

Visual Learning (Tutorial C2) Convolutional Networks

Instructor - Simon Lucey
RVSS - 2024



**AUSTRALIAN
INSTITUTE FOR
MACHINE LEARNING**

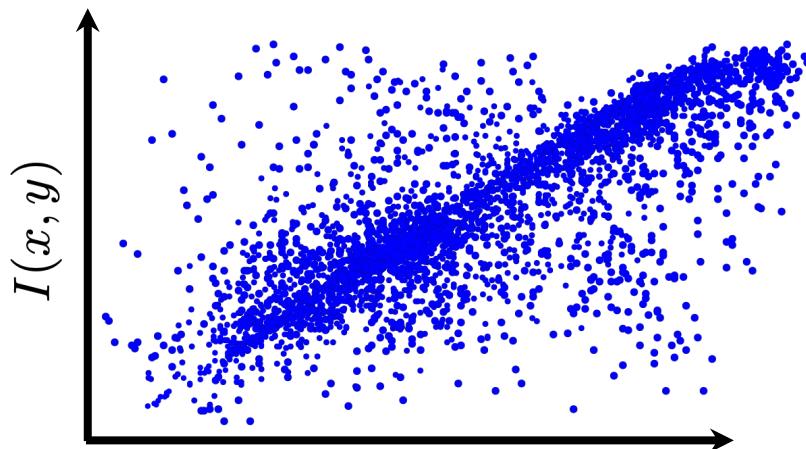
Today

- **ConvNets and AlexNet**
- Depthwise Convolution
- Skip/Residual Connections





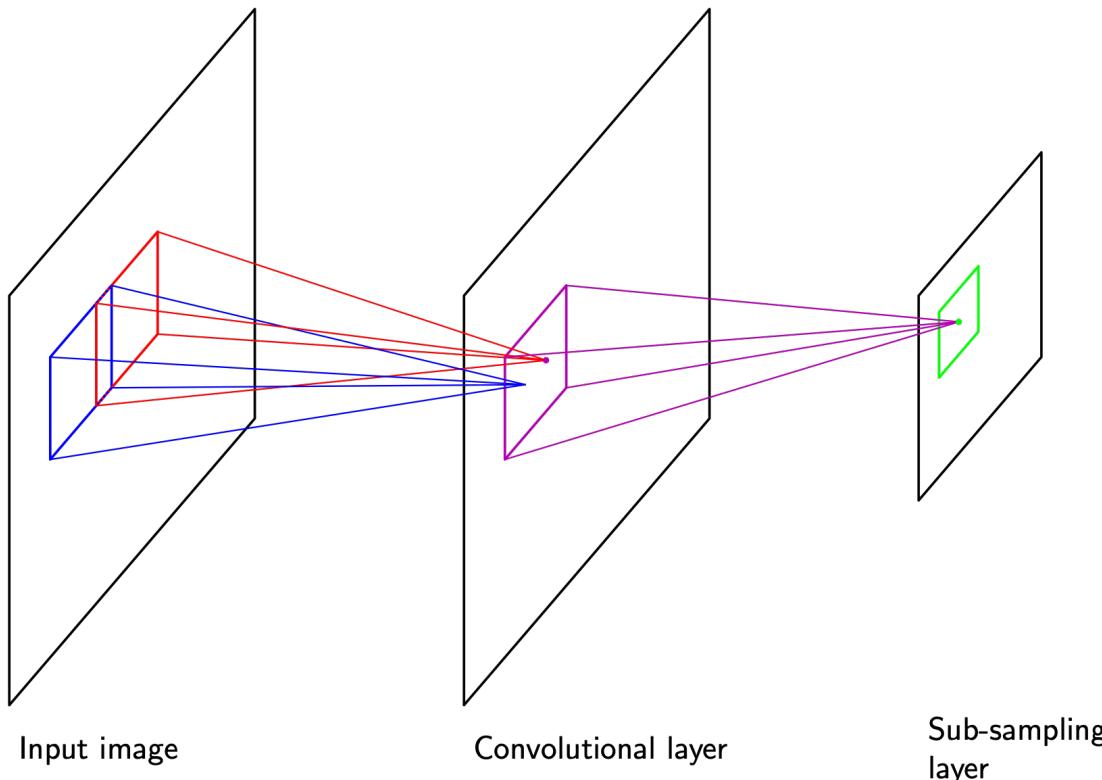
I



$I(x + \delta x, y + \delta y)$

Simoncelli & Olshausen 2001

Convolutional Neural Network

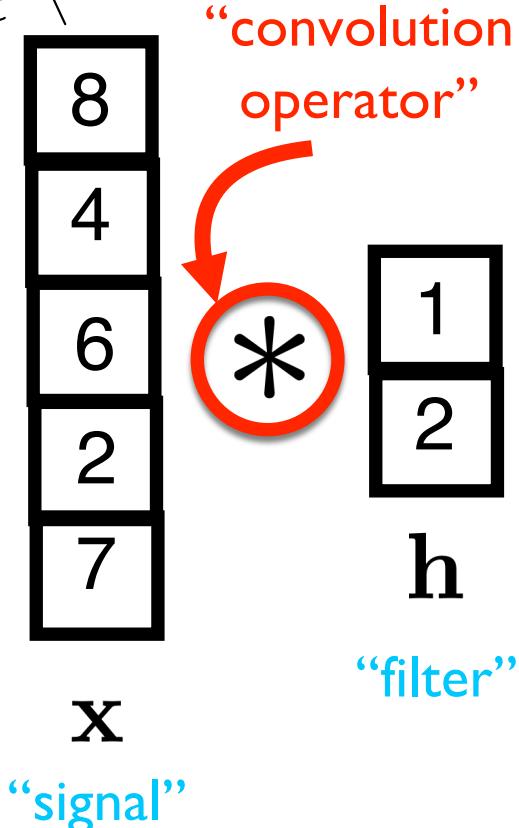


LeCun 1980

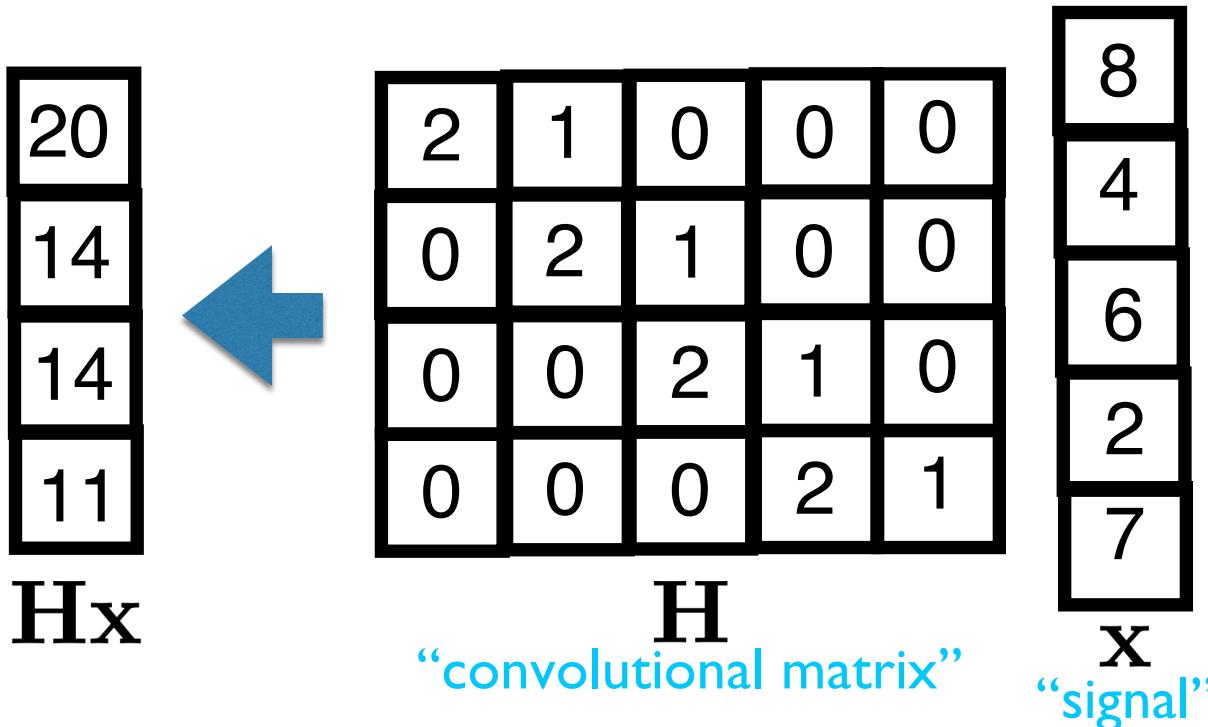
Reminder: Convolution

```
>>> from scipy.signal import \
    convolve as conv
>>> conv(x,h,'valid')
array([20, 14, 14, 11])
```

