

Cornell Bowers C-IS

College of Computing and Information Science

Modern Convolutional Neural Networks

CS4782: Intro to Deep Learning

Review: Convolutional Neural Networks (CNNs)

✓ Convolutions

Maintain spatial relation between pixels
Reduce number of parameters through weight sharing

✓ Pooling

Captures key information from across different areas of the feature maps
Together with convolutions allows for translational invariance

✓ BatchNorm

Increases speed and stability of training



input image

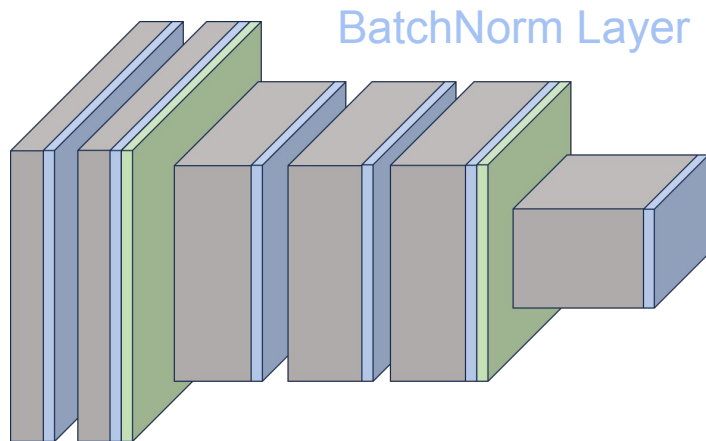
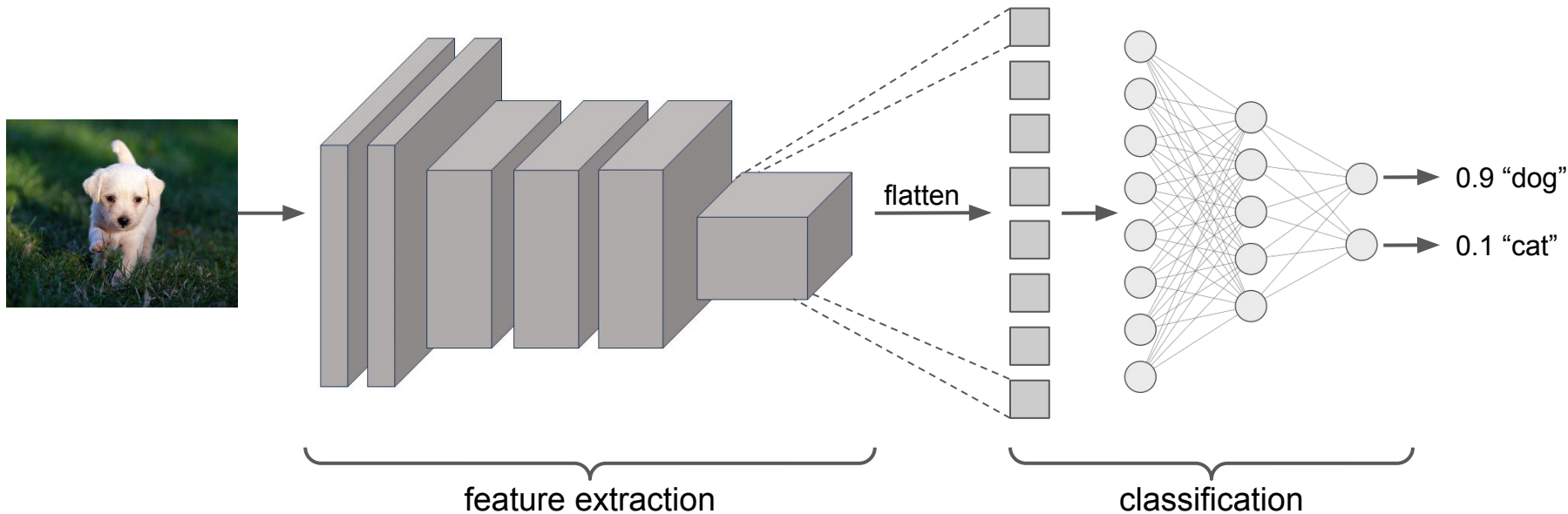


Image Classification

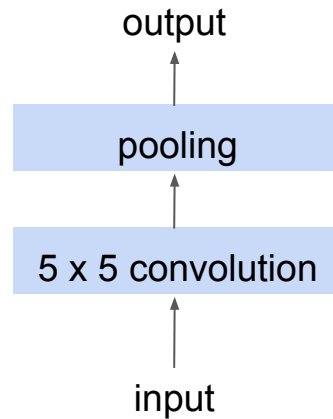
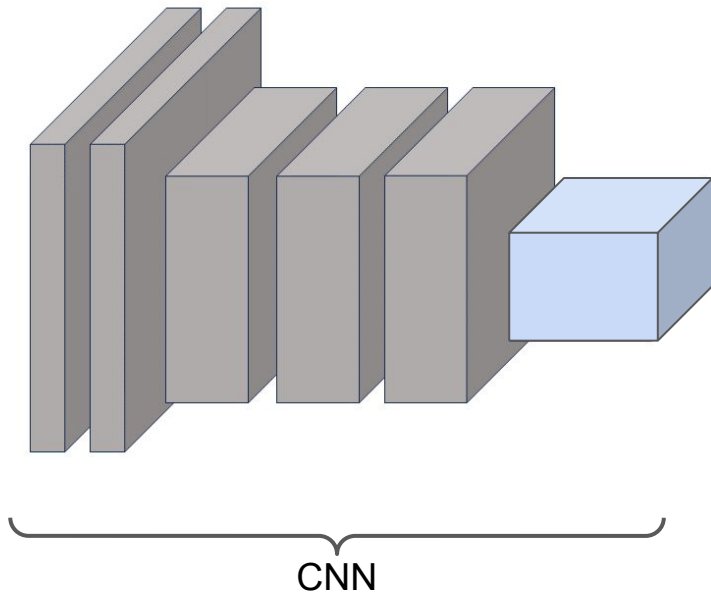
- Important: Everything is differentiable!
- Can calculate gradient of the loss with backpropagation
 - Train with SGD/Adam/etc.
 - Learn convolutional filters and classification head end-to-end!



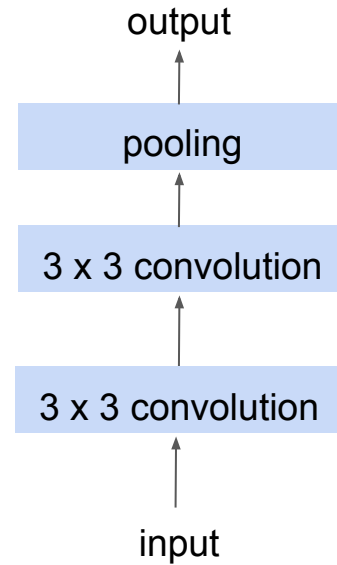
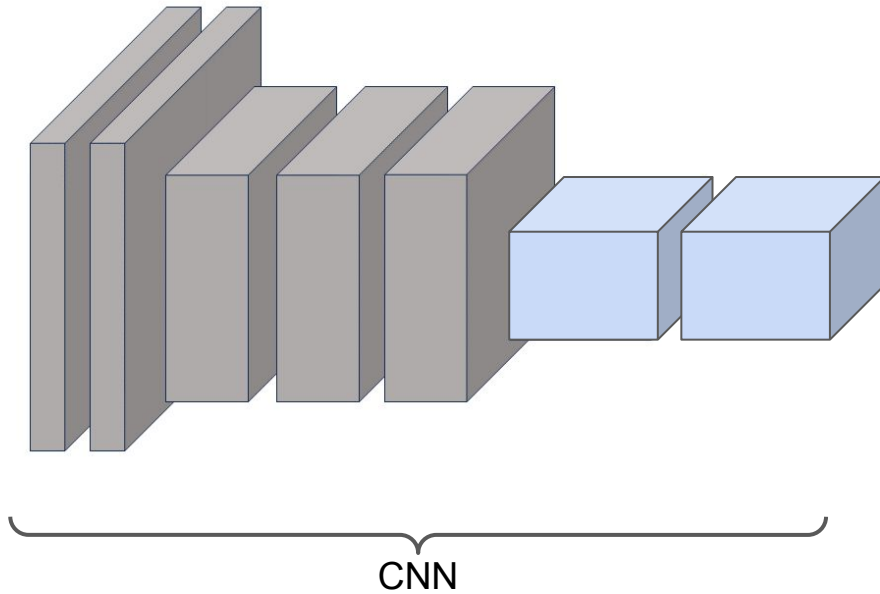
Discuss: Padding

- Given a 5x5 feature map and a 3x3 convolution:
 - How much padding do I need to maintain the spatial size of the feature map (i.e., 5x5)?
- What about when using a 5x5 convolution?

