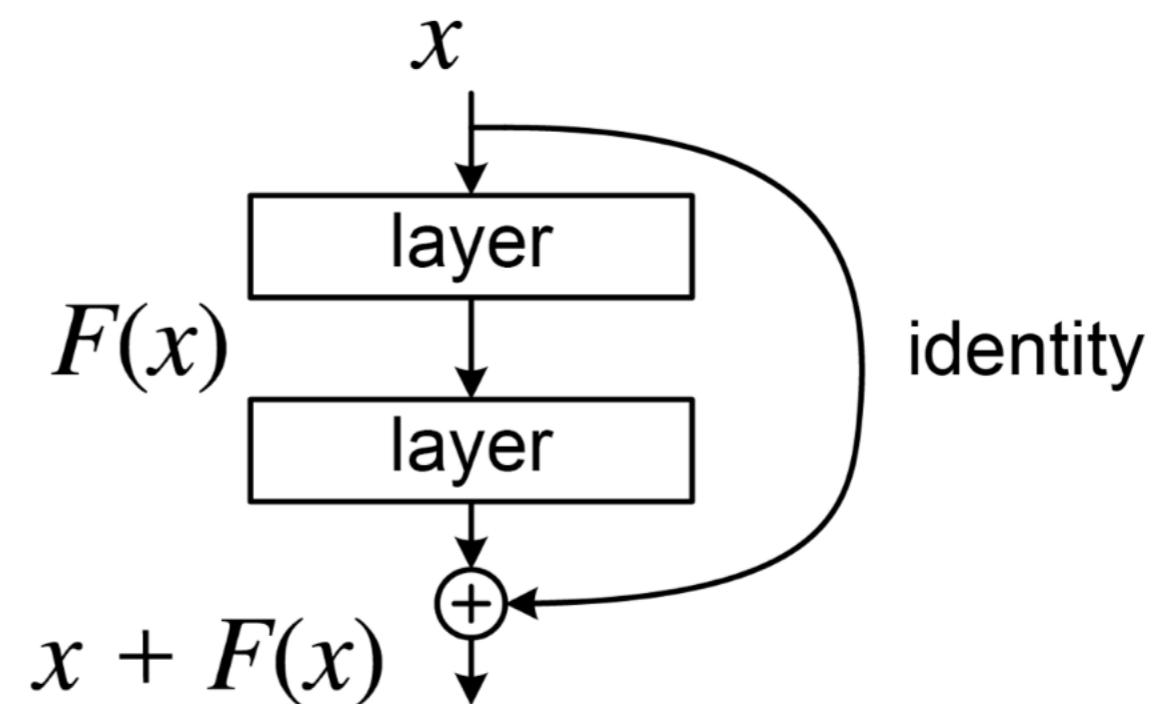


Day - 20, Oct - 22, 2024

CNN continuing — Resnet

Residual Neural Networks (Resnet)

introduced novel concept
('skip connections'). It allows
information to flow directly
from earlier layers to later ones,



