

Famous architectures



András Horváth

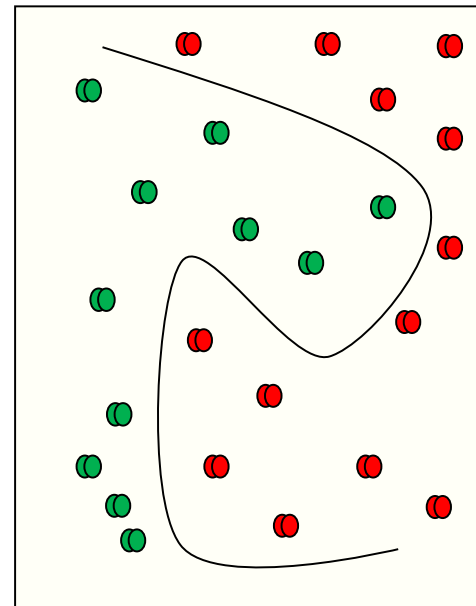
Budapest, 2018.12.03



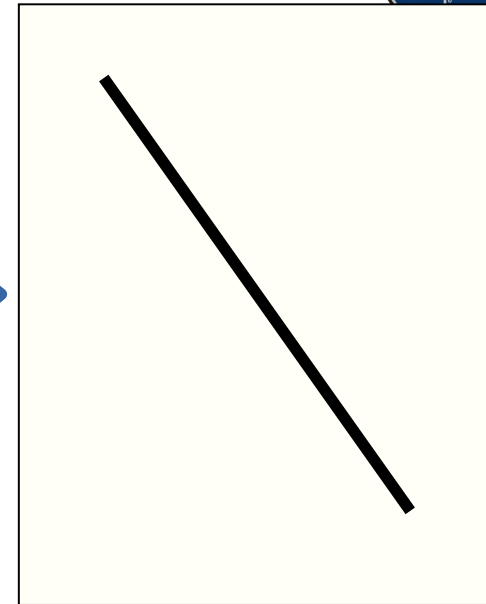
Neural Networks

- Classification - decision
- FNN, SVM – linear classification
- Is X larger than a limit? $X > k$?
- Finding a good feature representation:
 - Meaningful
 - Sparse - low dimensions
 - Ensures easy separation

Finding the representation with the help of machine learning



Input space



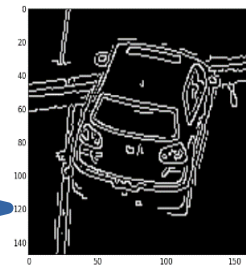
Feature space



Input Image

Convolution
Kernel

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



Feature Image

Convolutional neural networks

- A network of simple processing elements

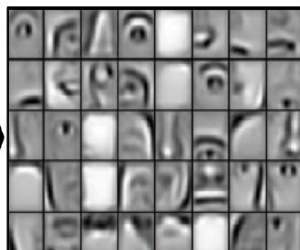
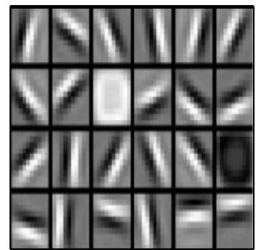
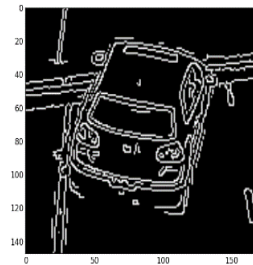
- Elements:

- Convolution



Convolution
Kernel

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



Low layers

Middle layers

High layers

- ReLU



Thresholding
all values
below zero

- Pooling

1	0	2	3
4	6	6	8
3	1	1	0
1	2	2	4



6	8
3	4

Selection of
the
maximal
response in
an area

Convolutional networks

Ok, but how many layers do we need?

How many features should be in each layer?

What should be the network architecture?

Convolutional networks

Ok, but how many layers do we need?

How many features should be in each layer?

What should be the network architecture?

These are called hyper-parameters:

Along with: non-linearity type, batch-norm, dropout etc.

Convolutional networks

Ok, but how many layers do we need?

How many features should be in each layer?

What should be the network architecture?

These are called hyper-parameters:

Along with: non-linearity type, batch-norm, dropout etc.

We can use a network which performed fairly well on an other dataset

It will probably work well on our task too

Large scale visual recognition challenge



Large scale visual recognition challenge

ImageNet:

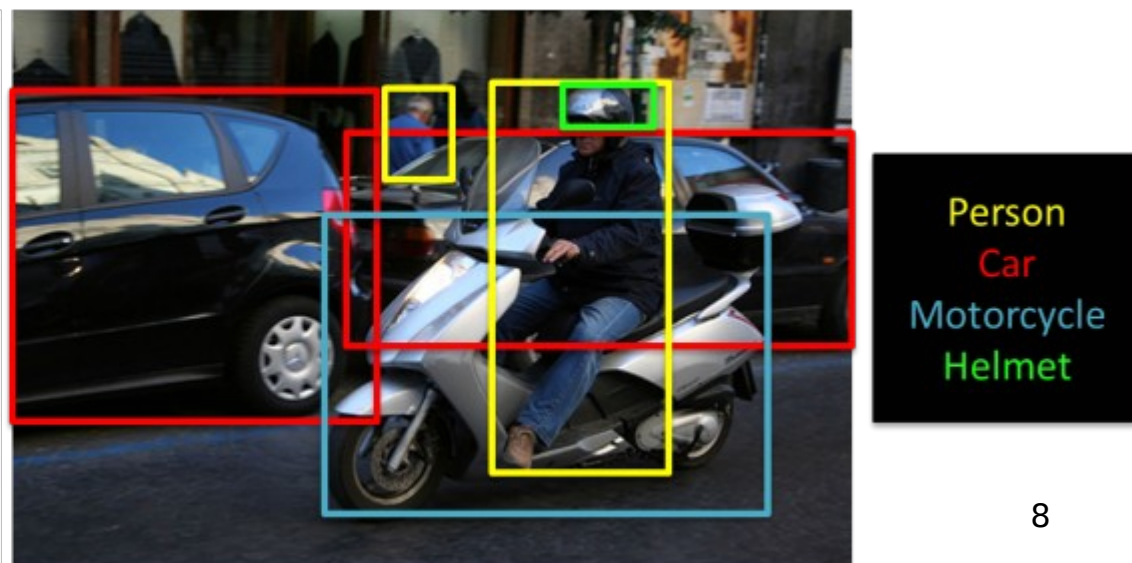
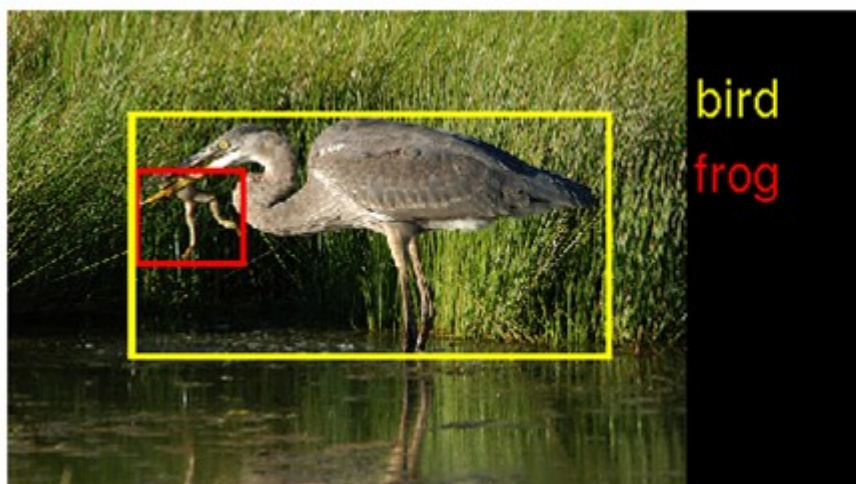
Over 15 million images, more than 22k categories

Detection and Classification



Detection for 200 fully labeled categories. 0.73 % mAP

Classification for 1000 categories. 0.97%



Large scale visual recognition challenge

Andrej Karpathy – „the human benchmark”



Shown 13 samples from each category

94.9% accuracy



The screenshot displays the ImageNet dataset interface. On the left, a large image of a hot dog is shown. Below it, a search bar contains the text "hotdog, hot dog, red hot". To the right of the search bar, there are buttons for "Show answer" and "Show google prediction". Below the search bar, there are two rows of predicted labels: "hotdog, hot dog, red hot" and "cheeseburger". Below these labels, there are two rows of predicted labels: "hotdog, hot dog, red hot" and "cheeseburger". Below these labels, there are two rows of predicted labels: "hotdog, hot dog, red hot" and "cheeseburger".

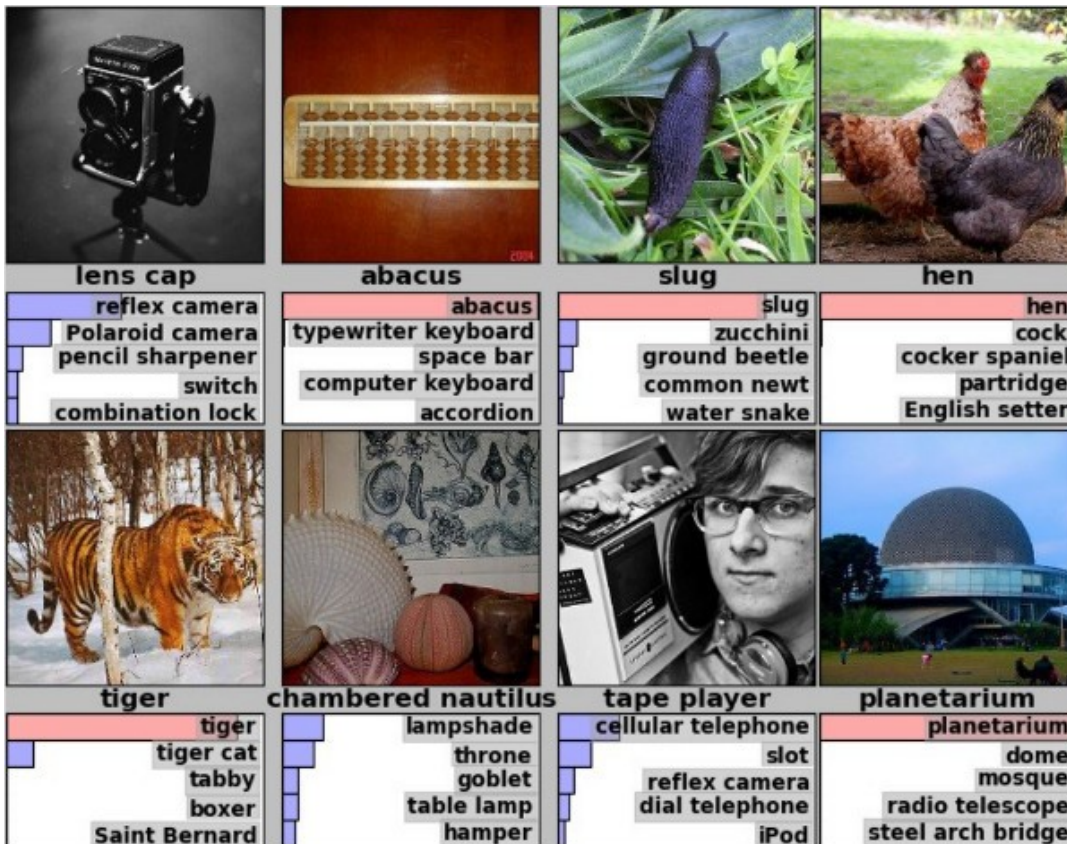
The main part of the interface shows a grid of 13 samples from each category. The categories are listed on the left: consomme, snack food, sandwich, hotdog, hot dog, red hot, hamburger, beefburger, burger, cheeseburger, course, entree, main course, plate, and dessert, sweet, afters, frozen dessert. Each category has a corresponding row of 13 sample images.

Large scale visual recognition challenge

Andrej Karpathy



Shown 13 samples from each category



Large scale visual recognition challenge

Andrej Karpathy



Shown 14 samples from each category

Can be tried out online:

<https://cs.stanford.edu/people/karpathy/ilstvrc/>



Alexnet

Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton (2012)

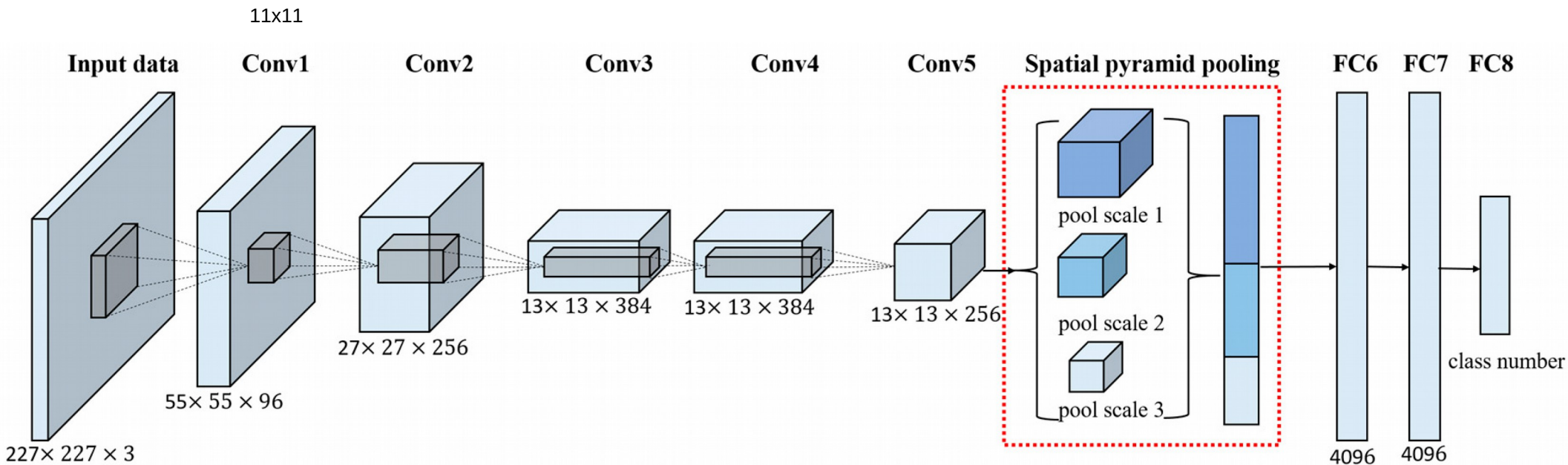
Trained whole ImageNet (15 million, 22,000 categories)

Used data augmentation (image translations, horizontal reflections, and patch extractions)

Used ReLU for the nonlinearity functions (Decreased training time compared to tanh) - Trained on two GTX 580 GPUs for six days

Dropout layers

2012 marked the first year where a CNN was used to achieve a top 5 test error rate of 15.4% (next best entry was with error of 26.2%)



VGG - 16/19

Karen Simonyan and Andrew Zisserman of the University of Oxford, 2014 Visual Geometry Group

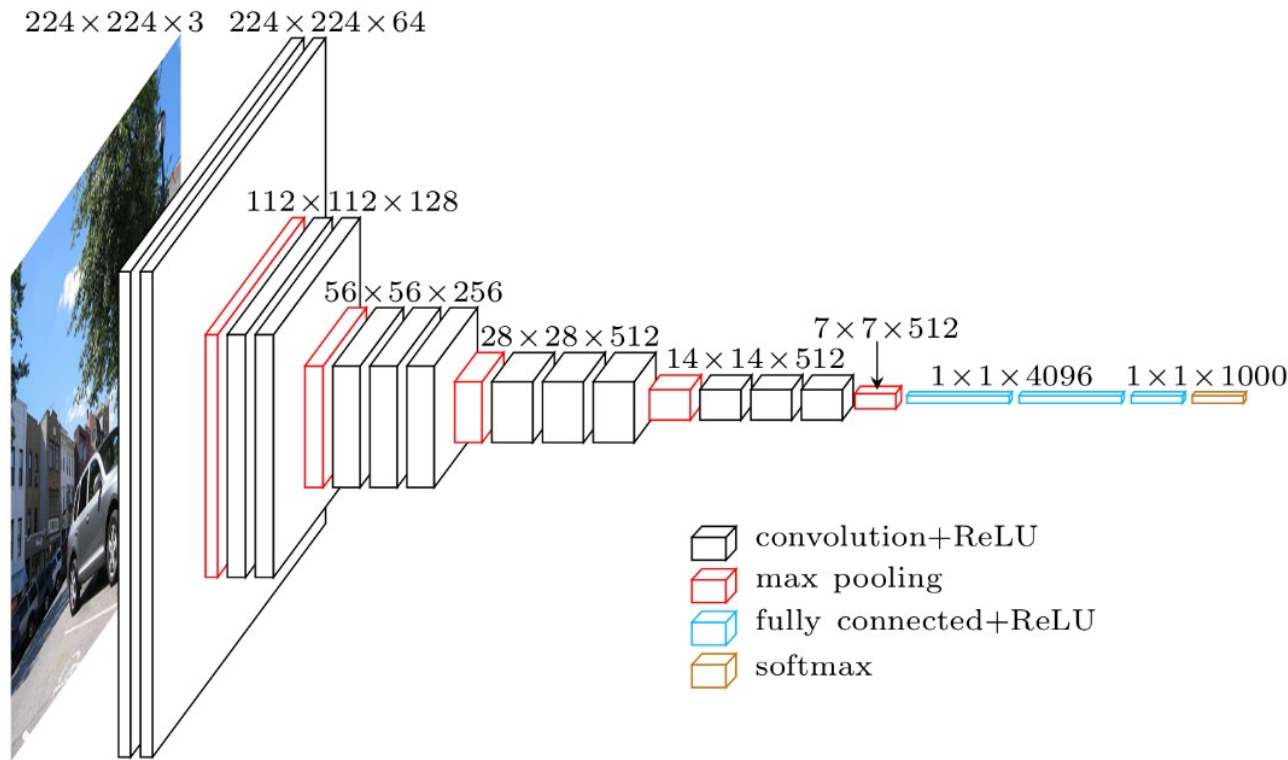
As the spatial size of the input volumes at each layer decrease (result of the conv and pool layers), the depth of the volumes increase due to the increased number of filters as you go down the network.

Shrinking spatial dimensions but growing depth

3x3 filters with stride and pad of 1, along with 2x2 maxpooling layers with stride 2

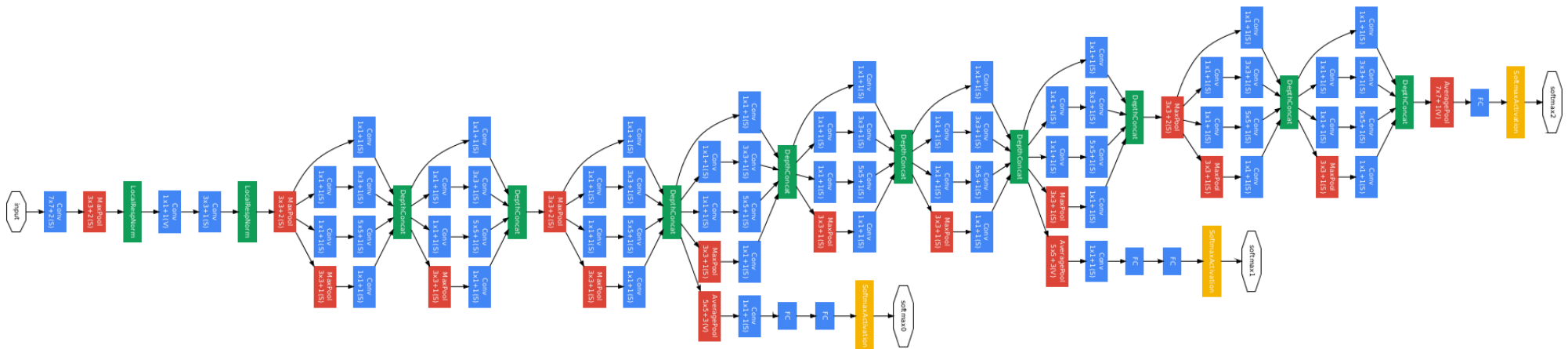
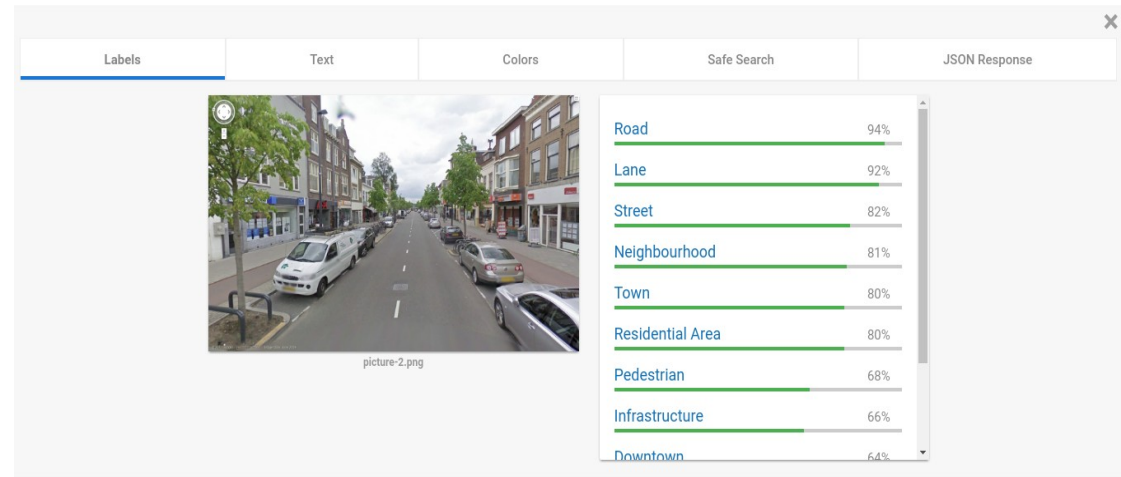
7.3% error rate

Simple architecture, still the
swiss knife of deep learning



Google - Inception architecture

- GoogLeNet:
- 22/42 layers (9 inception_v3 layers)
- 5 million free parameters
- ~1.5B operations/evaluations
- Demo: <https://cloud.google.com/vision/>