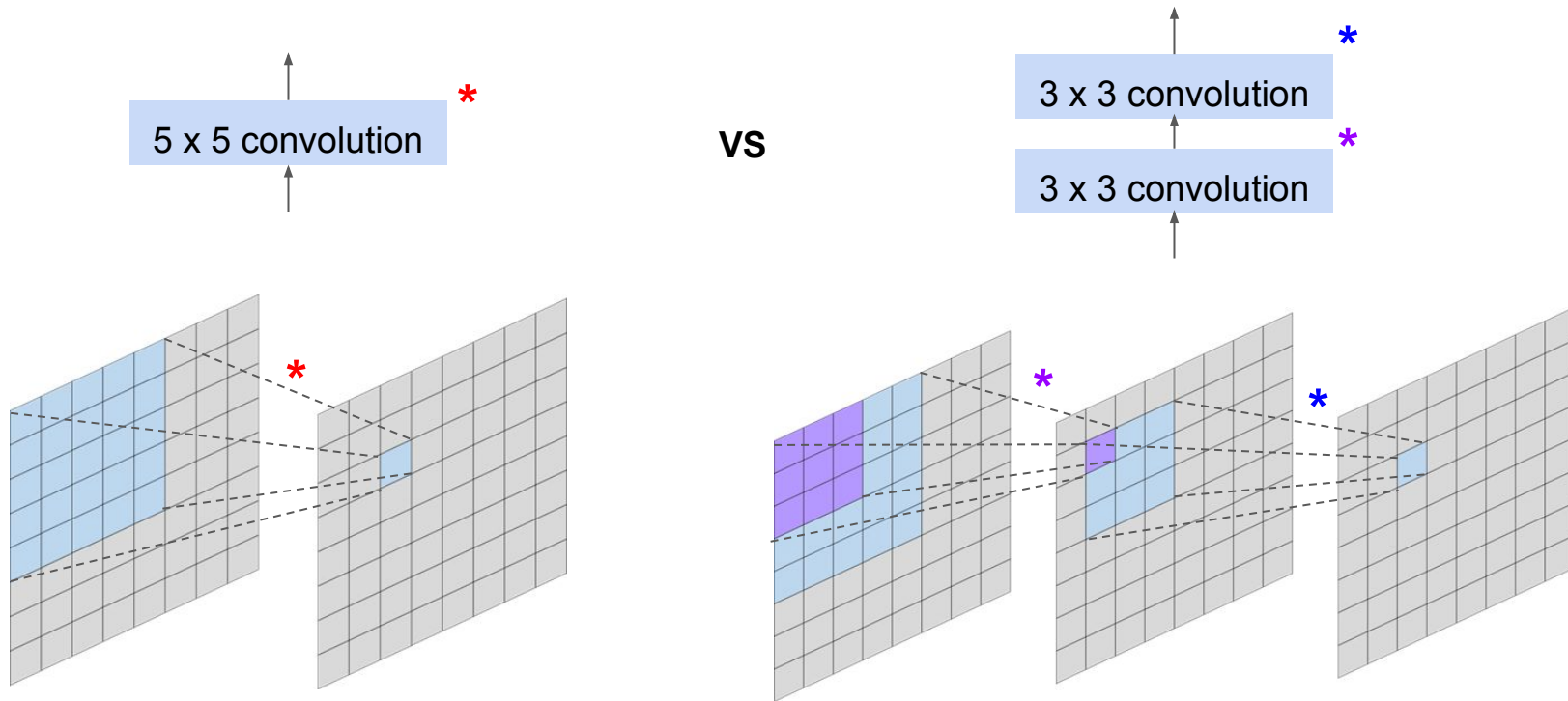
Deeper CNN Architectures



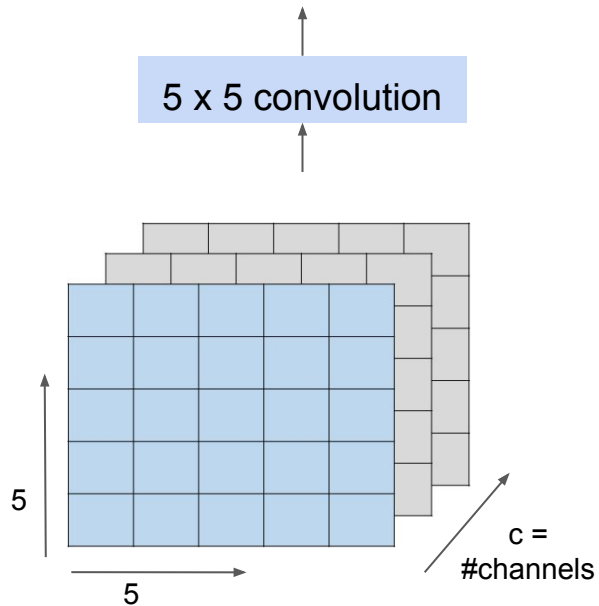
Deeper CNN Architectures



Deeper CNN Architectures

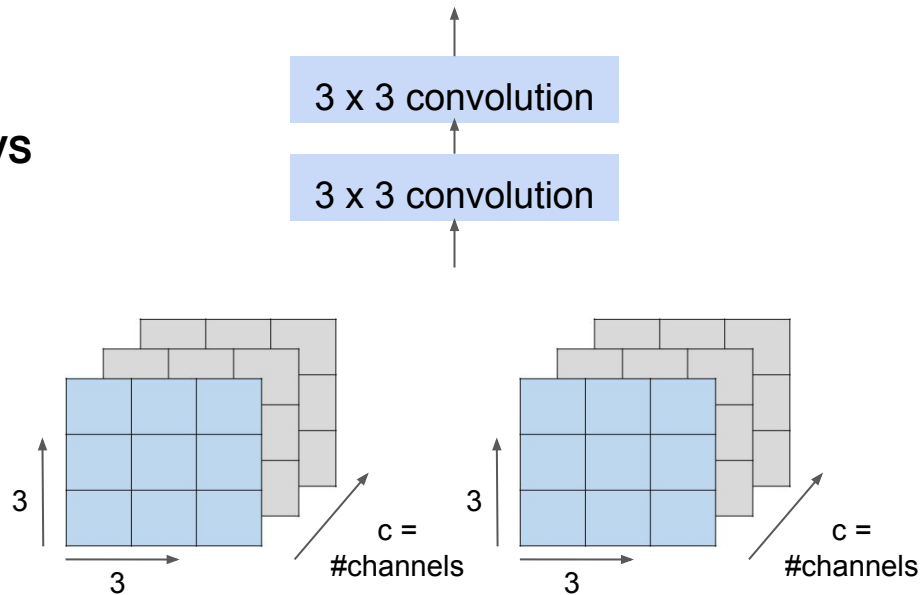


Deeper CNN Architectures



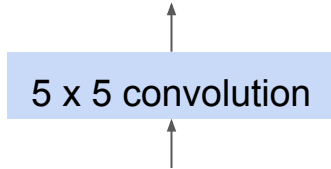
$$5 * 5 * c^2 = 25c^2 \text{ parameters}$$

VS

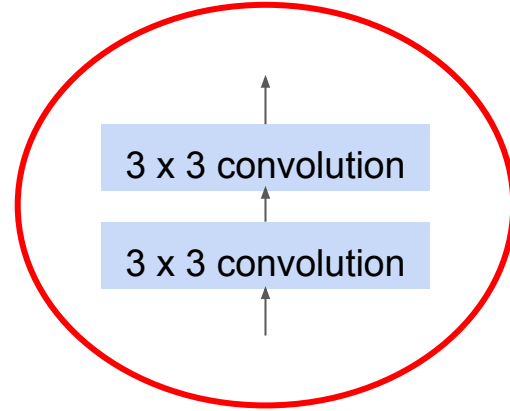


$$2 * 3 * 3 * c^2 = 18c^2 \text{ parameters}$$

Deeper CNN Architectures

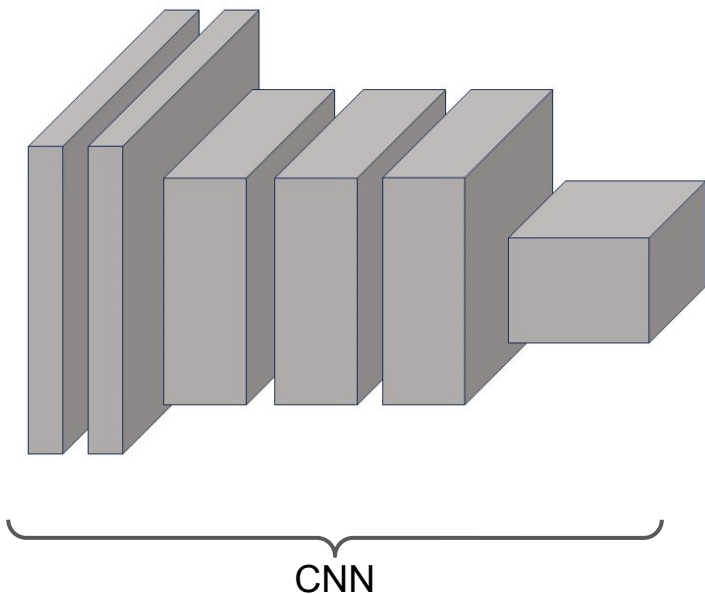


VS

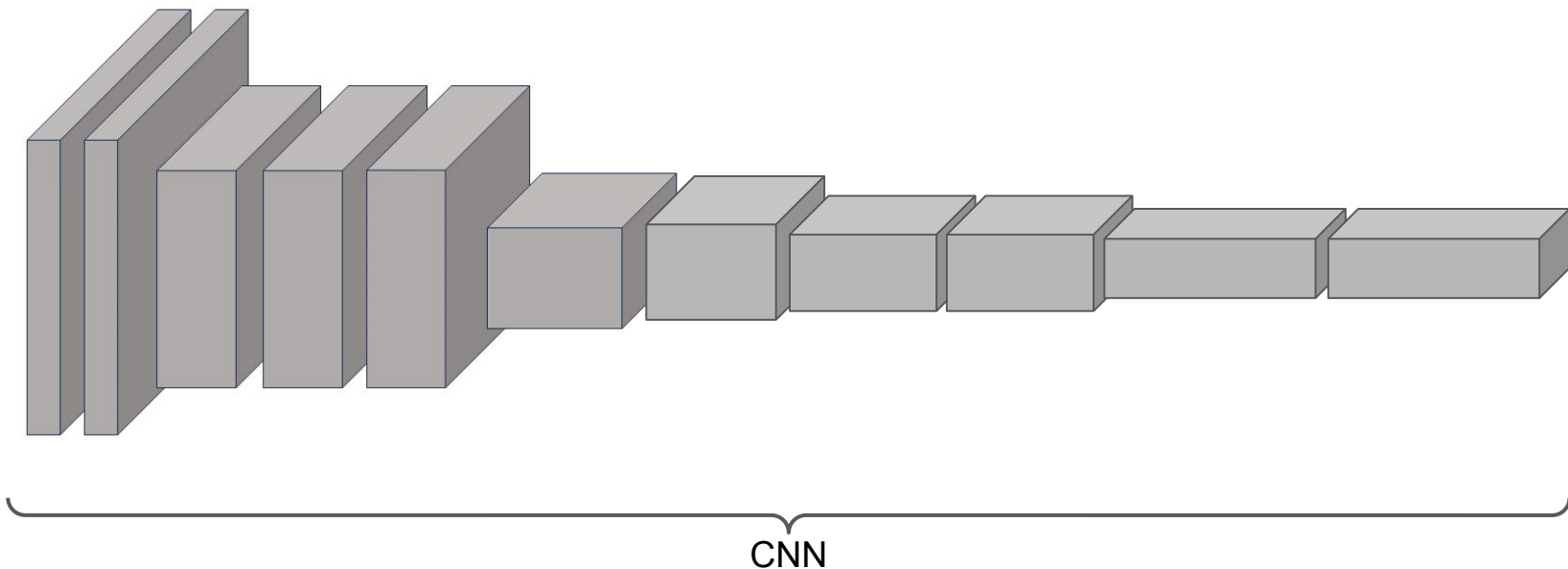


Performed better!

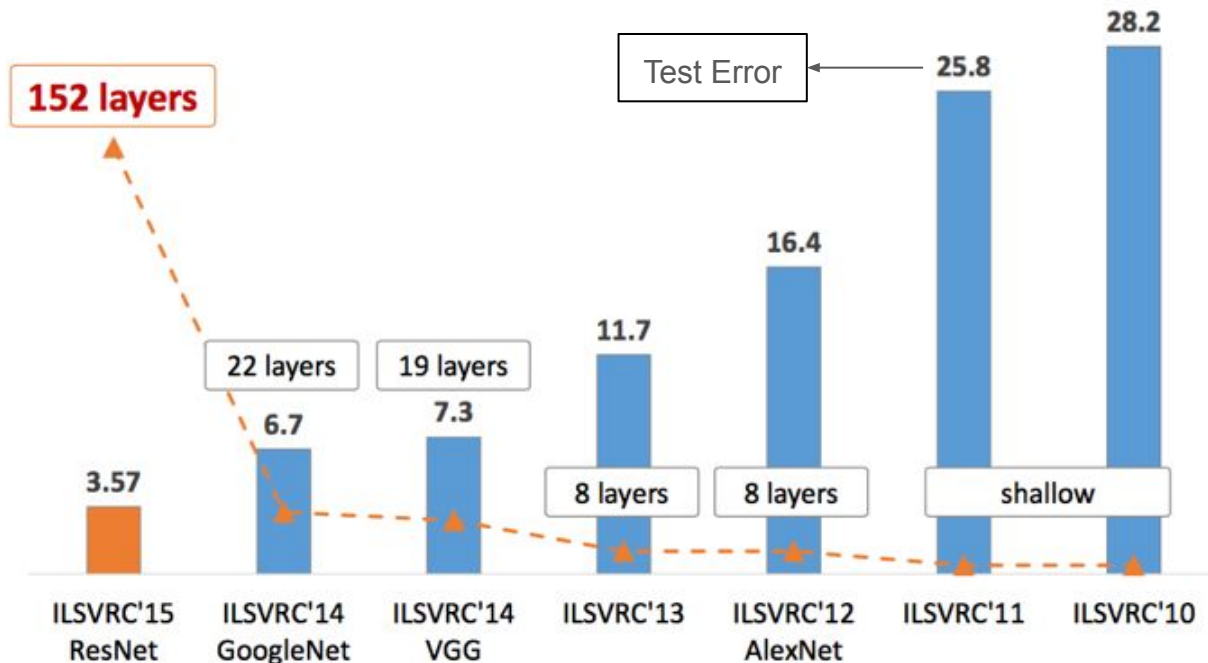
Deeper == better



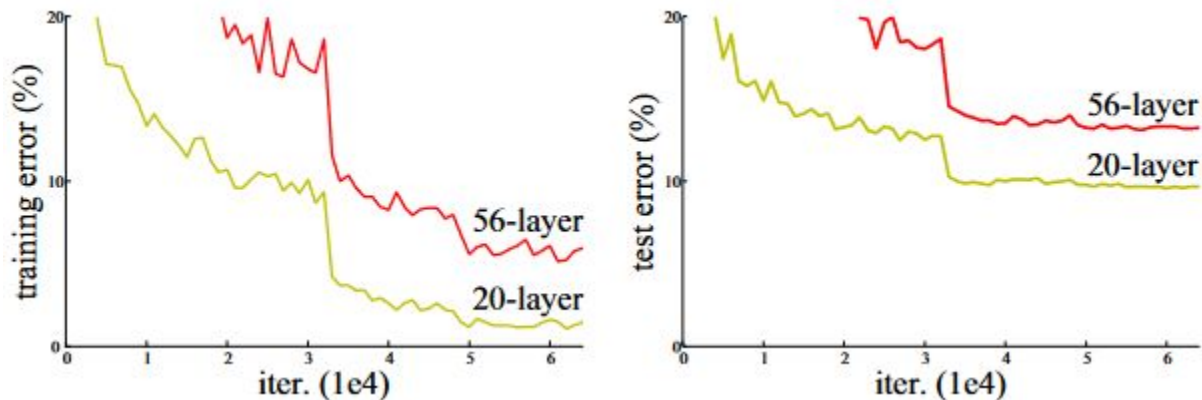
Deeper == better



ImageNet Classification Challenge: Deeper == better



Deeper == better?



