# E9 205 Machine Learning for Signal Procesing
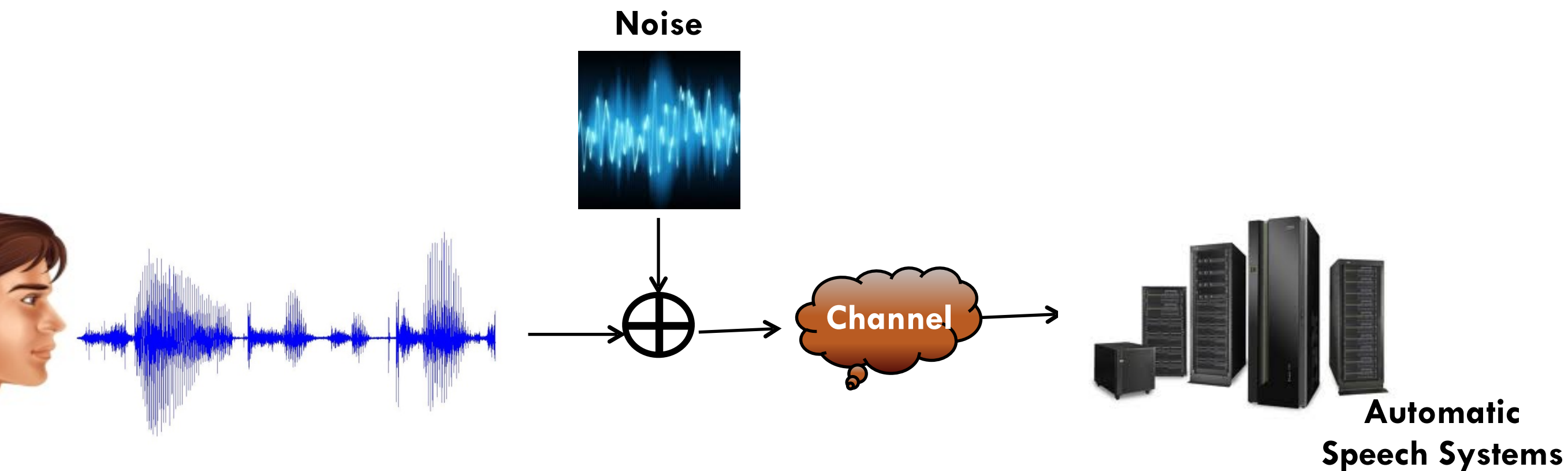
**Deep Learning for Audio and Vision**

