

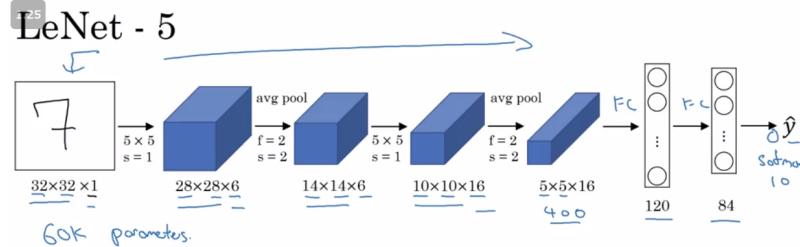
case studies

• classic network

- three inspiring neural network

• LeNet-5

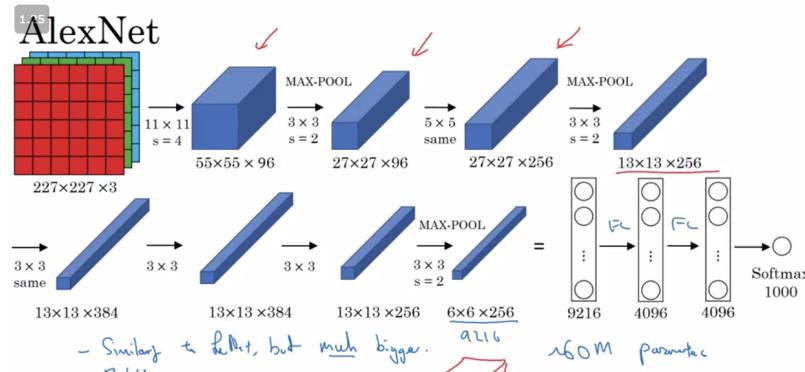
- publish in 1998
- task: identify handwritten digits in a $32 \times 32 \times 1$ gray image.



- $N_h, N_w \downarrow N_c \uparrow$
- come up with **the popular architecture**: Conv ==> Pool ==> Conv ==> Pool ==> FC ==> FC ==> softmax
 - #when in 1998 softmax isn't widely used, the author used another method to make multi-task classification.
 - #at that time researchers always used average pooling.
 - basic idea: N_h, N_w go down \downarrow , and N_c go up \uparrow

• AlexNet

- publish in 2012
- task: ImageNet challenge which classifies images into 1000 classes.

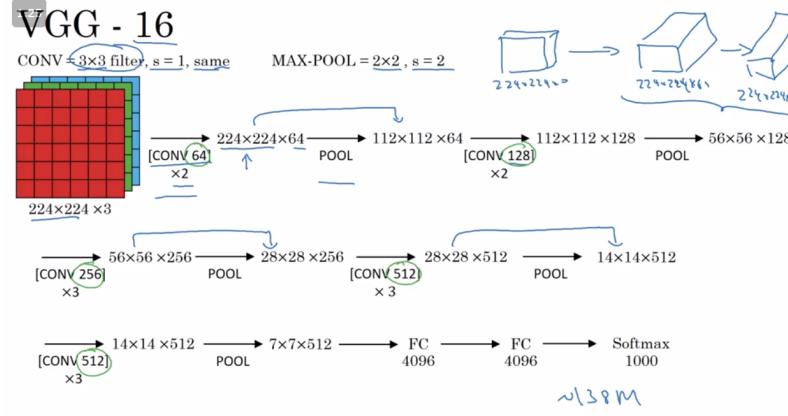


- [Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks] ↪ Andrew Ng
- The original paper contains Multiple GPUs and Local Response normalization (LRN).
 - Multiple GPUs were used because the GPUs were not so fast back then.
 - Researchers proved that Local Response normalization doesn't help much
 - convinced deep learning is important in computer vision.