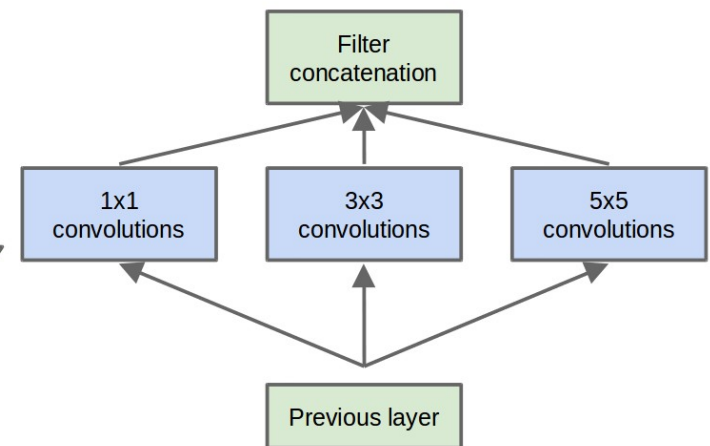
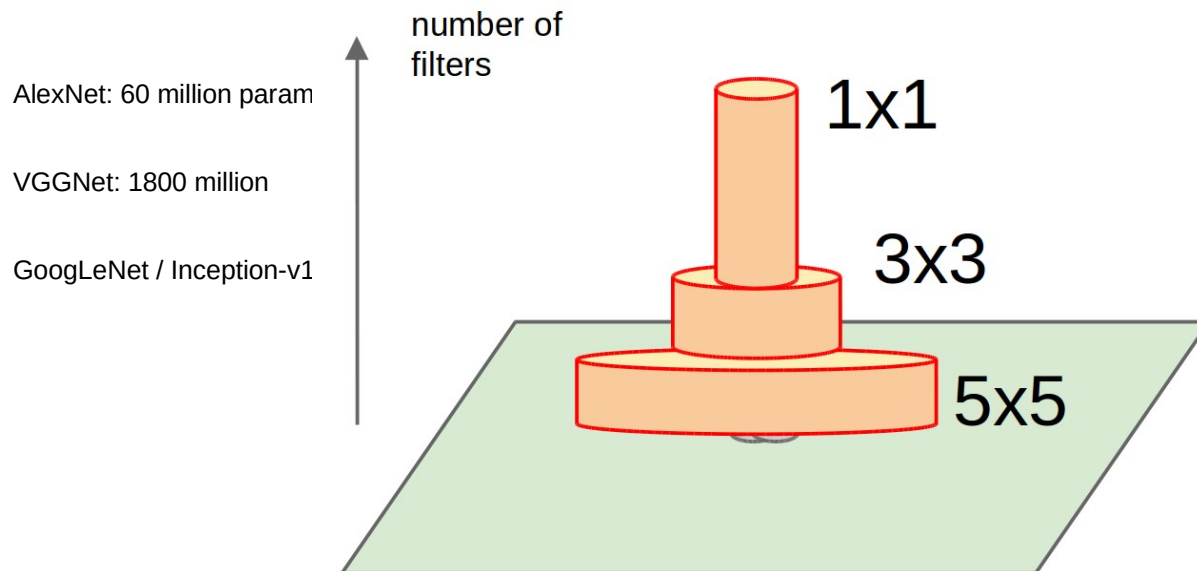


Inception

Google, Christian Szegedy

2014 with a top 5 error rate of 6.7%

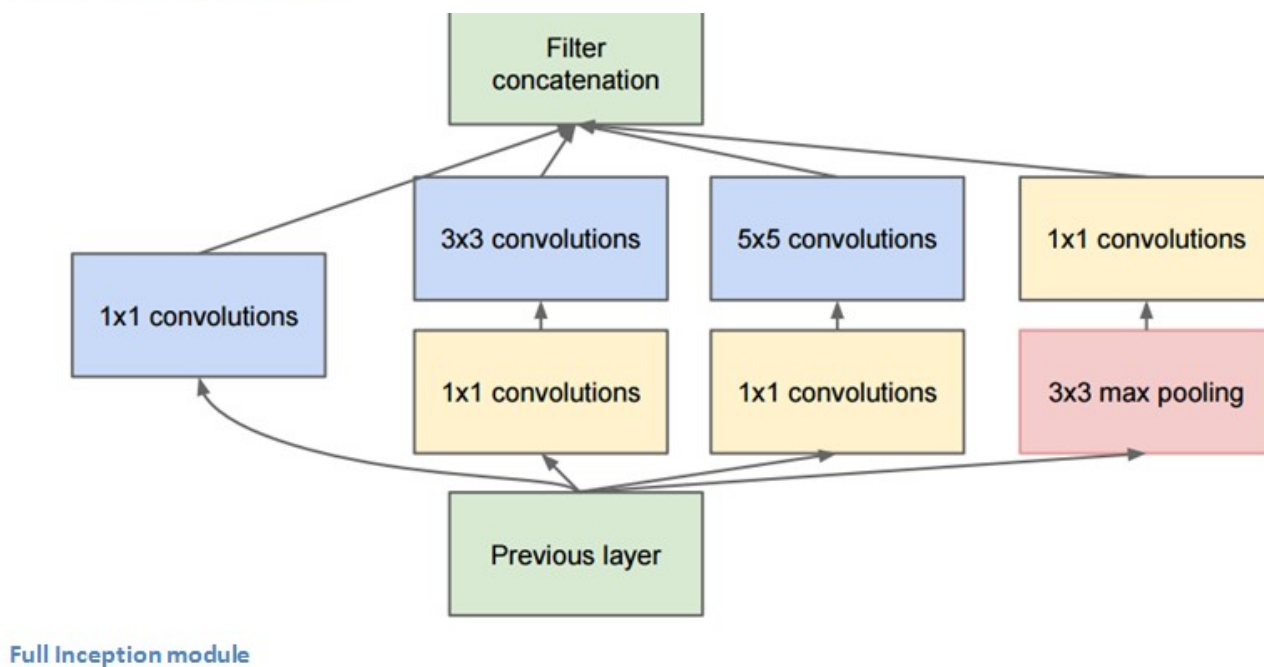
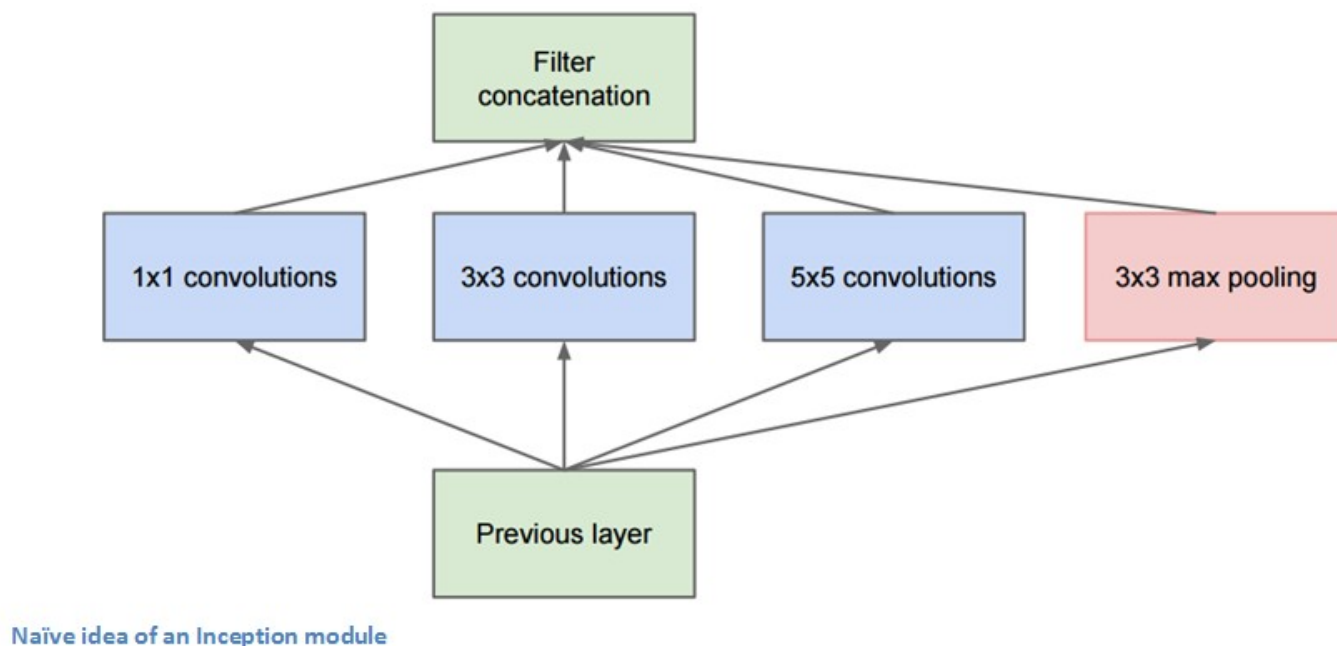
This can be thought of as a “pooling of features” because we are reducing the depth of the volume, similar to how we reduce the dimensions of height and width with normal maxpooling layers.





Rethinking Inception

Squeezing the number of channels for each kernel



Rethinking Inception

Larger convolutions were substituted by series of 3x3 convolutions

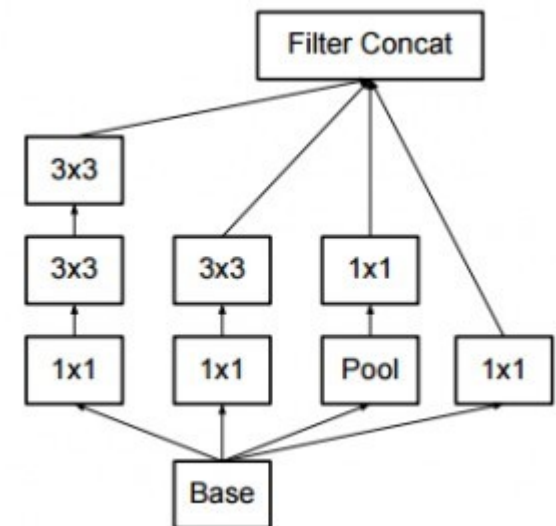
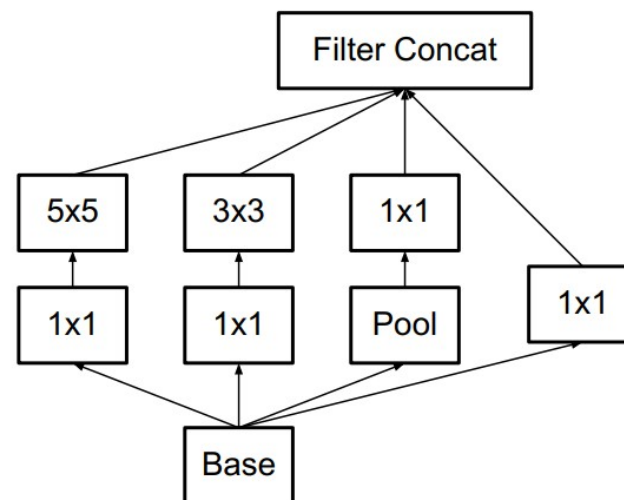
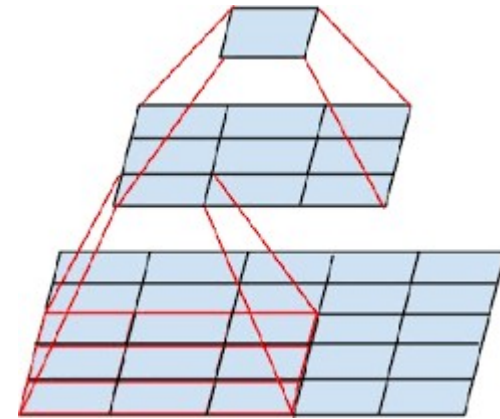


Figure 4. Original Inception module as described in [20].

Rethinking Inception

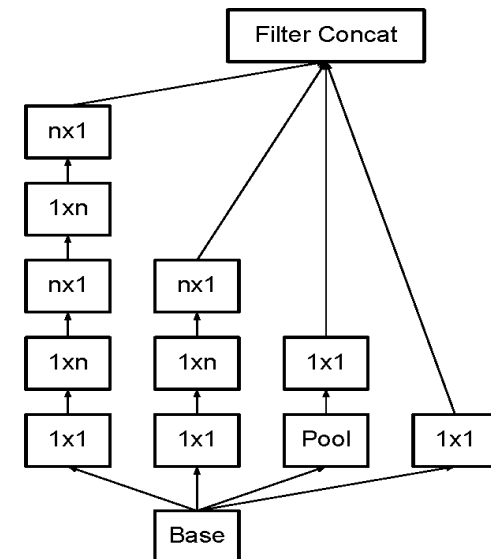
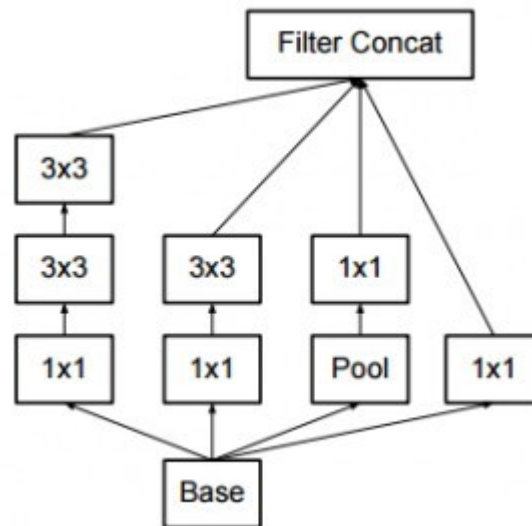
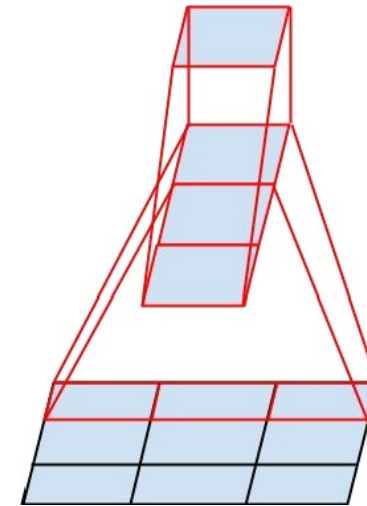
Larger convolutions were substituted by series of 3x3 convolutions

2D convolution were substituted by two 1D convolutions

AlexNet: 60 million parameters

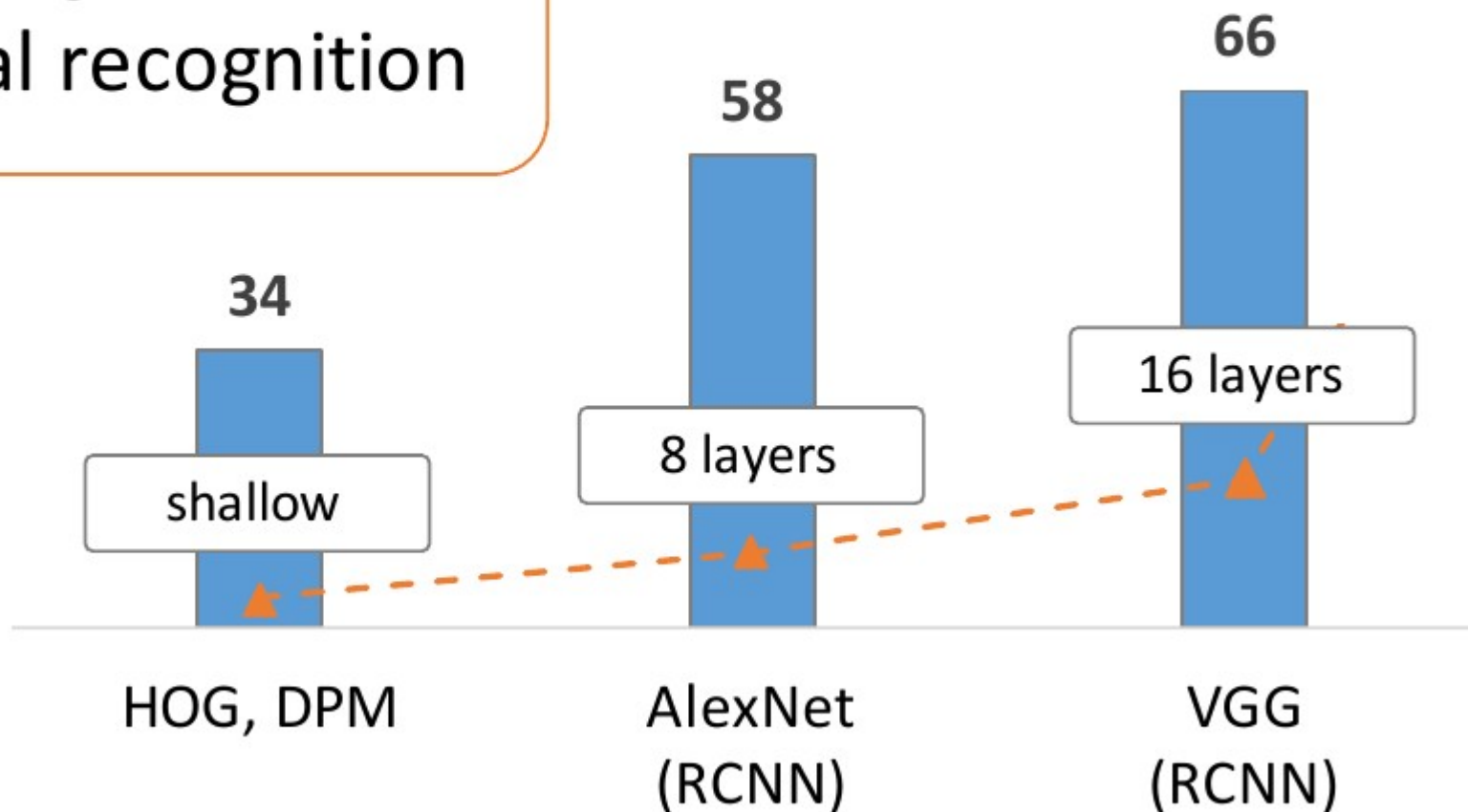
VGGNet :180 million parameters

GoogLeNet / Inception-v1: 7 million parameters



Revolution of Depth

Engines of
visual recognition



PASCAL VOC 2007 **Object Detection** mAP (%)

How deep could/should a network be?

History of network depth

Before 2012: four layers

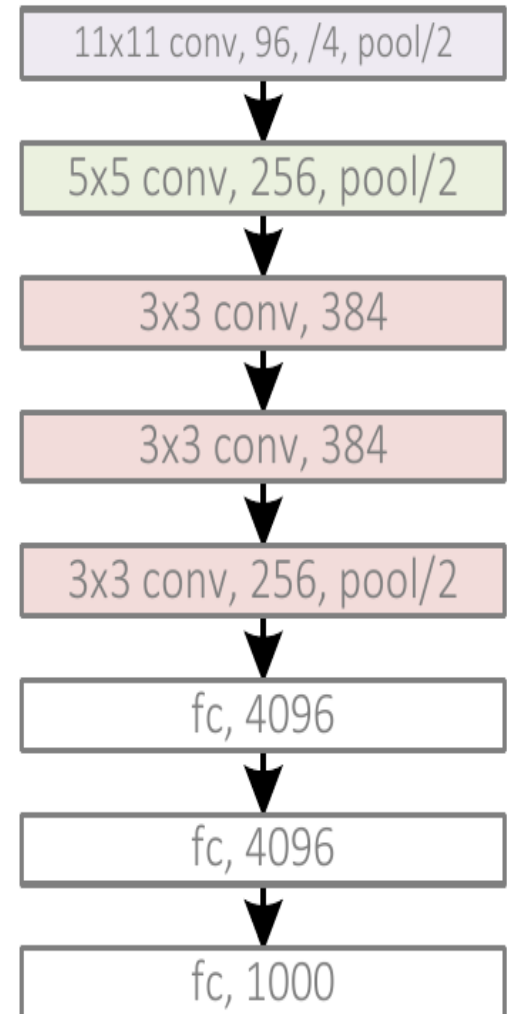
How deep could/should a network be?

History of network depth

Before 2012: four layers

2012: 8 layers

AlexNet, 8 layers
(ILSVRC 2012)





How deep could/should a network be?

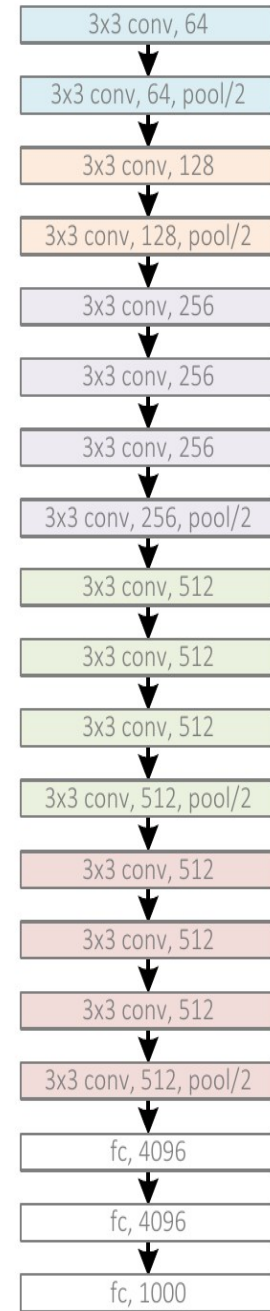
VGG, 19 layers
(ILSVRC 2014)

History of network depth

Before 2012: four layer

2012: 8 layers

2016: 22 layers



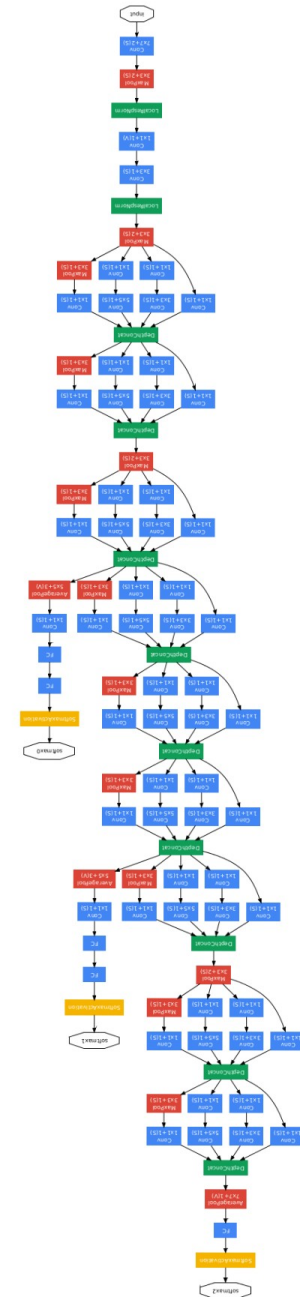
How deep could/should a network be?

History of network depth

Before 2012: four layer

2012: 8 layers

2016: 19-22 layers



How deep could/should a network be?

History of network depth

Before 2012: four layer

2012: 8 layers

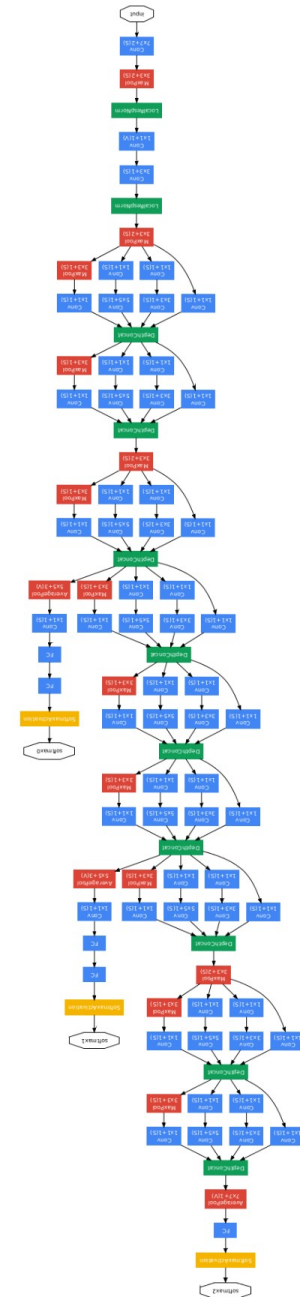
2016: 19-22 layers

Deeper network:

Possibility to approximate more complex functions



Higher number of parameters



How deep could/should a network be?