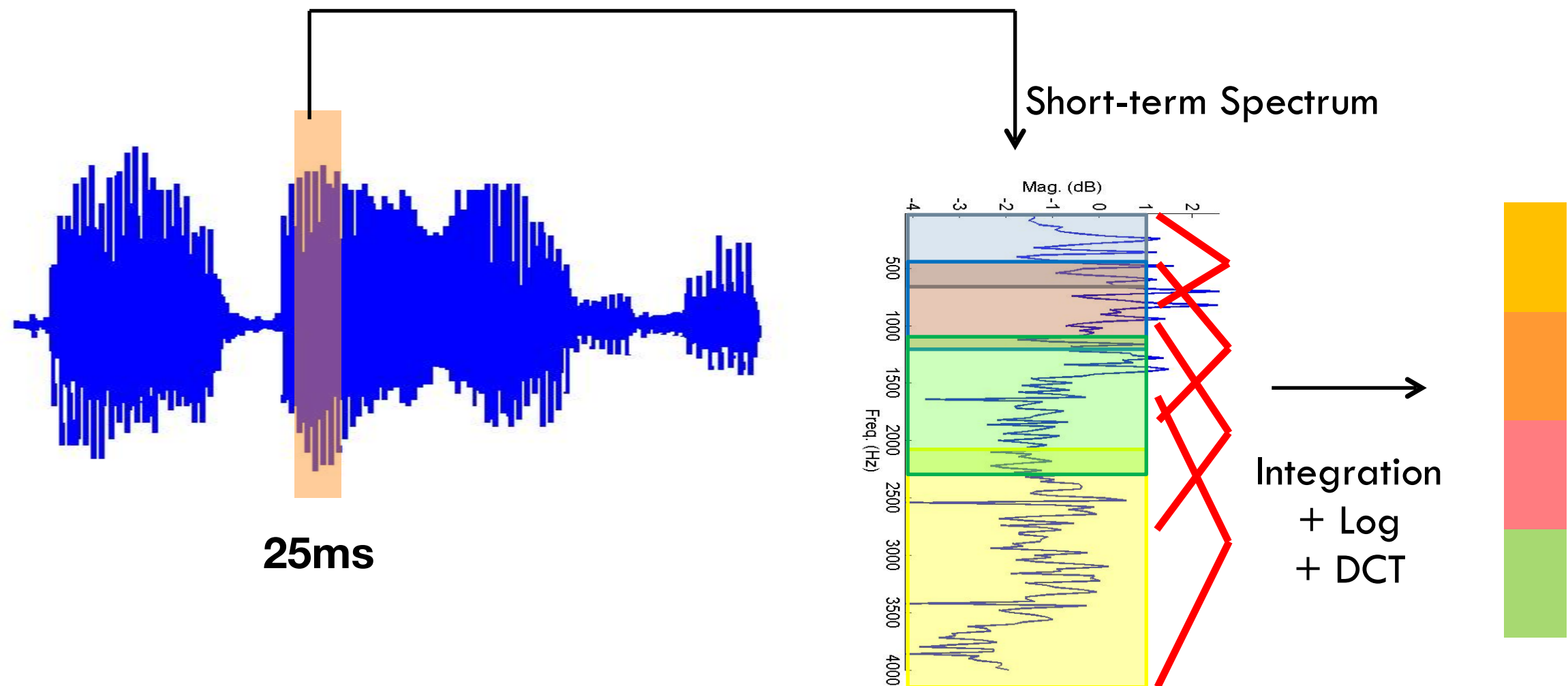20-11-2019

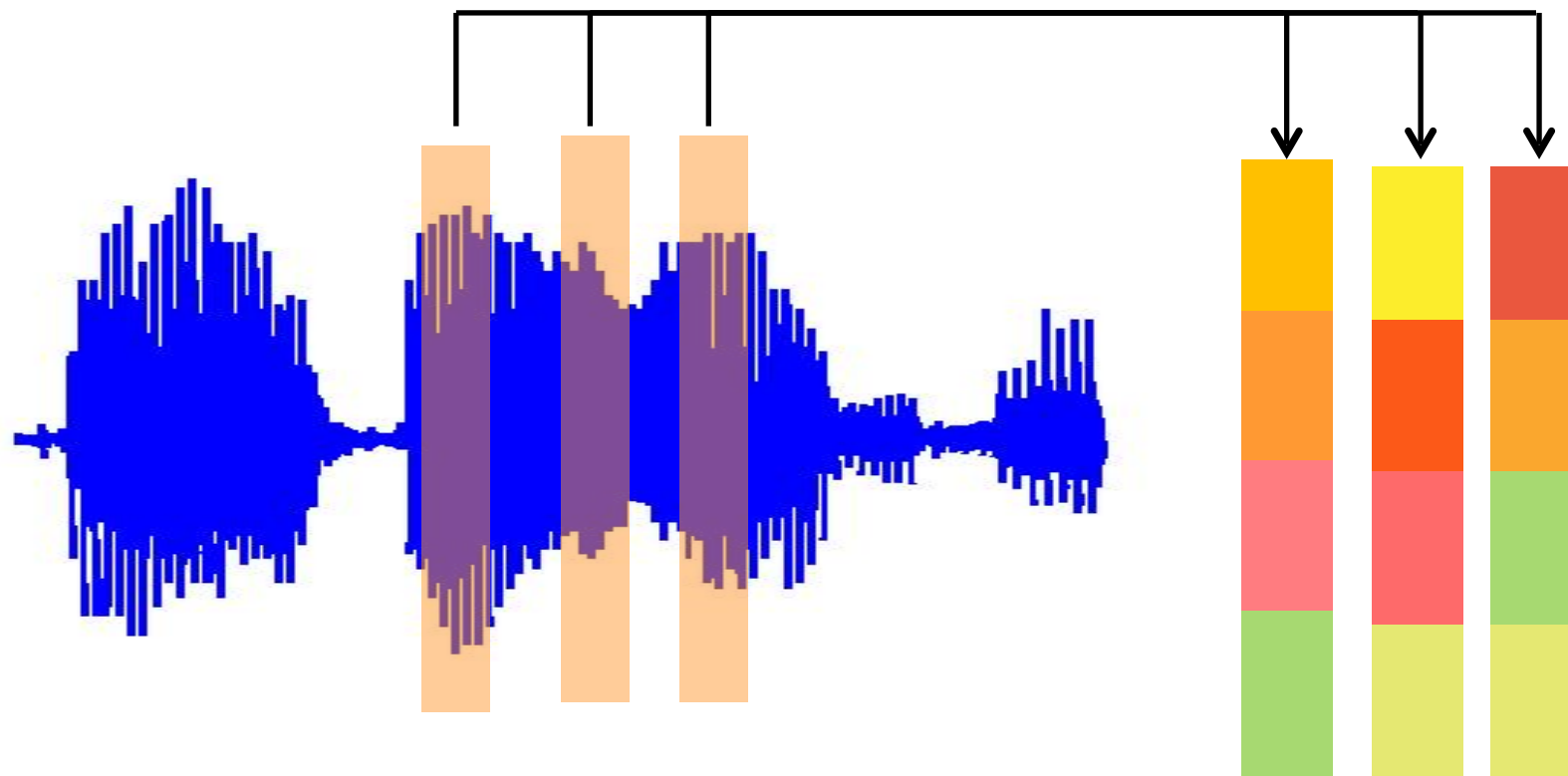# Speech Recognition



Noise

Channel

Automatic Speech Systems

# Signal Modeling

- Short-term spectra integrated in