# Semi-supervised Semantic Segmentation via Strong-weak Dual-branch Network

Wenfeng Luo[1] and Meng Yang*[1,2]

[1] School of Data and Computer Science, Sun Yat-sen University, Guangzhou, China
`luowf5@mail2.sysu.edu.cn`
[2] Key Laboratory of Machine Intelligence and Advanced Computing (SYSU),
Ministry of Education `yangm6@mail.sysu.edu.cn`
*Corresponding author

**Abstract.** While existing works have explored a variety of techniques to push the envelop of weakly-supervised semantic segmentation, there is still a significant gap compared to the supervised methods. In real-world application, besides massive amount of weakly-supervised data there are usually a few available pixel-level annotations, based on which semi-supervised track becomes a promising way for semantic segmentation. Current methods simply bundle these two different sets of annotations together to train a segmentation network. However, we discover that such treatment is problematic and achieves even worse results than just using strong labels, which indicates the misuse of the weak ones. To fully explore the potential of the weak labels, we propose to impose separate treatments of strong and weak annotations via a strong-weak dual-branch network, which discriminates the massive inaccurate weak supervisions from those strong ones. We design a shared network component to exploit the joint discrimination of strong and weak annotations; meanwhile, the proposed dual branches separately handle full and weak supervised learning and effectively eliminate their mutual interference. This simple architecture requires only slight additional computational costs during training yet brings significant improvements over the previous methods. Experiments on two standard benchmark datasets show the effectiveness of the proposed method.

**Keywords:** Semi-supervised, Strong-weak, Semantic Segmentation

## 1 Introduction

Convolutional Neural Networks (CNNs) [17, 30, 11] have proven soaring successes on the semantic segmentation problem. Despite their superior performance, these CNN-based methods are data-hungry and rely on huge amount of pixel-level annotations, whose collections are labor-intensive and time-consuming. Hence researchers have turned to *weakly-supervised learning* that could exploit weaker forms of annotation, thus reducing the labeling costs. Although numerous works [15, 36, 13, 19] have been done on learning segmentation models from weak supervisions, especially per-image labels, they still trail the accuracy of their fully-supervised counterparts and thus are not ready for real-world applications.
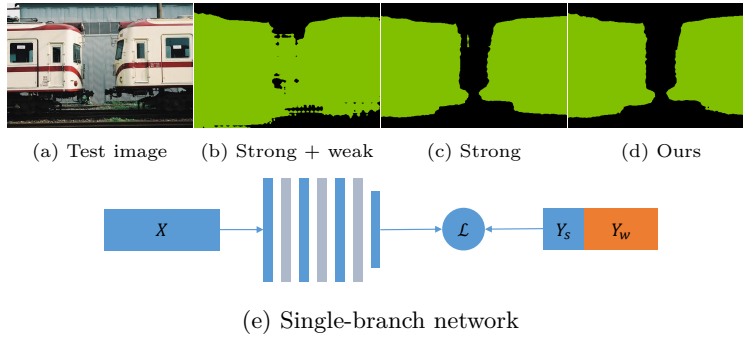
(a) Test image      (b) Strong + weak      (c) Strong      (d) Ours



(e) Single-branch network

**Fig. 1.** (a) Sample test image; (b) Result using both strong and weak annotations in a single-branch network; (c) Result using only the strong annotations; (d) Result using our strong-weak dual-branch network; (e) Single-branch network adopted by previous methods [26, 35, 19]; Images (a)(b)(c)(d) in the first row visually demonstrate that using extra weak annotations brings no improvement over only using the strong annotations when a single-branch network is employed. See Fig.7 for more visual comparisons.

In order to achieve good accuracy while still keeping the labeling budget in control, we focus on tackling a more practical problem under semi-supervised setting, where a combination of strongly-labeled (pixel-level masks) and weakly-labeled (image-level labels) annotations are utilized. However, previous methods (WSSL [26], MDC [35] and FickleNet[19]) simply scratch the surface of semi-supervised segmentation by exploring better weakly-supervised strategies to extract more accurate initial pixel-level supervisions, which are then mixed together with strong annotations to learn a segmentation network, as in Fig.1(e). However, we discover that the simple combination of strong and weak annotations with equal treatment may weaken the final performance.
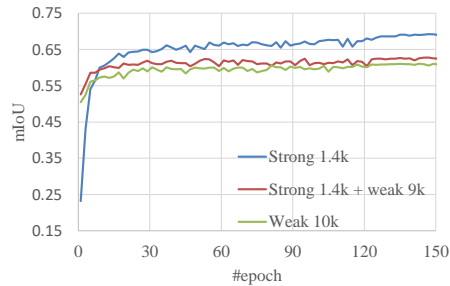


**Fig. 2.** Segmentation performance of the conventional single-branch network on *val* set using different training data. Here, the DSRG [13] is used to estimate the weak supervisions. The single-branch network supervised by the 1.4k strong data achieved much better result than that by the extra 9k weak annotations.

To further analyze the roles of strong and weak annotations in conventional single-branch network, we compare the segmentation performance on PASCAl VOC val set using different training data in Fig.2. When trained on small amount of strong data (1.4k in our experiments), performance of the segmentation network is not as low as people would expect. On the contrary, our implementation achieves a peak mIoU of 68.9% using only 1.4k strong annotations and it is already much better than other methods WSSL [26](64.6%), MDC [35](65.7%), FickleNet [19](65.8%) exploiting extra 9k weak annotations. Moreover, when simply bundling the strong and weak annotations to train a single segmentation network, the performance is not better than that using only the strong ones (quantatively shown in Fig.2 and visually shown in Fig.1(a)(b)(c)).

