al. use a random walk to address the issues of poor boundary localization and spatially fragmented predictions. Then Jiang et al. further conduct a sequence of random walks to approximate a stationary of the diffusion process. A random walk operation in the network can be defined as:

$$f(x)^L = \alpha P f(x)^{L-1} + f(x)^{L-1}, \quad (2)$$

where $f(x)^{L-1}$ is the neural representation of image x in layer $L-1$ and $f(x)^L$ is the neural representation after random walk in layer L . α is the weight parameter learned during training. The most key component of random walk is the probabilistic transition matrix P , whose unit p_{ij} measures the similarity between i -th and j -th elements of neural representations. Besides, all the units are positive, with every row of the matrix is summed to 1. Inner product, embedded gaussian (Wang et al. 2018), and diffusion distance (Sun and Xu 2019) have been widely used to compute the similarity between neural representations.

Self-Supervision

Consistency on neural representation is the property that should be concerned in almost all the learning-based tasks. Consistency would be helpful for detection, tracking, and definitely for segmentation. For this property, the consistent loss has been widely adopted. Since no ground truth required, this loss is more popular for unsupervised and semi-supervised tasks, and usually defined as the difference between neural representations of image and its transform:

$$ss(x, \phi) = l(T_\phi(f(x)), f(t_\phi(x))), \quad (3)$$

where t_ϕ denotes the transform operation on x with ϕ parameter, while T_ϕ corresponds to the transform operation on $f(x)$ (t_ϕ and T_ϕ are pair of corresponding transforms for self-supervision). l is the metric to define how the difference is measured. This kind of consistent loss is also referred to as self-supervision, which we denote as $ss(x, \phi)$. Self-supervision usually has to conduct two feed-forward processes, Laine and Aila 2016 propose temporal ensemble, which saves previous neural representations to avoid multiple feed-forward and facilitate the computation. Then, the mean teacher (Tarvainen and Valpola 2017) proposed to prepare a teacher network rather than save the auxiliary information. The self-supervision has been used for semi-supervised semantic segmentation (Mittal, Tatarchenko, and Brox 2019).

Method

The network of the proposed method is illustrated in Fig. 2, including two modules (ResNet backbone to extract features, and Similarity Measurement Module to compute the probabilistic transition matrix), one specific process (the random walk process) and three loss functions (the common cross-entropy loss, the maximum-entropy loss and self-supervised loss). In the following subsections, we will introduce the main components and discuss the insight.

Similarity Measurement Module

Similarity Measurement Module computes the distance between any pairs of neural representation elements and forms the probabilistic transition matrix. The module is illustrated in Fig. 2. There are several choices for the distance definition, such as Euclidean distance, inner product (Wang et al. 2018), and diffusion distance (Sun and Xu 2019). To preserve the property of the transition matrix (positive unit and each row sums to 1) and compute efficiency, in this work, we use the inner product with softmax operation to form the transition matrix. The probabilistic transition matrix P is defined as:

$$P = \text{softmax}(f(x)^{L-1T} f(x)^{L-1}), \quad (4)$$

where $f(x)^{L-1}$ is the neural representation of input x in layer $L-1$ and has flattened to $MN \times C$ dimension (M : height, N : width, C : feature channels). Thus, $f(x)^{L-1T} f(x)^{L-1}$ will produce a matrix of dimension $MN \times MN$. With softmax in the horizontal direction, we generate the adequate probabilistic transition matrix P .

Embedded Random Walk

We embed the random walk process in the computation flow of the neural network. The process is defined as Eq. 2 (with P as Eq. 4) and learned α to control the degree of random walk, and conducted on the final layer right before the classifier.

Through the random walk process, the i -th element in the neural representation of layer L , $f(x)_i^L$, equals to the weighted sum of all the elements in layer $L-1$ with the weight defined by the i -th row of P :

$$f(x)_i^L = \alpha \sum_{j=1}^{MN} P_{ij} f(x)_j^{L-1} + f(x)_i^{L-1}. \quad (5)$$

More similar $f(x)_j^{L-1}$ to $f(x)_i^{L-1}$, more large the P_{ij} is, and more important $f(x)_j^{L-1}$ to the $f(x)_i^L$.

