

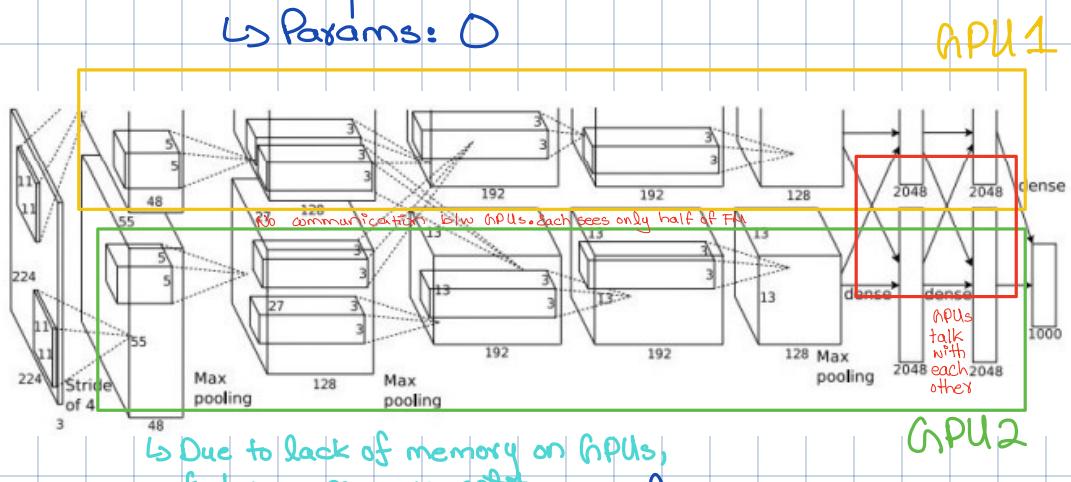
CNN Architecture

1) AlexNet: Error: 16.4%. (2012)

CONN1 → 96 11x11 filters at stride 4
 ↳ Output Volume: 55x55x96
 ↳ Parameters: $(11 * 11 * 3) * 96 = 35k$

NORM1
 CONV2
 MAXPOOL2
 NORM2
 CONV3
 CONV4
 CONV5
 MAXPOOL3
 FC6
 FC7
 FC8

MAXPOOL1 → 3x3 filters with stride 2
 ↳ Output Volume: 27x27x96
 ↳ Params: 0



[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)

Details:

↳ Used ReLU

↳ Used Norm layers

↳ Heavy data augmentation

↳ Dropout 0.5

↳ Batch size 128

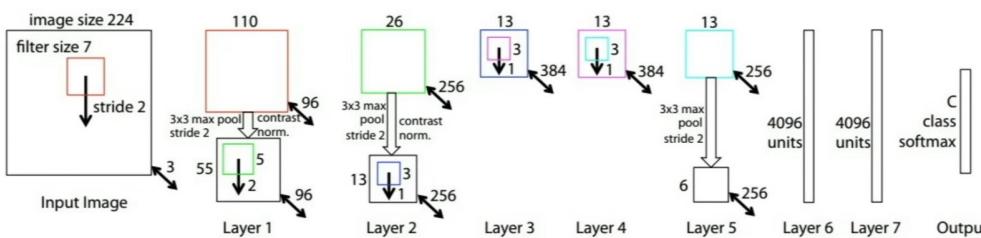
↳ SGD Momentum 0.9

↳ LR: 1e-2, reduced by 10 when val accuracy plateaus

↳ L2 weight decay 5e-4

↳ 7 CNN Ensemble: 18.2% → 15.4%

ZFNet: Error: 11.7 (2013)



