# On Regularized Losses
# for Weakly-supervised CNN Segmentation

Meng Tang[†]     Federico Perazzi[*]     Abdelaziz Djelouah[*]
Ismail Ben Ayed[‡]     Christopher Schroers[*]     Yuri Boykov[†]

[†]Computer Science, University of Waterloo, Canada
[*]Disney Research, Zürich, Switzerland
[‡]École de Technologie Supérieure, University of Quebec, Canada

**Abstract.** Minimization of regularized losses is a principled approach to weak supervision well-established in deep learning, in general. However, it is largely overlooked in semantic segmentation currently dominated by methods mimicking full supervision via "fake" fully-labeled training masks (proposals) generated from available partial input. To obtain such full masks the typical methods explicitly use standard regularization techniques for "shallow" segmentation, e.g. graph cuts or dense CRFs. In contrast, we integrate such standard regularizers directly into the loss functions over partial input. This approach simplifies weakly-supervised training by avoiding extra MRF/CRF inference steps or layers explicitly generating full masks, while improving both the quality and efficiency of training. This paper proposes and experimentally compares different losses integrating MRF/CRF regularization terms. We juxtapose our regularized losses with earlier proposal-generation methods using explicit regularization steps or layers. Our approach achieves state-of-the-art accuracy in semantic segmentation with near full-supervision quality.

## 1   Introduction

We advocate *regularized loss* functions for weakly-supervised training of semantic CNN segmentation. The use of unsupervised loss terms acting as regularizers on the output of deep-learning architectures is a principled approach to exploit structure similarity of partially labeled data [1,2]. Surprisingly, this general idea was largely overlooked in weakly-supervised CNN segmentation where current methods often introduce computationally expensive MRF/CRF layers or post-processing inference steps generating "fake" full masks from partial input.

We propose to use (relaxations of) MRF/CRF terms directly inside the loss avoiding explicit guessing of full training masks. This approach follows well-established ideas for weak supervision in deep learning [1,2] and continues our recent work [3] that proposed the integration of standard objectives in shallow[1] segmentation directly into loss functions. While [3] is entirely focused on the *normalized cut loss* motivated by a popular balanced segmentation criterion [4],

---

[1] In this paper "shallow" refers to standard segmentation methods unrelated to CNNs.

we now study a different class of *regularized losses* including (relaxations of) standard MRF/CRF potentials. While they are common as shallow regularizers [5,6,7,8] or as trainable layers [9], they were never used directly as losses.

We propose and evaluate several new losses motivated by MRF/CRF potentials and their combination with balanced partitioning criteria [10]. Such losses can be adapted to many forms of weak (or semi-) supervision based on diverse existing MRF/CRF formulations for interactive graph cut segmentation. But, the scope of this paper is limited to training with partial (user scribble) masks where regularized losses combined with cross entropy over the partial masks achieve the state-of-the-art close to full-supervision quality.

Besides basic Potts model [5], we use popular fully connected pairwise CRF potentials of Krähenbühl and Koltun [8], often referred to as *dense CRF*. In conjunction with CNNs dense CRFs have become the de-facto choice for semantic segmentation in the contexts of fully [11,12,9] and weakly/semi [13,14,15] supervised learning. For instance, DeepLab [11] popularized dense CRF as a post-processing step. In fully supervised setting, integrating the unary scores of a CNN classifier and the pairwise potentials of dense CRF achieve competitive performances [12]. This is facilitated by fast mean-field inference techniques for dense CRF based on high-dimensional filtering [16].

Weakly supervised semantic segmentation is commonly addressed by mimicking full supervision via synthesizing fully-labeled training masks (proposals) from the available partial inputs [15,14,17]. These schemes typically iterate two steps: CNN training and proposal generation via regularization-based shallow interactive segmentation, e.g. graph cut [17] or dense CRF mean-field inference [15,14]. In contrast, our approach avoids explicit inference steps by integrating shallow regularizers directly into the loss functions. Section 3 makes some interesting connections between proposal-generation and our regularized losses.

For simplicity, this paper uses a very basic quadratic relaxation of discrete MRF/CRF potentials, even though there are many alternatives, e.g. TV-based [18] and convex formulations [19,20], $L_p$ relaxations [21], LP and other relaxations [22,23]. Evaluation of different relaxations in the context of regularized weak supervision losses is left for future work. Our main contributions are:

- We propose and evaluate several *regularized losses* for weakly supervised CNN segmentation based on Potts [5], dense CRF [8], and kernel cut [10] regularizers (Sec.2). Our approach avoids explicit inference steps as in proposal-based methods. This continues the study of losses motivated by standard shallow segmentation energies started in [3] with *normalized cut loss*.
- We show that iterative proposal-generation schemes for weak supervision, which alternate CNN learning and mean-field inference, can be viewed as an approximate alternating direction optimization of regularized losses (Sec.3).
- Comprehensive experiments (Sec.4) with our regularized weakly supervised losses show (1) state-of-the-art performance for weakly supervised CNN segmentation reaching near full-supervision accuracy and (2) better quality and efficiency than proposal generating methods or normalized cut loss [3]. Alternating schemes (proposal generation) give higher loss at convergence.

## 2 Our Regularized Semi-supervised Losses

This section introduces our regularized losses for weakly-supervised segmentation. In general, the use of regularized losses is a well-established approach in semi-supervised deep learning [1,2]. We advocate this principle for semantic CNN segmentation, propose specific shallow regularizers for such losses, and discuss their properties.

Assuming image $I$ and its *partial* ground truth labeling or mask $Y$, let $f_\theta(I)$ be the output of a segmentation network parameterized by $\theta$. In general, CNN training with our joint regularized loss corresponds to optimization problem of the following form

$$\min_\theta \ell(f_\theta(I), Y) \; + \; \lambda \cdot R(f_\theta(I)) \tag{1}$$

where $\ell(S, Y)$ is a ground truth loss and $R(S)$ is a regularization term or regularization loss. Both losses have argument $S = f_\theta(I) \in [0,1]^{|\Omega| \times K}$, which is $K$-way softmax segmentation generated by a network. Using cross entropy over partial labeling as the ground truth loss, we have the following joint *regularized semi-supervised loss*

$$\sum_{p \in \Omega_\mathcal{L}} H(Y_p, S_p) \; + \; \lambda \cdot R(S) \tag{2}$$

where $\Omega_\mathcal{L} \subset \Omega$ is the set of labeled pixels and $H(Y_p, S_p) = -\sum_k -Y_p^k \log S_p^k$ is the cross entropy between network predicted segmentation $S_p \in [0,1]^K$ (a row of matrix $S$ corresponding to point $p$) and ground truth labeling $Y_p \in \{0,1\}^K$.

