



# DEEP VISION IN ACTION

From classification and detection to segmentation

---

Duy V. Huynh, (NAIR) NewAI-VN Research

May 15, 2020

# Outline

---

1. Recap: CNN Backbone Architecture
2. Segmentation and Detection



don't worry about it if you don't  
understand

## Recap: CNN Backbone Architecture

---

# LeNet<sup>12</sup>

- Convolutional: locally connected, weight-sharing
- Subsampling
- Fully-connected layer
- Train by back-propagation

⇒ Basic components of modern  
ConvNets!

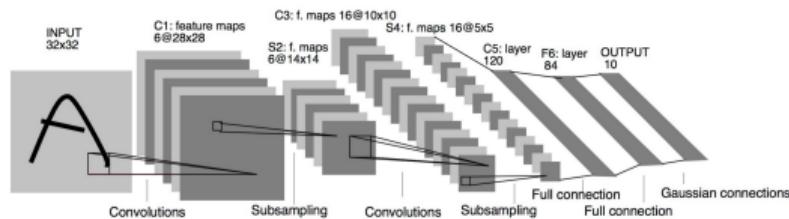


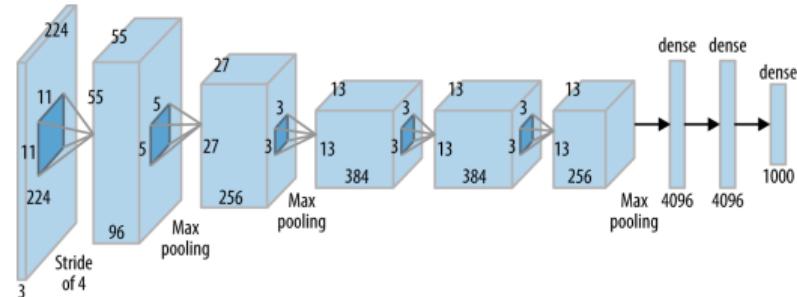
Figure 1: LeNet Architecture

<sup>1</sup>Yann Lecun et al. “Gradient-based learning applied to document recognition”. In: *Proceedings of the IEEE*. 1998, pp. 2278–2324.

<sup>2</sup>Le Cun et al. “Handwritten Digit Recognition with a Back-Propagation Network”. In: *Advances in Neural Information Processing Systems*. Morgan Kaufmann, 1990, pp. 396–404.

LeNet-style backbone, and:

- **ReLU** [Nair & Hinton 2010]
  - Revolution of deep learning.
  - Accelerate training, better than backprop + tanh.
- **Dropout** [Hinton et al 2012]
  - In-network ensembling.
  - Reduce overfitting



**Figure 2:** AlexNet Architecture

---

<sup>3</sup>Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks". In: *Advances in Neural Information Processing Systems*. P. 2012.

# VGG Net<sup>4</sup>

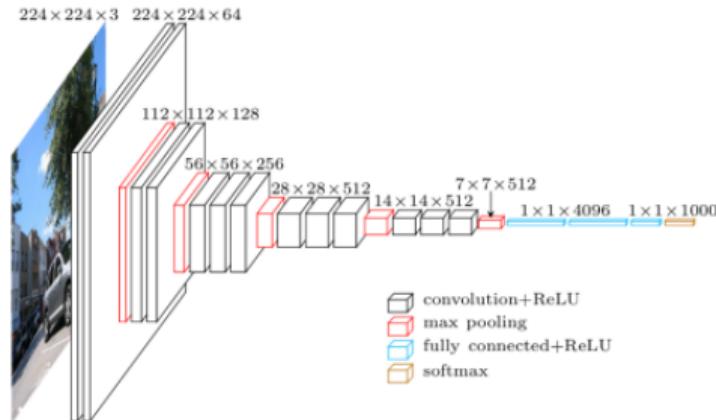


Figure 3: VGG-16 architecture

Modularized design:

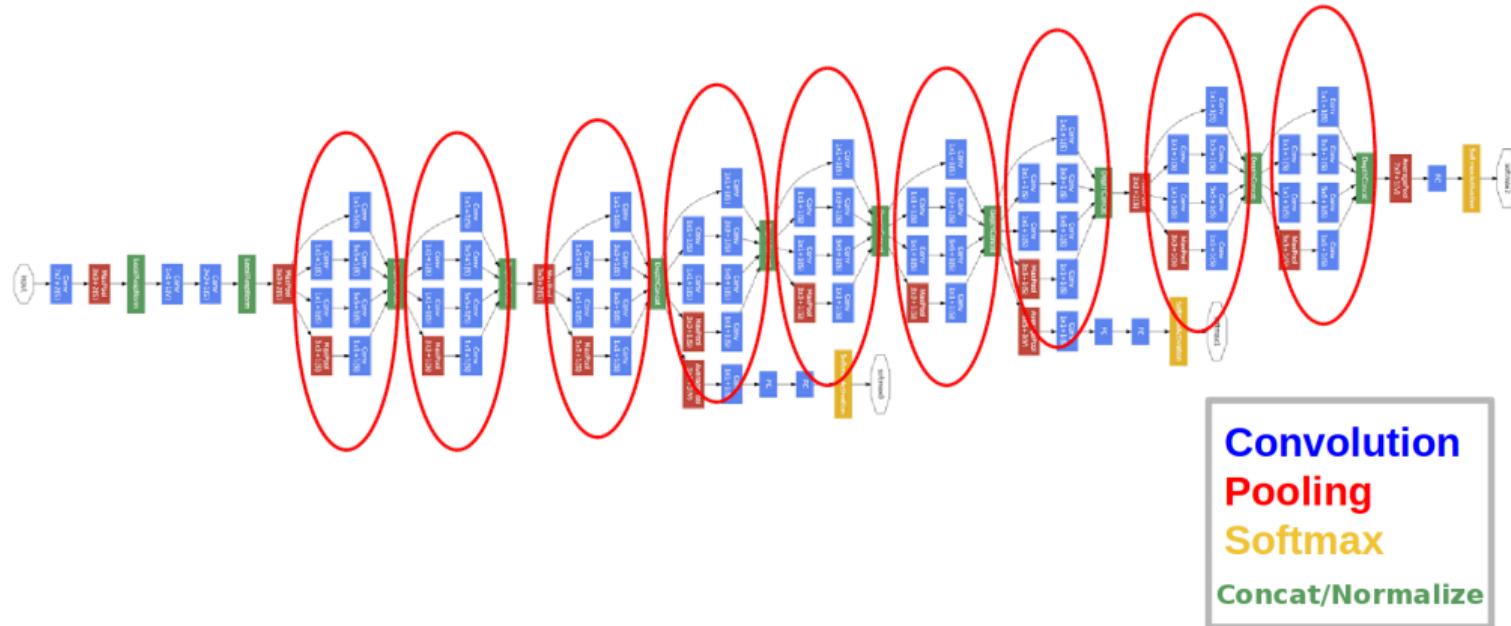
- $3 \times 3$  Conv = 1 module
- Stack the same module
- Same computation for each module  
( $1/2$  spatial size  $\Rightarrow$   $2x$  filters)

Stage-wise training

- VGG-11  $\Rightarrow$  VGG-13  $\Rightarrow$  VGG-16

<sup>4</sup>Karen Simonyan and Andrew Zisserman. *Very deep convolutional networks for large-scale image recognition.* 2014.

# GoogleNet/Inception<sup>5</sup>



**Figure 4:** GoogleNet/Inception architecture

<sup>5</sup>Christian Szegedy et al. *Going deeper with convolutions.* 2014.

## $1 \times 1$ convolution<sup>6</sup>

$$\begin{array}{ccc} \text{Input Tensor: } & \text{Filter: } & \text{Output Tensor: } \\ \begin{array}{c} \text{6} \times 6 \times 32 \\ \text{cube} \end{array} & \begin{array}{c} \text{*} \\ \text{1} \times 1 \times 32 \\ \text{cube} \end{array} & = \\ & & \begin{array}{c} \text{6} \times 6 \times 1 \\ \text{matrix} \end{array} \end{array}$$

$$\begin{array}{ccc} \text{Input Tensor: } & \text{Filter: } & \text{Output Tensor: } \\ \begin{array}{c} \text{6} \times 6 \times 32 \\ \text{cube} \end{array} & \begin{array}{c} \text{*} \\ \text{1} \times 1 \times 32 \\ \text{n_c filters} \end{array} & = \\ & & \begin{array}{c} \text{6} \times 6 \times \text{n_c} \\ \text{matrix stack} \end{array} \end{array}$$

**Figure 5:**  $1 \times 1$  conv with 1 filter

**Figure 6:**  $1 \times 1$  conv with n filter

---

<sup>6</sup>Min Lin, Qiang Chen, and Shuicheng Yan. *Network in Network*. 2013.

## Why $1 \times 1$ convolution?

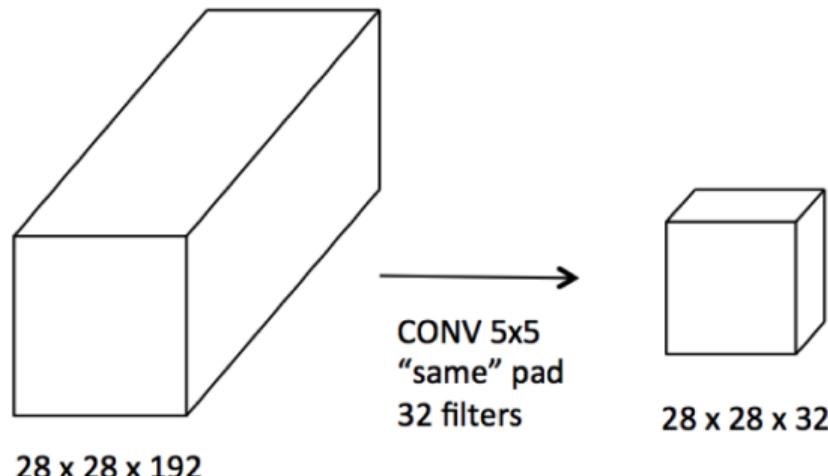
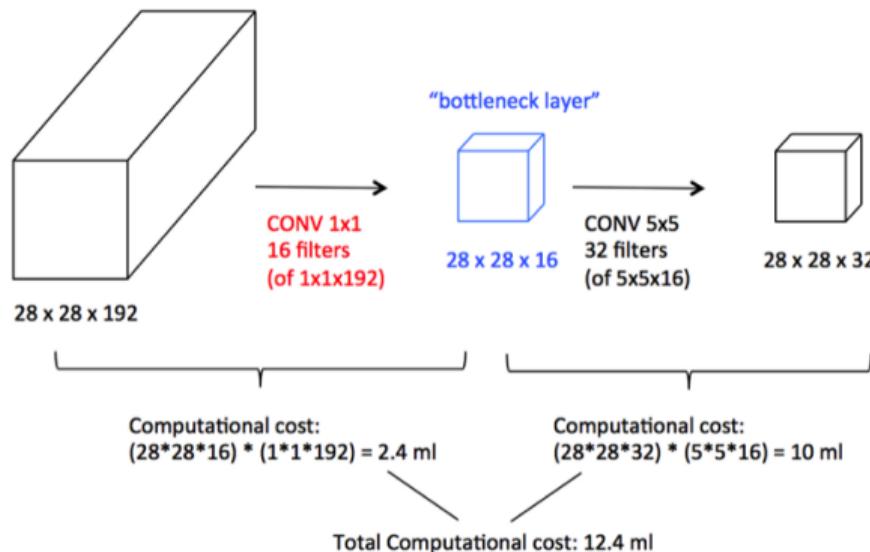


Figure 7:  $5 \times 5$  filter  $\times 32$

$\implies$  Computational cost:  $(28 \times 28 \times 32) \times (5 \times 5 \times 192) = 120.4M$

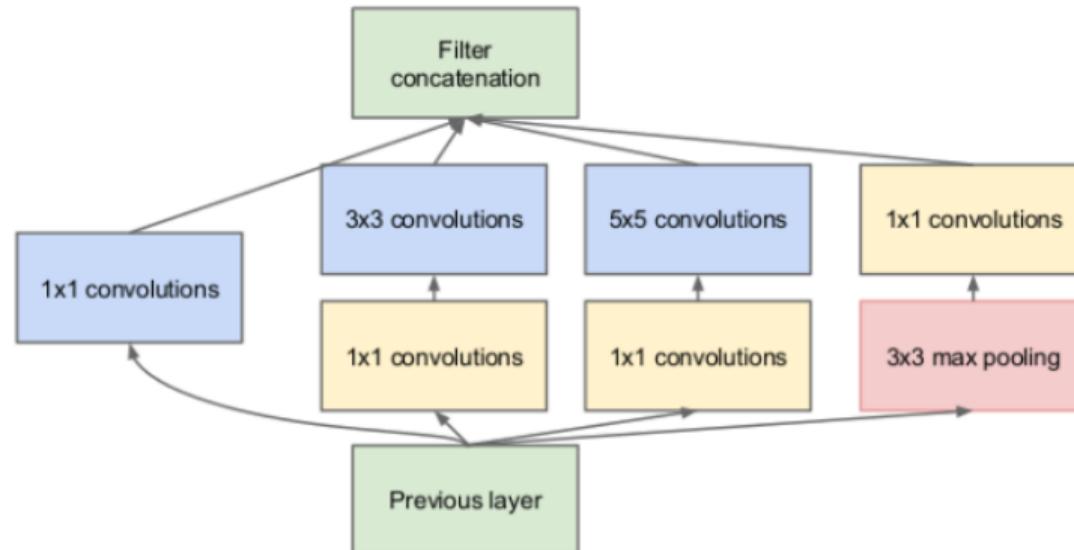
# Why $1 \times 1$ convolution?



