# Thanks to:

Varsha Kishore
Justin Lovelace
Anissa Dallmann
Stephanie Ginting

# Clarification: Dropout

In each forward pass, randomly set some neurons to zero.

The probability of ~~keeping~~ a neuron is a hyperparameter; p=0.5 is common.
zeroing
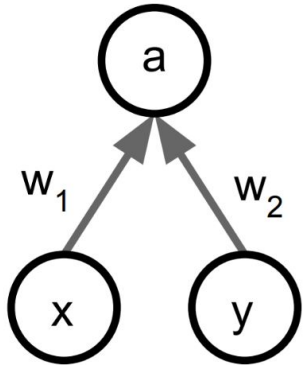


[Srivastava et al. 2014]

# Clarification: Dropout During Test Time

Need to re-scale activations so they are the same (in expectation) during training and testing

Consider a single neuron.

At test time we have: $E[a] = w_1 x + w_2 y$

During training we have:

$$E[a] = \frac{1}{4}(w_1 x + w_2 y) + \frac{1}{4}(w_1 x + 0y)$$
$$+ \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2 y)$$
$$= \frac{1}{2}(w_1 x + w_2 y)$$

At test time, **multiply** by **(1 - p)**

# Review: Image Classification



flatten

0.9 "dog"

0.1 "cat"
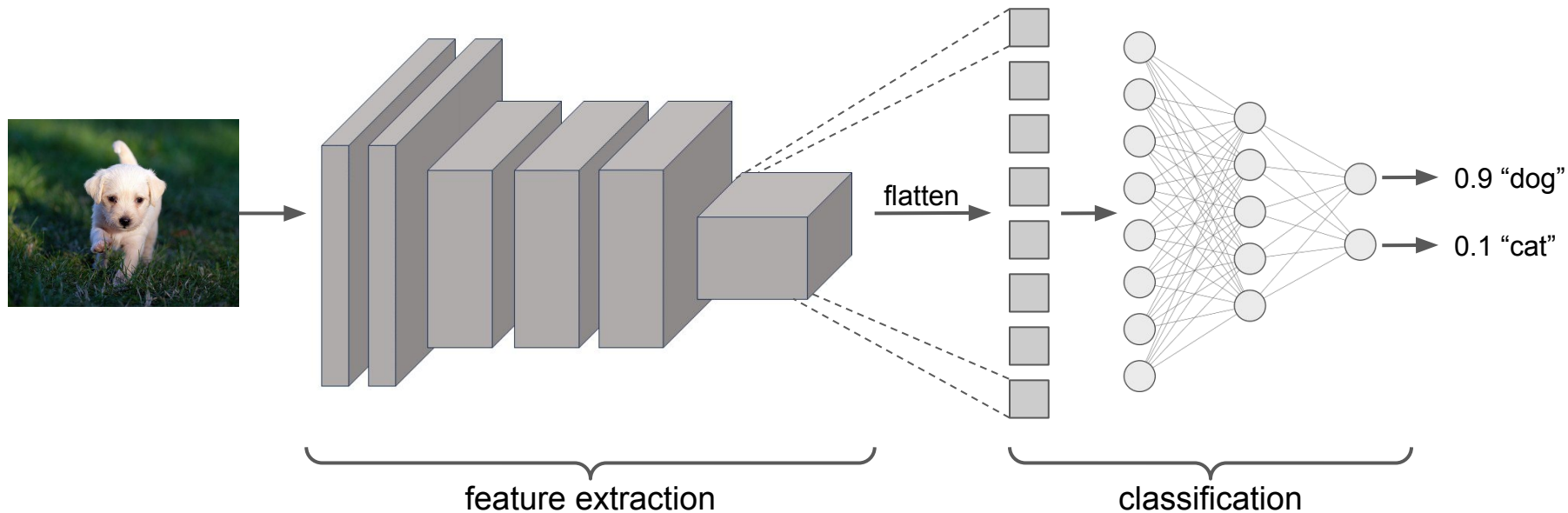
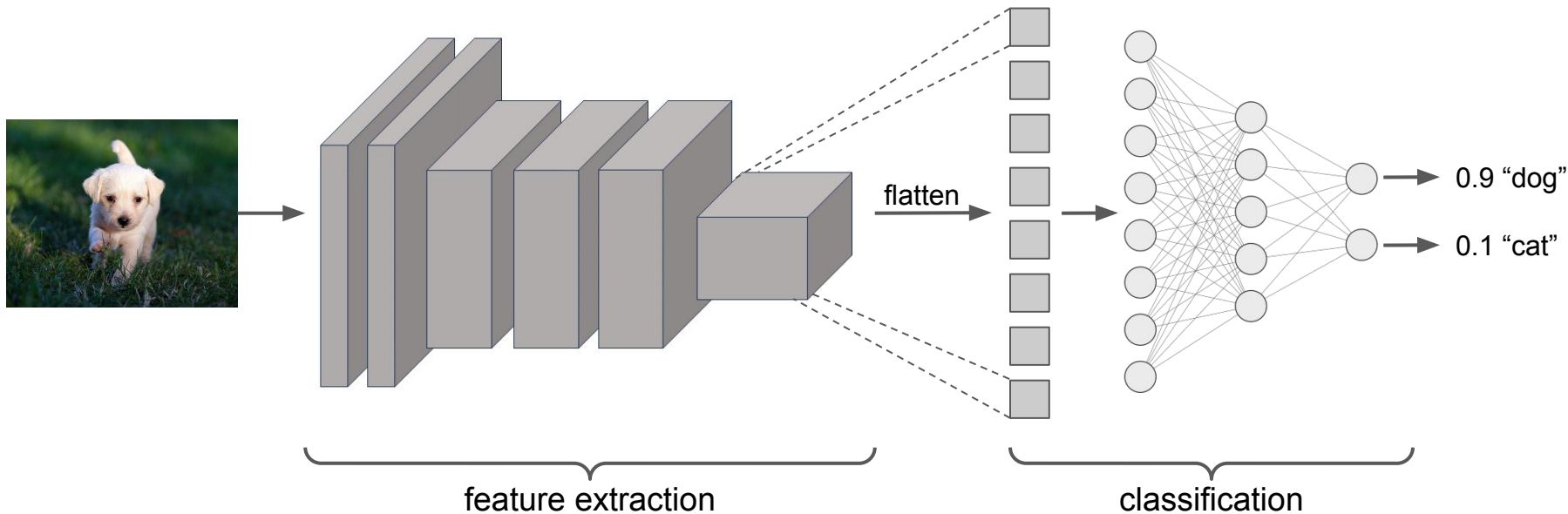feature extraction

classification

# 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!



feature extraction        classification

# Deeper CNN Architectures



feature extraction

classification

0.9 "dog"

0.2 "cat"

flatten

# Deeper CNN Architectures



5 x 5 convolution

**vs**

3 x 3 convolution

3 x 3 convolution

**Performed better!**

# Deeper == better



CNN

# Deeper == better



CNN

# Deeper == better?



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

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

# Discuss: How can a larger network achieve a higher training error?



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 <<< 1, gradient "vanishes" as the entire multiplication goes towards 0
    - => Weights not updated properly

# ImageNet Classification Challenge: Deeper == better



[Nguyen, Kien & Fookes, Clinton & Ross, Arun & Sridharan, Sridha. (2017). Iris Recognition with Off-the-Shelf CNN Features: A Deep Learning Perspective. IEEE Access. PP. 1-1. 10.1109/ACCESS.2017.2784352. ]

# GoogLeNet/Inception Net



Goal: given a fixed computational budget, design the best network

=> Deeper networks with computational...

In this paper, we will focus on an efficient deep neural network architecture for computer vision, codenamed Inception, which derives its name from the Network in network paper by Lin et al [12] in conjunction with the famous "we need to go deeper" internet meme [1]. In our case, the word

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

# Inception Module



Inception module = main
building blocks

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

# Inception Module

Still expensive!