History of network depth

Before 2012: four layer

2012: 8 layers

2015: 19-22 layers

Deeper network:



Possibility to approximate more complex functions

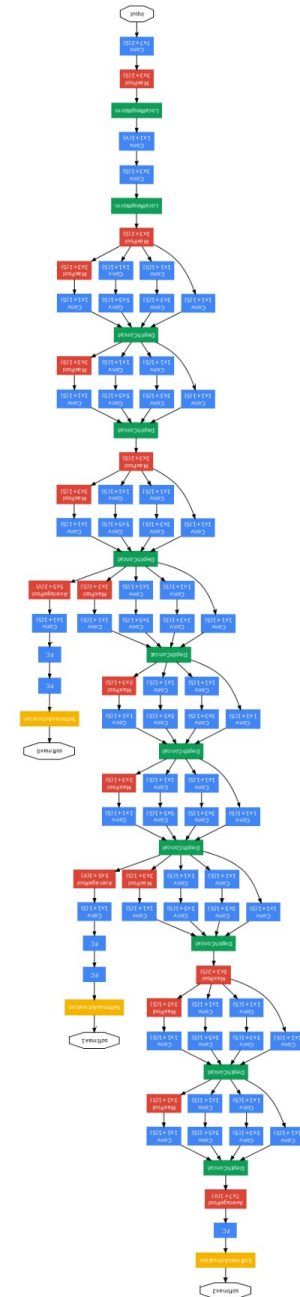


Higher number of parameters

There are no convolutional networks with more than 30 layers. Why?

The amount of transferred data is decreased from layer to layer

Training becomes difficult



Is a deeper network always better?

A deeper network would have higher approximation power

Higher number of parameters (both advantageous and disadvantageous)

Difficult to train the network

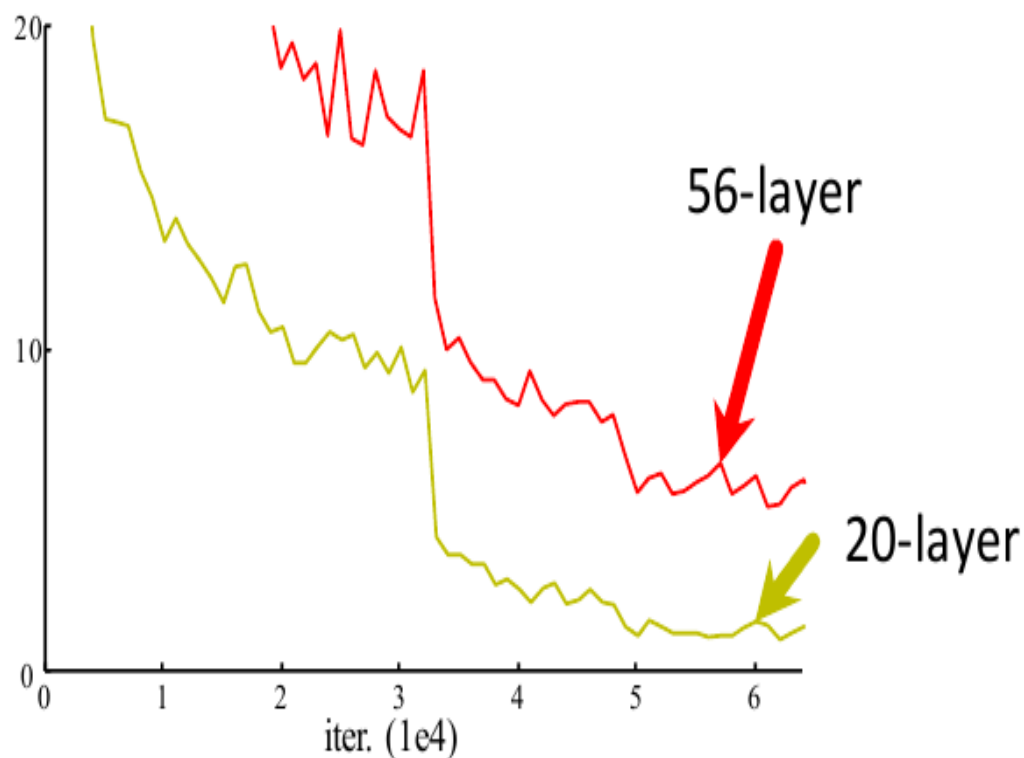
Is a deeper network always better?

A deeper network always has the potential to perform better, but training becomes difficult

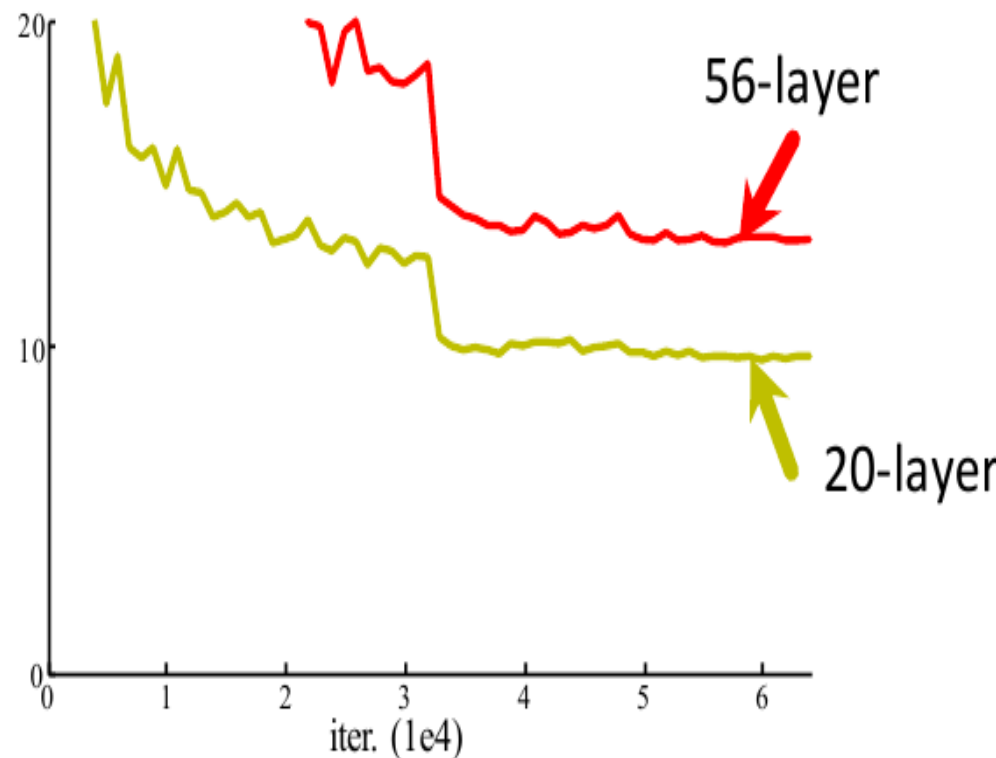
After a given depth, the same network with the same training on the same data, usually performs worse

CIFAR-10

train error (%)



test error (%)

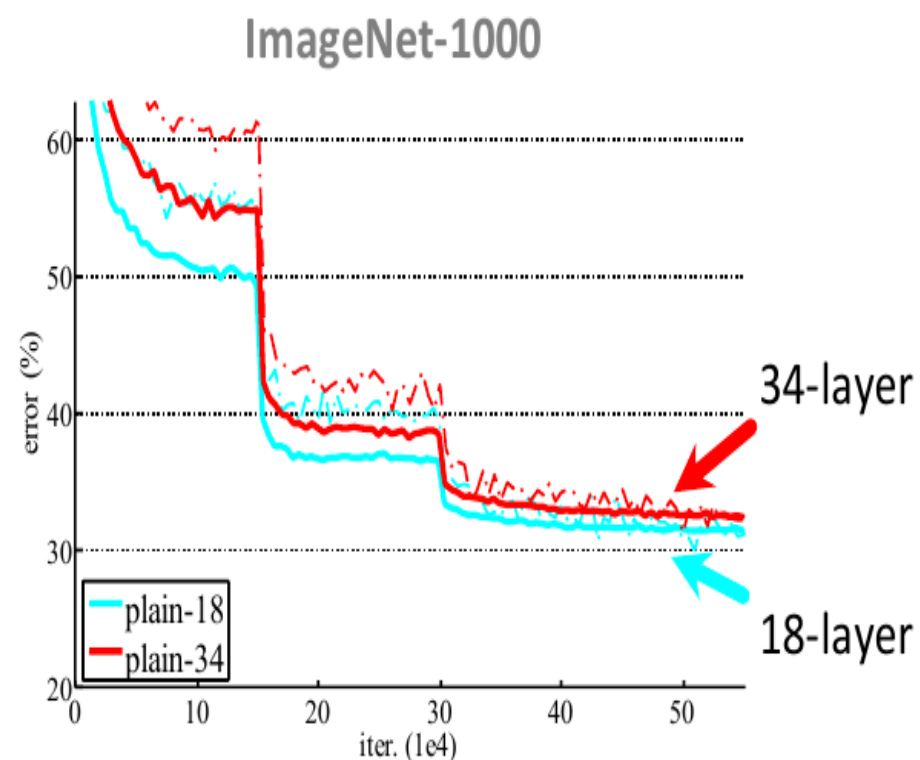
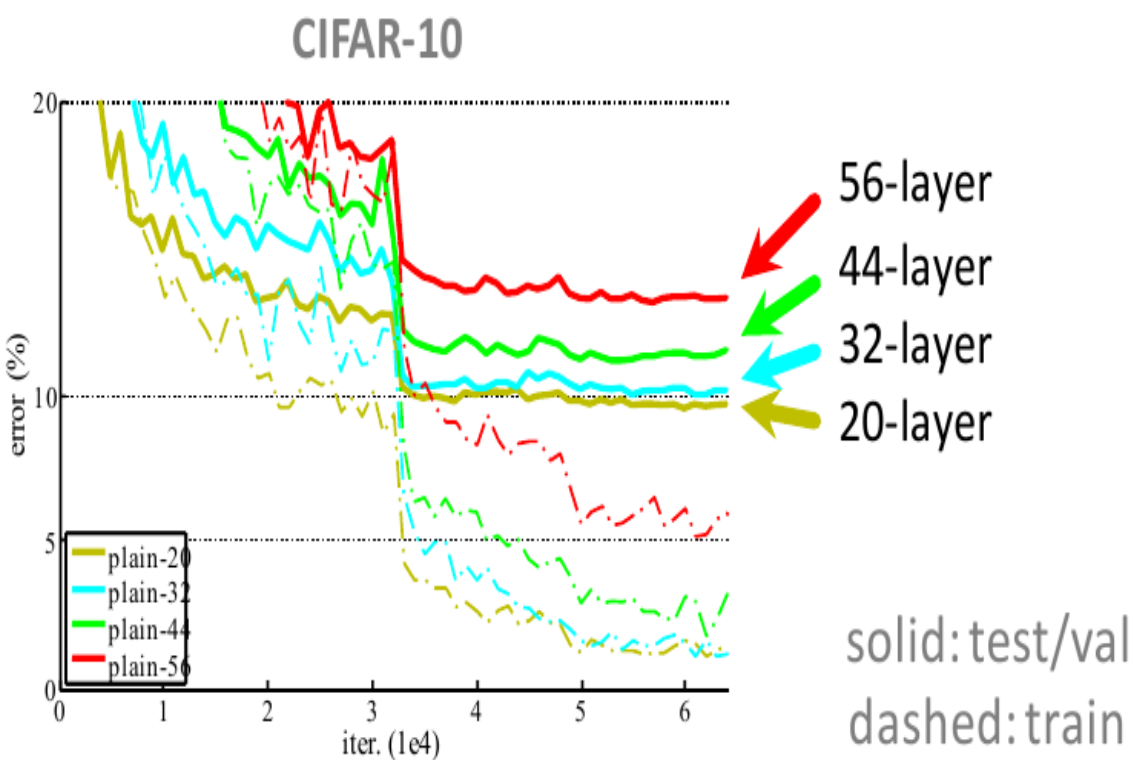


Is a deeper network always better?

A deeper network always have the potential to perform better, but training becomes difficult

We can not just simply stack convolutional layers to increase accuracy

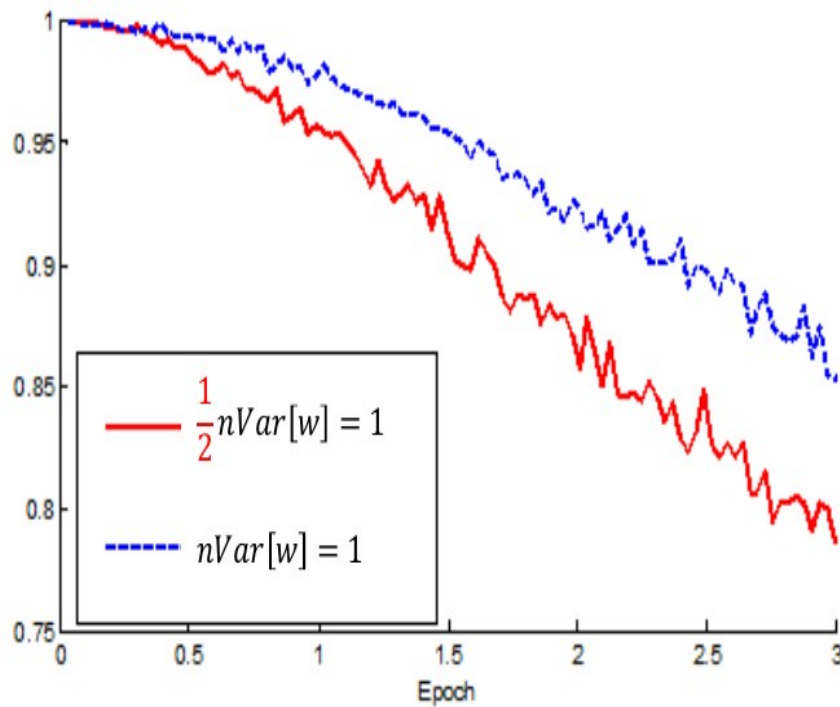
Example: stacking 3x3 convolution layers



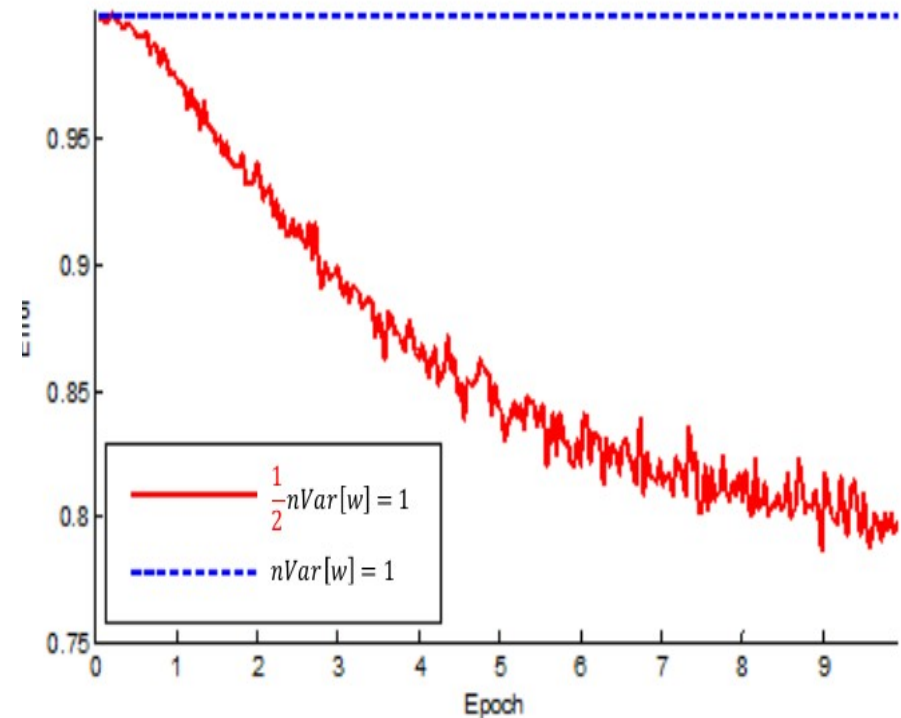
Is a deeper network always better?

A deeper network always have the potential to perform better, but training becomes difficult

22-layer ReLU net:
good init converges faster



30-layer ReLU net:
good init is able to converge



Is a deeper network always better?

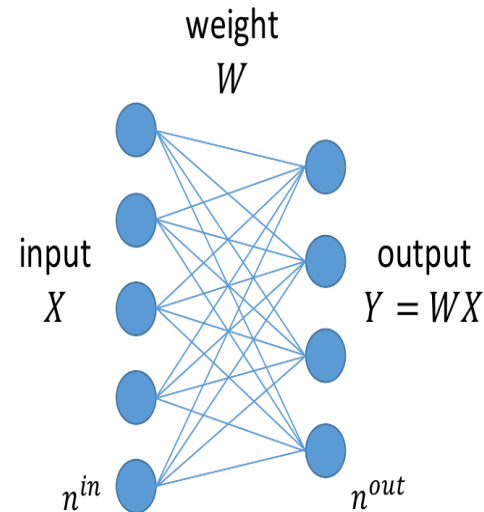
A deeper network would have higher approximation power

Higher number of parameters (both advantageous and disadvantageous)

Difficult to train the network:

Vanishing gradients

Initialization



If:

- Linear activation
- x, y, w : independent

Then:

1-layer:

$$Var[y] = (n^{in} Var[w]) Var[x]$$

Multi-layer:

$$Var[y] = \left(\prod_d n_d^{in} Var[w_d] \right) Var[x]$$

LeCun et al 1998 "Efficient Backprop"

Glorot & Bengio 2010 "Understanding the difficulty of training deep feedforward neural networks"

Is a deeper network always better?

A deeper network would have higher approximation power

Higher number of parameters (both advantageous and disadvantageous)

Difficult to train the network:

Vanishing gradients

The number of layers has an exponential impact

Initialization

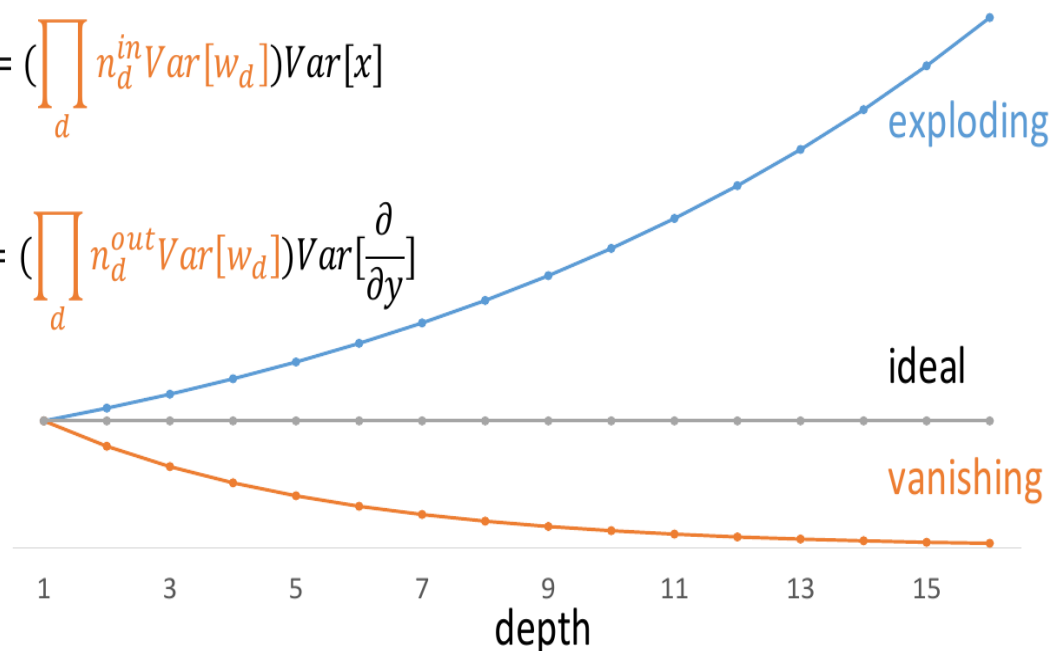
Both forward (response) and backward (gradient) signal can vanish/explode

Forward:

$$Var[y] = \left(\prod_d n_d^{in} Var[w_d] \right) Var[x]$$

Backward:

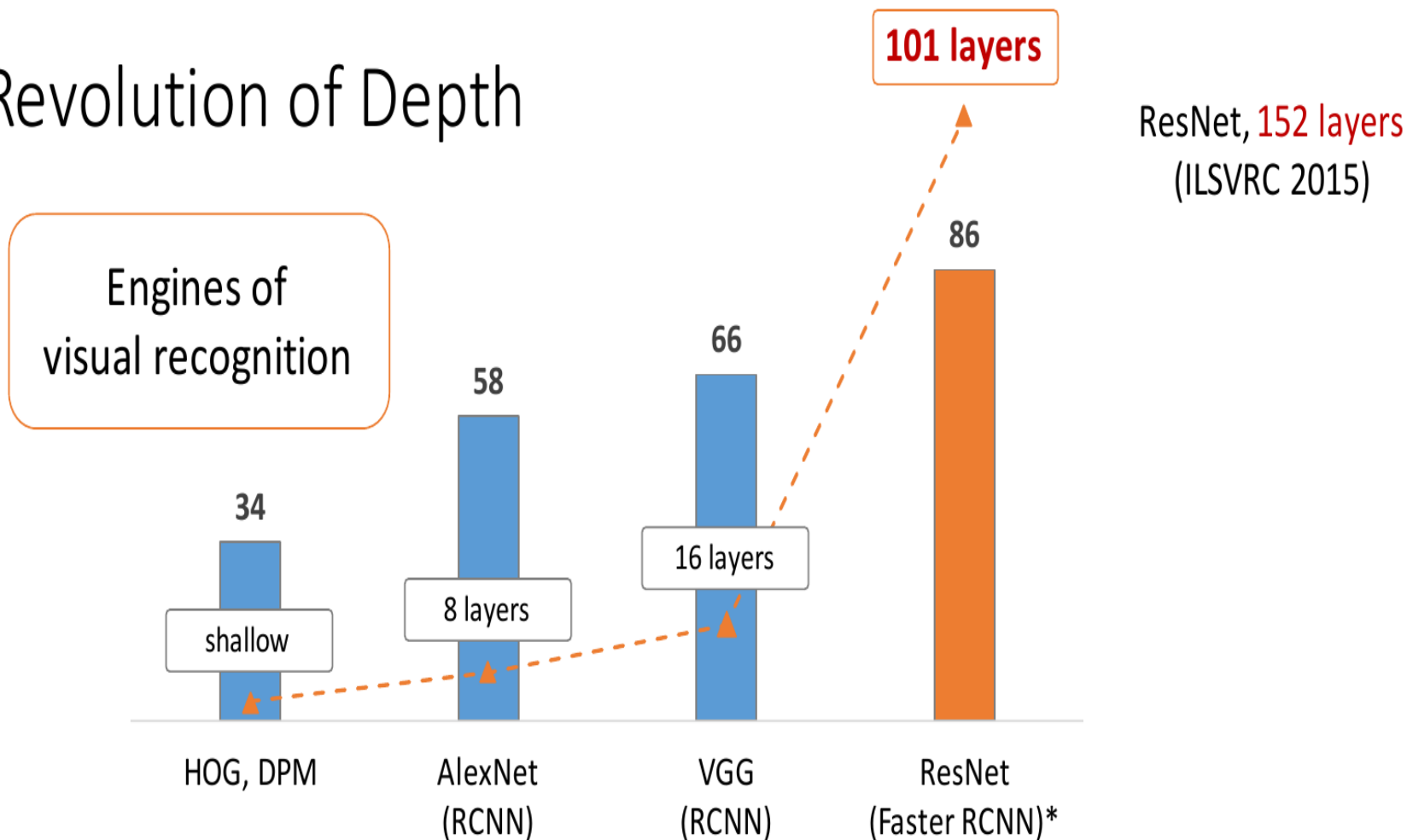
$$Var\left[\frac{\partial}{\partial x}\right] = \left(\prod_d n_d^{out} Var[w_d] \right) Var\left[\frac{\partial}{\partial y}\right]$$



How deep could a network be?