With this process, each element in neural representation have relationships with all the other elements. Thus, the scribble annotations will affect not only the labeled elements but also unlabeled elements. When the cross-entropy loss is applied, the best solution (achieve lower loss) will be the uniform predictions within each object, considering the random walk process's constraint. In Fig. 1(b), we visualize the neural representations (after the sum of absolute values in the channel dimension) before and after the random walk process. As can be seen, the neural representations become

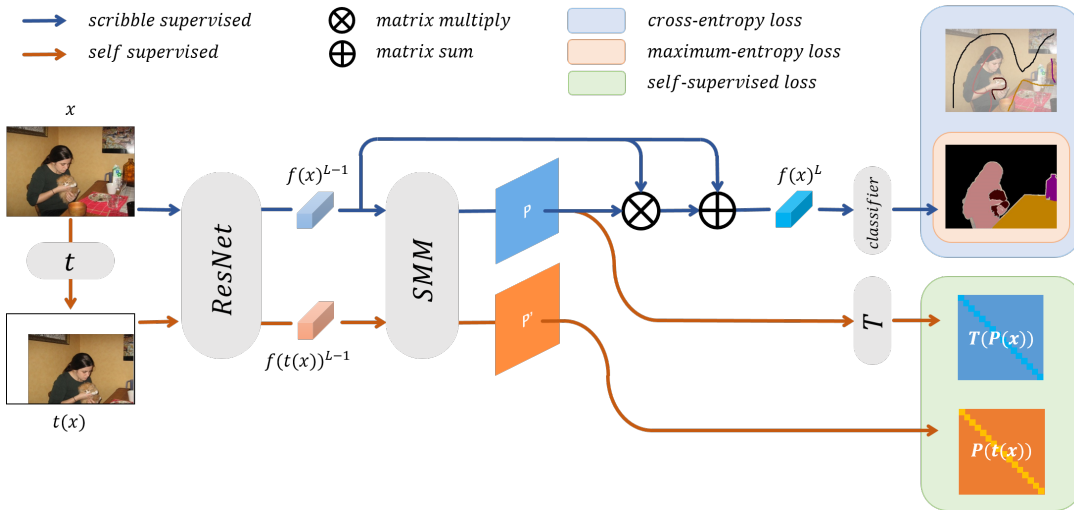


Figure 2: Network Architecture and Pipeline. We use blue flow to represent the scribble-supervised training and orange flow to represent self-supervised training. Given an image and its transform, we pass them to ResNet backbone to extract neural representations, from which the similarity measurement module (SMM) computes transition matrices. Then, a random walk is carried out on the neural representation of the original image. The results are used for classifying semantic. Simultaneously, the self-supervised loss is set between the transition matrices to realize the self-supervision on neural eigenspace. During inference, only blue flow is activated.

uniform within each semantic region after the random walk. It verifies our assumption in the introduction. It is worth noting that we have not imposed supervision on P , but P has gained semantic similarity knowledge from the results. It is the embedded random walk process with the scribble-annotations guide the formulation of P to produce uniform predictions.

Self-Supervision Loss

The self-supervision loss computes the difference between the neural representations of the image and its transform. There are several issues that need to be considered when applying self-supervision loss in this work: (1) where the self-supervision is involved; (2) how the self-supervision loss is calculated; (3) what kinds of the transform will be used. We address these issues in the following.

Self-Supervision on Eigenspace In our work, the typical choice of neural representation for self-supervision is $f(x)^L$, and thus the self-supervision loss will be

$$ss(x, \phi) = l(T_\phi(f(x)^L), f(t_\phi(x))^L). \quad (6)$$

However, as for semantic segmentation tasks with self-supervision, we argue that directly calculating loss on the whole neural representation is not necessary and may not optimal. When the image is distorted heavily after the transform, some parts of its neural representation will change greatly, so minimizing Eq. 6 will be hard and even ambiguous.

The transition matrix P could also be defined as $P = D^{-1}W$, where W is the affinity matrix, and D is the degree matrix. The eigenspace of P and the one of normalized Laplacian matrix L have close relationships, considering $L = D^{-1}(D - W)$. It can be proved that $\Lambda_P = 1 - \Lambda_L$

and $U_P = U_L$ (Λ denotes diagonal matrix with eigenvalues as entries, U denotes matrix with eigenvectors as columns). According to Von Luxburg 2007; Jiang, Vasconcelos, and Peng 2015; Jiang et al. 2019, columns of U_L have the capability to distinguish the main parts of the images. So, U_P will also inherit this property.