**Figure 8:** "Bottleneck" transfer

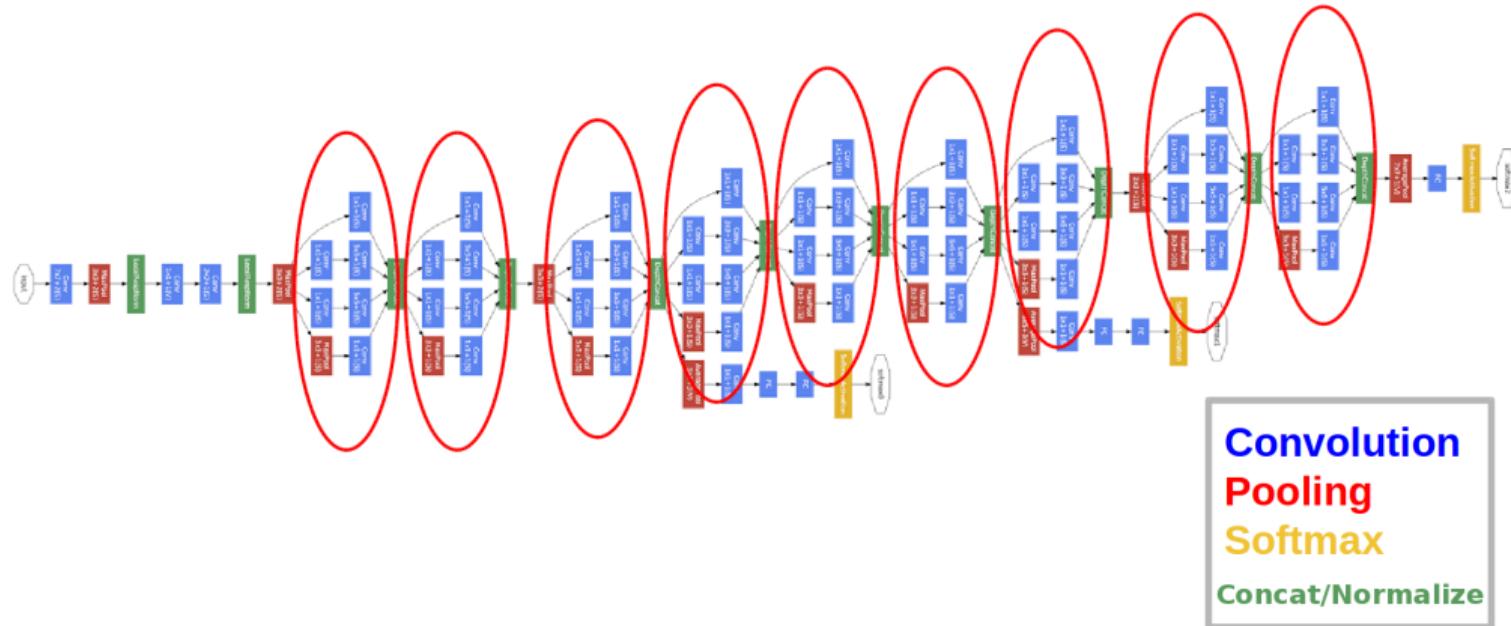
⇒  $1 \times 1$  reduces by nearly 10 times the total computational cost.

# Inception Module



**Figure 9:** Inception module

# GoogleNet/Inception<sup>7</sup>



**Figure 10:** GoogleNet/Inception architecture

<sup>7</sup>Christian Szegedy et al. *Going deeper with convolutions.* 2014.

## Conclusions

- Multiple "branches"
  - e.g.,  $1 \times 1$  conv,  $3 \times 3$  conv,  $5 \times 5$  conv, pool, ...
- Shorcuts
  - stand-alone  $1 \times 1$ , merged by concatenation method.
- Bottlenek
  - Reduce dimensionality by  $1 \times 1$  conv.

---

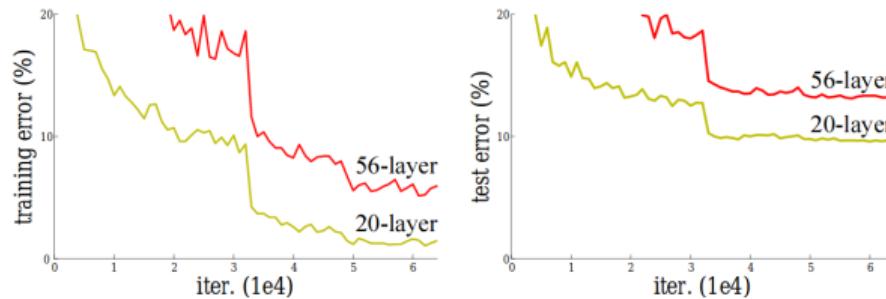
<sup>8</sup>Christian Szegedy et al. *Going deeper with convolutions*. 2014.



I WAS WINNING  
IMAGENET

UNTIL A  
DEEPER MODEL  
CAME ALONG

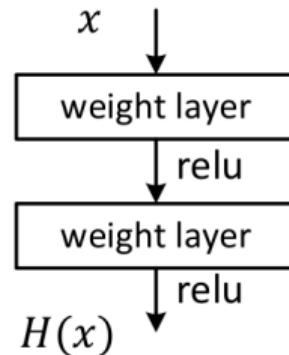
## Stacking more and more layers?<sup>9</sup>



**Figure 11:** Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer plain networks

<sup>9</sup>Kaiming He et al. "Deep Residual Learning for Image Recognition". In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778.

# Deep Residual Learning<sup>10</sup>

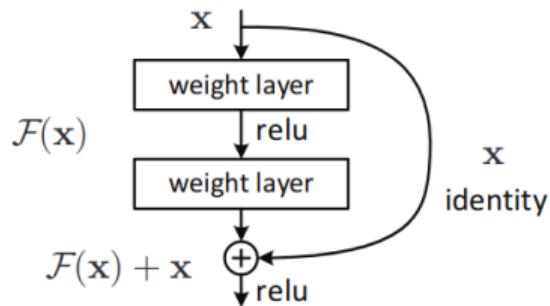


$H(x)$  is any desired mapping, hope the small subnet fit  $H(x)$ .

**Figure 12:** "Plain" connection in "plain" net

<sup>10</sup>Kaiming He et al. "Deep Residual Learning for Image Recognition". In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778.

# Deep Residual Learning<sup>11</sup>



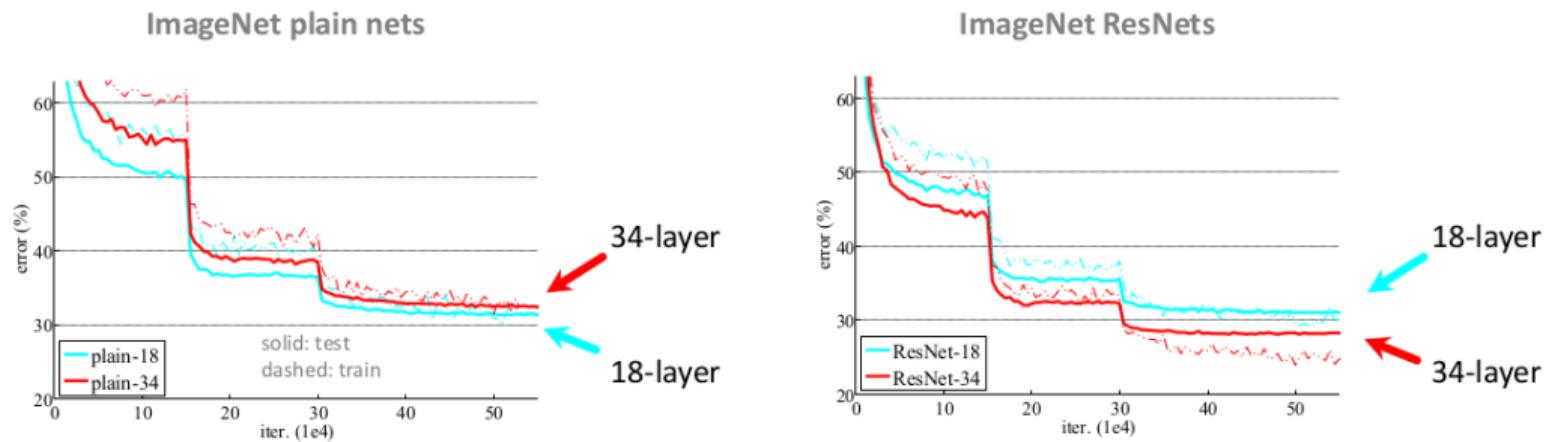
**Figure 13:** Residual connection in residual net

$H(x)$  is any desired mapping, ~~hope the small subnet fit  $H(x)$~~  hope the small subnet fit  $F(x)$  let  $H(x) = F(x) + x$ .

- If identity were optimal, easy to set weights as 0.
- If optimal mapping is closer to identity, easier to find small fluctuations.

<sup>11</sup>Kaiming He et al. "Deep Residual Learning for Image Recognition". In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778.

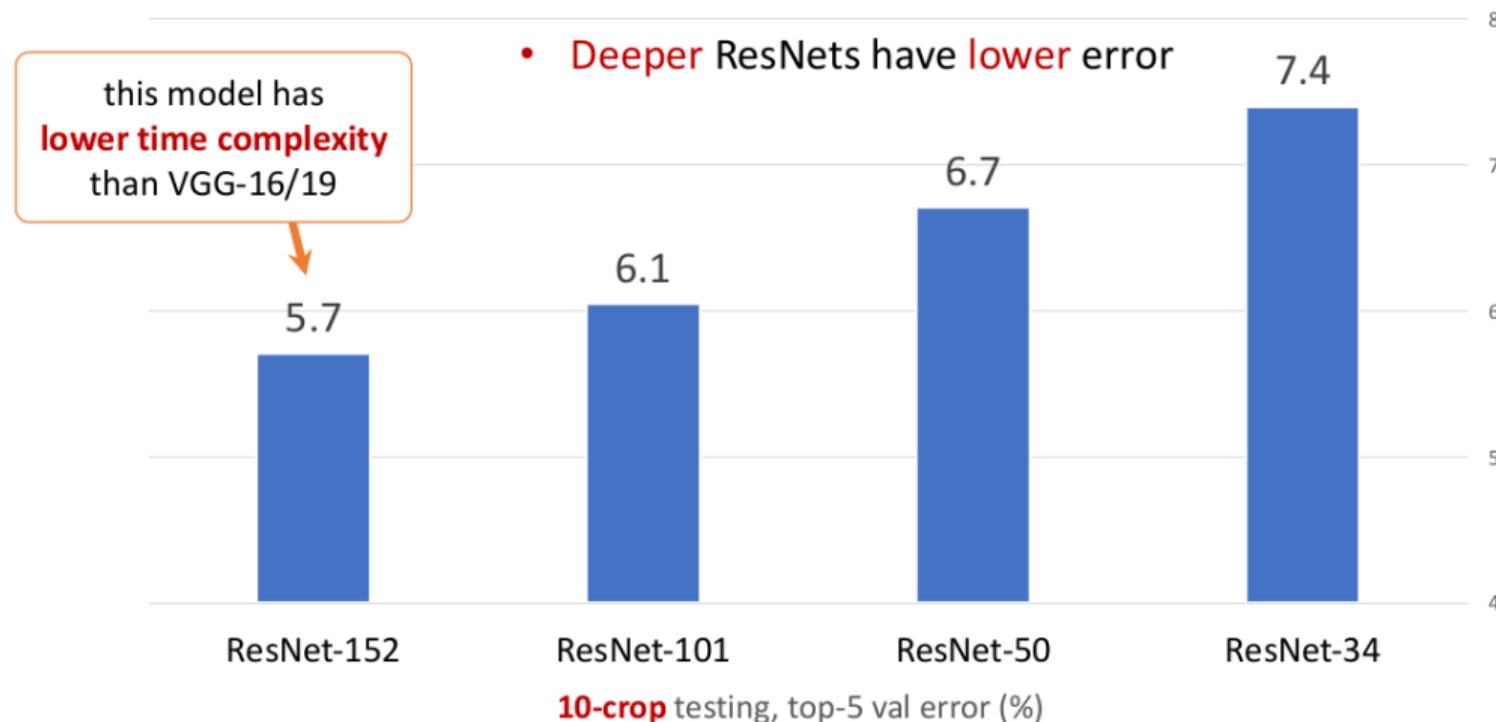
# ImageNet Experiments



**Figure 14:** Training on ImageNet

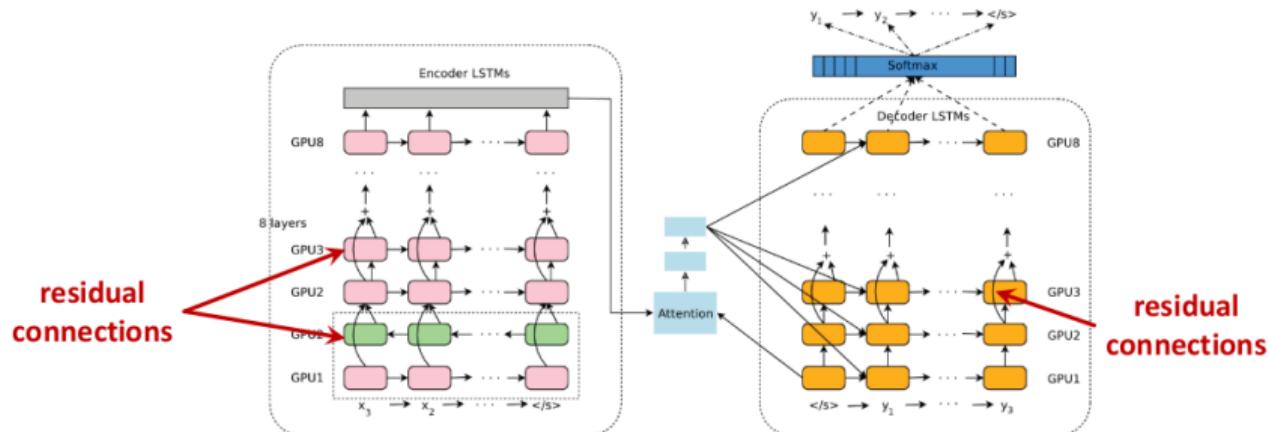
- Deep ResNets can be trained without any difficulties.
- Deeper ResNets have lower training error, and also lower test error.