Based on the above observations, it can be concluded that such treatment underuses the weakly-supervised data and thus introduces limited improvement, or even worse, downgrading the performance achieved by using only the strong annotations. We further point out two key issues that are previously unnoticed concerning the semi-supervised setting:

1. *sample imbalance*: there are usually much more weak data than the strong ones, which could easily result in overfitting to the weak supervisions.
2. *supervision inconsistency*: the weak annotations $Y_w$ are of relatively poor quality compared to the strong ones $Y_s$ and thus lead to poor performance.

To better jointly use the strong and weak annotations, we propose a novel method of strong-weak dual-branch network, which is a single unified architecture with parallel strong and weak branches to handle one type of annotation data (Fig.3). To fully exploit the joint discrimination of strong and weak annotations, the parallel branches share a common convolution backbone in exchange for supervision information of different level without competing with each other. The shared backbone enables the free flow of the gradient and the parallel branches can discriminate between the accurate and noisy annotations. Moreover, the dual branches are explicitly learned from strong and weak annotations separately, which can effectively avoid the affect of sample imbalance and supervision inconsistency. This simple architecture boosts the segmentation performance by a large margin while introducing negligible overheads. State-of-art performance has been achieved by the proposed strong-weak network under semi-supervised setting on both PASCAL VOC and COCO segmentation benchmarks. Remarkably, it even boosts the fully-supervised models when both branches are trained with strong annotations on PASCAL VOC.

The main contributions of our paper are three-folds:

1. We for the first time show that segmentation network trained under mixed strong and weak annotations achieves even worse results than using only the strong ones.
2. We reveal that sample imbalance and supervision inconsistency are two key obstacles in improving the performance of semi-supervised semantic segmentation.
3. We propose a simple unified network architecture to address the inconsistency problem of annotation data in the semi-supervised setting.

## 2   Related works

In this section, we briefly review weakly-supervised and semi-supervised visual learning, which are most related to our work. Although the idea of multiple branches to capture various context has already been explored in many computer vision tasks, we here highlight the primary difference between previous methods and this work.

**Weakly-supervised semantic segmentation** To relieve the labeling burden on manual annotation, many algorithms have been proposed to tackle semantic segmentation under weaker supervisions, including points [2], scribbles [21, 32], and bounding boxes [7, 31]. Among them, per-image class labels are most frequently explored to perform pixel-labeling task since their collections require the least efforts, only twenty seconds per image [2]. Class Activation Map (CAM) [39], is a common method to extract from classification network a sparse set of object seeds, which are known to concentrate on small discriminative regions. To mine more foreground pixels, a series of methods have been proposed to apply the erasing strategy, either on the original image [36] or high-level class activations [12]. Erasing strategy is a form of strong attention [37] which suppresses selective responsive regions and forces the network to find extra evidence to support the corresponding task. Some other works [25, 4, 34] also proposed to incorporate saliency prior to ease the localization of foreground objects.

Recently, Huang *et al.*[13] proposed Deep Seeded Region Growing (DSRG) to dynamically expand the discriminative regions along with the network training, thus mining more integral objects for segmentation networks. And Lee *et al.*[19] further improved the segmentation accuracy of DSRG by replacing the original CAM with stochastic feature selection for seed generation.

Despite the progress on weakly-supervised methods, there is still a large performance gap (over 10%) from their full-supervised versions [5, 6], which indicates that they are unsuitable for the real-world applications.

**Semi-supervised learning** In general, semi-supervised learning [40] addresses the classification problem by incorporating large amount of extra unlabeled data besides the labeled samples to construct better classifiers. Besides earlier methods, like semi-supervised Support Vector Machine [3], many techniques have been proposed to integrate into deep-learning models, such as Temporal Ensembling [18], Virtual Adversarial Training [24] and Mean Teacher [33].

In this paper, we focus on such *semi-supervised learning* setting on semantic segmentation problem, where the training data are composed of a small set of finely-labeled data and large amount of coarse annotations, usually estimated from a weakly-supervised methods. In this configuration, current models [27, 20, 35, 19] usually resort to the sophistication of weakly-supervised models to provide more accurate proxy supervisions and then simply bundle both sets of data altogether to learn a segmentation network. They pay no special attention to coordinating the usage of weak annotations with the strong ones. Such treatment, ignoring the annotation inconsistency, overwhelms the handful yet vital minority and consequently produces even worse results compared to using only the fine data.

**Multi-branch network** Networks with multiple parallel branches have been around for a long time and proven their effectiveness in a variety of vision-related tasks. Object detection models [9, 28, 22] usually ended with two parallel branches, one for the classification and the other for localization. In addition, segmentation networks, such as Atrous Spatial Pyramid Pooling (ASPP) [5] and Pyramid Scene Parsing (PSP) [38] network, explored multiple parallel branches to capture richer context to localize objects of different sizes. Unlike the above works, we instead utilize parallel branches to handle different types of annotation data.

## 3 Methods

As aforementioned, the proxy supervisions estimated by weakly-supervised methods are of relatively poor quality in contrast to manual annotations. For finely-labeled and weakly-labeled semantic segmentation task, a natural solution for different supervision is to separately train two different networks, whose outputs are then aggregated by taking the average (or maximum). Although this simple ensemble strategy is likely to boost the performance, it is undesirable to maintain two copies of network weights during both training and inference. Besides, separate training prohibits the exchange for supervision information. To enable information sharing and eliminate the sample imbalance and supervision inconsistency, we propose a dual-branch architecture to handle different types of supervision, eliminating the necessity of keeping two network copies. Fig.3 presents an overview of the proposed architecture.

