

Semantic Image Segmentation With Deep
Convolution Nets And Fully Connected CRFs



ABSTRACT

Deep Convolutional Neural Networks (DCNNs) have recently shown state of the art performance in high level vision tasks, such as image classification and object detection. This work brings together methods from DCNNs and probabilistic graphical models for addressing the task of **pixel-level classification** (also called "semantic image segmentation"). We show that responses at the final layer of DCNNs are not sufficiently localized for accurate object segmentation. This is due to the very invariance properties that make DCNNs good for high level tasks. We overcome this poor localization property of deep networks by combining the responses at the final DCNN layer with a fully connected Conditional Random Field (CRF). Qualitatively, our "DeepLab" system is able to localize segment boundaries at a level of accuracy which is beyond previous methods. Quantitatively, our method sets the new state-of-art at the PASCAL VOC-2012 semantic image segmentation task, reaching 71.6% IOU accuracy in the test set. We show how these results can be obtained efficiently: Careful network re-purposing and a novel application of the 'hole' algorithm from the wavelet community allow dense computation of neural net responses at 8 frames per second on a modern GPU.

픽셀단위로 분류를 하는 방법을 사용한다.

DCNN에서 기존에 사용하는 Segmentation Image Segmentation (71.6%) 정밀도가 높은
증명하고 동시에 pixel 단위로 분류하기 위해 "DeepLab" system을 소개하였다.

- 사용한 기법으로
- CRF
- Hole Algorithm이 소개된다.

There are two technical hurdles in the application of DCNNs to image labeling tasks: **signal downsampling**, and **spatial ‘insensitivity’** (invariance). The first problem relates to the reduction of signal resolution incurred by the repeated combination of max-pooling and downsampling ('striding') performed at every layer of standard DCNNs (Krizhevsky et al., 2013; Simonyan & Zisserman, 2014; Szegedy et al., 2014). Instead, as in Papandreou et al. (2014), we employ the ‘atrous’ (with holes) algorithm originally developed for efficiently computing the undecimated discrete wavelet transform (Mallat, 1999). This allows efficient dense computation of DCNN responses in a scheme substantially simpler than earlier solutions to this problem (Giusti et al., 2013; Sermanet et al., 2013).

The second problem relates to the fact that obtaining object-centric decisions from a classifier requires invariance to spatial transformations, inherently limiting the spatial accuracy of the DCNN model. We boost our model’s ability to capture fine details by employing a fully-connected Conditional Random Field (CRF). Conditional Random Fields have been broadly used in semantic segmentation to combine class scores computed by multi-way classifiers with the low-level information captured by the local interactions of pixels and edges (Rother et al., 2004; Shotton et al., 2009) or superpixels (Lucchi et al., 2011). Even though works of increased sophistication have been proposed to model the hierarchical dependency (He et al., 2004; Ladicky et al., 2009; Lemitsky et al., 2011) and/or high-order dependencies of segments (Delong et al., 2012; Gonfaus et al., 2010; Kohli et al., 2009; Chen et al., 2013; Wang et al., 2015), we use the fully connected pairwise CRF proposed by Krähenbühl & Koltun (2011) for its efficient computation, and ability to capture fine edge details while also catering for long range dependencies. That model was shown in Krähenbühl & Koltun (2011) to largely improve the performance of a boosting-based pixel-level classifier, and in our work we demonstrate that it leads to state-of-the-art results when coupled with a DCNN-based pixel-level classifier.

The three main advantages of our “DeepLab” system are (i) **speed**: by virtue of the ‘atrous’ algorithm, our dense DCNN operates at 8 fps, while Mean Field Inference for the fully-connected CRF requires 0.5 second, (ii) **accuracy**: we obtain state-of-the-art results on the PASCAL semantic segmentation challenge, outperforming the second-best approach of Mostajabi et al. (2014) by a margin of 7.2% and (iii) **simplicity**: our system is composed of a cascade of two fairly well-established modules, DCNNs and CRFs.

