



CS F425: Deep Learning

19

CNN Le-Net, ResNet DenseNet, U-Net



Dr. Kamlesh Tiwari
Assistant Professor, Department of CSIS,
BITS Pilani, Pilani Campus, Rajasthan-333031 INDIA
Mar 03, 2023 **ON-CAMPUS** Campus @ BITS-Pilani [Jan-May 2023]

<http://ktiwari.in/dl>

AlexNet ²

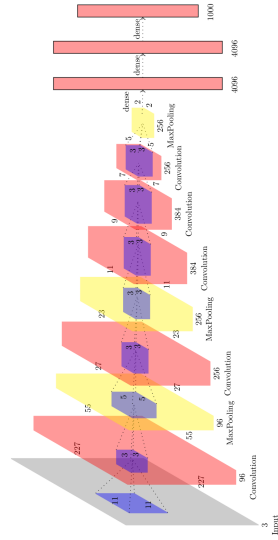
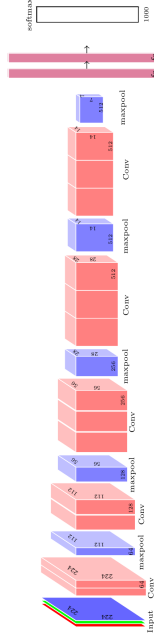


Image size $227 \times 227 \times 3 \rightarrow [96 F=11, S=4, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow [256 F=5, S=1, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow [384 F=3, S=1, P=0] \rightarrow [384 F=3, S=1, P=0] \rightarrow [256 F=3, S=1, P=0] \rightarrow [256 F=3, S=1, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow FC 4096 \rightarrow FC 1000$

² Cite: 128212, **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

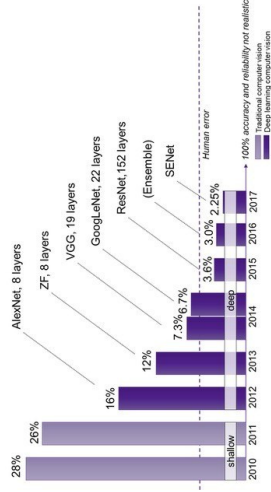
VGG16 ⁴



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

⁴ Cite 96789 **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 ILSVRC¹



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

¹ Imagenet large scale visual recognition challenge <http://www.image-net.org/challenges/LSVRC/>

ZFNet ³

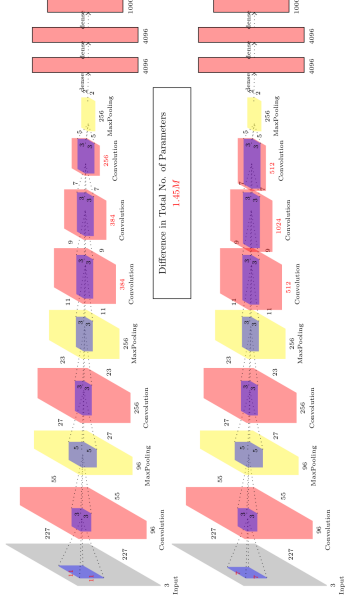
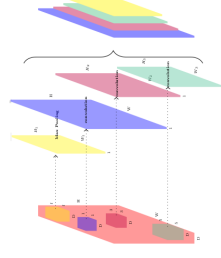


Image size $227 \times 227 \times 3 \rightarrow [69 F=7, S=4, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow [256 F=5, S=1, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow [512 F=3, S=1, P=0] \rightarrow [1024 F=3, S=1, P=0] \rightarrow [612 F=3, S=1, P=0] \rightarrow [MaxPool F=3, S=2] \rightarrow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

³ Cite: 18566 **Visualizing and understanding convolutional networks**, Zeiler, Matthew D and Fergus, Rob, pages 818–833, European conference on computer vision (ECCV) Springer-2014