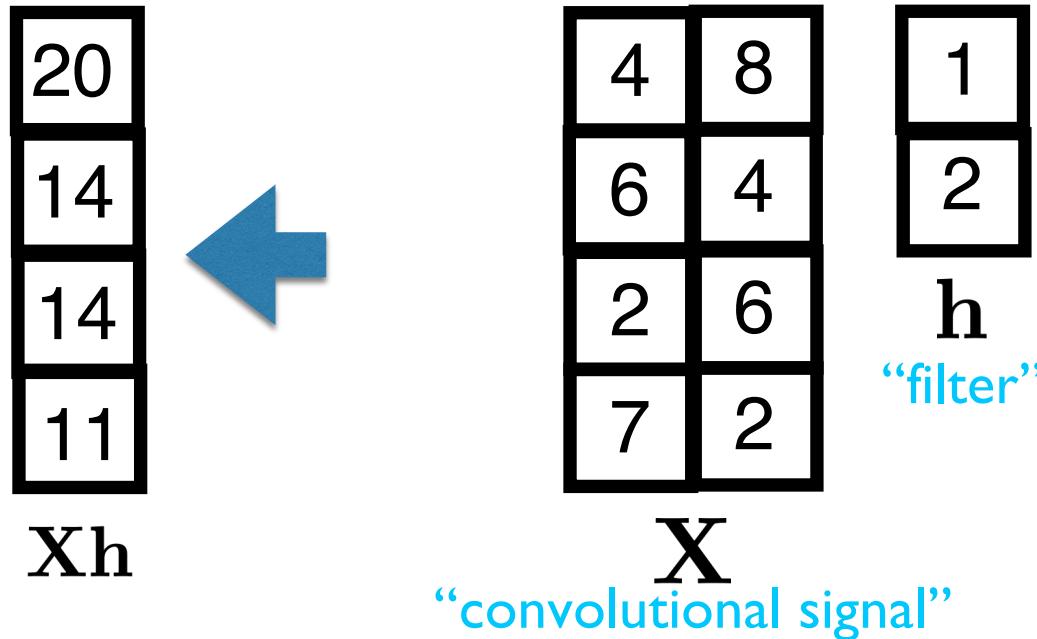
$$\frac{\partial[x * h]}{h^T} = ???$$



Reminder: Convolution



Reminder: Convolution



$$\frac{\partial[x * h]}{h^T} = \frac{\partial[Xh]}{h^T} = X^T$$

Efficiency of Convolution

Input size: 320 by 280

Filter size: 2 by 1

Output size: 319 by 280

	2	$319*280*320*28 > 8e9$	$2*319*280 = 178,640$
	$319*280*3 = 267,960$	$> 16e9$	Same as convolution (267,960)

Vectorizing 2D-Convolution

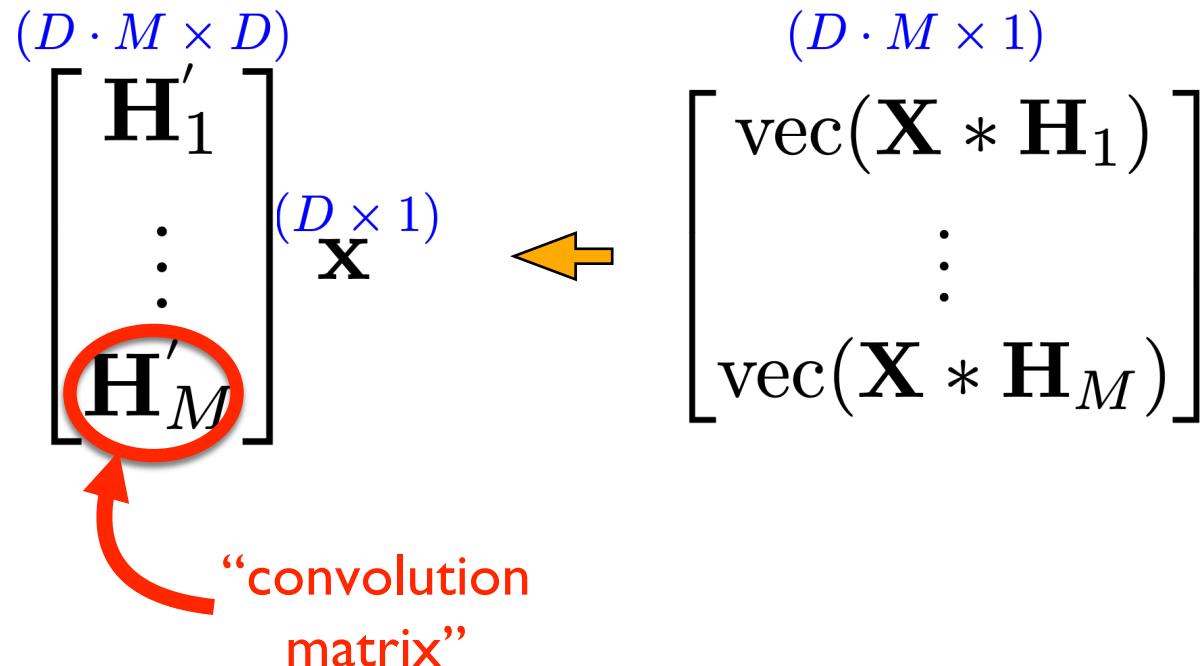
$$\text{vec}(\mathbf{Y}) = \text{vec}(\mathbf{X} * \mathbf{H})$$

```
def forward(self, x):
    # Flattens the image like structure into vector
    x = torch.flatten(x, start_dim=1)
```

Vectorizing 2D-Convolution

$$\text{vec}(\mathbf{Y}) = \mathbf{H}' \text{vec}(\mathbf{X})$$

Multiple Filters



Removing Redundancy - Striding

Stride: [4, 4]

207	245	77	21	247	211	240	1
219	41	58	179	161	154	184	98
215	145	187	71	251	249	65	100
192	2	189	247	166	63	232	213
105	94	66	190	156	61	89	145
159	154	87	184	101	105	72	71
192	111	6	94	60	70	65	226
175	120	210	226	80	183	168	184
134	56	36	240	159	178	76	135
239	244	199	9	132	104	188	185
245	210	78	199	0	92	9	246
5	121	187	122	107	47	12	119
230	171	135	36	82	54	65	37
61	140	79	19	161	96	127	187
56	223	46	6	180	186	142	244
28	20	61	2	178	187	98	220

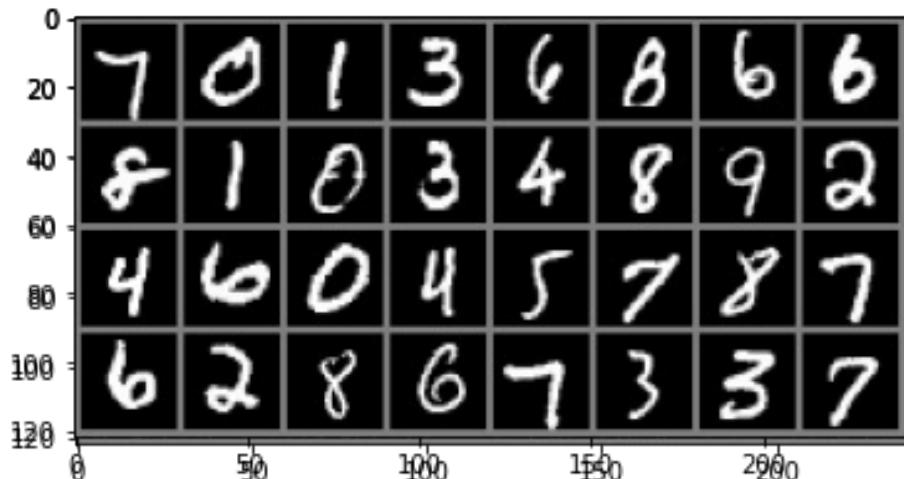
Let's have another play!!!