**Notation** Let the training images $X = (x_i)_{i \in [n]}$ be divided into two subsets: the images $X_s = \{(x_1^s, m_1^s), ..., (x_t^s, m_t^s)\}$ with strong annotations provided by the dataset and images $X_w = \{(x_1^w, m_1^w), ..., (x_k^w, m_k^w))\}$, the supervisions of which are estimated from a proxy ground-truth generator $G$:

$$m_i^w = G(x_i^w) \tag{1}$$

The proxy generator $G$ may need some extra information, such as class labels, to support its decision, but we can leave it for general discussion.

The rest of this section is organized as follows. Section 3.1 discusses in depth why training a single-branch network is problematic. Section 3.2 elaborates on the technical details of the proposed strong-weak dual-branch network.

### 3.1 Oversampling Doesn't Help with Single-branch Network

Previous works [27, 20, 35, 19] focus on developing algorithm to estimate more accurate initial supervision, but they pay no special attention on how to co-ordinate the strong and weak annotations. Notably, there are quite a few estimated masks of relatively poor quality (as shown Fig.4) when image scenes become more complex. Equal treatment biases the gradient signal towards the incorrect weak annotations since they are in majority during the computation
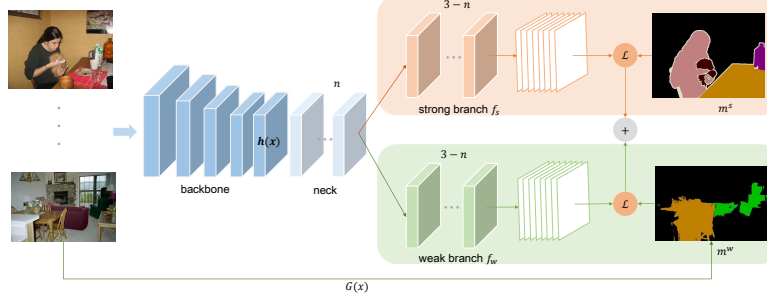
**Fig. 3.** Overview of the proposed dual-branch network. The proposed architecture consists of three individual parts: backbone, neck module and two parallel branches that share an identical structure but differ in the training annotations. The hyper-parameter $n$ controls the number (i.e.,3-$n$) of individual convolutional layers existing in the parallel branches.



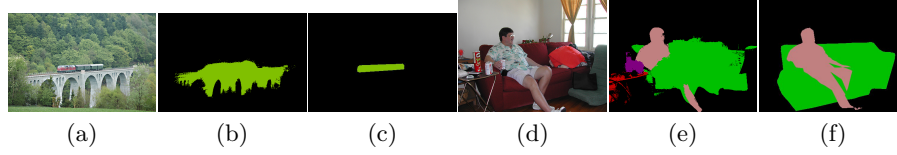| (a) | (b) | (c) | (d) | (e) | (f) |

**Fig. 4.** Two inaccurate weak annotations estimated by DSRG [13]. (a)(d): sample training images; (b)(e): masks estimated by DSRG; (c)(f): ground truth. In (b), the *train* mask expands to the background due to color similarity. In (e), large portions of the *human* body are misclassified.

of the training loss. Consequently, it offsets the correct concept learned from the strong annotations and therefore leads to performance degradation. In addition, we also conduct experiments via oversampling the strong annotations. The results shows that oversampling does improve the final segmentation accuracy steadily (62.8%→65.9%) as more strong annotations are duplicated, but it still fails to outperform the result (68.9%) using only the strong annotations. In conclusion, oversampling does not help with the single-branch network either. See the supplementary material for detailed results.

### 3.2 Strong-weak dual-branch network

**Network architecture** The strong-weak dual-branch network consists of three individual parts: convolutional backbone, neck module and two parallel branches with identical structure. Since our main experiments centre around the VGG16 network [30], we here give a detailed discussion of the architecture based on VGG16.

   **Backbone** The backbone is simply the components after removing the fully-connected layers. As in [5], the last two pooling layers are dropped and the

dilation rates of the subsequent convolution layers are raised accordingly to obtain features of output stride 8.

**Neck module** The neck module is a series of convolution layers added for better adaptation of the specific task. It could be shared between or added separately into subsequent parallel branches. Let $n$ be the number of convolution layers in the neck module shared by different supervision. Although the design of common components is simple, the backbone and the first $n$-layer neck module can effectively learn the joint discrimination from the full supervision and weak supervision. The total number of convolution layers in the neck and subsequent branch is fixed, but the hyper-parameter $n \in [0, 3]$ offers greater flexibility to control the information sharing. When $n$ is 0, each downstream branch has its own neck module. We denote the network and its output up until the neck module as $Z = h(X) \in R^{H \times W \times K}$.

**Strong-weak branches** These two parallel branches have the same structure while differ in the training annotations they receive. The strong branch is supervised by the fine annotation $X_s$, while the weak branch is trained by the coarse supervisions $X_w$. The way of separately processing different supervision is quite new because existing semi-supervised semantic segmentation methods adopt a single-branch network and current multi-branch network has never dealt with different types of annotation. The branches $f(Z; \theta_s)$ and $f(Z; \theta_w)$ are governed by independent sets of parameters. For brevity, we will omit the parameters in our notation and simply write $f_s(Z)$ and $f_w(Z)$. The normal cross entropy loss has the following form:

$$\mathcal{L}_{ce}(s, m) = -\frac{1}{|m|} \sum_c \sum_{u \in m_c} \log s_{u,c} \tag{2}$$

where tensor $s$ is the network outputs, $m$ is the annotation mask and $m_c$ denotes the set of pixels assigned to category $c$. Then the data loss of our method is:

$$
\begin{aligned}
s^s &= f_s(h(x^s)) \\
s^w &= f_w(h(x^w)) \\
\mathcal{L}_{data} &= \mathcal{L}_{ce}\left(s^s, m^s\right) + \mathcal{L}_{ce}\left(s^w, m^w\right)
\end{aligned}
\tag{3}
$$

We emphasize that all the loss terms are equally weighted so no hyper-parameter is involved.

