

BITS F464: Machine Learning

37

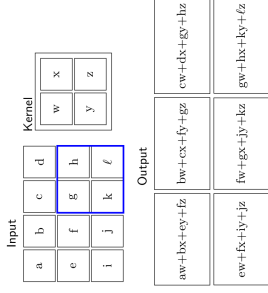
Convolution Neural Network (CNN)



Dr. Kamlesh Tiwari
 Assistant Professor, Department of CSIS,
 BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA
 April 28, 2021 **ONLINE** (Campus @ BITS-Pilani Jan-May 2021)

<http://kti.wari.in/ml>

Image is a 2D signal



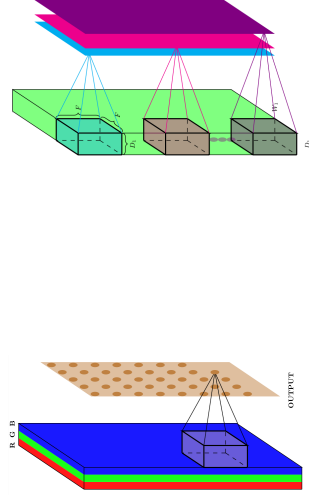
$$S_{ij} = (I * K)_{ij} = \sum_{a=0}^{m-1} \sum_{b=0}^{n-1} I_{i+a,j+b} K_{a,b}$$

General formula

$$S_{ij} = (I * K)_{ij} = \sum_{a=-\frac{K}{2}}^{\frac{K}{2}} \sum_{b=-\frac{K}{2}}^{\frac{K}{2}} I_{i-a,j-b} K_{a,b}$$

Filter in 3D

- It would refer to volume. Assume the filter extends to the depth
- Output is 2D



- Apply multiple filters to get multiple feature maps

CNN Motivation

Suppose you have a **noisy** sensor to have a measurement

How could you estimate the actual reading

- 1 Average of some readings
- 2 in local neighborhood
- 3 What about weighted average?

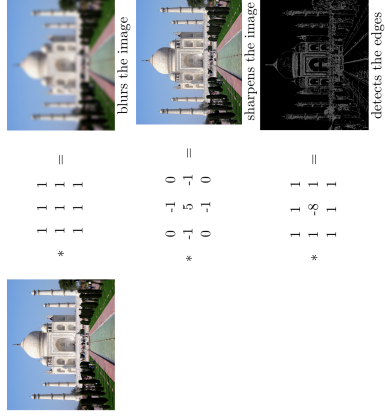
$$s_i = \sum_{l=0}^{\infty} x_{i-l} w_{-l}$$

w_{-6}	w_{-5}	w_{-4}	w_{-3}	w_{-2}	w_{-1}	w_0			
0.01	0.01	0.02	0.02	0.04	0.4	0.5			
X	1.00	1.10	1.20	1.40	1.70	1.80	2.10	2.20	2.40
S						1.80			

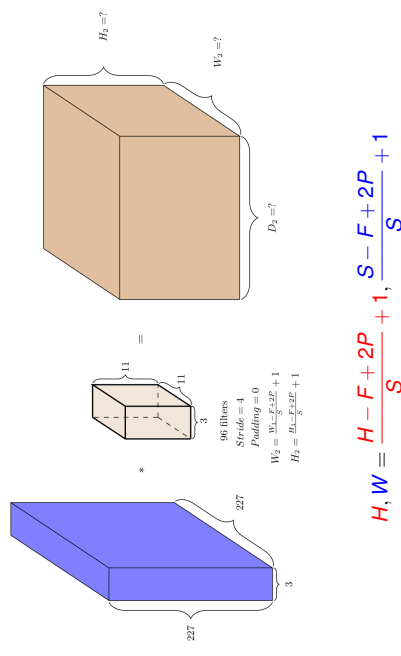
$$s_6 = x_6 w_0 + x_5 w_{-1} + x_4 w_{-2} + x_3 w_{-3} + x_2 w_{-4} + x_1 w_{-5} + x_0 w_{-6}$$

Here things are in one dimension only.