Some things to try!!!

- What happens to performance if you use a linear function?
- What happens when you permute the pixels?

```
from numpy.random import permutation
idx_permute = torch.from_numpy(permutation(784))
transform = transforms.Compose([transforms.ToTensor(),
                               transforms.Lambda(lambda x: x.view(-1)[idx_permute].view(1, 28, 28) ),
                               transforms.Normalize((0.5,), (0.5,)),
                               ])
```



Adding Striding - CNN

```
class ConvNet(nn.Module):
    def __init__(self):
        super(ConvNet, self).__init__()
        self.conv1 = nn.Conv2d(1, 6, 5, 2)
        self.conv2 = nn.Conv2d(6, 16, 5, 2)
        self.fc1 = nn.Linear(256, 84)
        self.fc2 = nn.Linear(84, 10)

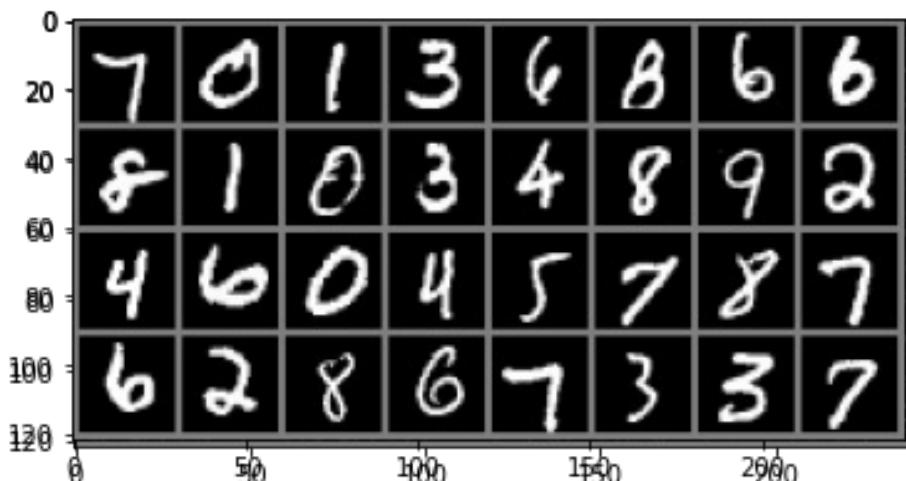
    def forward(self, x):
        # Input goes to convolution so no need to
        x = self.conv1(x)
        x = F.relu(x)
        x = self.conv2(x)
        x = F.relu(x)
        x = torch.flatten(x, start_dim=1)
        x = self.fc1(x)
        x = F.relu(x)
        x = self.fc2(x)
```

$$\eta(\mathbf{W}^{(1)}\eta(\mathbf{W}^{(0)}\mathbf{x} + \boldsymbol{\delta}^{(0)}) + \boldsymbol{\delta}^{(1)})$$

Permuting Pixels - CNN

- What happens when you permute the pixels in a CNN?

```
from numpy.random import permutation
idx_permute = torch.from_numpy(permutation(784))
transform = transforms.Compose([
    transforms.ToTensor(),
    transforms.Lambda(x: x.view(-1)[idx_permute].view(1, 28, 28)),
    transforms.Normalize((0.5,), (0.5,)),
])
```



ImageNet Classification with Deep Convolutional Neural Networks

Alex Krizhevsky

University of Toronto

kriz@cs.utoronto.ca

Ilya Sutskever

University of Toronto

ilya@cs.utoronto.ca

Geoffrey E. Hinton

University of Toronto

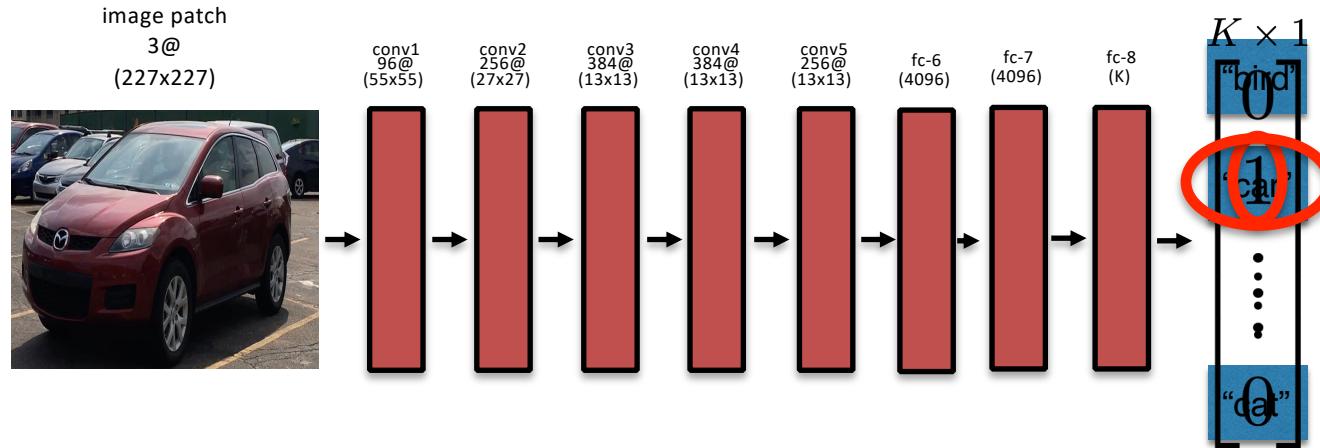
hinton@cs.utoronto.ca

Abstract

We trained a large, deep convolutional neural network to classify the 1.2 million high-resolution images in the ImageNet LSVRC-2010 contest into the 1000 different classes. On the test data, we achieved top-1 and top-5 error rates of 37.5% and 17.0% which is considerably better than the previous state-of-the-art. The neural network, which has 60 million parameters and 650,000 neurons, consists of five convolutional layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. To make training faster, we used non-saturating neurons and a very efficient GPU implementation of the convolution operation. To reduce overfitting in the fully-connected layers we employed a recently-developed regularization method called “dropout” that proved to be very effective. We also entered a variant of this model in the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%, compared to 26.2% achieved by the second-best entry.

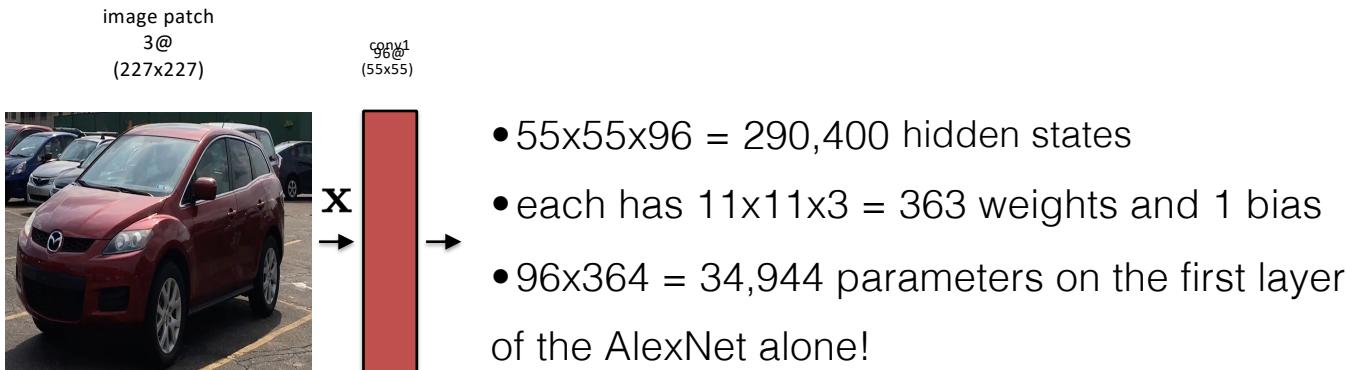
AlexNet

- AlexNet won the ILSVRC-2012 competition and achieved a winning top-5 test error rate of 15.3%. (Second best was 26.2 %).
- Network has 25 layers, but only 8 layers with learnable weights.
 - 5 convolutional weights.
 - 3 fully connected weights.