56 layer CNN has higher training and test error than 20 layer CNN on CIFAR-10 dataset for image classification

Deeper != better

- Long training times
- Vanishing gradient problem
 - Recall backpropagation to update weights

$$\frac{\partial z}{\partial z_i} = \frac{\partial z}{\partial z_{n-1}} \frac{\partial z_{n-1}}{\partial z_{n-2}} \cdots \frac{\partial z_{i+1}}{\partial z_i}$$

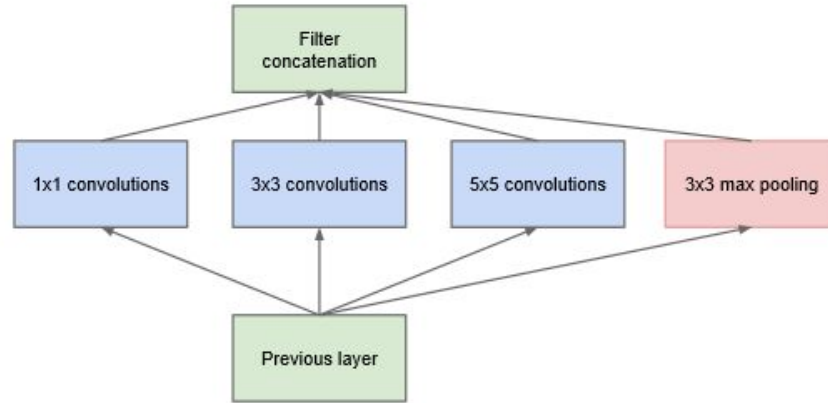
- If each term $\lll 1$, gradient “vanishes” as the entire multiplication goes towards 0
- => Weights not updated properly

GoogLeNet/Inception Net

Goal: given a fixed computational budget, optimize the depth and width of the network

=> Deeper networks with computational efficiency

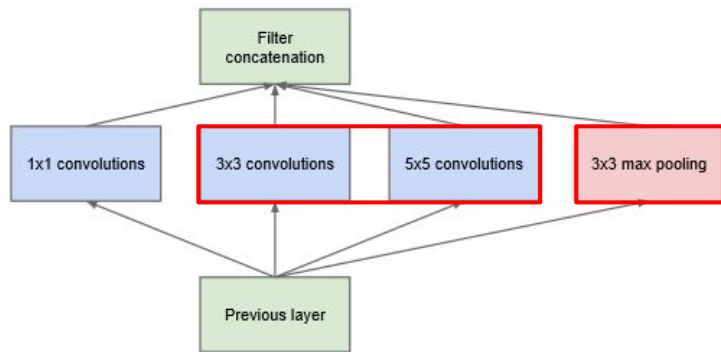
Inception Module



Inception module = main
building blocks

Inception Module

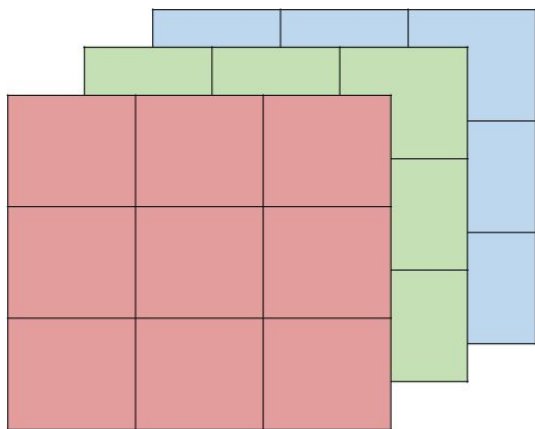
Still expensive!



- 3x3 and 5x5 convolutions have large number of operations
- Output of pooling layer increases the output channel dimension when concatenated