We visualize several eigenvectors of P in Fig. 1(c). As can be seen, compared with original neural representation, the eigenvectors of P are more powerful to distinguish the main parts from others and neglect some details, though P is also computed from neural representation. Based on the above analysis, in this work, we propose to set the self-supervision on the eigenspace of P :

$$ss(x, \phi) = l(T_\phi(U_P(x)), U_P(t_\phi(x))) + l(T_\phi(\Lambda_P(x)), \Lambda_P(t_\phi(x))). \quad (7)$$

Soft Eigenspace Self-Supervision Eq. 7 requires explicit eigendecomposition, which is time-consuming, especially within the deep neural network context. Though there are some approximation methods (Dang et al. 2018; Wang et al. 2019b; Sun and Xu 2019) proposed, their efficiency and stability are still far from satisfactory. To this end, we develop soft eigenspace self-supervision, which avoids explicit eigendecomposition. Firstly, in view of the fact that the matrix's trace is equal to the sum of its eigenvalues, we measure the consistency on the Λ by computing the difference on the trace of P , $tr(P)$. Secondly, given the consistency on the Λ , we propose to measure the consistency on the P to obtain consistent U indirectly. In other words, the soft eigenspace self-supervision loss is defined as:

$$ss_P(x, \phi) = l_1(T_\phi(P(x)), P(t_\phi(x))) + \gamma * l_2(T_\phi(tr(P(x))), tr(P(t_\phi(x)))) \quad (8)$$

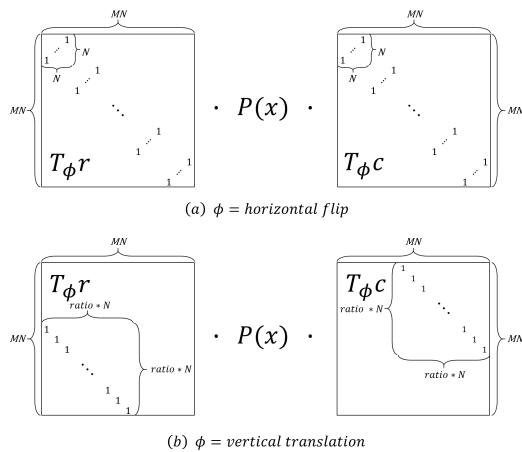


Figure 3: Visualization of predefined computing matrices for self-supervision on eigenspace.

where $P(x)$ denotes P for image x , $tr(P(x))$ is the trace of $P(x)$. Since $P(x)$ is the probabilistic transition matrix, we use Kullback-Leibler Divergence as l_1 to measure the difference. l_2 is defined as L_2 norm. γ is the weight to control the two terms.

Transform Operation and Computing Matrix In this work, we have linear transforms, including horizontal flip and translation, $\phi \in (\text{horizontal flip, translation})$. Compared with the transform effect on the neural representation, any transform will lead to a complex change on P and complicate the computation. However, since all the transform is linear, the probabilistic transition matrix after transform can be expressed as the multiplication of the original P with the predefined computing matrices, to facilitate the Eq. 8 computation. $T_\phi(P(x))$ can be defined as:

$$T_\phi(P(x)) = T_\phi r \cdot P(x) \cdot T_\phi c, \quad (9)$$

where $T_\phi r$ and $T_\phi c$ are predefined computing matrices for transform ϕ . In Fig. 3, we visualize computing matrices for horizontal flip and vertical translation when using soft eigenspace self-supervision. Please check supplementary for the detail definitions.

Maximum-entropy Loss

To force the network producing a more convinced prediction, we further minimize the maximum-entropy loss on the final prediction, which is defined as

$$E(s) = -\frac{1}{HW} \sum_{i,j} \sum_c (s(i,j,c) \cdot \log(s(i,j,c))), \quad (10)$$

where s represents the final prediction, and is of the size $H \times W \times C$ (C is the number of categories). $s(i,j,c)$ represents the probability that the pixel at position (i,j) of the image belongs to the c -th category.

Table 1: Ablation study on random walk, and operation and location of self-supervision.

Random Walk	Self-Supervision		mIoU
	operation	location	
	-	-	64.4
✓	-	-	67.6
✓	flip	$f(x)^{L-1}$	69.8
✓	flip	$f(x)^L$	70.1
✓	flip	Eigenspace	70.5
✓	translation	$f(x)^{L-1}$	70.5
✓	translation	$f(x)^L$	70.5
✓	translation	Eigenspace	70.8
✓	random	$f(x)^{L-1}$	70.2
✓	random	$f(x)^L$	70.3
✓	random	Eigenspace	71.2