Residual networks provide an answer to these questions

Revolution of Depth



PASCAL VOC 2007 **Object Detection** mAP (%)

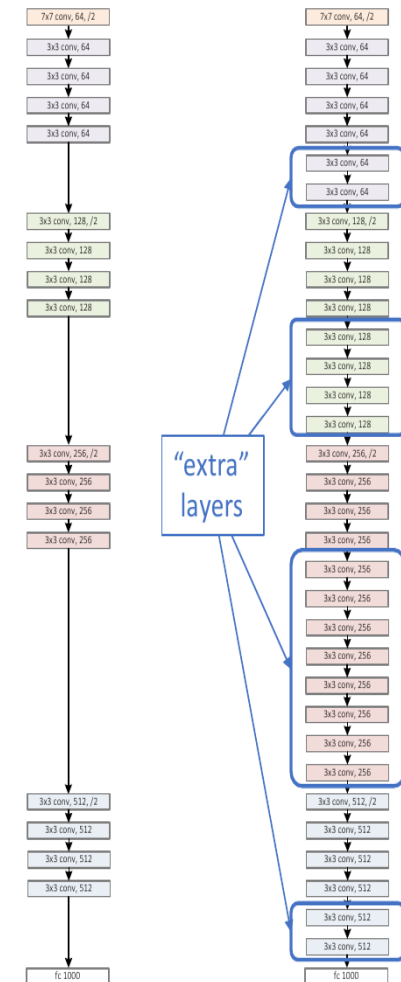
How could we create deeper networks?

A deeper network always have the potential to perform better, but training becomes difficult

How could we ensure that additional layers will not decrease accuracy (might even increase it)?

Let's start with a shallow model (18 layers) and add some extra layers (which we hope could increase accuracy)

a shallower
model
(18 layers)



How could we create deeper networks?

A deeper network always have the potential to perform better, but training becomes difficult

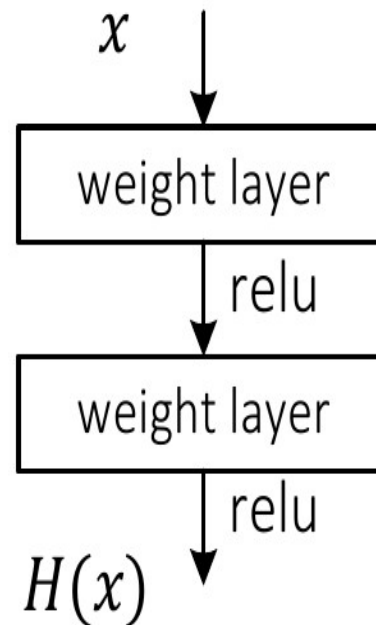
How could we ensure that additional layers will not decrease accuracy (might even increase it)?

Let's start with a shallow model (18 layers) and add some extra layers (which we hope could increase accuracy)

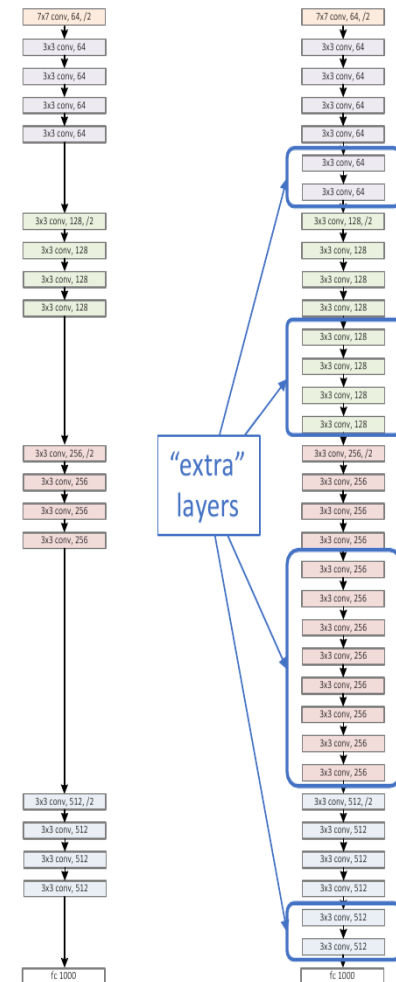
Our aim is to add “useful” operations $H(x)$

The problem is that $H(x)$ can ruin our accuracy because vanishing gradients, overfit - extra parameters

any two stacked layers



a shallower model
(18 layers)



How could we create deeper networks?

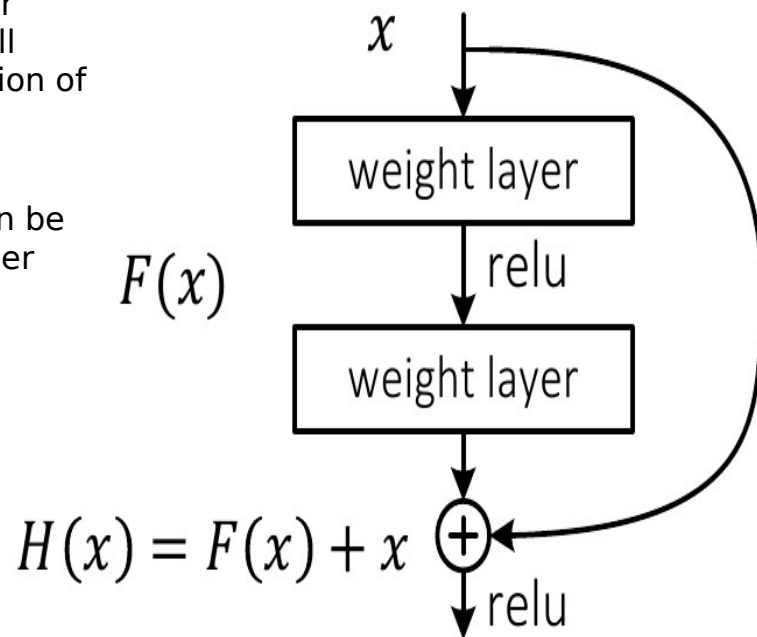
A deeper network always have the potential to perform better, but training becomes difficult

How could we ensure that additional layers will not decrease accuracy (might even increase it)?

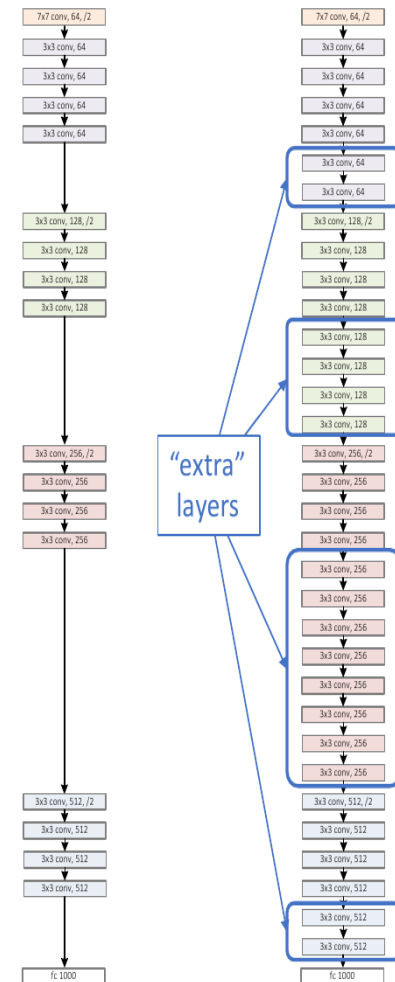
The trick is to use residual connection and as a starting point $F(x)$ could be zero, and $H(x)$ becomes the identity mapping

So $H(x)$ will not change our performance, gradients will remain, because the addition of x

Our accuracy will no be decreased, and might even be increased if we find a proper $F(x)$



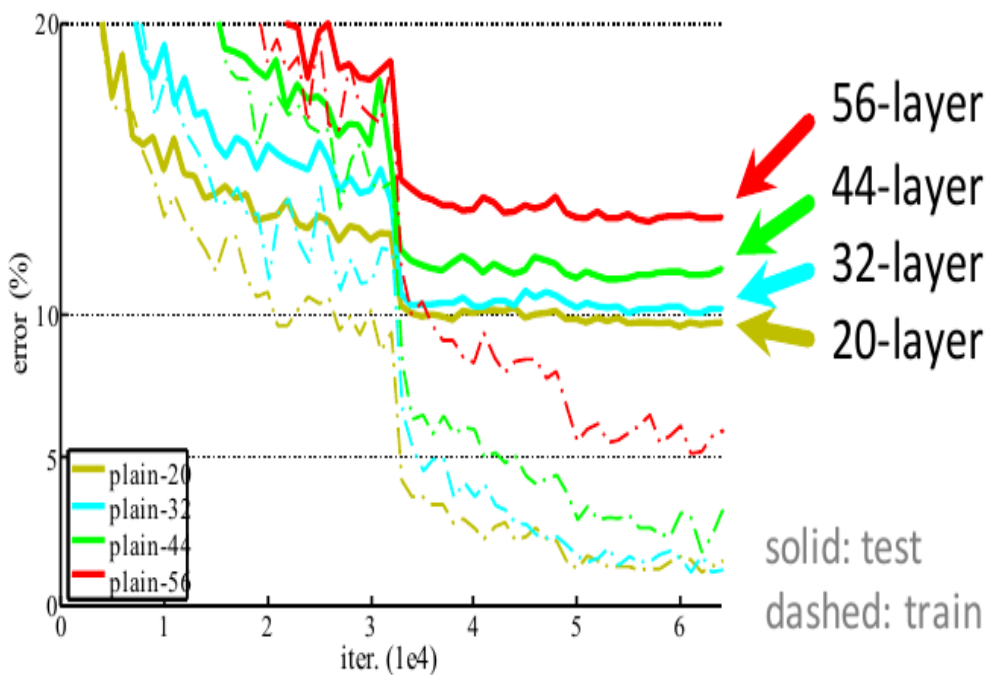
a shallower model
(18 layers)



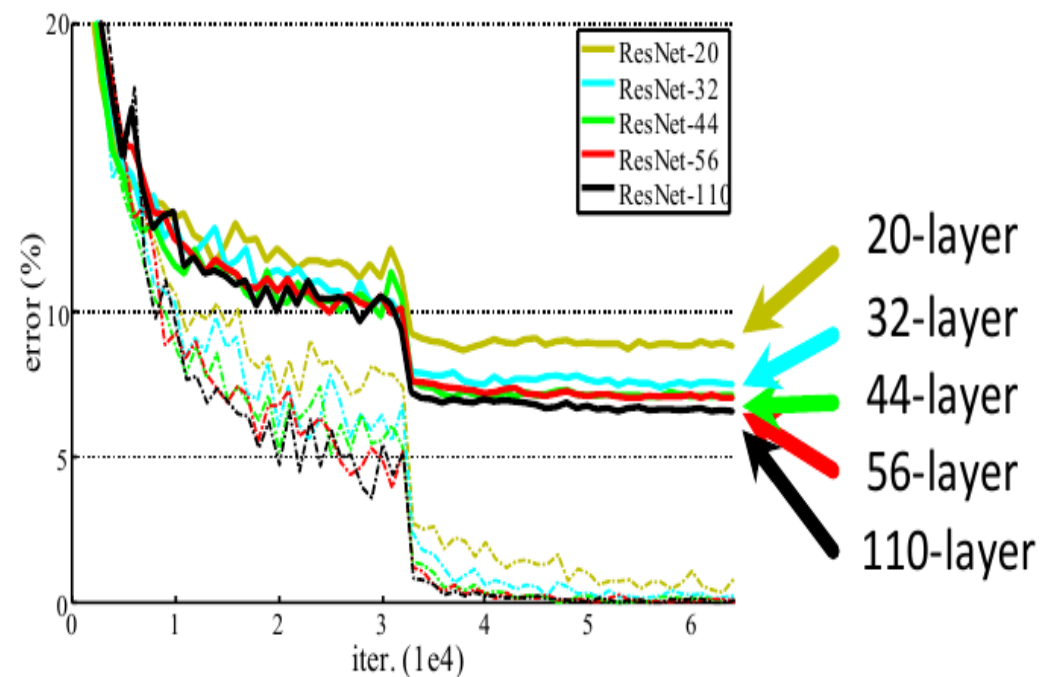
Residual networks

Results: Deeper residual networks result higher accuracy

CIFAR-10 plain nets

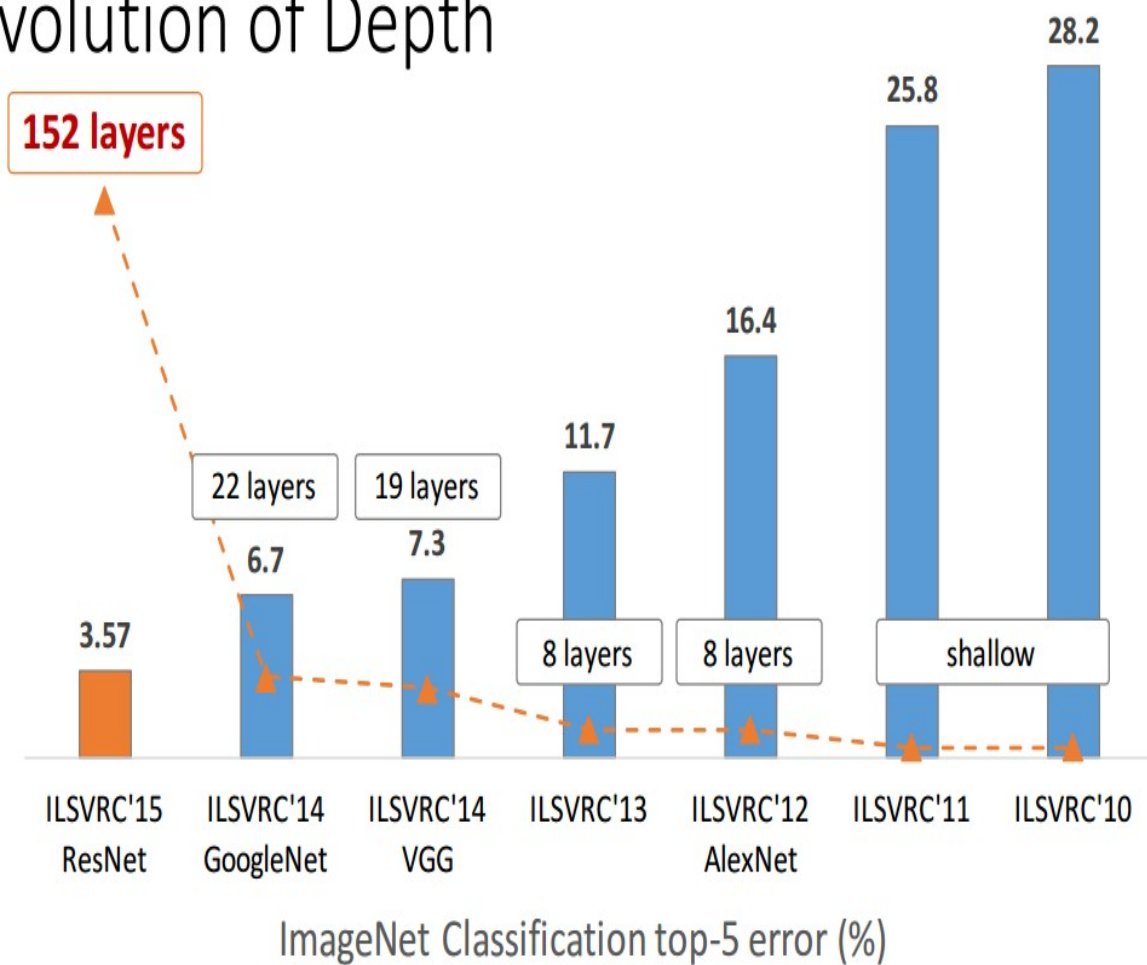


CIFAR-10 ResNets

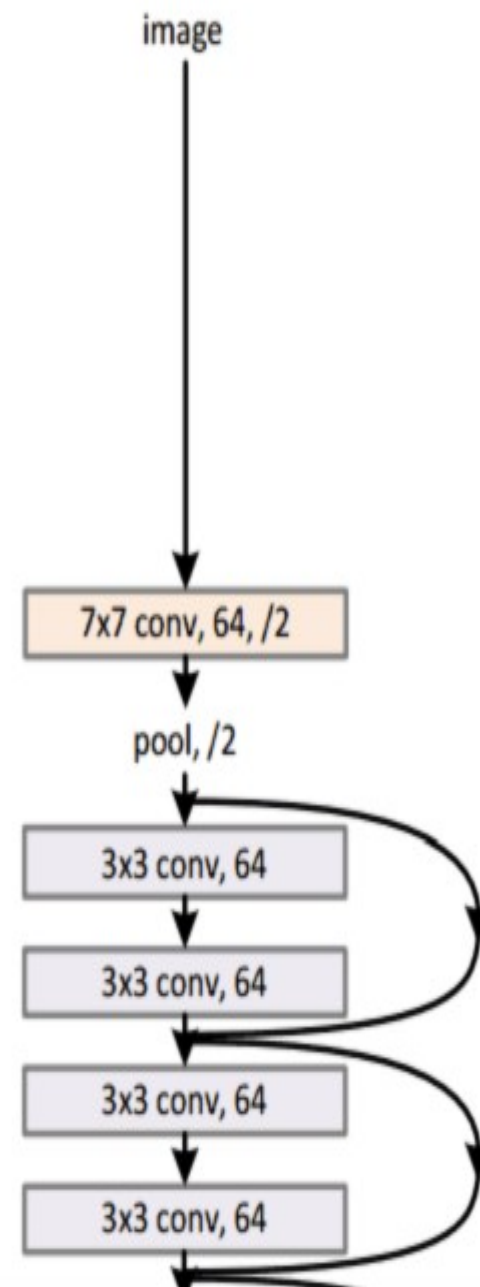


Results with ResNets

Revolution of Depth



34-layer residual



Results with ResNets

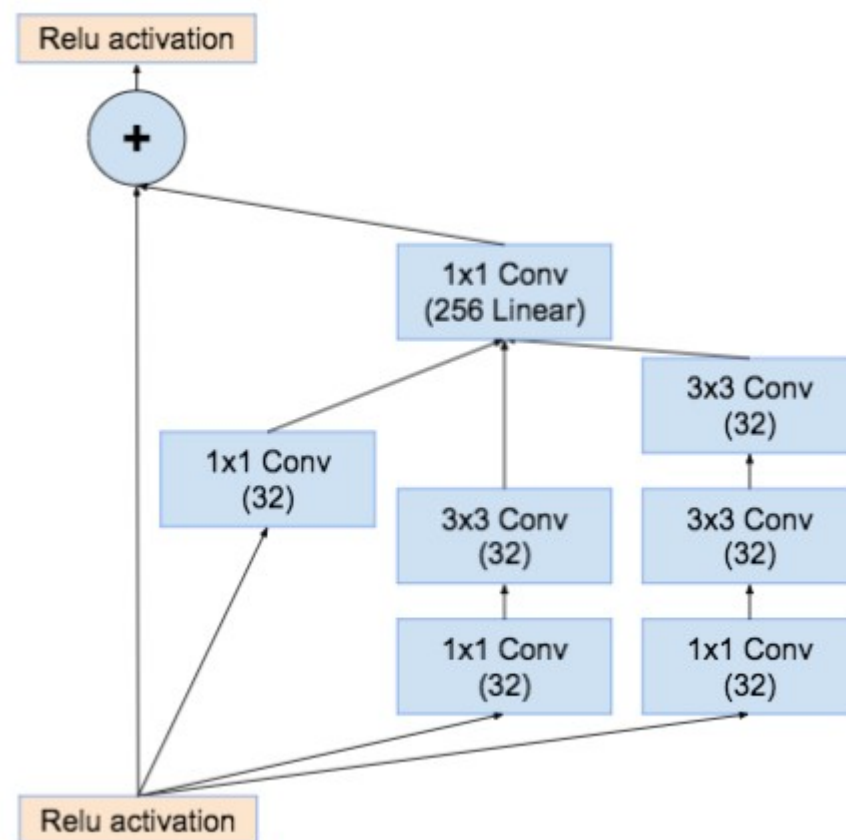
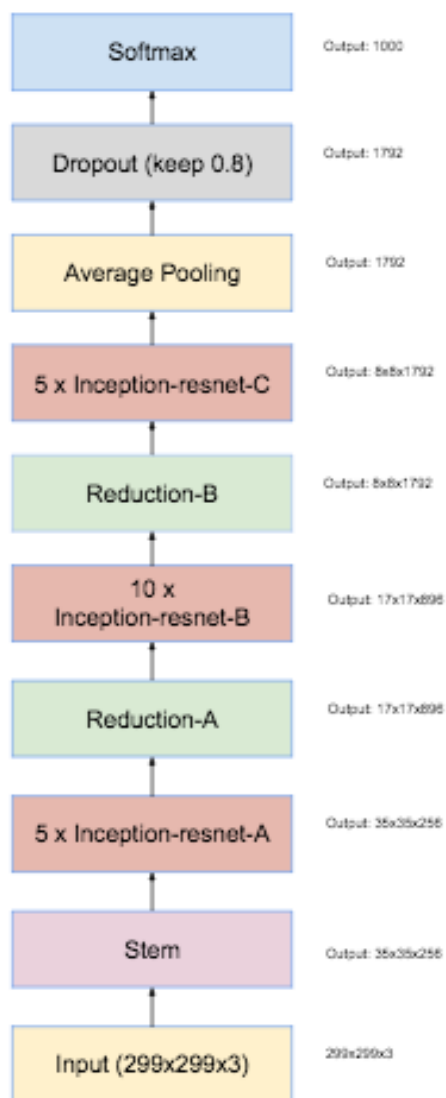
ResNets had the lowest error rate at most competitions since 2015