In principle, any differentiable function $R(S)$ can be used as a loss. This paper studies (relaxations of) regularizers from shallow segmentation as loss functions. Section 2.1 details our MRF/CRF loss and its implementation. In Section 2.2, we propose *kernel cut loss* combining CRF with normalized cut terms and justify this combination.

### 2.1 Potts/CRF Losses

Assuming that segmentation variables $S_p$ are restricted to binary class indicators $S_p \in \{0,1\}^K$, the standard Potts model [5] could be represented via Iverson brackets $[\cdot]$, as on the left hand side below

$$\sum_{p,q \in \Omega} W_{pq} \, [S_p \neq S_q] \; = \; \sum_{p,q \in \Omega} W_{pq} \, \|S_p - S_q\|^2, \tag{3}$$

where $W = [W_{pq}]$ is a matrix of pairwise discontinuity costs or an *affinity matrix*. The right hand side above is a particularly straightforward quadratic relaxation of the Potts model that works for relaxed $S_p \in [0,1]^K$ corresponding to a typical soft-max output of CNNs. In fact, this quadratic function is very common in the general context of regularized weakly supervised losses in deep learning [1].

As discussed in the introduction, this relaxation is not unique [18,19,20,21,22]. We use slightly different quadratic relaxation of the Potts model

$$R_{CRF}(S) = \sum_k S^{k'} W (\mathbf{1} - S^k) \tag{4}$$

expressed in terms of support vectors for each label $k$, i.e. columns of the segmentation matrix $S^k \in [0,1]^{|\Omega|}$. For discrete segment indicators (4) gives the cost of a cut between segments, same as the Potts model on the left hand side of (3), but it differs from the relaxation of the right hand side of (3).

The affinity matrix $W$ can be sparse or dense. Sparse $W$ commonly appears in the context of boundary regularization and edge alignment in shallow segmentation [6]. With dense Gaussian kernel $W_{pq}$ (4) is a relaxation of DenseCRF [24]. The implementation details including fast computation of the gradient (11) for CRF loss with dense Gaussian kernel is described in Sec. 4.

## 2.2   Kernel Cut Loss

Besides the CRF loss (4), we also propose its combination with normalized cut loss [3] where each term is a ratio of a segment's cut cost (Potts model) over the segment's weighted size (normalization)

$$R_{NC}(S) = \sum_k \frac{S^{k'} \hat{W}(1 - S^k)}{d' S^k}, \qquad (5)$$

where $d = \hat{W}\mathbf{1}$ are node degrees. Note that the affinity matrix $\hat{W}$ for normalized cut can be different from $W$ in CRF (4). The combined *kernel cut loss* is simply a linear combination of (4) and (5)

$$R_{KC}(S) = \sum_k S^{k'} W(\mathbf{1} - S^k) \;\; + \;\; \gamma \sum_k \frac{S^{k'} \hat{W}(1 - S^k)}{d' S^k} \qquad (6)$$

which is motivated by *kernel cut* shallow segmentation [10] with complementary benefits of balanced normalized cut partitioning and object boundary regularization or edge alignment as in Potts model. While the kernel cut loss is a high-order objective, its gradient (12) can be efficiently implemented, see Sec. 4.

This paper compares experimentally CRF, normalized cut and kernel cut losses for weakly supervised segmentation. In our experiments, the best weakly supervised segmentation is achieved with kernel cut loss.

Note that standard normalized cut and CRF objectives in shallow segmentation require fairly different optimization techniques (e.g. spectral relaxation or graph cuts), but the standard gradient descent approach for optimizing losses during CNN training allows significant flexibility in including different regularization terms, as long as there is a reasonable relaxation.

## 3   Connecting proposals generation and loss optimization

The main stream of weakly-supervised methods generate segmentation proposals and train with such "fake" ground truth [17,25,26,13,14,27]. In fact, many off-line shallow interactive segmentation techniques can be used to propagate labels and generate masks, e.g. graph cuts [6,7], random walker [28,21], etc. However,

training is vulnerable to mistakes in the proposals. While alternating proposal generation and network training [17] may improve the quality of the proposals, errors reinforce themselves in such self-taught learning scheme [29]. Rather than training networks to fit potential errors, our regularized semi-supervised loss framework is more direct and principled [29,1].

In this section, we show that proposal methods can be viewed as *approximate* alternating direction method (ADM) for optimization [30], which does not account directly for network variables $\theta$ in the ADM splitting. This optimization insight suggests that expressing very popular regularization terms, for instance, dense CRF, explicitly in terms of the network variables and performing direct back-propagation could be a better optimization alternative to the existing proposal generation methods, in both the quality of the obtained solutions and efficiency. Our optimization results confirm this, e.g. see the CRF loss plot in Fig. 3 and the training times in Table 3.

We consider proposal-generation schemes iterating between two steps, **network training** and **proposal generation**. Then alternation can happen either when training converges or online for each batch. At each iteration, the first step learns the network parameters $\theta$ from a given (fixed) ground-truth proposal $\tilde{X}$ computed at the previous iteration. This amounts to updating the K-way softmax segmentation $S$ to $\tilde{S} \equiv f_{\tilde{\theta}}(I)$ by minimizing the following proposal-based cross entropy with respect to parameters $\theta$ via standard back-propagation:

$$\tilde{\theta} = \arg \min_{\theta} \sum_{p \in \Omega_{\mathcal{L}}} H(Y_p, S_p) + \sum_{p \in \Omega_{\mathcal{U}}} H(\tilde{X}_p, S_p) \quad \text{for} \quad S \equiv f_{\theta}(I) \qquad (7)$$

where $\tilde{X}_p \in [0,1]^K$ are the ground-truth proposals for unlabeled pixels $p \in \Omega_{\mathcal{U}}$. Mask $\tilde{X}_p$ is constrained to be equal to $Y_p$ for labeled pixels $p \in \Omega_{\mathcal{L}}$. The second step fixes the network output $\tilde{S}$ and finds the next ground-truth proposal by minimizing regularization functionals that are standard in shallow segmentation:

$$\min_{X \in [0,1]^{|\Omega| \times K}} \sum_{p \in \Omega_{\mathcal{U}}} H(X_p, \tilde{S}_p) + \lambda R(X) \qquad (8)$$

where $X_p \in [0,1]^k$ denotes latent pixel labels within the probability simplex. Note that for fixed $\tilde{S}$ the cross entropy terms $H(X_p, \tilde{S}_p)$ in (8) are unary potentials for $X$. When $R$ corresponds to dense CRF, optimization of (8) is facilitated by fast mean-field inference techniques [8,31] significantly reducing the computational times via parallel updates of variables $X_p$ and high-dimensional filtering [16]. Appendix A shows that mean-field algorithms can be equivalently interpreted as a *convex-concave* approach to optimizing the following objective

$$\min_{X \in [0,1]^{|\Omega| \times K}} \sum_{p \in \Omega_{\mathcal{U}}} H(X_p, \tilde{S}_p) + \lambda R(X) - \sum_{p \in \Omega_{\mathcal{U}}} H(X_p) \qquad (9)$$

combining (8) and negative entropies $H(X_p) = -\sum_k X_p^k \log X_p^k$ that act as a simplex *barrier* for variables $X_p$. This yields closed-form independent (parallel) updates of variables $X_p$, while ensuring convergence under some conditions[2].

---

[2] Parallel updates are guaranteed to converge for concave CRF models, e.g. Potts [24].

**Proposition 1.** *Proposal methods alternating steps* (9) *and* (7) *can be viewed as approximate* alternating direction method *(ADM)*[3] *[30] for optimizing our regularized loss* (2) *using the following decomposition of the problem:*

$$\min_{\theta,X\in[0,1]^{|\Omega|\times K}} \sum_{p\in\Omega_{\mathcal{L}}} H(Y_p,S_p) + \lambda R(X) + \sum_{p\in\Omega_{\mathcal{U}}} KL(X_p|S_p) \qquad (10)$$

*where KL denotes the Kullback-Leibler divergence.*

*Proof.* The link between (10) and (9) comes directly from the following relation between the KL divergence and the entropies: $KL(X_p|S_p) = H(X_p,S_p) - H(X_p)$.

Instead of optimizing directly regularized loss (2) with respect to network parameters, proposal methods splits the optimization problem into two easier subproblems in (10). This is done by replacing the network softmax outputs $S_p$ in the regularization by latent distributions $X_p$ (the proposals) and minimizing a divergence between $S_p$ and $X_p$, which is KL in this case. This is conceptually similar to the general principles of ADM [30], except that the splitting is not done directly with respect the variables of the problem (i.e., parameters $\theta$) but rather with respect to network outputs $S$. This can be viewed as an *approximate* ADM scheme, which does not account directly for variables $\theta$ in the ADM splitting. Note that the method in [13] generates proposals via dense CRF layer, but their approach slightly deviates from the described ADM scheme since they also back-propagate through this layer[4]. But, as we show in Table 3, such backpropagation does not help and can be dropped. Moreover, our direct optimization of regularized losses makes such proposal generating layers (or procedures) entirely redundant. Our approach gives simpler and more efficient training avoiding expensive iterative inference [13] and obtaining better performance.

## 4   Experiments

Sec. 4.1 is the main experimental result of this paper. For weakly-supervised segmentation with scribbles [17], we train with different regularized losses. The experiments cover our proposed CRF loss, high-order normalized cut loss in [3] and kernel cut loss, as discussed in Sec. 2. We show that combining CRF (4) with normalized cut (5) *a la* KernelCut [10] yields the best performance.

In Sec. 4.2, using direct loss and using generated proposals for training are compared. In the light of the technical connection of the two schemes from optimization perspective in Sec. 3, we also evaluate how "regularized" are the segmentations obtained by computing the regularization energy. Besides for scribbles, we also utilize our regularized loss framework for image-level labels based

---

[3] In its basic form, *alternating direction method* transforms problem $\min_x f(x) + g(x)$ into $\min_{x,y} f(x) + g(y)$ s.t $x = y$ and alternates optimization over $x$ and $y$. This may work if optimizing $f$ and $g$ seperately is easier than the original problem.

[4] Cross-entropy loss $H(X(S),S)$ in [13] uses CRF layer proposal $X(S)$ generated from network output $S$. Dependence of $X$ on $S$ motivates back-propagation for this layer.

supervision and compare to SEC [13], a recent method based on proposal generation. Our method achieved the state-of-the-art for weakly supervised segmentation with scribbles or image-level labels.

We also investigate if regularization loss will facilitate fully or semi-supervised segmentation with unlabeled images. Some preliminary results are given in Sec. 4.3 for these extensions.

**Dataset** Most experiments are on the PASCAL VOC12 segmentation dataset. For all method, we train with the augmented dataset of 10,582 images. The scribble annotations for these training images are from [17]. Following standard protocol, mean intersection-over-union (mIoU) over the 21 classes is evaluated on the *val* set that contains 1,449 images. For image-level label supervision, our experiment setup and dataset is the same as that used in [13].

**Implementation details** Our implementation is based on DeepLab v2 [11]. We follow the learning rate strategy in DeepLab v2 [5] for the baseline with full supervision. For our method with regularized loss, we first train with partial cross entropy loss only for the seeds. Then we fine-tune with extra regularized losses of different types for the same number of iterations. Our CRF and normalized cut regularization losses are defined at full image resolution. If the network outputs shrinked labeling, which is typical, the labeling is interpolated to original resolution before feeding into the loss layer.