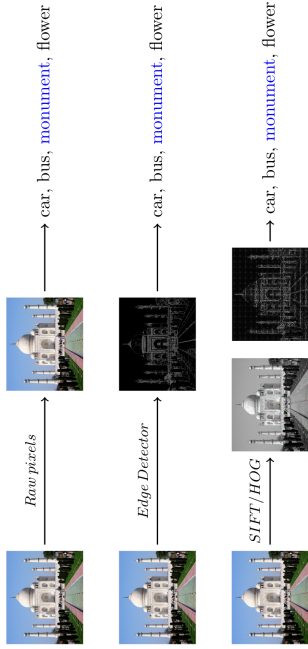
Image is a 2D signal



Filters, Padding and Stride

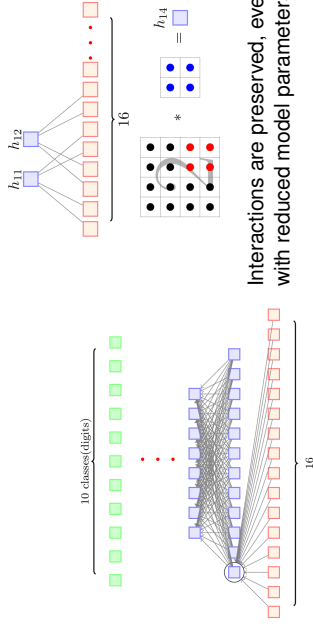


Classification Pipeline

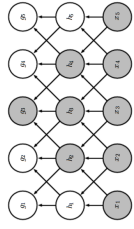


- Where is learning? (hand craft features, then learn weights for classification). One can see feature as a convolution.

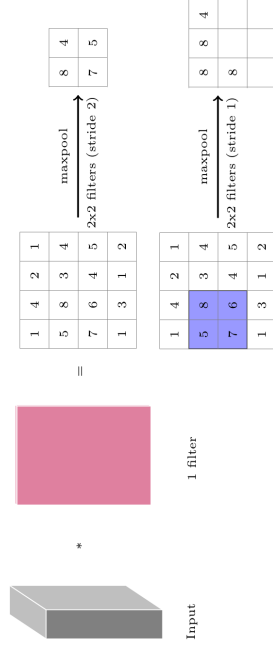
CNN has sparse connectivity with respect to NN



Interactions are preserved, even with reduced model parameters,



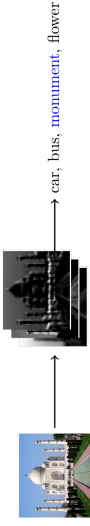
Pooling (max, min, average)



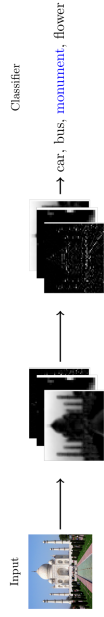
Training CNN? backpropagation!

Automate feature kernel discovery

Instead of handcrafted kernels, learn filters



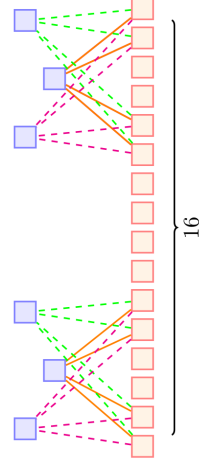
(Why not multiple?) learn multiple layers of kernels/filters



Treating these kernels as parameters and learning them.

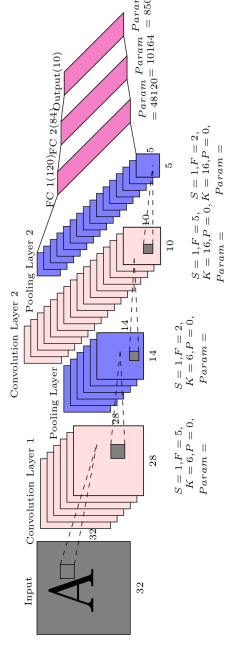
Weight sharing in CNN

One place you are extracting edge other place something else? So we do not want the kernel to be different for different portions of the image.