## ImageNet Experiments



**Figure 15:** Training more and more deeper model on ImageNet

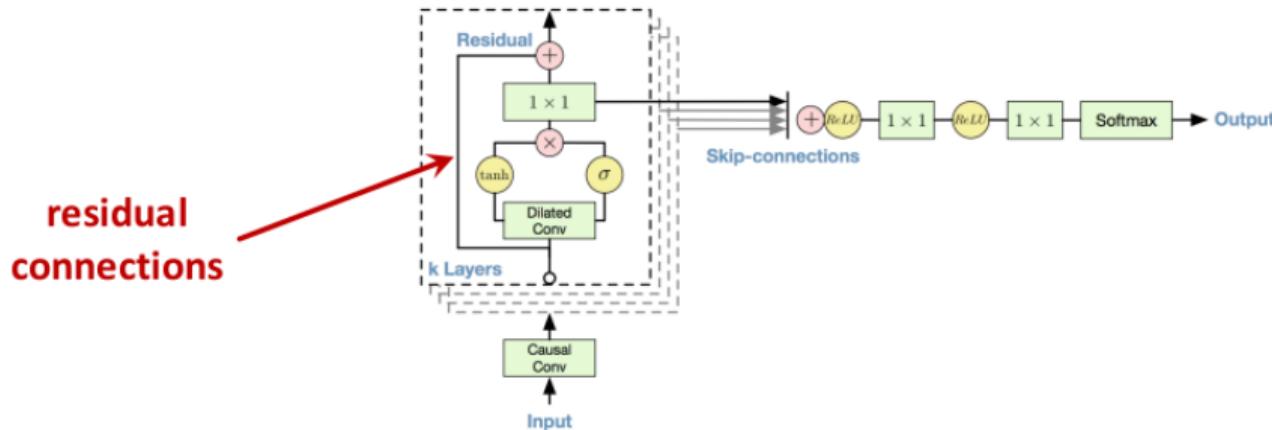
# Residual connection beyond Computer Vision



**Figure 16:** Neural Machine Translation (NMT)<sup>12</sup>: 8-layer LSTM

<sup>12</sup>Yonghui Wu et al. "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation". In: *CoRR* abs/1609.08144 (2016).

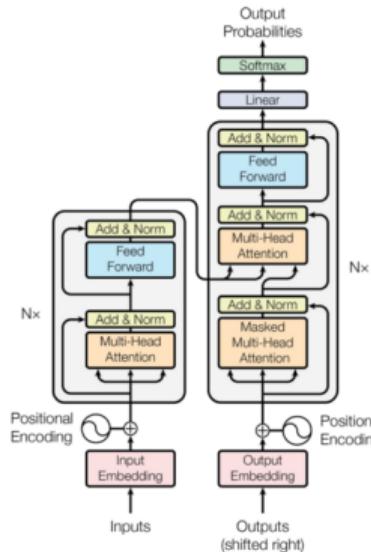
# Residual connection beyond Computer Vision



**Figure 17:** Speech synthesis (WaveNet)<sup>13</sup>: Residual CNNs on 1-d sequence

<sup>13</sup>Aäron van den Oord et al. “WaveNet: A Generative Model for Raw Audio”. In: SSW. 2016.

# Residual connection beyond Computer Vision



**Figure 18:** Transformer<sup>14</sup>:

<sup>14</sup>Ashish Vaswani et al. "Attention Is All You Need". In: *NIPS*. 2017.

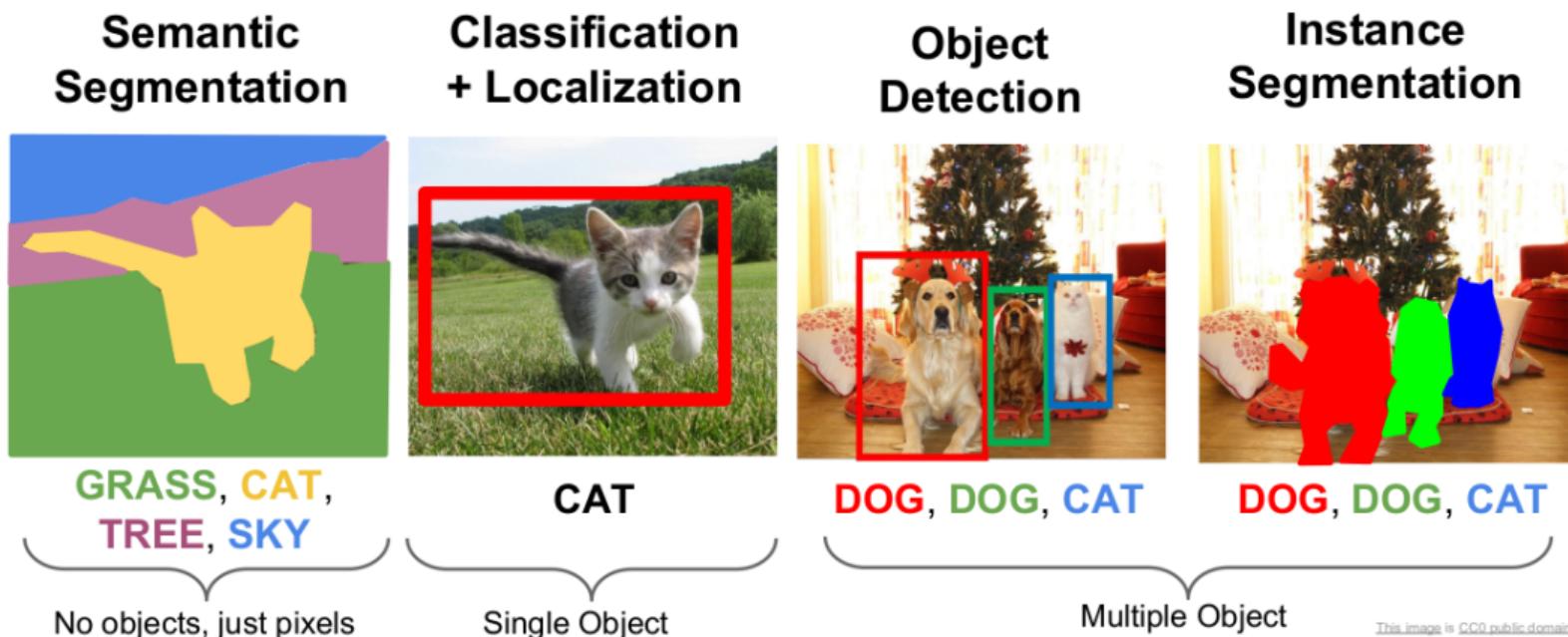
## More architectures

- **Inception-ResNet**<sub>[Szegedy et al 2017]</sub>: Inception as transformation + residual connection.
- **DenseNet**<sub>[Huang et al CVPR 2017]</sub>: Densely connected shortcuts w/ concat.
- **Xception**<sub>[Chollet CVPR 2017]</sub>, **MobileNets**<sub>[Howard et al 2017]</sub>: DepthwiseConv
- **ShuffleNet**<sub>[Zhang et al 2017]</sub>: More group/DepthwiseConv + shuffle
- **ResNeXt**<sub>[He et al CVPR 2017]</sub>: Residual connection, multi branch, group-conv
- ...

## **Segmentation and Detection**

---

# Computer Vision is not only Classification



# Computer Vision is not only Classification

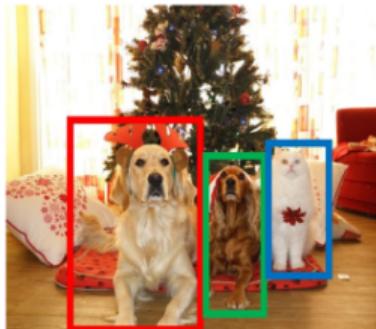
## Semantic Segmentation



GRASS, CAT,  
TREE, SKY

No objects, just pixels

## 2D Object Detection



DOG, DOG, CAT

Object categories +  
2D bounding boxes

## 3D Object Detection

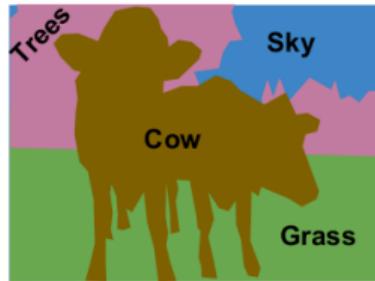
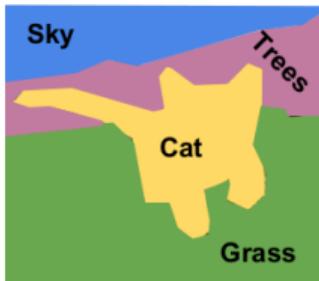


Car

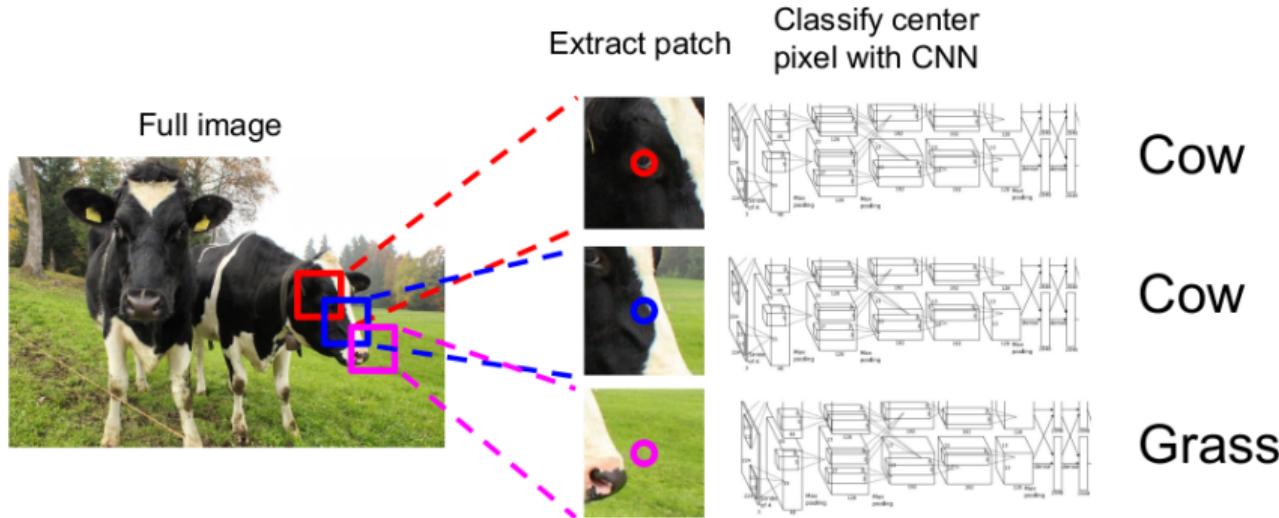
Object categories +  
3D bounding boxes

# Semantic Segmentation

- Label each pixel in the image with a category label.
- Don't differentiate instances, only care about pixels.