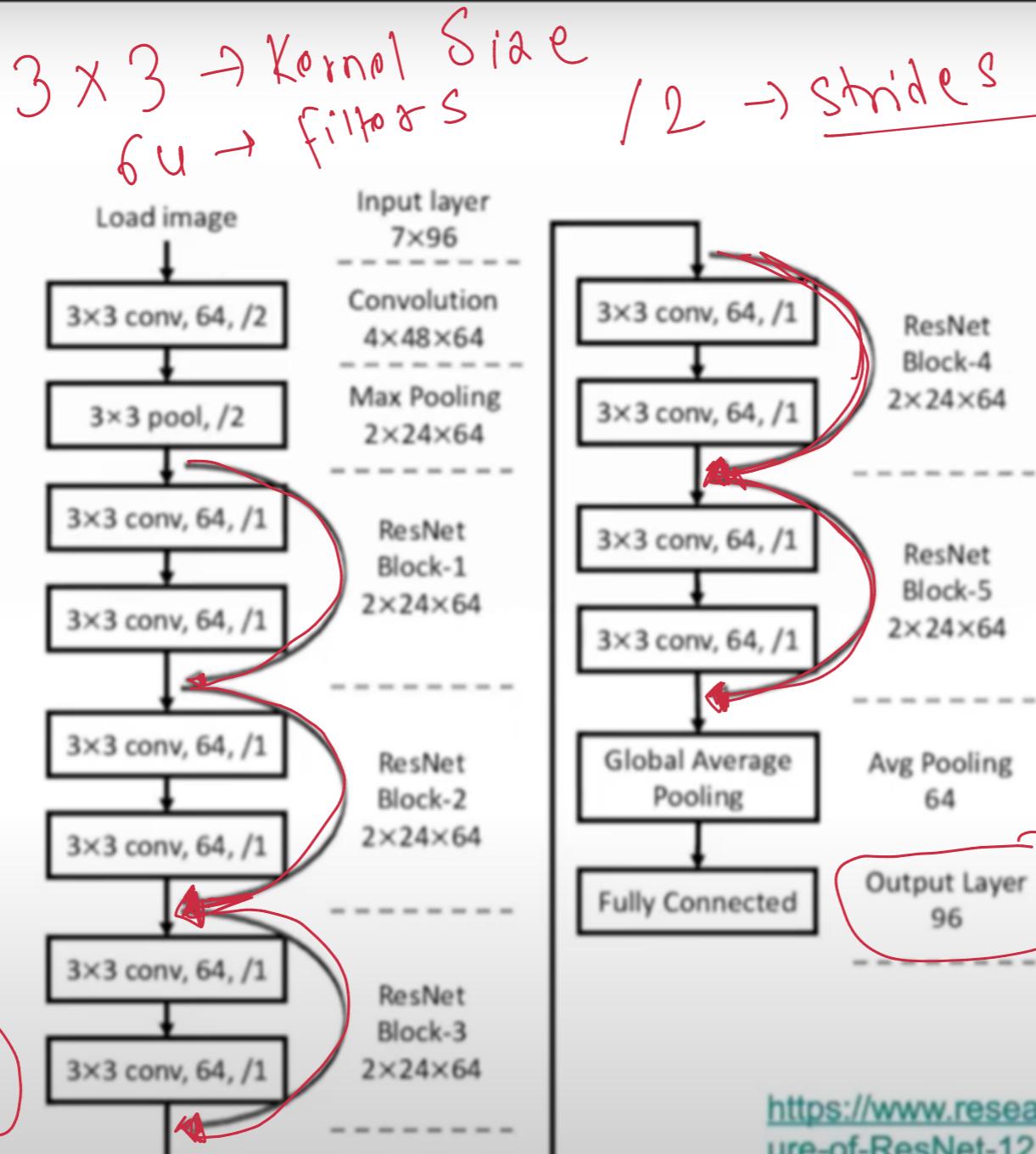
by passing intermediate layers. Direct path, Residual Blocks, Identity function,
and Solves Vanishing gradient problems.

Resnet

?

skip connections -> convolution

$4 \times 48 \times 64$
 $4 \rightarrow \text{rows}$
 $48 \rightarrow \text{columns}$
 $64 \rightarrow \text{filters}$
 $\rightarrow 64 \text{ feature Map}$



$3 \times 3 \rightarrow \text{Kernel Size}$
 $64 \rightarrow \text{filters}$
 $/2 \rightarrow \text{strides}$

Identity function
As Shortcut.)



https://www.researchgate.net/figure/The-structure-of-ResNet-12_fig1_329954455

HP 7×96
OIP -
 \rightarrow OIP - 96 classification

Identity function — input x to the function $f(x) \rightarrow x$ 'output'

Resnet

→ Skip Connections, layers (shortcut-connections)

