# SegNet

Presenter: Khakim Akhunov
22.10.2019, NTNU

# SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation

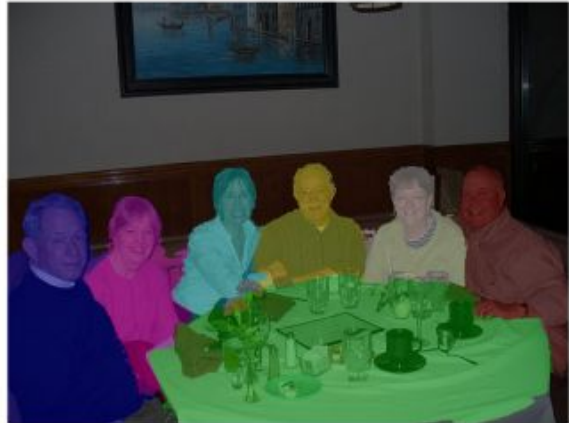Vijay Badrinarayanan, Alex Kendall, Roberto Cipolla, Senior Member, IEEE

10 Oct 2016

# What Semantic Segmentation is?

The process of assigning a label to every pixel in the image.

Semantic segmentation treats multiple objects of the same class as a single entity.

# Example of various Scene Understanding tasks



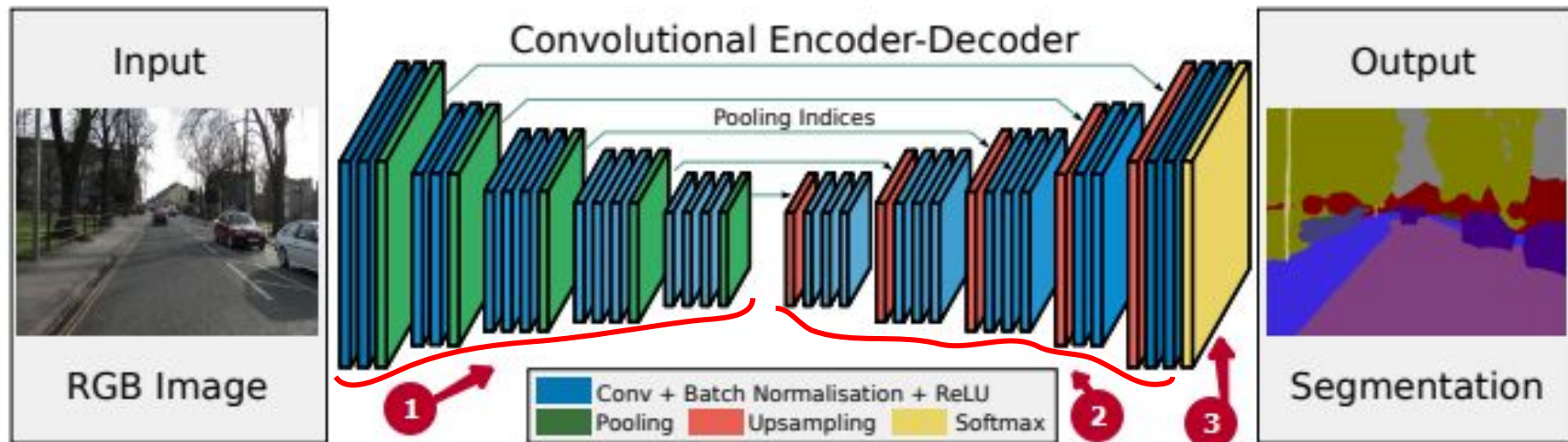| Object Detection | Semantic Segmentation | Instance Segmentation |
|---|---|---|
| Tags: Person, Dining Table | A group of people sitting at a table | Q: What were the people doing? A: Eating dinner |
| Image Classification | Image Captioning | Visual Question-Answering |

# What SegNet is?

Novel and practical deep fully convolutional neural network architecture for semantic pixel-wise segmentation.

SegNet = (Encoder + Decoder) + Pixel-Wise Classification layer

It is primarily motivated by road scene understanding applications which require the ability to model appearance (road, building), shape (cars, pedestrians) and understand the spatial-relationship (context) between different classes such as road and side-walk.

