

U-Net



The typical use of convolutional networks is on classification tasks, where the output to an image is a single class label. However, in many visual tasks, especially in biomedical image processing, the desired output should include localization, i.e., a class label is supposed to be assigned to each pixel. Moreover, thousands of training images are usually beyond reach in biomedical tasks.

↳ biomedical task를 해결하기 위해서는 output이 기존의 binary classification class label을 출력하는 것과 더불어 "위치정보"도 포함할 수 있어야 한다.

또한, 근본적으로 biomedical task를 해결하기 위해서는 주변, 동일개의 주변 데이터를 적용해야 하기 때문에 해당 부분의 처리는 data augmentation을 통해 강화해 가지고 논문 이후의 부분에 등장한다.

**Abstract.** There is large consent that successful training of deep networks requires many thousand annotated training samples. In this paper, we present a network and training strategy that relies on the strong use of data augmentation to use the available annotated samples more efficiently. The architecture consists of a contracting path to capture context and a symmetric expanding path that enables precise localization. We show that such a network can be trained end-to-end from very few images and outperforms the prior best method (a sliding-window convolutional network) on the ISBI challenge for segmentation of neuronal structures in electron microscopic stacks. Using the same net-

→ ① 주변 사용  
Window et 알리  
end-to-end 학습  
HTB

→ 두 번째 네트워크의 구조는 설계하는 과정에서 알맞은 결과만 V-Net은 두 개의 path로 구성된다.