mel frequency bands followed by log compression + DCT – mel frequency cepstral coefficients (MFCC) [Davis and Mermelstein, 1979].
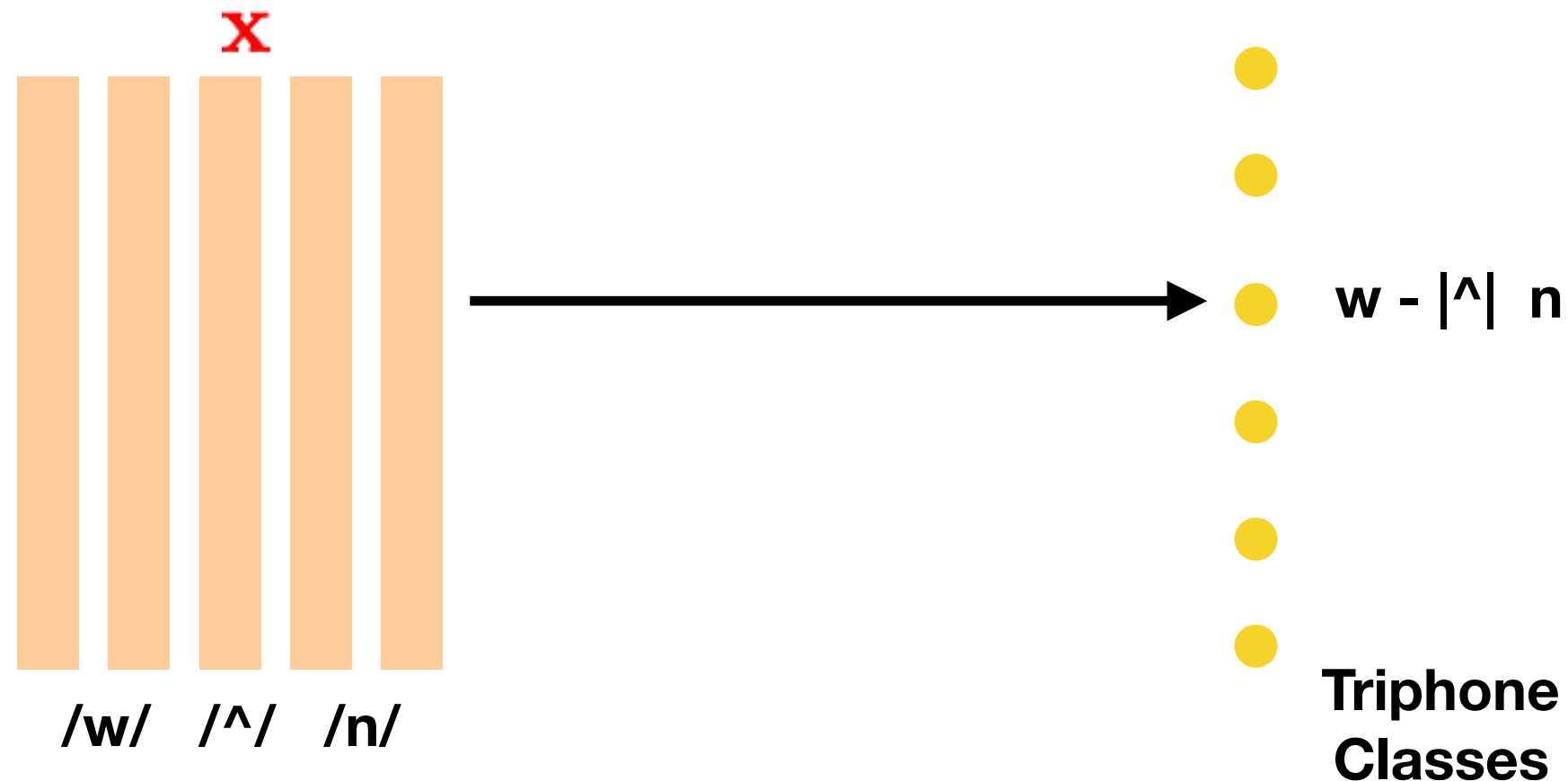
# Mel Frequency Cepstral Coefficients