GoogLeNet ⁵

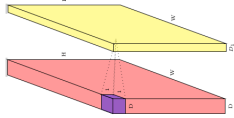
- Recall scale invariance in SIFT
- Multiple filters of different size is a good idea
- With $W \times H \times D$ input and $F \times F \times D$ filter and $S = 1$ and no padding, output is of size $(W - F + 1) \times (H - F + 1)$
- Each value needs $F \times F \times D$ computation



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

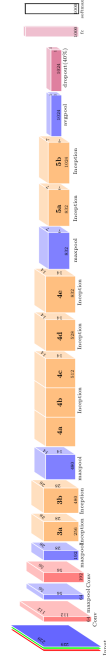
⁵ Cite: 47942, **Going deeper with convolutions**, Szegedy, Christian and Liu, Wei and Jia, Yangqing and Sermanet, Pierre and Reed, Scott and Anguelov, Dragomir and Erhan, Dumitru and Vanhoucke, Vincent and Rabinovich, Andrew. In: Proceedings of the IEEE conference on computer vision and pattern recognition pages 1–9, CVPR-2015

1 × 1 convolution



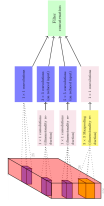
- 1×1 is $1 \times 1 \times D$
- They produce one output plane
- By using D_1 such 1×1 convolution output becomes $F \times F \times D_1$
- We have $D_1 < D$

GoogLeNet



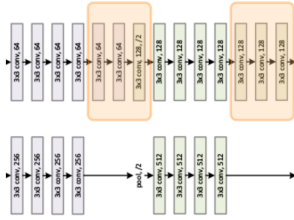
- Input is RGB $224 \times 224 \times 3$
- Each inception module have very specific configuration.

- (3a) $192 \times 28 \times 28$ 64 96 128 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
- (4c) $512 \times 14 \times 14$ 128 128 256 24 64 64
- (4d) $512 \times 14 \times 14$ 128 128 256 24 64 64
- (5a) $512 \times 14 \times 14$ 256 140 224 32 128 128
- (5b) $832 \times 7 \times 7$ 256 160 320 32 128 128
- (5c) $832 \times 7 \times 7$ 384 192 384 48 124 128



ResNet ⁶

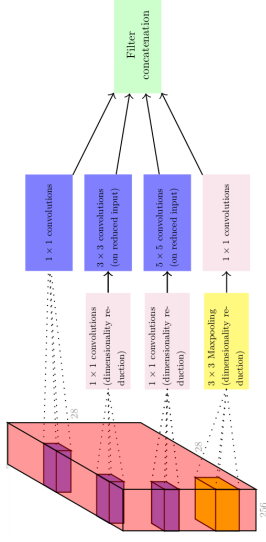
If a shallow neural network works well. What would happen if we add more layers?



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

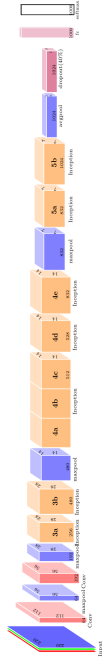
⁶ Cite: 32871, **Deep residual learning for image recognition**, He, Kaiming and Zhang, Xiangyu and Ren, Shaoqing and Sun, Jian, in: IEEE conference on computer vision and pattern recognition, pages 770-776, CVPR-2016