VGG-16

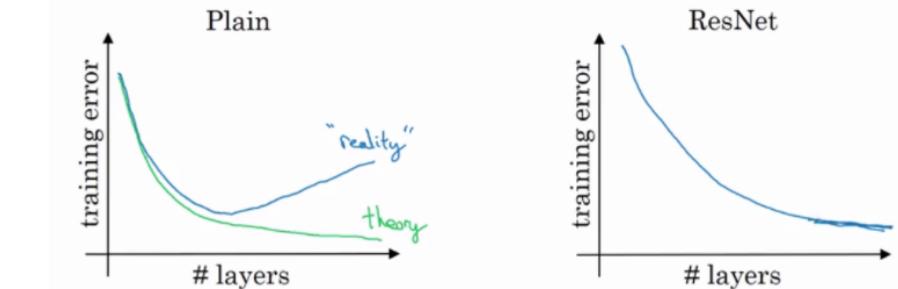
- publish in 2015
- modern one with bigger architecture
- simple NN architecture** → Focus on having only these blocks:

- CONV = 3 X 3 filter, s = 1, same
- MAX-POOL = 2 X 2 , s = 2



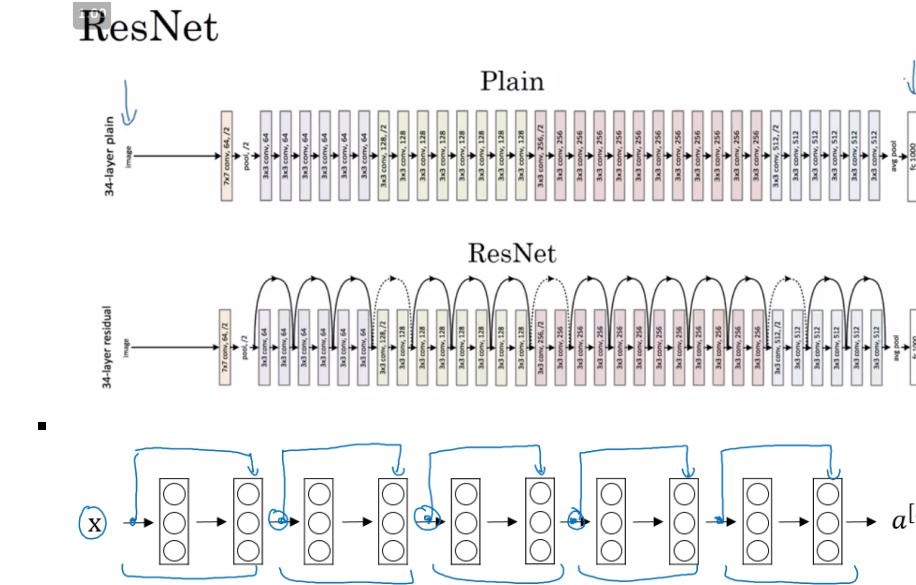
ResNet (Residual Network)

- design for **deeper** neural network.



- In **theory**, as NN grow deeper the training error would go down.
- However in **practice**, because of the vanishing and exploding gradients problems the performance of the network suffers as it goes deeper.
- ResNet can solve such a problem →

- architecture



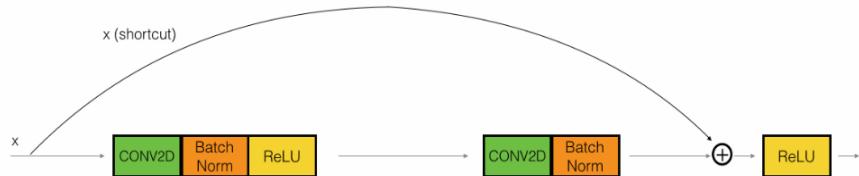


Figure 3: Identity block. Skip connection "skips over" 2 layers.