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

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

# Remember: 1x1 convolutions

**input**

**filters**

**output**



*

*1x1x64*

* 32 filters

=

Channel

Channel

*56x56x64*

*56x56x32*

# Discuss: Impact of Dimension Reduction

Assume you have an input feature map with 256 channels/features.

Compare the parameter counts from:

1. 3x3 conv with 256 filters

2. 1x1 conv with 64 filters → 3x3 conv with 64 filters → 1x1 conv with 256 filters

# Discuss: Impact of Dimension Reduction

Assume you have an input feature map with 256 channels/features.

Compare the parameter counts from:

1.  3x3 conv with 256 filters
    3*3*256*256 = ~590k parameters

2.  1x1 conv with 64 filters → 3x3 conv with 64 filters → 1x1 conv with 256 filters

# Discuss: Impact of Dimension Reduction

Assume you have an input feature map with 256 channels/features.

Compare the parameter counts from:
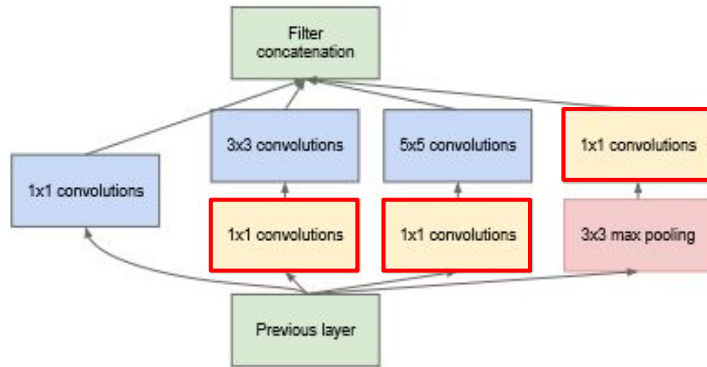
1. 3x3 conv with 256 filters
   3*3*256*256 = ~590k parameters

2. 1x1 conv with 64 filters → 3x3 conv with 64 filters → 1x1 conv with 256 filters
   1*1*256*64 + 3*3*64*64 + 1*1*64*256 = ~70k parameters

# Inception Module

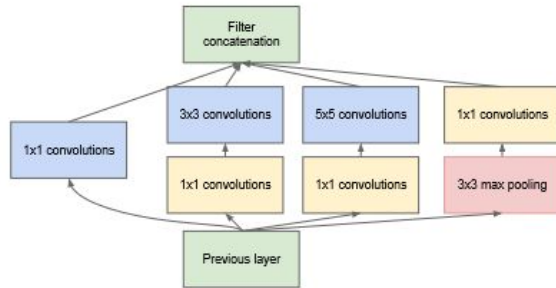Solution: Inception module with dimension reduction

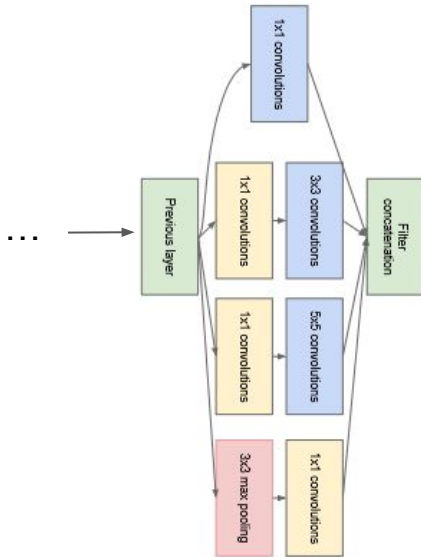

- "Bottleneck" with 1x1 convolutions to reduce dimensions

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

# 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.]

# 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.]
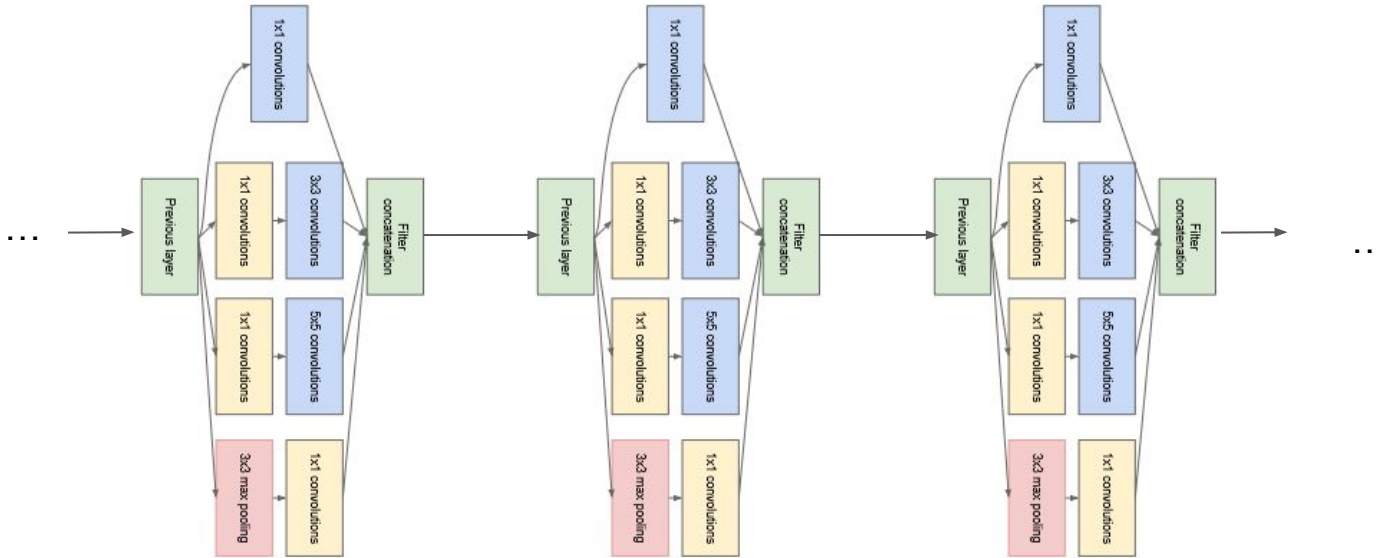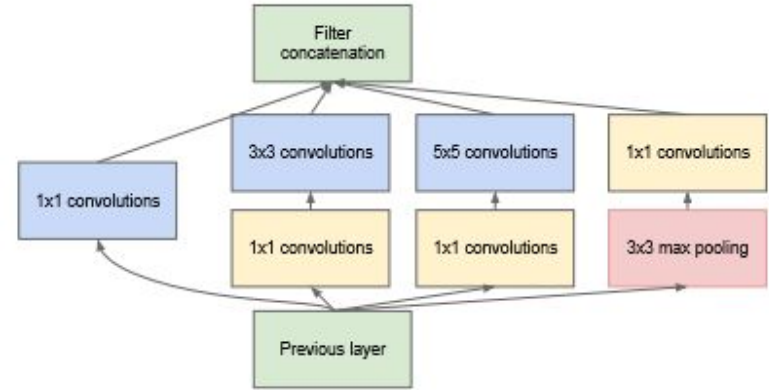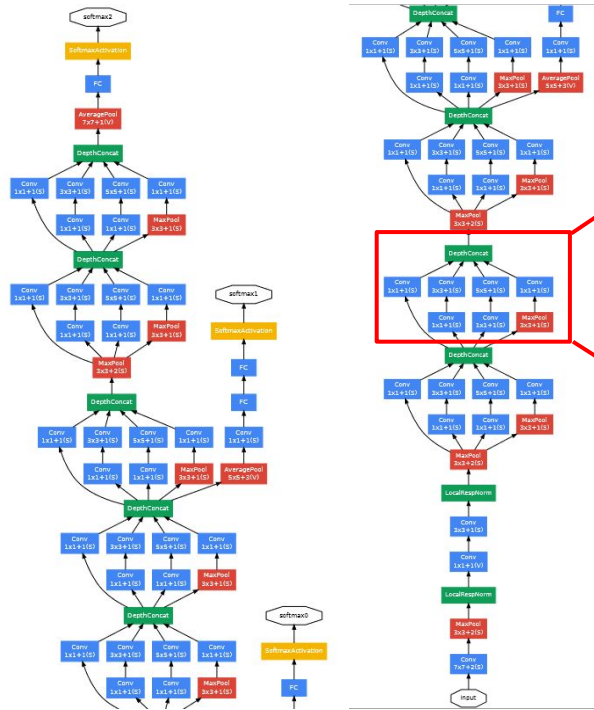
# GoogLeNet Architecture