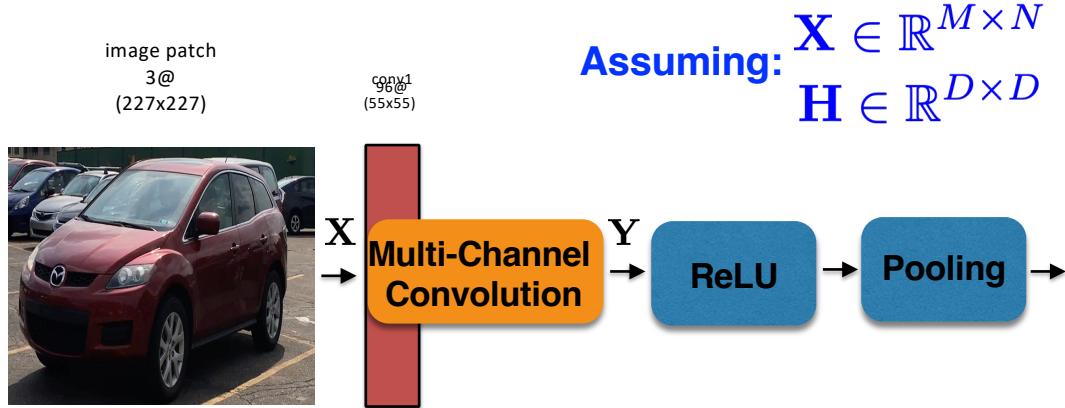


AlexNet Drawbacks

- AlexNet gave amazing performance but had some drawbacks.
- Fundamentally, it is extremely expensive to evaluate & learn.



Reminder: Multi-Channel Convolution



Assuming: $\mathbf{X} \in \mathbb{R}^{M \times N}$
 $\mathbf{H} \in \mathbb{R}^{D \times D}$

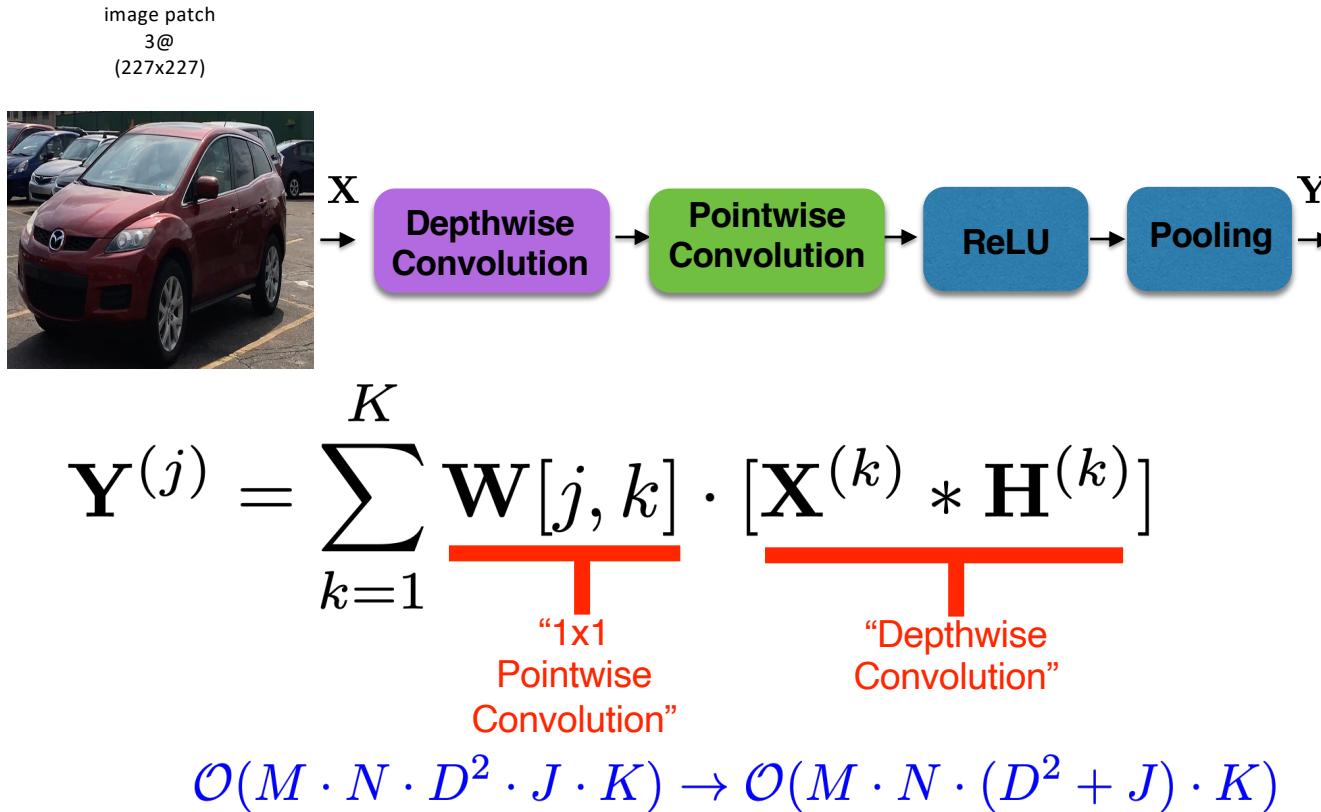
$$\mathbf{Y}^{(j)} = \sum_{k=1}^K \mathbf{X}^{(k)} * \mathbf{H}^{(j,k)}$$

Naive computational cost - $\mathcal{O}(M \cdot N \cdot D^2 \cdot K \cdot J)$

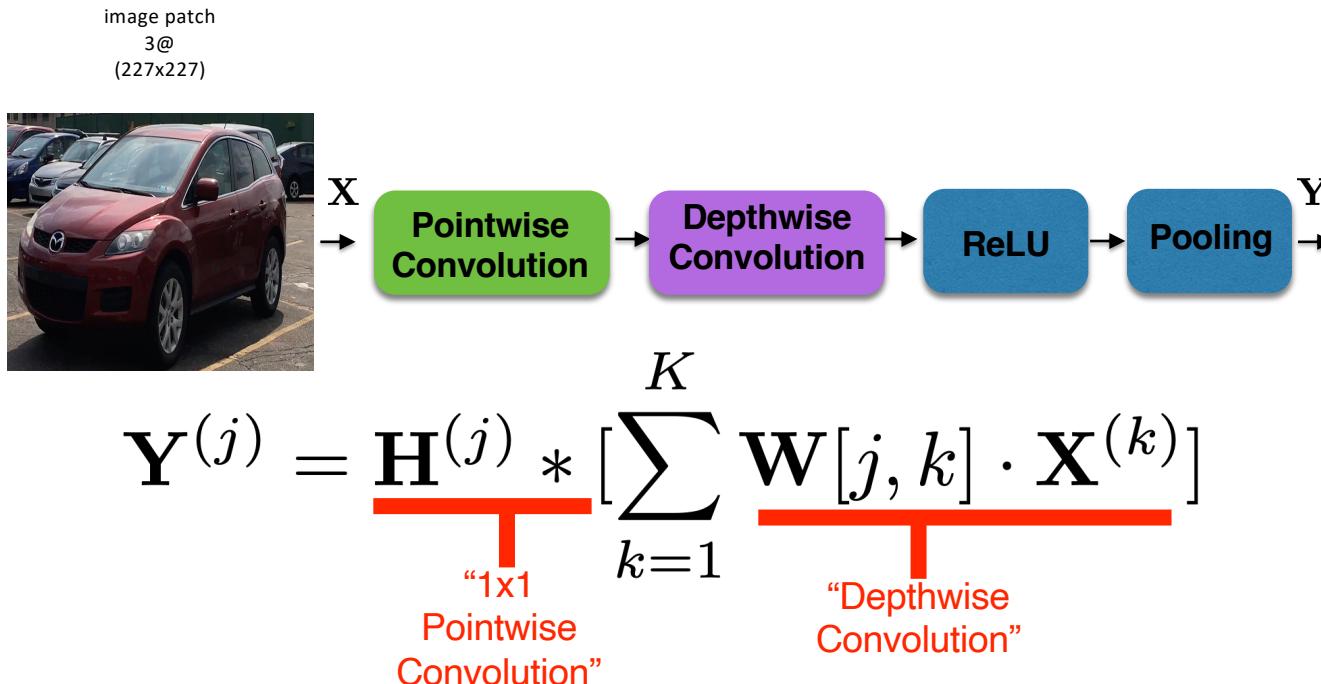
Today

- ConvNets and AlexNet
- **Depthwise Convolution**
- Skip/Residual Connections

Depthwise Convolution



Depthwise Convolution



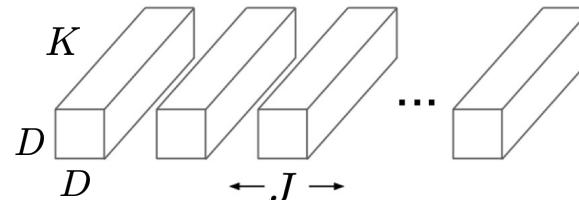
$$\mathcal{O}(M \cdot N \cdot D^2 \cdot J \cdot K) \rightarrow \mathcal{O}(M \cdot N \cdot (D^2 + J) \cdot K)$$

