# Fully Convolutional Netwoks for Semantic Segmentation

PR-SMARCLE
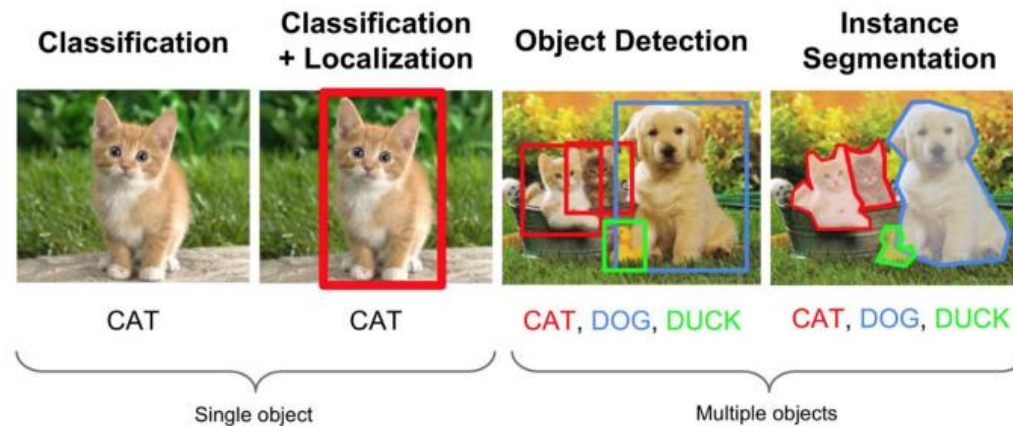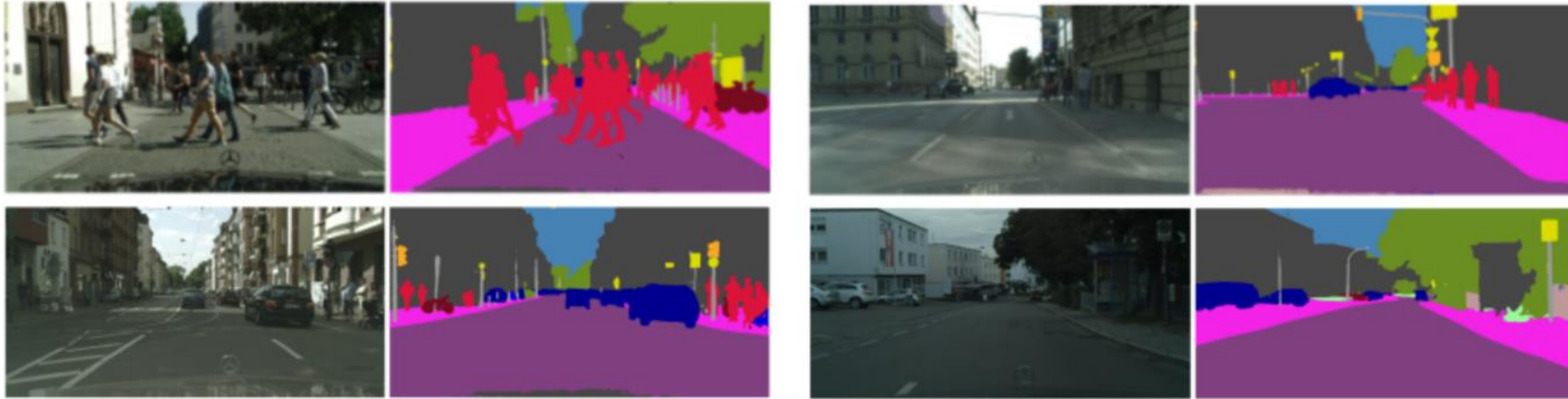
김찬영

SMARCLE

# Content

- Introduction to Semantic Segmentation

- Metric - IoU (Intersection over Union)

- Fully Convolutional Networks

- Encoder (Downsampling)

- Decoder (Upsampling)

- Results

- Key Points

SMARCLE

# Introduction to Semantic Segmentation

• **Classification** (분류): 인풋에 대해서 하나의 레이블을 예측하는 작업.
VGG, ResNet, Inception 등의 모델

• **Object Localization/Detection** (검출): 물체의 레이블을 예측하면서 그 물체가 어디에 있는
지 정보를 제공. 물체가 있는 곳에 네모를 그리는 등으로 표현
YOLO, R-CNN 등의 모델