Experiment

Implementation

The whole pipeline is shown in Fig. 2. We use pre-trained ResNet (He et al. 2016) with dilation (Chen et al. 2017) as the backbone to extract initial neural representations. The total loss in our work is defined as:

$$L = \sum_{p \in \Omega_{\mathcal{L}}} c(s_p, y_p) + \omega_1 * E(s) + \omega_2 * ss_P(x, \phi), \quad (11)$$

where ω_1 and ω_2 are the predefined weights. All the training images are randomly scaled (0.5 to 2), rotated (-10 to 10), blurred, and flipped for data augmentation, then cropped to 465×465 before feeding to the network. And the immediate output of *ResNet* ($f(x)^{L-1}$) is of spatial dimension 29×29 . The training process has two steps. Firstly, only the cross-entropy loss is used to train the network because the network may not perform well initially, and in this time, self-supervision will not bring benefits but prevent the optimization. After the network got reasonable performance, the whole Eq. 11 (scribble supervise and self supervise) is activated. The process has been visualized in Fig. 2. The two steps training is also used in (Tang et al. 2018a,b).

Experiment setting

Datasets We mainly compare with others on the common scribble-annotated dataset, *scribblesup* (Lin et al. 2016). This dataset is revised from pascal VOC (Everingham et al. 2015) and has every existing object in each image labeled by at least one scribble. However, our method does not need this hypothesis. To better demonstrate the advantage of the proposed method, we further proposed two variants of *scribblesup*. The first one is *scribble-drop*, where every object in images may drop (*i.e.* delete) all the scribble annotations with a defined possibility. The second one is *scribble-shrink*, where every scribble in the image is shrunk randomly (even to a point). All the images in two datasets are identical to the *scribblesup*'s, as well as the training and validation partition. In our experiments, many settings of drop and shrink rate are tested. Fig. 4 shows several representative samples of *scribble-drop* and *scribble-shrink*.

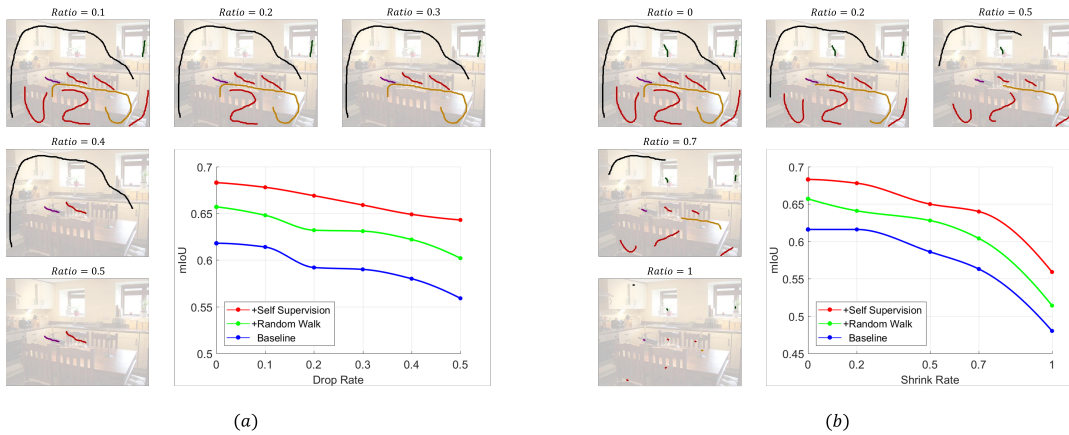


Figure 4: (a) A representative sample of *scribble-drop* with the scribble drop rate from 0.1 to 0.5, and the mIoU scores on different settings. (b) A representative sample of *scribble-shrink* with the scribble shrink rate from 0 to 1 (point), and the mIoU scores on different settings. (Zoom in for better visualization)

Table 2: Performance on the validation set of pascal VOC. The supervision types (Sup.) indicate: \mathcal{P} –point, \mathcal{S} –scribble, \mathcal{B} –bounding box, \mathcal{I} –image-level label, and \mathcal{F} –full label.

Method	Sup.	Backbone	wo/ CRF	w/ CRF
What’sPoint	\mathcal{P}	<i>VGG16</i>	46.0	-
SDI	\mathcal{B}	<i>ResNet101</i>	-	69.4
BCM	\mathcal{B}	<i>ResNet101</i>	-	70.2
CIAN	\mathcal{I}	<i>ResNet101</i>	64.1	67.3
FickleNet	\mathcal{I}	<i>ResNet101</i>	64.9	-
SCE	\mathcal{I}	<i>ResNet101</i>	64.8	66.1
DeepLabV2	\mathcal{F}	<i>ResNet101</i>	76.4	77.7
scribblesup	\mathcal{S}	<i>VGG16</i>	-	63.1
RAWKS	\mathcal{S}	<i>ResNet101</i>	59.5	61.4
NCL	\mathcal{S}	<i>ResNet101</i>	72.8	74.5
KCL	\mathcal{S}	<i>ResNet101</i>	73.0	75.0
BPG-PRN	\mathcal{S}	<i>ResNet101</i>	71.4	-
ours-ResNet50	\mathcal{S}	<i>ResNet50</i>	71.9	73.6
ours-ResNet101	\mathcal{S}	<i>ResNet101</i>	73.4	75.8