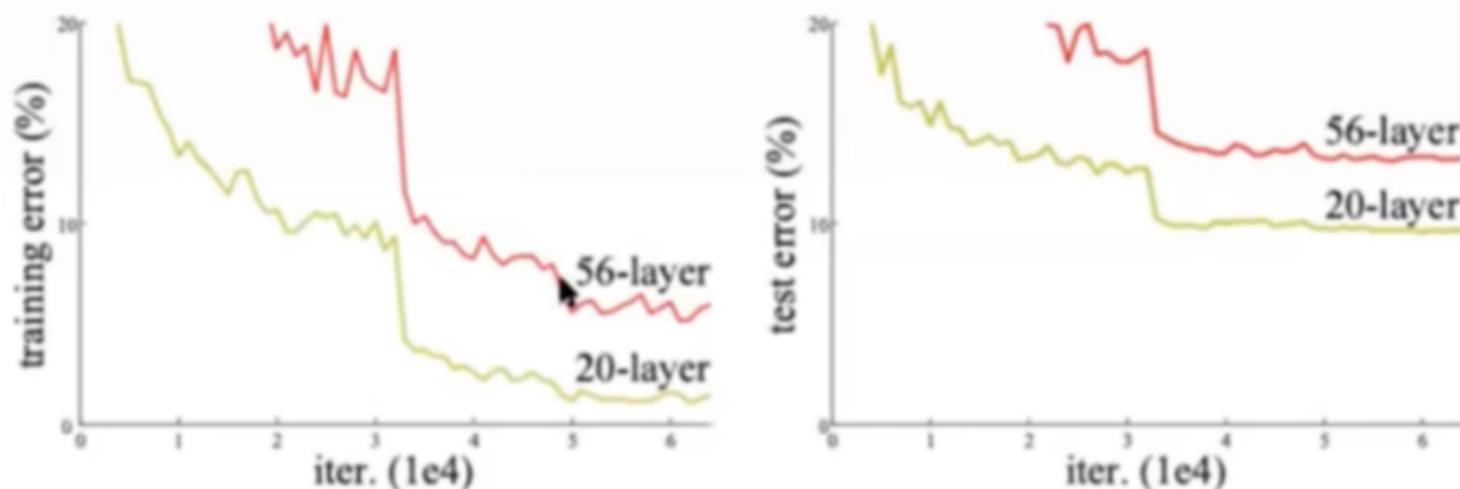


Figure 1. Training error (left) and test error (right) on CIFAR-10 with 20-layer and 56-layer “plain” networks. The deeper network has higher training error, and thus test error. Similar phenomena on ImageNet is presented in Fig. 4.

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



THIRD NEPAL WINTER SCHOOL AI, 2021

NepAI Applied Mathematics and
Informatics Institute for research



Nripesh Parajuli

nepalschool.naamii.org.np

As we add more layers to the deep Neural Networks then the problems of the ANN becomes unstabilize so we can use Resnet.

Figure 1 (18-layers No Resnet) (34-layers Resnet)

Fig 1.

Resnet-18 (18 layers)

Resnet

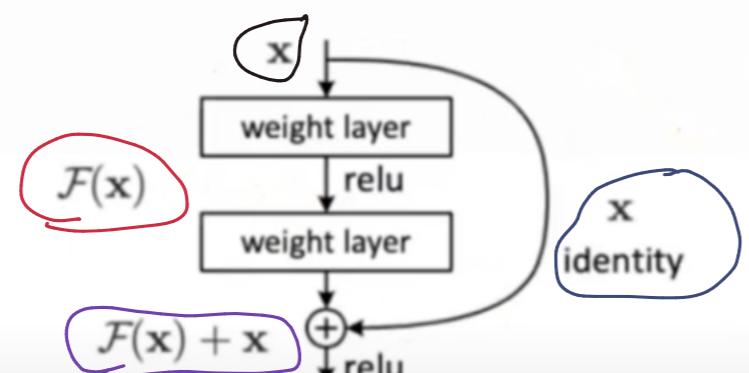


Figure 2. Residual learning: a building block.

- Assuming that the original desired mapping is $H(x)$:
 - Before: fit $F(x) = H(x) - x$
 - Resnet: fit $H(x) = F(x) + x$
- The hypothesis was that it's easier to fit the residual mapping than the original.
- In extreme case, easier to map to 0 than map to 1.



Fig -2 .

Resnet-34

perform better than

Resnet-18

due to
Resnet +
Applied .

Red \rightarrow 34

Blue \rightarrow 18

Fig 1 .

Resnet-34 has higher error than Resnet-18

plots Comparing the training and validation error curves for different ResNet Architecture -