① DCNN의 단점

- Signal downsampling

• 반복되는 pooling layer + downsampling이기 때문에

• 대상의 세밀한 경계선을 보완하는게 어렵다.

• 이를弥补하기 위해 atrous network를

제작하였다.

• 이를 통해서 더 dense한 이미지의 결과를 얻을 수 있다.

• 하지만 매우 느리다.

- Spatial insensitivity

• 대상의 위치에 따라 학습된 특징이 적용되면서 결과가 달라진다.

Moving towards works that lie closer to our approach, several other researchers have considered the use of convolutionally computed DCNN features for dense image labeling. Among the first have been Farabet et al. (2013) who apply DCNNs at multiple image resolutions and then employ a segmentation tree to smooth the prediction results; more recently, Hariharan et al. (2014a) propose to concatenate the computed inter-mediate feature maps within the DCNNs for pixel classification, and Dai et al. (2014) propose to pool the inter-mediate feature maps by region proposals. Even though these works still employ segmentation algorithms that are decoupled from the DCNN classifier's results, we believe it is advantageous that segmentation is only used at a later stage, avoiding the commitment to premature decisions.

More recently, the segmentation-free techniques of (Long et al., 2014; Eigen & Fergus, 2014) directly apply DCNNs to the whole image in a sliding window fashion, replacing the last fully connected layers of a DCNN by convolutional layers. In order to deal with the spatial localization issues outlined in the beginning of the introduction, Long et al. (2014) upsample and concatenate the scores from inter-mediate feature maps, while Eigen & Fergus (2014) refine the prediction result from coarse to fine by propagating the coarse results to another DCNN.

The main difference between our model and other state-of-the-art models is the combination of pixel-level CRFs and DCNN-based 'unary terms'. Focusing on the closest works in this direction, Cogswell et al. (2014) use CRFs as a proposal mechanism for a DCNN-based reranking system, while Farabet et al. (2013) treat superpixels as nodes for a local pairwise CRF and use graph-cuts for discrete inference; as such their results can be limited by errors in superpixel computations, while ignoring long-range superpixel dependencies. Our approach instead treats every pixel as a CRF node, exploits long-range dependencies, and uses CRF inference to directly optimize a DCNN-driven cost function. We note that mean field had been extensively studied for traditional image segmentation/edge detection tasks, e.g., (Geiger & Girosi, 1991; Geiger & Yuille, 1991; Kokkinos et al., 2008), but recently Krähenbühl & Kohl (2011) showed that the inference can be very efficient for fully connected CRF and particularly effective in the context of semantic segmentation.

DCNN의 classifier는 단계별로 inter-mediate feature map을 통해 classification, segmentation, labeling을 한다.

단계별 segmentation은 단계별로 이루어지며, 전처리하는 과정에서 일부를 빼거나 추가하는 과정이다.

spatial localization은 DCNN의 유통성을 통해 가능해지며, 기존 연구는

① upsample and concatenate scores from inter-mediate feature maps

② refine the prediction result from coarse to fine by propagating coarse results to another DCNN

단계별 segmentation은 단계별로 이루어지며, 그 다른 방법으로

combination of pixel-level CRFs & DCNN based 'unary terms'이다.

모든 픽셀은 CRF node인 경우에 따라, 다른 CRF의 결과물을 고려해 DCNN의 의해 계산되는 값을 보정하는 방식

A

DCNN (이경우는 VGG-16)은 어떻게 변형하는데 적용되는지 살펴

Herein we describe how we have re-purposed and finetuned the publicly available Imagenet pre-trained state-of-art 16-layer classification network of (Simonyan & Zisserman, 2014) (VGG-16) into an efficient and effective dense feature extractor for our dense semantic image segmentation system.

VGG-16 model은 classification model, backbone으로 활용하여 이미지를 통해 semantic segmentation을 위한 오류를 얻게된다.