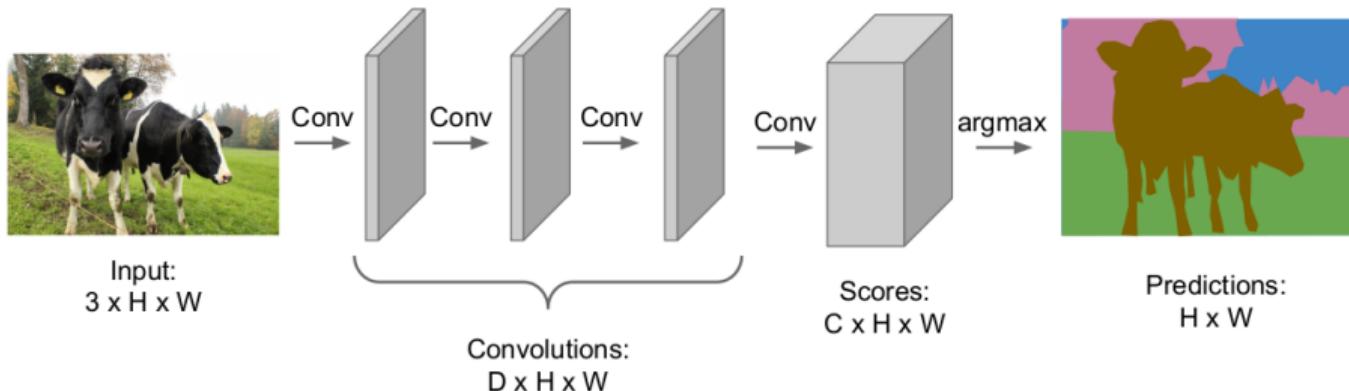


# Semantic Segmentation Idea: Sliding Window



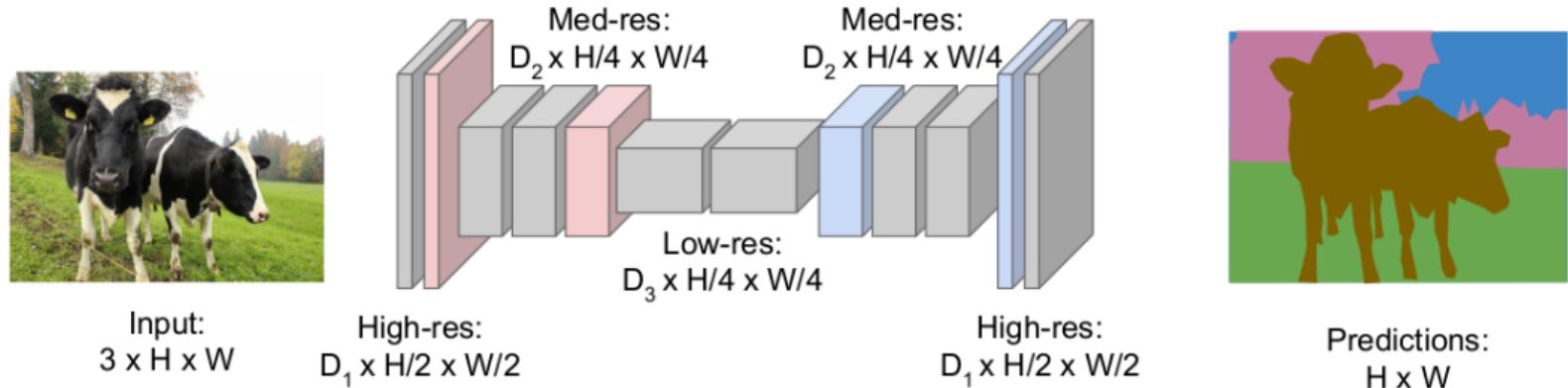
⇒ **Problem:** Computational cost is very expensive! Not reusing shared features between overlapping patches.

## Semantic Segmentation Idea: Fully Convolutional



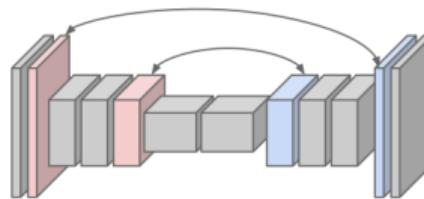
- Make predictions for **all pixels at once**
- ⇒ **Problem** again: Convolutions at original image resolution is still very expensive.

## Semantic Segmentation Idea: Fully Convolutional



**Figure 19:** Fully Convolutional Networks for Semantic Segmentation

## Semantic Segmentation Idea: Fully Convolutional



- Design network as a bunch of convolutional layers, with **down-sampling** and **up-sampling** inside<sup>15</sup>.
- **Down-sampling**: Pooling, strided convolution  $\implies$  **Up-sampling, Un-pooling ???**

---

<sup>15</sup>Evan Shelhamer, Jonathan Long, and Trevor Darrell. "Fully Convolutional Networks for Semantic Segmentation". In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 3431–3440.

# Naive Idea for Up-sampling

## Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4

5	6
7	8

Output: 2 x 2

Rest of the network

## Max Unpooling

Use positions from pooling layer

1	2
3	4

Input: 2 x 2

0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

Output: 4 x 4

# Naive Idea for Up-sampling

## Max Pooling

Remember which element was max!

1	2	6	3
3	5	2	1
1	2	2	1
7	3	4	8

Input: 4 x 4

5	6
7	8

Output: 2 x 2

Rest of the network

## Max Unpooling

Use positions from pooling layer

1	2
3	4

Input: 2 x 2

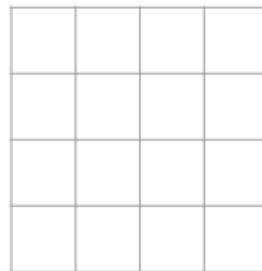
0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

Output: 4 x 4

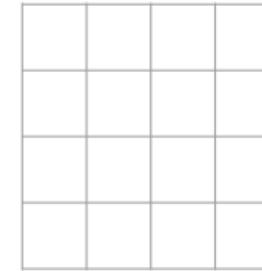
- Nonlearnable  $\implies$  Not a very good idea.

## Learnable Up-sampling: Transpose Convolution

**Recall:** Typical  $3 \times 3$  convolution, stride 1 pad 1



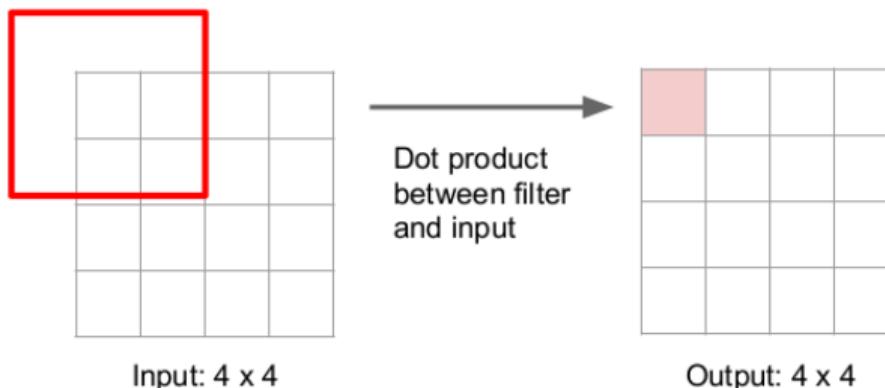
Input:  $4 \times 4$



Output:  $4 \times 4$

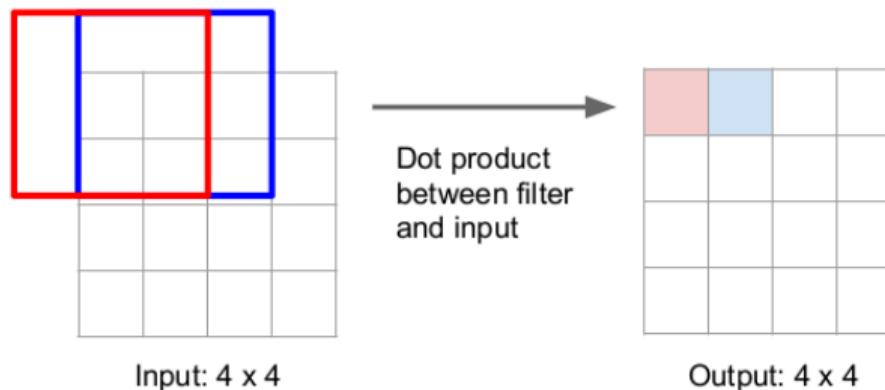
## Learnable Up-sampling: Transpose Convolution

**Recall:** Normal  $3 \times 3$  convolution, stride 1 pad 1



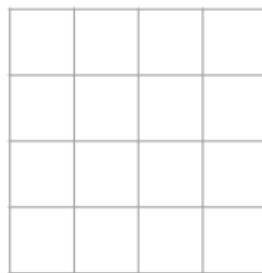
## Learnable Up-sampling: Transpose Convolution

**Recall:** Normal  $3 \times 3$  convolution, stride 1 pad 1

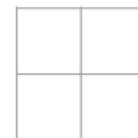


## Learnable Up-sampling: Transpose Convolution

**Recall:** Normal  $3 \times 3$  convolution, stride 2 pad 1



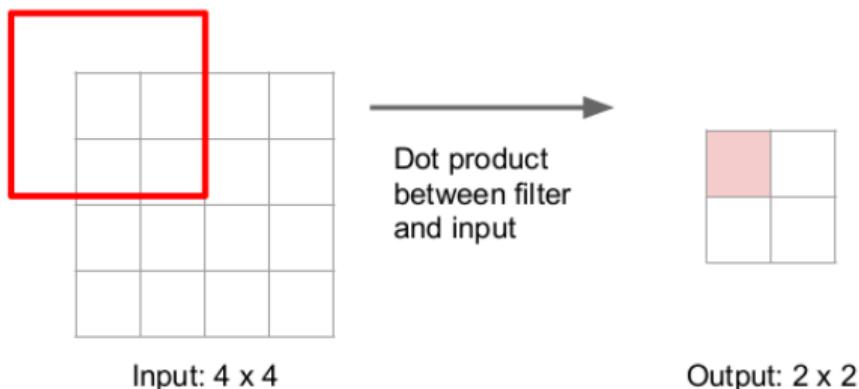
Input:  $4 \times 4$



Output:  $2 \times 2$

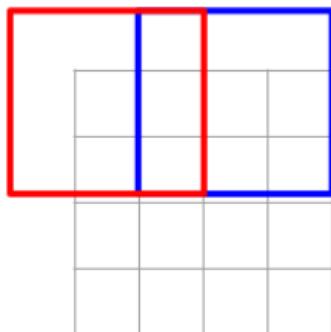
## Learnable Up-sampling: Transpose Convolution

**Recall:** Normal  $3 \times 3$  convolution, stride 2 pad 1



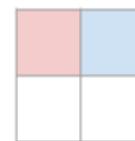
# Learnable Up-sampling: Transpose Convolution

**Recall:** Normal  $3 \times 3$  convolution, stride 2 pad 1



Input:  $4 \times 4$

Dot product  
between filter  
and input



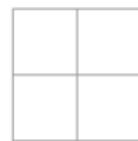
Output:  $2 \times 2$

Filter moves 2 pixels in  
the input for every one  
pixel in the output

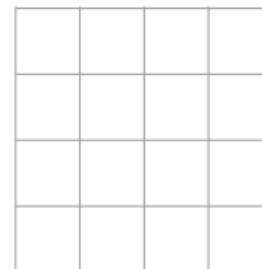
Stride gives ratio between  
movement in input and  
output

## Learnable Up-sampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1



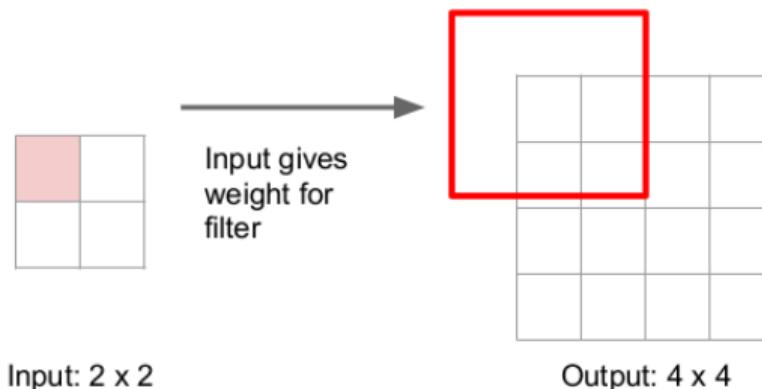
Input: 2 x 2



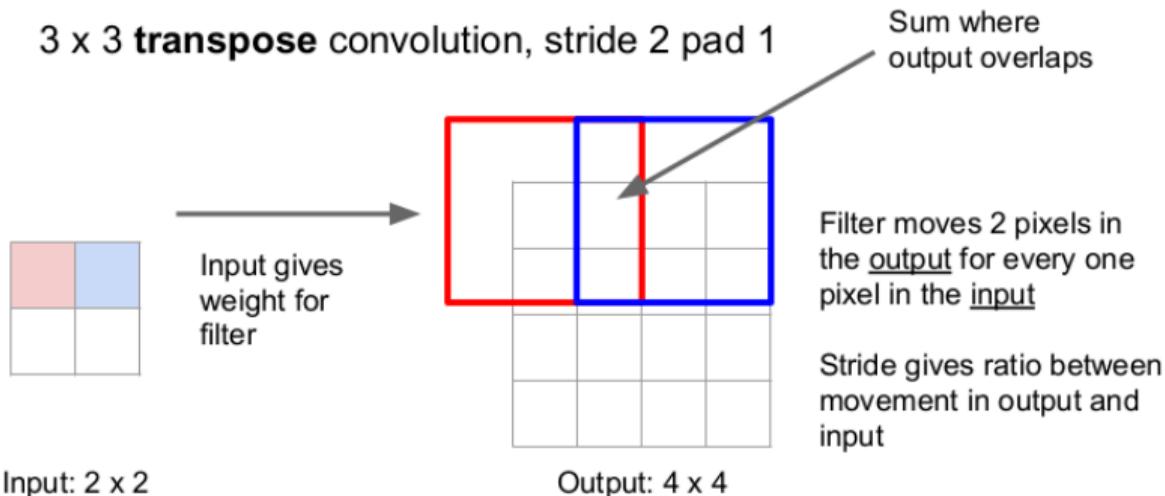
Output: 4 x 4

## Learnable Up-sampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

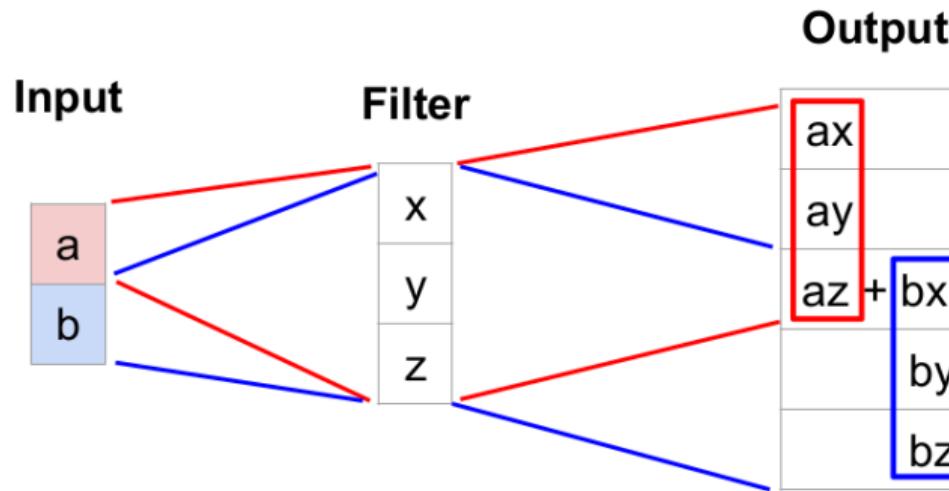


## Learnable Up-sampling: Transpose Convolution



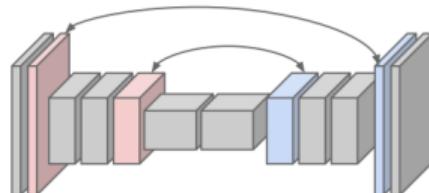
⇒ Transpose convolution a.k.a. Deconvolution, Upconvolution, Fractionally strided convolution, Backward convolution, ..