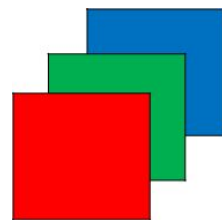
Slight Detour: 1x1 convolutions

input



3x3x3

filter

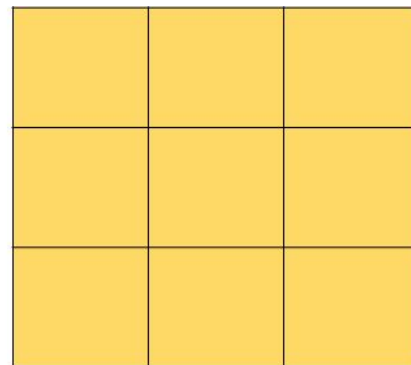


1x1x3

*

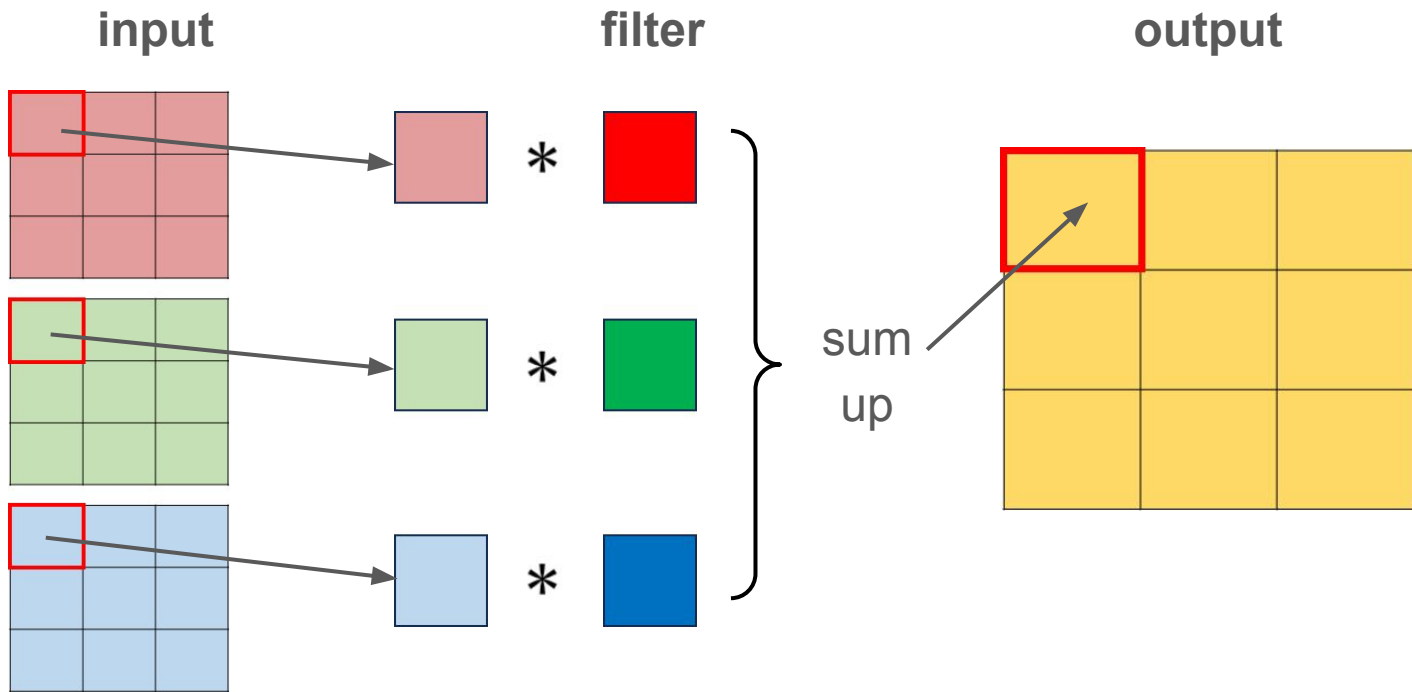
=

output

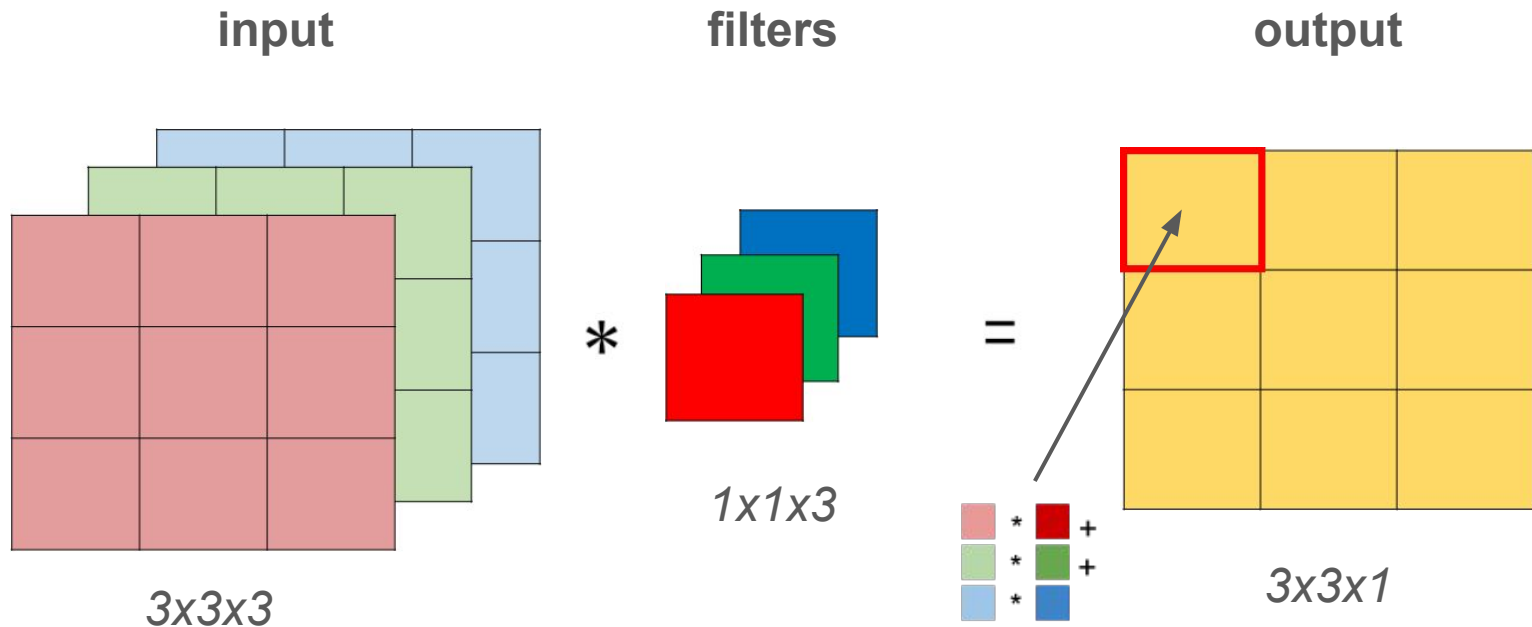


3x3x1

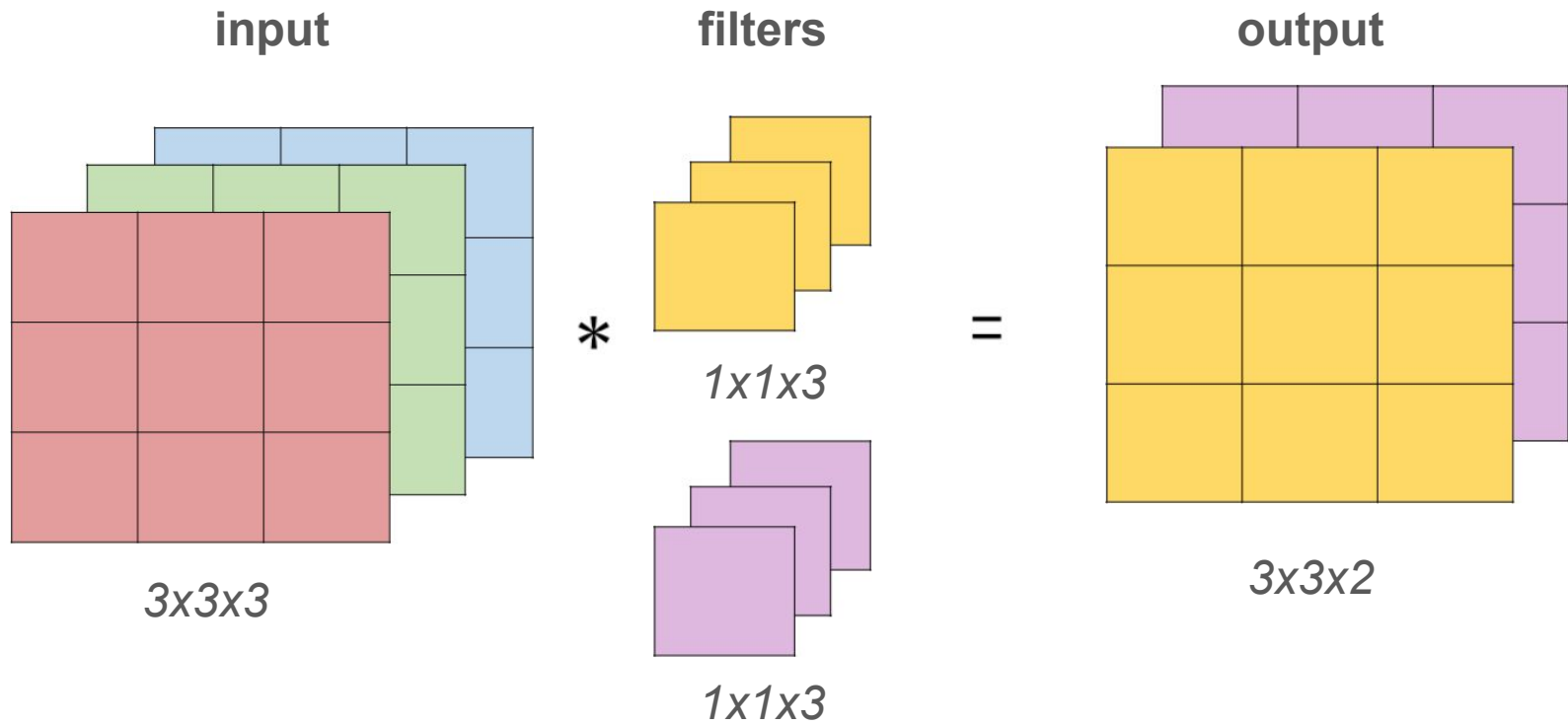
Slight Detour: 1x1 convolutions



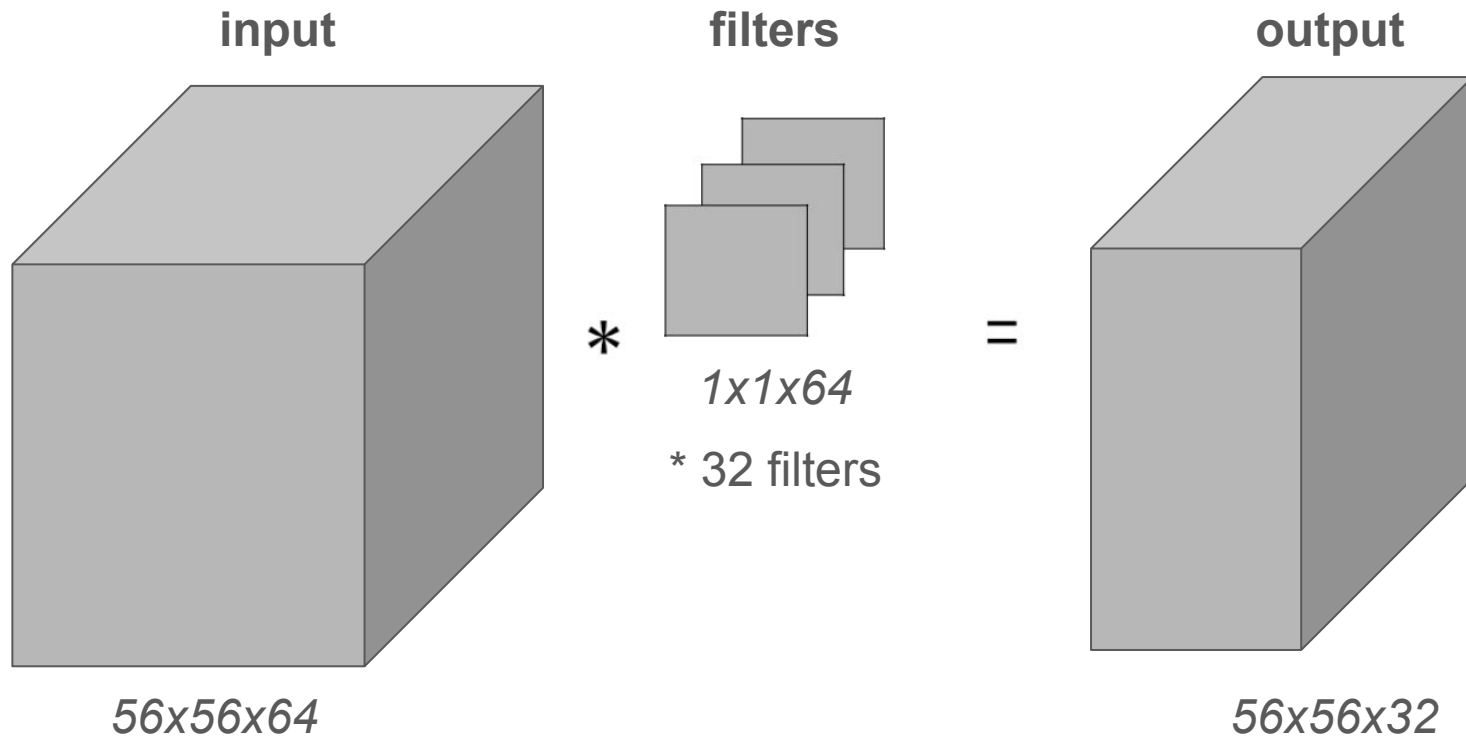
Slight Detour: 1x1 convolutions



Slight Detour: 1x1 convolutions

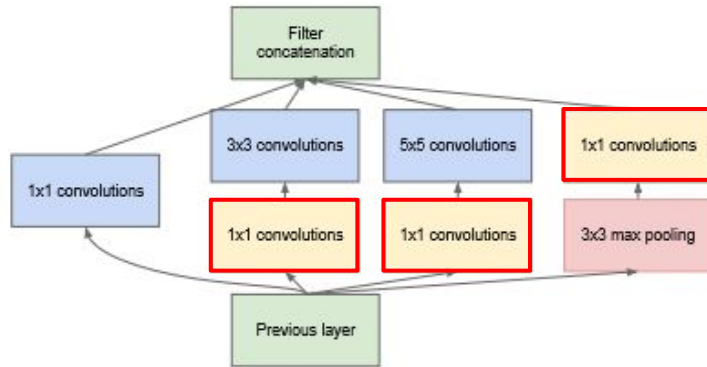


Slight Detour: 1x1 convolutions



Inception Module

Solution: Inception module with dimension reduction



- “Bottleneck” with 1x1 convolutions to reduce dimensions

Discuss: Impact of Dimension Reduction

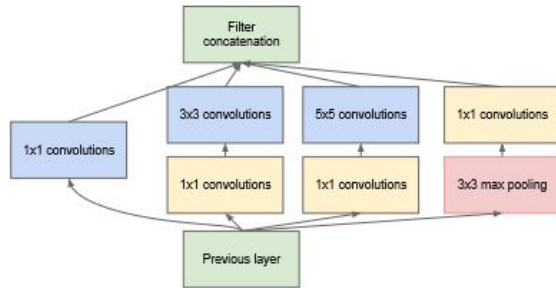
Assume you have an input feature map with 256 dimensions.

Compare the parameter counts from:

1. 3x3 conv with 256 filters
2. 1x1 conv with 64 filters \rightarrow 3x3 conv with 64 filters \rightarrow 1x1 conv with 256 filters