MobileNet

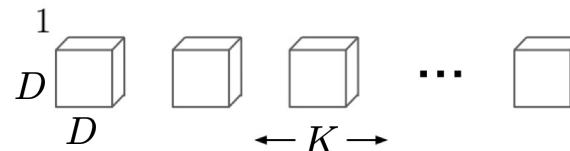
- MobileNet: (a) is divided into (b) and (c).
- Computation reduction

$$\frac{b+c}{a} = \frac{D^2 + J}{D^2 \cdot J}$$

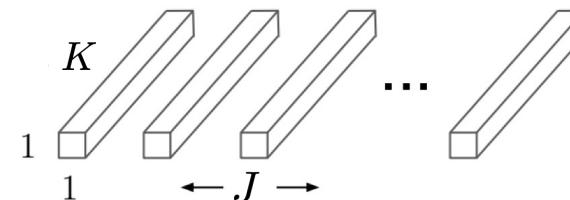
- D = 3: 8~9 times less computation.



(a) Standard convolution filters



(b) Depthwise convolution filters



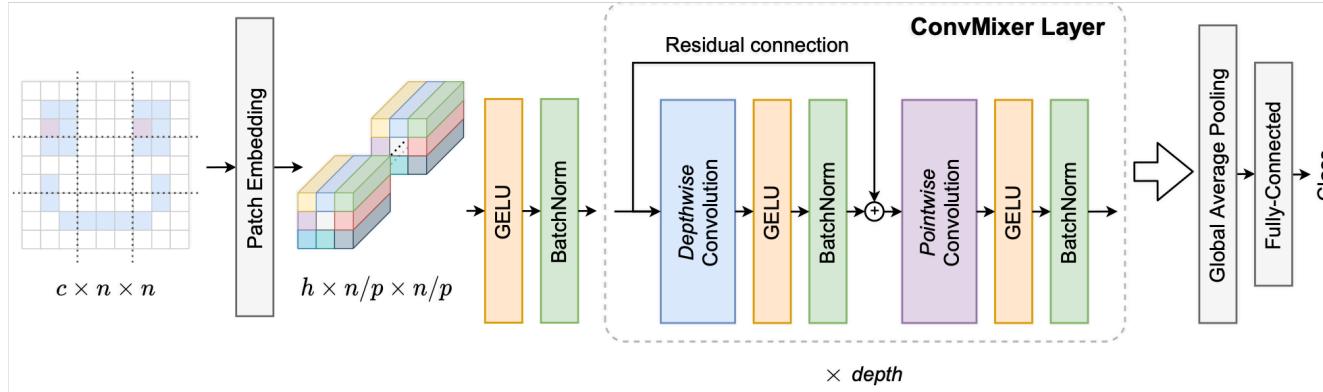
(c) 1x1 pointwise convolution

MobileNet

- MobileNet with similar accuracy but less computation and fewer parameters than VGG16 and GoogleNet
- $\alpha = 0.5$, input 160×160 : better than AlexNet while being 45 times smaller and 9.4 times less compute than AlexNet

Model	ImageNet	Million	Million
	Accuracy	Mult-Adds	Parameters
1.0 MobileNet-224	70.6%	569	4.2
GoogleNet	69.8%	1550	6.8
VGG 16	71.5%	15300	138
0.50 MobileNet-160	60.2%	76	1.32
SqueezeNet	57.5%	1700	1.25
AlexNet	57.2%	720	60

ConvMixer - Patches are all you need



```
1 def ConvMixer(h, depth, kernel_size=9, patch_size=7, n_classes=1000):
2     Seq, ActBn = nn.Sequential, lambda x: Seq(x, nn.GELU(), nn.BatchNorm2d(h))
3     Residual = type('Residual', (Seq,), {'forward': lambda self, x: self[0](x) + x})
4     return Seq(ActBn(nn.Conv2d(3, h, patch_size, stride=patch_size)),
5               *[Seq(Residual(ActBn(nn.Conv2d(h, h, kernel_size, groups=h, padding="same"))),
6                     ActBn(nn.Conv2d(h, h, 1))) for i in range(depth)],
7               nn.AdaptiveAvgPool2d((1,1)), nn.Flatten(), nn.Linear(h, n_classes))
```

Let's Have a Play

```
class ConvNet(nn.Module):
    def __init__(self):
        super(ConvNet, self).__init__()

        # Two convolution layers I am writing the first one
        # First convolutional layer takes single chennel images (batch_size specify the number of images) as input
        # We have 5x5 convolutions
        # We have 6 convolutional filter to produce output size 6*28*28 for a single training sample.
        # structure is : nn.conv2d(number of input channels, number of filters, conv kernel size, stride = 1)
        self.conv1 = nn.Conv2d(1, 6, 5, 1)
        # Note that Nparam 1*6*5*5 = 150 (+ 5 for bias per output).

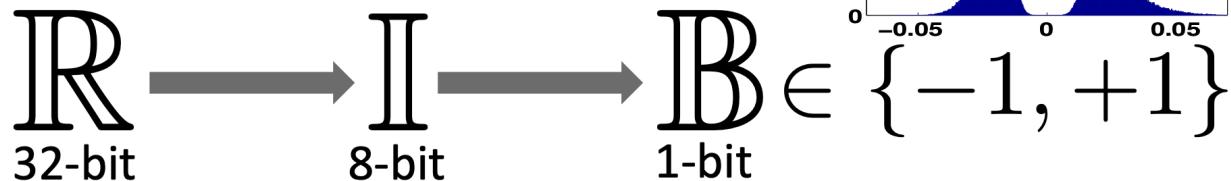
        #####
        # TODO: add another layer called self.conv2, 5x5 convolutions 16 filters in total.
        #####
        #self.conv2 = nn.Conv2d(6, 16, 5, 1)
        self.mlp2 = nn.Conv2d(6, 16, 1, 1)
        self.conv2 = nn.Conv2d(16, 16, 5, 1, groups=16)

        # max max fully connected layer: maximum input size - output size
```

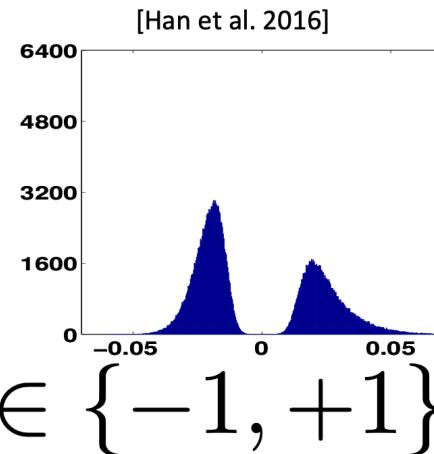
Trends - Lower Precision

Reducing Precision

- Saving Memory
- Saving Computation



$\{-1,+1\}$	$\{0,1\}$
MUL	XNOR
ADD, SUB	Bit-Count (popcount)