1st places in all five main tracks

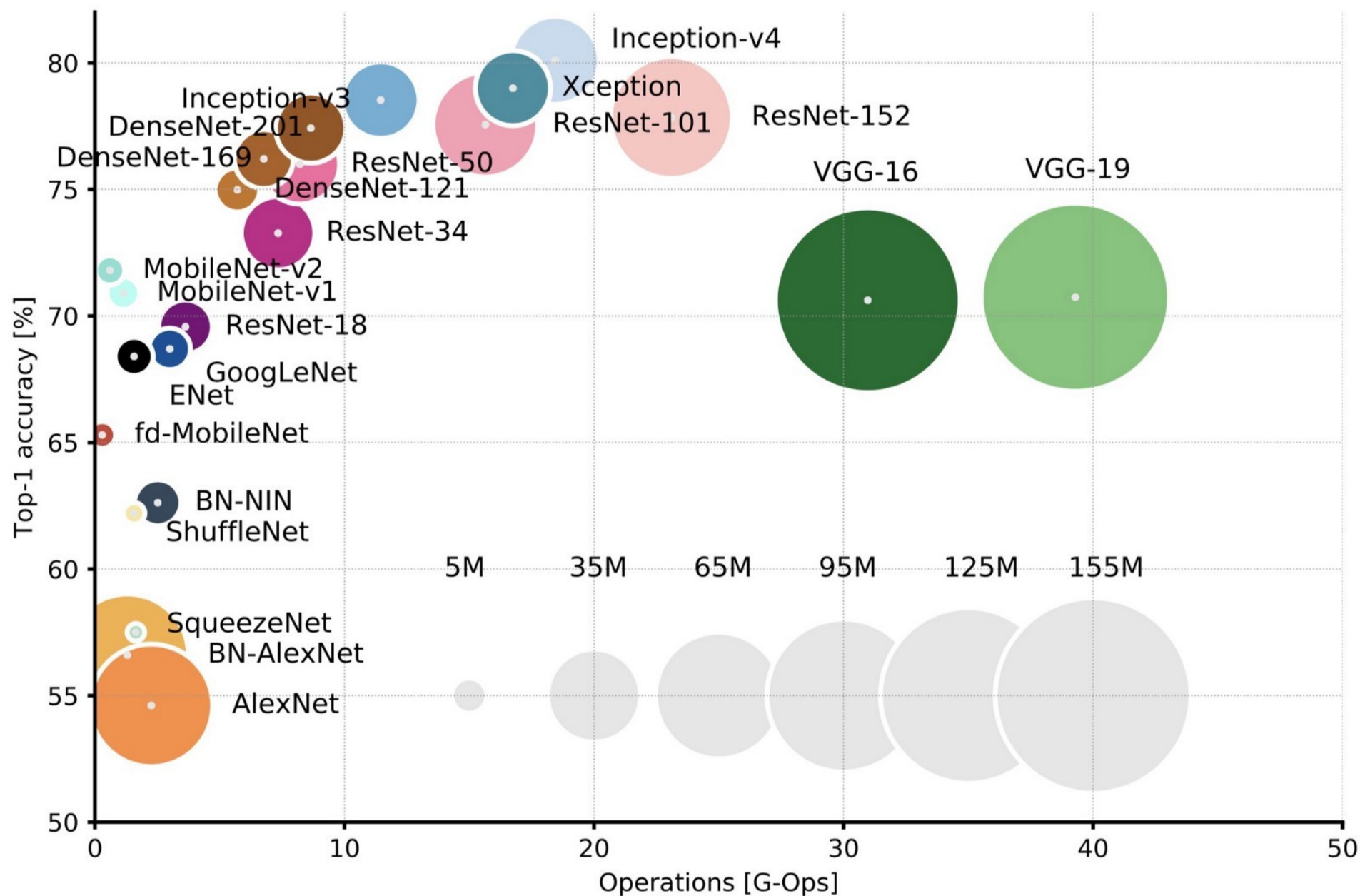
- ImageNetClassification: “Ultra-deep” 152-layer nets
- ImageNetDetection: **16%** better than 2nd
- ImageNetLocalization: **27%** better than 2nd
- COCO Detection: **11%** better than 2nd
- COCO Segmentation: **12%** better than 2nd

GoogleNet Inception v4

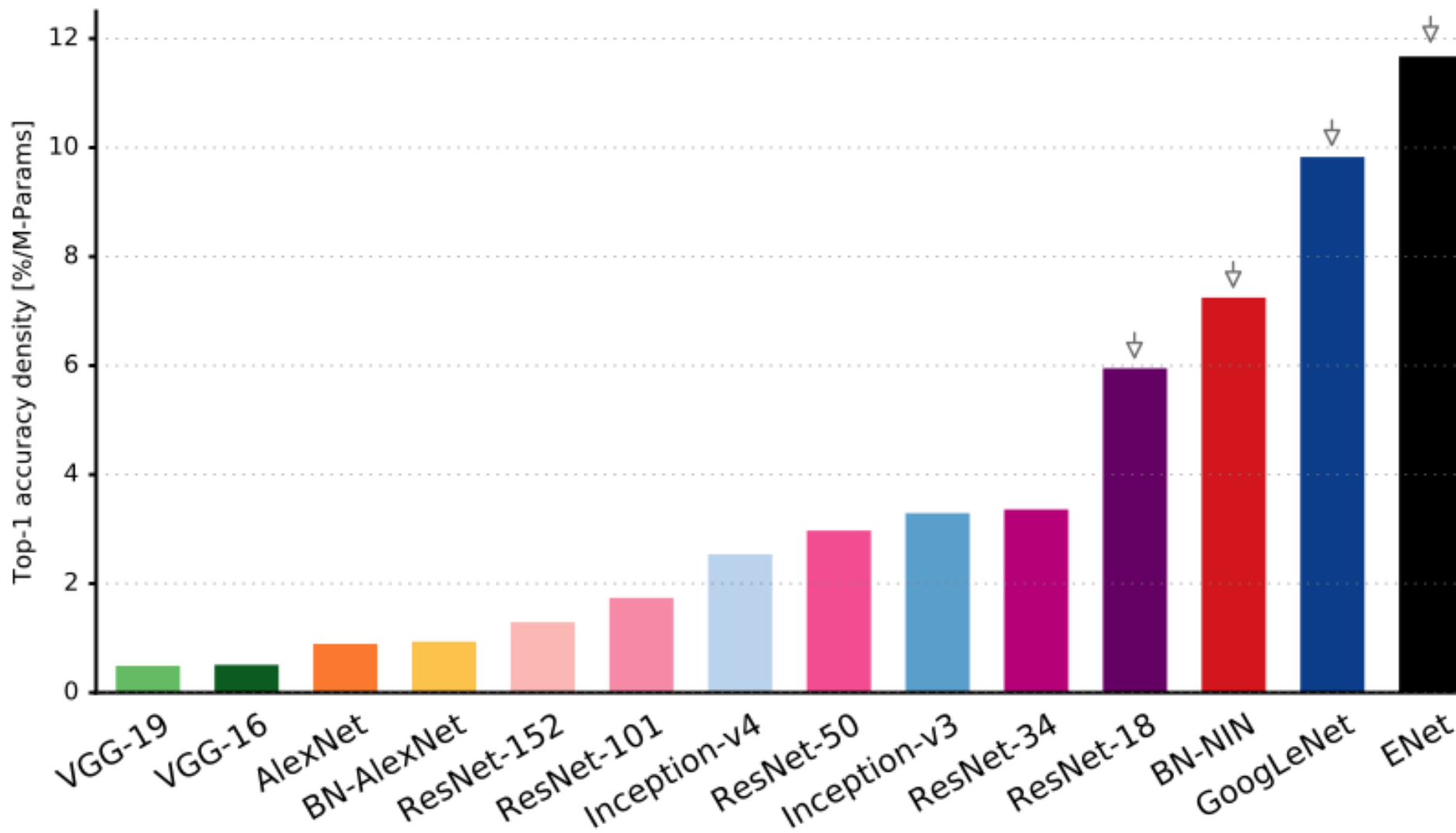
Inception architecture applied to residual networks



Efficiency of Neural Networks



Efficiency of Neural Networks

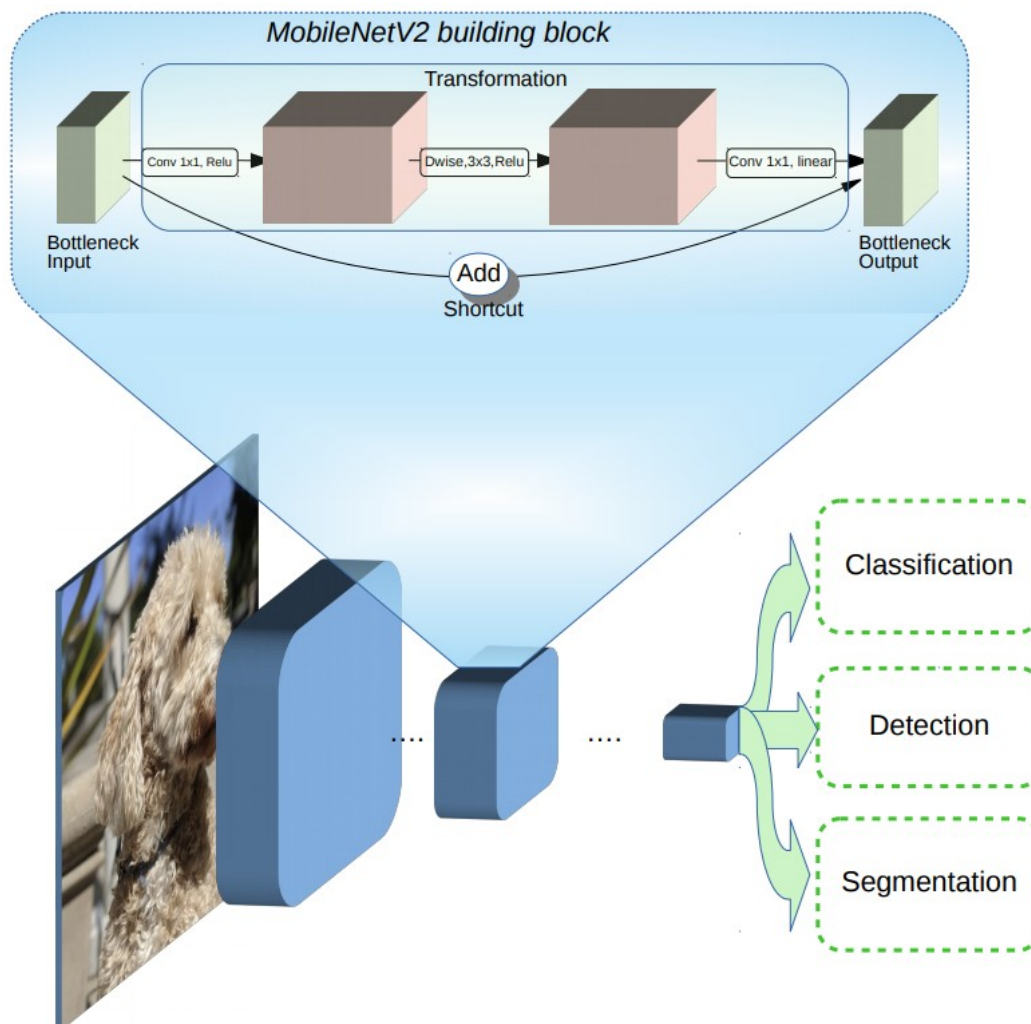


MobileNet

In this architecture depths are squeezed before each operation

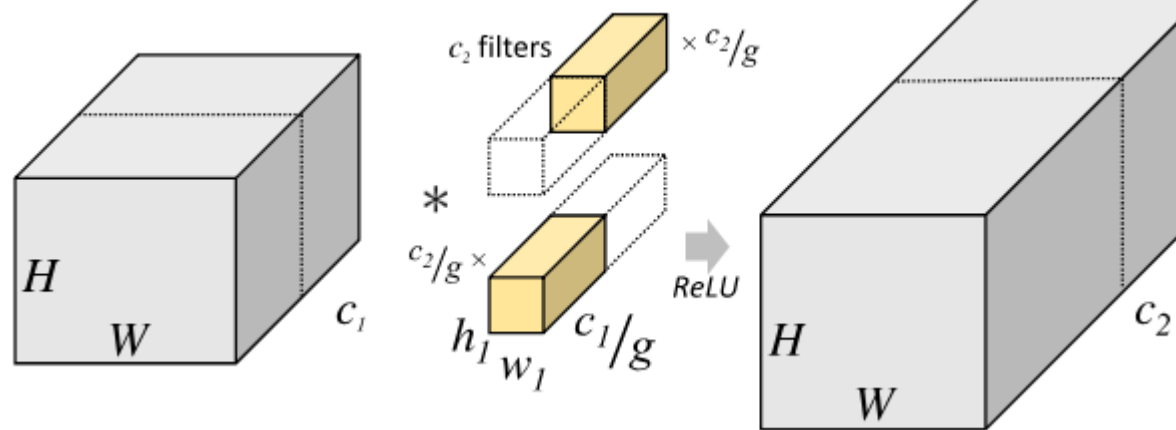
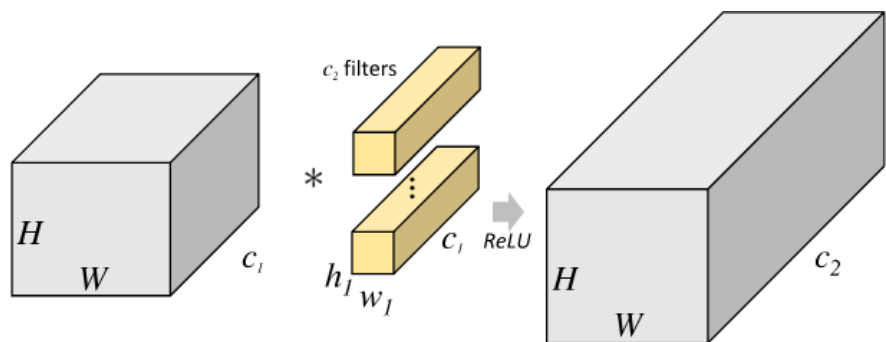
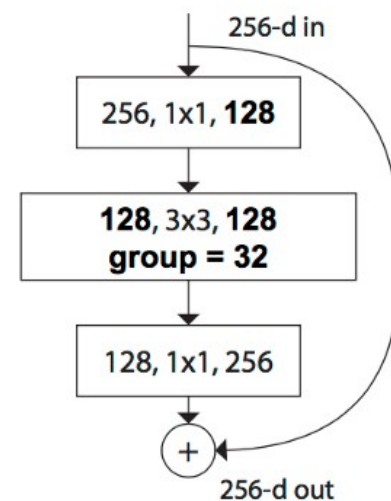
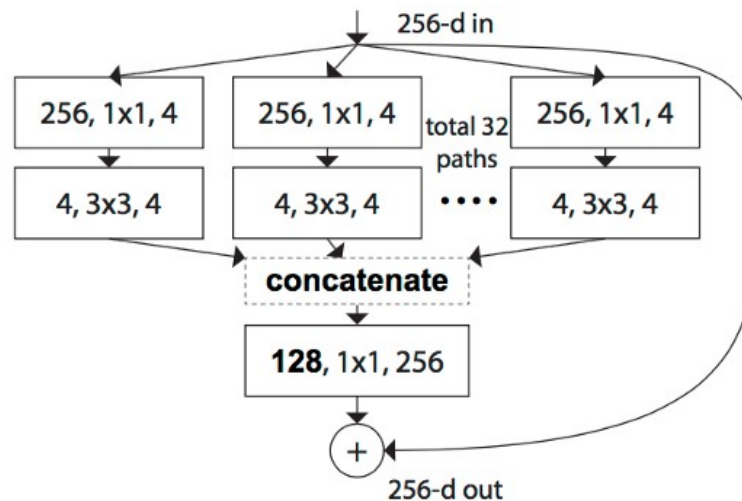
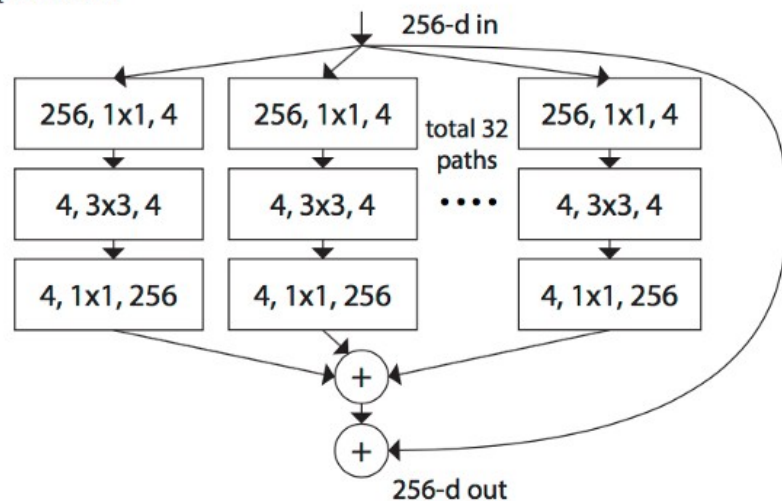
In a squeezed architecture we will use a linear approximation of 128 feature maps, using 16 independent feature maps
And expanded by 1x1 convolution

From the linear combination of these elements the new maps are created

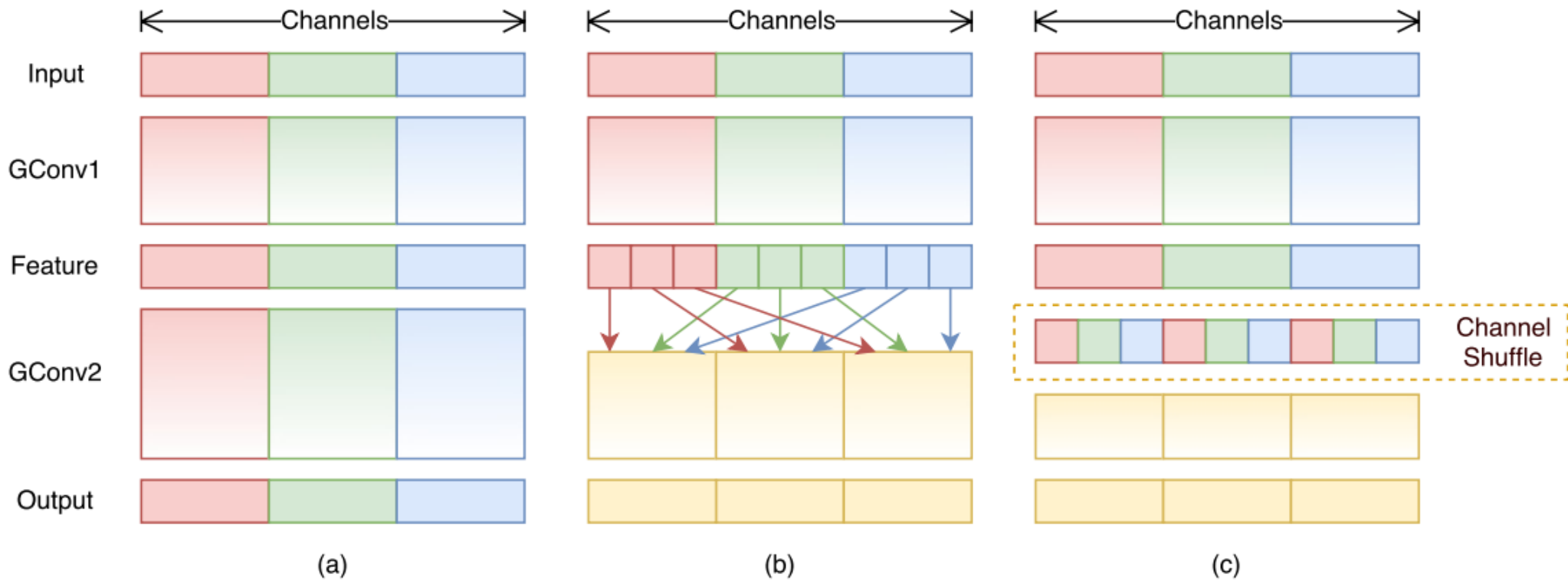


ResNext

equivalent



ShuffleNet



SqueezeNet

In this architecture depths are squeezed before each operation

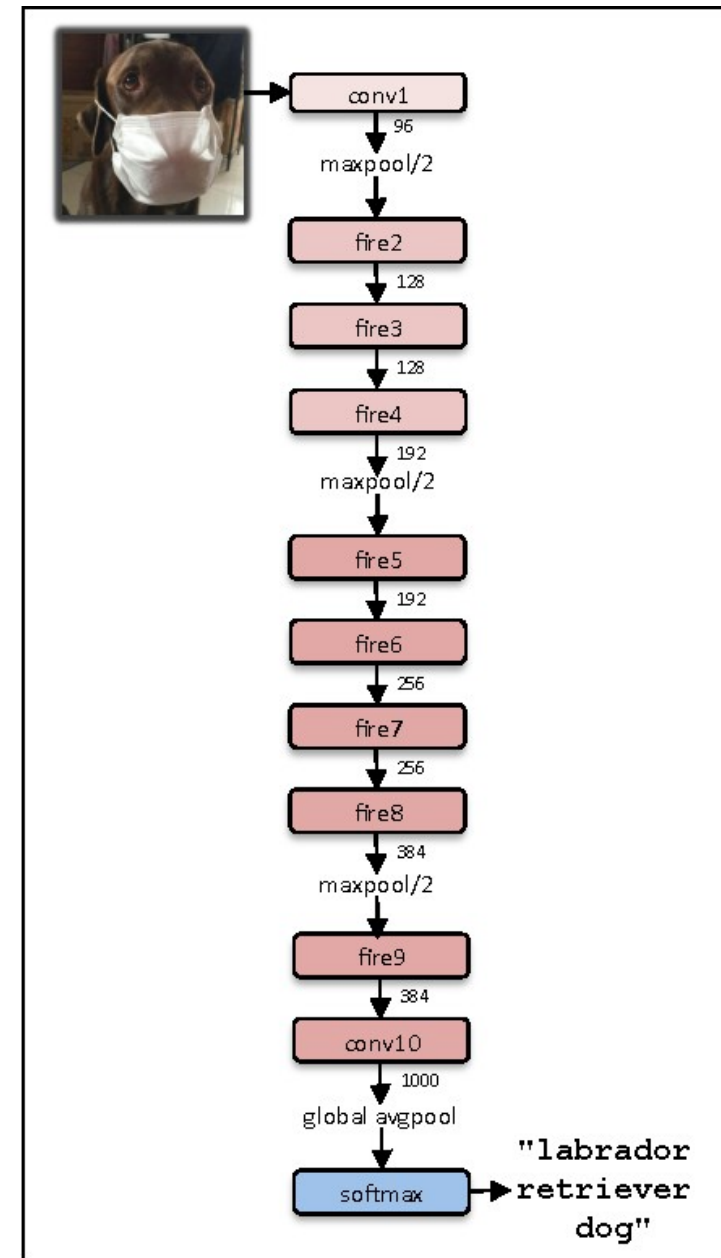
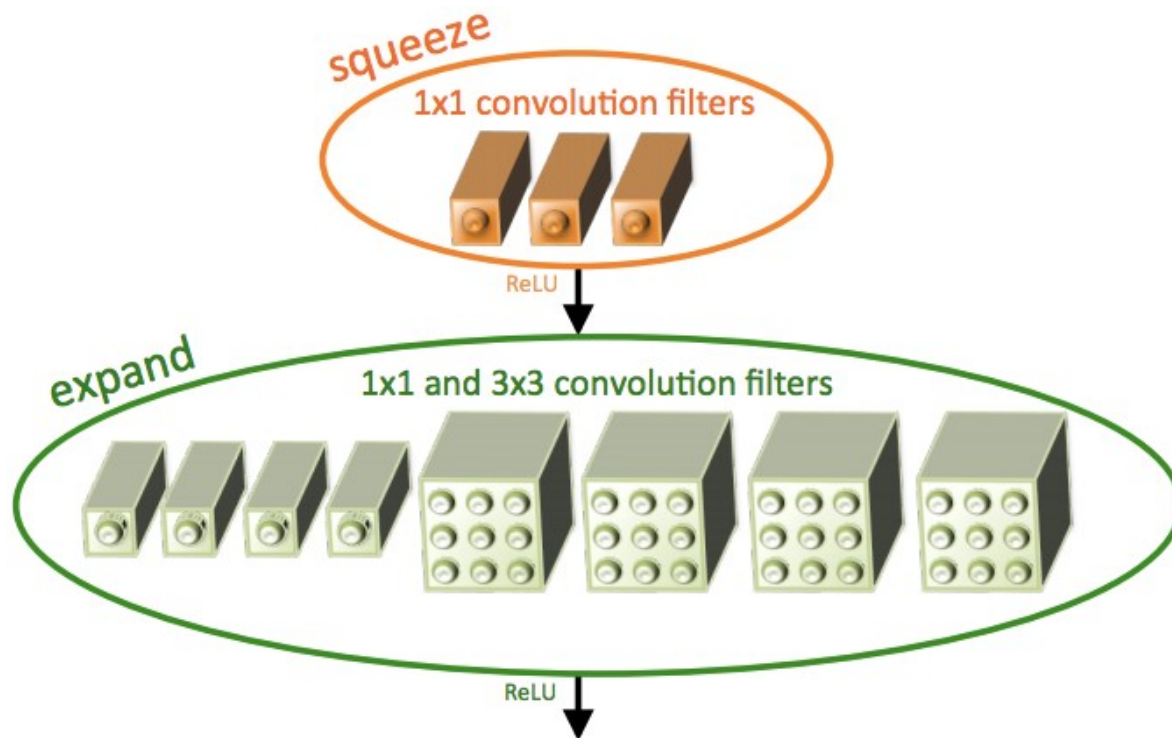