3.1 EFFICIENT DENSE SLIDING WINDOW FEATURE EXTRACTION WITH THE HOLE ALGORITHM

Dense spatial score evaluation is instrumental in the success of our dense CNN feature extractor. As a first step to implement this, we convert the fully-connected layers of VGG-16 into convolutional ones and run the network in a convolutional fashion on the image at its original resolution. However this is not enough as it yields very sparsely computed detection scores (with a stride of 32 pixels). To compute scores more densely at our target stride of 8 pixels, we develop a variation of the method previously employed by Giusti et al. (2013); Sermanet et al. (2013). We skip subsampling after the last two max-pooling layers in the network of Simonyan & Zisserman (2014) and modify the convolutional filters in the layers that follow them by introducing zeros to increase their length (2x in

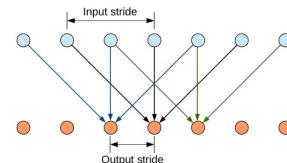
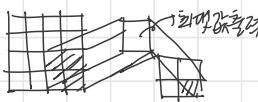


Figure 1: Illustration of the hole algorithm in 1-D, when $\text{kernel_size} = 3$, $\text{input_stride} = 2$, and $\text{output_stride} = 1$.

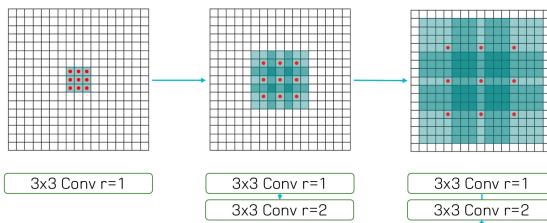
- pooling 층으로 인한 원본 이미지를 풀어 segmentation을
이용하여 Fully Connected Layer의 연결을 제거하는
방법을 학습해 왔다.
pooling 층에서 출력되는 Global Feature는 해당 객체의
각 부분을 각각의 미리학습된 (multi-scale) 이미지를 classification하는
방법이다.



pooling 층에서 풀어놓은
downsampling 층으로
특성화 단계를 줄인다.

- Dense prediction을 활용하는
① up convolution ② multi scale input을 사용하는
② pooling 층을 활용하는 두 가지 방법은
어디서 차이가 있는지?

$\cdot \text{optimizer} = \text{torch.optim.SGD}()$
 $\text{param} = [\text{f}_\theta \text{ for } \text{f}_\theta \in \text{model.parameters}()]$
 $\text{lr} = 0.005, \text{momentum} = 0.9,$
 $\text{weight_decay} = 0.0005$



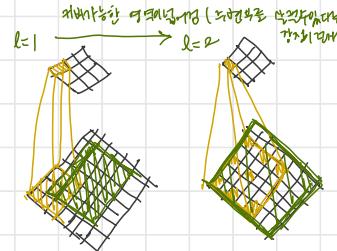
<Dilated Convolution>

$$(F * g_k)(p) = \sum_{s,t} F(s)g_k(t)$$

$\rightarrow l=1$: standard convolution

$\rightarrow l>1$: dilated convolution.

rate를 적용하는 경우
설정하는 경우
설정하는 경우



.: 디아이레이팅 된 후 풀어놓는다.
pooling 층에서는 풀어놓고 원래의 resolution을 놓는다.
(=downsampling)
↓
다음 층에 넣어서

3.2 CONTROLLING THE RECEPTIVE FIELD SIZE AND ACCELERATING DENSE COMPUTATION WITH CONVOLUTIONAL NETS

Another key ingredient in re-purposing our network for dense score computation is explicitly controlling the network's receptive field size. Most recent DCNN-based image recognition methods rely on networks pre-trained on the Imagenet large-scale classification task. These networks typically have large receptive field size: in the case of the VGG-16 net we consider, its receptive field is 224×224 (with zero-padding) and 404×404 pixels if the net is applied convolutionally. After converting the network to a fully convolutional one, the first fully connected layer has 4,096 filters of large 7×7 spatial size and becomes the computational bottleneck in our dense score map computation.