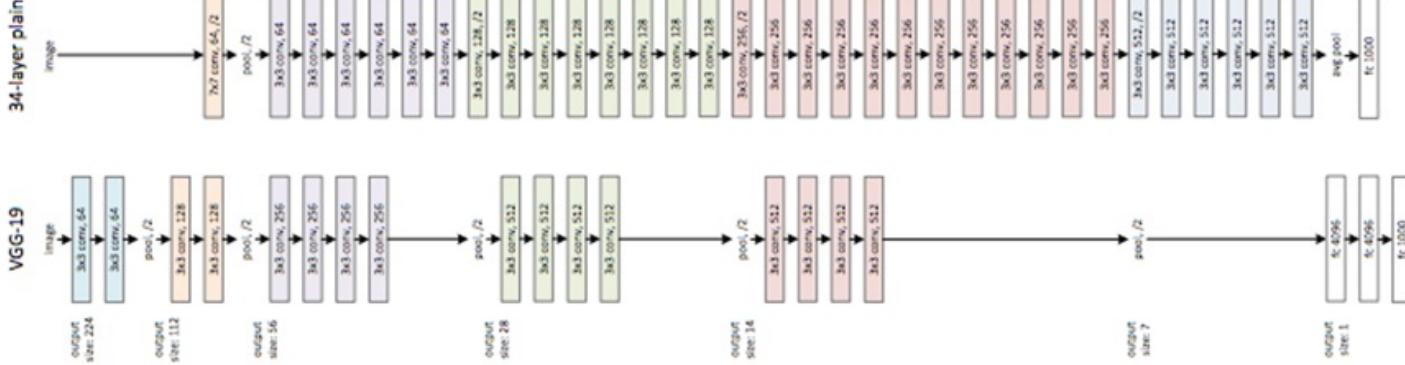


Today

- ConvNets and AlexNet
- Depthwise Convolution
- **Skip/Residual Connections**

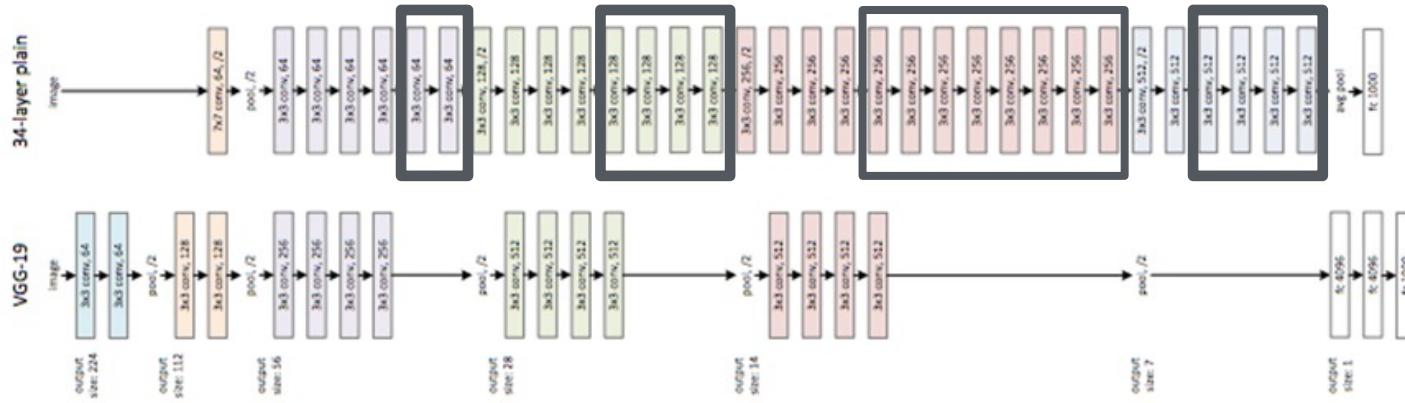
Neural Net : Deeper is Better?

- Very deep network -- training error increases.



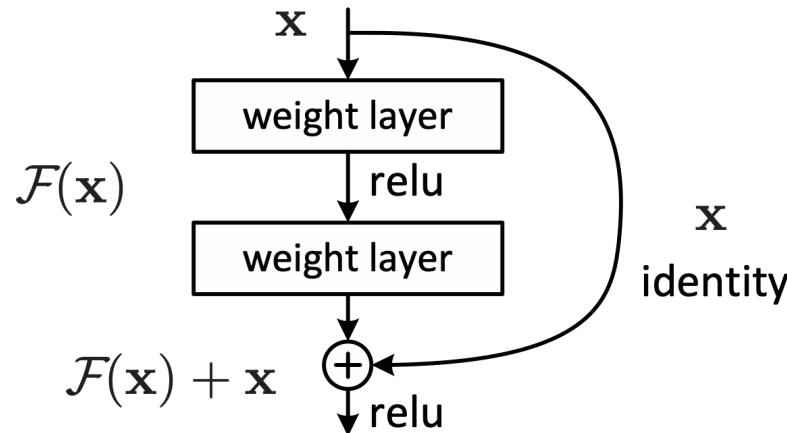
Neural Net : Deeper is Better?

- Very deep network -- training error increases.
 - What if I just learn to do NOT A THING in some layer?
 - 34 layer net then is just a 19 layer net with less loss!

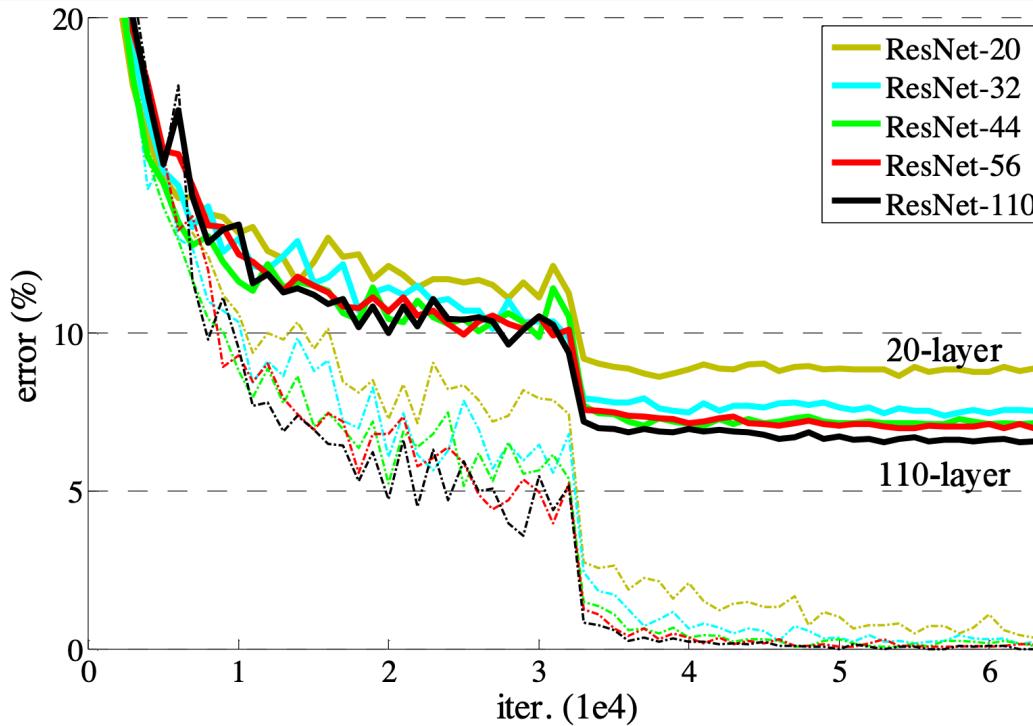


How Deep?

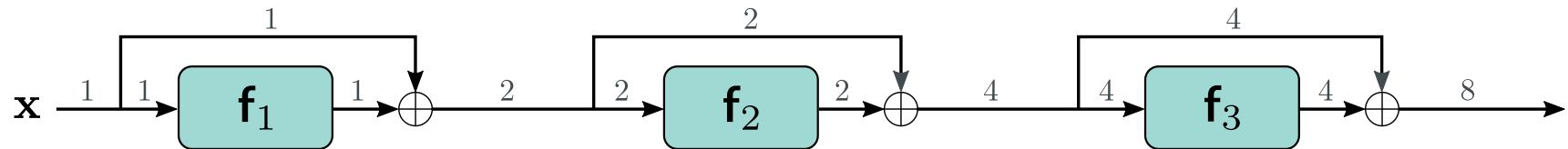
- Recent work has suggested that network depth is crucial for good performance (e.g. ImageNet).
- Counter intuitively, naively trained deeper networks tend to have higher train error than shallow networks.
- Innovation of residual learning has greatly helped with this.