### 3.3   How does the dual-branch network help?

During training, we need to construct a training batch with the same amount of strong and weak images. As in semi-supervised semantic segmentation, there are usually much more weak annotations than the strong ones. Consequently, the strong data have been looped through several times before the weak data are exhausted for the first pass, which essentially performs oversampling of the strong data. In this way, the strong data make a difference during training and thus mitigate the effect of sample imbalance.

Our dual-branch network imposes separate treatments on the strong and weak annotations and therefore prevents direct interference of different supervision information, so the supervision inconsistency can be well eliminated. Meanwhile, the coarse ones leave no direct influence on the strong branch, which determines the final prediction. Nevertheless, the extra weak annotations provide approximate location of objects and training them on a separate branch introduces regularization into the underlying backbone to some extent, hence improving the network's generalization capability.

### 3.4   Implementation detail

**Training** Here we introduce an efficient way to train the strong-weak dual-branch network. A presentation of the processing details can be found in Fig.5. During training, a batch of $2n$ images $X = [(x_1^s, m_1^s), ..., (x_1^w, m_1^w)...]$ are sampled, with the first half from $X_s$ and the second from $X_w$. Since the number of weak annotations is usually much bigger than that of strong ones, we are essentially performing an oversampling of the strong annotations $X_s$. For the image batch $X \in R^{2n \times h \times w}$, we make no distinction of the images and simply obtain the network logits in each branch, namely $S^s, S^w \in R^{2n \times h \times w}$, but half of them (in color gray Fig.5) have no associated annotations and are thus discarded. The remaining halves are concatenated to yield the final network output $S = [S^s[1 : n], S^w[n + 1 : 2n]]$, which are then used to calculate the cross entropy loss irrespective of the annotations employed. We find that this implementation eases the training and inference processes.



**Fig. 5.** The images are first forwarded through the network and half of the outputs (in color gray) are dropped before they are concatenated to compute the final loss (only the batch dimension is shown).

**Inference** When the network is trained, the weak branch is no longer needed since the information from weak annotations has been embedded into the convolution backbone and the shared neck module. So at inference stage, only the strong branch is utilized to generate final predictions.

## 4 Experiments

### 4.1 Experimental setup

**Dataset and evaluation metric** The proposed method is evaluated on two segmentation benchmarks, PASCAL VOC [8] and COCO dataset [23]. **PASCAL VOC**: There are 20 foreground classes plus 1 background category in PASCAL VOC dataset. It contains three subsets for semantic segmentation task, *train* set (1464 images), *val* set (1449 images) and *test* set (1456 images). As a common practice, we also include the additional annotations from [10] and end up with a *trainaug* set of 10582 images. For semi-supervised learning, we use the *train* set as the strong annotations and the remaining 9k images as weak annotations. We report segmentation results on both *val* and *test* set. **COCO**: We use the train-val split in the 2016 competition, where 118k images are used for training and the remaining 5k for testing. We report the segmentation performance on the 5k testing images.

The standard interaction-over-union (IoU) averaged across all categories is adopted as evaluation metric for all the experiments.

**Proxy supervision generator** $G$ To verify the effectiveness of the proposed architecture, we choose the recently popular weakly-supervised method, Deep Seeded Region Growing (DSRG) [13], as the proxy supervision generator $G$. We use the DSRG model before the retraining stage to generate proxy ground truth for our experiments. Further details could be found in the original paper.

**Training and testing settings** We use the parameters pretrained on the 1000-way ImageNet classification task to initialize our backbones (either VGG16 or ResNet101). We use Adam optimizer [14] with an initial learning rate of 1e-4 for the newly-added branches and 5e-6 for the backbone. The learning rate is decayed by a factor of 10 after 12 epochs. The network is trained under a batch size of 16 and a weight decay of 1e-4 for 20 epochs. We use random scaling and horizontal flipping as data augmentation and the image batches are cropped into a fixed dimension of $328 \times 328$.

In test phase, we use the strong branch to generate final segmentation for the testing images. Since fully-connected CRF [16] brings negligible improvements when the network predictions are accurate enough, we do not apply CRF as post refinement in our experiments.

### 4.2 Ablation study

To provide more insight into the proposed architecture, we conduct a series of experiments on PASCAL VOC using different experimental settings concerning different network architecture and training data. We use VGG16 as backbone unless stated otherwise.

**Using only 1.4k strong data** To obtain decent results with only 1.4k strong annotations, it is important to perform the same number of iterations (instead of epochs) to let the network converge. When trained enough amount

of iterations, we could achieve a mIoU of 68.9%, which is much better than the 62.5% reported in FickleNet [19].