Figure 2. The SqueezeNet architecture

SqueezeNet

In this architecture depths are squeezed before each operation

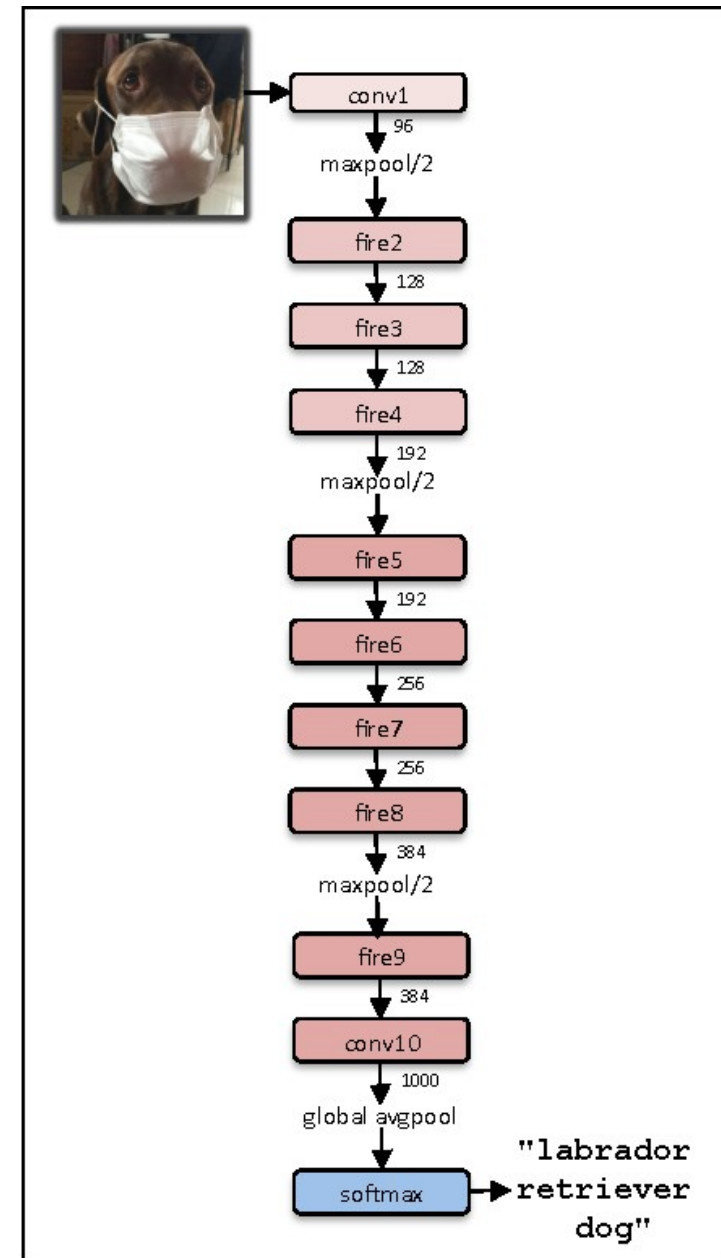
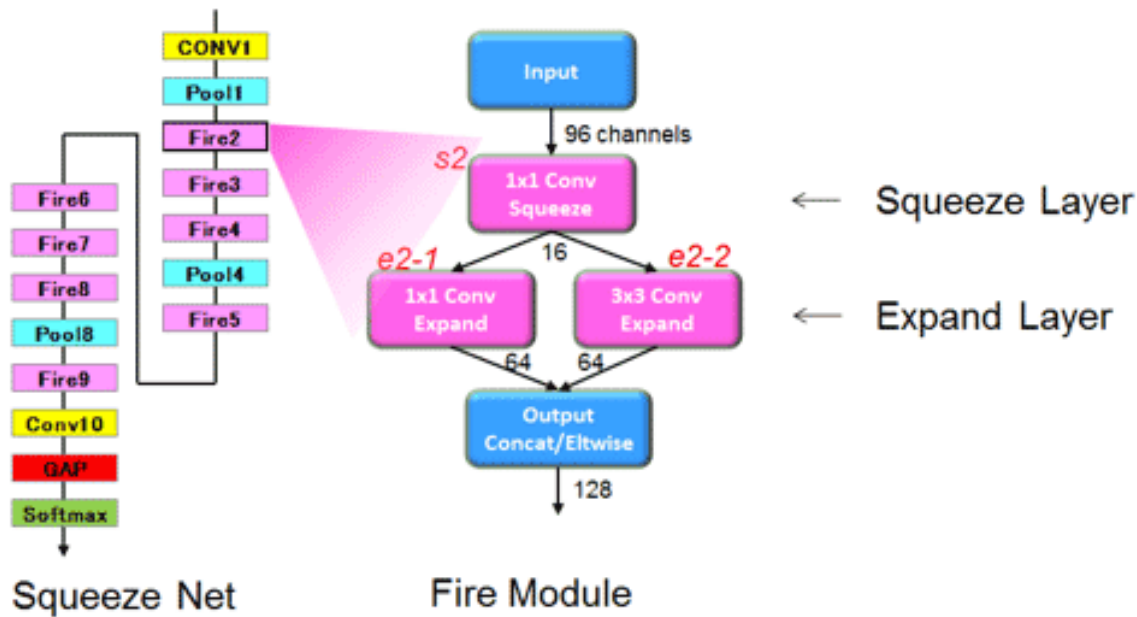


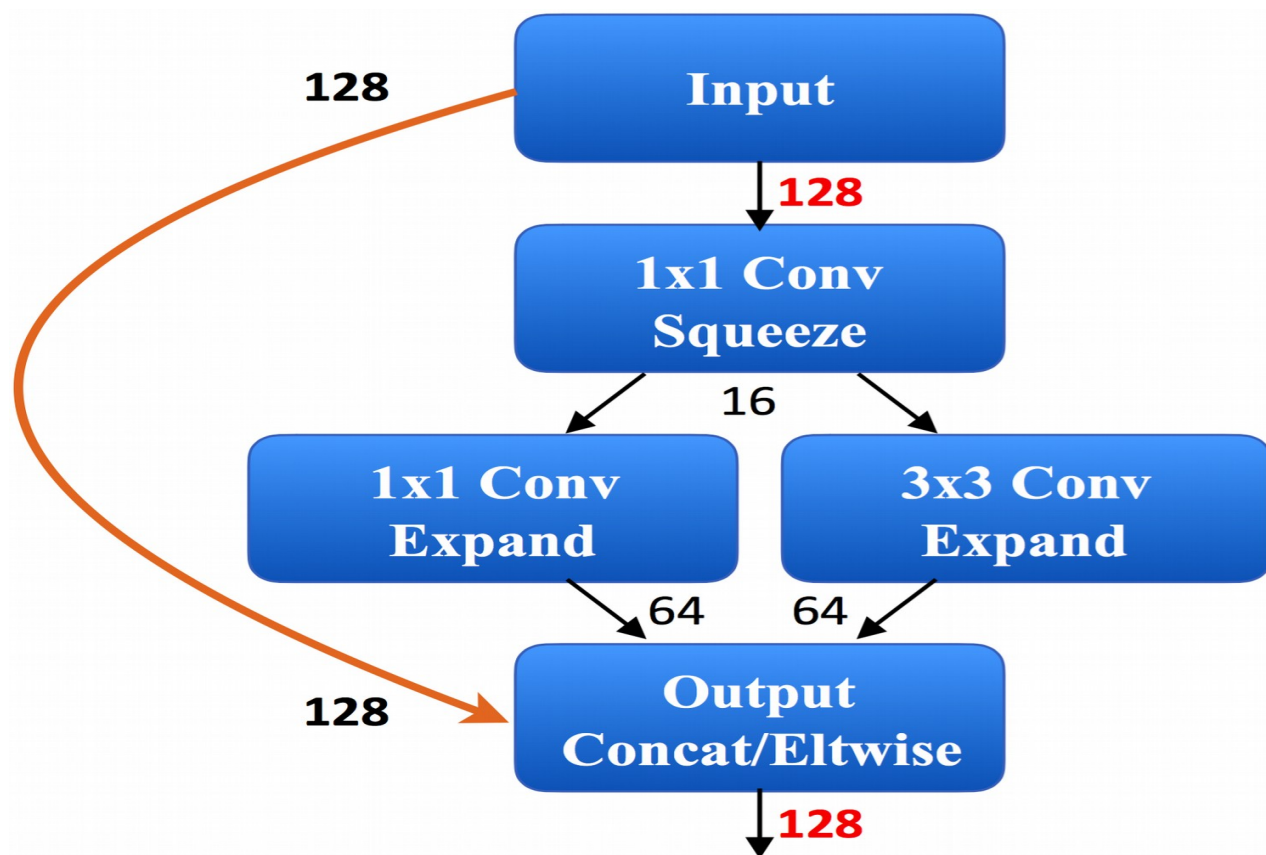
Figure 2. The SqueezeNet architecture

SqueezeNet

In this architecture depths are squeezed before each operation

In a squeezenet architecture we will use a linear approximation of 128 feature maps, using 16 independent feature maps

From the linear combination of these elements the new maps are created



Neural networks for regression

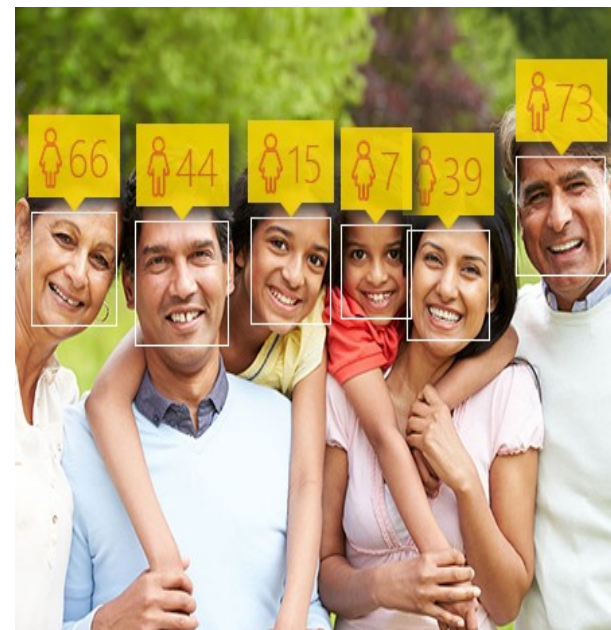
Age estimation

The output is not discrete classes or pixels, but continuous values

The network structure can remain the same but a different loss function and differently annotated dataset is needed.

Hard to interpret the error in common tasks.

E.G: Age estimation on images:



Neural networks for regression

Multiple object detection on a single image

Classification is good for a single object (can be extended for k objects – top k candidates)

How could we detect objects in general, when the number of objects is unknown

Classification



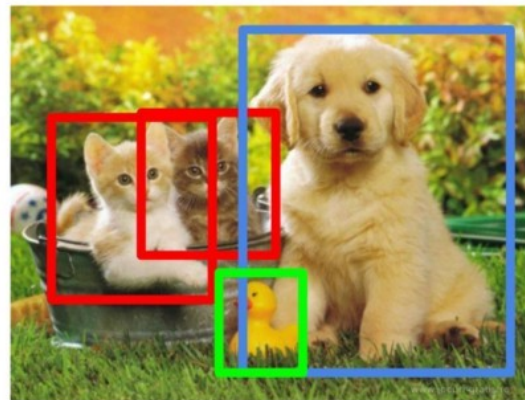
CAT

**Classification
+ Localization**



CAT

Object Detection



CAT, DOG, DUCK

**Instance
Segmentation**



CAT, DOG, DUCK

Single object

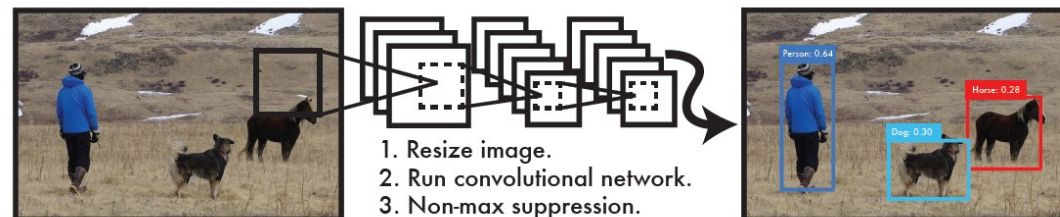
Multiple objects

Traditional method

Sliding window over the image

We might have objects in different scales

Slidign windowds in different scales, aspect ratios



Results a heat map → detect the objects: non-maximum suppression



Object detection as regression

Single shot object detector SSD (2016 March)

You Only Look Once YOLO (2016 May)

Classification



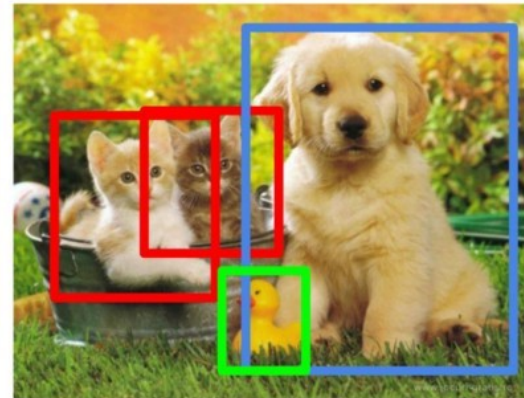
CAT

**Classification
+ Localization**



CAT

Object Detection



CAT, DOG, DUCK

**Instance
Segmentation**



CAT, DOG, DUCK

Single object

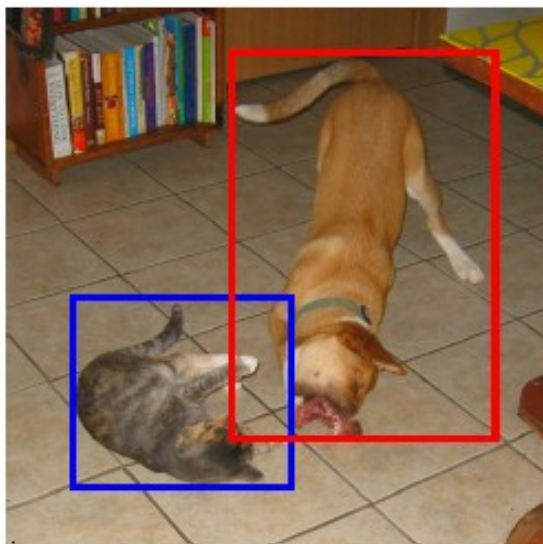
Multiple objects

SSD

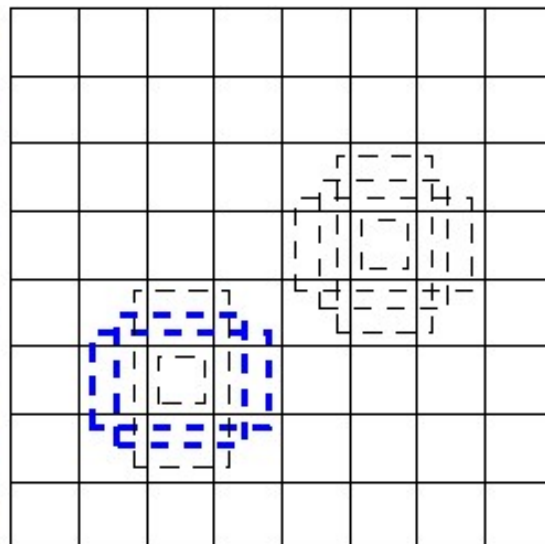
Single shot object detector SSD (2016 March)

Has a fixed resolution and the last feature maps (with different scales) can be considered as maps of bounding boxes

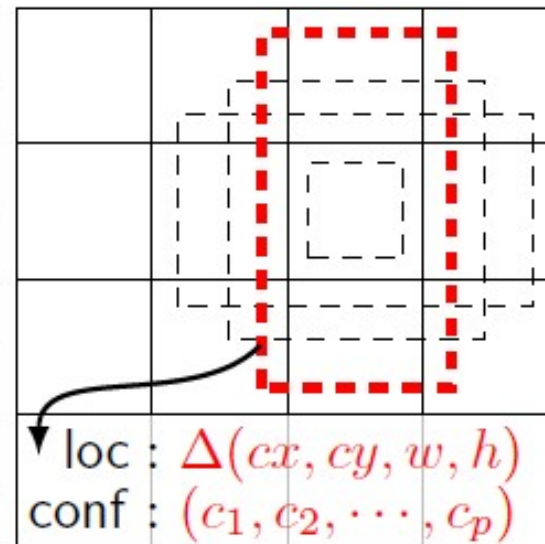
On these maps each pixel represent bounding boxes. A high pixel value represent high probability of the centerpoint of a detected object.



(a) Image with GT boxes



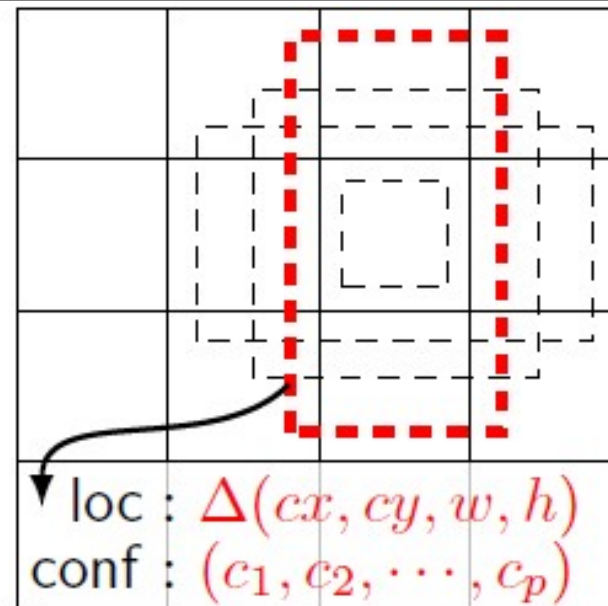
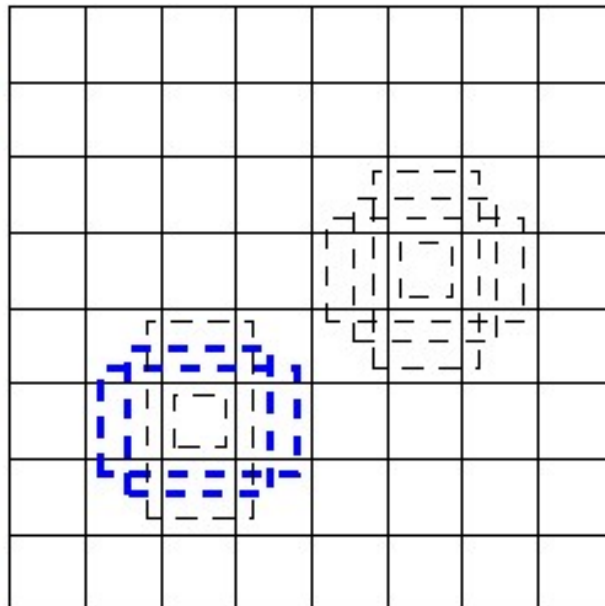
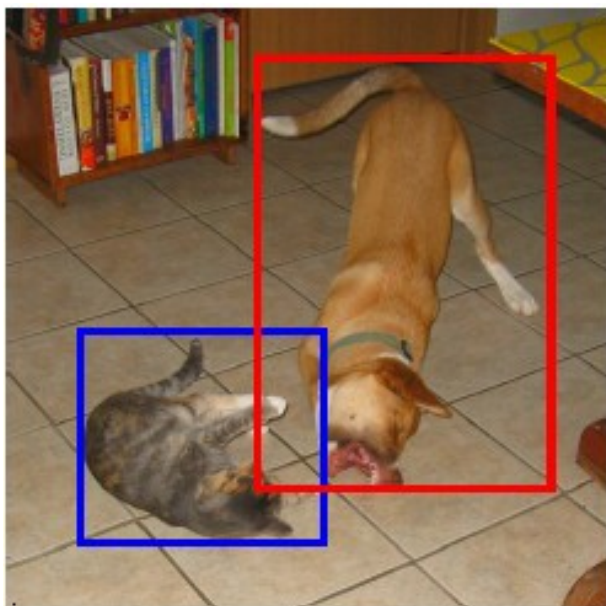
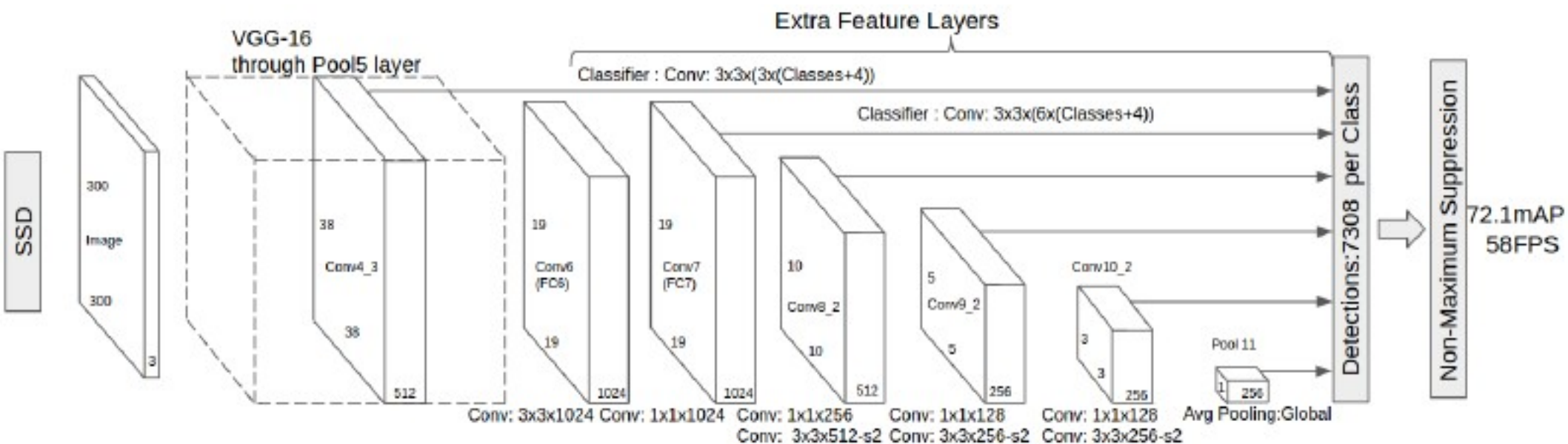
(b) 8×8 feature map



loc : $\Delta(cx, cy, w, h)$
conf : (c_1, c_2, \dots, c_p)