Resnet

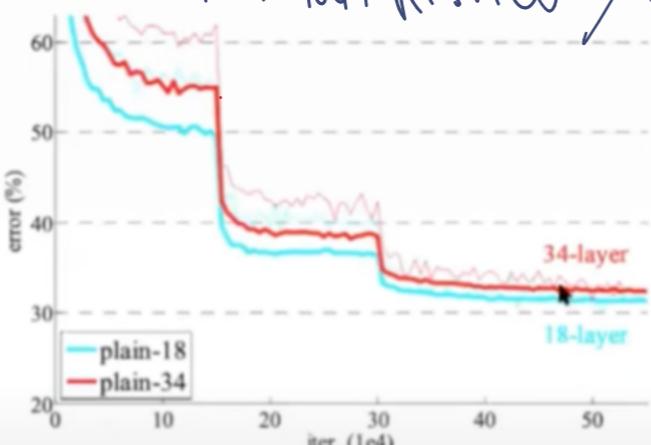


Fig 1
without Resnet X

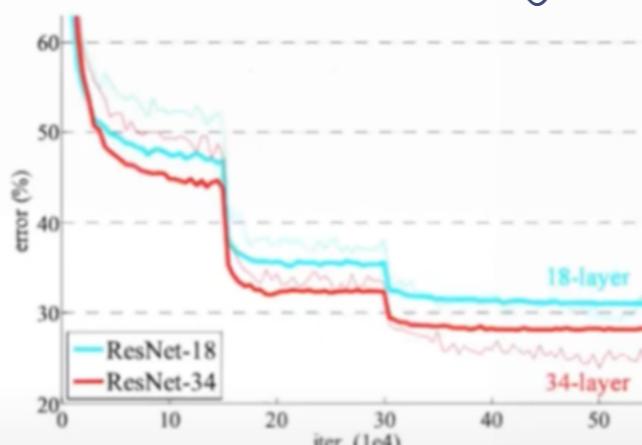


fig 2 .
with Resnet ✓

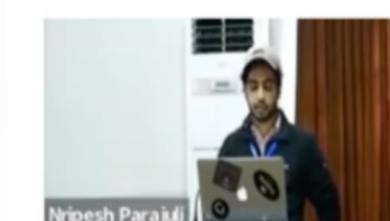


Figure 4. Training on ImageNet. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

GoogleNet

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

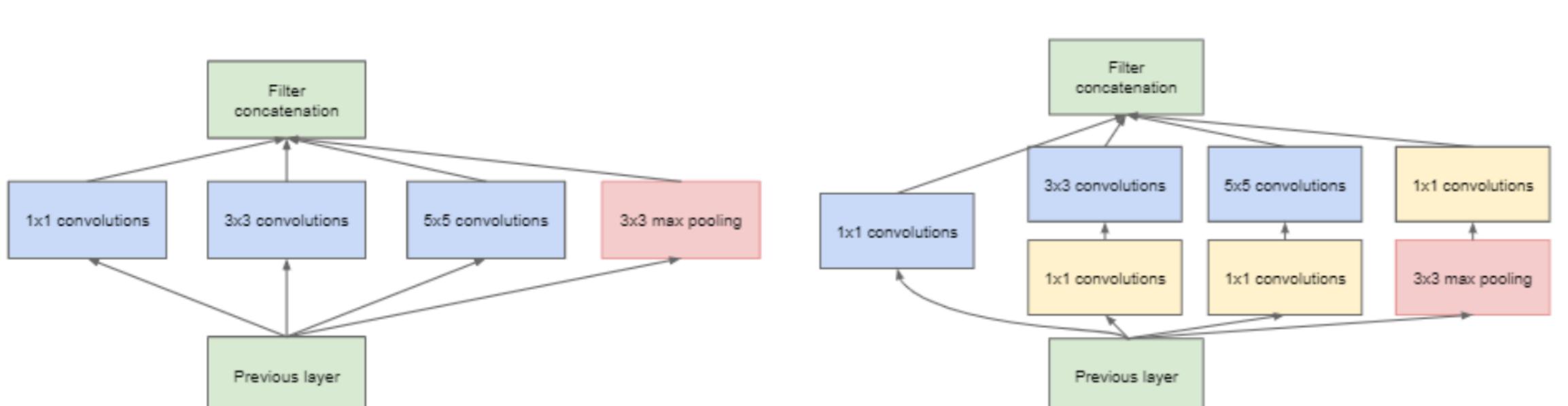


- Does dimensionality reduction on each layer and does convolutions at different resolutions to build sparse but more diverse/representative features.
- “visual information should be processed at various scales and then aggregated so that the next stage can abstract features from different scales simultaneously”



GoogleNet (Inception Module) allows the network to choose between several

Convolutional filter sizes in each layer that is beneficial for capturing both local and global features efficiently.