- MFCC processing repeated for every short-term frame yielding a sequence of features. Typically 25ms frames with 10ms hop in time.
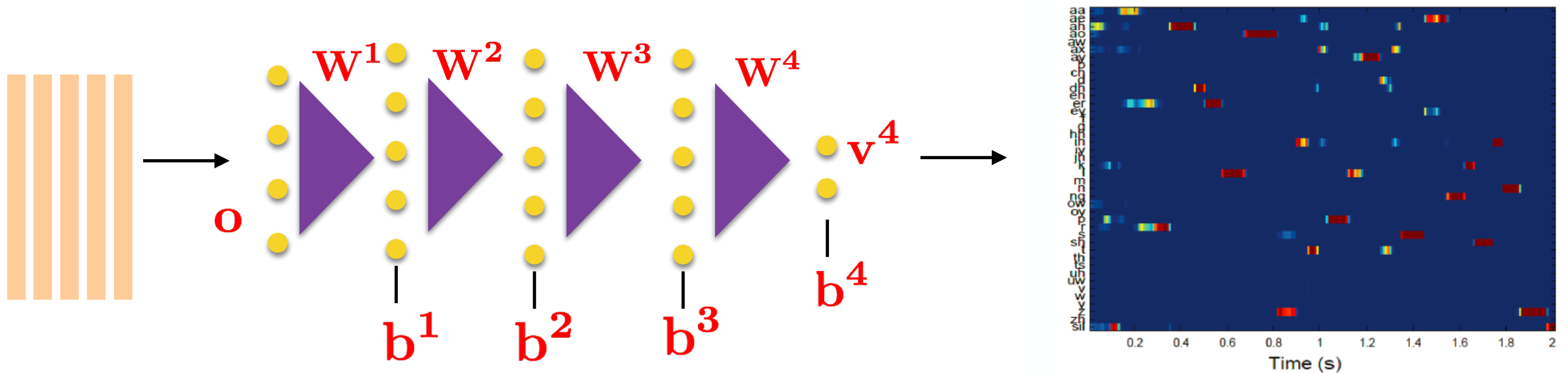
# Speech Recognition

- Map the features to phone class. Using phone labelled data.

**x**

/w/  /^/  /n/

w - |^|  n

**Triphone Classes**

- Classical machine learning - train a classifier on speech training data that maps to the target phoneme class.

LEAP

# Back to Speech Recognition

# Back to Speech Recognition

Mapping Speech Features to Phonemes to words
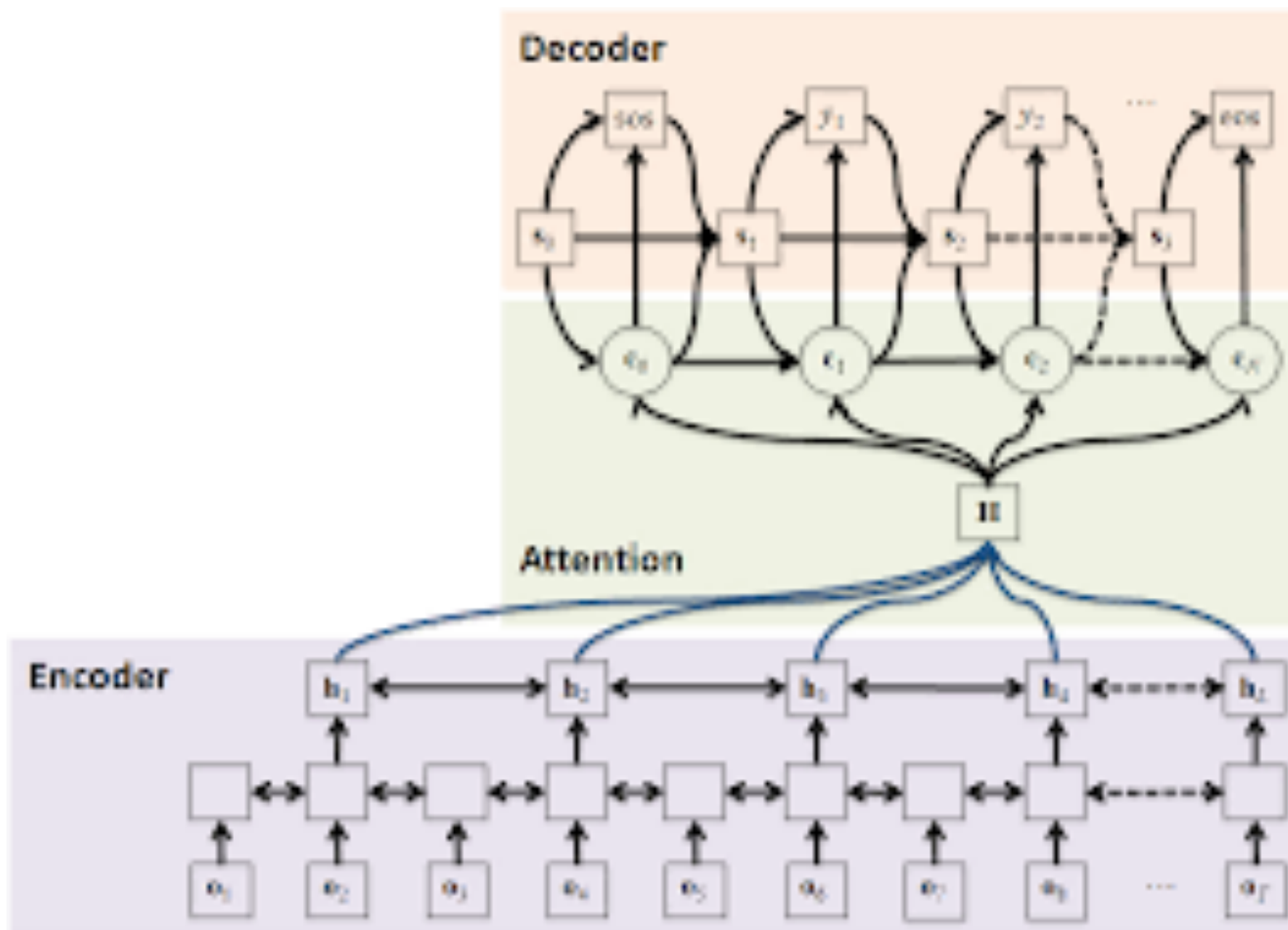
# State of Progress