(c) 4×4 feature map

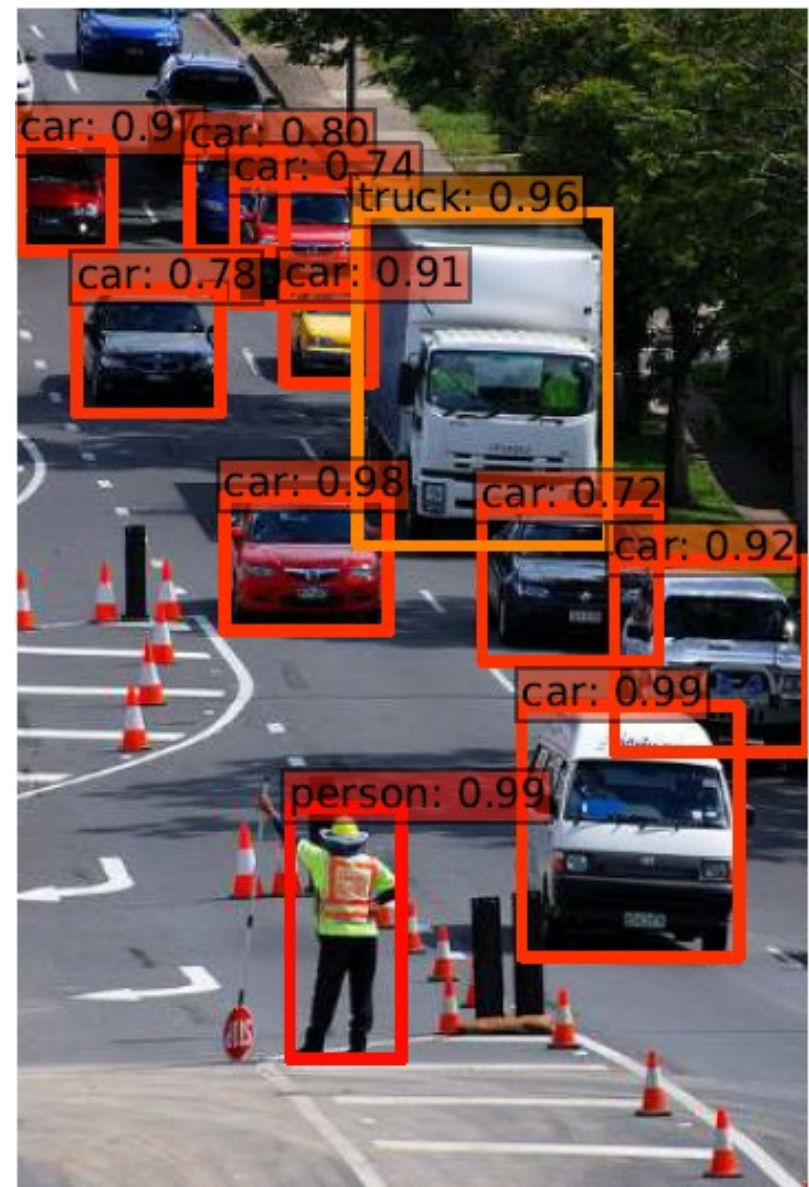
SSD architecture



Loss function for bounding boxes

$$L(x, c, l, g) = \frac{1}{N} (L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

- The **overall objective loss function** is a weighted sum of the localization loss and the confidence loss(conf)
- N: the number of matched default boxes
- l: predicted boxes
- g: the ground truth box
- x=1 denotes some certain default box is matched to a ground truth box



R-CNN

Region proposal CNN network

Separate the problem of object detection and classification

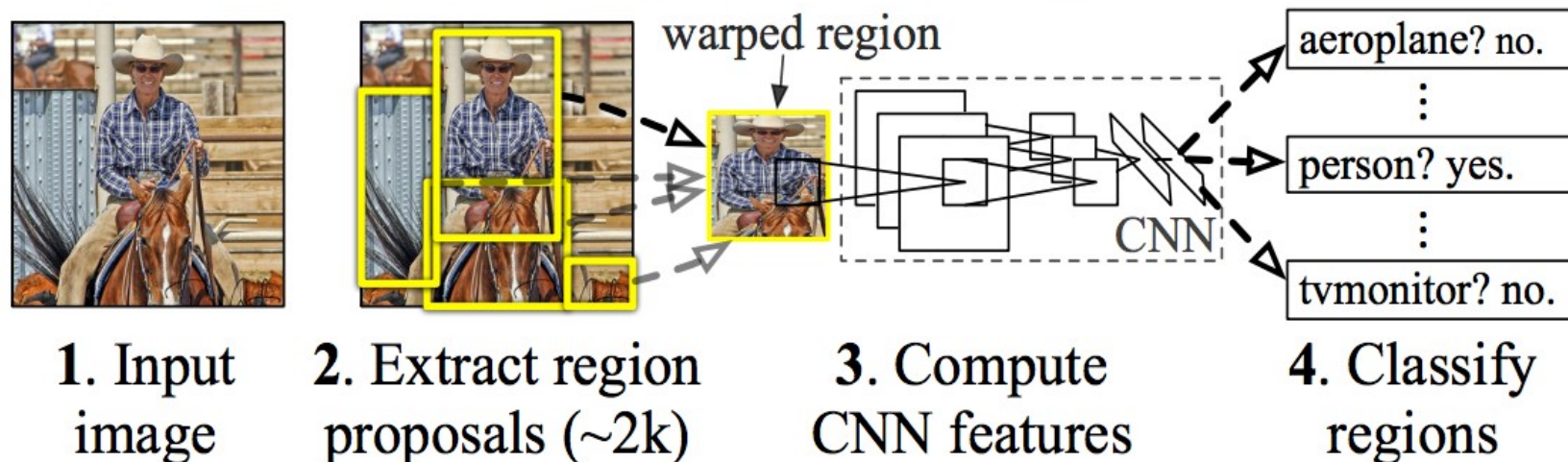
It consists of three modules.

The first generates category-independent region proposals. These proposals define the set of candidate detection available to detector.

The second module is a large convolutional neural network that extracts a fixed-length feature vector from each region.

The third module is a set of class-specific linear SVMs

R-CNN: Regions with CNN features





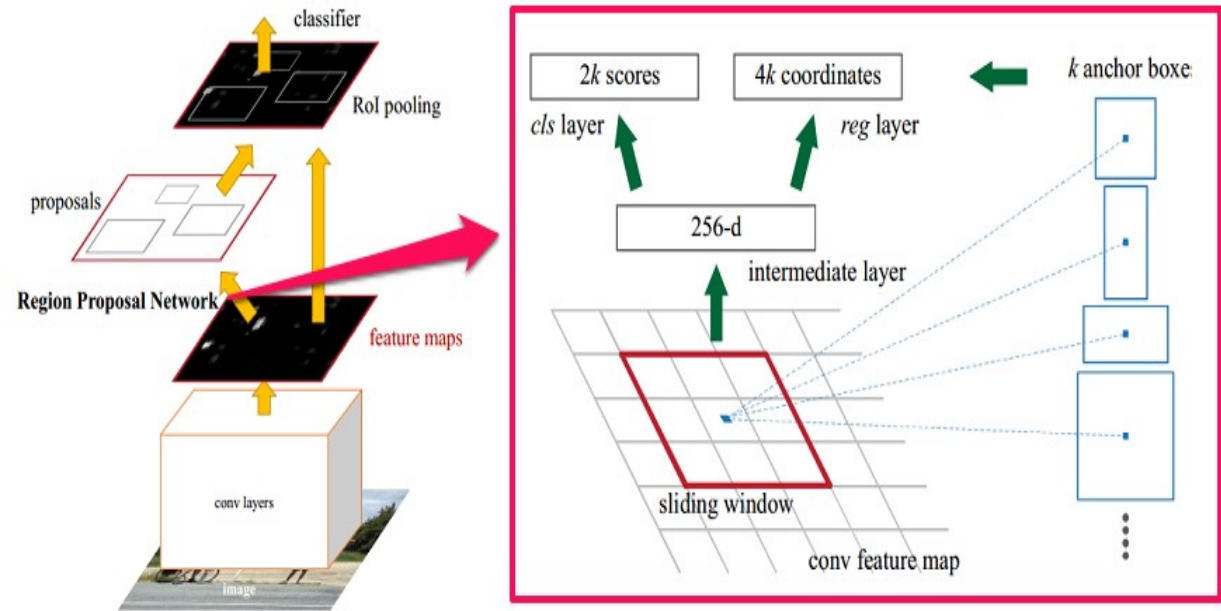
R-CNN

Region proposal from a network

Step 3 and 4 are standard CNN implementations

Extra layers for region proposals

Possible region refinement at the end



R-CNN: *Regions with CNN features*

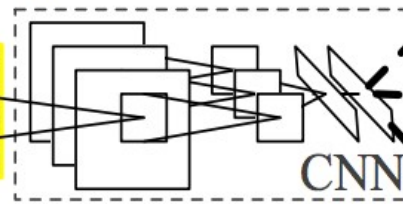


1. Input image



2. Extract region proposals (~2k)

warped region



3. Compute CNN features

aeroplane? no.

person? yes.

tvmonitor? no.

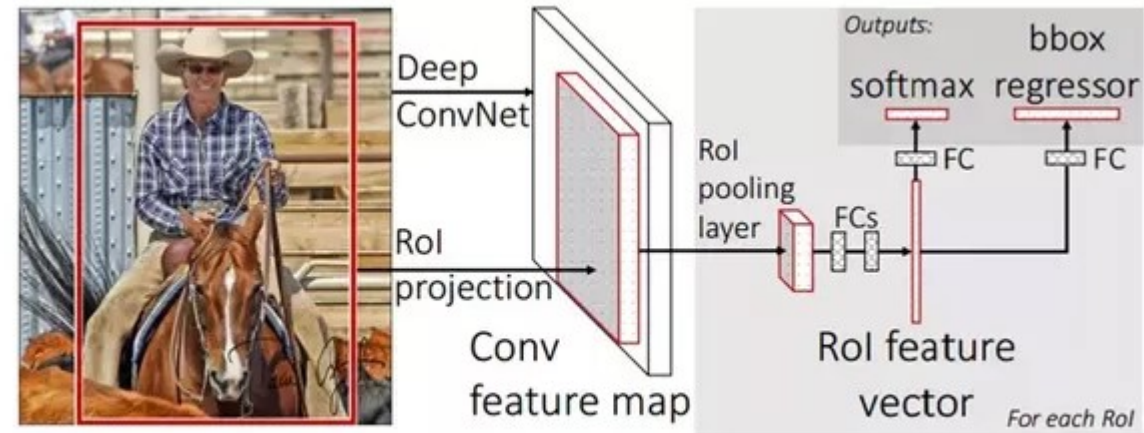
4. Classify regions

Fast R-CNN

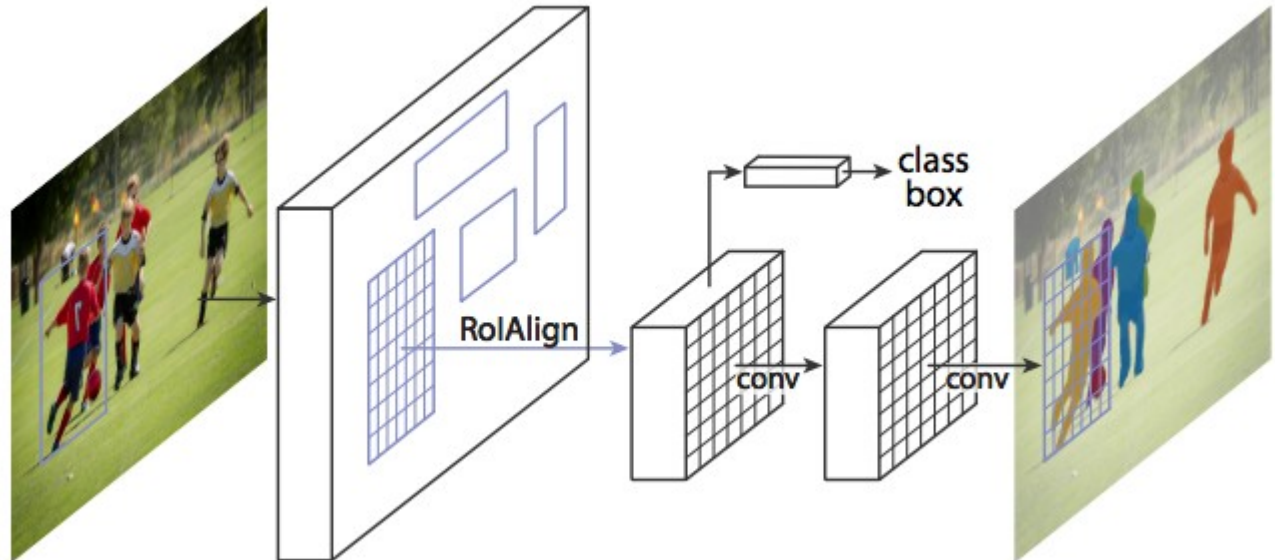
Fast R-CNN:

Speeding up by sharing computation of the conv layers between the proposal generation and classification

In this model, the image is first fed through a ConvNet, features of the region proposals are obtained from the last feature map of the ConvNet and lastly we have fully connected layers as well as our regression and classification heads.



Fast R-CNN workflow





YOLO, Detectnet

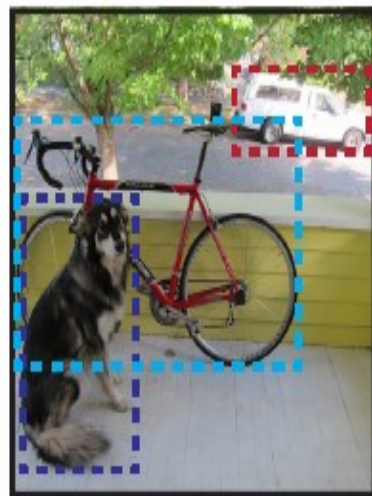
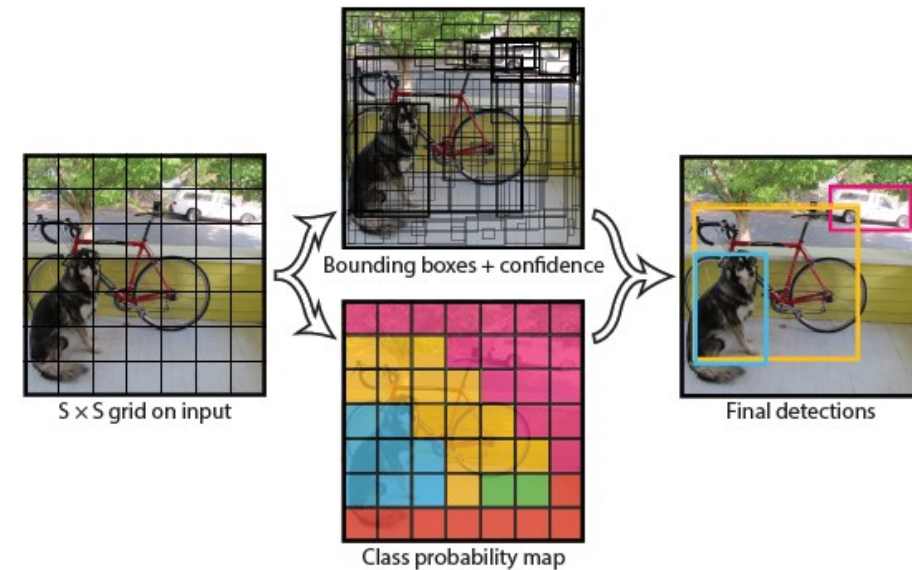
Models detection as a regression problem:

Divide the image into a grid and each cell can vote for the bounding box position of possible object.

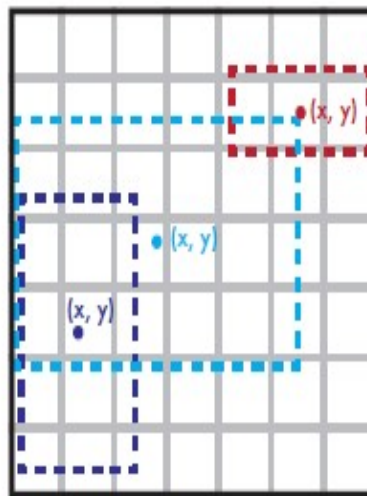
Boxes can have arbitrary sizes

Non-suppression on the boxes

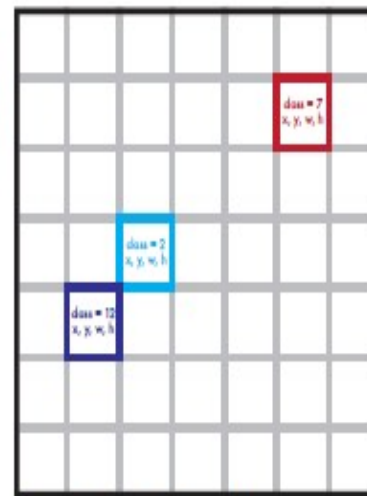
No need for scale search, the image is processed once and objects in different scales can be detected



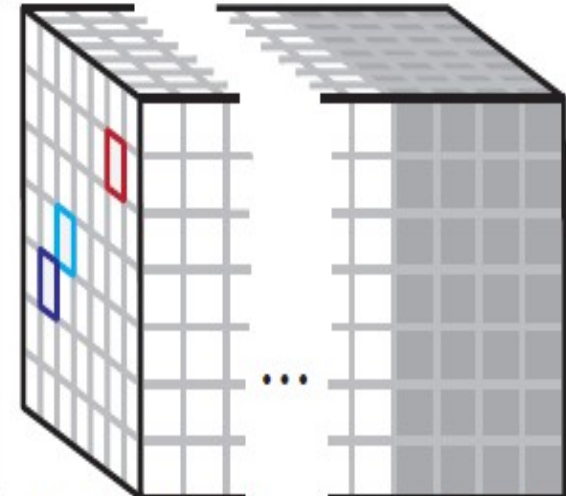
Resize The Image
And bounding boxes to 448 x 448.



Divide The Image
Into a 7 x 7 grid. Assign detections to grid cells based on their centers.



Train The Network
To predict this grid of class probabilities and bounding box coordinates.

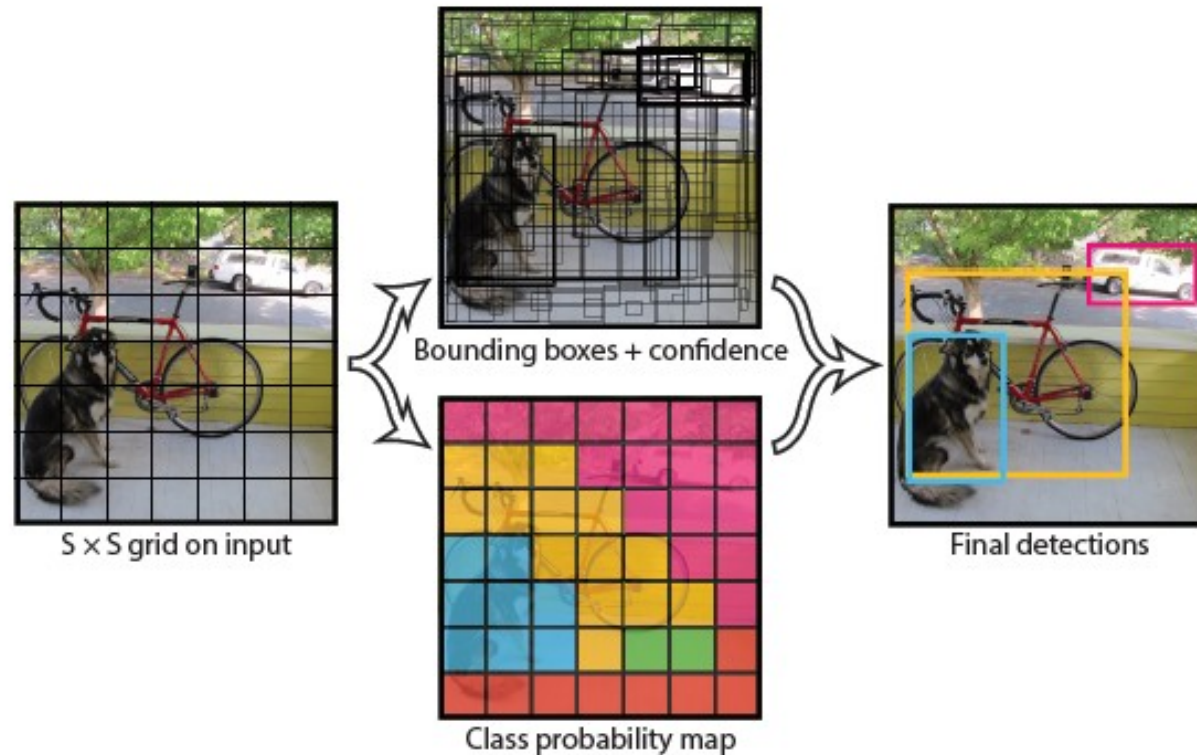


1st - 20th Channels:
Class probabilities
 $\text{Pr}(\text{Airplane}), \text{Pr}(\text{Bike})...$

Last 4 Channels:
Box coordinates
 x, y, w, h

Redmon, Joseph, et al. "You only look once: Unified, real-time object detection." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016.

How unified detection works?



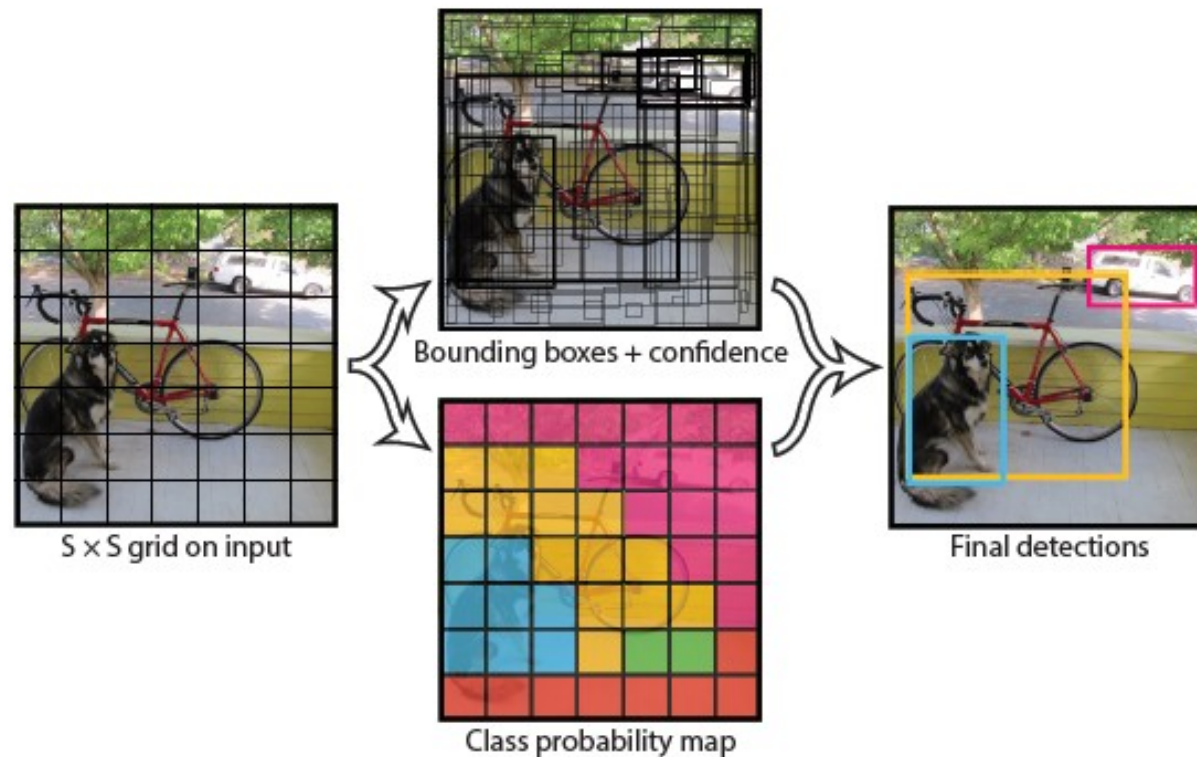
confidence scores: reflect how confident is that the box contains an object+how accurate the box is .

$$\Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

conditional class probabilities: conditioned on the grid cell containing an object

$$\Pr(\text{Class}_i | \text{Object})$$

How unified detection works?



$$\Pr(\text{Class}_i | \text{Object}) * \Pr(\text{Object}) * \text{IOU}_{\text{pred}}^{\text{truth}} = \Pr(\text{Class}_i) * \text{IOU}_{\text{pred}}^{\text{truth}}$$

- At test time, multiply the conditional class probabilities and the individual box confidence predictions
- giving class-specific confidence scores for each box
- Showing both the probability of that class appearing in the box and how well the predicted box fits the object

Pixel level segmentation

The expected output of the network is not a class, but a map representing the pixels belonging to a certain class.