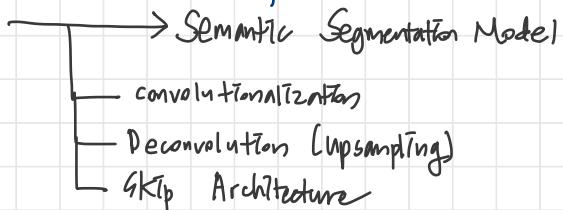
① contracting path  
— capture context

② expanding path  
— enables precise localization

# cf. FCN (Fully Connected Network)

convolution layer만 이용한 기계학  
학습 모델

(AlexNet, VGG16, GoogleNet)  
Image Classification Model

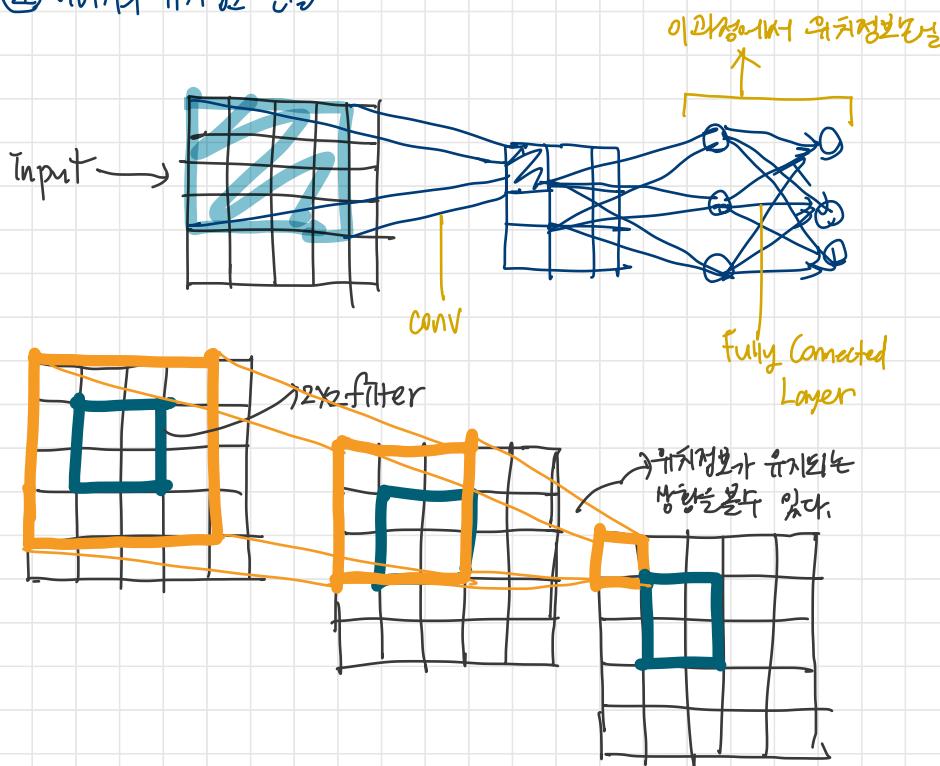


## 1) Convolutionalization

- 기존의 image classification 모델들은 흐리동에서 