**Claims of human parity using BLSTM based Models !!!**

# Moving to End-to-End


Text Output / Audio Features

# Image Processing

# Visual Graphics Group Network



224 x 224 x 3    224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

7 x 7 x 512

14 x 14 x 512

1 x 1 x 4096   1 x 1 x 1000

convolution+ReLU
max pooling
fully nected+ReLU
softmax

# ImageNet Task

1000 images in each of 1000 categories. In all, there are roughly 1.2 million training images, 50,000 validation images, and 150,000 testing images. ImageNet consists of variable-resolution images. Therefore, the images have been down-sampled to a fixed resolution of 224×224.

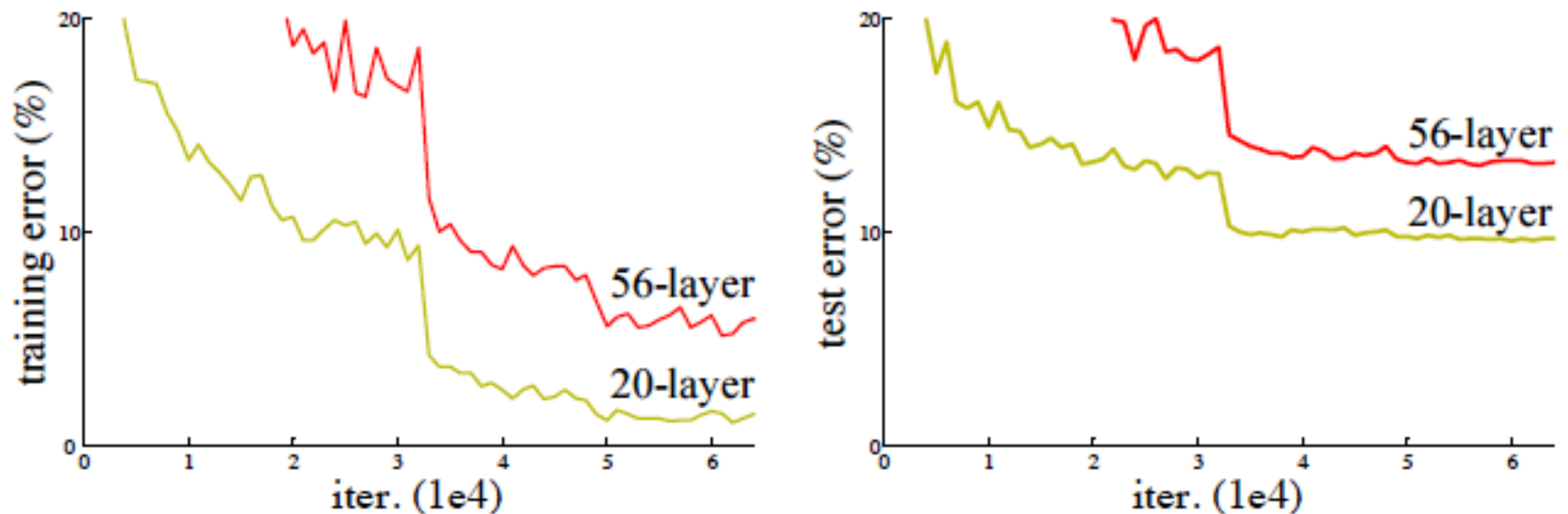| Method | top-1 val. error (%) | top-5 val. error (%) | top-5 test error (%) |
|---|---|---|---|
| VGG (2 nets, multi-crop & dense eval.) | **23.7** | **6.8** | **6.8** |
| VGG (1 net, multi-crop & dense eval.) | 24.4 | 7.1 | 7.0 |
| VGG (ILSVRC submission, 7 nets, dense eval.) | 24.7 | 7.5 | 7.3 |
| GoogLeNet (Szegedy et al., 2014) (1 net) | - | 7.9 | |
| GoogLeNet (Szegedy et al., 2014) (7 nets) | - | 6.7 | |
| MSRA (He et al., 2014) (11 nets) | - | - | 8.1 |
| MSRA (He et al., 2014) (1 net) | 27.9 | 9.1 | 9.1 |
| Clarifai (Russakovsky et al., 2014) (multiple nets) | - | - | 11.7 |
| Clarifai (Russakovsky et al., 2014) (1 net) | - | - | 12.5 |
| Zeiler & Fergus (Zeiler & Fergus, 2013) (6 nets) | 36.0 | 14.7 | 14.8 |
| Zeiler & Fergus (Zeiler & Fergus, 2013) (1 net) | 37.5 | 16.0 | 16.1 |
| OverFeat (Sermanet et al., 2014) (7 nets) | 34.0 | 13.2 | 13.6 |
| OverFeat (Sermanet et al., 2014) (1 net) | 35.7 | 14.2 | - |
| Krizhevsky et al. (Krizhevsky et al., 2012) (5 nets) | 38.1 | 16.4 | 16.4 |
| Krizhevsky et al. (Krizhevsky et al., 2012) (1 net) | 40.7 | 18.2 | - |

# Can we go deeper



Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer "plain" networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