# SegNet architecture



1 - encoder network
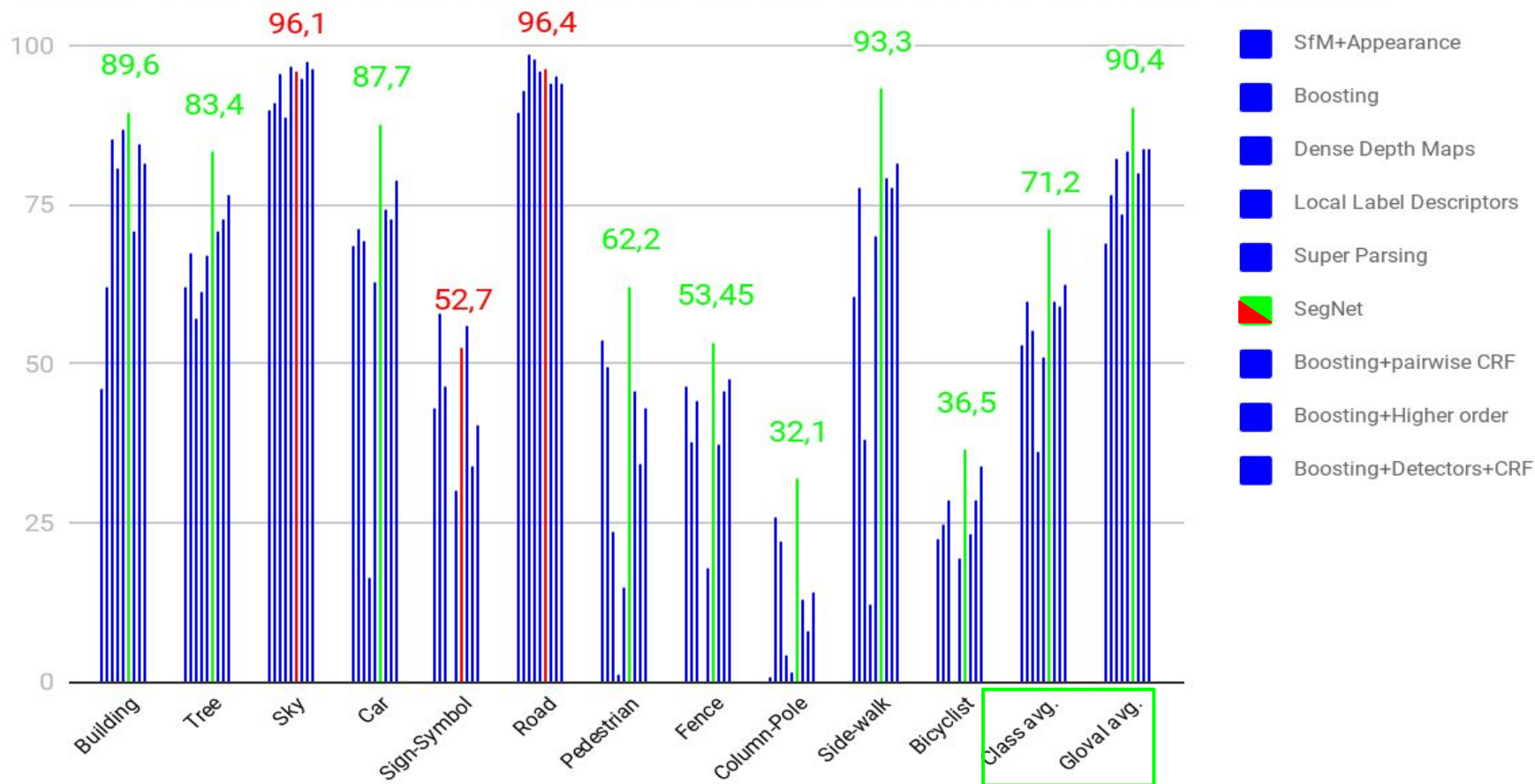2 - decoder network
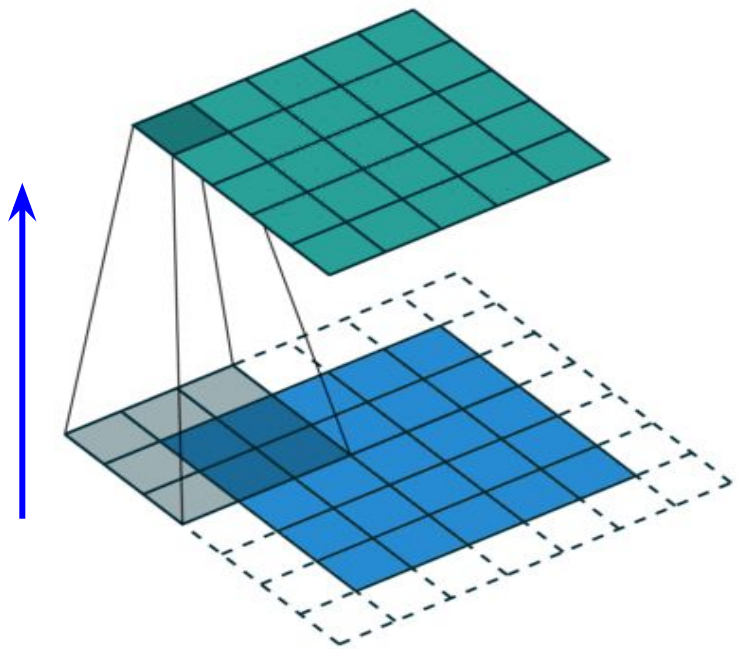3 - pixel-wise classification layer