이전에 입력받은 특징맵들을 (convolution을 적용된) fully connected layer에 넣어서 최종 분류를 진행한다.

그러나 이를

### ① 이미지의 위치정보 병합



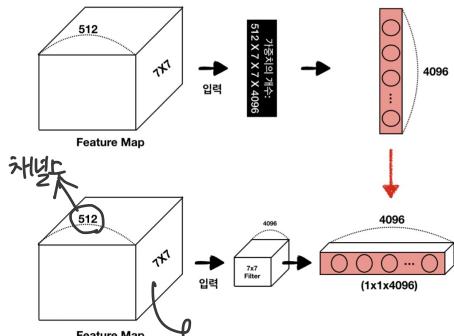
## ② 입력 이미지의 크기의 고정

fully-connected layer( fully connected)는 Dense layers를 사용한 것을 말한다.

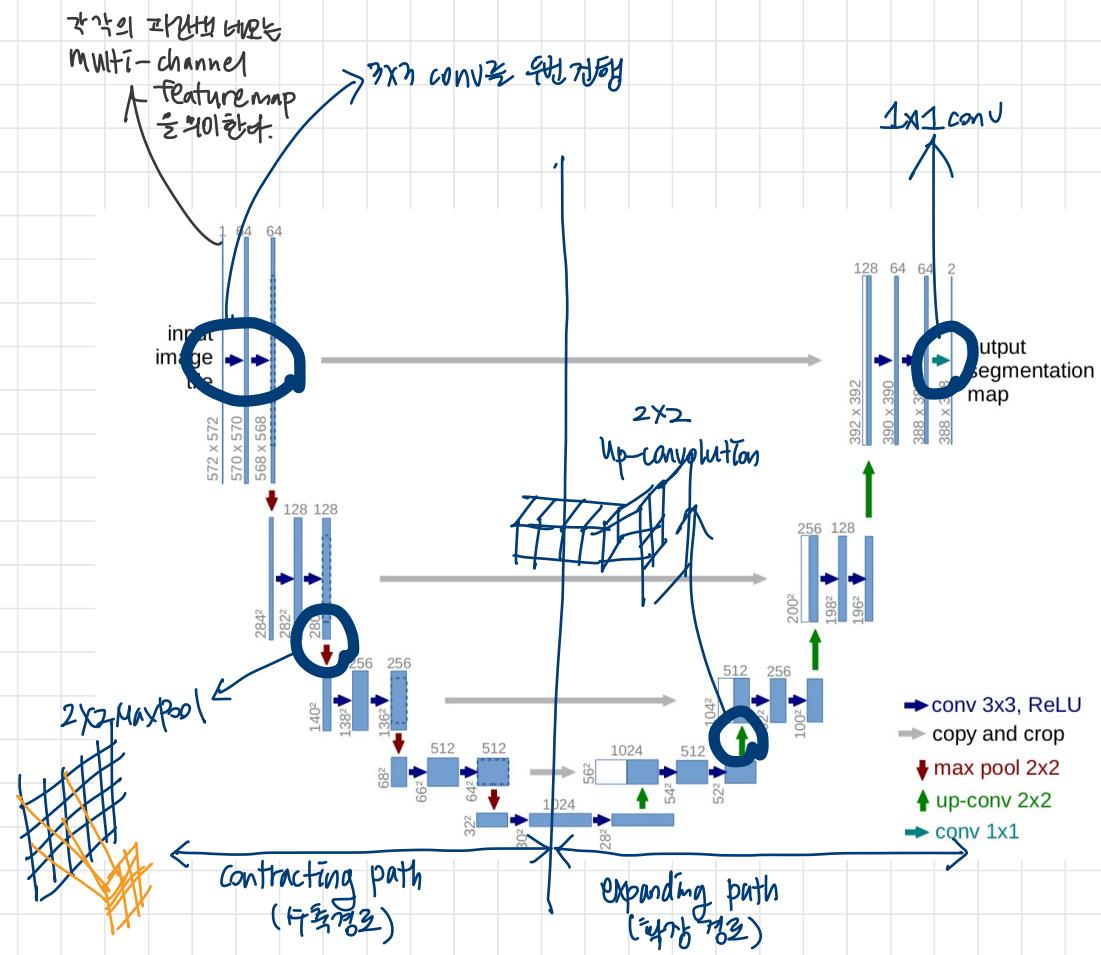
Dense-layer에는 가중치가 하나도 고려되어있지 않아서 fully-layer의 feature map의 크기 또한 고정이 되며, 따라서 입력 이미지의 크기로 고정된다.