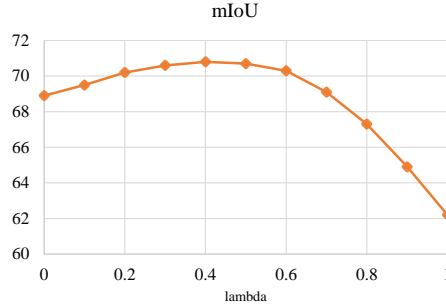   **Two separate networks** As aforementioned, the single-branch network



**Fig. 6.** The segmentation mIoU (%) with respect to different $\lambda$'s.

trained under the mixture of strong and weak annotations achieves no better performance than using only the strong ones. Therefore, it is natural to train two different networks on two sets of data since there exists an obvious annotation inconsistency. Specifically, we train two networks, the first supervised by the strong annotations and the second by extra weak annotations. Then their outputs are aggregated through the following equation:

$$F(x) = \lambda * F_w(x) + (1 - \lambda) * F_s(x) \tag{4}$$

where $F_w$ and $F_s$ denote the weak and strong network respectively. Fig.6 shows the segmentation accuracy under different $\lambda$ values. Simply training on the strong annotations yields an accuracy of 68.9%, 6.1% higher than the weak one. The result could be improved up to 70.8% with $\lambda$ equal to 0.4, a 1.9% boost over the strong network. However, separate networks double the computation overhead during both training and inference.

   **Single branch vs. dual branch** Our VGG16-based implementation of the DSRG method achieves mIoU of 57.0% and 60.1% after retraining, presented in Table 1. With a combination of 1.4k strong annotations and 9k weak annotations estimated from DSRG, the single-branch network only improves the segmentation accuracy by 2.7%. However, it already achieves much higher accuracy of 68.9% under only 1.4k strong annotations, which means the extra 9k weak annotations bring no benefits but actually downgrade the performance dramatically, nearly 6% drop. This phenomenon verifies our hypothesis that equal treatment of strong and weak annotations are problematic as large amount inaccurate weak annotations mislead the network training.

   We then train the proposed architecture with the strong branch supervised by the 1.4k strong annotations and the weak one by extra 9k weak annotations. This time the accuracy successfully goes up to 72.2%, a 3.3% improvement over

**Table 1.** Ablation experiments concerning network architectures and training data. Rows marked with "*" are results from the proposed dual-branch network and others are from the single-branch network.

| Backbone | Strong branch | Weak branch | mIoU |
|---|---|---|---|
| VGG16 | 10k weak | - | 57.0 |
| VGG16 | 10k weak (retrain) | - | 60.1 |
| VGG16 | 1.4k strong + 9k weak | - | 62.8 |
| VGG16 | 1.4k strong | - | 68.9 |
| VGG16 | 10k strong | - | 71.4 |
| VGG16*(w/o oversampling) | 1.4k strong | 1.4k strong + 9k weak | 66.4 |
| VGG16*(w oversampling) | 1.4k strong | 1.4k strong + 9k weak | 72.2 |
| VGG16*(w oversampling) | 1.4k strong | 10k strong | 73.9 |

the 1.4k single-branch model. Remarkably, this result is even better than training a single-branch model with 10k strong annotations, which implies that there is an inconsistency between the official 1.4k annotations and the additional 9k annotations provided by [10]. Based on this observation, we conduct another experiment on our dual-branch network with 1.4k strong annotations for the strong branch and 10k strong annotations for the weak branch. As expected, the accuracy is further increased by 1.7%. To see the impact of sample imbalance, we also conduct a experiment with our network without oversampling. And it achieves accuracy of 66.4%, up from mixed training result 62.8% but still worse than that of using only strong annotations. So we conclude that it's more effective to combine dual-branch model with oversampling training strategy.

### 4.3   Comparison with the state-of-arts

Table 2 compares the proposed method with current state-of-art weakly-and-semi supervised methods: SEC [15], DSRG [13], FickleNet [19], WSSL [26], Box-Sup [7], etc. For fair comparison, the result reported in the original paper is listed along with the backbone adopted.

The weakly-supervised methods are provided in the upper part of Table 2 as reference since many of them used relatively weak supervision, with FickleNet (61.9%) achieving the best performance among other baselines using only class labels. However, the elimination of the demand for pixel-level annotations results in significant performance drop, around 11% compared to their fully-supervised counterparts. There are some recent works exploring other weak supervisions, such as Normalized cut loss [32] and Box-driven method [31]. They improved the segmentation performance significantly with slightly increasing labeling efforts. Our method actually serves as an alternative direction by using a combination of strong and weak annotations to achieve excellent results.

The lower part of Table 2 presents results of the semi-supervised methods. DSRG and FickleNet used the same region growing mechanism to expand the original object seeds. As shown in the table, all previous methods achieved

**Table 2.** Segmentation results of different methods on PASCAL VOC 2012 *val* and *test* set. * - Result copied from the FickleNet paper

| Methods | Backbone | Val | Test |
|---|---|---|---|
| Supervision: 10k scribbles | | | |
| Scribblesup [21] | VGG16 | 63.1 | - |
| Normalized cut [32] | ResNet101 | 74.5 | - |
| Supervision: 10k boxes | | | |
| WSSL [26] | VGG16 | 60.6 | 62.2 |
| BoxSup [7] | VGG16 | 62.0 | 64.2 |
| Supervision: 10k class | | | |
| SEC [15] | VGG16 | 50.7 | 51.7 |
| AF-SS [36] | VGG16 | 52.6 | 52.7 |
| Multi-Cues [29] | VGG16 | 52.8 | 53.7 |
| DCSP [4] | VGG16 | 58.6 | 59.2 |
| DSRG [13] | VGG16 | 59.0 | 60.4 |
| AffinityNet [1] | VGG16 | 58.4 | 60.5 |
| MDC [35] | VGG16 | 60.4 | 60.8 |
| FickleNet [19] | VGG16 | 61.2 | 61.9 |
| Supervision: 1.4k pixel + 9k class | | | |
| DSRG [13]* | VGG16 | 64.3 | - |
| FickleNet [19] | VGG16 | 65.8 | - |
| WSSL [26] | VGG16 | 64.6 | 66.2 |
| MDC [35] | VGG16 | 65.7 | 67.6 |
| Ours | VGG16 | 72.2 | 72.3 |
| Ours | ResNet101 | **76.6** | **77.1** |

roughly the same and poor performance when learned under 1.4k pixel annotations and 9k class annotations, with the best accuracy 67.6% by MDC approach.