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.]
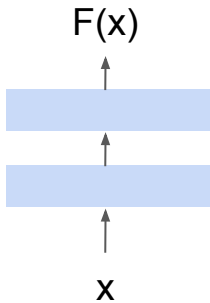
# The Entire GoogLeNet Architecture



Inception Module

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

# CNN Architectures

**"Plain" CNN**
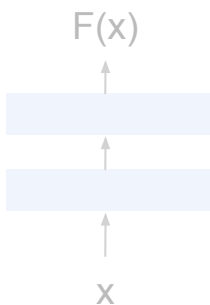
Simple connection from previous to next layer

F(x)

x

# CNN Architectures

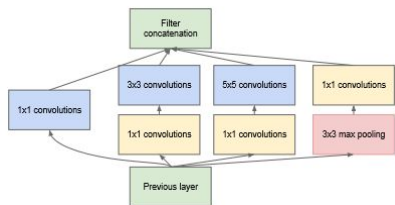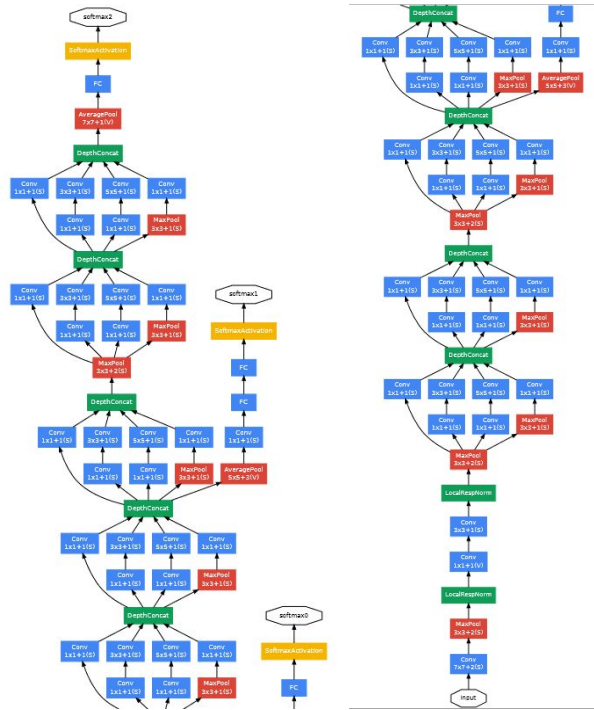| "Plain" CNN | GoogLeNet |
|:---:|:---:|
| Simple connection from previous to next layer | 1x1, 3x3, 5x5 convolutions and pooling between each layer |

# The Entire GoogleNet Architecture



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

[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?
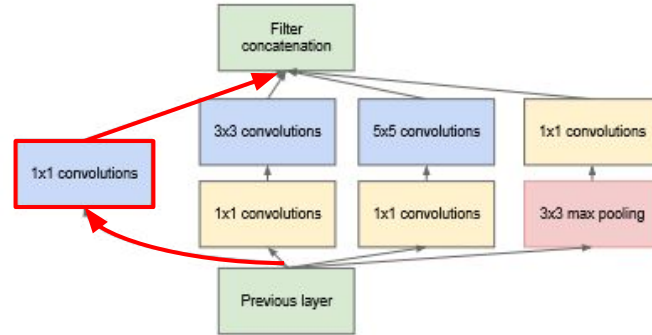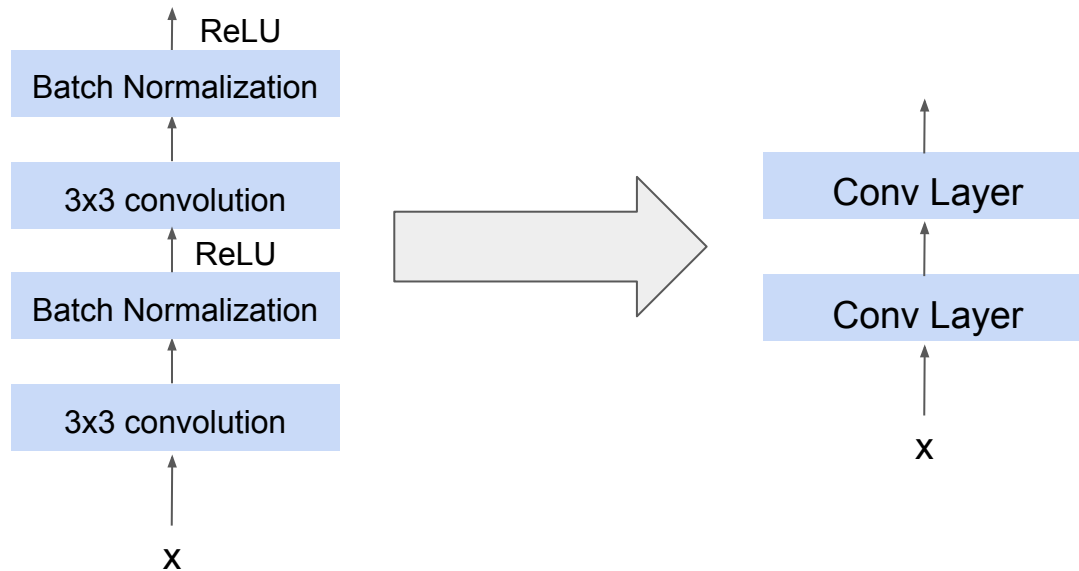


**Residual connections**

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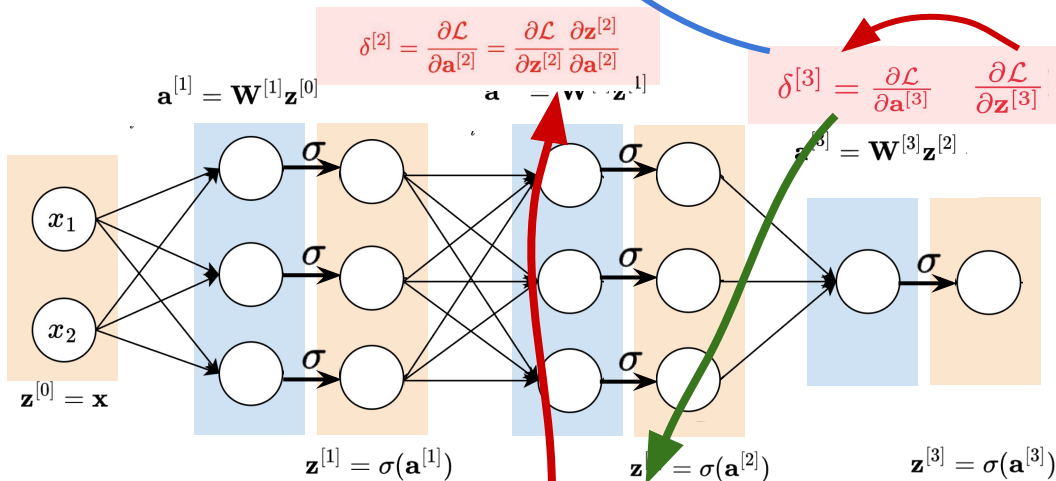
# Aside: Conv Layer Abstraction