# Deep Residual Learning for Image Recognition

Kaiming He      Xiangyu Zhang      Shaoqing Ren      Jian Sun

Microsoft Research

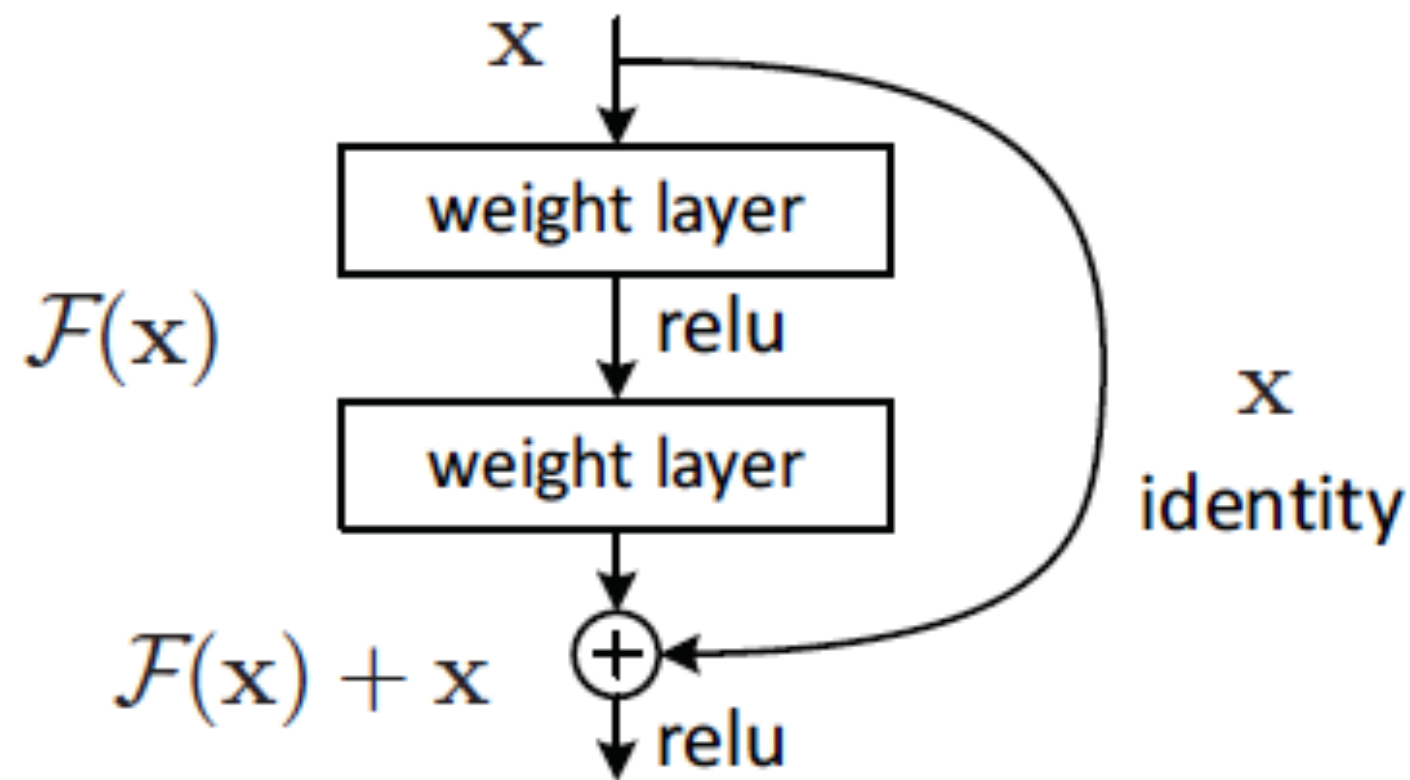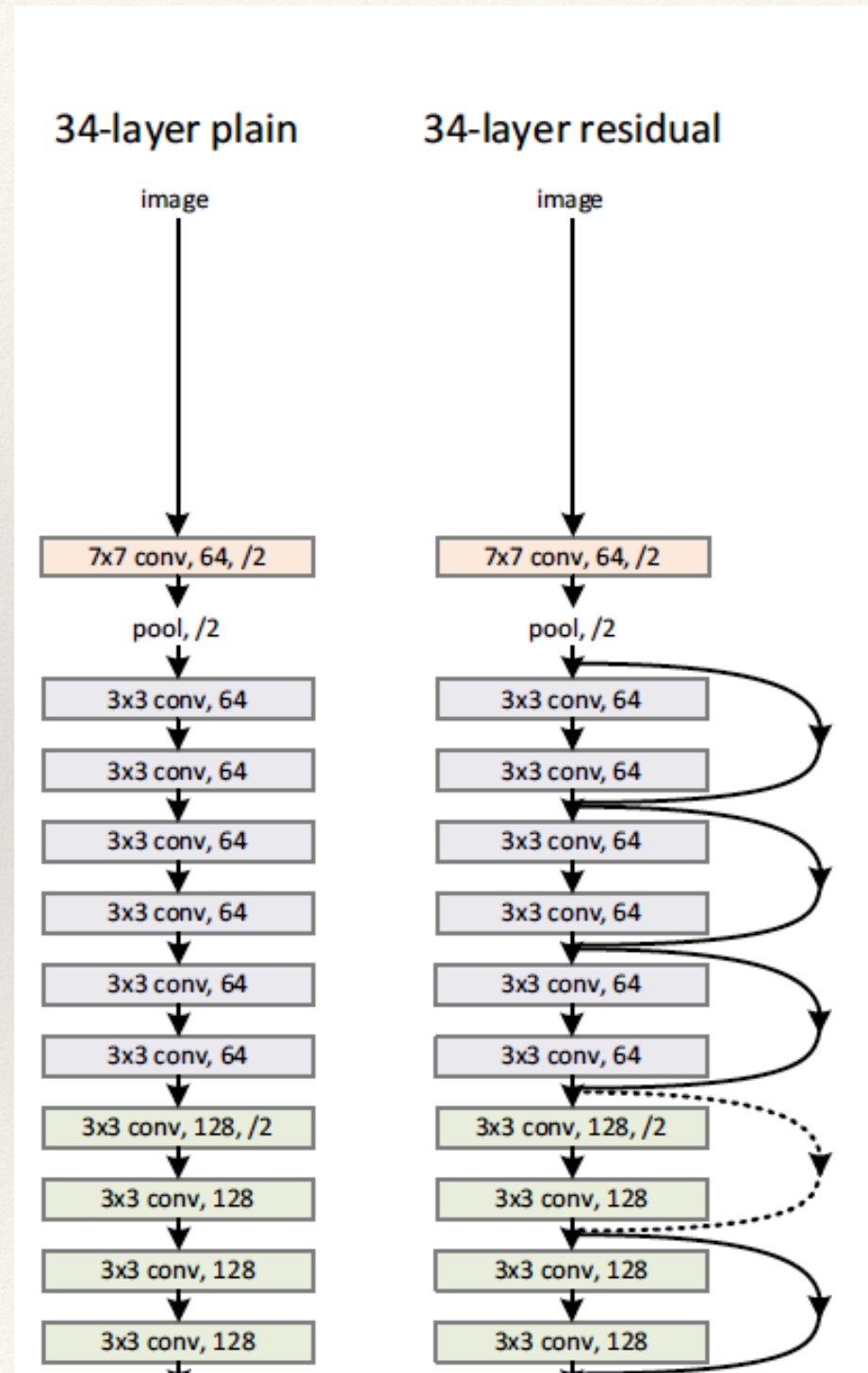{kahe, v-xiangz, v-shren, jiansun}@microsoft.com