그러나 segmentation의 경우에 목적이 원본 이미지를 각 픽셀 단위로 구분하고  
제작을 원할 때는 각 픽셀 단위로 위치별로의 정보가 더욱 중요하다.

\* 아래의 fclayer의 학습은 주로 이미지에서 모든 fclayers는  
Convolution Layer로 대체하였다. 즉 feature map을 coarse하는 방식이 관계된다.



↑ 각 feature map의 width와 height 같은 크기의  
filter를 사용해서 4096개의 가중치 개수를 유지할 수 있게 하면

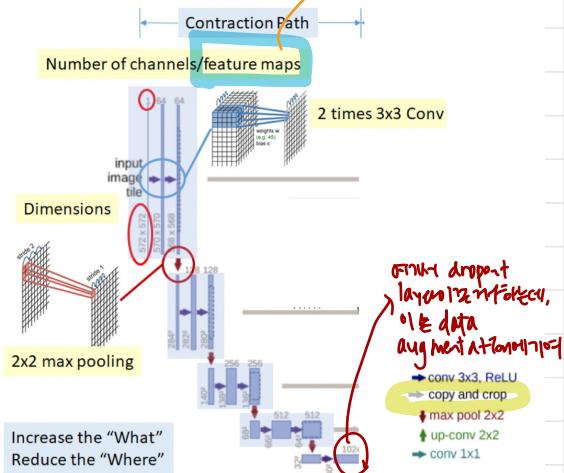


① ① ② ③ feature map의 context 활용

② DownSampling 과정은 차례로 feature map  
임

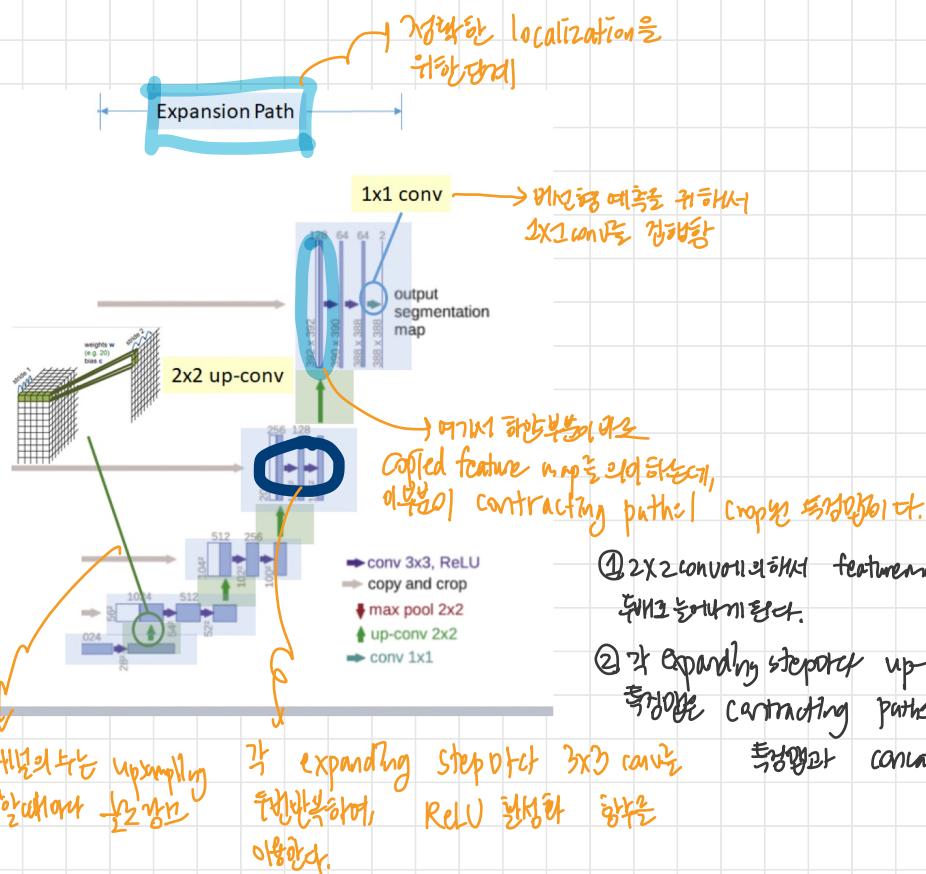
① 정확한 localization을 예상하기

② ③ ④ ⑤ layer의 결과를 concatenation 해  
: up-sampling은 전원 홉과의 feature map을  
정확히 맞춰주는 단계로 그 목적은 이미 된다



① 3x3 convolution layer에 padding 2로 특성맵의 크기를 감소한다.  
이제 padding = 1로 사용된다.

② 2x2 max pooling layer에서의 해상도 down sampling은 적용된다.  
padding = 2로 여기서 특성맵의 크기는 절반으로 줄어들게 된다.



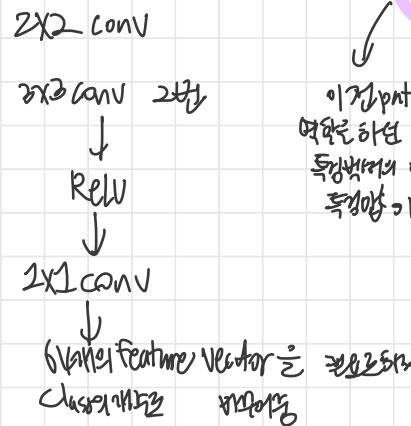
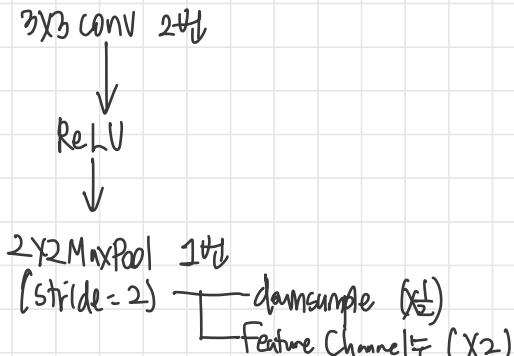
① 2x2 convolution의 해상도 feature map은 절반으로 늘어나게 된다.

② 각 expanding step마다 up-conv로  
특성맵과 contracting path cropped된  
특성맵과 concat

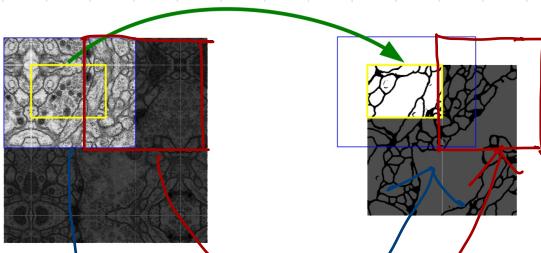
# ∴ 23개의 계층이 이루어진 FCN9201학.