# Main motivations behind SegNet

- **Retain boundary information** in the extracted image representation
- **Efficient** in terms of both **memory and computation time**
- **Able to train** end-to-end using **efficient** weight update technique

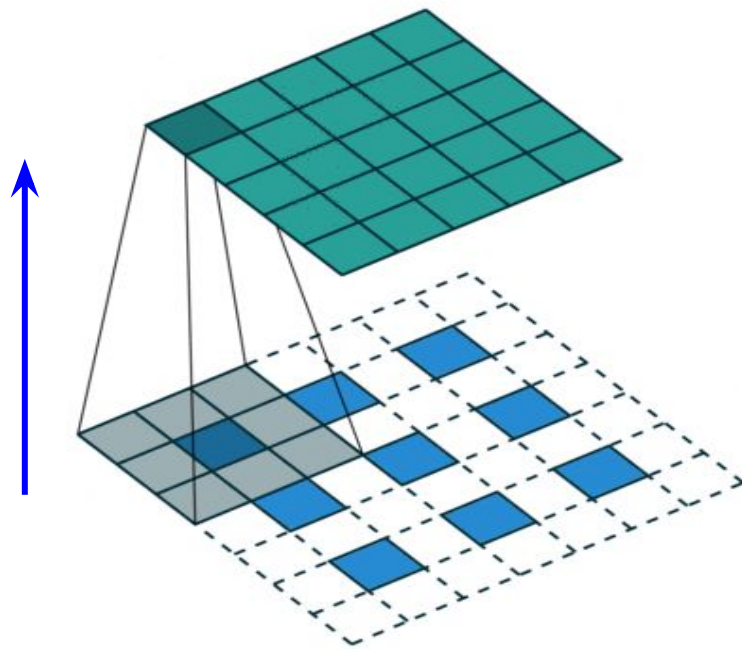# SegNet predictions on road scenes and indoor scenes