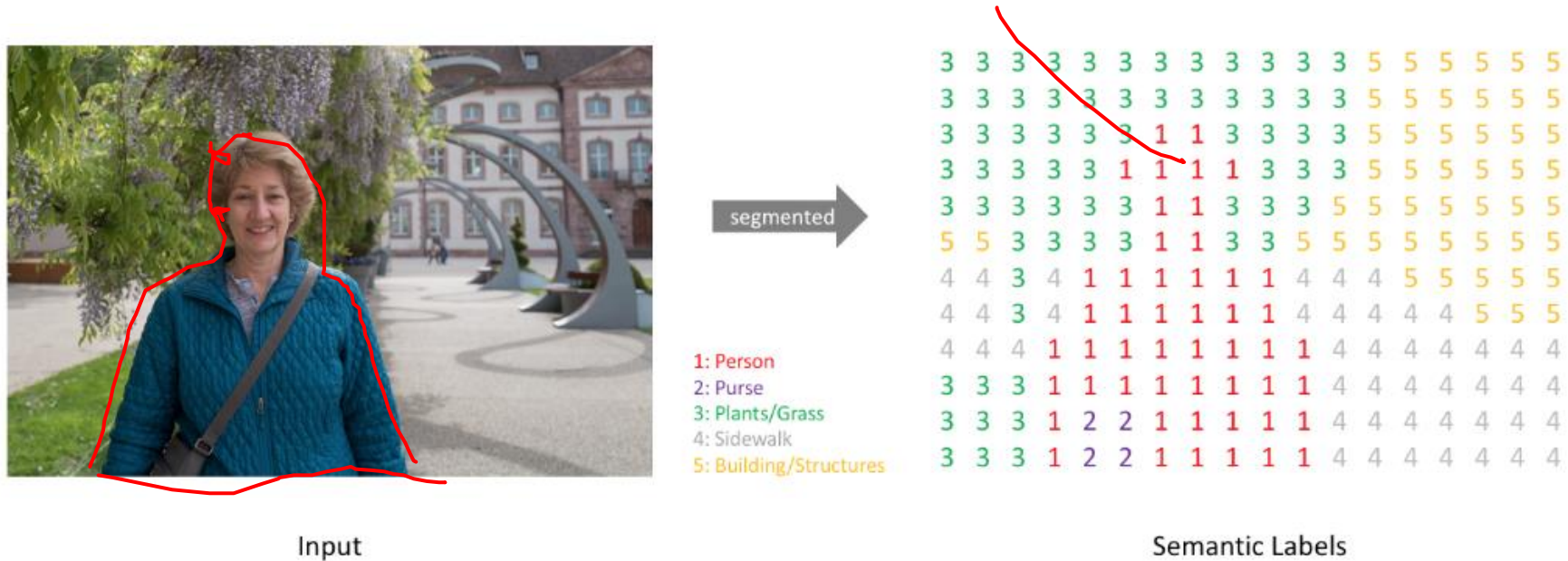• **Segmentation** (분할): 모든 픽셀의 레이블을 예측
FCN, SegNet, DeepLab 등의 모델

# Introduction to Semantic Segmentation



Task Target : 사진 또는 영상의 모든 픽셀을 해당하는 Class로 분류하는 것
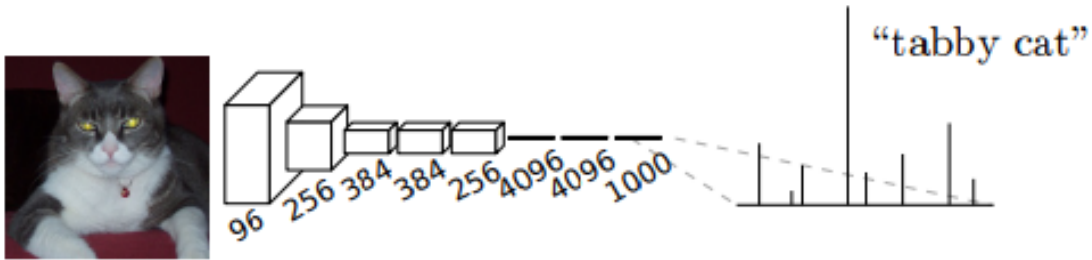
# Introduction to Semantic Segmentation