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

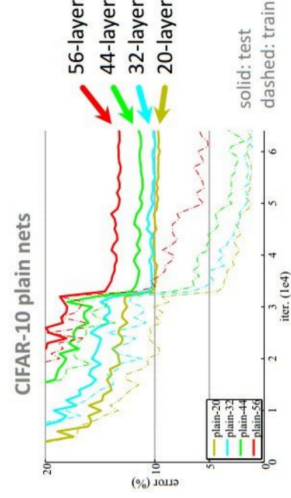
GoogLeNet



- 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%

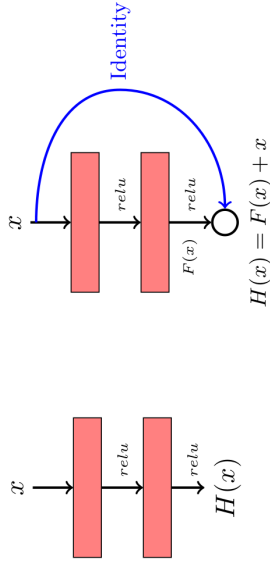
ResNet

But, in practice it was not happening



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

Let me tell this to the network

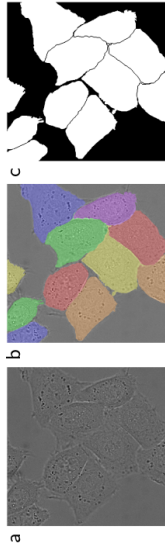


$$H(x) = F(x) + x$$

ResNet Hyper-parameters and Issues

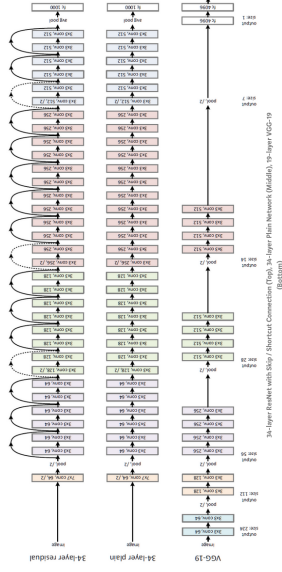
- Training takes huge time
- Batch Normalization
- Xavier/2 initialization
- SGD and momentum
- Small learning rate 0.1
- Mini-batch size 256
- Weight decay
- No Dropout

Segmentation



- ISBI challenge for segmentation of neuronal structures in electron microscopic stacks
- Works with very few training images (30/application) and touching boundary. Yield more precise segmentation
- Data augmentation is essential (mainly shift, rotation and elastic deformation)

ResNet Comparison 7

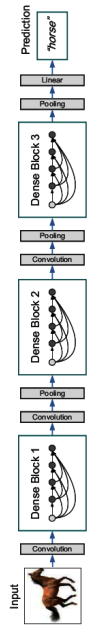
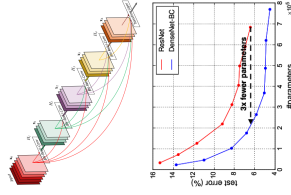


152-layer deep net was better than human. Only 3.6% error rate
ImageNet Classification ^a

^a Better than the 2nd best system
ImageNet Detection: 16% ImageNet Localization: 27% ImageNet Classification: 11% COCO Segmentation: 12%

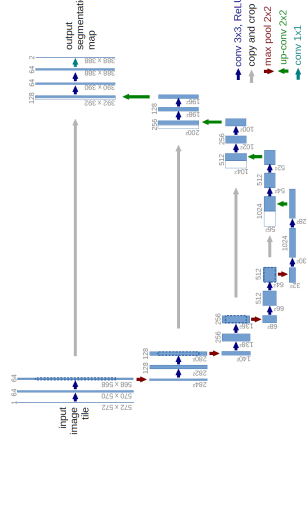
DenseNet 8

- Connect each layer to every other layer in feed-forward fashion
- It alleviates vanishing-gradient problem, strengthen feature propagation
- With feature reuse, substantially reduce number of parameters
- CIFAR10-3.6%, CIFAR100-19.6%,



⁸ Cite 32471, Huang, Gao and Liu, Zhuang and Van Der Maaten, Laurens and Weinberger, Kilian Q, **Densely connected convolutional networks**. In IEEE CVPR vision and pattern recognition, pp 4700–4708, 2017

U-Net⁹



- ISBI DIC-HeLa achieved 77.6% iou as compared to 46.0% second
- ISBI Cell tracking 2015, achieved 92% IoU as compared to 83%

second

⁹ Cite 57528 O. Ronneberger and P.Fischer, and T. Brox, *U-Net: Convolutional Networks for Biomedical Image Segmentation*, Medical Image Computing and Computer-Assisted Intervention (MICCAI), LNCS-8351, pages 234–241, Springer-2015

Thank You!

Thank you very much for your attention!