Hypertuned params of AlexNet

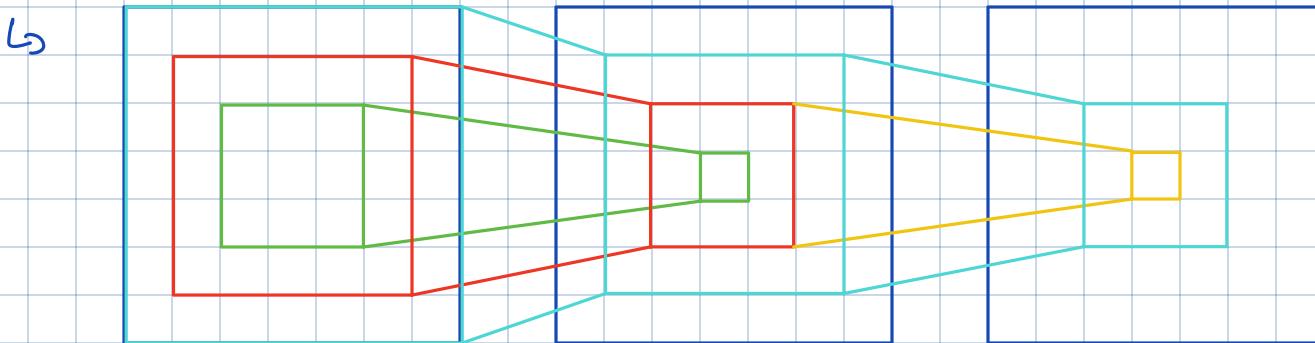
VhNet: Errors: 7.3% (2014)

↳ Filter size: 3×3

Q. Why use smaller filters?

↳ Stack of three 3×3 filters conv (stride 1) layers has the same effective receptive field as one 7×7 conv layer
 ↳ Maximum area that a CNN can capture! process

Q. What is the receptive field of three 3×3 conv (stride 1) layers?



↳ So, effective receptive field of 3 3×3 filter is same as one 7×7 filter

↳ Less params also. $3 \times (3 \times 3 \times C \times C) = 27C^2$ vs $7 \times 7 \times C \times C = 49C^2$

```

INPUT: [224x224x3]      memory: 224*224*3=150K  params: 0      (not counting biases)
CONV3-64: [224x224x64]   memory: 224*224*64=3.2M  params: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64]   memory: 224*224*64=3.2M  params: (3*3*64)*64 = 36,864
POOL2: [112x112x64]     memory: 112*112*64=800K  params: 0
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  params: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  params: (3*3*128)*128 = 147,456
POOL2: [56x56x128]       memory: 56*56*128=400K  params: 0
CONV3-256: [56x56x256]   memory: 56*56*256=800K  params: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256]   memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256]   memory: 56*56*256=800K  params: (3*3*256)*256 = 589,824
POOL2: [28x28x256]       memory: 28*28*256=200K  params: 0
CONV3-512: [28x28x512]   memory: 28*28*512=400K  params: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512]   memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512]   memory: 28*28*512=400K  params: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]       memory: 14*14*512=100K  params: 0
CONV3-512: [14x14x512]   memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]   memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512]   memory: 14*14*512=100K  params: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]          memory: 7*7*512=25K  params: 0
FC: [1x1x4096]            memory: 4096 params: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]            memory: 4096 params: 4096*4096 = 16,777,216
FC: [1x1x1000]            memory: 1000 params: 4096*1000 = 4,096,000

```

TOTAL memory: 24M * 4 bytes \sim 96MB / image (only forward! \sim 2 for bwd)

TOTAL params: 138M parameters



VGG16

VGG19

Details:

↳ 2nd Place

↳ No Local Response Norm

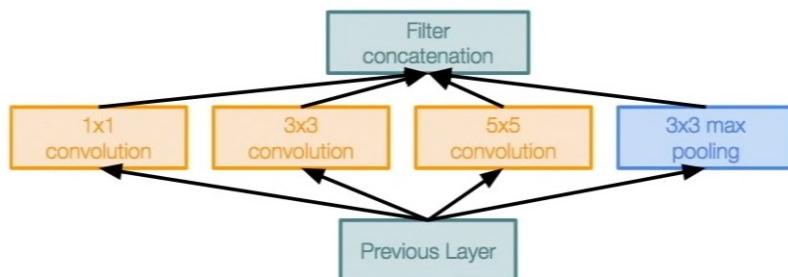
↳ VhG16 or VhG19

↳ Use ensemble for best results

↳ FC7 good feature representation

GoogleNet:

↳ Used Inception modules

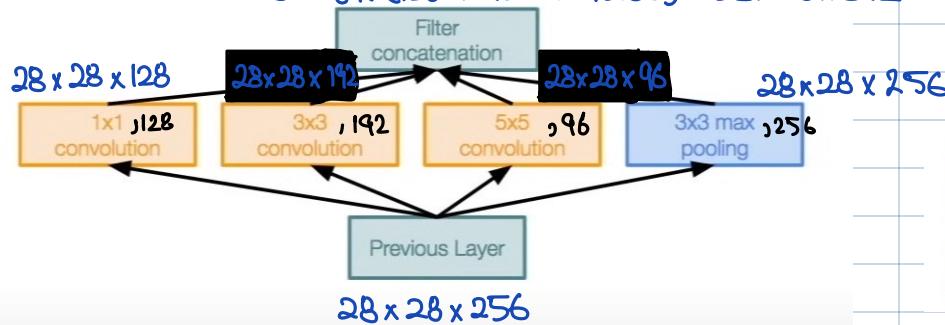


Naive Inception module

Q. What's the issue?

↳ Computational complexity

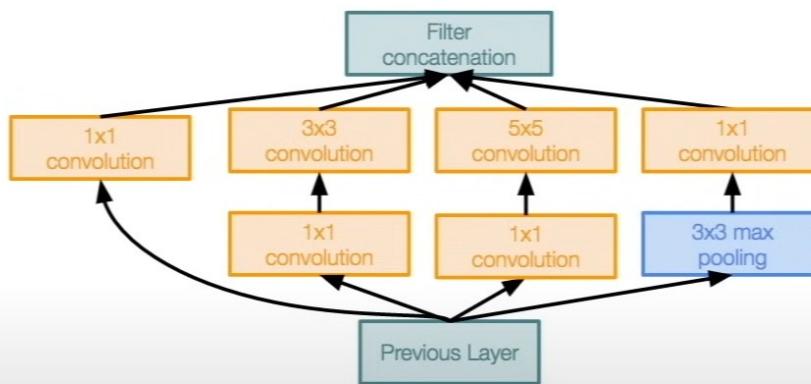
$$28 \times 28 \times (128 + 192 + 96 + 256) = 28 \times 28 \times 672$$



Naive Inception module

Solution:

↳ Project input depth to lowest depth using 1x1 convs



Inception module with dimension reduction

- ↳ Apply parallel filter operations on the input from previous layers.
- ↳ Multiple receptive field sizes for conv ($1 \times 1, 3 \times 3, 5 \times 5$)
- ↳ Pooling op (3×3)
- ↳ Concatenate all filters outputs depth wise

Conv Ops:

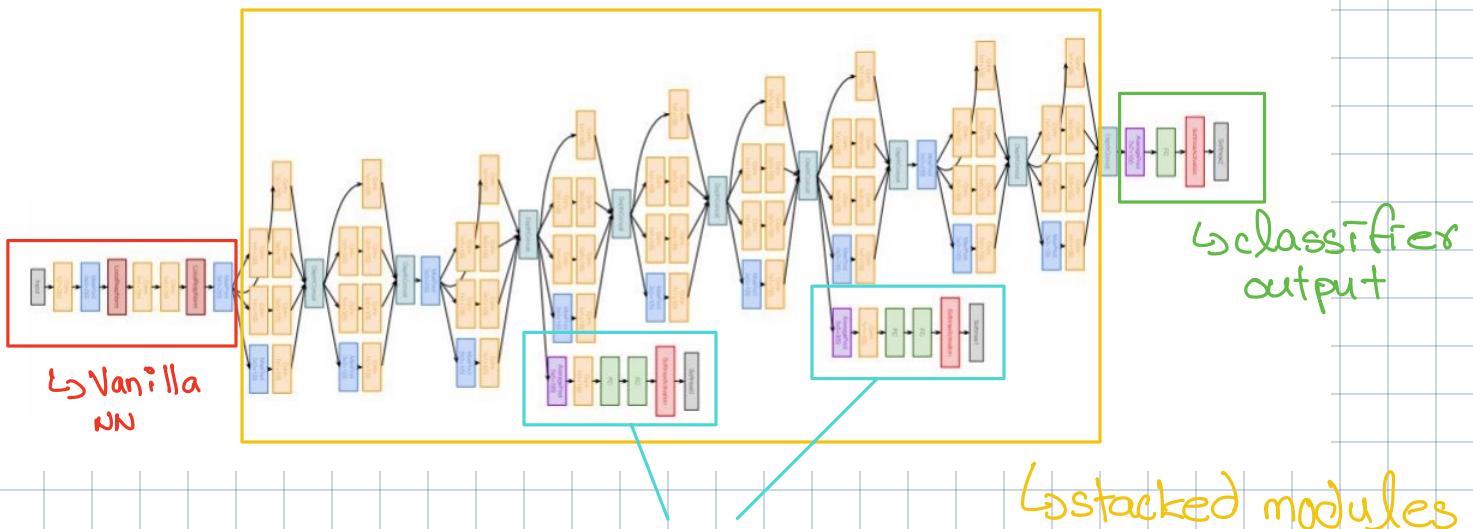
[1x1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$
 [3x3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 256$
 [5x5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 256$

Total: 854M ops

Conv Ops:

[1x1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$
 [1x1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$
 [1x1 conv, 128] $28 \times 28 \times 128 \times 1 \times 1 \times 256$
 [3x3 conv, 192] $28 \times 28 \times 192 \times 3 \times 3 \times 64$
 [5x5 conv, 96] $28 \times 28 \times 96 \times 5 \times 5 \times 64$
 [1x1 conv, 64] $28 \times 28 \times 64 \times 1 \times 1 \times 256$

Total: 358M ops



Details:

- ↳ 22 layers
- ↳ Efficient "Inception" module
- ↳ No FC layers
- ↳ 12x less params than Alex Net
- ↳ 6.7% error

Aux classification outputs. Mini NN.

Helps amplifying grads as the grads from end may end up becoming minimized when moving backward

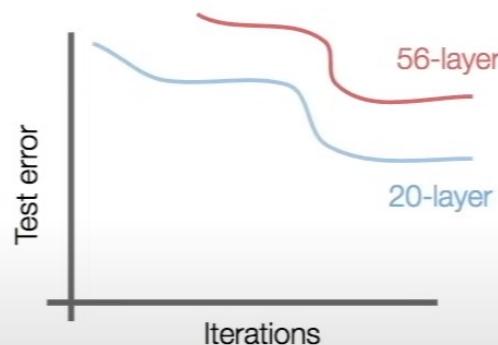
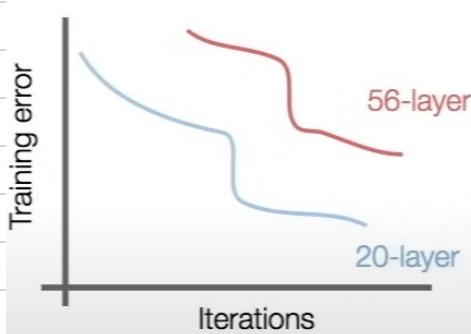
↳ stacked modules

ResNet: (2015)

- ↳ 152 layers model
- ↳ 3.57% error

Q. What happens when we continue stacking deeper layers on a plain convolutional neural network?

↳ No.