We have addressed this practical problem by spatially subsampling (by simple decimation) the first FC layer to 4×4 (or 3×3) spatial size. This has reduced the receptive field of the network down to 128×128 (with zero-padding) or 308×308 (in convolutional mode) and has reduced computation time for the first FC layer by 2 - 3 times. Using our Caffe-based implementation and a Titan GPU, the resulting VGG-derived network is very efficient: Given a 306×306 input image, it produces 39×39

작은 이미지로 처리
작은 계산량

CF-DNN

Deep Convolutional Neural Network

dense raw feature scores at the top of the network at a rate of about 8 frames/sec during testing. The speed during training is 3 frames/sec. We have also successfully experimented with reducing the number of channels at the fully connected layers from 4,096 down to 1,024, considerably further decreasing computation time and memory footprint without sacrificing performance, as detailed in Section 5. Using smaller networks such as Krizhevsky et al. (2013) could allow video-rate test-time dense feature computation even on light-weight GPUs.

- VGG's DCNN Normal Max-pooling layer 2개는 미리 학습한 내용을 사용함으로서 2개의 feature map은 모두 extra convolution을 이용해 더 넓은 receptive field를 확장할 수 있다.
Pooling을 적용할 때 같은 크기의 kernel을 이용한 convolution layer와 적용하지 않은 때는 차원마다 receptive field가 확장된다.
마지막 정확도를 위해 마지막 pooling layer는 각각 다른 atrous convolution과 함께 pooling을 한다.

다음은 예제.

- 그림에서 4x4의 16개의 뉴런을 만들면서 노드의 수를 줄이면서 (batch size 128개가 있다, 1024개)

첫번째 Fully Connected layers 4x4 / 3x3의 spatial size로 줄인다.



VGG-16 network의 pooling layer 274를
제거한 601 convolution layer로 바꾸면
diluted convolution을 적용

dilated convolution layers 320
789 미리 네트워크를 제거

4 DETAILED BOUNDARY RECOVERY: FULLY-CONNECTED CONDITIONAL RANDOM FIELDS AND MULTI-SCALE PREDICTION

4.1 DEEP CONVOLUTIONAL NETWORKS AND THE LOCALIZATION CHALLENGE

As illustrated in Figure 2, DCNN score maps can reliably predict the presence and rough position of objects in an image but are less well suited for pin-pointing their exact outline. There is a natural trade-off between classification accuracy and localization accuracy with convolutional networks: Deeper models with multiple max-pooling layers have proven most successful in classification tasks, however their increased invariance and large receptive fields make the problem of inferring position from the scores at their top output levels more challenging.

Recent work has pursued two directions to address this localization challenge. The first approach is to harness information from multiple layers in the convolutional network in order to better estimate the object boundaries (Long et al., 2014; Eigen & Fergus, 2014). The second approach is to employ a super-pixel representation, essentially delegating the localization task to a low-level segmentation method. This route is followed by the very successful recent method of Mostajabi et al. (2014).

In Section 4.2, we pursue a novel alternative direction based on coupling the recognition capacity of DCNNs and the fine-grained localization accuracy of fully connected CRFs and show that it is remarkably successful in addressing the localization challenge, producing accurate semantic segmentation results and recovering object boundaries at a level of detail that is well beyond the reach of existing methods.

(B) 디지털 이미지의 boundary는 디지털 사진 Localization Challenge는 디지털 사진

→ 예상한 결과는 fully connected CRF로 예상해서 문제를 해결하겠다 한다.