- **skip-connection**

- $a[1+2] = g(z[1+2] + a[1])$
- $= g(W[1+2] a[1+1] + b[1+2] + a[1])$
- if the dimension of $a[1]$ and $a[1+2]$ doesn't match
 - $a[1+2] = g(z[1+2] + ws * a[1])$ # The added ws make the dimensions equal
 - ws also can be a zero padding.
 - or use 1×1 convolutions to match (which we would talk below)

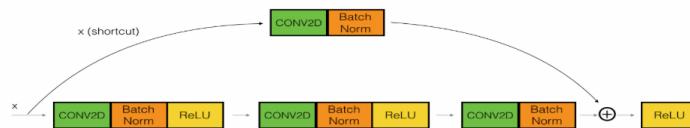
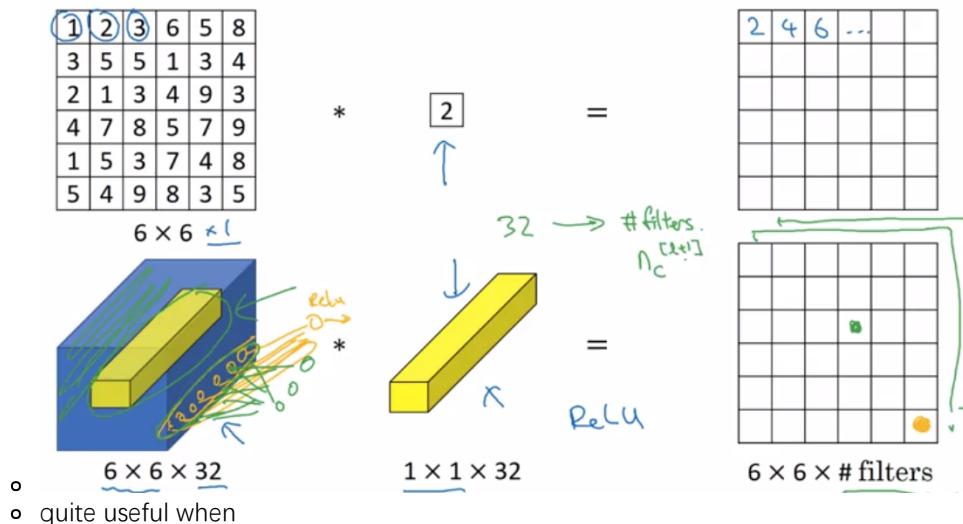


Figure 4: Convolutional block

- why it work?

- firstly, it would not hurt the plain NN performance
 - if using L2 regularization, assume that it lead W and b equal to 0.
 - then $g(W[1+2] a[1+1] + b[1+2] + a[1])$ would equal to $g(a[1])$
 - assume that we are using ReLU, then we have $g(a[1]) = a[1]$
- secondly, it's likely to improve NN performance by learning more features from the residual layer.

• Networks in networks and 1×1 Convolutions



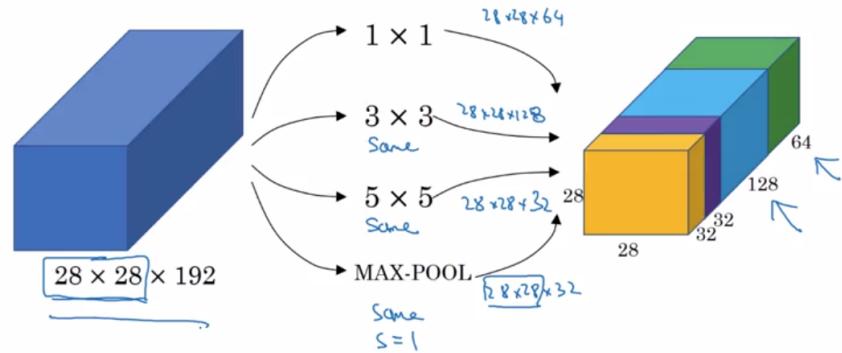
• Inception network

- fun part: where does its name come from?