The network architecture is illustrated in Figure 1. It consists of a contracting path (left side) and an expansive path (right side). The contracting path follows the typical architecture of a convolutional network. It consists of the repeated application of two  $3 \times 3$  convolutions (unpadded convolutions), each followed by a rectified linear unit (ReLU) and a  $2 \times 2$  max pooling operation with stride 2 for downsampling. At each downsampling step we double the number of feature channels. Every step in the expansive path consists of an upsampling of the feature map followed by a  $2 \times 2$  convolution ("up-convolution") that halves the number of feature channels, a concatenation with the correspondingly cropped feature map from the contracting path, and two  $3 \times 3$  convolutions, each followed by a ReLU. The cropping is necessary due to the loss of border pixels in every convolution. At the final layer a  $1 \times 1$  convolution is used to map each 64-component feature vector to the desired number of classes. In total the network has 23 convolutional layers.

To allow a seamless tiling of the output segmentation map (see Figure 2), it is important to select the input tile size such that all  $2 \times 2$  max-pooling operations are applied to a layer with an even x- and y-size.



Feature Map  
+  
Cropped Feature map



**Fig. 2.** Overlap-tile strategy for seamless segmentation of arbitrarily large images (here segmentation of neuronal structures in EM stacks). Prediction of the segmentation in the yellow area, requires image data within the blue area as input. Missing input data is extrapolated by mirroring

이미지의 유팽영역  
segmentation을 위한  
 tile에 대한 segmentation map을 구현하는  
방법의 일부를 표시합니다

## Overlap-tile Strategy

This strategy allows the seamless segmentation of arbitrarily large images by an overlap-tile strategy (see Figure 2). To predict the pixels in the border region of the image, the missing context is extrapolated by mirroring the input image. This tiling strategy is important to apply the network to large images, since otherwise the resolution would be limited by the GPU memory.

① Fully-connected Layer는 2차원 영상 입력 이미지의 크기에 제한이 있다.  
따라서 이미지를 여러 틀에 걸쳐 분할하여  
이미지 전체를无缝하게 overlap-tile 형태로 만들고자 한다.

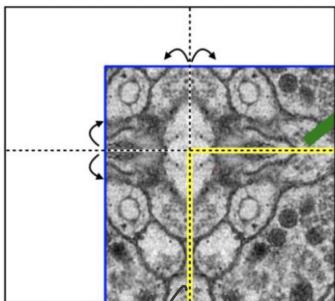
② 각각의 patch는 원본 이미지의 일부로 사용된다.

③ 이미지의 경계부분의 pixel에 대한 segmentation은 흔히 0.14%의  
padding을 적용하는데 이미지가 확장되는 경우에 이를  
extrapolation 기법을 사용한다.

To allow a seamless tiling of the output segmentation map (see Figure 2), it is important to select the input tile size such that all 2x2 max-pooling operations are applied to a layer with an even x- and y-size.



## Mirroring extrapolation



## Mirroring Extrapolation.

- ④ 모니터의 화면 해상도 pixel로 바꾸기

Segmentation은 이미지의 object를 segmentation하는 것  
사용하는 대상이 이미지의 object를 찾는 것  
mirroring을 이용한 extrapolation하는 것  
사용한다.

이미지인데, 이 표현이 이미지에 대해서  
기억에 반영되는 대상으로 확장된다

Another challenge in many cell segmentation tasks is the separation of touching objects of the same class; see Figure 3. To this end, we propose the use of a weighted loss, where the separating background labels between touching cells obtain a large weight in the loss function.

① 디지털 cell segmentation 문제들에서 객체들 사이의 (같은 class)接触하는 혹은 겹쳐있는 부분의 분리가 중요한 상황이다.