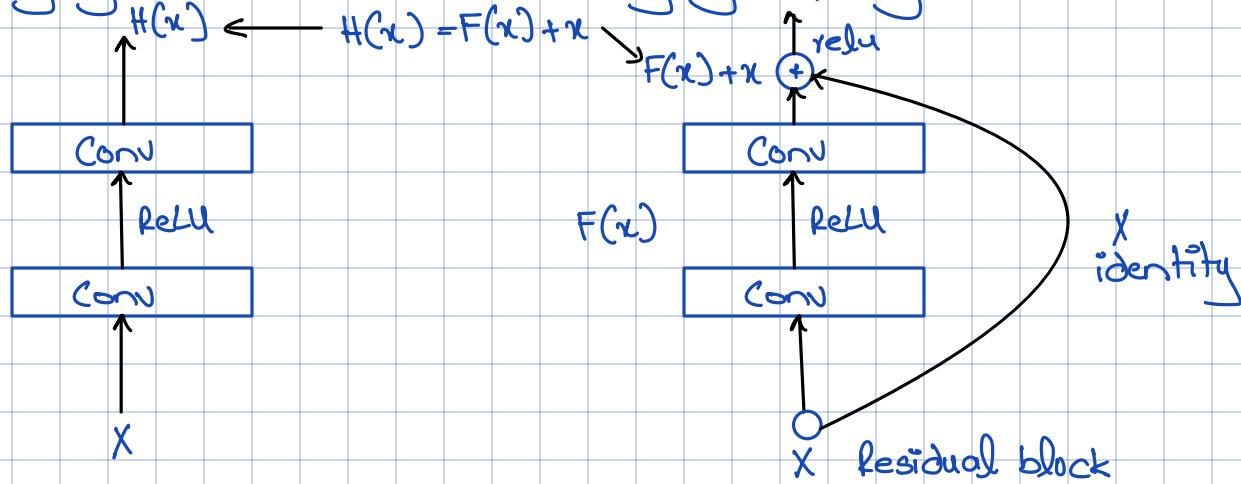
↳ Even though 56 layers model performs worse, it's NOT due to overfitting

↳ The problem is an optimization problem

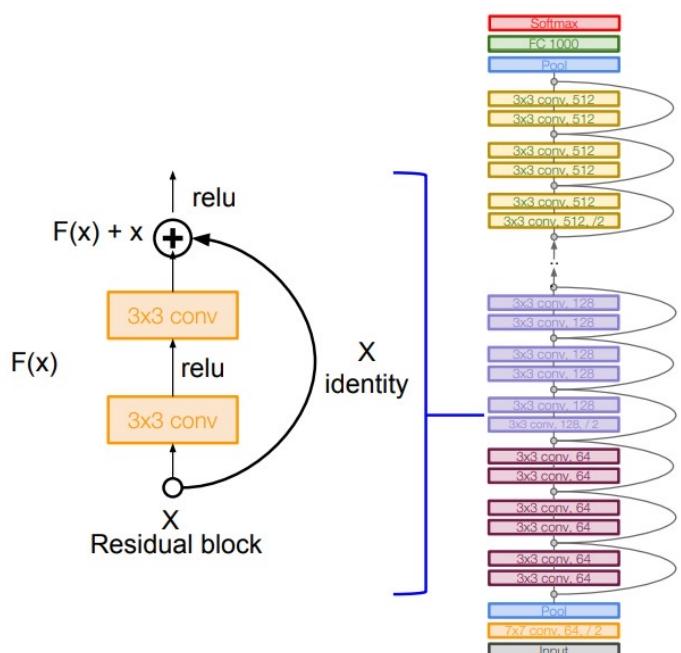
↳ A solution by construction is copying the learned layers from the shallower model and setting additional layers to identity mapping