Our method significantly outperforms all the weakly-and-semi supervised method by a large margin, with state-of-art 77.1% mIoU on the *test* set when ResNet101 backbone is adopted.

### 4.4   Visualization result

Fig.7 shows segmentation results of sample images from PASCAL VOC *val* set. As can be seen in the third column, weakly-supervised method (DSRG) generates segmentation maps of relatively poor quality and no improvement is visually significant if combined with 1.4k strong annotations. Our approach manages to remove some of the false positives in the foreground categories, as in the second and third examples. The last line demonstrates a failure case when neither approach is effective to generate correct prediction.

### 4.5   Results on COCO

To verify the generality of the proposed architecture, we conduct further experiments on the Microsoft COCO dataset, which contains a lot more images

**Table 3.** Per-class IoU on COCO val set. (a) Single-branch network using 20k strong annotations; (b) Single-branch network using 98k extra weak annotations; (c) Dual-branch network using 98k extra weak annotations.

| Cat. | Class | (a) | (b) | (c) |
|---|---|---|---|---|
| BG | background | 86.2 | 78.4 | 86.7 |
| P | person | 74.4 | 60.7 | 75.2 |
| Vehicle | bicycle | 54.2 | 48.4 | 55.3 |
| | car | 47.4 | 38.2 | 49.5 |
| | motorcycle | 70.4 | 63.7 | 70.6 |
| | airplane | 63.3 | 30.5 | 66.0 |
| | bus | 69.7 | 64.1 | 71.5 |
| | train | 67.2 | 46.7 | 69.8 |
| | truck | 43.3 | 36.4 | 45.2 |
| | boat | 42.5 | 26.1 | 41.9 |
| Outdoor | traffic light | 42.9 | 27.6 | 47.1 |
| | fire hydrant | 74.2 | 47.3 | 75.5 |
| | stop sign | 82.3 | 53.6 | 87.3 |
| | parking meter | 48.4 | 42.7 | 53.8 |
| | bench | 32.6 | 25.3 | 34.9 |
| Animal | bird | 56.6 | 33.9 | 62.0 |
| | cat | 76.7 | 65.1 | 77.5 |
| | dog | 68.7 | 60.6 | 69.0 |
| | horse | 64.4 | 50.0 | 66.2 |
| | sheep | 70.5 | 55.5 | 73.3 |
| | cow | 61.7 | 49.7 | 65.3 |
| | elephant | 79.9 | 67.6 | 81.2 |
| | bear | 79.7 | 60.4 | 81.7 |
| | zebra | 81.7 | 61.2 | 82.9 |
| | giraffe | 74.3 | 47.0 | 75.0 |
| Accessory | backpack | 11.4 | 2.5 | 12.6 |
| | umbrella | 57.9 | 44.3 | 59.1 |
| | handbag | 6.8 | 0.0 | 8.2 |
| | tie | 34.5 | 20.6 | 35.4 |
| | suitcase | 53.1 | 48.4 | 57.6 |
| Sport | frisbee | 48.2 | 39.4 | 50.8 |
| | skis | 14.6 | 5.3 | 11.8 |
| | snowboard | 37.8 | 15.6 | 39.1 |
| | sports ball | 27.0 | 13.3 | 29.7 |
| | kite | 32.1 | 23.7 | 36.2 |
| | baseball bat | 10.4 | 0.0 | 11.1 |
| | baseball glove | 28.4 | 0.0 | 37.6 |
| | skateboard | 32.0 | 20.4 | 31.6 |
| | surfboard | 43.7 | 32.2 | 44.5 |
| | tennis racket | 55.7 | 47.3 | 58.1 |
| | bottle | 39.7 | 33.0 | 39.4 |