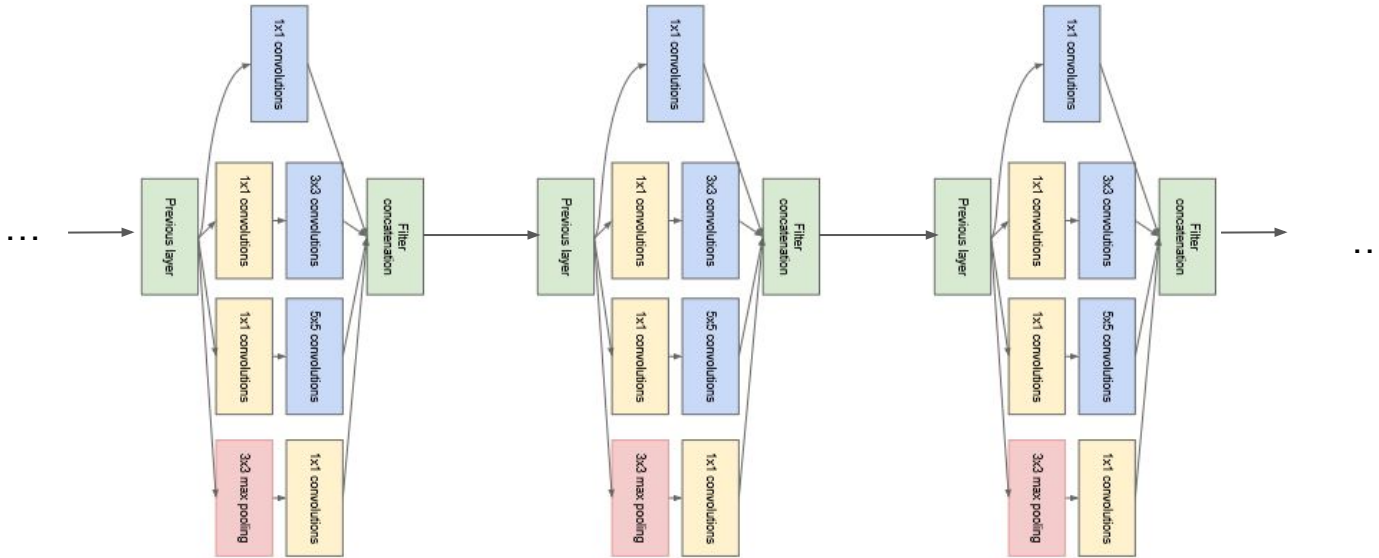
GoogLeNet Architecture

Key idea: stack inception modules together



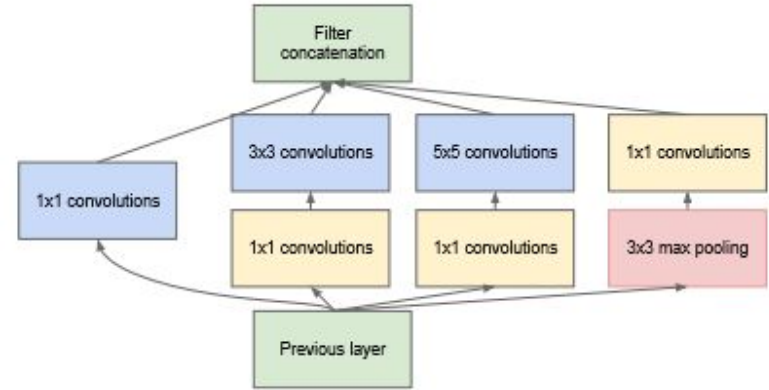
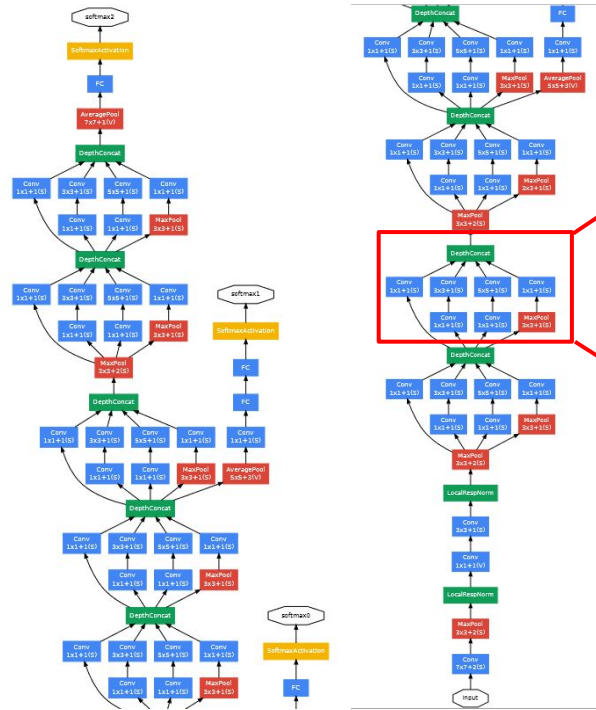
GoogLeNet Architecture

Key idea: stack inception modules together



[Szegedy, Christian, et al. "Going deeper with convolutions." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015.]

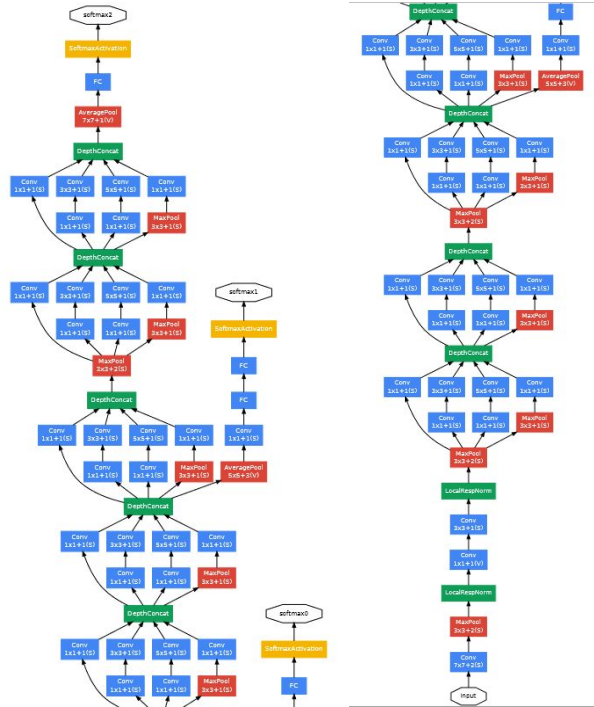
The Entire GoogLeNet Architecture



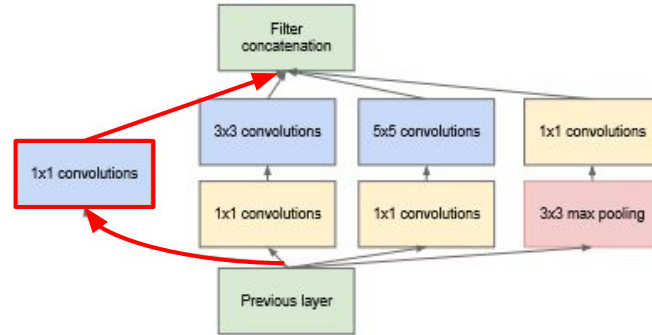
Inception Module

[Szegedy, Christian, et al. "Going deeper with convolutions." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015.]

The Entire GoogleNet Architecture



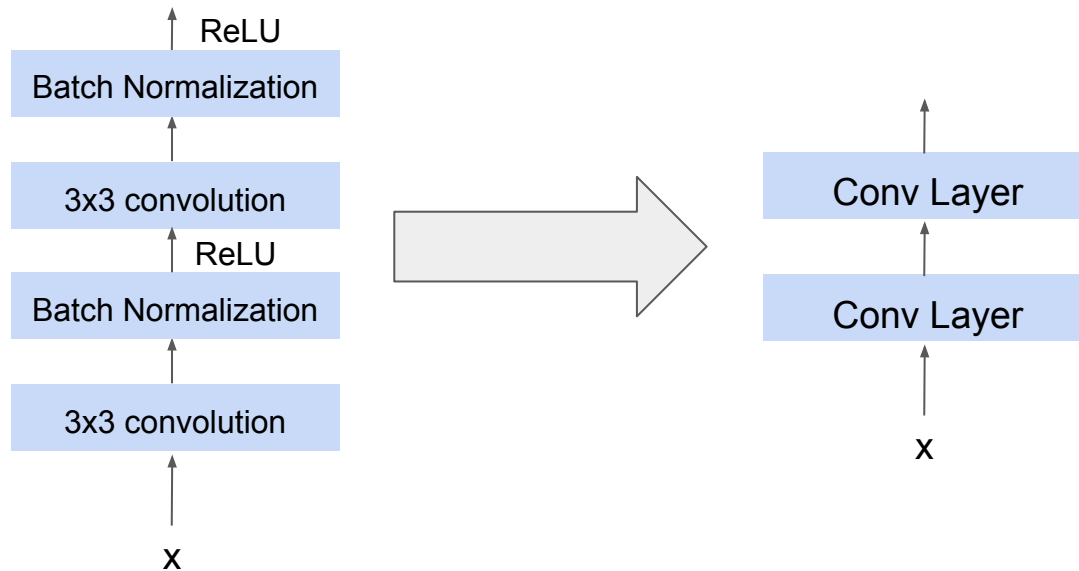
Very complicated - how exactly did this architecture solve the problem?



Residual connections

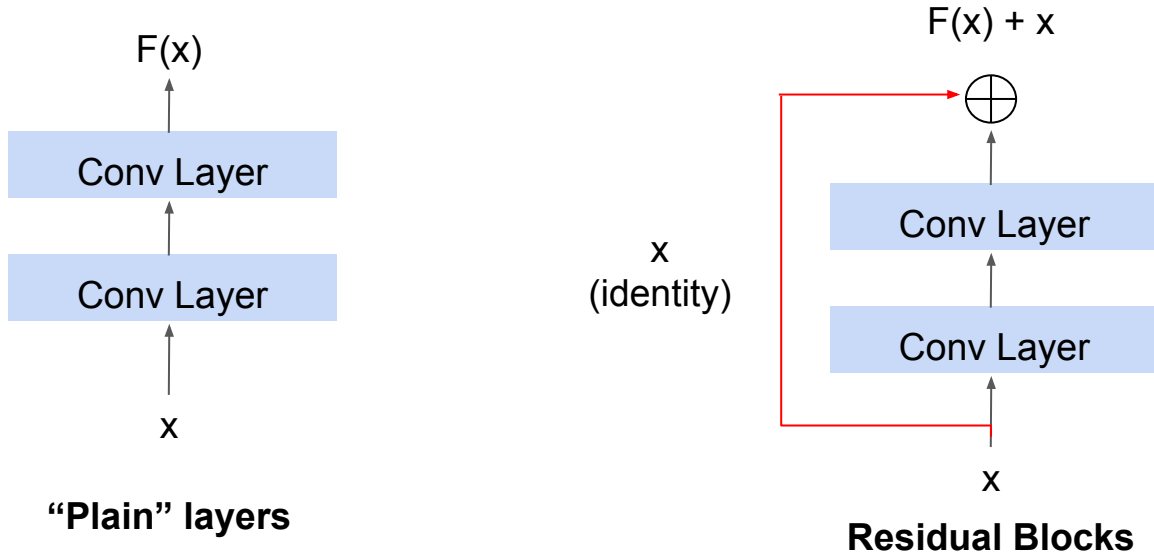
[Szegedy, Christian, et al. "Going deeper with convolutions." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015.]