How Deep?



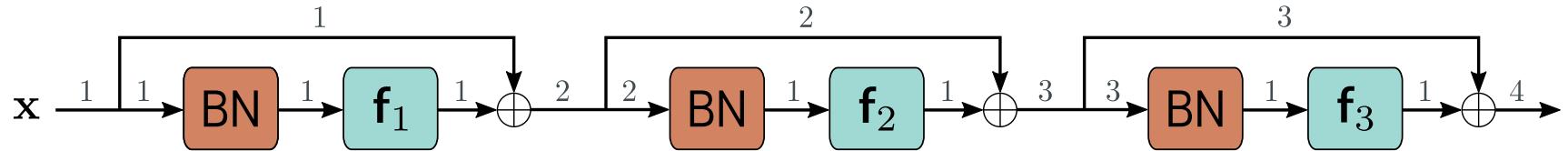
Initialization Problem



Try example:-

https://colab.research.google.com/github/udlbook/udlbook/blob/main/Notebooks/Chap11/11_3_Batch_Normalization.ipynb

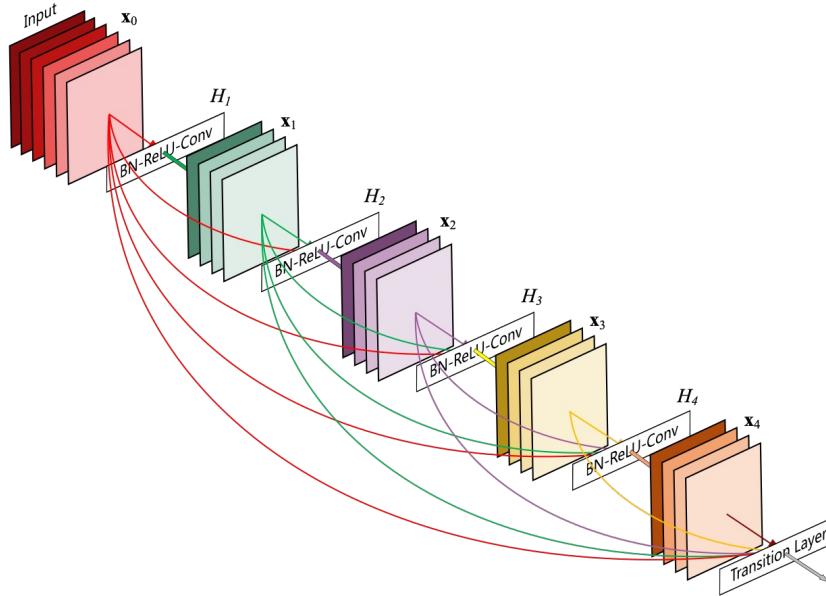
Batch Norm



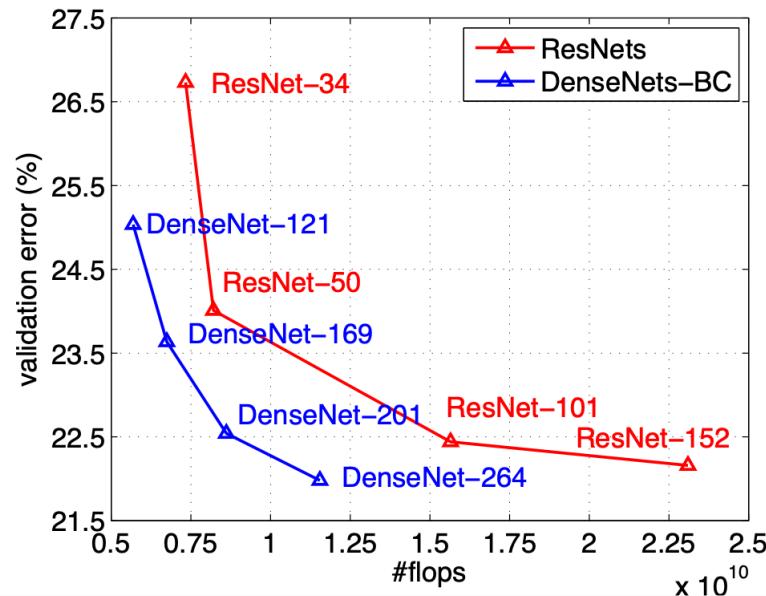
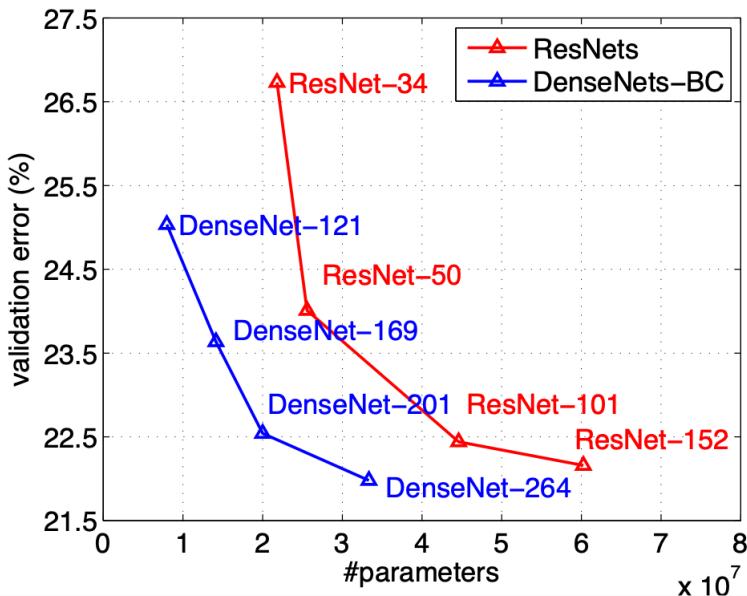
$$m_h = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} h_i \quad h_i \leftarrow \frac{h_i - m_h}{s_h + \epsilon}$$

$$s_h = \sqrt{\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} (h_i - m_h)^2}, \quad h_i \leftarrow \gamma h_i + \delta$$

Densely Connected Networks

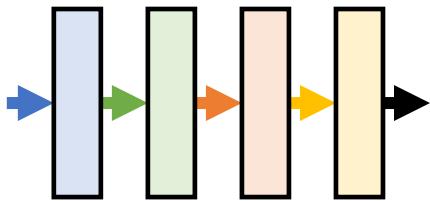


Densely Connected Networks

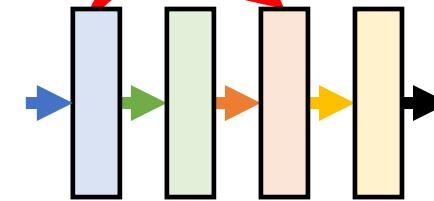


Network
Architecture

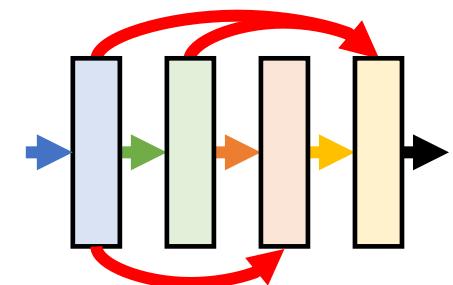
Alexnet



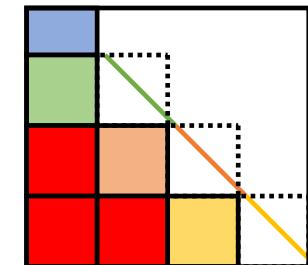
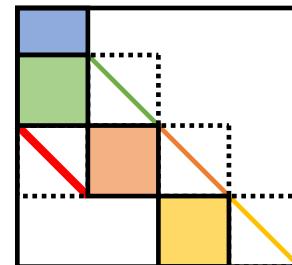
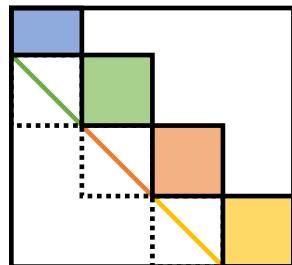
ResNet