# Backpropagation

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[3]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{W}^{[3]}}$$

$$= \delta^{[3]} (\mathbf{z}^{[2]})^T$$
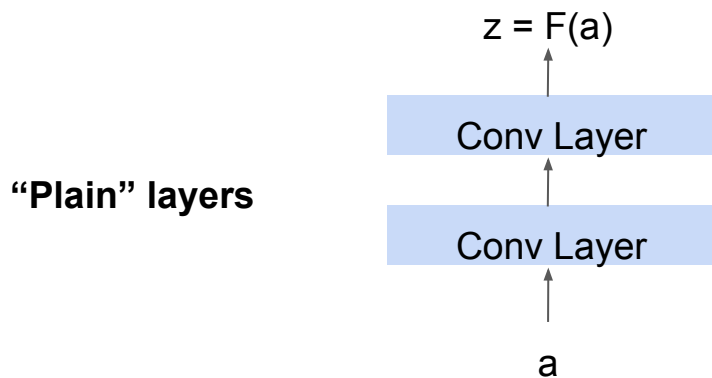
**Algorithm** Backward Pass through MLP (Detailed)

1: **Input:** $\{\mathbf{z}^{[1]}, \ldots, \mathbf{z}^{[L]}\}$, $\{\mathbf{a}^{[1]}, \ldots, \mathbf{a}^{[L]}\}$, loss gradient $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}}$

2: $\delta^{[L]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \frac{\partial \mathbf{z}^{[L]}}{\partial \mathbf{a}^{[L]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[L]}} \odot \sigma^{[L]\prime}(\mathbf{a}^{[L]})$     ▷ Error term

3: **for** $l = L$ **to** 1 **do**

4:     $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{W}^{[l]}} = \delta^{[l]} (\mathbf{z}^{[l-1]})^T$     ▷ Gradient of weights

5:     $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[l]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{b}^{[l]}} = \delta^{[l]}$     ▷ Gradient of biases

6:     $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l]}} \frac{\partial \mathbf{a}^{[l]}}{\partial \mathbf{z}^{[l-1]}} = (\mathbf{W}^{[l]})^T \delta^{[l]}$

7:     $\delta^{[l-1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[l-1]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[l-1]}} \frac{\partial \mathbf{z}^{[l-1]}}{\partial \mathbf{a}^{[l-1]}} = ((\mathbf{W}^{[l]})^T \delta^{[l]}) \odot \sigma^{[l-1]\prime}(\mathbf{a}^{[l-1]})$

8: **end for**

9: **Output:** $\frac{\partial \mathcal{L}}{\partial \mathbf{W}^{[1:L]}}$, $\frac{\partial \mathcal{L}}{\partial \mathbf{b}^{[1:L]}}$



$\delta^{[2]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} \frac{\partial \mathbf{z}^{[2]}}{\partial \mathbf{a}^{[2]}}$

$\delta^{[3]} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \quad \frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$

$\mathbf{a}^{[1]} = \mathbf{W}^{[1]}\mathbf{z}^{[0]}$

$\mathbf{a}^{[3]} = \mathbf{W}^{[3]}\mathbf{z}^{[2]}$

$\mathbf{z}^{[0]} = \mathbf{x}$

$\mathbf{z}^{[1]} = \sigma(\mathbf{a}^{[1]})$

$\mathbf{z}^{[2]} = \sigma(\mathbf{a}^{[2]})$

$\mathbf{z}^{[3]} = \sigma(\mathbf{a}^{[3]})$

$\mathcal{L}(\mathbf{z}^{[3]}, \mathbf{y})$

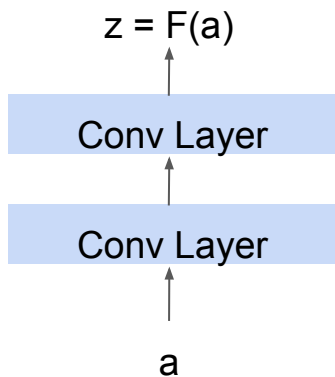We can directly compute $\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[3]}}$!

$\frac{\partial \mathcal{L}}{\partial \mathbf{z}^{[2]}} = \frac{\partial \mathcal{L}}{\partial \mathbf{a}^{[3]}} \frac{\partial \mathbf{a}^{[3]}}{\partial \mathbf{z}^{[2]}} = (W^{[3]})^T \delta^{[3]}$

# Backpropagation through "plain" conv layers

z = F(a)

Conv Layer

**"Plain" layers**

Conv Layer

a

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial a} = \frac{\partial L}{\partial z}\boxed{(F'(a))}$$
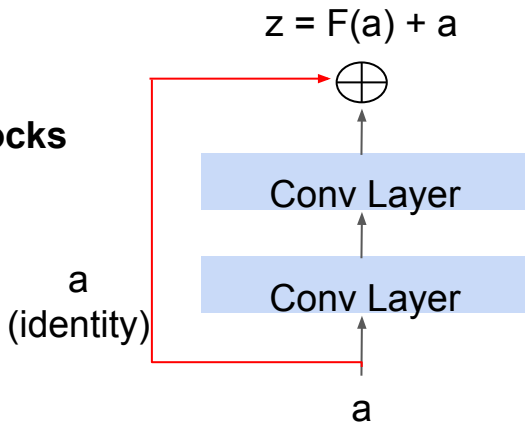
# Discussion: Backpropagation through Residual blocks

z = F(a)

Conv Layer

**"Plain" layers**

Conv Layer

a

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial a} = \frac{\partial L}{\partial z}\boxed{(F'(a))}$$
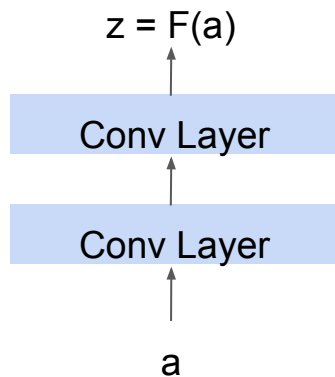
z = F(a) + a

$\oplus$

**Residual Blocks**

Conv Layer

a
(identity)

Conv Layer

a

$$\frac{\partial L}{\partial a} =$$

# Backpropagation through Residual blocks
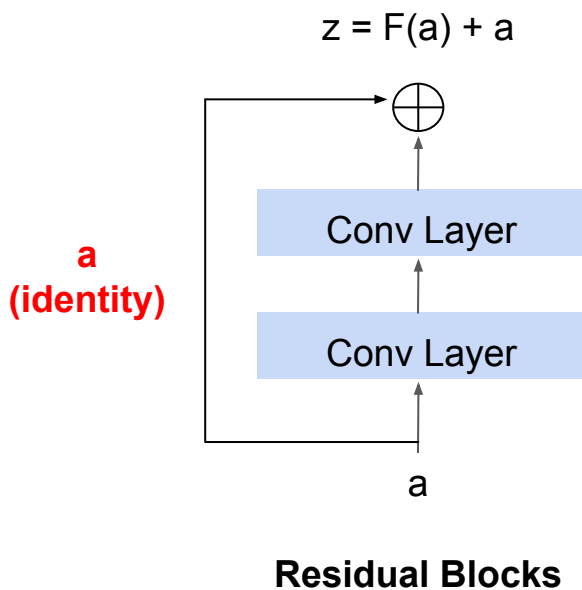
z = F(a)

**"Plain" layers**

Conv Layer

Conv Layer

a

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial a} = \frac{\partial L}{\partial z}\boxed{(F'(a))}$$

z = F(a) + a

**Residual Blocks**

$\oplus$

Conv Layer

Conv Layer

a
(identity)

a

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial a} = \frac{\partial L}{\partial z}\boxed{(1 + F'(a))}$$

# Residual Connections

z = F(a) + a

**a**
**(identity)**

Conv Layer

Conv Layer

a

**Residual Blocks**

Identity mapping

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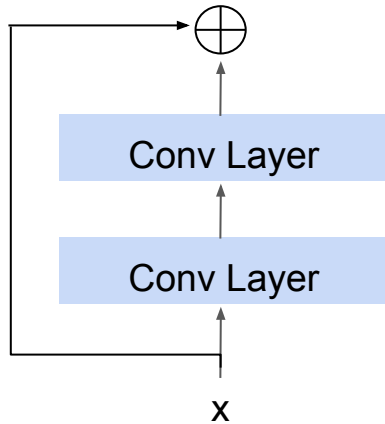- only learn *difference* of feature maps