Compared methods We compare with recently proposed scribble-supervised methods including scribblesup (Lin et al. 2016), RAWKS (Vernaza and Chandraker 2017), NCL (Tang et al. 2018a), KCL (Tang et al. 2018b) and BPG-PRN (Wang et al. 2019a), and also other weakly-supervised methods such as point supervised (What’sPoint (Bearman et al. 2016)), bounding-box supervised (SDI (Khoreva et al. 2017), BCM (Song et al. 2019)) and image-level-label supervised (CIAN (Fan et al. 2020), FickleNet (Lee et al. 2019), SCE (Chang et al. 2020)). Besides, full-label supervised method (DeepLabV2 (Chen et al. 2017)) is also compared. We use *mIoU* as the main metric to evaluate these methods and ours. When comparing with others, we mainly refer their reported scores if available.

Hyper-parameters The proposed method has 200 epochs of training, with the first 100 epochs have no self-supervised loss. For every step, eight images (batch size) are randomly selected to train the network with Adam (Kingma and Ba

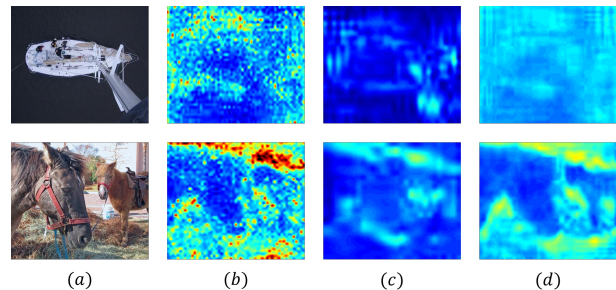


Figure 5: Visualization of variation by different self-supervision operations. (a) The input image, (b) The variation on $f(x)^{L-1}$, (c) The variation on $P(x)$, (d) The variation on $f(x)^L$. The first row shows the variation for the flip operation, while the second row is for the translation operation.

Table 3: Variation comparison under the same transform.

	$f(x)^{L-1}$	$f(x)^L$	$P(x)$
flip	25.3%	27.7%	11.1%
translation	7.7%	15.2%	6.9%

2014) optimizer, Sync-BatchNorm (Ioffe and Szegedy 2015) and learning rate as $1e-3$ for the first 100 epochs and $1e-4$ for the rest. The weights γ , ω_1 and ω_2 are set to be 0.01, 0.2 and 1, respectively. Besides, as the common way for semantic segmentation (Zhao et al. 2018), data augmentation is adopted during training. All the computations are carried out on NVIDIA TITAN RTX GPUs.

Ablation Study

In Tab. 1, we do ablation study w/wo random walk, w/wo self-supervision, and self-supervision on $f(x)^{L-1}$, $f(x)^L$ and $P(x)$ on *Scribblesup* dataset. We use *ResNet50* as the backbone without maximum-entropy loss and CRF to study