DenseNet



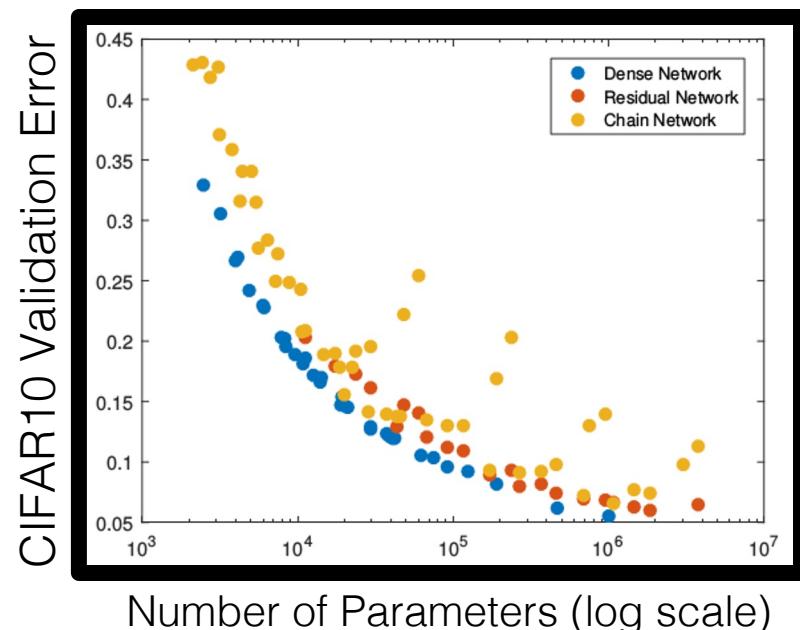
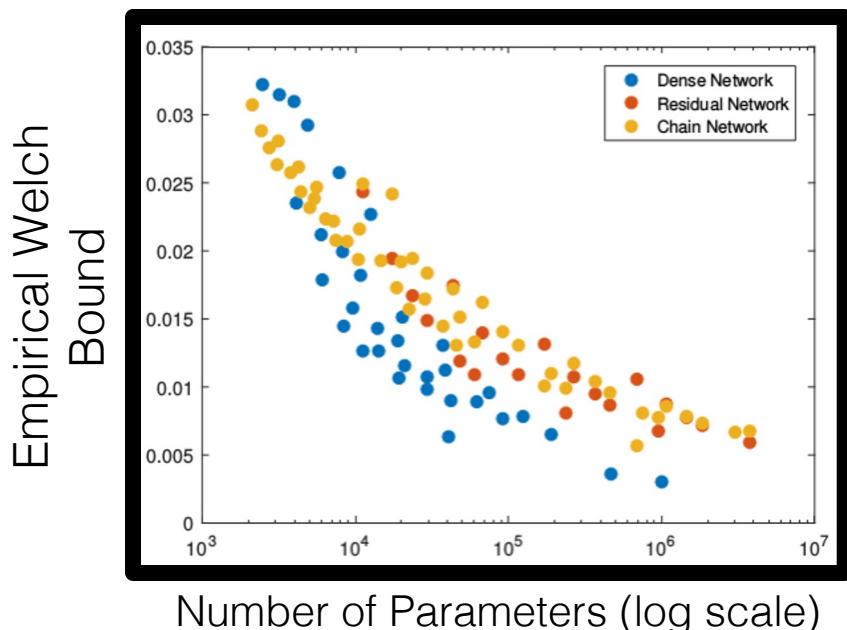
Induced
Dictionary
Structure



Gao Huang, Zhuang Liu, Kilian Weinberger, and Laurens van der Maaten.
“Densely Connected Convolutional Networks.” CVPR, 2017.

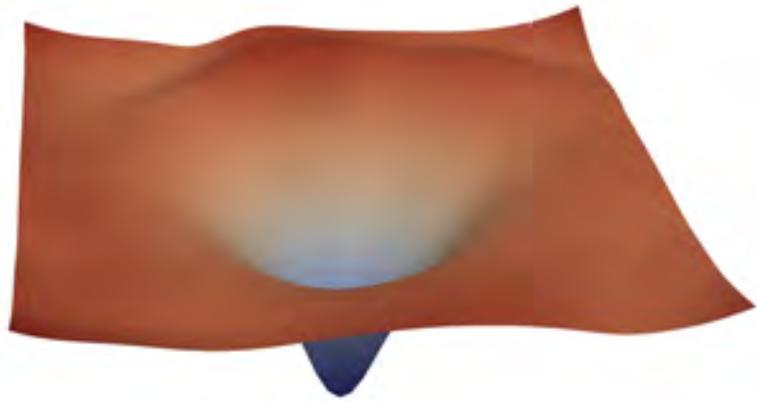
Experimental Results

- Dense skip connections induce dictionary structures with lower Welch bounds and give better performance with fewer parameters.



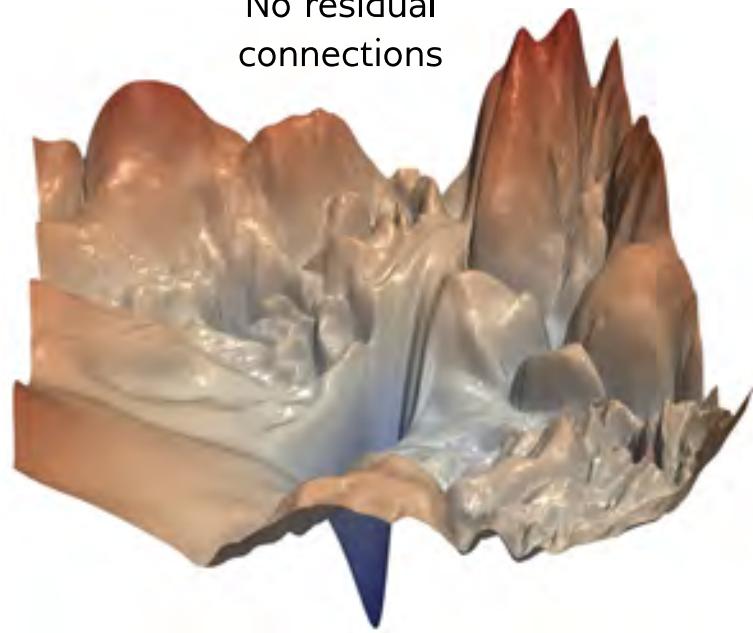
a)

Residual
connections

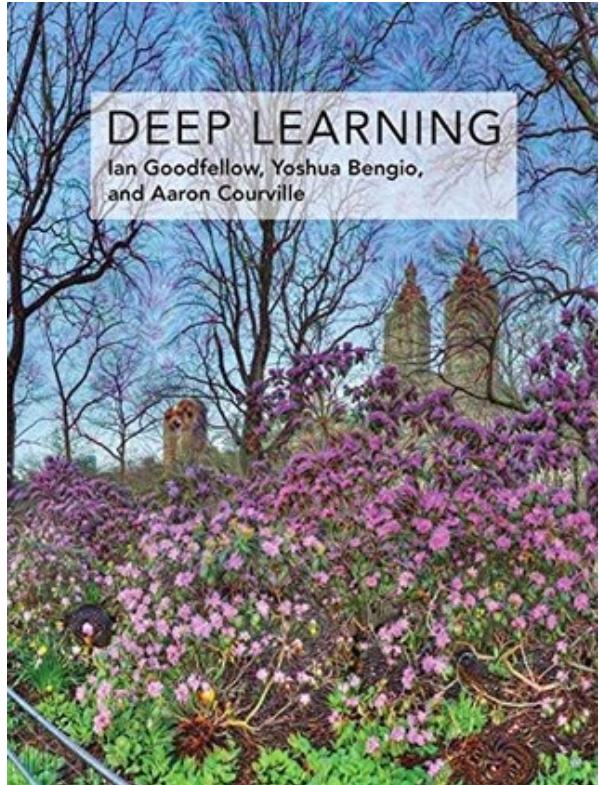


b)

No residual
connections



More to read...



DEEP LEARNING

Ian Goodfellow, Yoshua Bengio,
and Aaron Courville

SIMON J. D. PRINCE

Understanding Deep Learning