Additive component of identity

- alleviates vanishing gradients

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial a} = \frac{\partial L}{\partial z}(1 + F'(a))$$
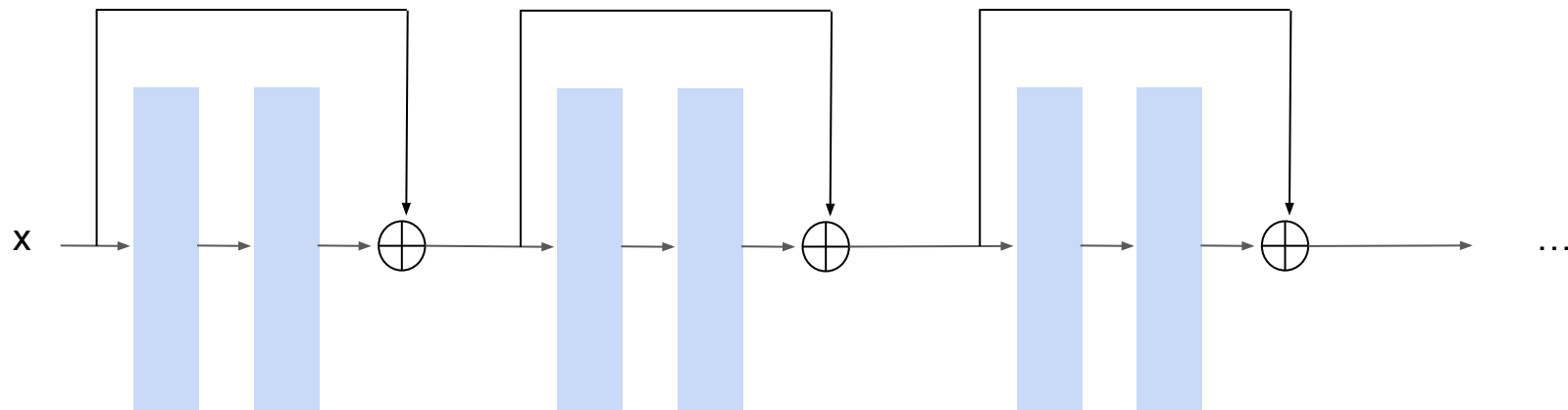
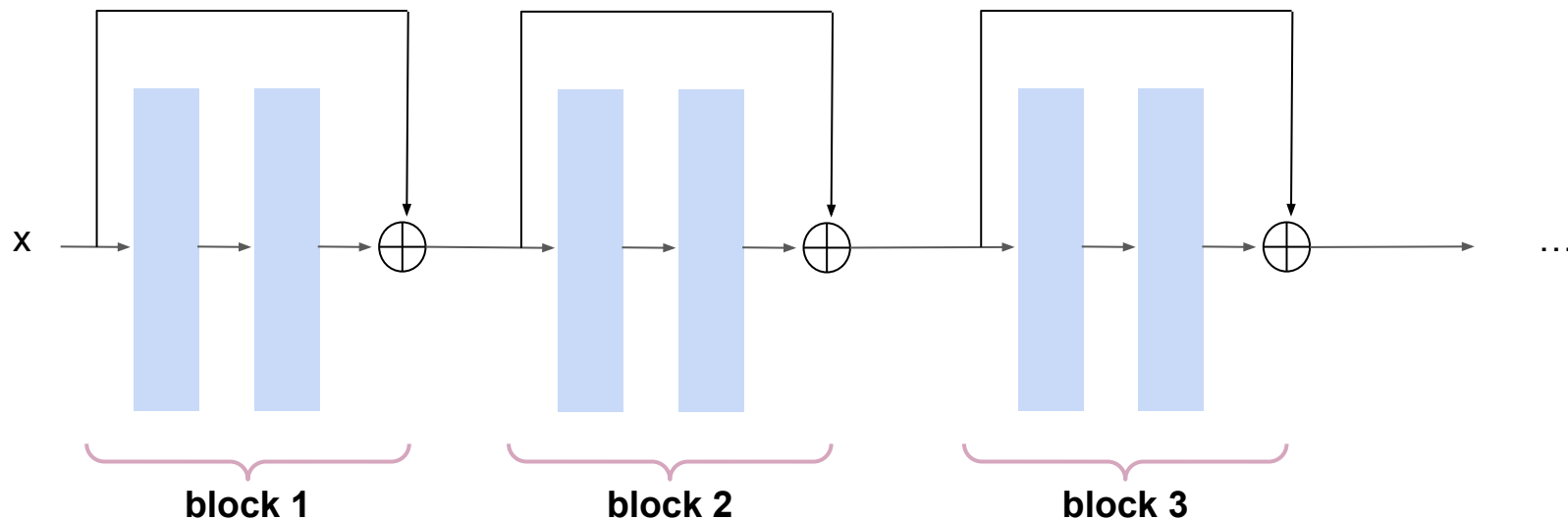[He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.]

# ResNet

Stack residual blocks together!



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

# ResNet

Stack residual blocks together!
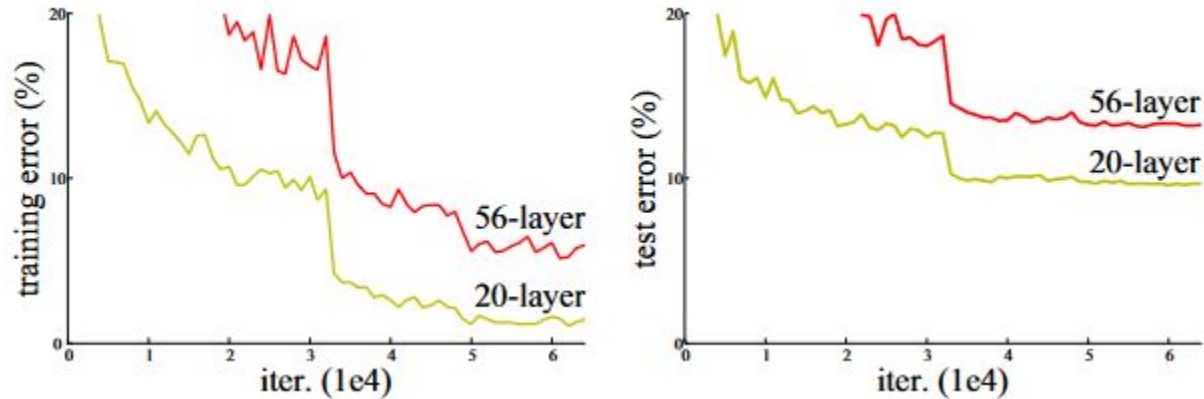


[He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.]

# ResNet

Stack residual blocks together!



block 1                    block 2                    block 3

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

# ResNet

Stack residual blocks together!