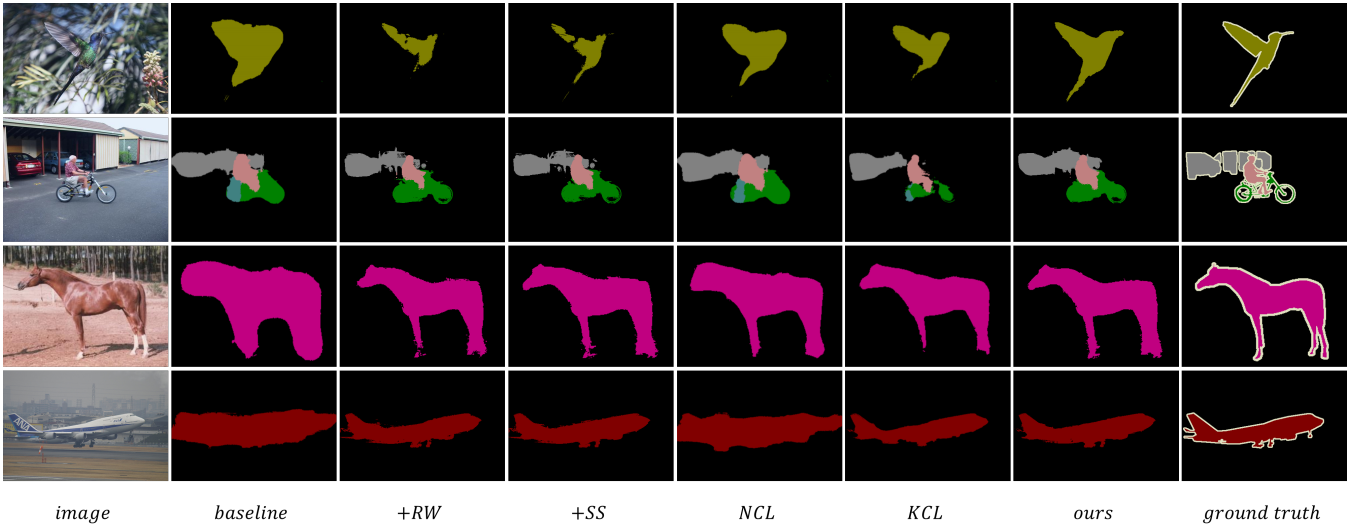


Figure 6: Visual comparison between proposed method and others on *scribblesup* dataset.

Table 4: The performance drop ratios compared to no drop and no shrink when only using baseline and gradually adding random walk and self-supervision.

drop rate	0.1	0.2	0.3	0.4	0.5	shrink rate	0.2	0.5	0.7	1
baseline	0.7%	4.2%	4.5%	6.2%	9.6%	baseline	0%	4.9%	8.6%	22.1%
+Random Walk	1.4%	3.8%	4.0%	5.3%	8.4%	+Random Walk	2.4%	4.4%	8.1%	21.8%
+Self-Supervision	0.7%	2.1%	3.5%	5.0%	5.9%	+Self-Supervision	0.7%	4.8%	6.3%	18.2%

the random walk and self-supervision more thoroughly. The first one is a fully convolutional neural network (Baseline). With random walk, the mIoU receives more than 3% promotion. After incorporating self-supervision (the last one), the performance further boosted by 3.6% increasing. As for the self-supervision operations, we observe that both are helpful and self-supervision on the eigenspace outperforms on the other locations consistently no matter what kind of operation applied. Besides, randomly combined on the eigenspace leads to the best performance, while other locations do not show such promotion. It is worth noting that eigenspace is also a kind of neural representation that lies in the middle between $f(x)^{L-1}$ and $f(x)^L$. The middle one always outperforms sides, indicating a general advantage of self-supervision on neural eigenspace.

In Tab. 3, we show the mean variations of $f(x)^{L-1}$, $f(x)^L$ and $P(x)$ under the same transform (no self-supervision applied yet). The variation is measured by the relative error defined as $|T_\phi(f) - f'| / (|T_\phi(f)| + |f'|)$ (f : feature, f' : feature after transform). As can be seen, the same transform will always lead to less change on $P(x)$. Considering the property of eigenvectors, we believe $P(x)$ is not sensitive to the trivial structures (whose semantic may be heavily changed after transform). In Fig. 5, we visualize the variation of two images by different self-supervision operations, respectively. It can be seen the variation on $P(x)$ are most distributed on the main objects or object boundaries, while $f(x)^{L-1}$ and $f(x)^L$ also highlight backgrounds with uneven distri-

bution. This phenomenon indicates that self-supervision on $P(x)$ will be more effortless and less confusing.

Quantitative Results

After introducing maximum-entropy loss, our method gets 71.9% mIoU with *ResNet50* and 73.4% with *ResNet101* on *Scribblesup* dataset. When comparing with others, we also report the performance with CRF as others. Tab. 2 lists all the scores of compared methods under different settings. In addition to the scribble-supervised methods, we also show methods with other type labels, including point, bounding-box, image-level-label and full-label. The proposed method reaches state-of-the-art performance compared with scribble and other weakly supervised methods and is even comparable to the full-label one. The reported full-label method (DeepLabV2) had been pre-trained on COCO dataset.

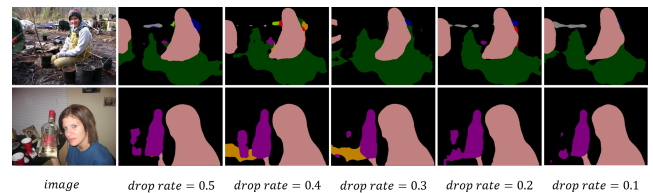


Figure 7: Results of proposed method on *scribble-drop* dataset with different drop rate.