Input

Semantic Labels

1: Person
2: Purse
3: Plants/Grass
4: Sidewalk
5: Building/Structures

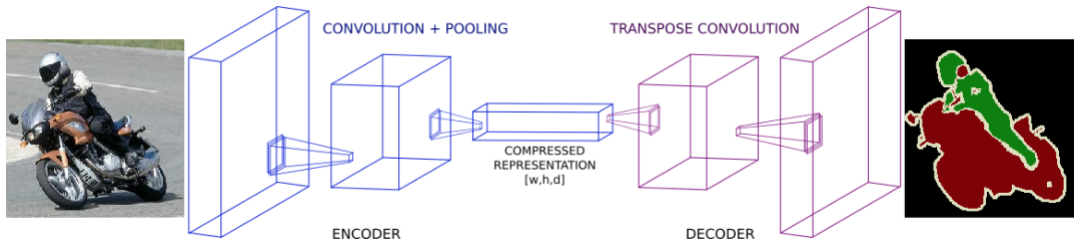Input : (Height X Width X 3) or (Height X Width X 1)

Output : Segmenatation Map

SMARCLE

# Introduction to Semantic Segmentation



일반적인 CNN의 구조
-> 마지막 Fully Connected Layer에 의해 위치정보를 잃음



Semantic Segmentation 모델들은 Encoder Decoder형태
를 사용해 위치 정보를 잃지 않음

대부분의 Semantic Segmentation 방법론은
오늘 리뷰하는 FCN을 기반으로 함

SMARCLE

# Metric - IoU (Intersection over Union)



Score = $\dfrac{\text{Area of overlap}}{\text{Area of union}}$

Sample IoU scores

| 0.905 | 0.532 | 0.391 | 0.143 | 0.0 |

- pixel accuracy: $\sum_i n_{ii} / \sum_i t_i$
- mean accuraccy: $(1/n_{cl}) \sum_i n_{ii}/t_i$
- mean IU: $(1/n_{cl}) \sum_i n_{ii} / \left( t_i + \sum_j n_{ji} - n_{ii} \right)$
- frequency weighted IU:
  $(\sum_k t_k)^{-1} \sum_i t_i n_{ii} / \left( t_i + \sum_j n_{ji} - n_{ii} \right)$

SMARCLE

# Fully Convolutional Networks



Figure 1. Fully convolutional networks can efficiently learn to make dense predictions for per-pixel tasks like semantic segmentation.

FCN's Overall Network Structure

Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2015.

# Fully Convolutional Networks



Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.

Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2015.

# Encoder

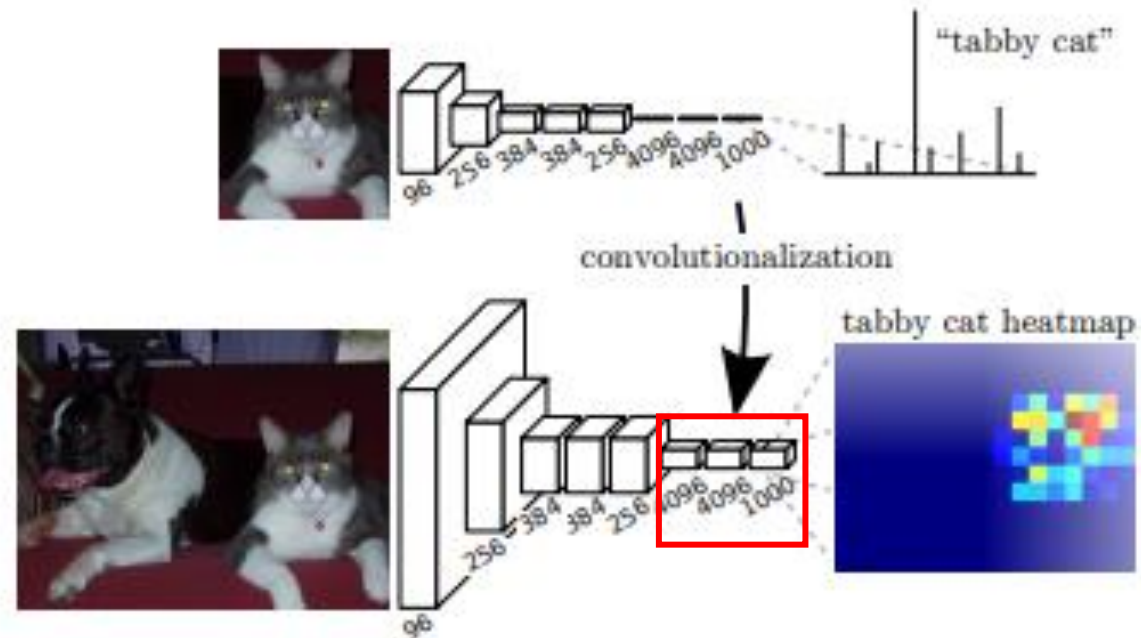Adapting classifiers for dense prediction