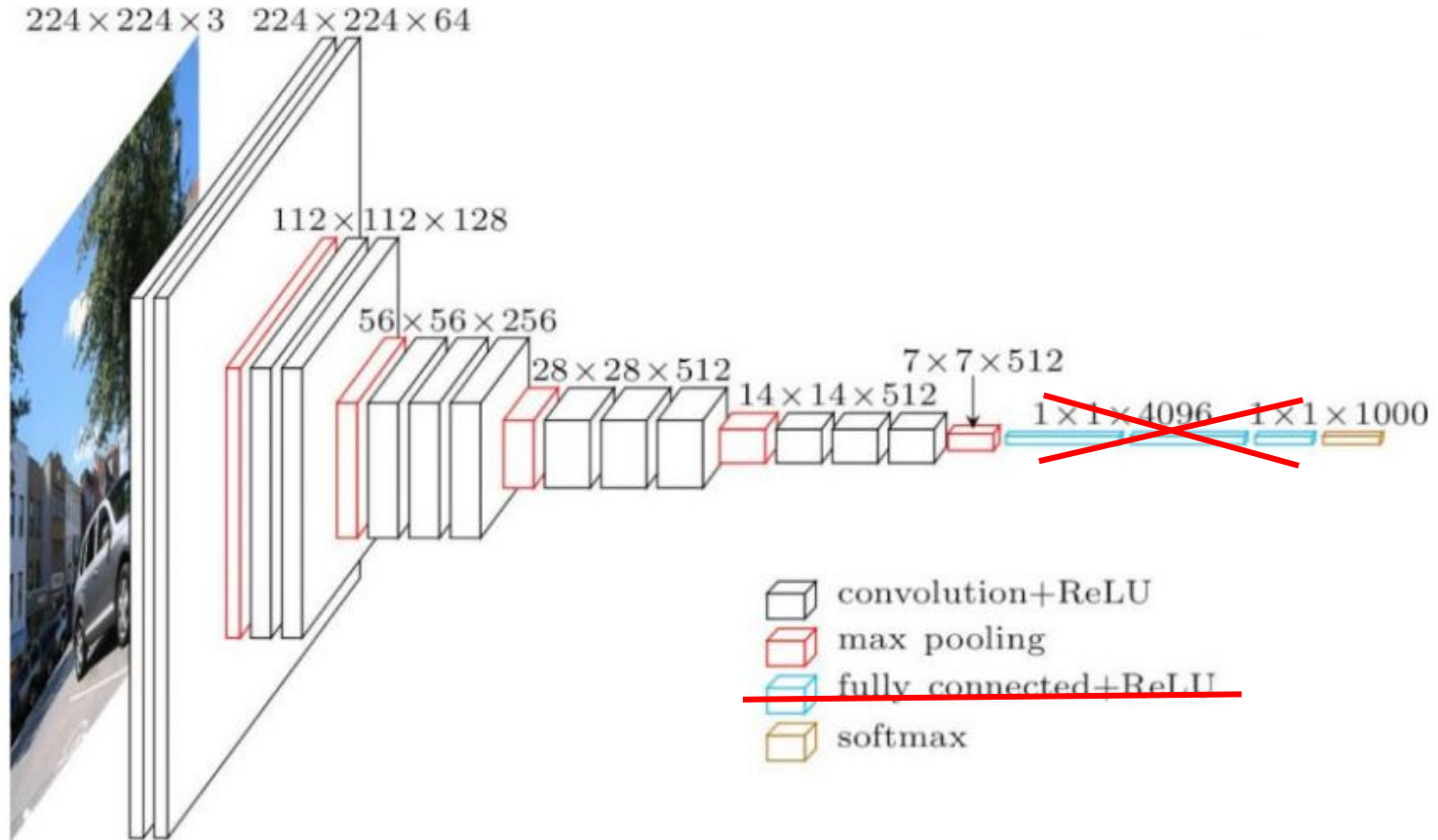Quantitative comparisons of SegNet with traditional methods on the CamVid 11 road class segmentation problem
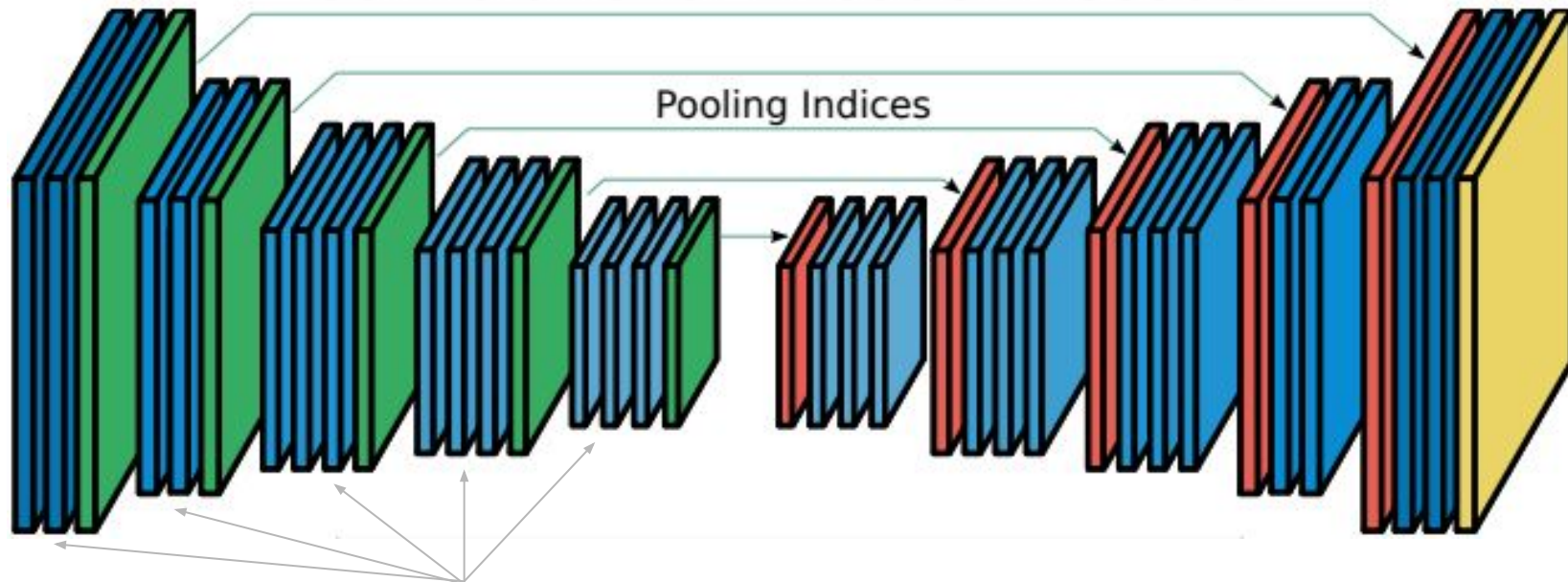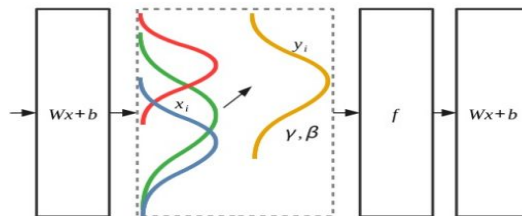
Encoding