We choose dense Gaussian kernel over RGBXY channels for $R_{CRF}(S)$, $R_{NC}(S)$ and $R_{KC}(S)$. As hyper-parameter, the Gaussian bandwidth is optimized via validation for DenseCRF, normalized cut and kernel cut. As is also mentioned in [3], naive forward and backward pass of such fully-connected pairwise or high-order loss layer would be prohibitively slow ($O(|\Omega|^2)$ for $|\Omega|$ pixels). For example, to implement $R_{CRF}(S)$ (4) as a loss, we need to compute its gradient w.r.t. $S^k$ during backpropagation,

$$\frac{\partial R_{CRF}(S)}{\partial S^k} = -2WS^k. \qquad (11)$$

For DenseCRF where $W$ is fully connected Gaussian, computing the gradient (11) becomes a standard Bilateral filtering problem, for which many fast methods were proposed [16,32]. We implement our loss layers using fast Gaussian filtering [16], which is also utilized in the inference of DenseCRF [8,9]. Using the same fast filtering component, we can also computer the following gradient (12) of our Kernel Cut loss (6) in linear time. Note that our CRF and KC loss layer is much faster than CRF inference layer [13,9] since no iterations is needed.

$$\frac{\partial R_{KC}(S)}{\partial S^k} = -2WS^k + \gamma \frac{S^{k'}\hat{W}S^k d}{(d'S^k)^2} - \gamma \frac{2\hat{W}S^k}{d'S^k}. \qquad (12)$$

## 4.1  Comparison of regularized losses

Tab. 1 summaries the results with different regularized losses. Here we report both result with or without CRF post-processing on various networks. The base-

---

| | Weak | | | | Full |
|---|---|---|---|---|---|
| | CE only | w/ NC [3] | w/ CRF | w/ KernelCut | |
| DeepLab-MSc-largeFOV | 56.0 (8.1) | 60.5 (3.6) | 63.1 (1.0) | **63.5 (0.6)** | 64.1 |
| DeepLab-MSc-largeFOV+CRF | 62.0 (6.7) | 65.1 (3.6) | 65.9 (2.8) | **66.7 (2.0)** | 68.7 |
| DeepLab-VGG16 | 60.4 (8.4) | 62.4 (6.4) | 64.4 (4.4) | **64.8 (4.0)** | 68.8 |
| DeepLab-VGG16+CRF | 64.3 (7.2) | 65.2 (6.3) | 66.4 (5.1) | **66.7 (4.8)** | 71.5 |
| DeepLab-ResNet101 | 69.5 (6.1) | 72.8 (2.8) | 72.9 (2.7) | **73.0 (2.6)** | 75.6 |
| DeepLab-ResNet101+CRF | 72.8 (4.0) | 74.5 (2.3) | **75.0 (1.8)** | **75.0 (1.8)** | 76.8 |

Table 1: mIOU on PASCAL VOC2012 $val$ set. Our flexible framework allows various types of regularization losses for weakly supervised segmentation, e.g. normalized cut, CRF or their combinations (KernelCut [10]) as joint loss. We achieved the state-of-the-art with scribbles. In () shows the offset to the result with full masks.

lines are with cross entropy losses of full labeled masks or partial seeds ignoring unlabeled region. We choose the weight of the regularization term to achieve the best validation accuracy. The state-of-the-art of scribble-based segmentation is from prior work [3] with extra normalized cut loss. Consistently over different networks, using the proposed CRF loss outperforms that with normalized cut loss. Our best result is obtained when combining both normalized cut loss and DenseCRF loss. Clearly, utilization of CRF loss and KernelCut loss reduce the gap toward the full supervision baseline. With DeepLab-MSc-largeFOV followed by CRF post processing, using KernelCut regularized loss achieved mIOU of 66.7%, while previous best is 65.1% with normalized cut loss [3]. Our result with scribbles approaches **97.6%** of the quality of that with full supervision, yet only **3%** of all pixels are scribbled. This paper pushes the limit of weakly supervised segmentation.

To get some intuition about these losses and their regularization effect, we visualize their gradient w.r.t. segmentation $\frac{\partial R(S)}{\partial S}$ in Fig. 1. Note that the *sign* of gradients indicates whether to encourage or discourage certain labeling. The color coded gradients clearly show evidence toward better color clustering /edge alignment/ object separation with regularized loss. The gradients of different losses are slightly different. Since kernel cut is the combination of normalized cut with CRF, then its gradient is the sum of that of each.

Fig. 2 shows some qualitative examples with different losses. Results with regularized loss is better than that without. Besides, the segmentation with kernel cut loss have better edge alignment compared to that with normalized cut loss. This is because of the extra pairwise CRF loss. The effect of CRF loss and normalized cut loss is different. Our Kernel Cut loss combines the benefit of both regional color clustering (normalized cut) and pairwise regularization (DenseCRF). By combining both we can achieve better segmentation regularization.

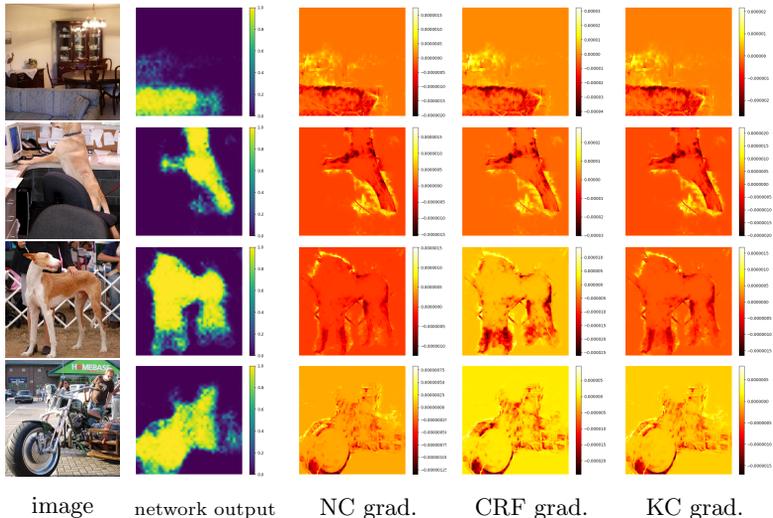| image | network output | NC grad. | CRF grad. | KC grad. |

Fig. 1: Visualization of the gradient for different losses. The negative (positive) gradients are coded in red (yellow). For example, negative gradients on the sofa drives the network to predict "sofa" for these pixels. Also note how the dog pops out in the gradient map.

### 4.2 Direct loss vs proposal generation

Here we compare our direct loss and proposal generation methods (Sec. 3) in weakly supervised setting mainly focusing on scribbles. Proposals can be generated *offline* or *online*. One straightforward proposal method is to treat GrabCut output as "fake" ground truth for training. ScribbleSup [17] refines GrabCut output using network predicted segmentation as unary potentials. The proposals are updated but are generated offline. By online proposal generation, we let network output go through a CRF inference layer during training at each iteration. The loss for proposal generation is the cross entropy between the input and output of the CRF inference layer, see Sec.3. A recent work that generates proposals online for tag-based weakly-supervised segmentation is SEC [13].