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

# Full ResNet Architecture



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

# Recall: How can a larger network achieve a higher training error?



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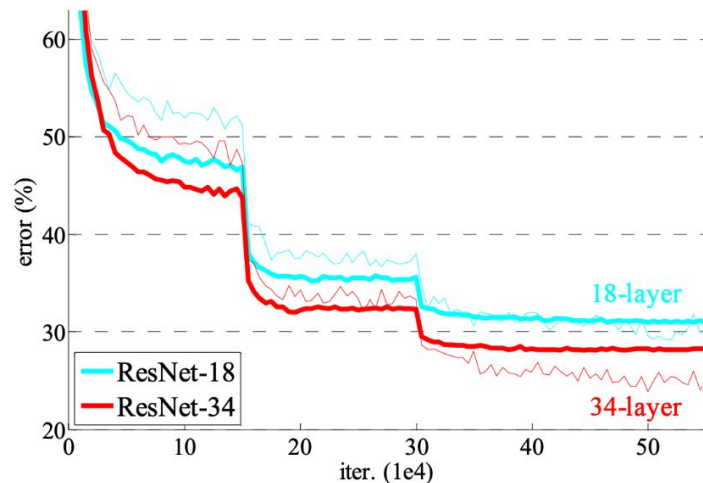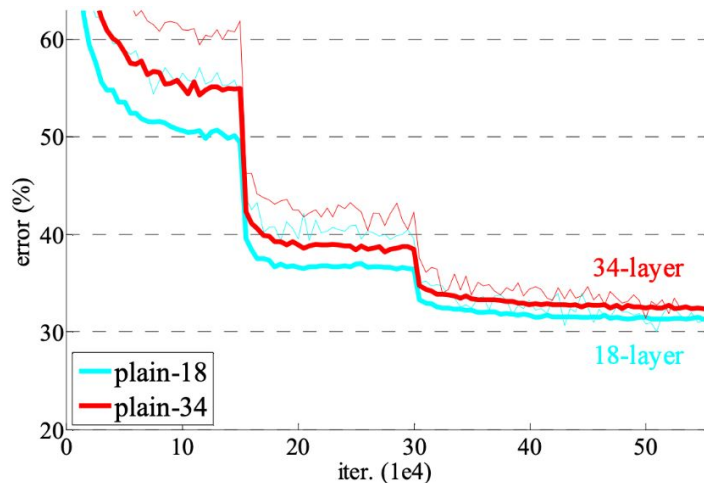
# Deeper == better

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.

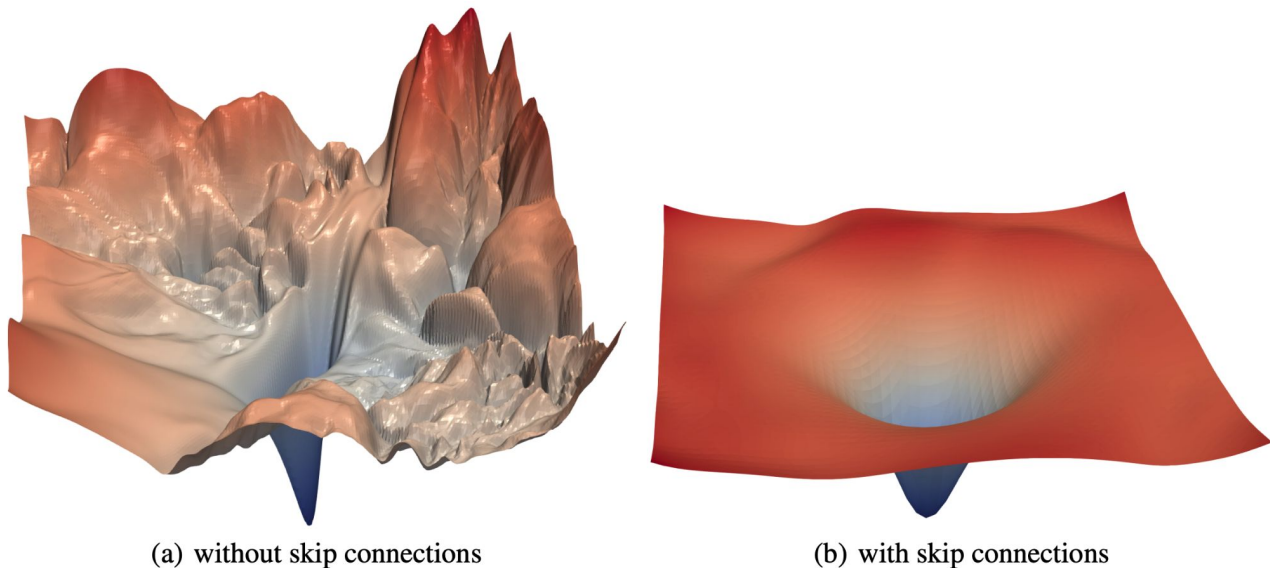# Visualizing the Effect of Skip Connections

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.
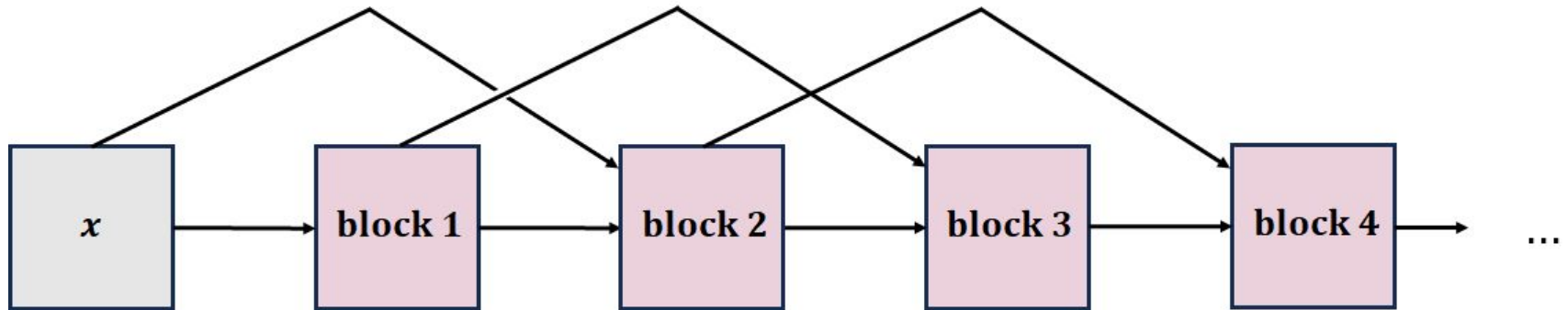
[Li, Hao, et al. "Visualizing the loss landscape of neural nets." Advances in neural information processing systems 31 (2018).]

# Stochastic Depth

Still have long training times! Solution: stochastic depth

[Huang, Gao, et al. "Deep networks with stochastic depth." *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14.* Springer International Publishing, 2016.]

# Stochastic Depth

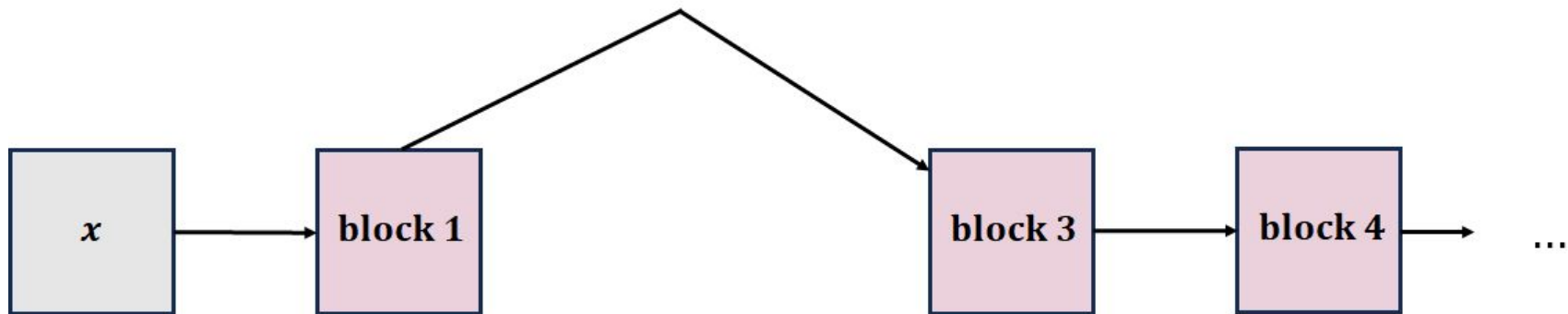During training, randomly drop Residual Blocks using skip connections

Like dropout but with residual blocks instead of individual neurons



[Huang, Gao, et al. "Deep networks with stochastic depth." *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14.* Springer International Publishing, 2016.]

# Stochastic Depth

During training, randomly drop Residual Blocks using skip connections

Like dropout but with residual blocks instead of individual neurons



[Huang, Gao, et al. "Deep networks with stochastic depth." *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14.* Springer International Publishing, 2016.]