Decoding

# VGG16 architecture



$224 \times 224 \times 3$  $224 \times 224 \times 64$

$112 \times 112 \times 128$

$56 \times 56 \times 256$

$28 \times 28 \times 512$

$14 \times 14 \times 512$

$7 \times 7 \times 512$

$1 \times 1 \times 4096$  $1 \times 1 \times 1000$

convolution+ReLU

max pooling

fully connected+ReLU

softmax

# Convolutional Encoder-Decoder

Pooling Indices

3x3conv

$1{\times}1$ $1{\times}0$ $1{\times}1$ 0 0
$0{\times}0$ $1{\times}1$ $1{\times}0$ 1 0
$0{\times}1$ $0{\times}0$ $1{\times}1$ 1 1
0 0 1 1 0
0 1 1 0 0

4

$wx+b$ → $x_i$ → $y_i$, $\gamma, \beta$ → $f$ → $wx+b$

Batch Normalization

ReLU

$R(z) = max(0, \ z)$

ReLU

# Convolutional Encoder-Decoder

Pooling Indices

| 1 | 3 | 2 | 9 |
|---|---|---|---|
| 7 | 4 | 1 | 5 |
| 8 | 5 | 2 | 3 |
| 4 | 2 | 1 | 4 |

| 7 | 9 |
|---|---|
| 8 | |

## 2x2 Max-pooling

# Convolutional Encoder-Decoder

Pooling Indices

| $a$ | 0 | 0 | 0 |
|---|---|---|---|
| 0 | 0 | $b$ | 0 |
| 0 | 0 | 0 | $d$ |
| $c$ | 0 | 0 | 0 |