Table 2 compares our direct loss method to proposal generation variants above. We used the public implementation of SEC's *constrain-to-boundary loss*[6] that combines explicit dense CRF proposal layer and cross entropy loss between the proposal and network output. We report the results for SEC*, our adaptation of tag-based SEC to weak-supervision with scribbles from [17]. We find that (frequent) online proposal updates give better results than those with fixed proposals. Compared to our direct loss method, (online) proposal generation gives inferior segmentation accuracy over different networks, see Table 2.

We further evaluate online proposal generation. Figure 3 compares it to our regularized loss method in terms of segmentation accuracy and obtained loss

---

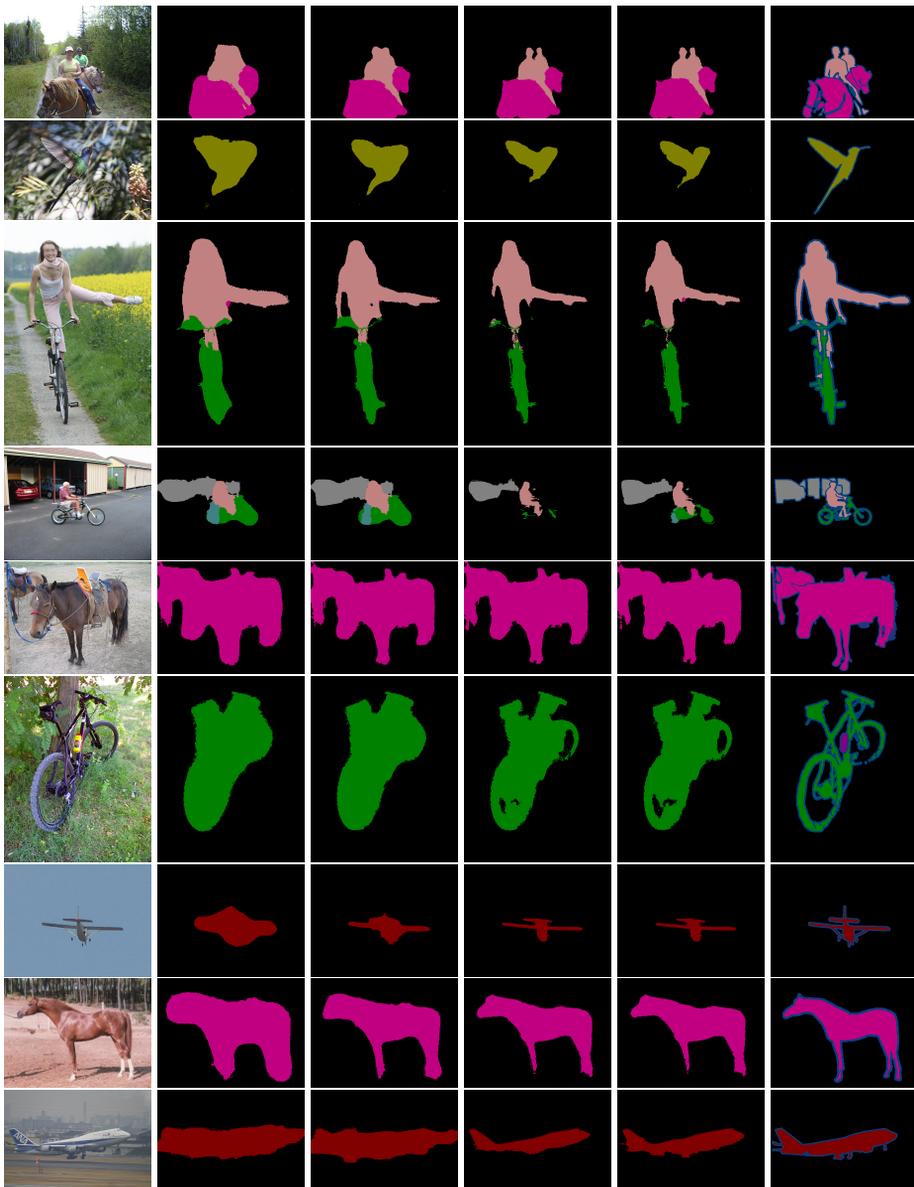[6] https://github.com/kolesman/SEC

image      CE loss only  w/ NC loss [3]  w/ CRF loss   KernelCut loss  ground truth

Fig. 2: Examples on PASCAL VOC *val* set. Kernel cut as regularization loss gives qualitatively better result than that with normalized cut loss. We found kernel cut results to have better edge alignment.

values. Even though the proposal generation scheme indirectly minimizes our

| | Weak | | | | Full |
|---|---|---|---|---|---|
| | proposal generation | | | direct loss | |
| | GrabCut (one time) | ScribbleSup (iterative) | SEC* (online) | CRF loss | |
| DeepLab-MSc-largeFOV | 55.5 | n/a | 61.3 | **63.1** | 64.1 |
| DeepLab-MSc-largeFOV+CRF | 59.7 | 63.1 | 65.4 | **65.9** | 68.7 |
| DeepLab-VGG16 | 59.0 | n/a | 63.4 | **64.4** | 68.8 |
| DeepLab-ResNet101 | 63.9 | n/a | 72.5 | **72.9** | 75.6 |

Table 2: Results using weak supervision (scribbles). The baseline is training with interactive GrabCut output. ScribbleSup [17] alternates between GrabCut and CNN training, but the proposals are generated offline. It helps to have frequent online proposal updates at each iteration of training as in SEC*, our adaptation of tag-based SEC [13] to weak supervision with scribbles in [17]. The best (quality and speed) training is based on simple direct loss optimization avoiding proposal generations. This comparison uses the same dense CRF Gaussian bandwidths.

regularized loss, such training scheme gives higher loss values than those obtained with our direct loss minimization. Also, direct loss minimization gives higher mIOUs for the training and validation.