## Learnable Up-sampling: 1D Example



- Output contains copies of the filter weighted by the input, summing at where it overlaps in the output.
- Need to crop one pixel from output to make output exactly 2x input.

## Semantic Segmentation Idea: Fully Convolutional

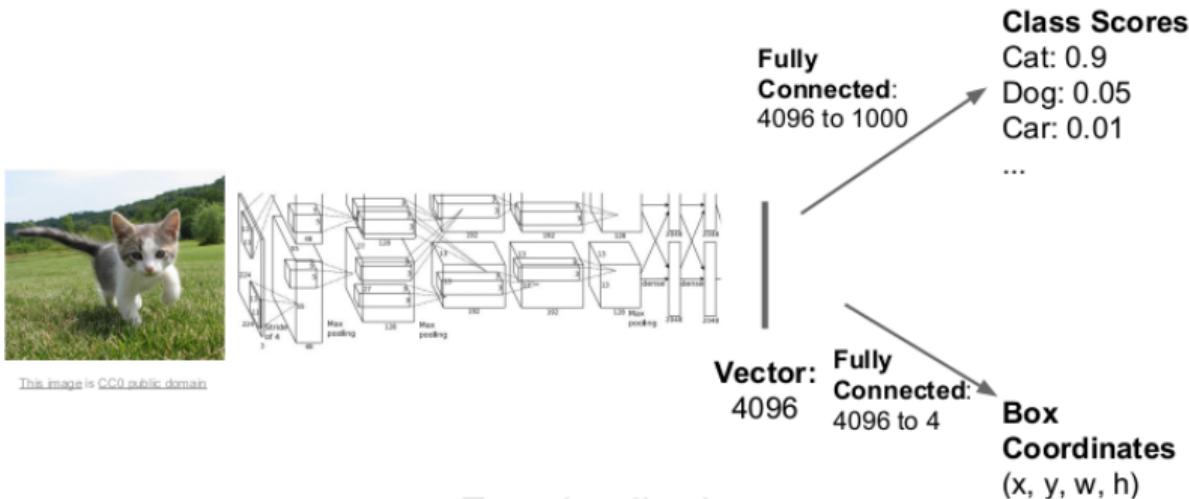


- Design network as a bunch of convolutional layers, with **down-sampling** and **up-sampling** inside<sup>16</sup>.
- **Down-sampling:** Pooling, strided convolution  $\implies$  **Up-sampling:** Un-pooling, strided transpose convolution

---

<sup>16</sup>Evan Shelhamer, Jonathan Long, and Trevor Darrell. "Fully Convolutional Networks for Semantic Segmentation". In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 3431–3440.

# Classification and Localization



**Figure 20:** Single Object Detection

## Classification and Localization

---

How would we train such a single object detection network that produces  $4 + n$  values?

- Predicting the class of the object ( $n$  class probabilities) is a **classification** problem.
  - Predicting the four coordinates for the bounding box is a **regression** problem.
- ⇒ Need a loss function that combines these two problems.

# Classification and Localization

---

- Localization loss

$$\mathcal{L}_{loc} = \sum_{i=1}^4 |x_{pred,i} - x_{truth,i}|$$

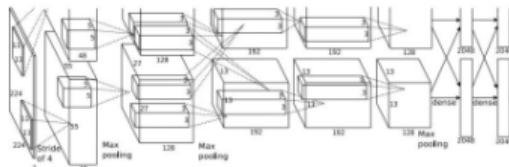
- Confidence loss

$$\mathcal{L}_{conf} = \frac{-1}{N} (\sum_{i=1}^N y_{truth,i} \cdot \log(y_{pred,i}))$$

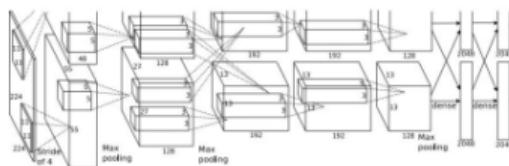
⇒ Combined loss function

$$\mathcal{L} = \mathcal{L}_{loc} + \alpha \cdot \mathcal{L}_{conf}$$

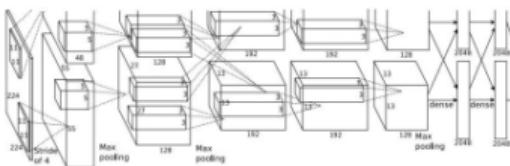
# Classification and Localization



CAT:  $(x, y, w, h)$  **4 numbers**



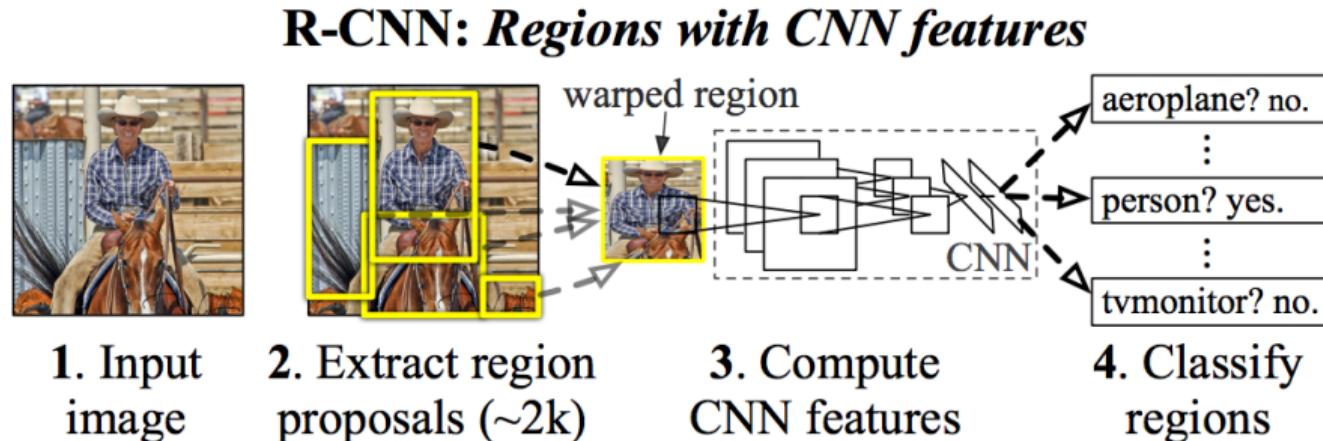
DOG:  $(x, y, w, h)$   
DOG:  $(x, y, w, h)$  **16 numbers**  
CAT:  $(x, y, w, h)$



DUCK:  $(x, y, w, h)$  **Many**  
DUCK:  $(x, y, w, h)$  **numbers!**

...

# Multiple Object Detection: Region Proposal Methods



**Figure 21:** Region Convolutional Neural Network<sup>17</sup>

<sup>17</sup>Ross Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2014.

# Multiple Object Detection: Region Proposal Methods

## Problems with R-CNN

- Extracting 2,000 regions for each image based on selective search.
- Extracting features using CNN for every image region. Suppose we have  $N$  images, then the number of CNN features will be  $N \times 2,000$ .
- The entire process of object detection using R-CNN has three models:
  - CNN for feature extraction.
  - Linear SVM classifier for identifying objects.
  - Regression model for tightening the bounding boxes.

⇒ **Very slow**, it takes around 40-50 seconds to make predictions for each image.

# Multiple Object Detection: Region Proposal Methods

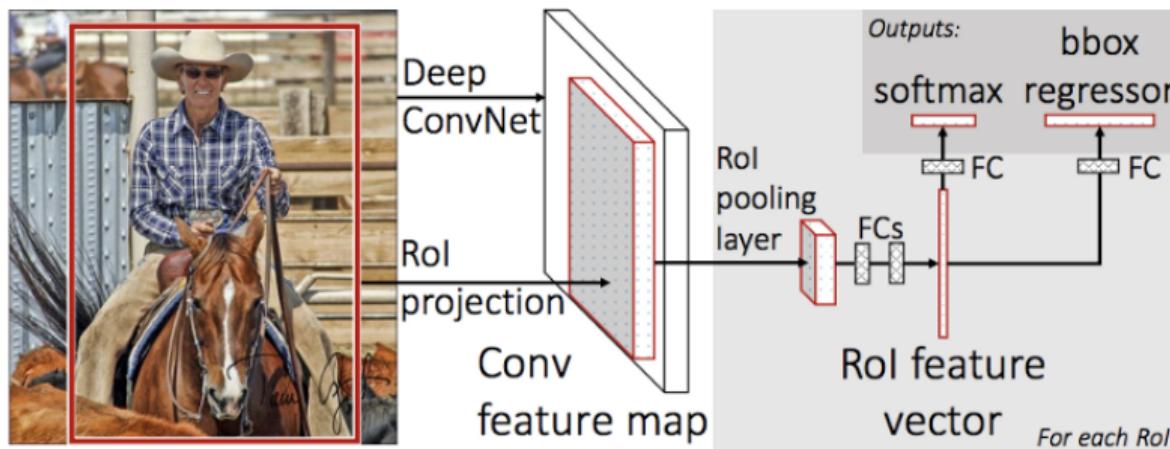


Figure 22: Fast R-CNN<sup>18</sup>

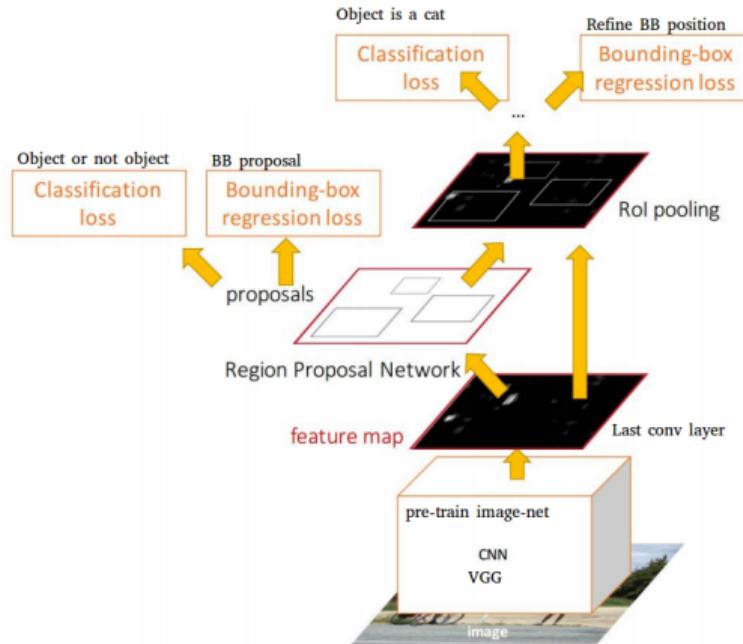
<sup>18</sup>Ross B. Girshick. "Fast R-CNN". In: *2015 IEEE International Conference on Computer Vision (ICCV)* (2015), pp. 1440–1448.