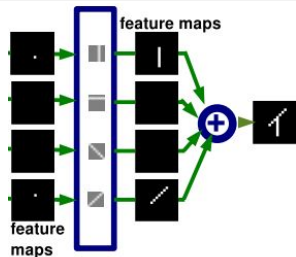| $a$ | $b$ |
|---|---|
| $c$ | $d$ |

Max-pooling Indices

Upsampling

# Using pooling indices for upsampling

# Convolutional Encoder-Decoder

Pooling Indices



Decoder Filter Bank

feature maps

feature maps

Batch Normalization

$Wx+b$

$x_i$

$y_i$

$\gamma, \beta$

$f$

$Wx+b$

# Convolutional Encoder-Decoder



Pooling Indices

Softmax classifier
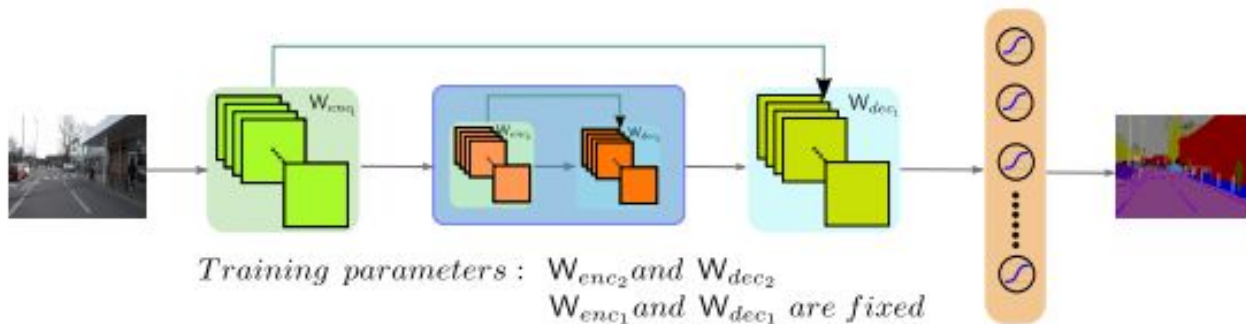
$$y_k = \frac{\exp(a_k)}{\sum_i \exp(a_i)}$$

# Training the SegNet



Training parameters : $W_{enc_1}$ and $W_{dec_1}$

(a)

$$y_k = \frac{\exp(a_k)}{\sum_i \exp(a_i)}$$

Training parameters : $W_{enc_2}$ and $W_{dec_2}$
$W_{enc_1}$ and $W_{dec_1}$ are fixed
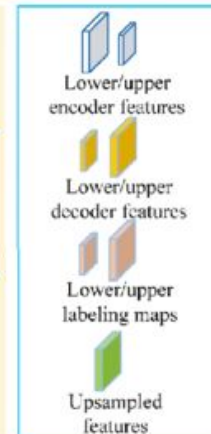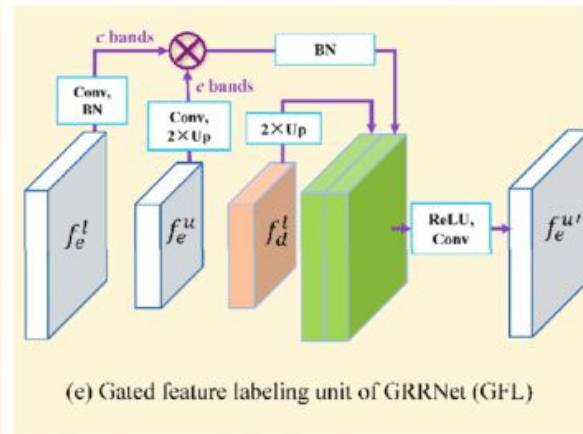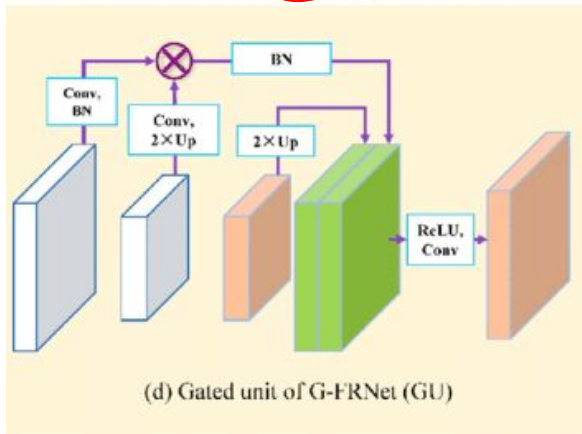
(b)
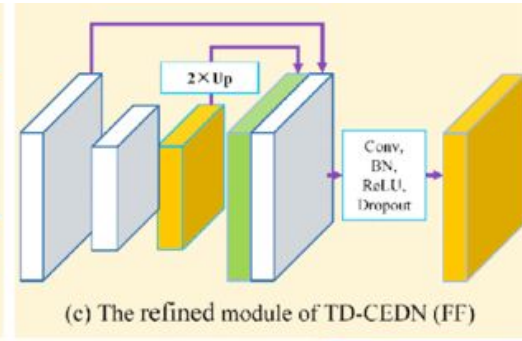
$$y_k = \frac{\exp(a_k)}{\sum_i \exp(a_i)}$$