Creation of a labeled dataset (handmade pixel level mask) is a tedious task

More complex architectures are needed (compared to classification)

Popular architectures (Sharpmask, U-NET ...)



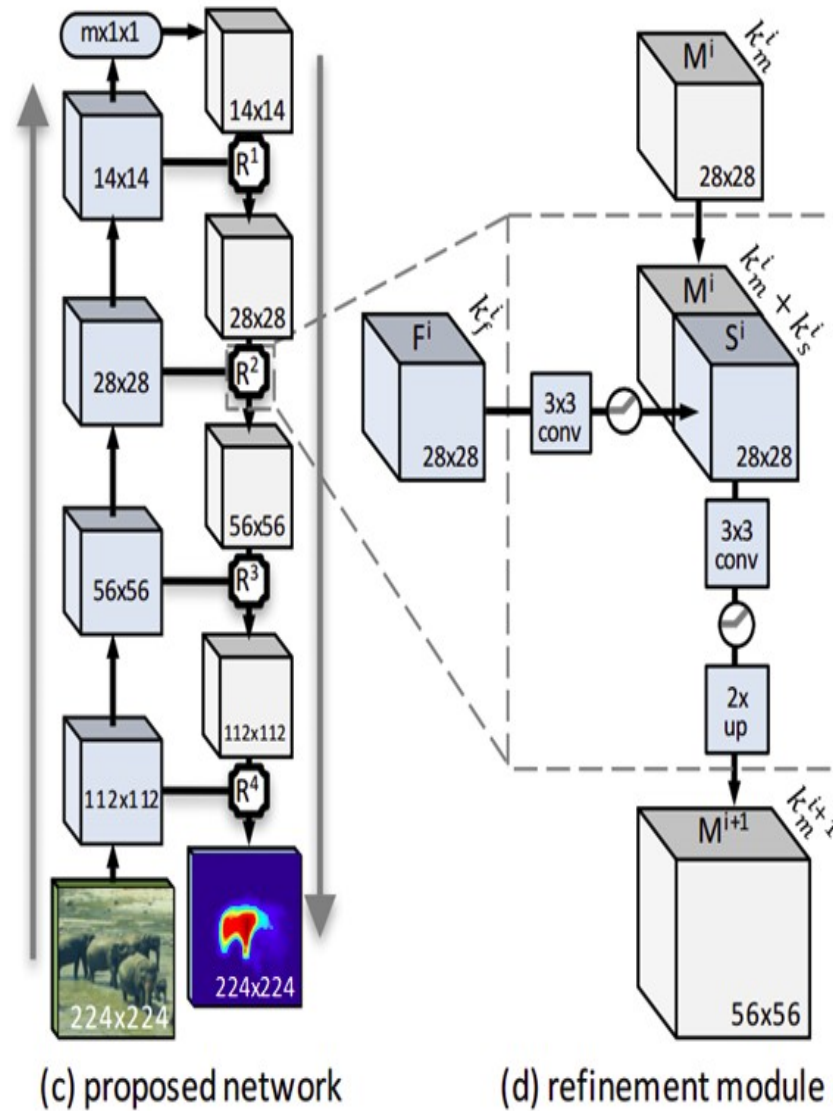
sky tree road grass water bldg mntn fg obj.

SharpMask: *Learning to Refine Object Segments*. Pedro O. Pinheiro, Tsung-Yi Lin, Ronan Collobert, Piotr Dollár (ECCV 2016)



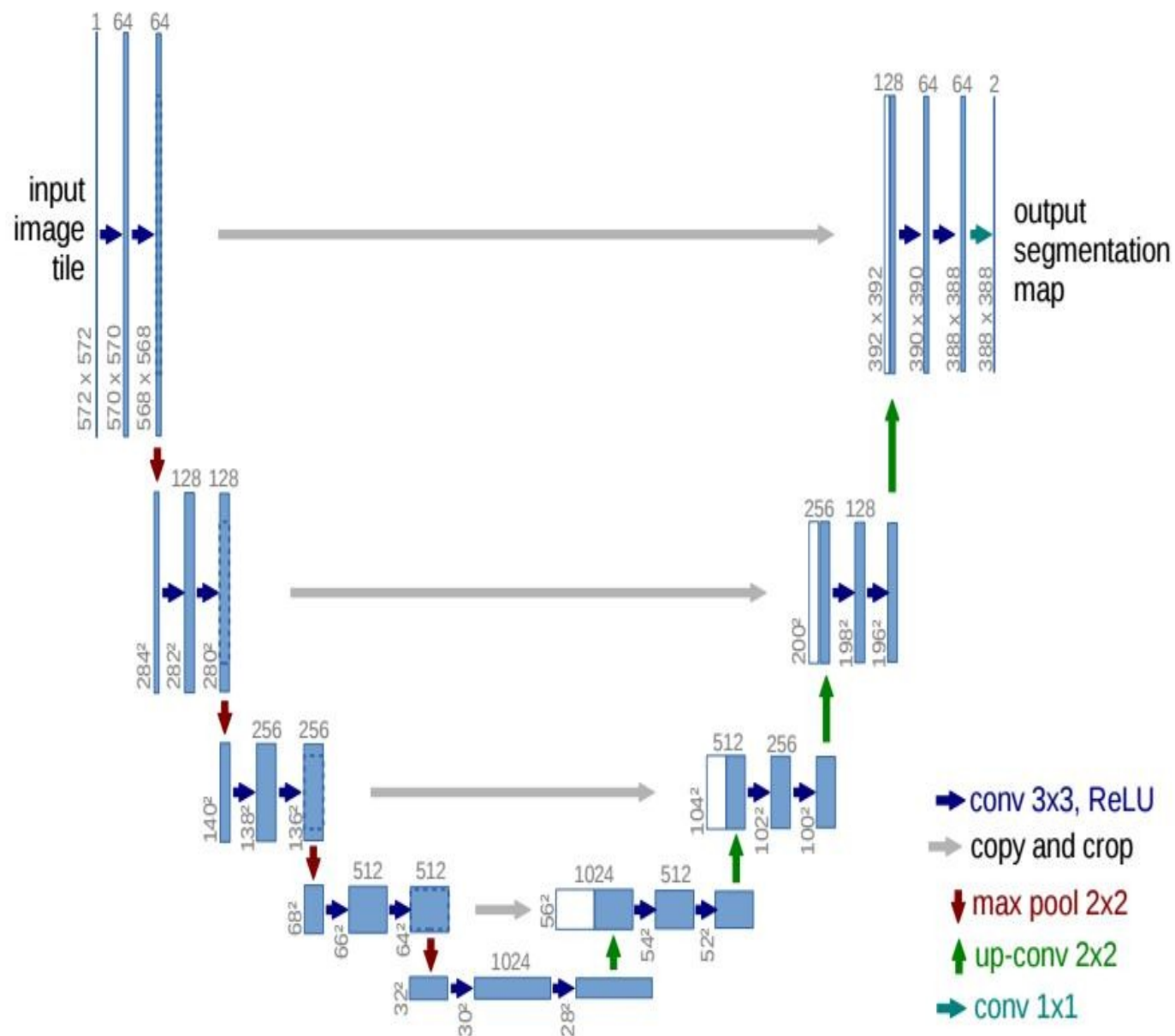
SEMANTIC IMAGE SEGMENTATION WITH DEEP CONVOLUTIONAL NETS AND FULLY CONNECTED CRFS Liang-Chieh Chen et al. ICLR 2015

Sharpmask



SharpMask network architecture

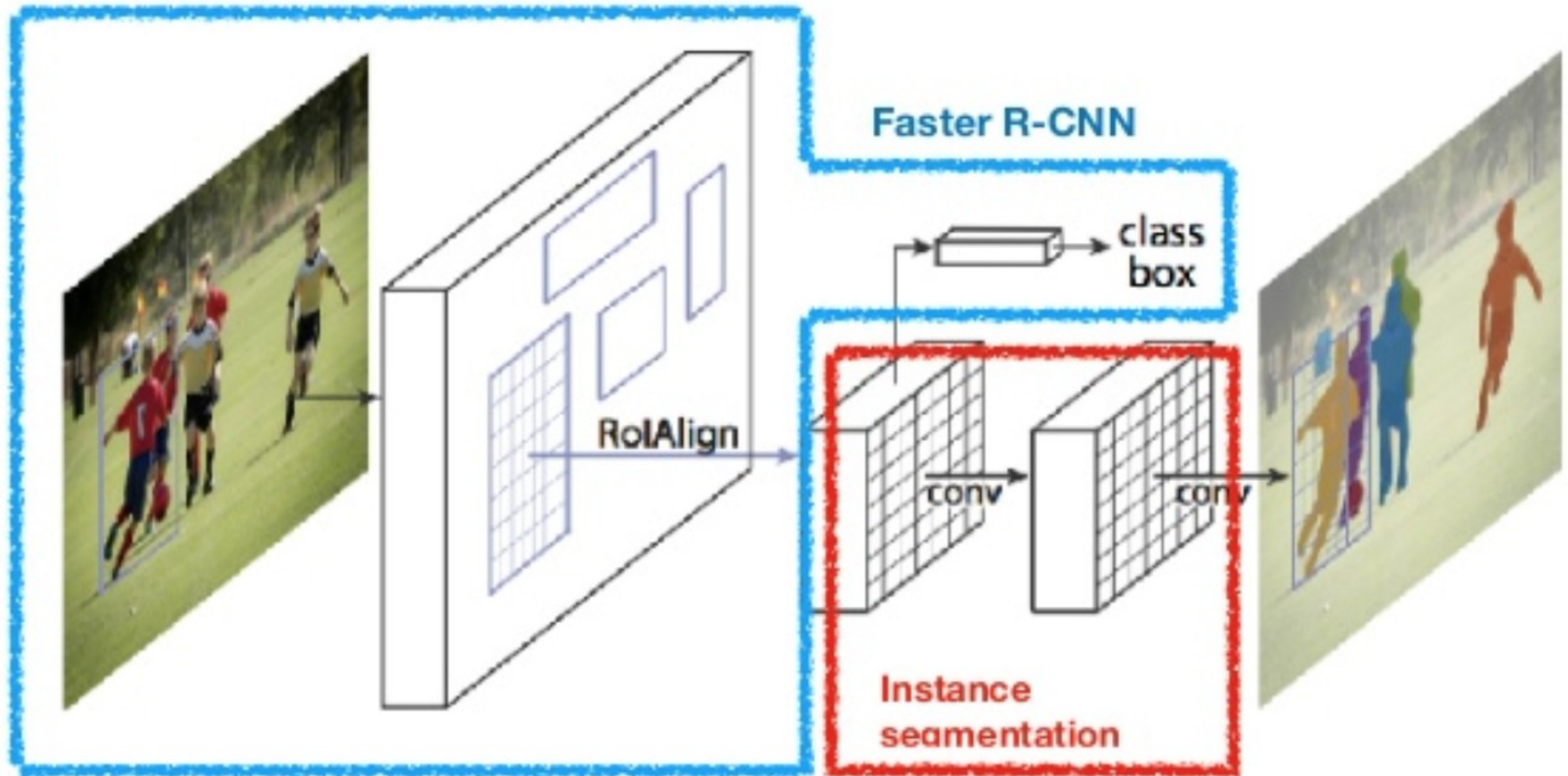
U-net



Mask RCNN, RetinaNet

These networks generate bounding boxes and semantic segmentation maps simultaneously

They can be trained on images having labels for only one or both types of output

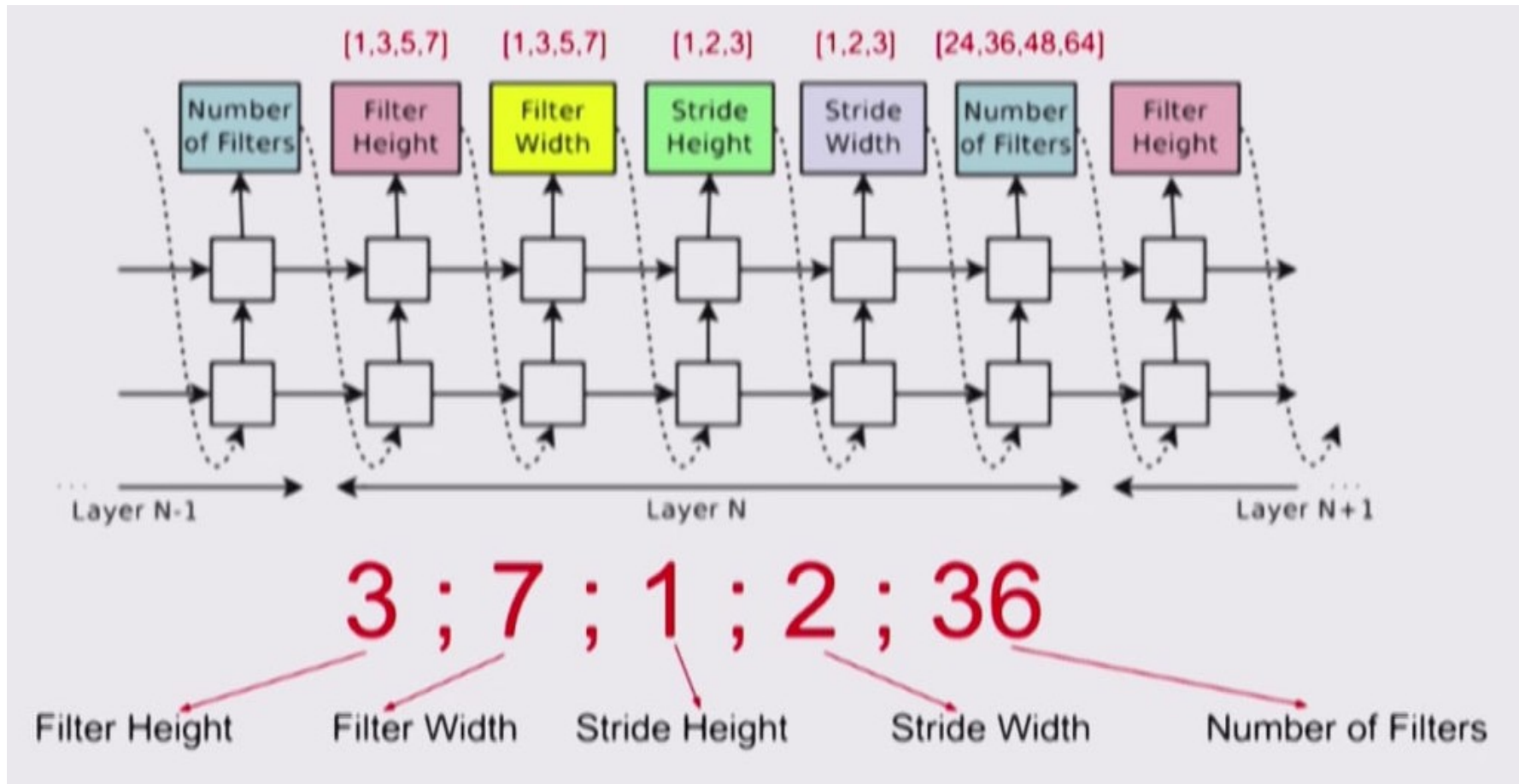


These networks generate bounding boxes and semantic segmentation maps simultaneously



Starting from scratch

Neural architecture search:
 Networks can be described as a series of
 operations
 As series of words → text



Starting from scratch

Neural architecture search:
 Networks can be described as a series of operations
 As series of words → text

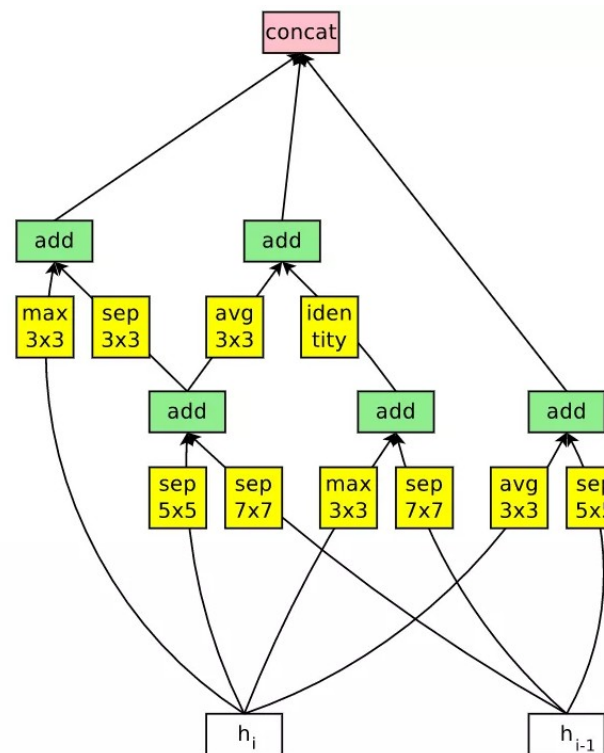
The parameters of each layer can be described as numbers
 The input(s)/output(s) of the layer can be IDs

The whole network can be described as a graph

```

layers {
  bottom: "conv1"
  top: "conv1"
  name: "relu0"
  type: RELU
}
layers {
  bottom: "conv1"
  top: "ccc1"
  name: "ccc1"
  type: CONVOLUTION
  blobs_lr: 1
  blobs_lr: 2
  convolution_param {
    num_output: 96
    kernel_size: 1
    stride: 1
  }
}

```



Starting from scratch

Neural architecture search:
 Networks can be described as a series of operations
 As series of words \rightarrow text

The parameters of each layer can be described as numbers
 The input(s)/output(s) of the layer can be IDs

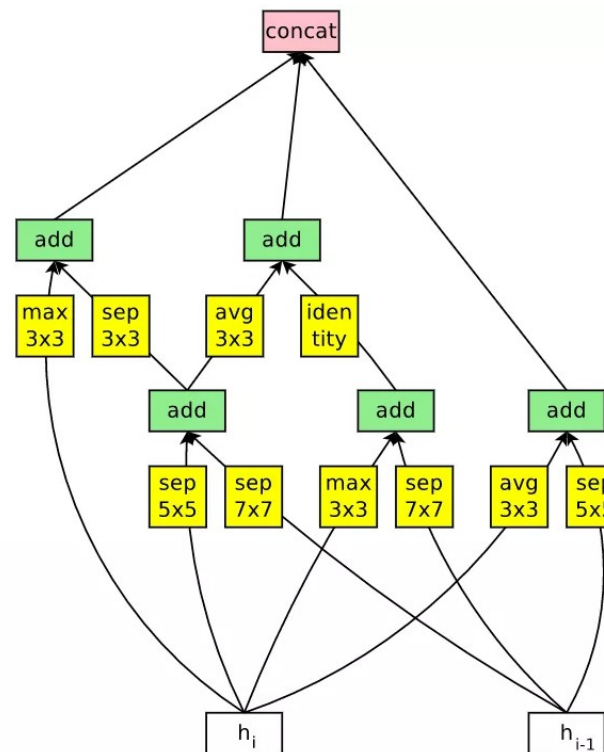
The whole network can be described as a graph

We have a problem space where we have text as an input and an accuracy number as an output

```

layers {
  bottom: "conv1"
  top: "conv1"
  name: "relu0"
  type: RELU
}
layers {
  bottom: "conv1"
  top: "conv1"
  name: "conv1"
  type: CONVOLUTION
  blobs_lr: 1
  blobs_lr: 2
  convolution_param {
    num_output: 96
    kernel_size: 1
    stride: 1
  }
}

```



Starting from scratch

Neural architecture search:
 Networks can be described as a series of operations
 As series of words \rightarrow text

The parameters of each layer can be described as numbers
 The input(s)/output(s) of the layer can be IDs

The whole network can be described as a graph

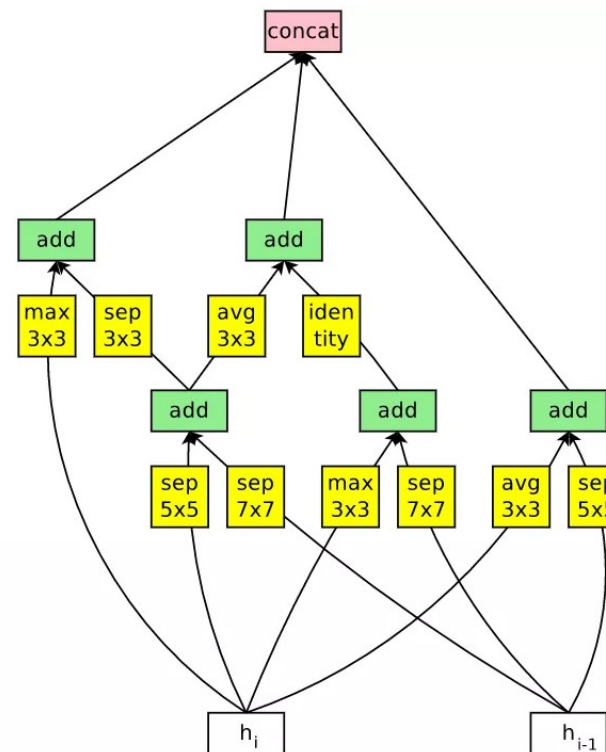
We have a problem space where we have text as an input and an accuracy number as an output

We can train an RNN for regression, which approximates the accuracy of a given network

```

layers {
  bottom: "conv1"
  top: "conv1"
  name: "relu0"
  type: RELU
}
layers {
  bottom: "conv1"
  top: "conv1"
  name: "conv1"
  type: CONVOLUTION
  blobs_lr: 1
  blobs_lr: 2
  convolution_param {
    num_output: 96
    kernel_size: 1
    stride: 1
  }
}

```



Starting from scratch

Neural architecture search:
 Networks can be described as a series of
 operations
 As series of words \rightarrow text

We can turn the problem around:

A recurrent network can be trained with
 reinforcement learning which can train a
 network with predefined accuracy on a
 given dataset.

This recurrent network will understand the
 effect of the elements on this dataset

Test accuracy On CIFAR-10:
 96.35%

Best pervious accuraccy:
 96.26

This architecture os also 1.05 times faster
 (less computations)