| Cat. | Class | (a) | (b) | (c) |
|---|---|---|---|---|
| Kitchenware | wine glass | 42.5 | 36.0 | 45.2 |
| | cup | 38.8 | 30.9 | 38.9 |
| | fork | 16.6 | 0.0 | 17.2 |
| | knife | 3.4 | 0.1 | 6.9 |
| | spoon | 5.9 | 0.0 | 5.4 |
| | bowl | 33.0 | 22.4 | 34.7 |
| Food | banana | 62.4 | 53.1 | 63.3 |
| | apple | 36.6 | 29.8 | 37.3 |
| | sandwich | 44.3 | 35.1 | 46.0 |
| | orange | 55.3 | 50.3 | 57.9 |
| | broccoli | 49.9 | 37.3 | 53.3 |
| | carrot | 34.4 | 31.8 | 37.0 |
| | hot dog | 38.8 | 36.0 | 39.8 |
| | pizza | 74.8 | 68.6 | 76.6 |
| | donut | 49.4 | 48.6 | 53.9 |
| | cake | 45.6 | 40.6 | 45.3 |
| Furniture | chair | 24.4 | 12.3 | 25.2 |
| | couch | 41.0 | 20.5 | 42.6 |
| | potted plant | 23.4 | 15.5 | 24.5 |
| | bed | 46.9 | 38.2 | 50.4 |
| | dining table | 34.8 | 9.2 | 35.0 |
| | toilet | 61.5 | 45.3 | 62.7 |
| Electronics | tv | 49.9 | 22.5 | 52.5 |
| | laptop | 56.2 | 40.6 | 57.4 |
| | mouse | 38.5 | 0.7 | 35.5 |
| | remote | 37.8 | 25.7 | 30.9 |
| | keyboard | 44.3 | 35.9 | 47.1 |
| | cell phone | 44.1 | 36.8 | 42.3 |
| Appliance | microwave | 47.2 | 32.8 | 44.6 |
| | oven | 42.9 | 29.7 | 47.9 |
| | toaster | 0.0 | 0.0 | 0.0 |
| | sink | 40.0 | 30.5 | 42.4 |
| | refrigerator | 55.5 | 34.8 | 57.3 |
| Indoor | book | 29.9 | 16.4 | 29.0 |
| | clock | 57.5 | 16.4 | 59.5 |
| | vase | 45.8 | 30.1 | 43.9 |
| | scissors | 57.1 | 34.1 | 56.4 |
| | teddy bear | 64.4 | 57.9 | 66.1 |
| | hair drier | 0.0 | 0.0 | 0.0 |
| | toothbrush | 13.2 | 8.5 | 17.6 |
| | **mean IoU** | **46.1** | **33.4** | **47.6** |

(118k) and semantic categories (81 classes), thus posing a challenge even for fully-supervised segmentation approaches. We randomly select 20k images as our strong set and the remaining 98k images as the weak set, whose annotations are estimated from the DSRG method. This splitting ratio is roughly the same compared to PASCAL VOC experiments. We report per-class IoU over all 81 semantic categories on the 5k validation images. As shown in Table 3, with 20k strong annotations, the single branch network achieves an accuracy of 46.1%. When we bring in extra 98k weak annotations estimated by DSRG, the
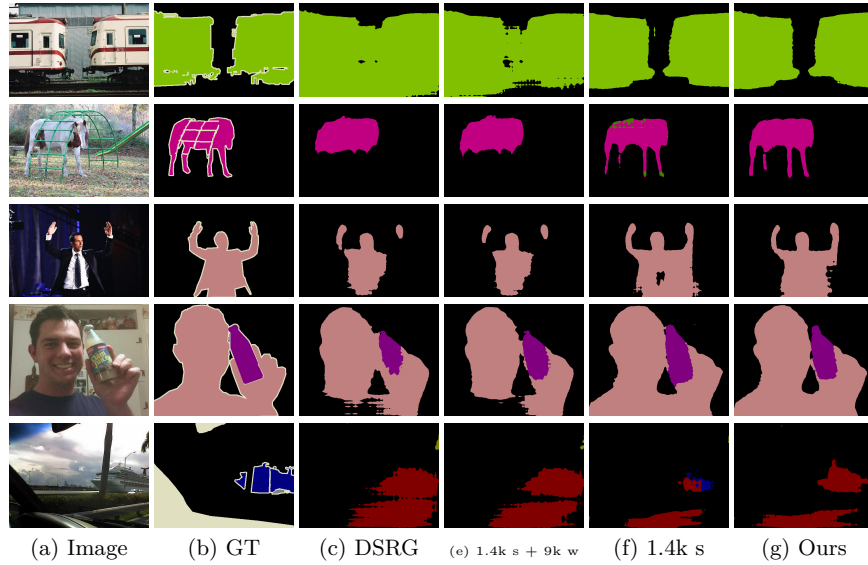
|   (a) Image   |   (b) GT   |   (c) DSRG   |   (e) 1.4k s + 9k w   |   (f) 1.4k s   |   (g) Ours   |

**Fig. 7.** Demonstration of sample images. (a) Original images; (b) Ground truth; (c) DSRG; (e) Mixing 1.4k s + 9k w for training; (f) 1.4k strong annotations; (g) Ours under 1.4k s + 9k w.

performance downgrades by 12.7%, down to only 33.4%, which again verifies our hypothesis. Using our dual-branch network, the performance successfully goes up to 47.6%, which means our approach manages to make use of the weak annotations.

## 5   Conclusion

We have addressed the problem of semi-supervised semantic segmentation where a combination of finely-labeled masks and coarsely-estimated data are available for training. Weak annotations are cheap to obtain yet not enough to train a segmentation model of high quality. We propose a strong-weak dual-branch network that has fully utilized the limited strong annotations without being overwhelmed by the bulk of weak ones. It manages to eliminate the learning obstacles of sample imbalance and supervision inconsistency. Our method significantly outperforms the weakly-supervised and almost reaches the accuracy of fully-supervised models. We think semi-supervised approaches could serve as an alternative to weakly-supervised methods by retaining the segmentation accuracy while still keeping labeling budget in control.