To overcome these limitations of short-range CRFs, we integrate into our system the fully connected CRF model of Krähenbühl & Koltun (2011). The model employs the energy function

$$E(x) = \sum_i \theta_i(x_i) + \sum_{ij} \theta_{ij}(x_i, x_j) \quad (1)$$

where x is the label assignment for pixels. We use as unary potential $\theta_i(x_i) = -\log P(x_i)$, where $P(x_i)$ is the label assignment probability at pixel i as computed by DCNN. The pairwise potential is $\theta_{ij}(x_i, x_j) = \mu(x_i, x_j) \sum_{m=1}^M w_m \cdot k^m(f_i, f_j)$, where $\mu(x_i, x_j) = 1$ if $x_i \neq x_j$, and zero otherwise (i.e., Potts Model). There is one pairwise term for each pair of pixels i and j in the image no matter how far from each other they lie, i.e., the model's factor graph is fully connected. Each k^m is the Gaussian kernel depends on features (denoted as f) extracted for pixel i and j and is weighted by parameter w_m . We adopt bilateral position and color terms, specifically, the kernels are

$$w_1 \exp\left(-\frac{\|p_i - p_j\|^2}{2\sigma_\alpha^2} - \frac{\|I_i - I_j\|^2}{2\sigma_\beta^2}\right) + w_2 \exp\left(-\frac{\|p_i - p_j\|^2}{2\sigma_\gamma^2}\right) \quad (2)$$

where the first kernel depends on both pixel positions (denoted as p) and pixel color intensities (denoted as I), and the second kernel only depends on pixel positions. The hyper parameters σ_α , σ_β and σ_γ control the "scale" of the Gaussian kernels.

Crucially, this model is amenable to efficient approximate probabilistic inference (Krähenbühl & Koltun, 2011). The message passing updates under a fully decomposable mean field approximation $b(x) = \prod_i b_i(x_i)$ can be expressed as convolutions with a Gaussian kernel in feature space. High-dimensional filtering algorithms (Adams et al., 2010) significantly speed-up this computation resulting in an algorithm that is very fast in practice, less than 0.5 sec on average for Pascal VOC images using the publicly available implementation of (Krähenbühl & Koltun, 2011).

CRF

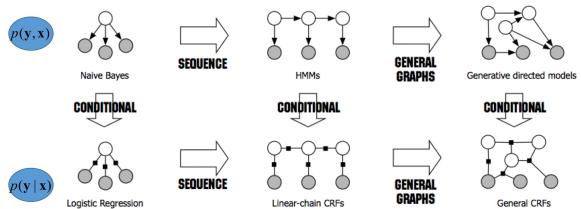
- Graph Model

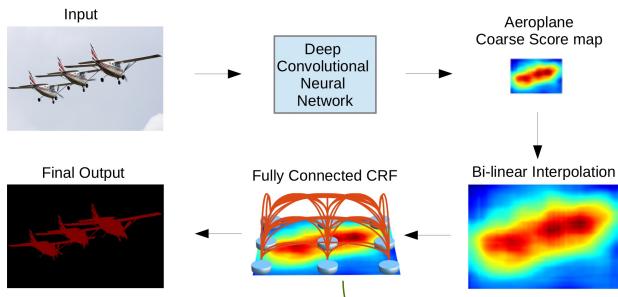
- X node = label of pixel

- X nodes = latent variable = Pixel (normal RGBB)

- X node helps in Modeling

- posteriori 확률 계산하는 모델 확률





Fully Connected CRFs

Maximize Posterior.

$$P(X|I) = \frac{1}{Z(I)} \exp(-\sum \phi_c(x_i|I))$$

(label) (Image) Normalization

• 노드는 각 계층마다 서로 연결되어
가장 먼저 CRF가 입력된 예제를 통해 살펴보자.
다음 예제에서 약간 다른 MeanField Approximation을 적용.