# Comparison of decoder variants

| Variant | Params (M) | Storage multiplier | Inference time (ms) |
|---|---|---|---|
| Fixed upsampling | | | |
| Bilinear-Interpolation | 0.625 | 0 | 24.2 |
| Upsampling using max-pooling indices | | | |
| SegNeg-Basic | 1.425 | 1 | 52.6 |
| SegNeg-Basic-Encoder Addition | 1.425 | 64 | 53.0 |
| SegNeg-Basic-SingleChannelDecoder | 0.625 | 1 | 33.1 |
| Learning to upsample | | | |
| FCN-Basic | 0.65 | 11 | 24.2 |
| FCN-Basic-NoAddition | 0.65 | n/a | 23.8 |
| FCN-Basic-NoDimReduction | 1.625 | 64 | 44.8 |
| FCN-Basic-NoAddition-NoDimReduction | 1.625 | 0 | 43.9 |

# Summary of different decoders analysis

- The best performance is achieved when encoder feature maps are stored in full
- Compressed forms of encoder feature maps can be stored and used for decoding to meet memory constraints
- Larger decoders increase performance for a given encoder netrowork

# Schematic representation of different architectures



(a) The refined module of SegNet

(b) The refined module of Res-U-Net

(c) The refined module of TD-CEDN (FF)

(d) Gated unit of G-FRNet (GU)

(e) Gated feature labeling unit of GRRNet (GFL)

Lower/upper encoder features

Lower/upper decoder features

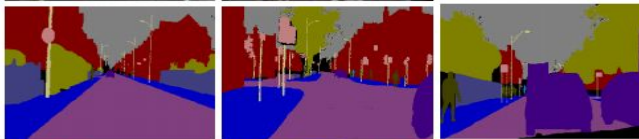Lower/upper labeling maps

Upsampled features

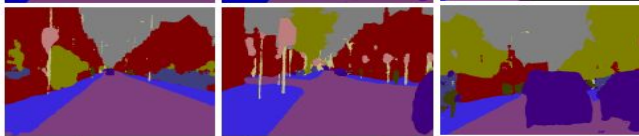# Road scene segmentation

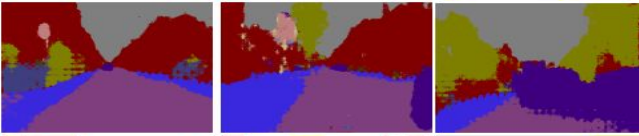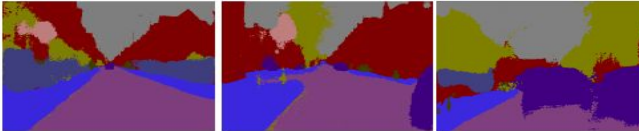

Test samples

Ground Truth

SegNet

DeepLab-LargeFOV

FCN

DeconvNet

# Indoor scene segmentation



Test samples

Ground Truth

SegNet

DeepLab-LargeFOV

FCN (learnt deconv)

DeconvNet

# Conclusion

- **SegNet is more efficient** compared to other architectures since it **only stores the max-pooling indices** of the feature maps and uses them in its decoder network **to achieve good performance**

- On large and well known datasets **SegNet performs competitively**, achieving high scores for road scene understanding

- **End-to-end learning** of deep segmentation architectures is a **harder challenge**

# References

1) http://mi.eng.cam.ac.uk/projects/segnet/
2) https://www.cyberailab.com/home/segnet-an-image-segmentation-neural-network
3) http://www.robots.ox.ac.uk/~tvg/publications/2017/CRFMeetCNN4SemanticSegmentation.pdf
4) https://neurohive.io/en/popular-networks/vgg16/

# Questions?