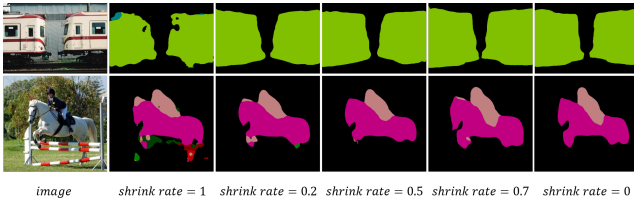


Figure 8: Results of proposed method on *scribble-shrink* dataset with different shrink rate.

It should be noted that some methods, such as Lin et al. 2016; Vernaza and Chandraker 2017, require every existing object in the image labeled. However, ours does not have this limit. To evaluate how the performance affected when the scribble randomly dropped and shrank, we further prepare two datasets, *scribble-drop* and *scribble-shrink* modified from *scribblesup* as described before. We conduct experiments with different drop and shrink rate and show the mIoU scores in Fig. 4. Besides, in Tab. 4, we also show the performance drop ratios under different drop and shrink rates. The methods used are identical to the ones in Tab. 1. It can be seen that the proposed method demonstrates certain robustness when the drop rate and shrink rate increase, even all the scribbles were shrunk to points (point-supervised).

Qualitative Results

We show visual comparison in Fig. 6, Fig. 7 and Fig. 8. Fig. 6 presents results of *NCL*, *KCL* and ours on *scribblesup* dataset. With the proposed random walk (RW) and self-supervision on eigenspace (SS), the results are gradually refined, and the complete method outperforms others clearly and shows significant promotion over the baseline. Fig. 7 and Fig. 8 further demonstrate our results on *scribble-drop* and *scribble-shrink*. Here no CRF using for better evaluation. It can be seen some details are missing when the annotations for training are gradually shrunk, but the main parts are preserved well. As for the random drop, our method shows promising robustness. When each scribble was dropped with 50% probability during training, the prediction does not degrade much. (The results and scores in this section are all from the validation set. Please check supplementary for more results.)

Efficiency

We analyze the efficiency after adding random walk and self-supervision in Tab. 5. Since the self-supervision branch shares the same trainable-parameters with the main branch, no extra parameters are needed. Besides, the self-supervision output will act as the target, so no intermediate features need to store. Finally, self-supervision is only conducted during training. Consequently, in our implementation, the time and memory cost by random walk and self-supervision is acceptable.

Conclusion

In this work, we present a scribble-supervised semantic segmentation method. The insight is to guide the network to

Table 5: Trainable-parameters (P), memory (M), and inference speed (S) statistics when only using baseline and changes after gradually adding random walk and self-supervision, respectively.

	P	M	S
Baseline	23.61 M	1090.82 MB	5.87 it/s
+Random Walk	+0.31 M	+9.3 MB	-0.07 it/s
+Self-supervision	+0 M	+2.7 MB	-0 it/s

produce uniform and consistent predictions by embedding a random walk process on the neural representation and imposing a self-supervision on the neural eigenspace. Thoroughly ablation studies and intermediate visualization have verified the effectiveness of proposed components. Finally, the complete method reaches state-of-the-art performance compared with others, even is comparable to the full-label supervised ones. Moreover, the proposed method shows robustness to the data with randomly dropped or shrank labels.

References

- Bearman, A.; Russakovsky, O.; Ferrari, V.; and Fei-Fei, L. 2016. What’s the point: Semantic segmentation with point supervision. In *European conference on computer vision*, 549–565. Springer.
- Bertasius, G.; Torresani, L.; Yu, S. X.; and Shi, J. 2017. Convolutional random walk networks for semantic image segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 858–866.
- Chang, Y.-T.; Wang, Q.; Hung, W.-C.; Piramuthu, R.; Tsai, Y.-H.; and Yang, M.-H. 2020. Weakly-Supervised Semantic Segmentation via Sub-Category Exploration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8991–9000.
- Chen, L.-C.; Papandreou, G.; Kokkinos, I.; Murphy, K.; and Yuille, A. L. 2017. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence* 40(4): 834–848.
- Dang, Z.; Moo Yi, K.; Hu, Y.; Wang, F.; Fua, P.; and Salzmann, M. 2018. Eigendecomposition-free training of deep networks with zero eigenvalue-based losses. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 768–783.
- Everingham, M.; Eslami, S. A.; Van Gool, L.; Williams, C. K.; Winn, J.; and Zisserman, A. 2015. The pascal visual object classes challenge: A retrospective. *International journal of computer vision* 111(1): 98–136.
- Fan, J.; Zhang, Z.; Tan, T.; Song, C.; and Xiao, J. 2020. CIAN: Cross-Image Affinity Net for Weakly Supervised Semantic Segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the*