# Stochastic Depth

Another benefit: robustness/mitigating overfitting



[Huang, Gao, et al. "Deep networks with stochastic depth." *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14.* Springer International Publishing, 2016.]

# Stochastic Depth

Increases training loss, but… decreases te



**Fig. 3.** Test error on CIFAR-1 data augmentation, correspond

**Fig. 5.** With stochastic depth, the 1202-layer ResNet still significantly improves over the 110-layer one.

[Huang, Gao, et al. "Deep networks with stochastic depth." *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14.* Springer International Publishing, 2016.]

# CNN Architectures

| "Plain" CNN | GoogLeNet |
|---|---|
| Simple connection from previous to next layer | 1x1, 3x3, 5x5 convolutions and pooling between each layer |

# CNN Architectures

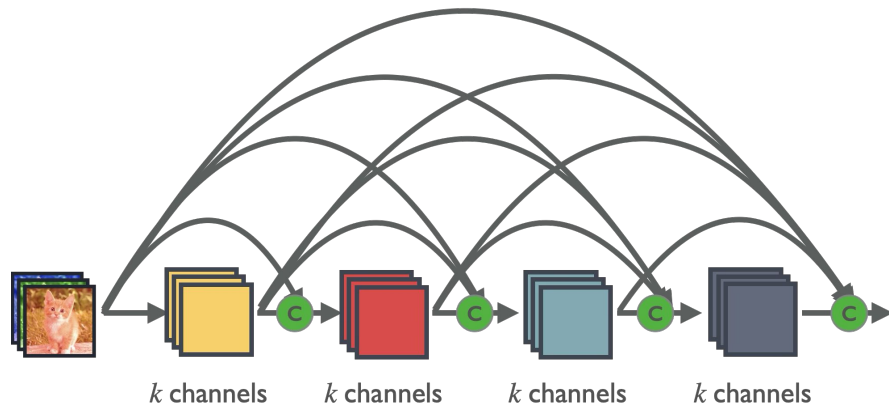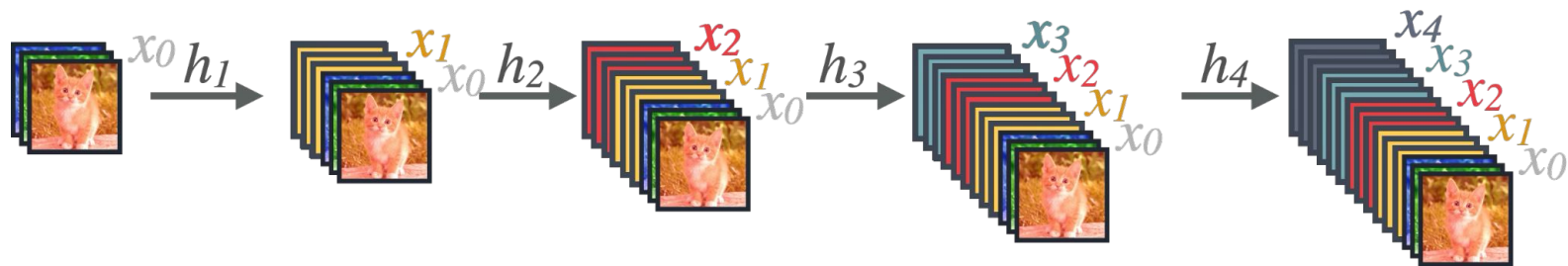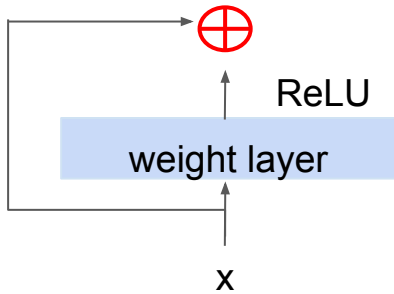| "Plain" CNN | GoogLeNet | ResNet |
|:---:|:---:|:---:|
| Simple connection from previous to next layer | 1x1, 3x3, 5x5 convolutions and pooling between each layer | Skip connections<br><br>Add output of previous layer to next layer |

# From ResNets to DenseNets

**ResNet**

⊕ : Element-wise addition

**DenseNet**

$k$ channels    $k$ channels    $k$ channels    $k$ channels

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

# DenseNets

Feature concatenation

# Dense Blocks

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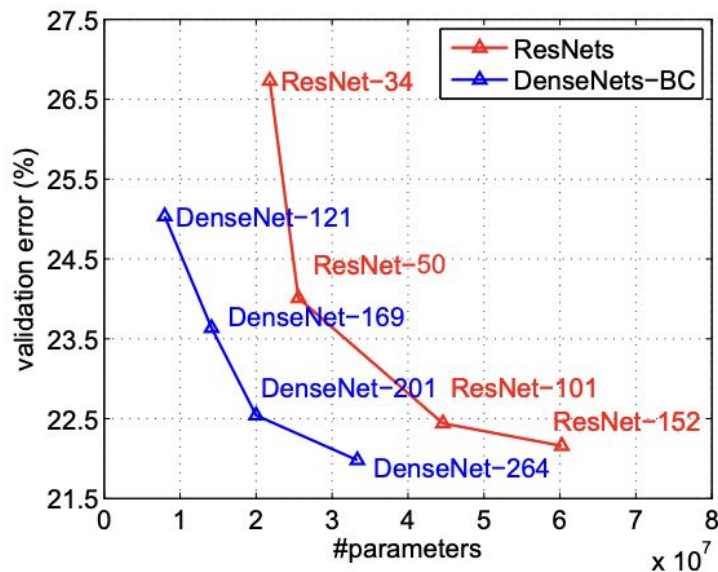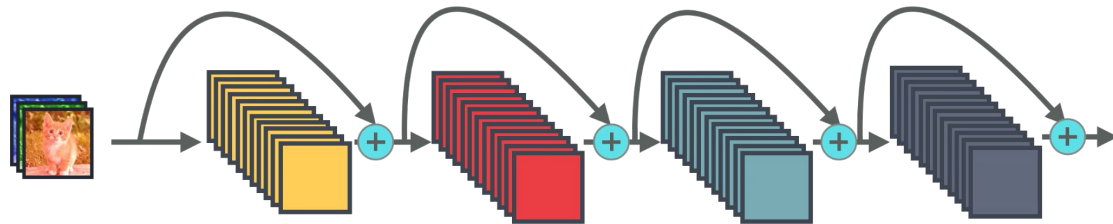


**Residual Blocks**

**Dense Blocks**

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

**Figure 3:** Comparison of the DenseNets and ResNets top-1 error rates (single-crop testing) on the ImageNet validation dataset as a function of learned parameters (*left*) and FLOPs during test-time (*right*).
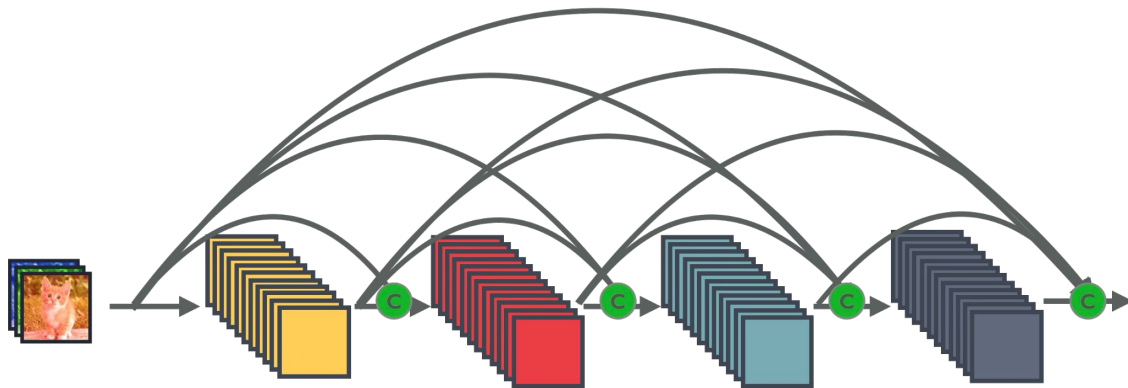
Discussion: What design choices might allow a ~100-layer DenseNet to have fewer parameters than a ~100-layer ResNet?