Solution:

↳ Use network layers to fit a residual mapping instead of directly trying to fit a desired underlying mapping



↳ Use layers to fit residual $F(x) = H(x) - x$ instead of $H(x)$ directly



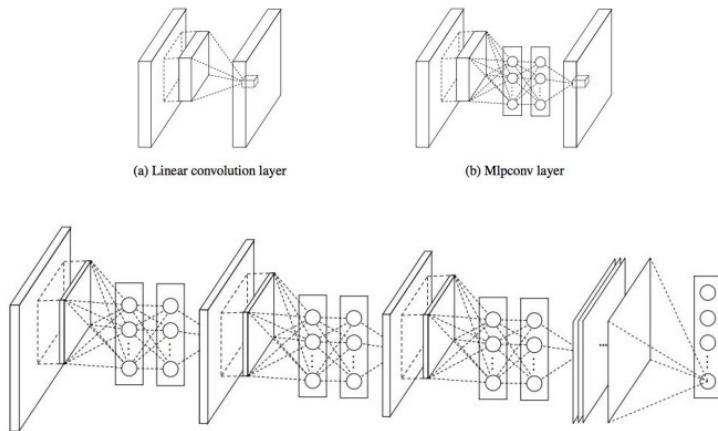
- ↳ Stack residual blocks
- ↳ Every residual block has two 3×3 conv layers
- ↳ Periodically, double no. of filters
- ↳ ↓ downsample spatially using stride 2
- ↳ Additional conv layer at the beginning
- ↳ No FC
- ↳ Total depths: 34, 50, 152
- ↳ For ResNet 50+, use 1×1 conv layers for efficiency

Training ResNet:

- ↳ Batch Norm after every CONV layer
- ↳ Xavier initialization
- ↳ SGD + momentum (0.9)
- ↳ LR: 0.1, / by 10 when VE plateaus
- ↳ Batch: 256
- ↳ Weight decay 1e-5
- ↳ No dropout

Network in Network:

↳ Convolves a more complex hierarchical filter

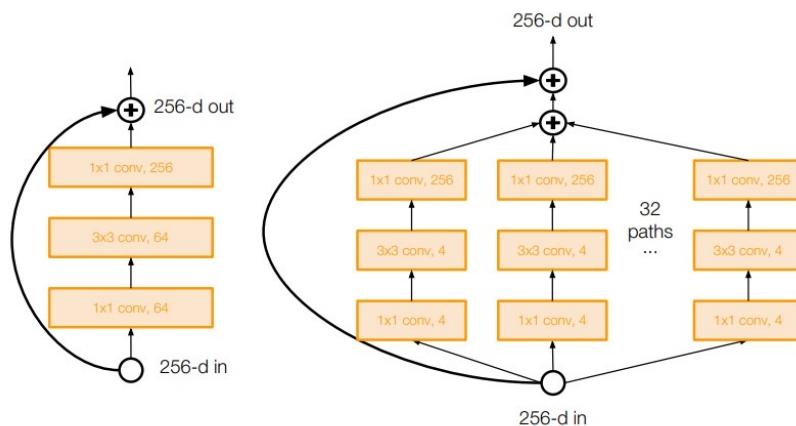


Wider Residual Networks:

- ↳ Argued residuals were important & NOT depth of network
- ↳ They used wider residual blocks ($F \times k$ filters instead of F) i.e. More filters
- ↳ 50 layer wide ResNet outperforms 152 layers original ResNet

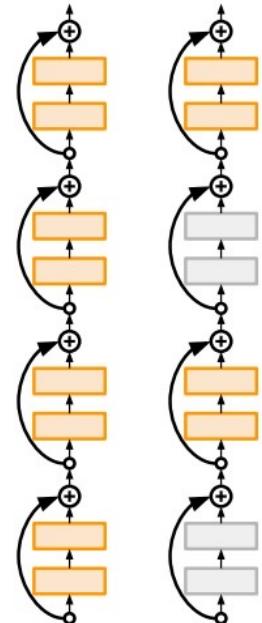
ResNeXt:

- ↳ Increases width of residual block through multiple parallel pathways "Cardinality"
- ↳ Parallel pathways similar in spirit to inception modules



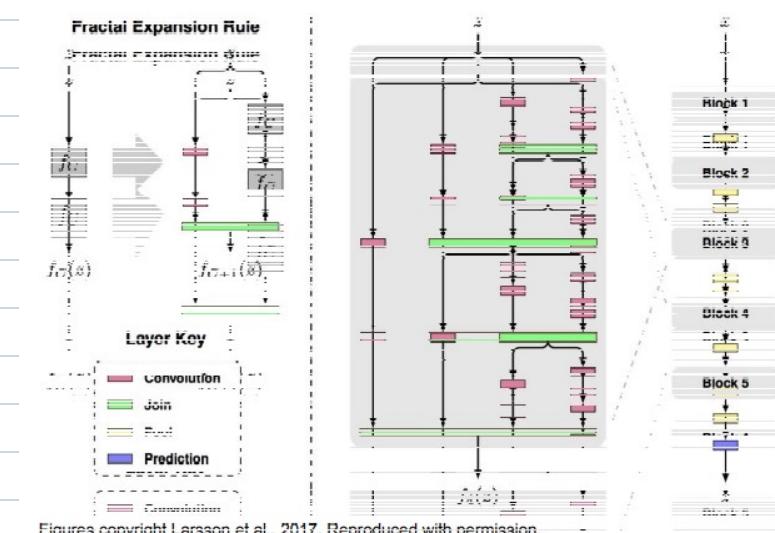
Stochastic Depth:

- ↳ Reduce vanishing grads
- ↳ Randomly drop a subset of layers
- ↳ Bypass with Identity function
- ↳ Use full deep networks at test time



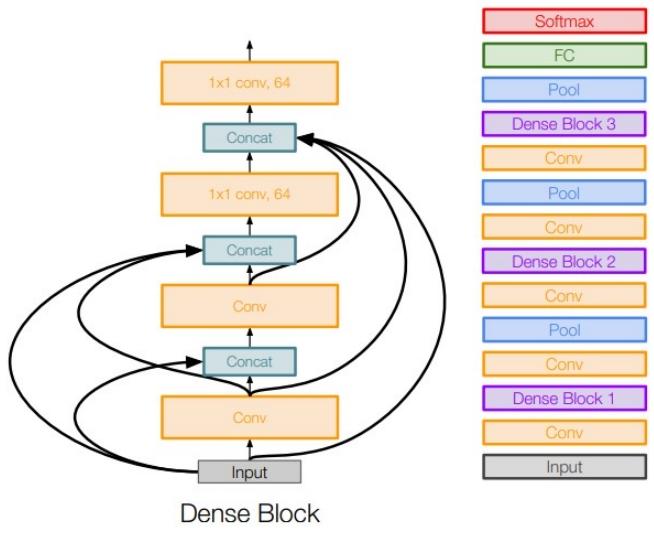
Fractal Net:

- ↳ key is transitioning from shallow to deep networks \wedge NOT residual representations
- ↳ Fractal architecture with both shallow \wedge deep pathways to output
- ↳ Trained with dropping out sub-paths
- ↳ Full network at test time



Dense Nets:

- ↳ Dense blocks where each layer is connected to every other layer in feed forward fashion.
- ↳ Alleviates vanishing grads, strengthens feature prop, encourages feature use.



SqueezeNet: → efficient architectures

- ↳ Fire modules consisting of a squeeze layer with 1×1 filters feeding an 'expand' layer with 1×1 & 3×3 filters
- ↳ AlexNet accuracy with $50 \times$ fewer params
- ↳ $500 \times$ smaller than AlexNet