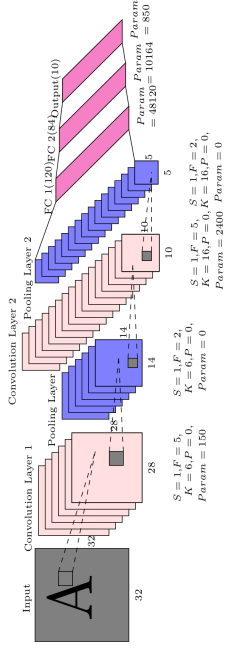
- Weight sharing in CNN makes the job of learning weights easier
- Multiple kernels help get different feature at the same level

Case-Study: LeNet-5 for handwritten character recognition¹



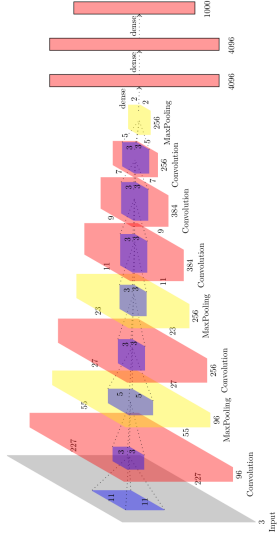
¹Yann Lecun and Leon Bottou and Yoshua Bengio and Patrick Haffner, Gradient-based learning applied to document recognition, pp 2278-2324, IEEE-1998

LeNet-5 for handwritten character recognition ²



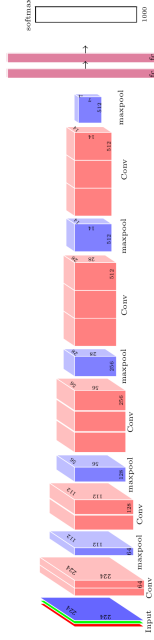
² Yann Lecun and Leon Bottou and Yoshua Bengio and Patrick Haffner, *Gradient-based learning applied to document recognition*, pp 2278–2324, IEEE: 1998

AlexNet 4



⁴ Clje: 50086, *Imagenet classification with deep convolutional neural networks*, Krizhevsky, Alex and Sutskever, Ilya and Hinton, Geoffrey E. In: Advances in neural information processing systems pages: 1097–1105, NIPS:2012 → [MaxPool F=3, S=2] → FC:4096 → Softmax [8944 F=3, S=1, P=0] → [384 F=3, S=1, P=0] → [638 F=3, S=1, P=0] → [MaxPool F=3, S=2] → FC:4096 →

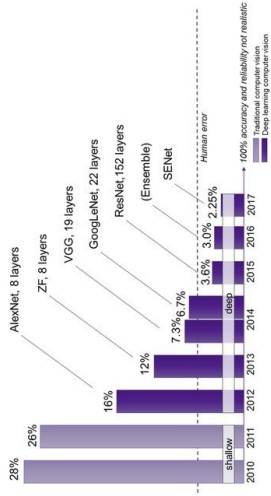
VGG16 6



- ➊ Kernel size is always 3×3
- ➋ 16M parameters in pre-FC and 122 in FC. First FC layer is huge
- ➌ Layers represents abstract representation and can be reused (FC or Conv)

⁶ Clje 29548 *Very deep convolutional networks for large-scale image recognition*, Simonyan, Karen and Zisserman, Andrew, pages 818–833, arXiv preprint arXiv:1409.1556 - In: ICLR-2015

ImageNet ILSVRC³



- (2009) 22K category, 14M images
- Challenge 1000 class, 1431167 images
- HoG, LBP, SVM ...