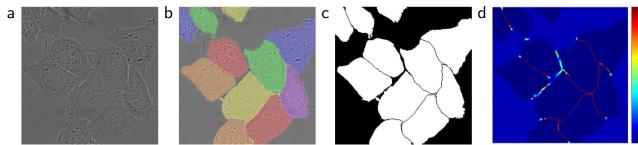


Fig. 3. HeLa cells on glass recorded with DIC (differential interference contrast) microscopy. (a) raw image. (b) overlay with ground truth segmentation. Different colors indicate different instances of the HeLa cells. (c) generated segmentation mask (white: foreground, black: background). (d) map with a pixel-wise loss weight to force the network to learn the border pixels.

따라서 객체간의 분리를 위해  
가중치가 부여된 부분을 적용한다.

즉, 비슷한 높이(label에 대해서)  
(接触하는 cell 사이의 높은  
통에 해당하는 높이)는  
가중치를 주면서 다른 cell은  
인자화될 필요를 한다.

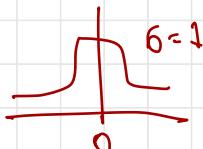
※ momentum = 0.99 (단위)

↳ 높은값은 학습에서 이전경험에서 학습에 사용한 sample이  
현재 optimization step는 학습에서 사용된다.

※ SGD를 사용 (학습적 경사 하강법) 해서 맵을 제작하고 이를 학습에 사용한다.

※ Convolution Layer의 padding을 확장해 output이 더 넓어, 단위

※ Gaussian Distribution을 이용해서 학습과 동시에 초기화



The energy function is computed by a pixel-wise soft-max over the final feature map combined with the cross entropy loss function. The soft-max is defined as  $p_k(\mathbf{x}) = \exp(a_k(\mathbf{x})) / \left( \sum_{k'=1}^K \exp(a_{k'}(\mathbf{x})) \right)$  where  $a_k(\mathbf{x})$  denotes the activation in feature channel  $k$  at the pixel position  $\mathbf{x} \in \Omega$  with  $\Omega \subset \mathbb{Z}^2$ .  $K$  is the number of classes and  $p_k(\mathbf{x})$  is the approximated maximum-function. I.e.  $p_k(\mathbf{x}) \approx 1$  for the  $k$  that has the maximum activation  $a_k(\mathbf{x})$  and  $p_k(\mathbf{x}) \approx 0$  for all other  $k$ . The cross entropy then penalizes at each position the deviation of  $p_{\ell(\mathbf{x})}(\mathbf{x})$  from 1 using

$$E = \sum_{\mathbf{x} \in \Omega} w(\mathbf{x}) \log(p_{\ell(\mathbf{x})}(\mathbf{x})) \quad (1)$$

*이 부분에서 학습 과정에서 크로스 엔트로피(판별기) softmax를 사용한다.*

where  $\ell : \Omega \rightarrow \{1, \dots, K\}$  is the true label of each pixel and  $w : \Omega \rightarrow \mathbb{R}$  is a weight map that we introduced to give some pixels more importance in the training.

We pre-compute the weight map for each ground truth segmentation to compensate the different frequency of pixels from a certain class in the training data set, and to force the network to learn the small separation borders that we introduce between touching cells (See Figure 3c and d).

The separation border is computed using morphological operations. The weight map is then computed as

$$w(\mathbf{x}) = w_c(\mathbf{x}) + w_0 \cdot \exp \left( -\frac{(d_1(\mathbf{x}) + d_2(\mathbf{x}))^2}{2\sigma^2} \right) \quad (2)$$

where  $w_c : \Omega \rightarrow \mathbb{R}$  is the weight map to balance the class frequencies,  $d_1 : \Omega \rightarrow \mathbb{R}$  denotes the distance to the border of the nearest cell and  $d_2 : \Omega \rightarrow \mathbb{R}$  the distance to the border of the second nearest cell. In our experiments we set  $w_0 = 10$  and  $\sigma \approx 5$  pixels.

① 각각의 class마다 다른 가중치를 이용해서 ground-truth에 대한 weightmap을 학습에 이용한다.