(a) Inception module, naïve version

(b) Inception module with dimension reductions

[1x1 is multiplied by constant]

#Auxiliary Classifier:

- Tackle gradient Vanishing problem
- provide additional gradient flow to the previous layers
- Incorporated at the intermediate layers.



Figure 3: GoogLeNet network with all the bells and whistles

GoogleNet

- Uses 7x7 convolutions after the input - which is rather unusual.
- Auxiliary classifiers:
 - The model had two additional outputs, trying to predict the same thing as the primary output.
 - These are added in the middle part of the network - in order to increase the discriminative power of the middle layers.
 - The losses are added to overall loss, but multiplied by 0.3
 - This helped in training the model better and not having bottlenecks in the middle layers.



GoogleNet on ImageNet

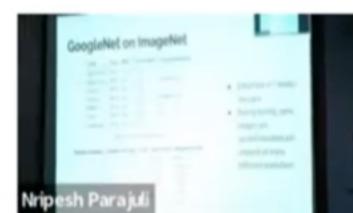
Team	Year	Place	Error (top-5)	Uses external data
SuperVision	2012	1st	16.4%	no
SuperVision	2012	1st	15.3%	Imagenet 22k
Clarifai	2013	1st	11.7%	no
Clarifai	2013	1st	11.2%	Imagenet 22k
MSRA	2014	3rd	7.35%	no
VGG	2014	2nd	7.32%	no
GoogLeNet	2014	1st	6.67%	no

Table 2: Classification performance

Number of models	Number of Crops	Cost	Top-5 error	compared to base
1	1	1	10.07%	base
1	10	10	9.15%	-0.92%
1	144	144	7.89%	-2.18%
7	1	7	8.09%	-1.98%
7	10	70	7.62%	-2.45%
7	144	1008	6.67%	-3.45%

Table 3: GoogLeNet classification performance break down

- Ensemble of 7 models are used.
- During testing, same images are resized/rescaled and cropped at many different resolutions.

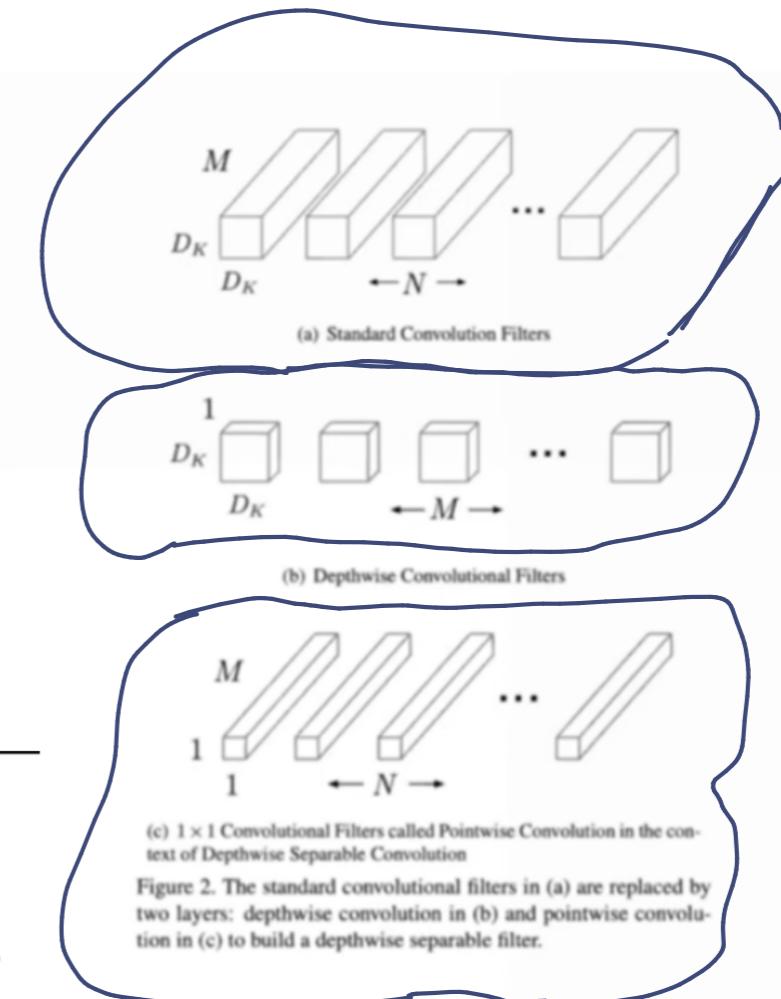


Mobilenet

more efficient
light weight
for mobile
devices'

Mobilenet

- Uses separable depthwise convolution where regular convolution - it is separated into depthwise convolution (3×3), followed by (1×1) convolution.
- This reduces number of parameters by about 8-9 times.