# Residual Blocks



Figure 2. Residual learning: a building block.

# Deep Networks with Residual Blocks
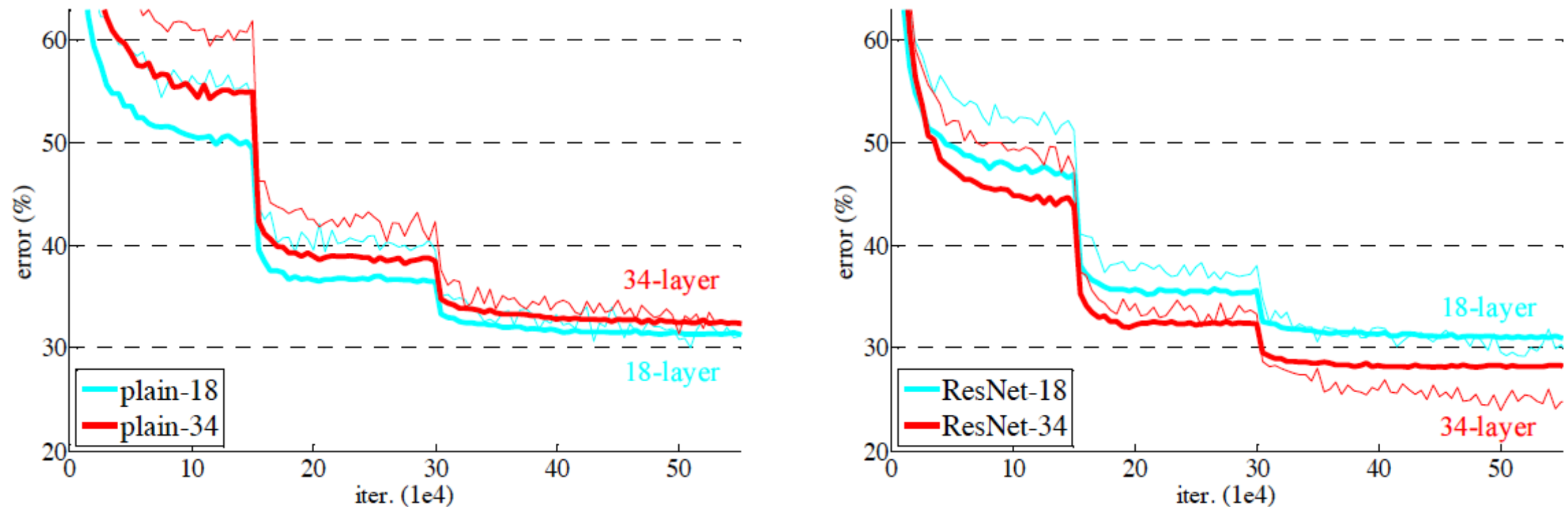
# Deep Networks with Residual Blocks



Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

# Results with ResNet

|  | plain | ResNet |
|---|---|---|
| 18 layers | 27.94 | 27.88 |
| 34 layers | 28.54 | **25.03** |

Table 2. Top-1 error (%, 10-crop testing) on ImageNet validation. Here the ResNets have no extra parameter compared to their plain counterparts. Fig. 4 shows the training procedures.
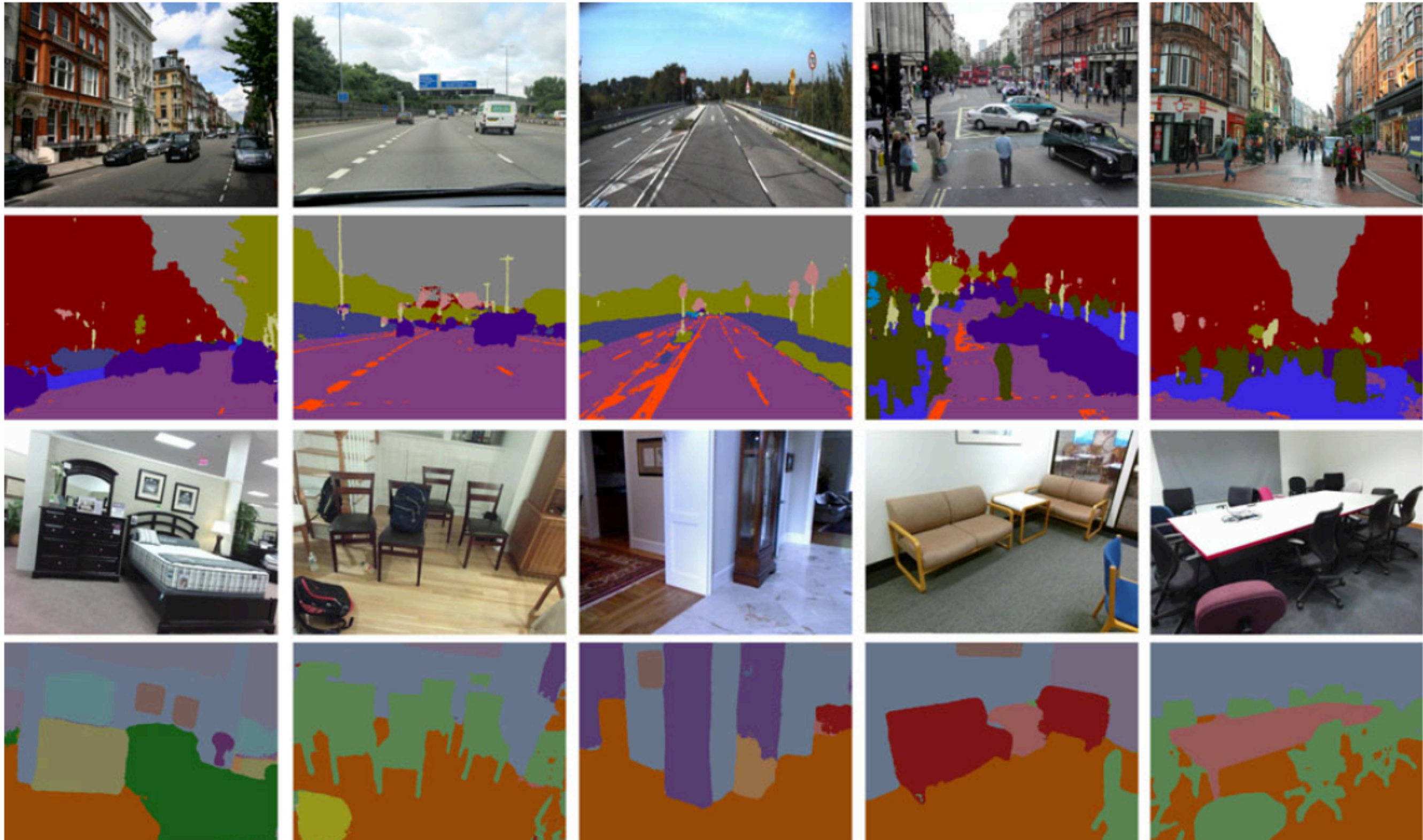
| ResNet-34 C | 24.19 | 7.40 |
|---|---|---|
| ResNet-50 | 22.85 | 6.71 |
| ResNet-101 | 21.75 | 6.05 |
| ResNet-152 | **21.43** | **5.71** |