As mentioned earlier, SEC [13] was originally focused on tag-based supervision and Table 3 reports some tests for that form of weak supervision. We compare SEC with its simplification replacing their constrain-to-boundary loss by our direct regularization loss. We train using different combinations of losses for supervision based on image-level labels/tags. Our CRF loss helps to improve training to 43.9% compared to 38.4% without it. There is only small improvement in segmentation mIOU when replacing constrain-to-boundary loss by CRF loss. However, the direct loss layer is several times faster than SEC integrating explicit proposal layer. The segmentation accuracy and overall training speed are also reported in Tab. 3. The results are for the DeepLab-largeFOV network since it is fast to train. We also tested a variant of SEC without (CRF) proposal layer back-propagation, which we show is redundant in practice.
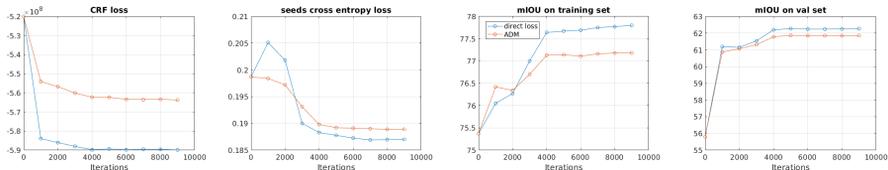


Fig. 3: Our direct loss scheme achieves better mIOU accuracy on training and val set. The CRF loss and seeding loss of our trained model are also less than that with proposal generation scheme. For fair comparison, our CRF loss and the CRF inference layer in proposal generation method have the same Gaussian kernel in this experiment.
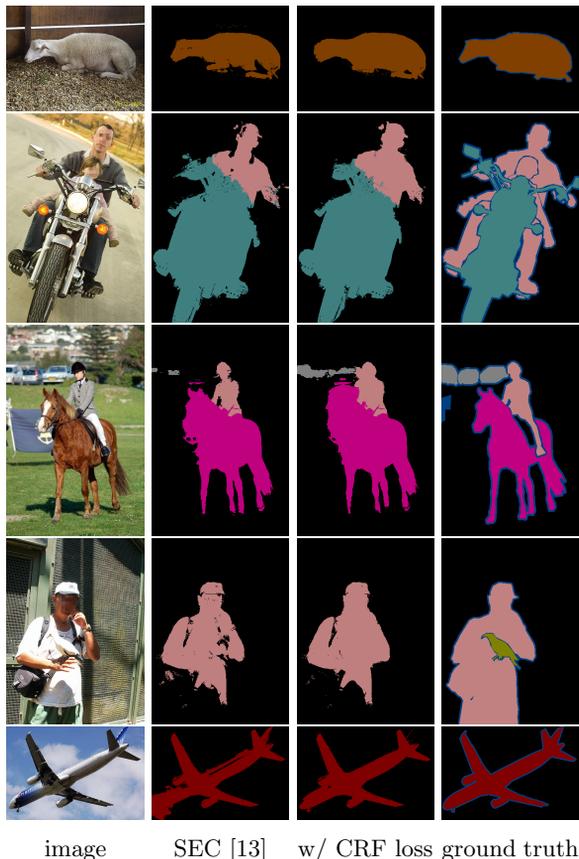
image      SEC [13]    w/ CRF loss ground truth

Fig. 4: Examples on PASCAL VOC *val* set for supervision with image-level labels (tags). We train using the seeding loss, expansion loss in SEC [13] and our CRF loss. Similar segmentation is obtained yet we avoid any iterative CRF inference and have the direct loss instead.

Fig. 4 shows testing examples for our method and SEC with image tags as supervision. Using direct loss rather than the constrain-to-boundary loss gives similar segmentation, while being faster to train since no inference is needed.

To see the limit of our algorithm with scribble supervision, we train with shortened scribbles visualized in Fig. 5. Note that with length zero, there is only one click or spot for each object. For different length ratios from zero to 100%, our direct loss method achieved much better segmentation than ScribbleSup [17], see Fig. 6. The improvement over ScribbleSup [17] is more significant for shorter scribbles or even clicks.

|  |  | include this loss? | | | |
|---|---|---|---|---|---|
| Losses | Seeding loss [13] | ✓ | ✓ | ✓ | ✓ |
|  | Expansion loss [13] | ✓ | ✓ | ✓ | ✓ |
|  | Constrain-to-boundary loss [13] |  | ✓ | ⋆ |  |
|  | Our direct CRF loss |  |  |  | ✓ |
| mIOU (%) | | 38.4 | 43.7 | 43.8 | 43.9 |
| Overall training time in s/batch | | 0.86 | 1.19 (0.33) | 1.19 (0.33) | 0.98 (0.12) |

Table 3: Tag-based weak supervision. We train with different combinations of the losses in SEC [13] and our CRF loss. Replacing the constrain-to-boundary loss in SEC [13] by direct CRF loss gives minor improvement in accuracy, but training with our direct loss is faster since no iterative CRF inference is needed. We also compare to a variant (⋆) of SEC without back-propagation of the CRF inference layer. Parenthesis (·) show the computational times for the constrain-to-boundary loss layer or our direct loss layer.
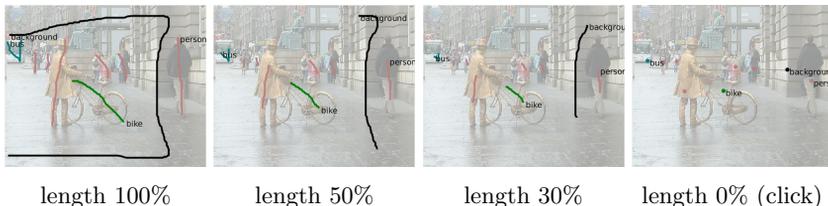


length 100%        length 50%        length 30%        length 0% (click)

Fig. 5: Similar to [17], we shorten the scribbles with different length ratios. With length zero (clicks) is the most challenging case for training.

## 4.3    Fully and semi supervised segmentation

We've demonstrated the usefulness of regularized loss for weakly supervised segmentation. Here we test if it also helps full supervision or semi-supervision with extra unlabeled images. For full supervision, we add NC loss on labeled masks besides the cross entropy loss. This experiment is on a simple saliency dataset [33] where color clustering is obvious and likely to help. As shown in Tab. 4, when we increase the weight of $R_{NC}(S)$, we indeed obtained segmentation that
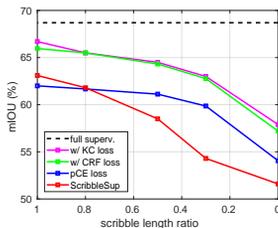


Fig. 6: mIOU on *val* set when train with shorter scribbles or clicks (length zero).

| NC loss weight | mIOU | cross entropy loss | NC loss |
|:---:|:---:|:---:|:---:|
| 0 | 89.85% | 0.106 | 0.536 |
| 0.05 | 89.53% | 0.108 | 0.526 |
| 0.1 | 89.38% | 0.110 | 0.517 |
| 0.2 | 89.39% | 0.112 | 0.509 |
| 0.5 | 88.75% | 0.125 | 0.485 |

Table 4: Negative effect of regularization loss for full supervision.

is more regularized. However, with extra regularization loss during training, the cross entropy loss got worse and mIOU decreased. The conclusion is that imposing regularized loss naively on labeled images doesn't help fully supervision segmentation. Empirical risk minimization is in some sense optimal for fully labeled data. Extra regularization loss steers the network in the wrong direction if the regularization doesn't totally agree with the ground truth. Reporting this result though negative helps to complete our investigation of regularized loss for fully, weakly and semi-supervised settings.