ResNet
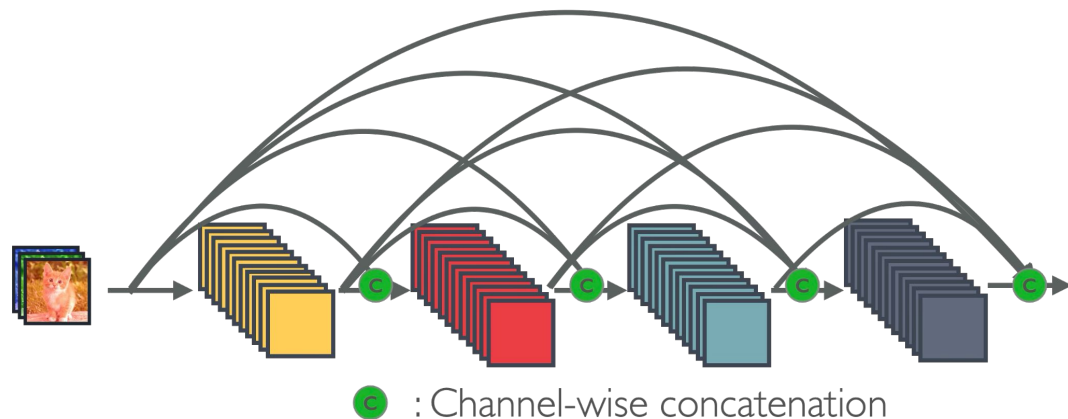
+ : Element-wise addition

DenseNet

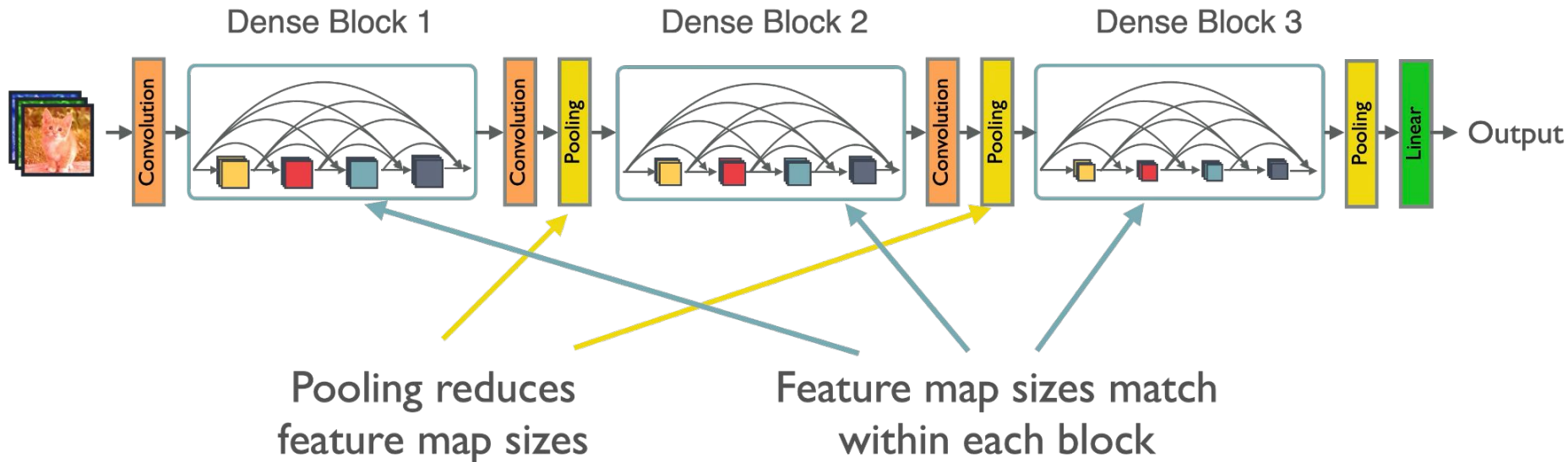c : Channel-wise concatenation

# Dense Connections

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

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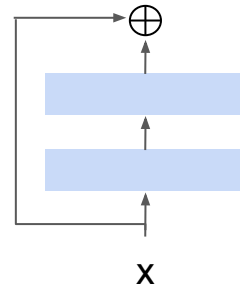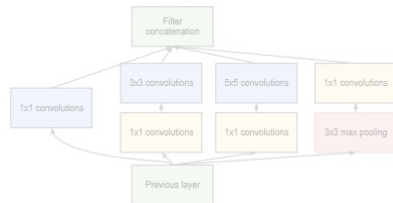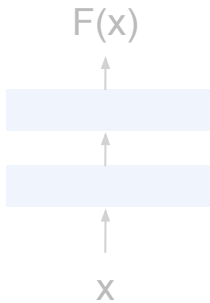- less parameters to learn
- alleviates vanishing gradient



c : Channel-wise concatenation

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

# DenseNets



Dense Block 1

Dense Block 2

Dense Block 3

Convolution

Pooling

Output

Pooling reduces
feature map sizes

Feature map sizes match
within each block

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

# CNN Architectures

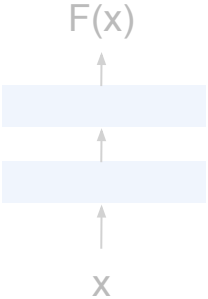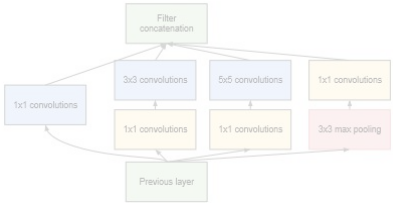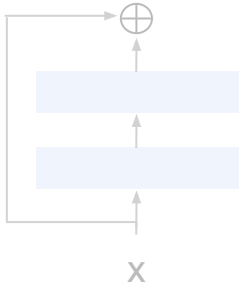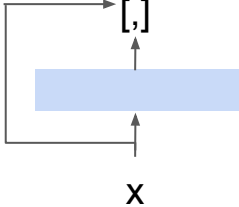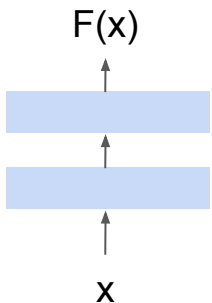| "Plain" CNN | GoogLeNet | ResNet |
|---|---|---|
| Simple connection from previous to next layer | 1x1, 3x3, 5x5 convolutions and pooling between each layer | Skip connections<br><br>Add output of previous layer to next layer |

# CNN Architectures

| "Plain" CNN | GoogLeNet | ResNet | DenseNet |
|---|---|---|---|
| Simple connection from previous to next layer | 1x1, 3x3, 5x5 convolutions and pooling between each layer | Skip connections<br><br>Add output of previous layer to next layer | Dense connections<br><br>Concatenate output of previous layer to next layer |

F(x)

x

Filter concatenation

3x3 convolutions    5x5 convolutions    1x1 convolutions

1x1 convolutions    1x1 convolutions    3x3 max pooling

1x1 convolutions

Previous layer

⊕

x

[,]

x

# Summary of Models

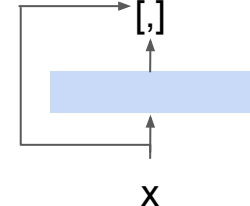| "Plain" CNN | Google Net | ResNet | DenseNet |
|---|---|---|---|
| Simple connection from previous to next layer | 1x1, 3x3, 5x5 convolutions and pooling between each layer | Skip connections<br><br>Add output of previous layer to next layer | Dense connections<br><br>Concatenate output of previous layer to next layer |

# 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