Minimize Energy

(= 흐름을 제어하는
최종 에너지)

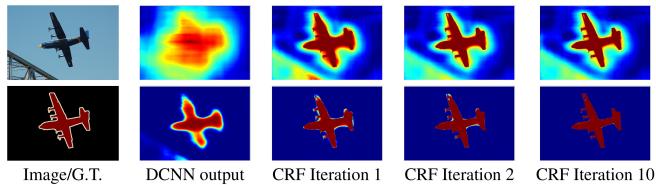
$$E(X|I) = \sum \phi_c(x_i|I)$$

$$E(x) = \sum \psi_c(x_i)$$

$$E(X) = \sum_i \psi_i(d_i) + \sum_{i,j} \psi_{ij}(d_i, d_j)$$

Unary Fully Connected

$$= \sum_i -\log P(d_i) + \sum_{i,j} \psi_{ij}(d_i, d_j)$$



↳ CRF는 주변에 대한 정보의 경계와
상호작용을 고려하는데 좋다.

Pairwise Term (Impose consistency of labeling) → defined over neighboring pixels

$$\psi_{ij}(d_i, d_j) = \left[\begin{array}{c} \psi(d_i, d_j) \\ W_1 \exp\left(-\frac{|p_i - p_j|^2}{2\sigma^2} - \frac{|I_i - I_j|^2}{2\sigma_f^2}\right) + W_2 \exp\left(-\frac{|p_i - p_j|^2}{2\sigma^2}\right) \end{array} \right]$$

$\psi(d_i, d_j) = \begin{cases} 1 & (d_i \neq d_j) \\ 0 & (\text{else}) \end{cases}$

∴ 퍼셀이 서로 비슷한 이미지인 경우 (유사한 색상 RGB 128, 128)

label들이 서로 다른 경우 energy가 증가하여 penalty를 얻게 된다.

* P_1, P_2 : 각각의 차이

* I_1, I_2 : 각각의 RGB

* W_1, W_2 : kernel weights

* $\theta_1, \theta_2, \theta_3$: Hyper parameter

정리 정리 정리 정리 정리 Mean Field Approximation 정리 정리.

$$P(X|I) = \frac{1}{Z(I)} \exp(-E(X))$$

정리 $Q(x) = \prod_i Q_i(x_i)$ 이 정리하고 $D_{KL}(Q||P)$ 은 확장화 하드를 허용한다.

$$Q_i(d_i = I) = \frac{1}{Z_i} \exp(-\psi_i(d_i) - \sum_j \mu(l, l') \sum_m \sum_s k^m(f_i, f_s) Q_s(l'))$$

↓

$$\psi_i(d_i, x_i) = \mu(d_i, d_s) \left(w_1 \exp\left(-\frac{|P_i - P_s|^2}{2\sigma_1^2} - \frac{|I_i - I_s|^2}{2\sigma_2^2}\right) + w_2 \exp\left(-\frac{|P_i - P_s|^2}{2\sigma_2^2}\right) \right)$$

↑ Smoothness kernel.

Update Rule

$$Q_i(d_i = l) = \frac{1}{Z_i} \exp\left(-\psi_i(d_i) - \sum_j \mu(l, l') \sum_m \sum_s k^m(f_i, f_s) Q_s(l')\right)$$

$$\text{정리 } Q_i(d_i = l) = \frac{1}{Z_i} \exp(-\psi_i(d_i))$$

수정함수 정리...

→ Gradient Kernel.

→ feature vectors for pixel i, j

MeanField Update rule for CRF

$$\tilde{Q}_i(l) = \sum_s k^m(f_i, f_s) Q_s(l) \quad \text{Message Passing}$$

→ weight of the m th kernel

$$Q_i(l) = \sum_n \mu(n) \tilde{Q}_i^n(l) \quad \text{Weighting (가중치 부여)}$$

→ label compatibility function

$$\hat{Q}_i(l) = \sum_{l'} \mu(l, l') Q_i(l) \quad \text{Compatibility Transform}$$

$$Q_i(l) = -(\psi_u(d_i) + \hat{Q}_i(l)) \quad \text{Adding Unary (Local Update)}$$

Normalization (softmax)

cf. 디자인 kernel approximation
마우스 커널 convolutional kernel
연산 학습에 적용하기 어려움
자연어 처리방법 이용하는데 어려움
제작 품질 고려해 kernel approximation