# Image Segmentation

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

Vijay Badrinarayanan, Alex Kendall [iD], and Roberto Cipolla, *Senior Member, IEEE*
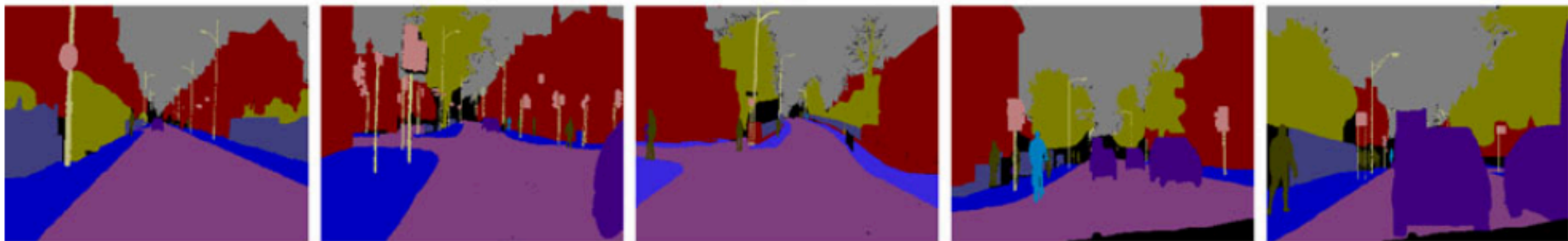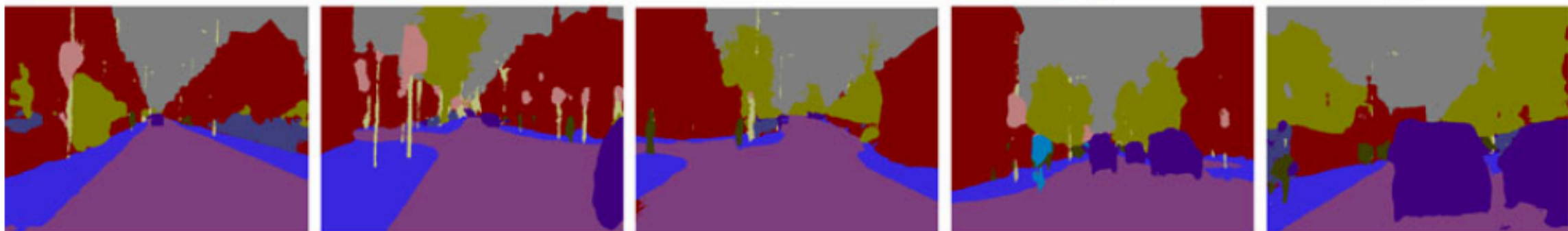
# SegNet Architecture
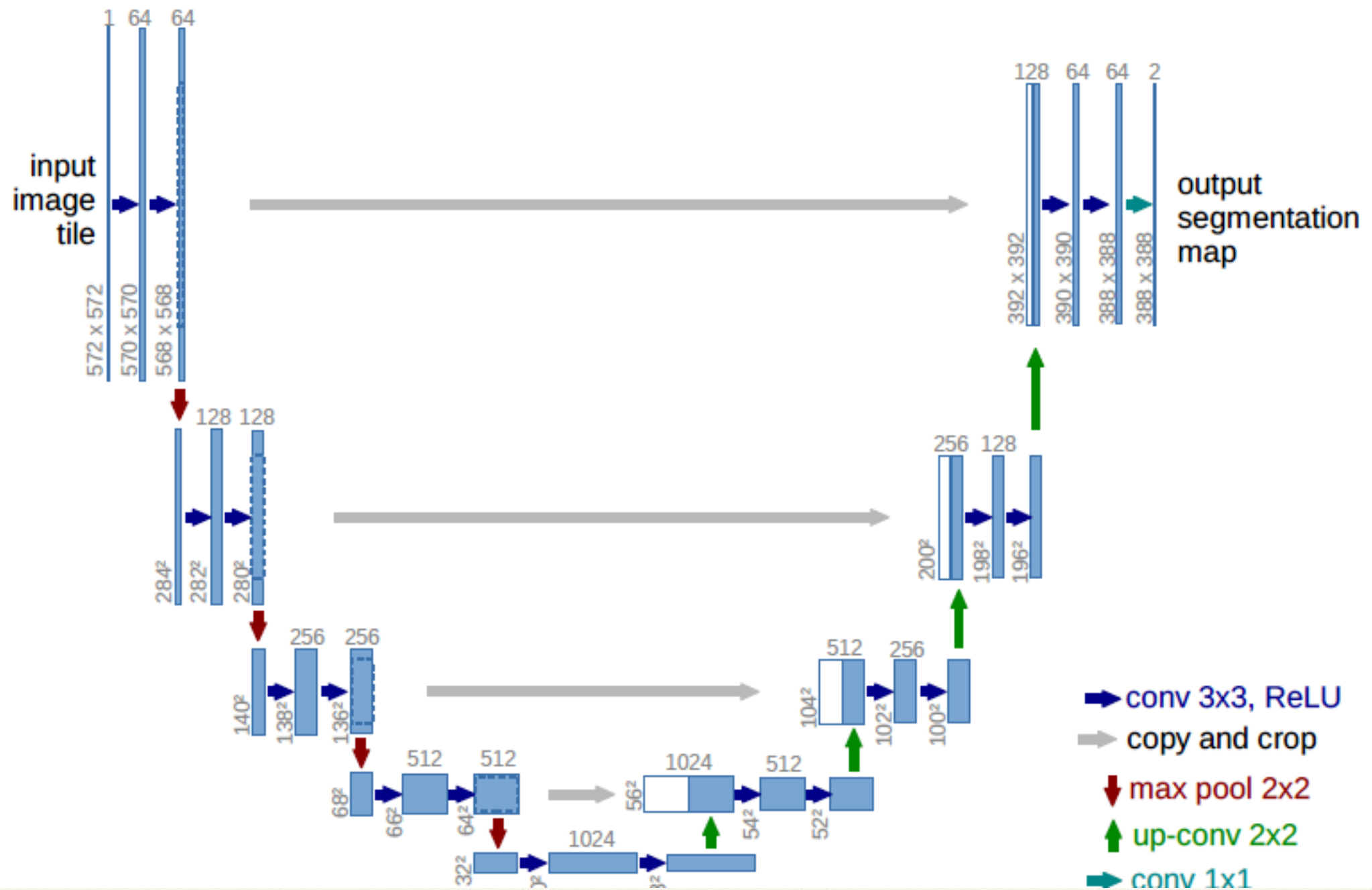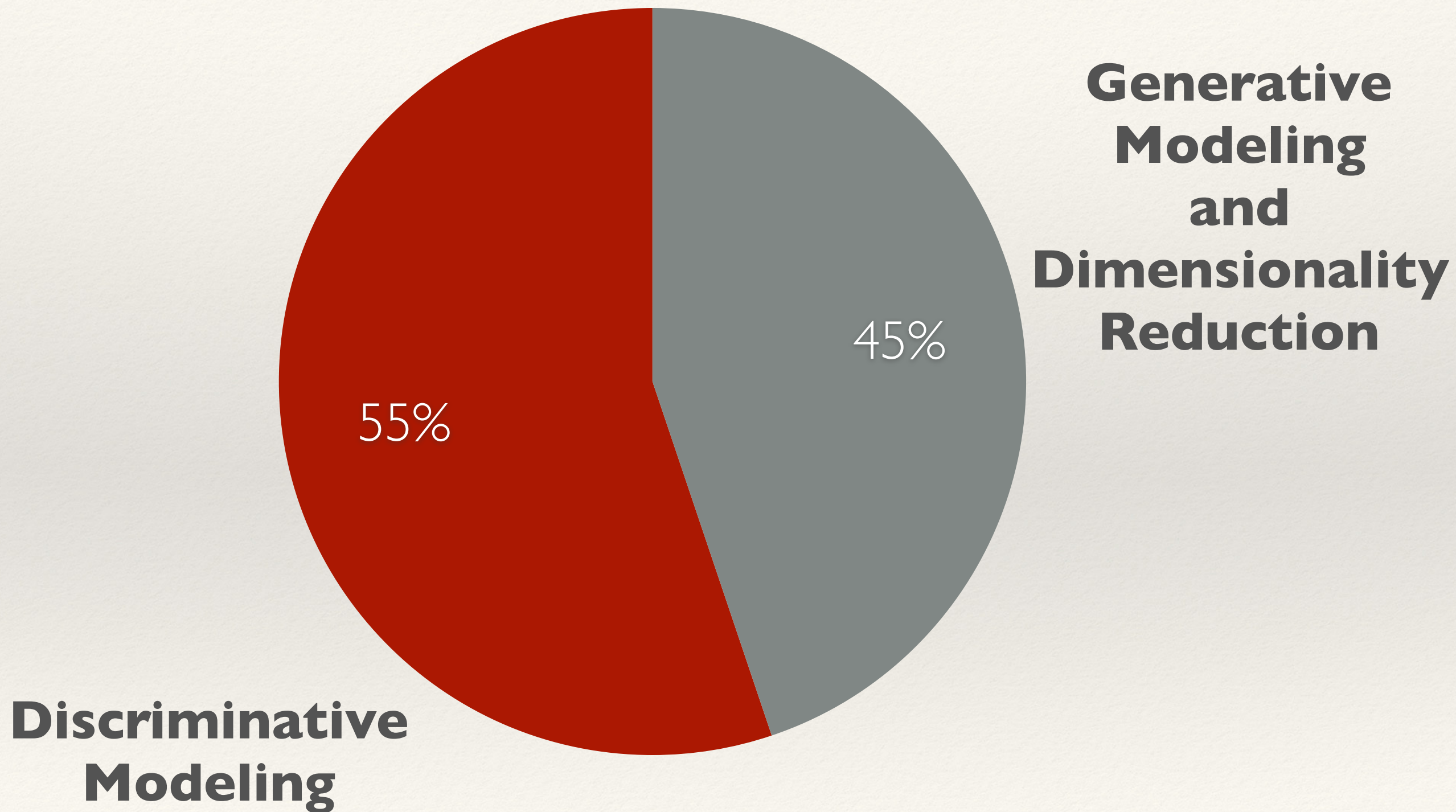
# Results from Segnet



Test samples

Ground Truth

SegNet

U-Net: Convolutional Networks for Biomedical Image Segmentation

# U-net

# Summary of the Course

# Generative Modeling and Dimensionality Reduction

When we started …

# Dates of Various Rituals

* 5 Assignments spread over 3 months (roughly one assignment every two weeks).

* September 1st week - project topic announcements.

* September 3rd week - 1st Midterm

* September 4th week - project topic and team finalization and proposal submission. [1 and 2 person teams].

* October 1st week - Project Proposal

* October 3rd week - 2nd MidTerm

* November 1st week - Project MidTerm Presentations.

* **December 1st week - Final Exams**

* **December 2nd week - Project Final Presentations.**

# Content Delivery