## Multiple Object Detection: Region Proposal Methods

---

⇒ Still have a **problem**, it also uses selective search as a proposal method to find the Regions of Interest, which is a slow and time consuming process. It takes around 2 seconds per image to detect objects, which is much better compared to RCNN

# Multiple Object Detection: Region Proposal Methods

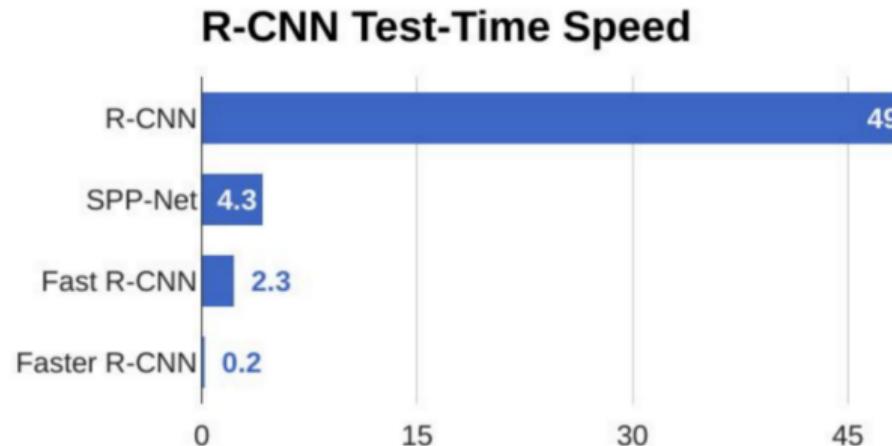


Insert **Region Proposal Network (RPN)** to predict proposals from features. Jointly train with 4 losses:

- RPN classify object/not object
- RPN regress box coordinates
- Final classification score (object classes)
- Final box coordinates

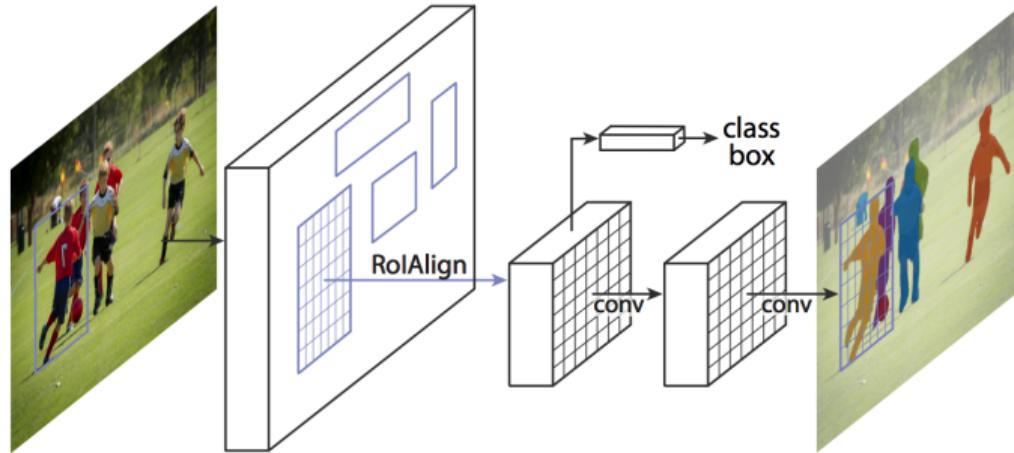
Figure 23: Faster R-CNN

## Multiple Object Detection: Region Proposal Methods



**Figure 24:** Experimentation results

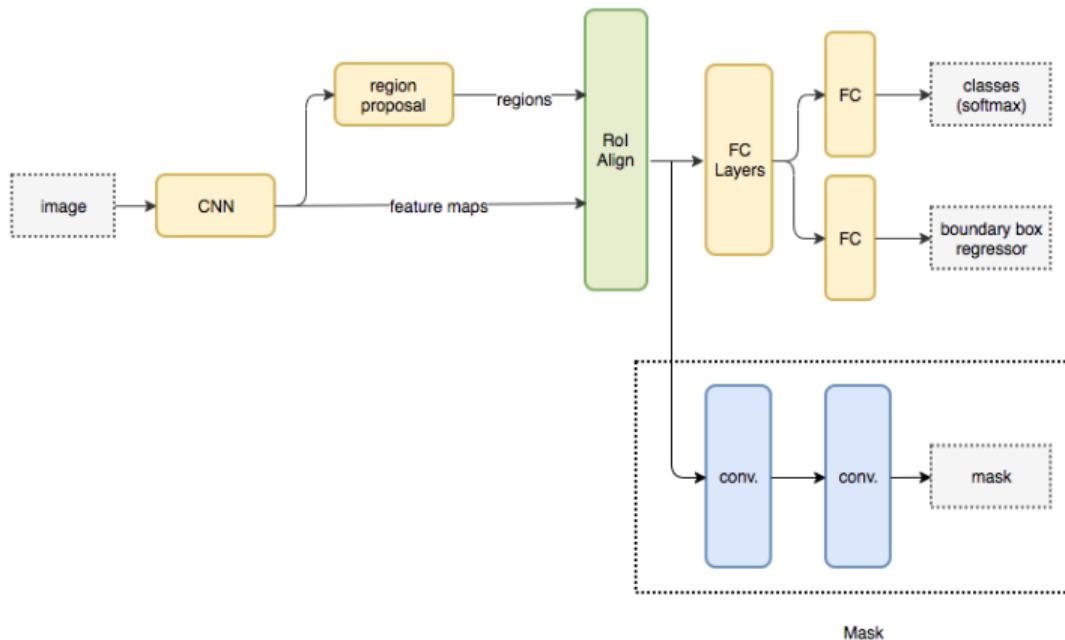
# Instance Segmentation



**Figure 25:** Mask R-CNN<sup>19</sup>

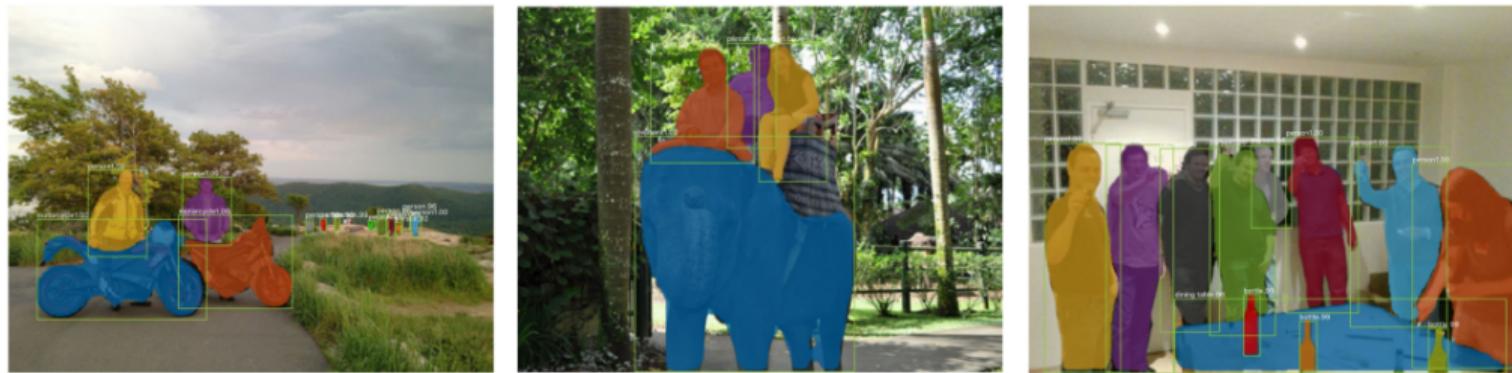
<sup>19</sup>Kaiming He et al. "Mask R-CNN". In: *2017 IEEE International Conference on Computer Vision (ICCV)* (2017), pp. 2980–2988.

# Instance Segmentation



**Figure 26:** Mask R-CNN architecture

# Instance Segmentation



**Figure 27:** Amazing results

# You Look Only Once

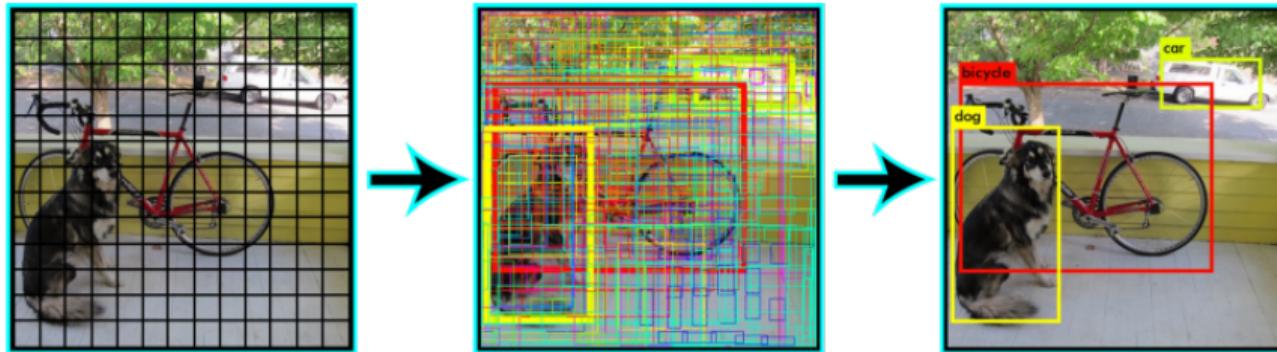
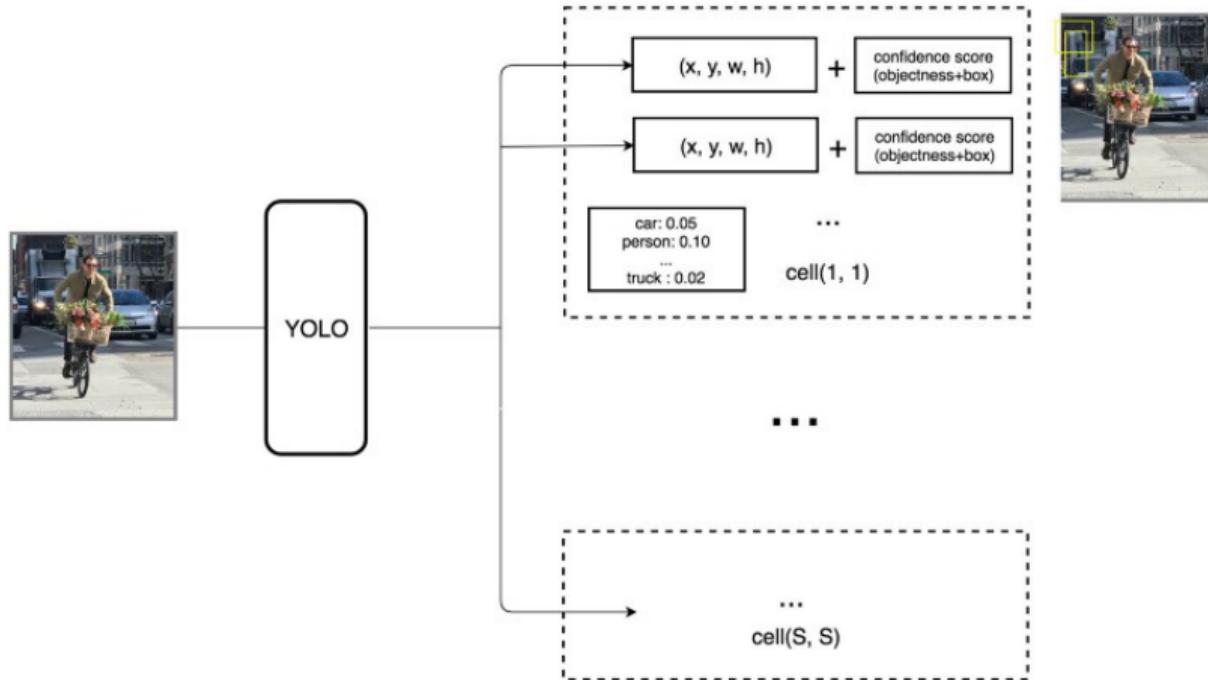


Figure 28: YOLO method motivation<sup>20</sup>

<sup>20</sup>Joseph Redmon et al. "You Only Look Once: Unified, Real-Time Object Detection". In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 779–788.

# You Look Only Once



YOLO divide input image into  $S \times S$  grid cells. It will predict:

- **B** boundary boxes and each box has one box confidence score (5 elements  $(x, y, w, h)$  and confidence score)
- **One** object only regardless of the number of boxes B
- **C** is a conditional class probabilities (one per class for the likeliness of the object class). The conditional class probability is the probability that the detected object belongs to a particular class.

⇒ Output tensor shape  $(S, S, B \times 5 + C)$

# You Look Only Once

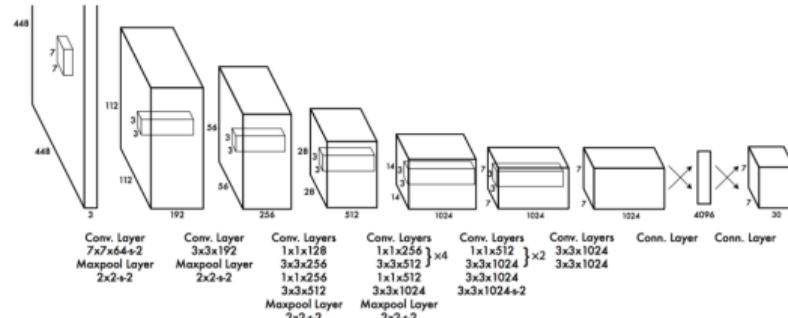


Figure 30: Network architecture design

YOLO has 24 convolutional layers followed by 2 fully connected layers (FC). Some convolution layers use  $1 \times 1$  reduction layers alternatively to reduce the depth of the features maps.

YOLO uses sum-squared error between the predictions and the ground truth to calculate loss. The loss function composes of:

- Classification loss
- Localization loss
- Confidence loss

# You Look Only Once

---

Final loss function:

$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ \left( \sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left( \sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} \left( C_i - \hat{C}_i \right)^2 \\ & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

## Intersection Over Union

$$\text{IoU} = \frac{\text{Area of Overlap}}{\text{Area of Union}}$$



## Non-maxima suppression algorithm

1. Sort the predictions by the confidence scores.
2. Start from the top scores, ignore any current prediction if we find any previous predictions that have the same class and  $\text{IoU} \geq 0.5$  with the current prediction.
3. Repeat step 2 until all predictions are checked.

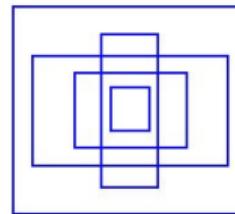
## YOLOv2: Accuracy Improvements

---

- Batch normalization: +2% mAP
- High-resolution classifier: +4% mAP
- Convolutional with Anchor Boxes
- Direct location prediction
- Fine-grained features
- Multi-scale training

## YOLOv2: Convolutional with Anchor Boxes

---



Instead of predicting 5 arbitrary boundary boxes, we predict offsets to each of the anchor boxes above.