Mobilenet

Table 8. MobileNet Comparison to Popular Models

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
GoogleNet	69.8%	1550	6.8
VGG 16	71.5%	15300	138

Table 9. Smaller MobileNet Comparison to Popular Models

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
0.50 MobileNet-160	60.2%	76	1.32
SqueezeNet	57.5%	1700	1.25
AlexNet	57.2%	720	60

- For comparison, resnet-152 had 79% accuracy (top-1 error). But is so much bigger and slower.



Mobilenet

- The model allows for reducing the size / latency further by using less channels (controlled by alpha), and by reducing input image resolution (controlled by α_0).
- Small reduction in number of channels and resolution, still gives good performance.

→ MobileNet are

Especially made for mobile and embedded visual applications.

→ Resnet can handle deep neural network solving vanishing gradient problem, skip connectionism.

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
Conv MobileNet	71.7%	4866	29.3
MobileNet	70.6%	569	4.2

Table 5. Narrow vs Shallow MobileNet

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
0.75 MobileNet	68.4%	325	2.6
Shallow MobileNet	65.3%	307	2.9

Table 6. MobileNet Width Multiplier

Width Multiplier	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
0.75 MobileNet-224	68.4%	325	2.6
0.5 MobileNet-224	63.7%	149	1.3
0.25 MobileNet-224	50.6%	41	0.5

Table 7. MobileNet Resolution

Resolution	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
1.0 MobileNet-192	69.1%	418	4.2
1.0 MobileNet-160	67.2%	290	4.2
1.0 MobileNet-128	64.4%	186	4.2



MobileNet

→ For mobile, embedded visual system & less computational devices

→ Depthwise Separable Convolutions, faster

(~ 4.2 million)

→ May sacrifice too much accuracy.

Resnet

→ Residual Blocks
(Skip Connections)

→ Can be very deep
Autote models for
classification but costly

(~ 26 million)

→ Deep learning visual task

GoogleNet

→ Inception Modules
($1 \times 1, 3 \times 3, 5 \times 5$)

→ Features various scales with fewer parameters (~ 5 million)

→ General vision task

From Textbook: # Knowledge Representation: Unification

Unification: Unification is a process in FOL that aims to make two logical expression identical by finding a substitution (set of variable replacements).

Example:
→ Expressions: $p(x, a)$ and $p(b, y)$
→ Unification: By substituting $x=b$ & $y=a$ these two expressions become identical
 $p(b, a)$ and $p(b, a)$ ✓ Substituting.

Lifting: Lifting refers to extending the inference process from propositional logic to first-order-logic; variables are involved, predicates involved.

It is the process of transforming a ground clause (a clause without variables) into the general clause (a clause with variables).

Generalizing Resolution, Avoiding Grounded Inference, Offfiling Inference to FOL
Learning Utzay Challenge Oct 3 - Nov 3 2024 12 of

and Role of unification in lifting are some points to be considered.

Example: Suppose you have two clauses:

$p(x) \rightarrow$ where x is variables

$\neg p(a) \rightarrow$ where a is constant

Before applying Inference (resolution), we must "lift" the process by Unifying $x=a$, making the clauses compatible.

(1) $p(x) \vee q(x)$

(2) $\neg p(a)$

(3) $p(a) \vee p(a)$

① Initial clauses:

① clause 1: $p(x) \vee q(x)$

② clause 2: $\neg p(a)$

② Unification: $x = a$

clause 1: $p(a) \vee q(a)$

clause 2: $\neg p(a)$

③ Applying Resolution: Clause 1: $p(a) \vee q(a)$
Clause 2: $\neg p(a)$

Resolution step:

$$(P(a) \vee Q(a))$$

Complementary literals:

$P(a)$ $Q(a)$

Result $\rightarrow Q(a)$. → holds true -

Machine vision

- Application-specific, real-time processing, decision-making, automation
- Heavy reliance on specialized hardware

Computer vision

- General-purpose, understanding & interpreting visual data
- ML, DL Algorithms, Object recognition