Figure 2. Transforming fully connected layers into convolution layers enables a classification net to output a heatmap. Adding layers and a spatial loss (as in Figure 1) produces an efficient machine for end-to-end dense learning.

| | FCN-AlexNet | FCN-VGG16 | FCN-GoogLeNet[4] |
|---|---|---|---|
| mean IU | 39.8 | **56.0** | 42.5 |
| forward time | 50 ms | 210 ms | 59 ms |
| conv. layers | 8 | 16 | 22 |
| parameters | 57M | 134M | 6M |
| rf size | 355 | 404 | 907 |
| max stride | 32 | 32 | 32 |

Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2015.

SMARCLE

# Decoder (Skip Architecture)
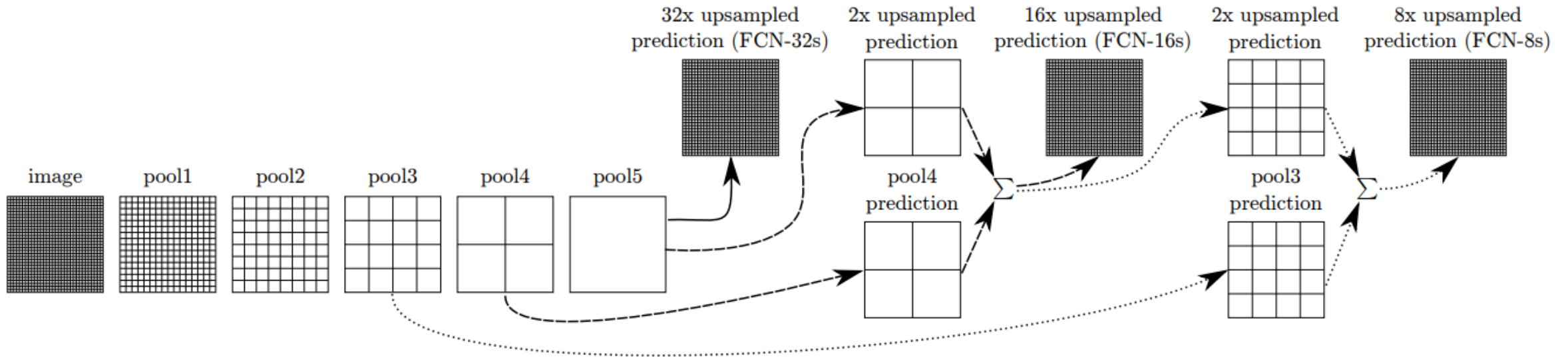


Table 2. Comparison of skip FCNs on a subset of PASCAL VOC2011 validation[7]. Learning is end-to-end, except for FCN-32s-fixed, where only the last layer is fine-tuned. Note that FCN-32s is FCN-VGG16, renamed to highlight stride.

| | pixel acc. | mean acc. | mean IU | f.w. IU |
|---|---|---|---|---|
| FCN-32s-fixed | 83.0 | 59.7 | 45.4 | 72.0 |
| FCN-32s | 89.1 | 73.3 | 59.4 | 81.4 |
| FCN-16s | 90.0 | 75.7 | 62.4 | 83.0 |
| **FCN-8s** | **90.3** | **75.9** | **62.7** | **83.2** |