Aside: Conv Layer Abstraction

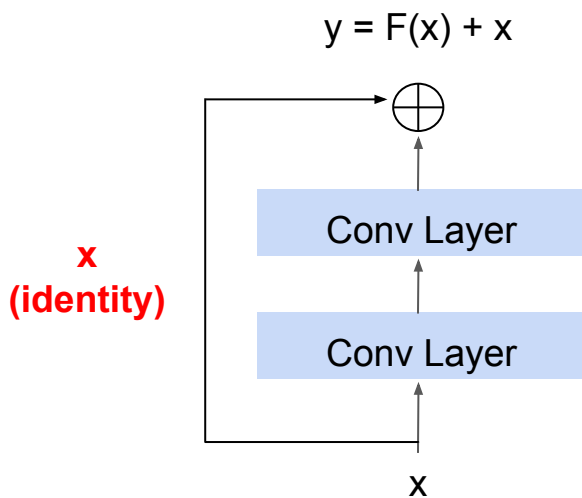


Residual Connections

aka skip connections - add an identity mapping to the output function



Residual Connections

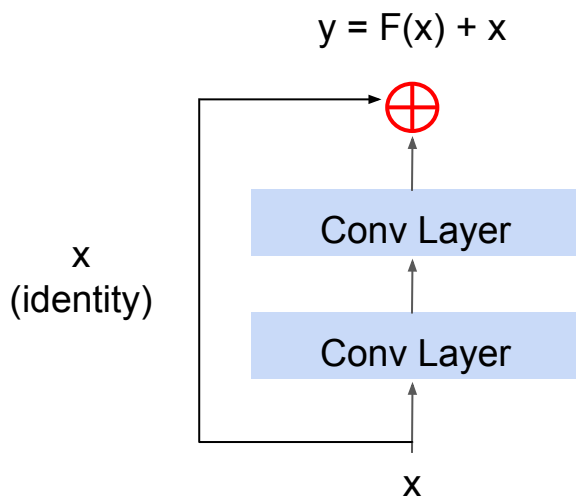


Residual Blocks

Identity mapping

- can propagate features forward
- only learn *difference* of feature maps

Residual Connections



Residual Blocks

Additive component of identity

- alleviates vanishing gradients

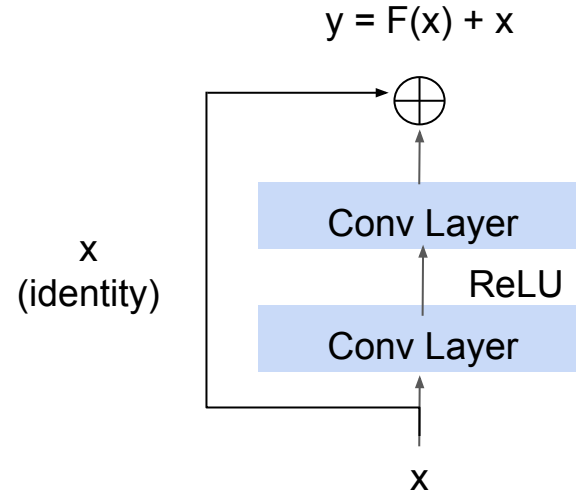
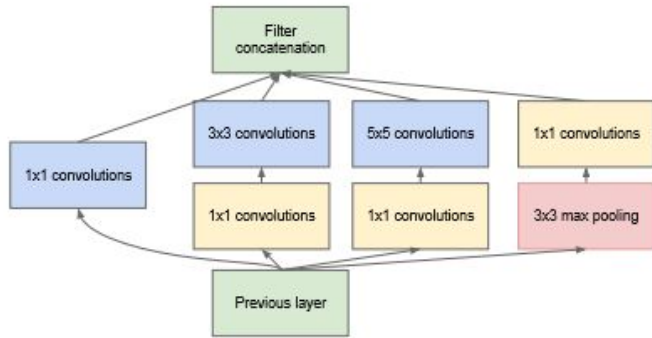
$$\frac{\delta L}{\delta x} = \frac{\delta L}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta L}{\delta y} (F'(x))$$

Plain

$$\frac{\delta L}{\delta x} = \frac{\delta L}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta L}{\delta y} (1 + F'(x))$$

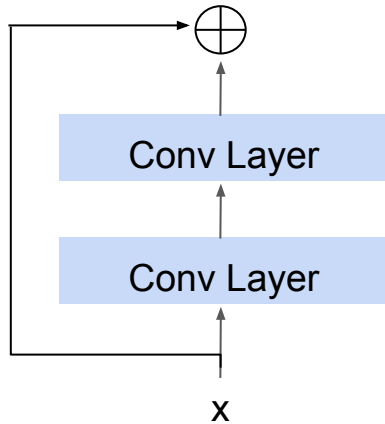
ResNet

Discuss: Spot (and explain) the difference



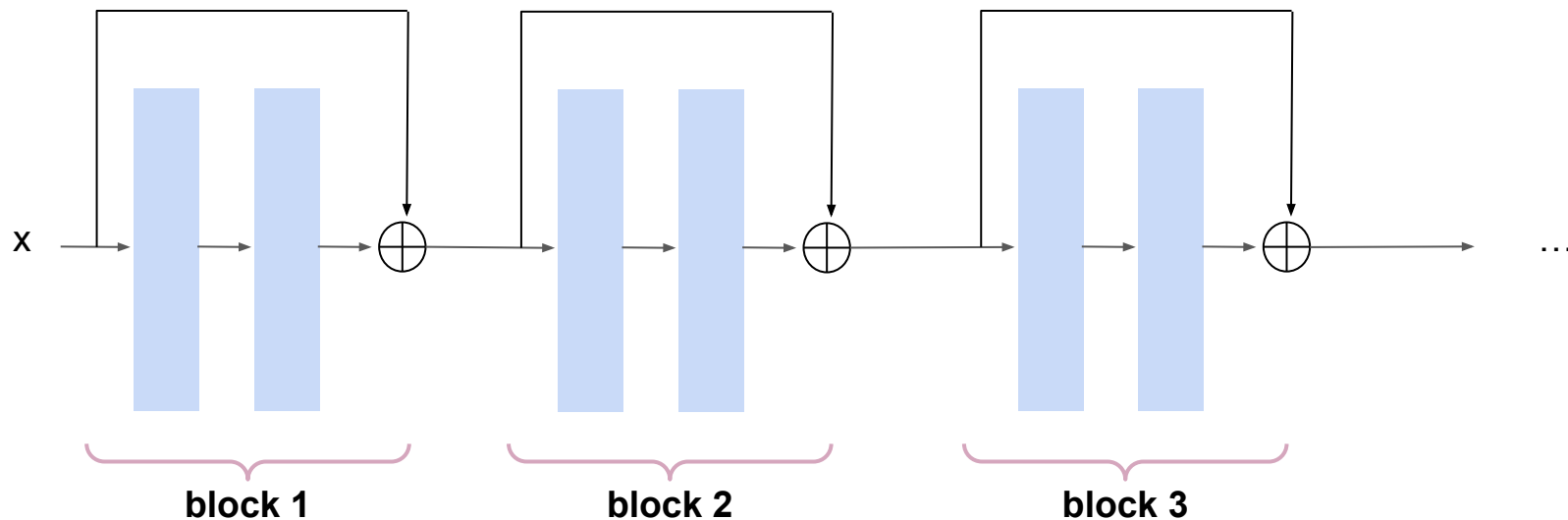
ResNet

Stack residual blocks together!



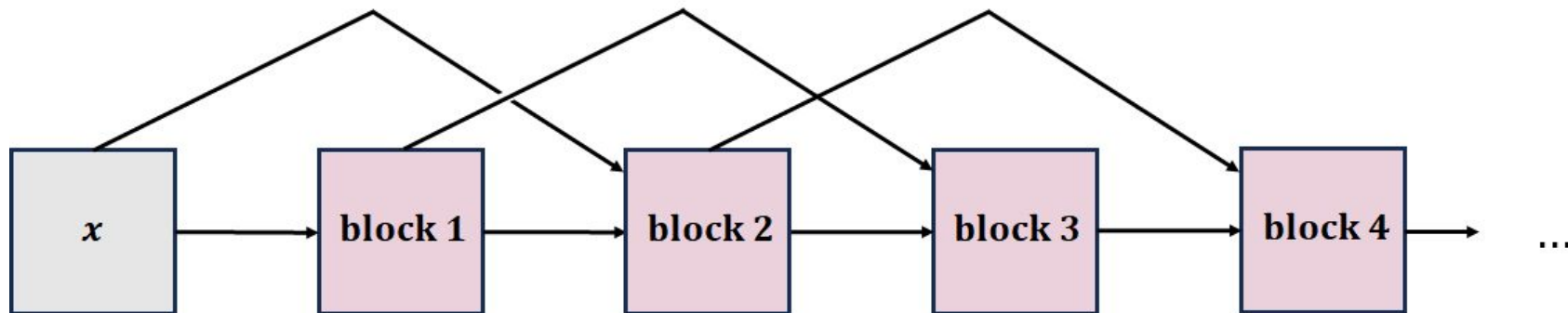
ResNet

Stack residual blocks together!



ResNet

Stack residual blocks together!

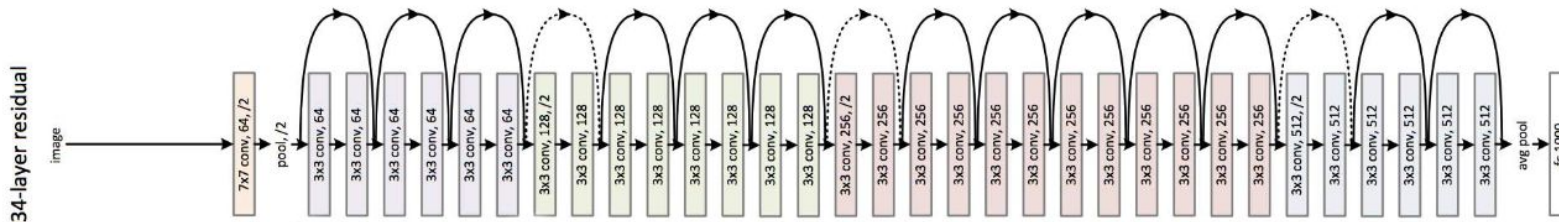


Full ResNet Architecture

“Plain” Network



ResNet



[He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.]

Deeper == better

Can train deeper models!

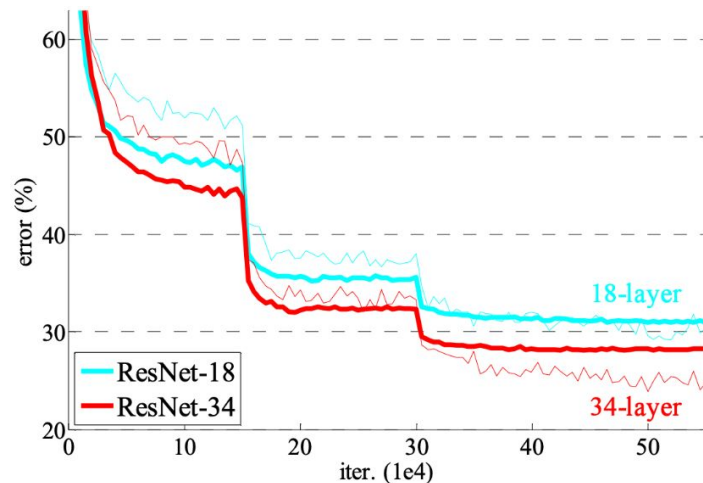
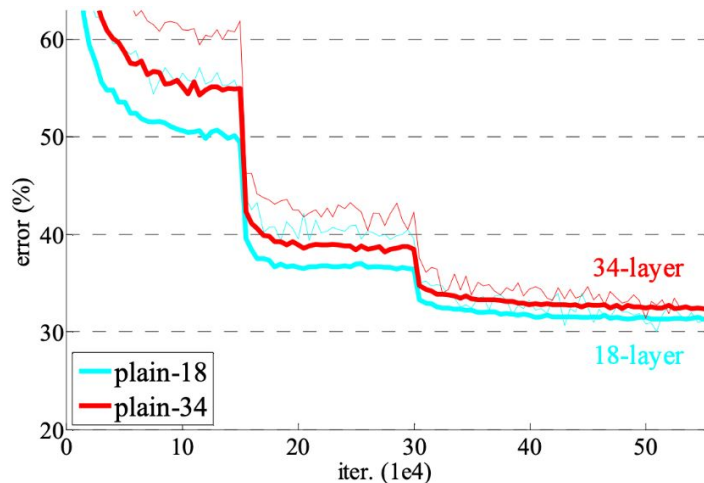
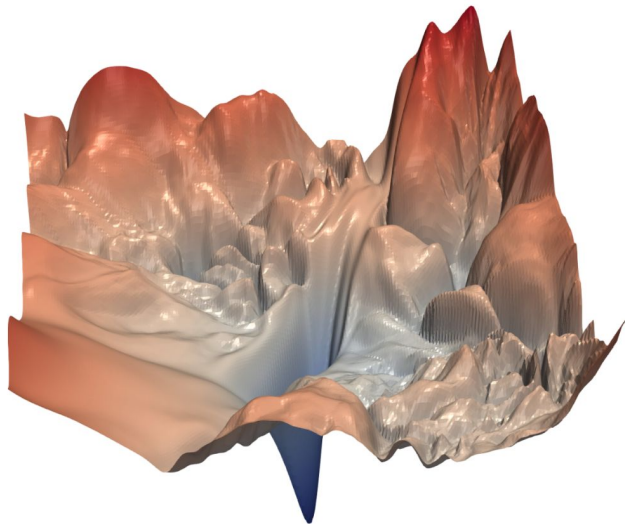


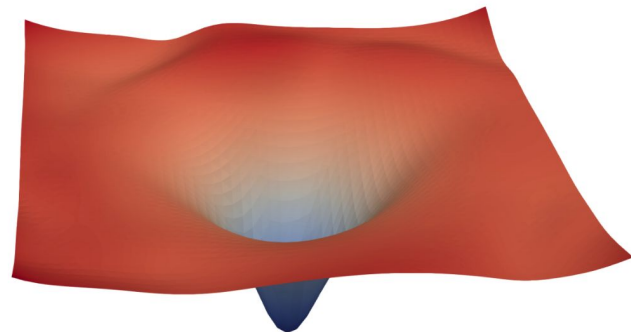
Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

Visualizing the Effect of Skip Connections

Makes optimization easier!