For training with both labeled images and unlabeled images, our joint losses include cross entropy on labeled images and regularization on unlabeled ones. The 11K labeled images are from PASCAL VOC 2012 and the 10K unlabeled ones are from VOC 2007. We train DeepLab-LargeFOV with different amount of labeled & unlabeled images, see Tab. 5. For the baseline that can only utilize labeled images, the performance degrades with less masks, as expected. For our framework, the labeled and unlabeled images are mixed and randomly sampled in each batch. We observed 0.7% 1.9% improvement with our regularized loss. Note that this result is highly preliminary and detailed analysis of overfitting, generalization property and comparison to recent semi-supervised segmentation [34] with extra unlabeled images will be our future work.

## 5   Conclusion and Future Work

*Regularized semi-supervised loss* is a principled approach to semi-supervised deep learning [1,2], in general. We utilize such principle for weakly supervised CNN segmentation. In particular, this paper is continuation of the study of losses motivated by standard shallow segmentation [3]. While [3] is entirely on normalized cut loss, in this paper we propose and evaluate several regularized loss

| training data | # of labeled images | 11K | 11K | 7K | 5K | 3K |
|:---:|:---|:---:|:---:|:---:|:---:|:---:|
| | # of unlabeled images | 10K | 0 | 4K | 6K | 8K |
| losses | cross entropy only | 63.5% | 63.5% | 61.5% | 60.1% | 57.6% |
| | cross entropy + CRF reg. | 64.6% | 63.5% | 63.4% | 61.8% | 58.3% |

Table 5: Our regularization loss $R_{CRF}(S)$ on unlabeled images help to improve semi-supervised segmentation.

for weakly-supervised CNN segmentation based on Potts/CRF [5,8], normalized cut [4] and KernelCut [10] regularizer. DenseCRF [8] is very popular as post-processing [11] or trainable layer [12] for CNN segmentation. We are the first to use a relaxed version of DenseCRF directly as part of the loss.

In contrast to our direct regularized loss approach, the main stream in weakly supervised segmentation rely on generating "fake" full masks from partial input and train a network to match the proposals [17,25,26,13,14,27]. Proposals can be pre-computed or iteratively updated. Some work even back-propagate the proposal generation step [25,13]. We show that proposal methods can be viewed as approximate alternating direction method (ADM) for optimization of our direct loss. Using direct loss gives better optimization while being more efficient than proposal generation scheme since no CRF inference is needed.

This paper pushes the limit of weakly-supervised segmentation. Comprehensive experiments (Sec.4) with our regularized weakly supervised losses show (1) state-of-the-art performance for weakly supervised CNN segmentation reaching near full-supervision accuracy and (2) better quality and efficiency than proposal generating methods or normalized cut loss [3]. Alternating schemes (proposal generation) give higher loss at convergence. Besides for weak supervision, we also report our preliminary results for full and semi-supervision with unlabeled images.

In principle, any differentiable loss function fits our regularized loss framework. Exploring other relaxations of CRF as losses [18,19,20,21,22,23] and corresponding efficient gradient computation is left for future work. Also it would be interesting to apply our CRF regularized loss framework for weakly-supervised computer vision problems other than segmentation.

## A   Mean-field inference for DenseCRF

Here we show that the iterative parallel mean-field inference [8] indeed minimizes (9) with pairwise DenseCRF regularizer and unary potentials $\tilde{S}_p$ (e.g. given by network).

$$E(X) = \sum_p H(X_p, \tilde{S}_p) + \lambda R_{CRF}(X) - \sum_p H(X_p).$$

For positive semidefinite affinity matrix $W$, e.g. with Gaussian Kernel,

$$R_{CRF}(X) = \sum_k X^{k'} W(\mathbf{1} - X^k) \stackrel{c}{=} -\sum_k X^{k'} W X^k$$

is concave[7]. Since the cross entropy $H(X_p, \tilde{S}_p)$ is linear and the negative entropy $-H(X_p)$ is convex w.r.t. $X_p$, the concave-convex procedure (CCCP) can

---

7 $\stackrel{c}{=}$ means up to an additive constant.

iteratively solve an approximation of $E(X)$ by linearizing the concave part at $\tilde{S}$.

$$A(X) = \sum_p H(X_p, \tilde{S}_p) + \lambda \langle X, \nabla R_{CRF}(X)|_{\tilde{S}} \rangle - \sum_p H(X_p).$$

$$= \sum_p H(X_p, \tilde{S}_p) - 2\lambda \sum_p \sum_k X_p^k \cdot [W\tilde{S}^k]_p - \sum_p H(X_p).$$

$$= -\sum_k X_p^k \cdot \log \tilde{S}_p^k - 2\lambda \sum_k X_p^k \cdot [W\tilde{S}^k]_p - \sum_p H(X_p).$$

KKT approach for minimizing $A(X)$ subject to probability simplex constraints $\sum_k X_p^k = 1$ yields the following optima,

$$\arg\min_X A(X) = \frac{1}{z_p} \exp\{-\log \tilde{S}_p^k - 2\lambda[W\tilde{S}^k]_p\}, \tag{13}$$

where $z_p$ is a normalization constant for softmax. (13) is exactly the mean-field update for dense CRF [8]. Note that the updates (13) is also justified in a similar way in [24] for convergent optimization of KL distance between factorial marginal distribution and Gibbs distribution induced by CRF. Our justification of (13) is different. We show alternative interpretation of mean-field updates (13) as minimizing CRF potential $R_{CRF}(X)$ plus negative entropy $-H(X)$.