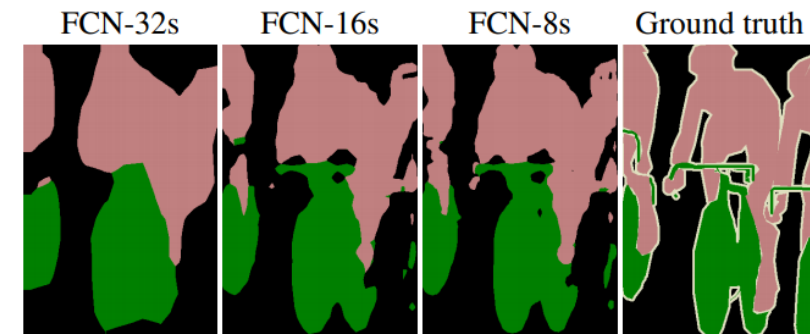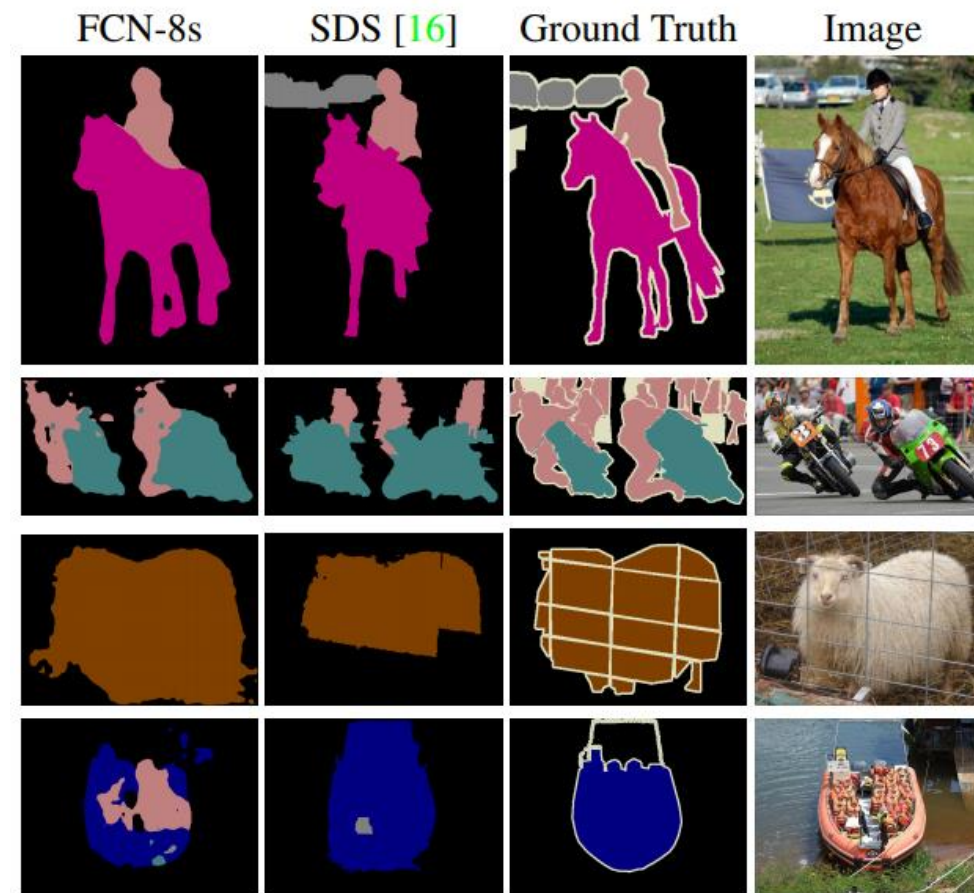Figure 4. Refining fully convolutional nets by fusing information from layers with different strides improves segmentation detail. The first three images show the output from our 32, 16, and 8 pixel stride nets (see Figure 3).

Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015.

# Results

## PASCAL VOC

Table 3. Our fully convolutional net gives a 20% relative improvement over the state-of-the-art on the PASCAL VOC 2011 and 2012 test sets, and reduces inference time.

|  | mean IU VOC2011 test | mean IU VOC2012 test | inference time |
|---|---|---|---|
| R-CNN [12] | 47.9 | - | - |
| SDS [16] | 52.6 | 51.6 | $\sim 50$ s |
| FCN-8s | **62.7** | **62.2** | $\sim$ **175 ms** |

Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015.

# Results

## NYUDv2

Table 4. Results on NYUDv2. *RGBD* is early-fusion of the RGB and depth channels at the input. *HHA* is the depth embedding of [14] as horizontal disparity, height above ground, and the angle of the local surface normal with the inferred gravity direction. *RGB-HHA* is the jointly trained late fusion model that sums RGB and HHA predictions.