## YOLOv3: An Incremental Improvement

---

- Class prediction: Replace softmax function with independent logistic classifiers (sum of output can be greater than 1 now).  $\implies$  Use binary cross-entropy loss for each label  $\implies$  Reduce computation complexity.
- Feature Pyramid Network (FPN)

# YOLOv3: Performance

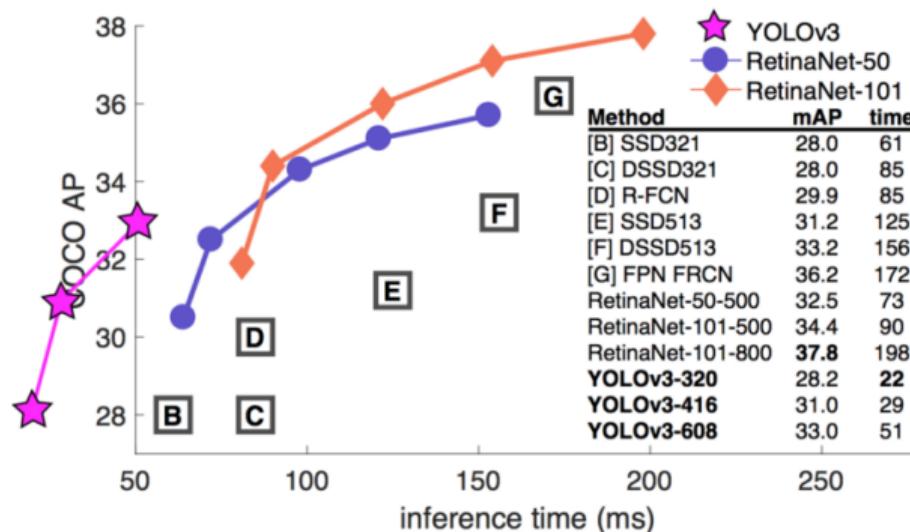
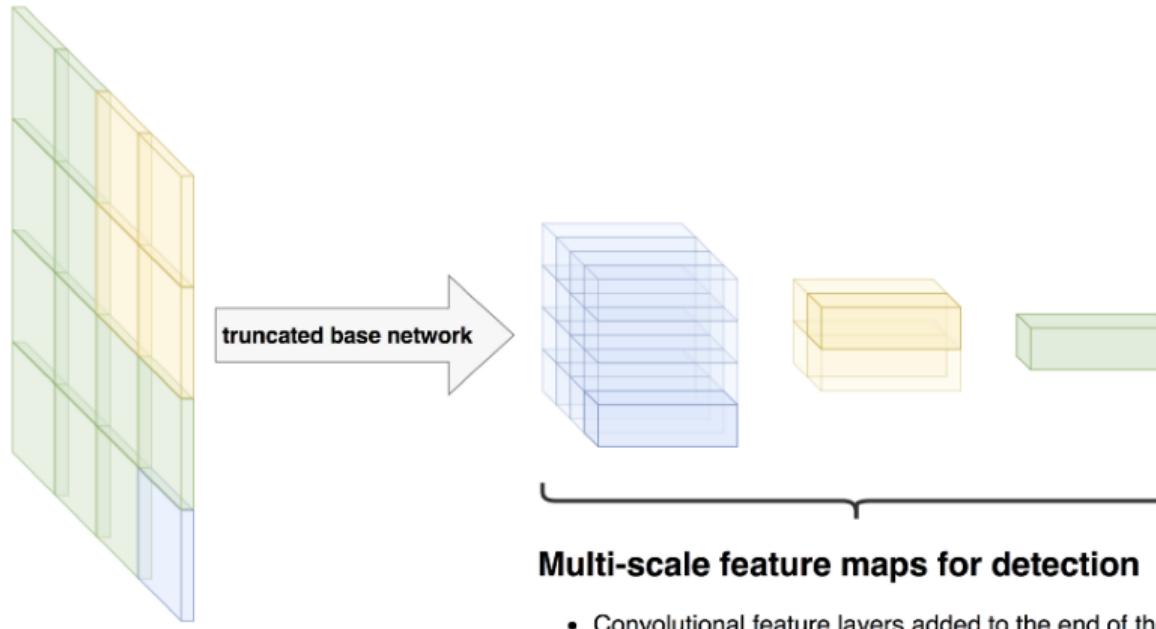


Figure 31: YOLOv3 Performance<sup>21</sup>

<sup>21</sup>Joseph Redmon and Ali Farhadi. "YOLOv3: An Incremental Improvement". In: *CoRR* abs/1804.02767 (2018).

# Single Shot Multibox Detector (SSD)



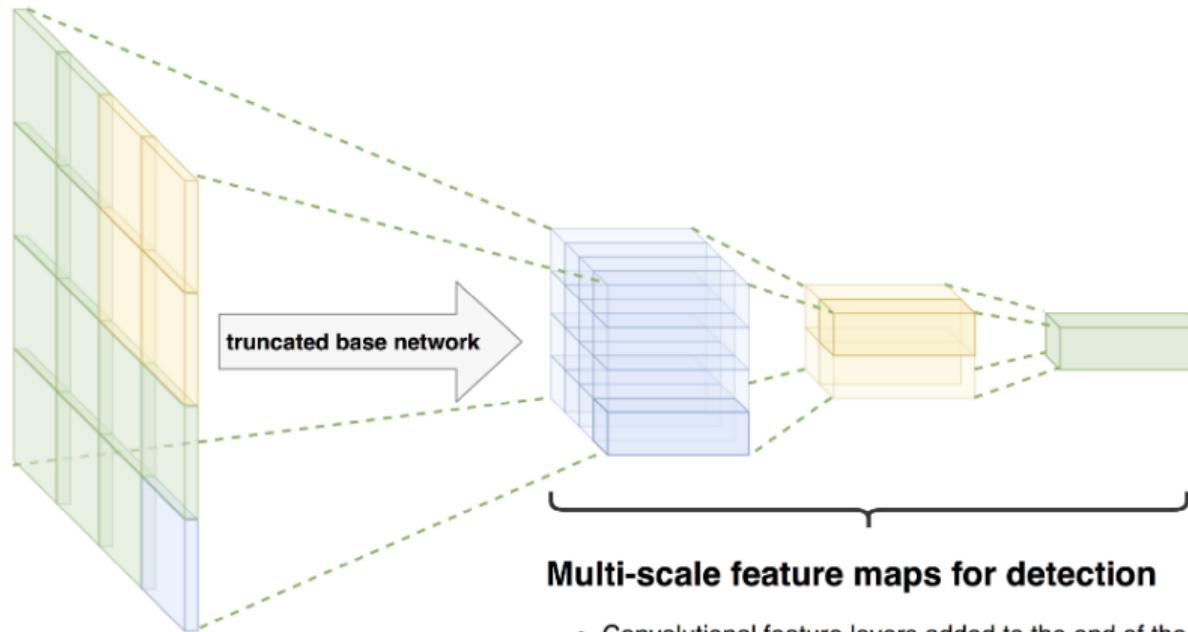
**Input image**

Shape i.e. [3, 224, 224]

## Multi-scale feature maps for detection

- Convolutional feature layers added to the end of the truncated base network
- Decrease in size progressively and allow predictions of detections at multiple scales
- Shapes in this example: [256, 4, 4], [256, 2, 2], and [256, 1, 1]

# Single Shot Multibox Detector (SSD)



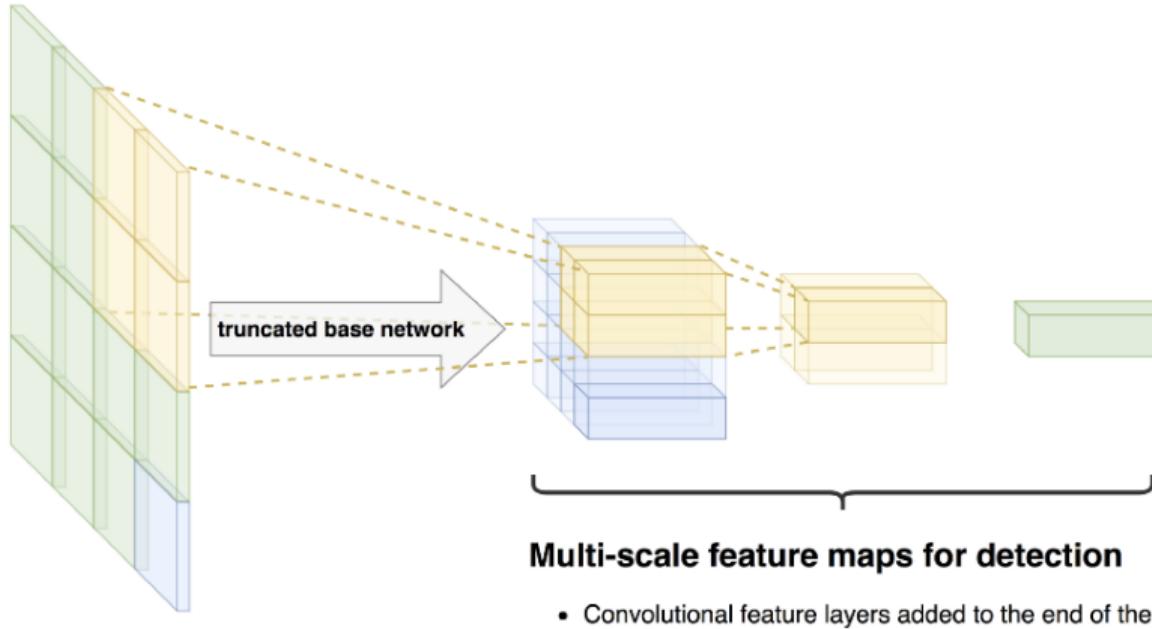
**Input image**

Shape i.e. [3, 224, 224]

## Multi-scale feature maps for detection

- Convolutional feature layers added to the end of the truncated base network
- Decrease in size progressively and allow predictions of detections at multiple scales
- Shapes in this example: [256, 4, 4], [256, 2, 2], and [256, 1, 1]

# Single Shot Multibox Detector (SSD)



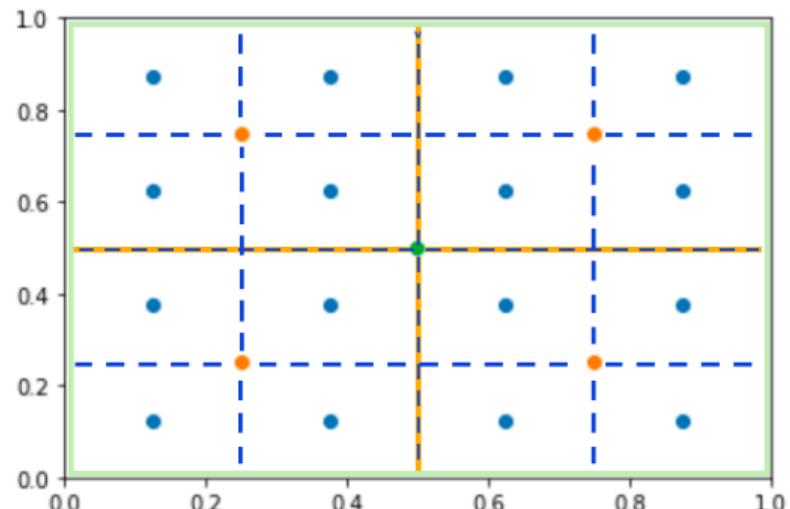
**Input image**

Shape i.e. [3, 224, 224]

## Multi-scale feature maps for detection

- Convolutional feature layers added to the end of the truncated base network
- Decrease in size progressively and allow predictions of detections at multiple scales
- Shapes in this example: [256, 4, 4], [256, 2, 2], and [256, 1, 1]

# Single Shot Multibox Detector (SSD)

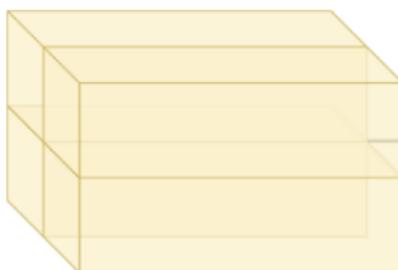


- Green dots:  $4 \times 4$  boxes
- Orange dots:  $2 \times 2$  boxes
- Blue dots:  $1 \times 1$  boxes

# Single Shot Multibox Detecto (SSD)

## Feature map for the 2x2 grid of default boxes

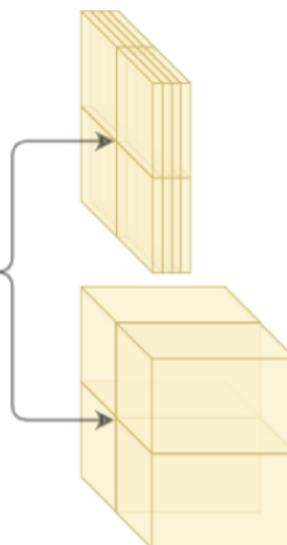
Shape i.e. [256, 2, 2]



## Convolutional predictors for detection

Predicted bounding box coordinates (offsets to the respective default box coordinates)

Shape [4, 2, 2]



Class probabilities for every default box in the 2x2 grid

Shape [ $n_{\text{classes}}$ , 2, 2]

Fed every one of the three feature maps is input to two more convolutional layers.

## Trade off between speed and accuracy

Let's summarize something:

- We defined several grids of differently sized default boxes that will allow us to detect objects at different scales in one single forward pass. Many previous architectures detected objects at different scales for example by passing differently sized versions of the same picture to the detector which is computationally more expensive.
  - For each default box in every grid, the network outputs  $n$  class probabilities and 4 offsets to the respective default box coordinates which give the predicted bounding box coordinates.
- ⇒ But still less accurate than two-stage methods.