# References

1. Ahn, J., Kwak, S.: Learning pixel-level semantic affinity with image-level supervision for weakly supervised semantic segmentation. In: CVPR (June 2018)
2. Bearman, A., Russakovsky, O., Ferrari, V., Li, F.: What's the Point: Semantic Segmentation with Point Supervision. ECCV (2016)
3. Bennett, K., Demiriz, A.: Semi-supervised support vector machines. In: NIPs. pp. 368–374. MIT Press, Cambridge, MA, USA (1999), http://dl.acm.org/citation.cfm?id=340534.340671
4. Chaudhry, A., Dokania, P.K., Torr, P., Toor, P.: Discovering class-specific pixels for weakly-supervised semantic segmentation. In: BMVC. vol. abs/1707.05821 (2017)
5. Chen, L., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. TPAMI **40**, 834–848 (2016)
6. Chen, L., Zhu, Y., Papandreou, G., Schroff, F., Adam, H.: Encoder-decoder with atrous separable convolution for semantic image segmentation. In: ECCV (2018)
7. Dai, J., He, K., Sun, J.: Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. ICCV pp. 1635–1643 (2015)
8. Everingham, M., Van Gool, L., Williams, C.K.I., Winn, J., Zisserman, A.: The PASCAL Visual Object Classes Challenge 2012 (VOC2012) Results. http://www.pascal-network.org/challenges/VOC/voc2012/workshop/index.html (2012)
9. Girshick, R.: Fast r-cnn. ICCV pp. 1440–1448 (2015)
10. Hariharan, B., Arbelaez, P., Bourdev, L., Maji, S., Malik, J.: Semantic contours from inverse detectors. In: ICCV (2011)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. CVPR pp. 770–778 (2015)
12. Hou, Q., Jiang, P., Wei, Y., Cheng, M.: Self-erasing network for integral object attention. In: NIPs (2018)
13. Huang, Z., Wang, X., Wang, J., Liu, W., Wang, J.: Weakly-supervised semantic segmentation network with deep seeded region growing. In: CVPR (June 2018)
14. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. CoRR **abs/1412.6980** (2014)
15. Kolesnikov, A., Lampert, C.: Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In: ECCV. vol. abs/1603.06098 (2016)
16. Krähenbühl, P., Koltun, V.: Efficient inference in fully connected crfs with gaussian edge potentials. In: NIPS (2011)
17. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: NIPs. pp. 1097–1105. Curran Associates Inc., USA (2012), http://dl.acm.org/citation.cfm?id=2999134.2999257
18. Laine, S., Aila, T.: Temporal ensembling for semi-supervised learning. In: ICLR. vol. abs/1610.02242 (2016)
19. Lee, J., Kim, E., Lee, S., Lee, J., Yoon, S.: Ficklenet: Weakly and semi-supervised semantic image segmentation using stochastic inference. In: CVPR (June 2019)
20. Li, K., Wu, Z., Peng, K., Ernst, J., Fu, Y.: Tell me where to look: Guided attention inference network. CVPR pp. 9215–9223 (2018)
21. Lin, D., Dai, J., Jia, J., He, K., Sun, J.: Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. CVPR pp. 3159–3167 (2016)
22. Lin, T., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. ICCV pp. 2999–3007 (2017)

23. Lin, T., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.: Microsoft coco: Common objects in context. In: ECCV (2014)
24. Miyato, T., Maeda, S., Koyama, M., Ishii, S.: Virtual adversarial training: A regularization method for supervised and semi-supervised learning. TPAMI **41**(8), 1979–1993 (Aug 2019). https://doi.org/10.1109/TPAMI.2018.2858821
25. Oh, S., Benenson, R., Khoreva, A., Akata, Z., Fritz, M., Schiele, B.: Exploiting saliency for object segmentation from image level labels. In: CVPR (2017), to appear
26. Papandreou, G., Chen, L., Murphy, K.P., Yuille, A.L.: Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In: ICCV. pp. 1742–1750 (Dec 2015). https://doi.org/10.1109/ICCV.2015.203
27. Papandreou, G., Chen, L., Murphy, K.P., Yuille, A.L.: Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In: ICCV. pp. 1742–1750. ICCV '15, IEEE Computer Society, Washington, DC, USA (2015). https://doi.org/10.1109/ICCV.2015.203, http://dx.doi.org/10.1109/ICCV.2015.203
28. Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks. TPAMI **39**, 1137–1149 (2015)
29. Roy, A., Todorovic, S.: Combining bottom-up, top-down, and smoothness cues for weakly supervised image segmentation. In: CVPR (July 2017)
30. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2015)
31. Song, C., Huang, Y., Ouyang, W., Wang, L.: Box-driven class-wise region masking and filling rate guided loss for weakly supervised semantic segmentation. In: CVPR (June 2019)
32. Tang, M., Djelouah, A., Perazzi, F., Boykov, Y., Schroers, C.: Normalized cut loss for weakly-supervised cnn segmentation. In: CVPR (June 2018)
33. Tarvainen, A., Valpola, H.: Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In: ICLR (2017)
34. Wang, X., You, S., Li, X., Ma, H.: Weakly-supervised semantic segmentation by iteratively mining common object features. In: CVPR (June 2018)
35. Wei, Y., Xiao, H., Shi, H., Jie, Z., Feng, J., Huang, T.: Revisiting dilated convolution: A simple approach for weakly- and semi-supervised semantic segmentation. In: CVPR (June 2018)
36. Wei, Y., Feng, J., Liang, X., Cheng, M.M., Zhao, Y., Yan, S.: Object region mining with adversarial erasing: A simple classification to semantic segmentation approach. In: CVPR (July 2017)
37. Zhang, J., Bargal, S.A., Lin, Z., Brandt, J., Shen, X., Sclaroff, S.: Top-down neural attention by excitation backprop. IJCV **126**, 1084–1102 (2016)
38. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. CVPR pp. 6230–6239 (2016)
39. Zhou, B., Khosla, A., A., L., Oliva, A., Torralba, A.: Learning Deep Features for Discriminative Localization. CVPR (2016)
40. Zhu, X.: Semi-supervised learning literature survey. Tech. Rep. 1530, Computer Sciences, University of Wisconsin-Madison (2005)