³ imagenet large scale visual recognition challenge <http://www.image-net.org/challenges/LSVRC/>

ZFNet 5

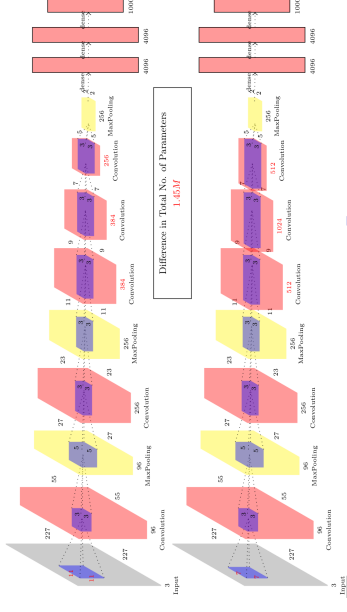
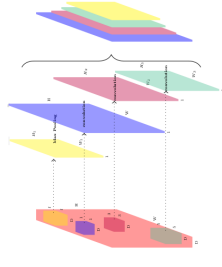


Image size $227 \times 227 \times 3 \rightarrow [69 F=7, S=4, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow [256 F=5, S=1, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow [512 F=3, S=1, P=0] \rightarrow [1024 F=3, S=1, P=0] \rightarrow [612 F=3, S=1, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow FC:4096 \rightarrow FC:1000$

⁵ Clje 7560 *Visualizing and understanding convolutional networks*, Zeiler, Matthew D and Fergus, Rob, pages 818–833, European conference on computer vision (ECCV) Springer-2014

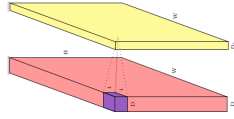
GoogLeNet 7

- Recall scale invariance in SIFT
- Multiple filters of different size is a good idea
- With $W \times H \times D$ input and $F \times F \times D$ filter and $S = 1$ and no padding, output is of size $(W - F + 1) \times (H - F + 1)$
- Each value needs $F \times F \times D$ computation



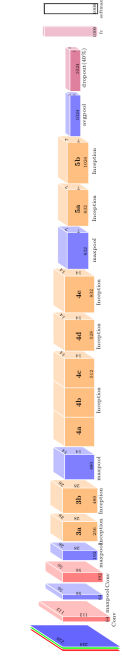
Can we reduce this computation a bit?
Idea is to have 1×1 computation

1 x 1 convolution



- 1×1 is $1 \times 1 \times D$
- They produce one output plane
- By using D_1 such 1×1 convolution output becomes $F \times F \times D_1$
- We have $D_1 < D$

GoogLeNet



- Input is RGB 229×229

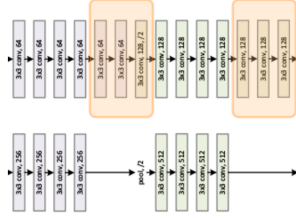
- Each inception module have very specific configuration.

- (3a) $192 \times 28 \times 28$ 64 96 126 16 32 32
- (3b) $256 \times 28 \times 28$ 28 128 192 32 96 61
- (4a) $48 \times 14 \times 14$ 192 96 208 16 48 96
- (4b) $512 \times 14 \times 14$ 160 112 224 24 64 64
- (5a) $512 \times 14 \times 14$ 128 128 256 24 64 64
- (5b) $512 \times 14 \times 14$ 128 128 256 24 64 64
- (6a) $528 \times 14 \times 14$ 256 160 320 32 128 128
- (6b) $832 \times 7 \times 7$ 256 160 320 32 128 128
- (7) $832 \times 7 \times 7$ 384 192 384 48 124 128



ResNet⁸

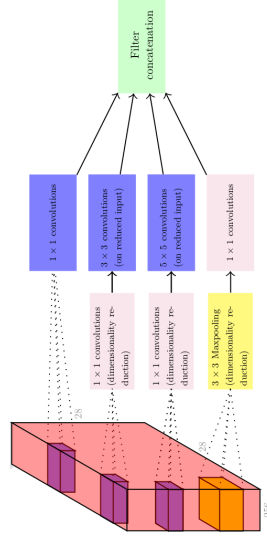
If a shallow neural network works well. What would happen if we add more layers?