(a) without skip connections



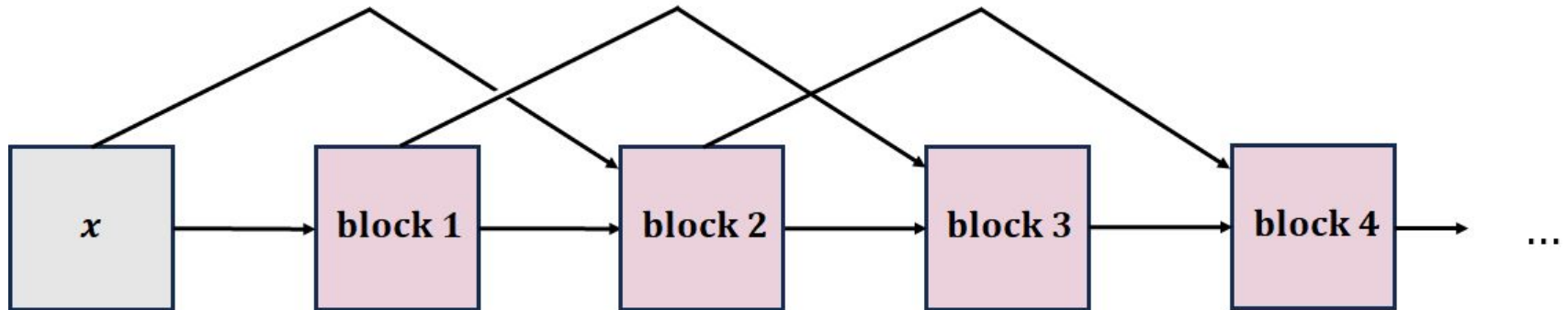
(b) with skip connections

Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

Stochastic Depth

During training, randomly drop Residual Blocks using skip connections

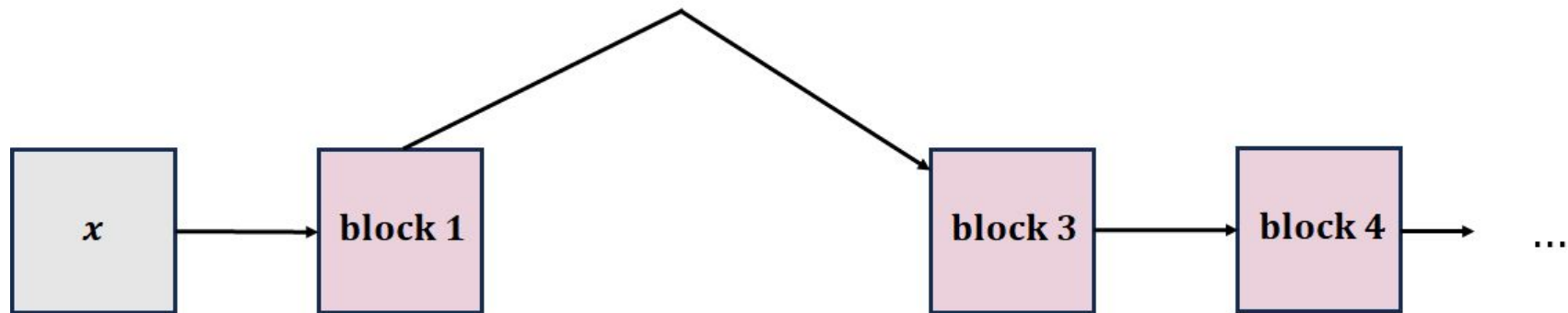
Like dropout but with residual blocks instead of individual neurons



Stochastic Depth

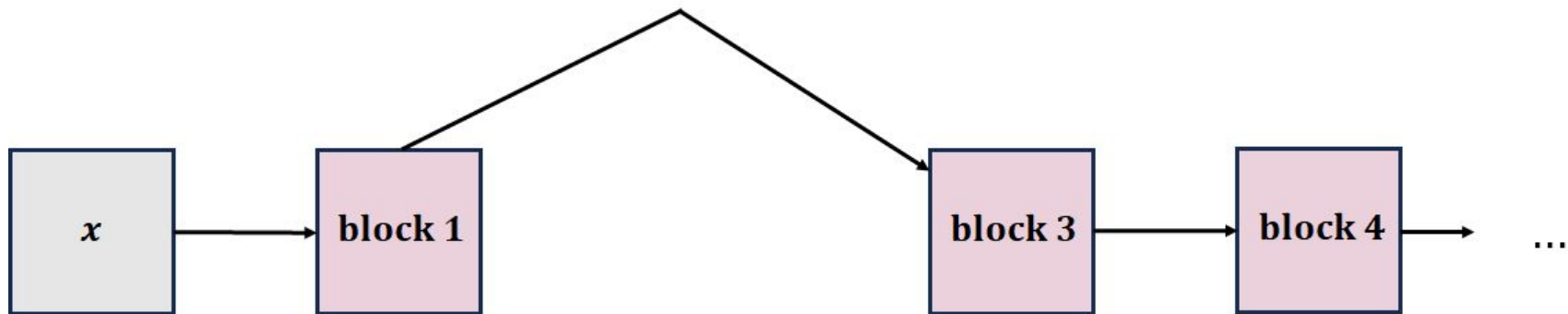
During training, randomly drop Residual Blocks using skip connections

Like dropout but with residual blocks instead of individual neurons



Stochastic Depth

Another benefit: robustness/mitigating overfitting



Stochastic Depth

Increases training loss, but... decreases test error!

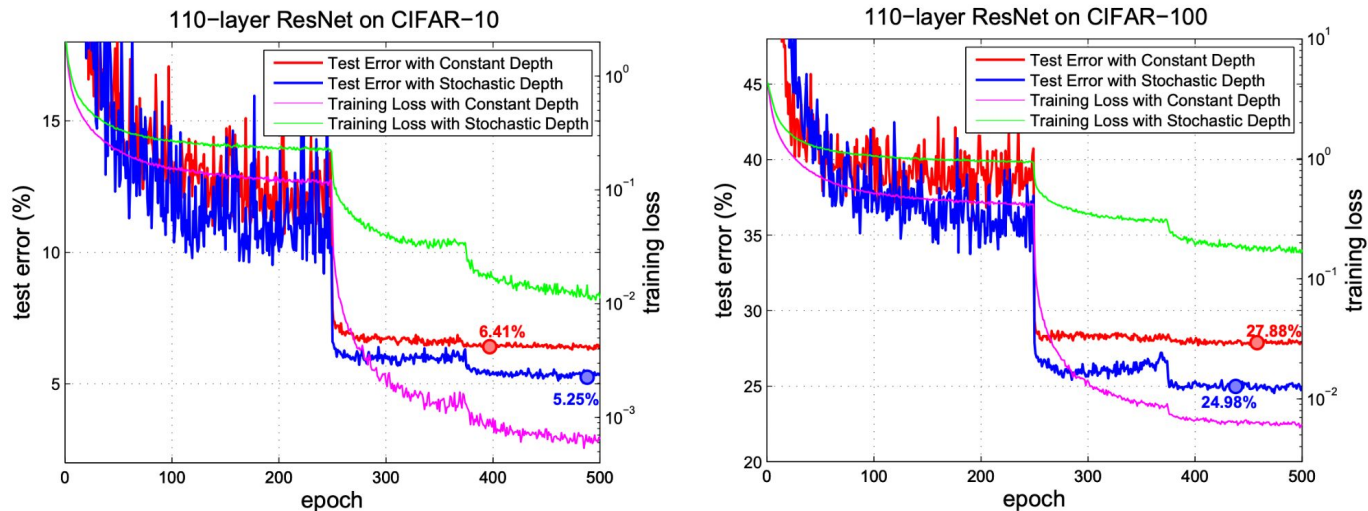


Fig. 3. Test error on CIFAR-10 (*left*) and CIFAR-100 (*right*) during training, with data augmentation, corresponding to results in the first two columns of Table 1.

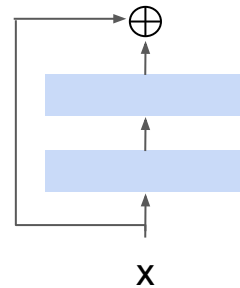
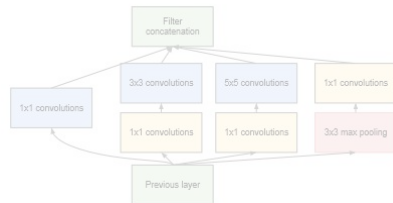
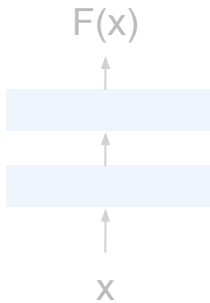
CNN Architectures

“Plain” CNN	GoogLeNet	ResNet
-------------	-----------	--------

Simple connection from previous to next layer

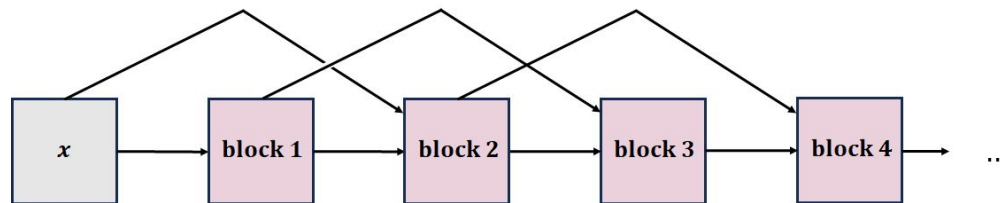
1x1, 3x3, 5x5 convolutions and pooling between each layer

Skip connections
Add output of previous layer to next layer

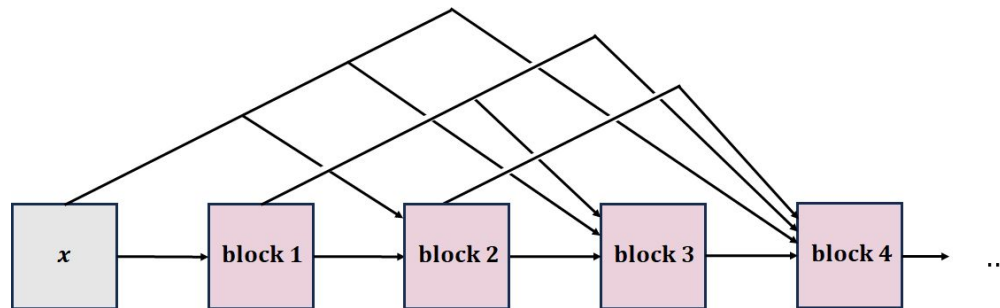


From ResNets to DenseNets

ResNet



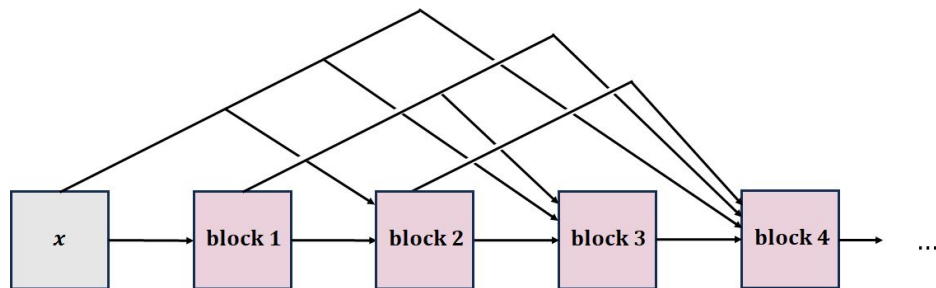
DenseNet



Dense Connections

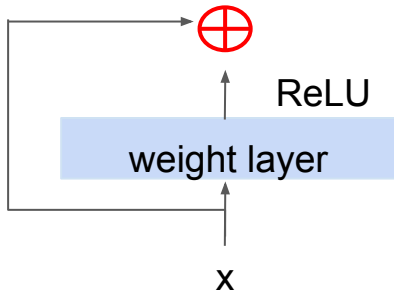
Each layer has access to every other layer before it, which:

- maximizes information flow
- allows for feature-map reuse
- less parameters to learn
- alleviates vanishing gradient

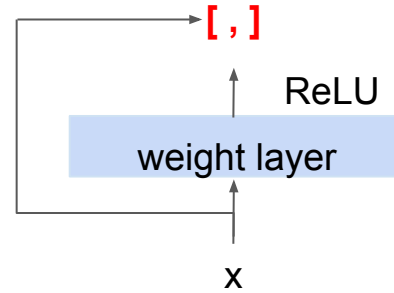


Dense Blocks

To create dense connections, dense blocks use the same structure as residual blocks, but concatenate (denoted by $[,]$) inputs instead of simply adding them



Residual Blocks

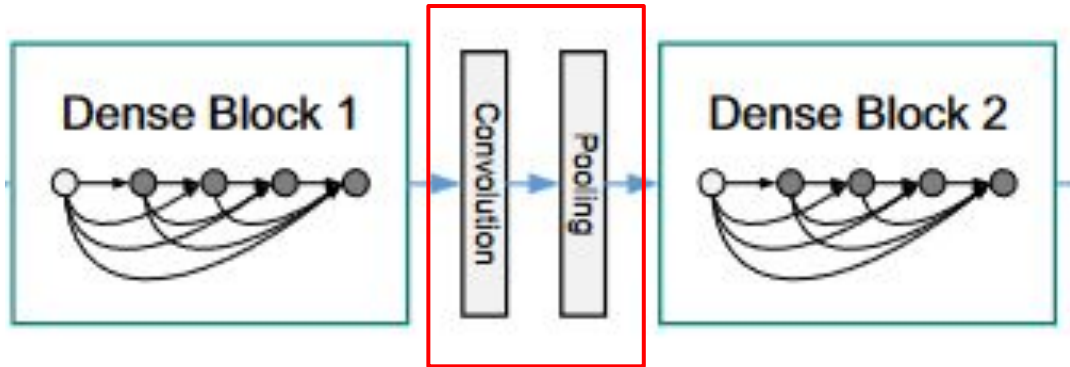


Dense Blocks

Transition Layers

Each dense block increases the number of dimensions

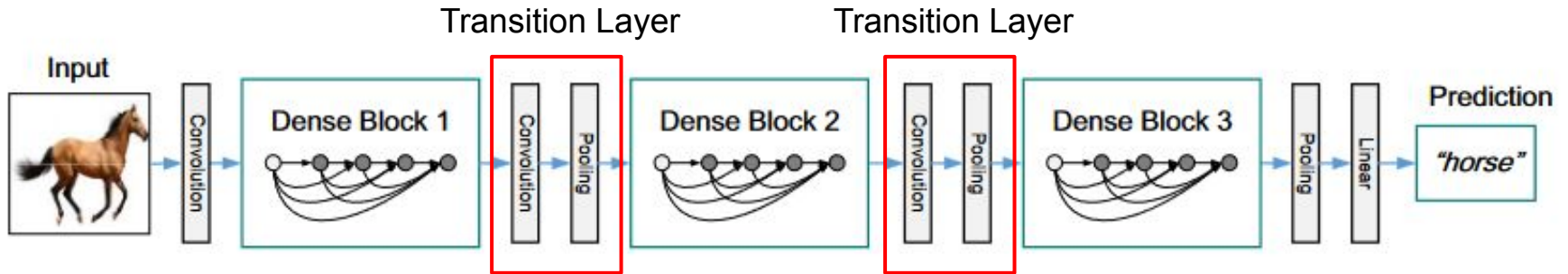
Maintain size of dimension with 1x1 convolutions and pooling (transition layer)



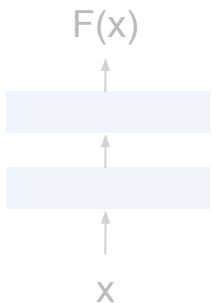
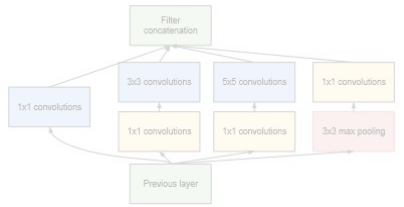
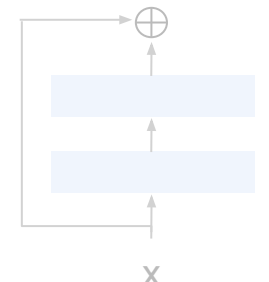
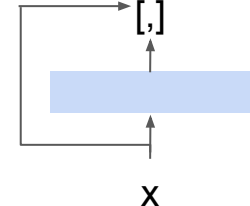
Transition Layer

DenseNet Architecture

Stack Dense Blocks together with transition layers in between each Dense Block



CNN Architectures

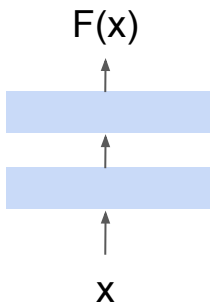
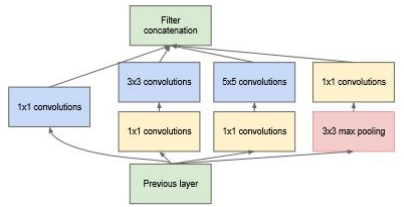
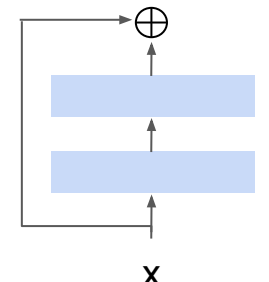
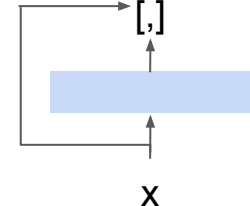
"Plain" CNN	GoogLeNet	ResNet	DenseNet
<p>Simple connection from previous to next layer</p>	<p>1x1, 3x3, 5x5 convolutions and pooling between each layer</p>	<p>Skip connections Add output of previous layer to next layer</p>	<p>Dense connections Concatenate output of previous layer to next layer</p>
 <p>The diagram shows a vertical flow starting with an input X at the bottom. An arrow points up to a light blue rectangular layer. A second arrow points up to another light blue rectangular layer. A final arrow points up to the output $F(x)$.</p>	 <p>The diagram illustrates the GoogLeNet architecture. It starts with a 'Previous layer' at the bottom. The path branches into several parallel operations: '1x1 convolutions', '3x3 convolutions', '5x5 convolutions', and '3x3 max pooling'. These paths then converge into '1x1 convolutions' and finally 'Filter concatenation' at the top.</p>	 <p>The diagram shows a ResNet block. An input X enters from the bottom. It splits into two paths: one goes directly to a circular addition node (\oplus), and the other goes through a light blue rectangular layer. The outputs of both paths are combined at the addition node.</p>	 <p>The diagram shows a DenseNet block. An input X enters from the bottom. It splits into two paths: one goes directly to a square concatenation node ($[,]$), and the other goes through a light blue rectangular layer. The outputs of both paths are combined at the concatenation node.</p>

Model Comparison - error rates

Method	Depth	Params	C10	C10+	C100	C100+	SVHN
Network in Network [22]	-	-	10.41	8.81	35.68	-	2.35
All-CNN [31]	-	-	9.08	7.25	-	33.71	-
Deeply Supervised Net [20]	-	-	9.69	7.97	-	34.57	1.92
Highway Network [33]	-	-	-	7.72	-	32.39	-
FractalNet [17]	21	38.6M	10.18	5.22	35.34	23.30	2.01
with Dropout/Drop-path	21	38.6M	7.33	4.60	28.20	23.73	1.87
ResNet [11]	110	1.7M	-	6.61	-	-	-
ResNet (reported by [13])	110	1.7M	13.63	6.41	44.74	27.22	2.01
ResNet with Stochastic Depth [13]	110	1.7M	11.66	5.23	37.80	24.58	1.75
	1202	10.2M	-	4.91	-	-	-
Wide ResNet [41]	16	11.0M	-	4.81	-	22.07	-
	28	36.5M	-	4.17	-	20.50	-
with Dropout	16	2.7M	-	-	-	-	1.64
ResNet (pre-activation) [12]	164	1.7M	11.26*	5.46	35.58*	24.33	-
	1001	10.2M	10.56*	4.62	33.47*	22.71	-
DenseNet ($k = 12$)	40	1.0M	7.00	5.24	27.55	24.42	1.79
DenseNet ($k = 12$)	100	7.0M	5.77	4.10	23.79	20.20	1.67
DenseNet ($k = 24$)	100	27.2M	5.83	3.74	23.42	19.25	1.59
DenseNet-BC ($k = 12$)	100	0.8M	5.92	4.51	24.15	22.27	1.76
DenseNet-BC ($k = 24$)	250	15.3M	5.19	3.62	19.64	17.60	1.74
DenseNet-BC ($k = 40$)	190	25.6M	-	3.46	-	17.18	-

[Huang, Gao, et al. "Densely connected convolutional networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2017.]

Summary of Models

“Plain” CNN	Google Net	ResNet	DenseNet
<p>Simple connection from previous to next layer</p>	<p>1x1, 3x3, 5x5 convolutions and pooling between each layer</p>	<p>Skip connections Add output of previous layer to next layer</p>	<p>Dense connections Concatenate output of previous layer to next layer</p>
 <p>The diagram shows a vertical flow starting with an input 'X' at the bottom. An upward arrow points to a blue rectangular layer. Another upward arrow points to a second blue rectangular layer. A final upward arrow points to the output 'F(x)'.</p>	 <p>The diagram shows a 'Previous layer' (green box) at the bottom. An upward arrow splits into four paths leading to: a blue box labeled '1x1 convolutions', a blue box labeled '3x3 convolutions', a blue box labeled '5x5 convolutions', and a yellow box labeled '1x1 convolutions'. Below these are two yellow boxes labeled '1x1 convolutions' and a red box labeled '3x3 max pooling'. Arrows from all these boxes point to a green box at the top labeled 'Filter concatenation'.</p>	 <p>The diagram shows an input 'X' at the bottom. An upward arrow splits into two paths. The main path goes through a blue rectangular layer. The second path goes around the layer and ends at a circle with a plus sign (⊕). An arrow from the top of the blue layer also points to this circle.</p>	 <p>The diagram shows an input 'X' at the bottom. An upward arrow splits into two paths. The main path goes through a blue rectangular layer. The second path goes around the layer and ends at a circle containing '[' and ',' ([,]). An arrow from the top of the blue layer also points to this circle.</p>

Summary

- Deep CNNs outperform shallow CNNs
- But...
 - Harder optimization problem!
- Residual (and dense) connections make training easier!
 - Can train networks with 100s of layers!
- Stochastic depth let's you train deeper networks faster
 - 1000+ layers!
- In general...
 - Build large networks as stacks of (many!) simple building blocks