# References

1. Weston, J., Ratle, F., Mobahi, H., Collobert, R.: Deep learning via semi-supervised embedding. In: Neural Networks: Tricks of the Trade. Springer (2012) 639–655
2. Goodfellow, I., Bengio, Y., Courville, A.: Deep learning. MIT press (2016)
3. Tang, M., Djelouah, A., Perazzi, F., Boykov, Y., Schroers, C.: Normalized Cut Loss for Weakly-supervised CNN Segmentation. In: IEEE conference on Computer Vision and Pattern Recognition (CVPR), Salt Lake City (June 2018)
4. Shi, J., Malik, J.: Normalized cuts and image segmentation. IEEE Transactions on Pattern Analysis and Machine Intelligence **22** (2000) 888–905
5. Boykov, Y., Veksler, O., Zabih, R.: Fast approximate energy minimization via graph cuts. IEEE transactions on Pattern Analysis and Machine Intelligence **23**(11) (November 2001) 1222–1239
6. Boykov, Y., Jolly, M.P.: *Interactive graph cuts* for optimal boundary & region segmentation of objects in N-D images. In: ICCV. Volume I. (July 2001) 105–112
7. Rother, C., Kolmogorov, V., Blake, A.: Grabcut - interactive foreground extraction using iterated graph cuts. In: ACM trans. on Graphics (SIGGRAPH). (2004)
8. Krahenbuhl, P., Koltun, V.: Efficient inference in fully connected CRFs with Gaussian edge potentials. In: NIPS. (2011)
9. Zheng, S., Jayasumana, S., Romera-Paredes, B., Vineet, V., Su, Z., Du, D., Huang, C., Torr, P.H.: Conditional random fields as recurrent neural networks. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1529–1537
10. Tang, M., Marin, D., Ayed, I.B., Boykov, Y.: Normalized Cut meets MRF. In: European Conference on Computer Vision (ECCV), Amsterdam, Netherlands (October 2016)

11. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. arXiv:1606.00915 (2016)
12. Arnab, A., Zheng, S., Jayasumana, S., Romera-Paredes, B., Larsson, M., Kirillov, A., Savchynskyy, B., Rother, C., Kahl, F., Torr, P.: Conditional random fields meet deep neural networks for semantic segmentation. IEEE Signal Processing Magazine (2017)
13. Kolesnikov, A., Lampert, C.H.: Seed, expand and constrain: Three principles for weakly-supervised image segmentation. In: European Conference on Computer Vision (ECCV), Springer (2016)
14. Papandreou, G., Chen, L.C., Murphy, K.P., Yuille, A.L.: Weakly-and semi-supervised learning of a deep convolutional network for semantic image segmentation. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1742–1750
15. Rajchl, M., Lee, M.C., Oktay, O., Kamnitsas, K., Passerat-Palmbach, J., Bai, W., Damodaram, M., Rutherford, M.A., Hajnal, J.V., Kainz, B., et al.: Deepcut: Object segmentation from bounding box annotations using convolutional neural networks. IEEE transactions on medical imaging **36**(2) (2017) 674–683
16. Adams, A., Baek, J., Davis, M.A.: Fast high-dimensional filtering using the per-mutohedral lattice. Computer Graphics Forum **29**(2) (2010) 753–762
17. Lin, D., Dai, J., Jia, J., He, K., Sun, J.: Scribblesup: Scribble-supervised convolutional networks for semantic segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 3159–3167
18. Chambolle, A., Darbon, J.: On total variation minimization and surface evolution using parametric maximum flows. International Journal of Computer Vision **84**(3) (April 2009) 288
19. Chan, T.F., Esedoglu, S., Nikolova, M.: Algorithms for finding global minimizers of image segmentation and denoising models. SIAM journal on applied mathematics **66**(5) (2006) 1632–1648
20. Pock, T., Chambolle, A., Cremers, D., Bischof, H.: A convex relaxation approach for computing minimal partitions. In: IEEE conference on Computer Vision and Pattern Recognition (CVPR). (2009)
21. Couprie, C., Grady, L., Najman, L., Talbot, H.: A unifying graph-based optimization framework. IEEE Transactions on Pattern Analysis and Machine Intelligence **33**(7) (July 2011) 1384–1399
22. Desmaison, A., Bunel, R., Kohli, P., Torr, P.H., Kumar, M.P.: Efficient continuous relaxations for dense crf. In: European Conference on Computer Vision, Springer (2016) 818–833
23. Thalaiyasingam, A., Desmaison, A., Bunel, R., Salzmann, M., Torr, P.H., Kumar, M.P.: Efficient linear programming for dense crfs. In: Conference on Computer Vision and Pattern Recognition. (2017)
24. Krähenbühl, P., Koltun, V.: Parameter learning and convergent inference for dense random fields. In: International Conference on Machine Learning (ICML). (2013)
25. Vernaza, P., Chandraker, M.: Learning random-walk label propagation for weakly-supervised semantic segmentation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Volume 3. (2017)
26. Khoreva, A., Benenson, R., Hosang, J., Hein, M., Schiele, B.: Simple does it: Weakly supervised instance and semantic segmentation. In: 30th IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2017), Honolulu, HI, USA (2017)

27. Dai, J., He, K., Sun, J.: Boxsup: Exploiting bounding boxes to supervise convolutional networks for semantic segmentation. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1635–1643
28. Grady, L.: Random walks for image segmentation. IEEE transactions on pattern analysis and machine intelligence **28**(11) (2006) 1768–1783
29. Chapelle, O., Schölkopf, B., Zien, A., eds.: Semi-Supervised Learning. MIT Press, Cambridge, MA (2006)
30. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. Foundations and Trends in Machine Learning **3**(1) (2011) 1–122
31. Baque, P., Bagautdinov, T.M., Fleuret, F., Fua, P.: Principled parallel mean-field inference for discrete random fields. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (2016)
32. Paris, S., Durand, F.: A fast approximation of the bilateral filter using a signal processing approach. International journal of computer vision **81**(1) (2009) 24–52
33. Cheng, M.M., Mitra, N.J., Huang, X., Torr, P.H.S., Hu, S.M.: Global contrast based salient region detection. IEEE TPAMI **37**(3) (2015) 569–582
34. Hung, W.C., Tsai, Y.H., Liou, Y.T., Lin, Y.Y., Yang, M.H.: Adversarial learning for semi-supervised semantic segmentation. arXiv preprint arXiv:1802.07934 (2018)