- IEEE conference on computer vision and pattern recognition*, 770–778.
- Ioffe, S.; and Szegedy, C. 2015. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. In *International Conference on Machine Learning*, 448–456.
- Jiang, P.; Gu, F.; Wang, Y.; Tu, C.; and Chen, B. 2018. Difnet: Semantic segmentation by diffusion networks. In *Advances in Neural Information Processing Systems*, 1630–1639.
- Jiang, P.; Pan, Z.; Tu, C.; Vasconcelos, N.; Chen, B.; and Peng, J. 2019. Super diffusion for salient object detection. *IEEE Transactions on Image Processing* 29: 2903–2917.
- Jiang, P.; Vasconcelos, N.; and Peng, J. 2015. Generic Promotion of Diffusion-Based Salient Object Detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- Khoreva, A.; Benenson, R.; Hosang, J.; Hein, M.; and Schiele, B. 2017. Simple does it: Weakly supervised instance and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 876–885.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Laine, S.; and Aila, T. 2016. Temporal ensembling for semi-supervised learning. *arXiv preprint arXiv:1610.02242*.
- Lee, J.; Kim, E.; Lee, S.; Lee, J.; and Yoon, S. 2019. Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 5267–5276.
- Lin, D.; Dai, J.; Jia, J.; He, K.; and Sun, J. 2016. Scribble-sup: Scribble-supervised convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3159–3167.
- Long, J.; Shelhamer, E.; and Darrell, T. 2015. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3431–3440.
- Mittal, S.; Tatarchenko, M.; and Brox, T. 2019. Semi-supervised semantic segmentation with high-and low-level consistency. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Nadler, B.; and Galun, M. 2007. Fundamental limitations of spectral clustering. In *Advances in neural information processing systems*, 1017–1024.
- Ng, A. Y.; Jordan, M. I.; and Weiss, Y. 2002. On spectral clustering: Analysis and an algorithm. In *Advances in neural information processing systems*, 849–856.
- Song, C.; Huang, Y.; Ouyang, W.; and Wang, L. 2019. Box-Driven Class-Wise Region Masking and Filling Rate Guided Loss for Weakly Supervised Semantic Segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Sun, J.; and Xu, Z. 2019. Neural Diffusion Distance for Image Segmentation. In *Advances in Neural Information Processing Systems*, 1443–1453.
- Tang, M.; Djelouah, A.; Perazzi, F.; Boykov, Y.; and Schroers, C. 2018a. Normalized cut loss for weakly-supervised cnn segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 1818–1827.
- Tang, M.; Perazzi, F.; Djelouah, A.; Ben Ayed, I.; Schroers, C.; and Boykov, Y. 2018b. On regularized losses for weakly-supervised cnn segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 507–522.
- Tarvainen, A.; and Valpola, H. 2017. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *Advances in neural information processing systems*, 1195–1204.
- Vernaza, P.; and Chandraker, M. 2017. Learning random-walk label propagation for weakly-supervised semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7158–7166.
- Von Luxburg, U. 2007. A tutorial on spectral clustering. *Statistics and computing* 17(4): 395–416.
- Wang, B.; Qi, G.; Tang, S.; Zhang, T.; Wei, Y.; Li, L.; and Zhang, Y. 2019a. Boundary Perception Guidance: A Scribble-Supervised Semantic Segmentation Approach. In *IJCAI*, 3663–3669.
- Wang, W.; Dang, Z.; Hu, Y.; Fua, P.; and Salzmann, M. 2019b. Backpropagation-friendly eigendecomposition. In *Advances in Neural Information Processing Systems*, 3162–3170.
- Wang, X.; Girshick, R.; Gupta, A.; and He, K. 2018. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 7794–7803.
- Yang, C.; Zhang, L.; Lu, H.; Ruan, X.; and Yang, M.-H. 2013. Saliency detection via graph-based manifold ranking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 3166–3173.
- Zhao, H.; Zhang, Y.; Liu, S.; Shi, J.; Change Loy, C.; Lin, D.; and Jia, J. 2018. Psanet: Point-wise spatial attention network for scene parsing. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 267–283.