## Trade off between Speed and Accuracy

---

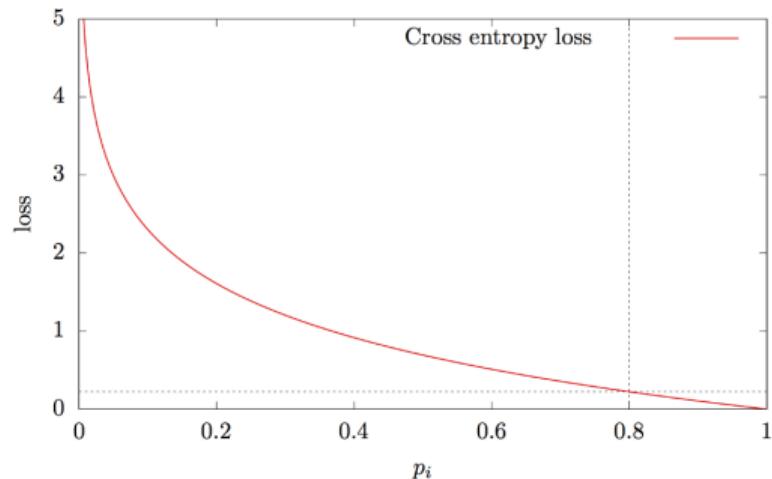
- The advantage of two-stage methods try to show off a **few region of interest** - ignore background.
- Meanwhile single-stage methods (like YOLO, SSD) evaluate roughly between ten to hundred thousand candidate locations (way more than the  $4 \times 4$ ,  $2 \times 2 + 1$  default boxes in our example here) and of course most of these boxes do not contain an object.  
     $\Rightarrow$  Imbalance class

# Trade off between Speed and Accuracy

- Cross-entropy loss function, where  $i$  is the index of class,  $y_i$  is the label (1 or 0), and  $p_i$  is the predicted probability that object belongs to class  $i$ .

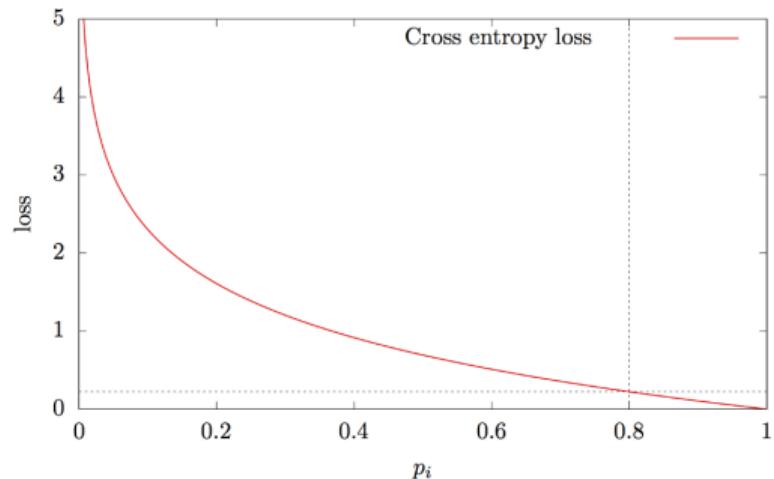
$$C(p, y) = - \sum_i y_i \log p_i$$

- Lets say a box contains background and the network is 80% sure that it actually is only background. In this case  $y_{background} = 1$ , all other  $y_i$  are 0 and  $p_{background} = 0.8$ .



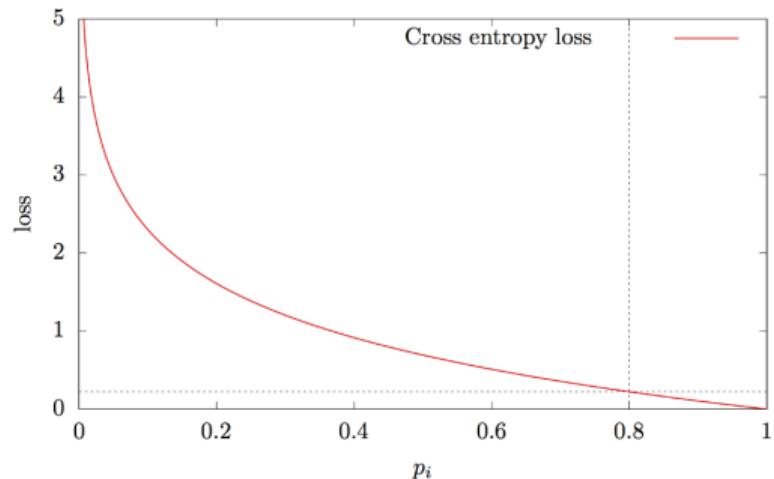
# Trade off between Speed and Accuracy

- We can see that at 80% certainty that the box contains only background, the loss is still  $\sim 0.22$ .
- Imagine: We have ten actual objects in the image and the network is not really sure to which class they belong so that their loss is i.e.  $\sim 3$ . This would give us around 30



## Trade off between Speed and Accuracy

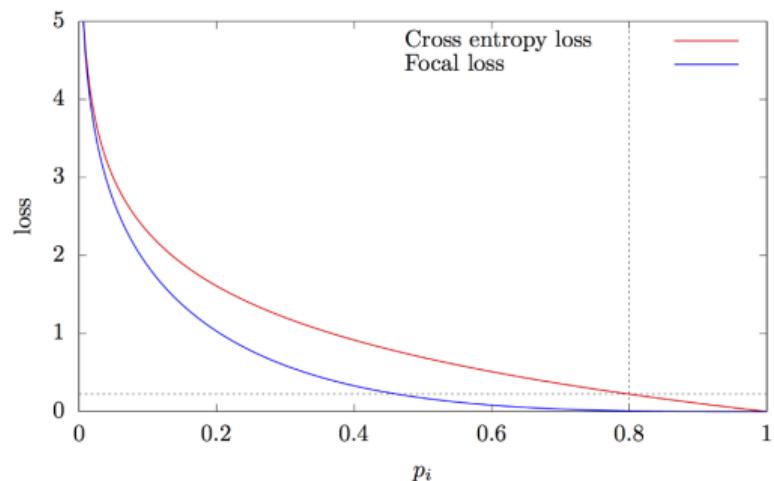
- All the other  $\sim 10000$  default boxes are background and the network is 80% sure that they are just background.  
⇒ This give us a loss of around:  
 $20 \times 0.22 = 2200$



## Focal Loss<sup>22</sup>

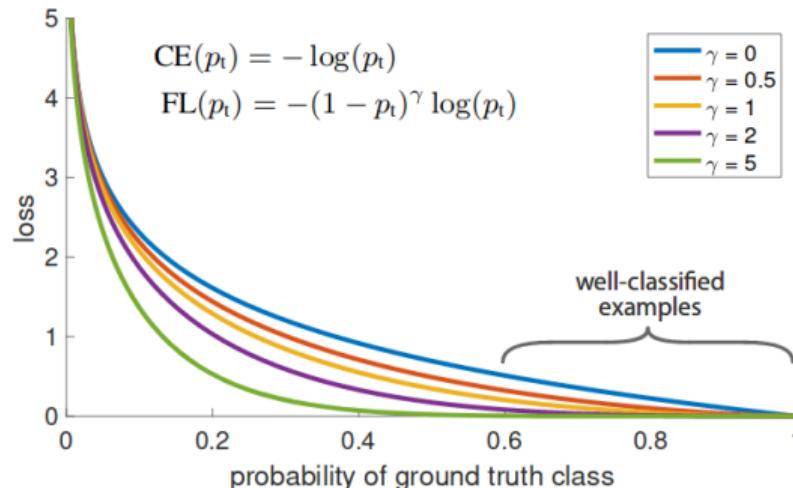
- Lin et al. (2017) propose a loss function so that the learning can focus on the few interesting cases. The authors called their loss function **Focal loss** and their architecture **RetinaNet**.

$$C(p, y) = - \sum_i y_i (1 - p_i)^\gamma \log p_i$$



<sup>22</sup>Tsung-Yi Lin et al. "Focal Loss for Dense Object Detection". In: *IEEE transactions on pattern analysis and machine intelligence* (2018).

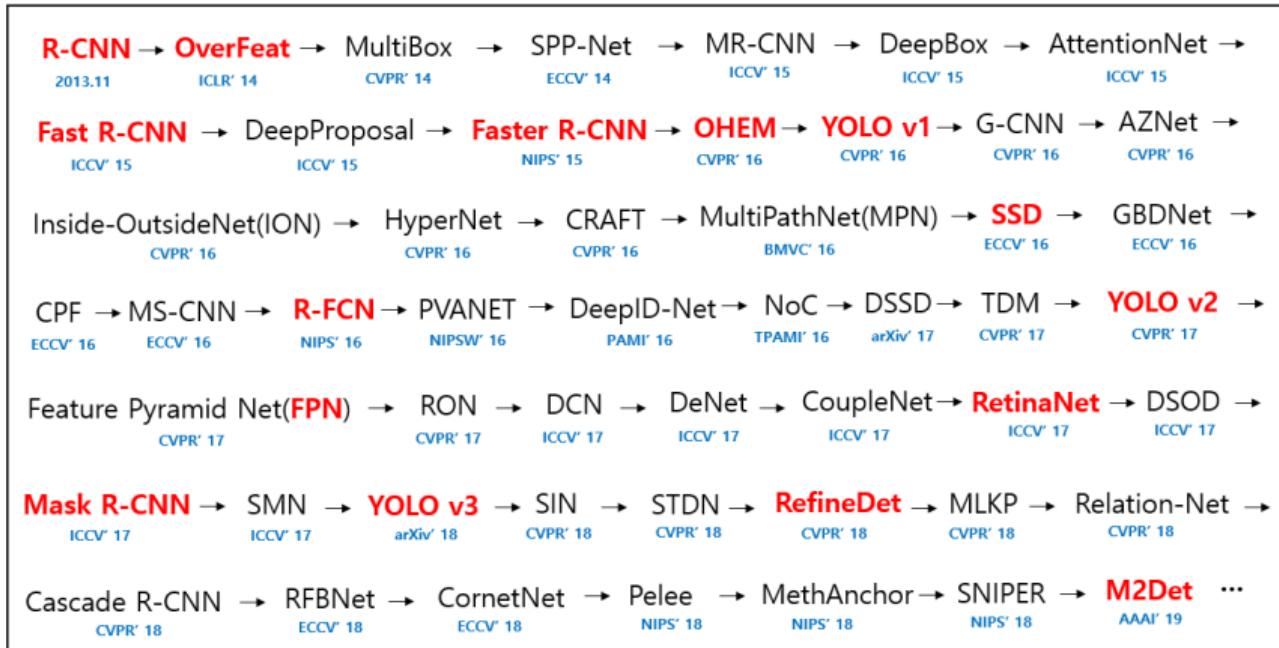
## Focal Loss



**Figure 32:** Focal loss with some  $\gamma$  factor<sup>23</sup>

<sup>23</sup>Tsung-Yi Lin et al. "Focal Loss for Dense Object Detection". In: *IEEE transactions on pattern analysis and machine intelligence* (2018).

## Remark



**Figure 33:** Deep Learning methods for Object detection history [Lee Hoseong 2019]

# Q&A

## References

---

-  Le Cun et al. "Handwritten Digit Recognition with a Back-Propagation Network". In: *Advances in Neural Information Processing Systems*. Morgan Kaufmann, 1990, pp. 396–404.
-  Ross B. Girshick. "Fast R-CNN". In: *2015 IEEE International Conference on Computer Vision (ICCV)* (2015), pp. 1440–1448.
-  Ross Girshick et al. "Rich feature hierarchies for accurate object detection and semantic segmentation". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2014.

## References ii

-  Kaiming He et al. "Deep Residual Learning for Image Recognition". In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 770–778.
-  Kaiming He et al. "Mask R-CNN". In: *2017 IEEE International Conference on Computer Vision (ICCV)* (2017), pp. 2980–2988.
-  Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks". In: *Advances in Neural Information Processing Systems*. P. 2012.
-  Yann Lecun et al. "Gradient-based learning applied to document recognition". In: *Proceedings of the IEEE*. 1998, pp. 2278–2324.
-  Min Lin, Qiang Chen, and Shuicheng Yan. *Network in Network*. 2013.

## References iii

-  Tsung-Yi Lin et al. “Focal Loss for Dense Object Detection”. In: *IEEE transactions on pattern analysis and machine intelligence* (2018).
-  Aäron van den Oord et al. “WaveNet: A Generative Model for Raw Audio”. In: *SSW*. 2016.
-  Joseph Redmon and Ali Farhadi. “YOLOv3: An Incremental Improvement”. In: *CoRR* abs/1804.02767 (2018).
-  Joseph Redmon et al. “You Only Look Once: Unified, Real-Time Object Detection”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 779–788.

## References iv

-  Evan Shelhamer, Jonathan Long, and Trevor Darrell. "Fully Convolutional Networks for Semantic Segmentation". In: *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 3431–3440.
-  Karen Simonyan and Andrew Zisserman. *Very deep convolutional networks for large-scale image recognition*. 2014.
-  Christian Szegedy et al. *Going deeper with convolutions*. 2014.
-  Ashish Vaswani et al. "Attention Is All You Need". In: *NIPS*. 2017.
-  Yonghui Wu et al. "Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation". In: *CoRR* abs/1609.08144 (2016).