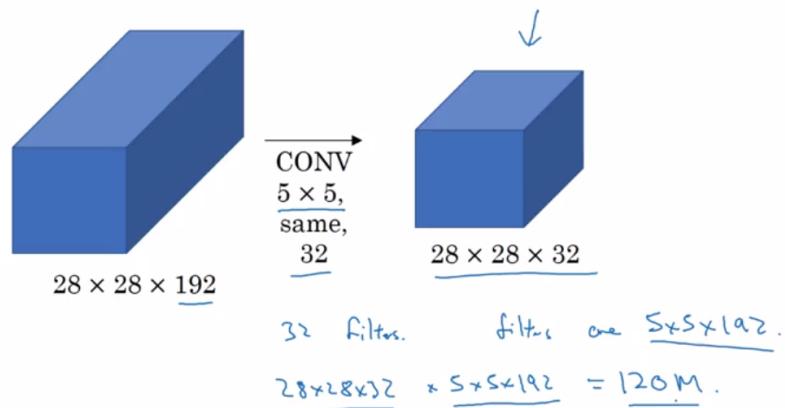
- <http://knowyourmeme.com/memes/we-need-to-go-deeper>
- "don't hesitate for what to use, we implement it all LoL."
- basic blocks below ↗

Motivation for inception network



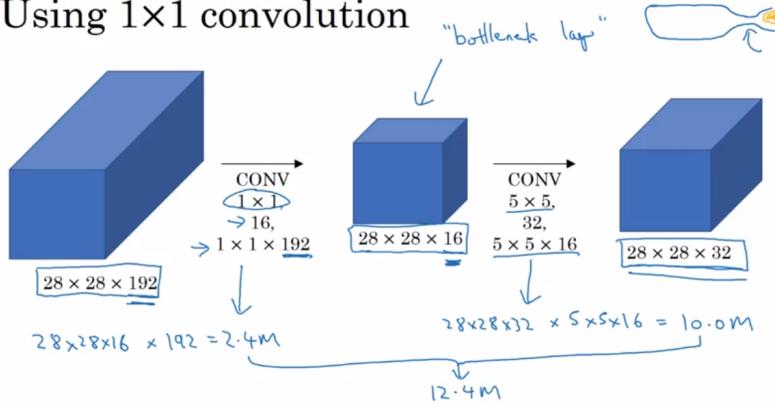
- trick: Using 1 X 1 convolution to reduce computational cost
 - A 1 x 1 Conv here is called **Bottleneck BN**
 - problem:

The problem of computational cost

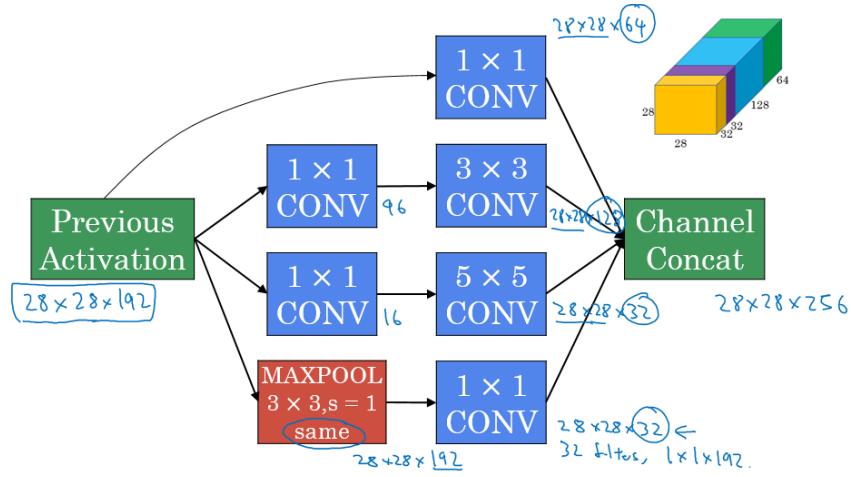


- using bottleneck layer

Using 1×1 convolution

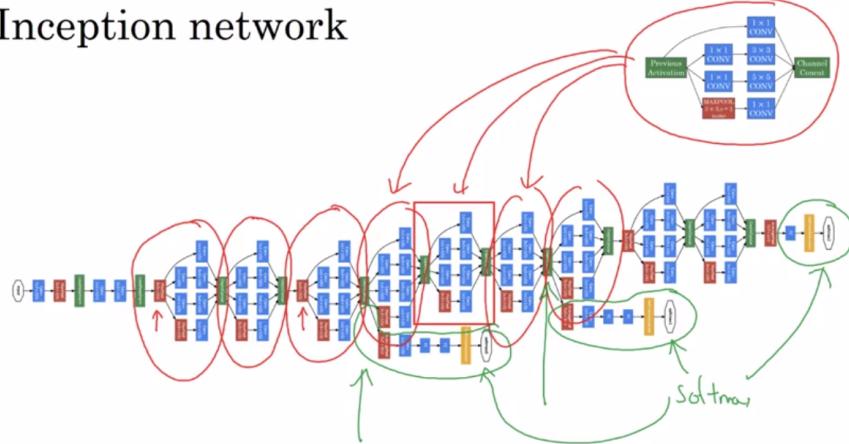


- so now basic model could be:



- architecture of inception NN in paper

Inception network

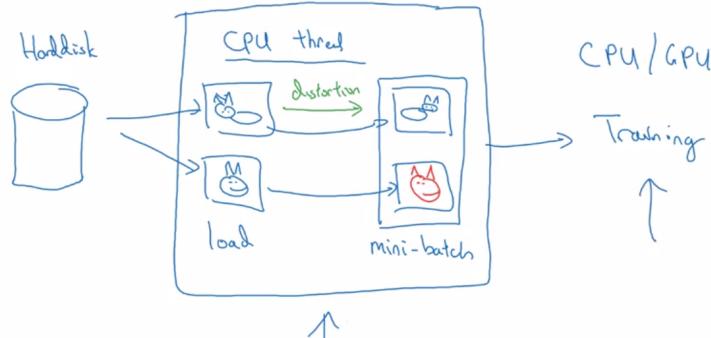


- Max-Pool block (unit in red) is used before the inception module to reduce the dimensions of the inputs.
- 3 Softmax branches :
 - make prediction from hidden layers
 - helps to ensure that the intermediate features are good enough to the network to learn
 - it turns out that softmax0 and softmax1 gives regularization effect.

practical advices

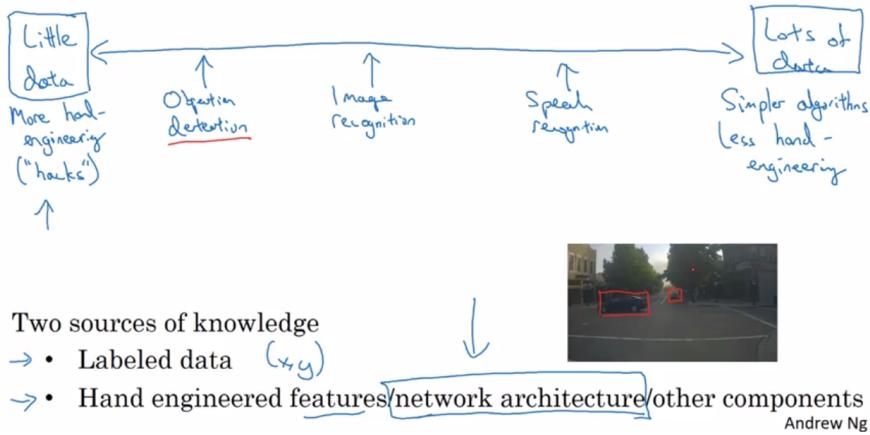
- using open-source implementation (like repo in Github)
- transfer learning
- data augmentation
 - methods:
 - rotation/mirroring/random cropping/shearing
 - color shifting **PCA color augmentation**
 - Ex. parallel threads of DL

Implementing distortions during training



- state of CV
 - "we still don't have enough data for CV tasks like Object Detection."

Data vs. hand-engineering



- tips (though improvement might be quite subtle)

Tips for doing well on benchmarks/winning competitions

Ensembling

3 - 15 networks



- Train several networks independently and average their outputs

Multi-crop at test time

- Run classifier on multiple versions of test images and average results

10-crop



- use open source code

Use open source code

- Use architectures of networks published in the literature
- Use open source implementations if possible
- Use pretrained models and fine-tune on your dataset