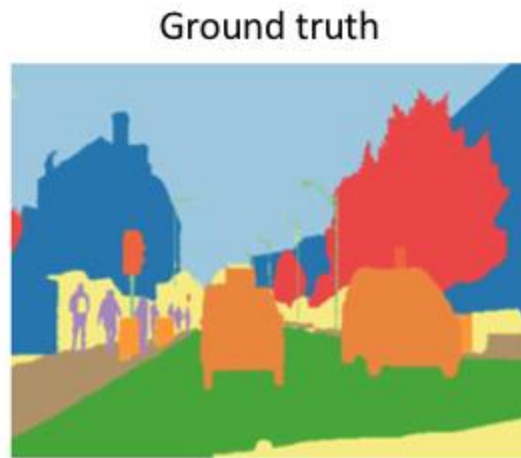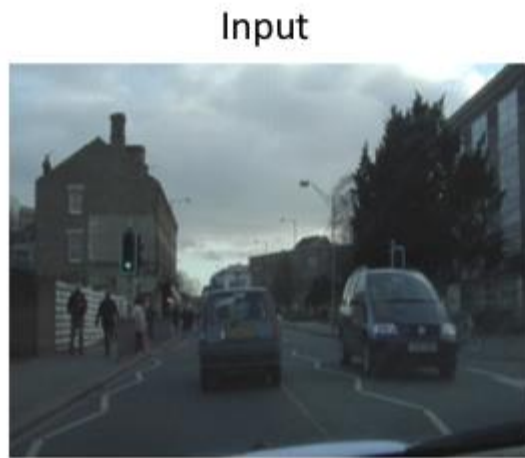|  | pixel acc. | mean acc. | mean IU | f.w. IU |
|---|---|---|---|---|
| Gupta *et al.* [14] | 60.3 | - | 28.6 | 47.0 |
| FCN-32s RGB | 60.0 | 42.2 | 29.2 | 43.9 |
| FCN-32s RGBD | 61.5 | 42.4 | 30.5 | 45.5 |
| FCN-32s HHA | 57.1 | 35.2 | 24.2 | 40.4 |
| FCN-32s RGB-HHA | 64.3 | 44.9 | 32.8 | 48.0 |
| FCN-16s RGB-HHA | **65.4** | **46.1** | **34.0** | **49.5** |

## SIFT Flow

Table 5. Results on SIFT Flow[10] with class segmentation (center) and geometric segmentation (right). Tighe [33] is a non-parametric transfer method. Tighe 1 is an exemplar SVM while 2 is SVM + MRF. Farabet is a multi-scale convnet trained on class-balanced samples (1) or natural frequency samples (2). Pinheiro is a multi-scale, recurrent convnet, denoted $RCNN_3$ ($\circ^3$). The metric for geometry is pixel accuracy.

|  | pixel acc. | mean acc. | mean IU | f.w. IU | geom. acc. |
|---|---|---|---|---|---|
| Liu *et al.* [23] | 76.7 | - | - | - | - |
| Tighe *et al.* [33] | - | - | - | - | 90.8 |
| Tighe *et al.* [34] 1 | 75.6 | 41.1 | - | - | - |
| Tighe *et al.* [34] 2 | 78.6 | 39.2 | - | - | - |
| Farabet *et al.* [8] 1 | 72.3 | 50.8 | - | - | - |
| Farabet *et al.* [8] 2 | 78.5 | 29.6 | - | - | - |
| Pinheiro *et al.* [28] | 77.7 | 29.8 | - | - | - |
| FCN-16s | **85.2** | **51.7** | 39.5 | 76.1 | **94.3** |

Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2015.

SMARCLE

# Key Points

1. Adopting Fully Convolutional Layer instead of Fully Connected Layer
   -> to maintain corresponding spatial dimension of input data

2. Defined Novel Skip Architecture
   -> to combine deep, coarse, semantic information and shallow, fine, appearance information



Input      Ground truth      Prediction

Long, Jonathan, Evan Shelhamer, and Trevor Darrell. "Fully convolutional networks for semantic segmentation." *Proceedings of the IEEE conference on computer vision and pattern recognition.* 2015.

SMARCLE

# 감사합니다

PR-SMARCLE

SMARCLE

김찬영