- Deep network should also work well (It would learn identity in new layers)

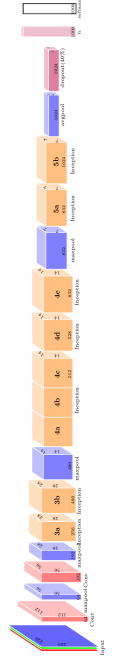
⁸ Cite: 32871, **Deep residual learning for image recognition**, He, Kaiming and Zhang, Xiangyu and Ren, Shaoqing and Sun, Jian, in: IEEE conference on computer vision and pattern recognition, pages 770-778, CVPR-2016

Inception Block: Multiple convolutions



- 1×1 convolution
- 1×1 convolution followed by 3×3
- 1×1 convolution followed by 5×5
- 3×3 maxpool followed by 1×1
- Appropriate padding is done to make things of same size

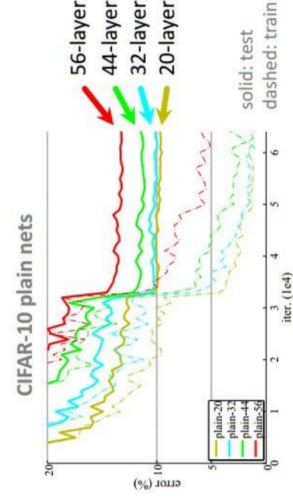
GoogLeNet



- VGGNET has $512 \times 7 \times 7$ size at pre-FC this was an issue to connect with 4096
- GoogLeNet applies a average pool. Gives 49 time reduction. has 1024 values only
- Dropout and connect to 1000
- 12 times less connections as compared to AlexNet
- 2 times more computation as compared to AlexNet
- Very high accuracy. Error reduced from 16% -to- 6.7%

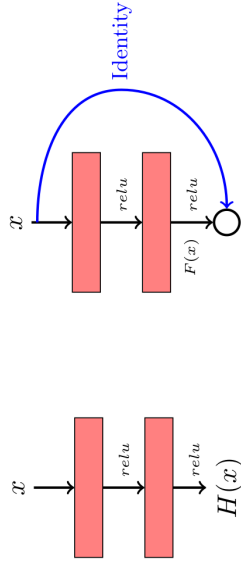
ResNet

But, in practice it was not happening



Why? Identity is one of the solution in large domain.

Let me tell this to the network

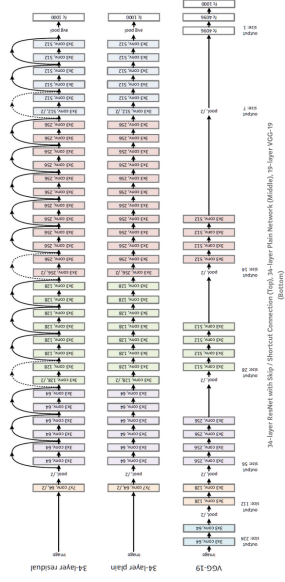


$$H(x) = F(x) + x$$

ResNet Hyper-parameters and Issues

- Training takes huge time
- Batch Normalization
- Xavier/2 initialization
- SGD and momentum
- Small learning rate 0.1
- Mini-batch size 256
- Weight decay
- No Dropout

ResNet Comparison ⁹



152-layer deep net was better than human. Only 3.6% error rate
ImageNet Classification ^a

^a Better than the 2nd best system
ImageNet Detection: 16% ImageNet Localization: 27% COCO Detection: 11% COCO Segmentation: 1.2%

⁹ ResNet, Winner of ILSVRC 2015 (Image Classification, Localization, Detection) Sik-Ho Tsang

Thank You!

Thank you very much for your attention!