

Video classification with keras

딥러닝스터디
2020-10-21
이예림

Video classification vs. Image classification

Video? A series of individual images!

Video classification = N images classification
(N = # of frames)

Video classification vs. Image classification

Video? A series of individual images!

~~Video classification = N images classification
(N = # of frames)~~

Video classification vs. Image classification

Video? A series of individual images!

Video classification = N images classification
(N = # of frames)

Subsequent frames are correlated with *semantic contents*.

> Temporal dimension > video classification results improval
(i.e., LSTM, RNN ...)

Video classification vs. Image classification

Video? A series of individual images!

Video classification = N images classification
(N = # of frames)

Subsequent frames are correlated with *semantic contents*.

> Temporal dimension > video classification results improval
(i.e., LSTM, RNN ...)

Rolling prediction averaging!

Image classification (CNN) > Video classification (Rolling averaging)

Video classification vs. Image classification

Image classification

1. Input : Image
2. Prediction
3. The largest corresponding probability >> The label

Video classification

1. Input : Frame
2. Loop over all frames
3. Prediction *individually* and *independently*
4. The largest corresponding probability >> The label

“Prediction flickering”

Video classification vs. Image classification

Image classification

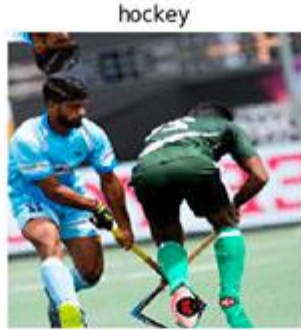
1. Input : Image
2. Prediction
3. The largest corresponding probability >> The label

Video classification

1. Input : Frame
2. Loop over all frames
3. Prediction *individually* and *independently*
4. A list of the last K predictions
5. The average of the last K predictions
6. The largest corresponding probability >> The label

"Smoothing out"

The sports classification dataset



- | | |
|--------------------|-------------------|
| 1.Swimming | 12.Gymnasium |
| 2.Badminton | 13.Weight lifting |
| 3.Wrestling | 14.Volleyball |
| 4.Olympic Shooting | 15.Table tennis |
| 5.Cricket | 16.Baseball |
| 6.Football | 17.Formula 1 |
| 7.Tennis | 18.Moto GP |
| 8.Hockey | 19.Chess |
| 9.Ice Hockey | 20.Boxing |
| 10.Kabaddi | 21.Fencing |
| 11.WWE | 22.Basketba |

The sports classification dataset



- | | |
|--------------------|-------------------|
| 1.Swimming | 12.Gymnasium |
| 2.Badminton | 13.Weight lifting |
| 3.Wrestling | 14.Volleyball |
| 4.Olympic Shooting | 15.Table tennis |
| 5.Cricket | 16.Baseball |
| 6.Football | 17.Formula 1 |
| 7.Tennis | 18.Moto GP |
| 8.Hockey | 19.Chess |
| 9.Ice Hockey | 20.Boxing |
| 10.Kabaddi | 21.Fencing |
| 11.WWE | 22.Basketba |

Transfer learning

1. Take a network pre-trained on a dataset.
2. Utilize the network to recognize image/object categories it was not trained on.

Advantages?

- Train a network with a new dataset > Cost + time loss, Not working
- Hundred of parameters

Two steps of transfer learning

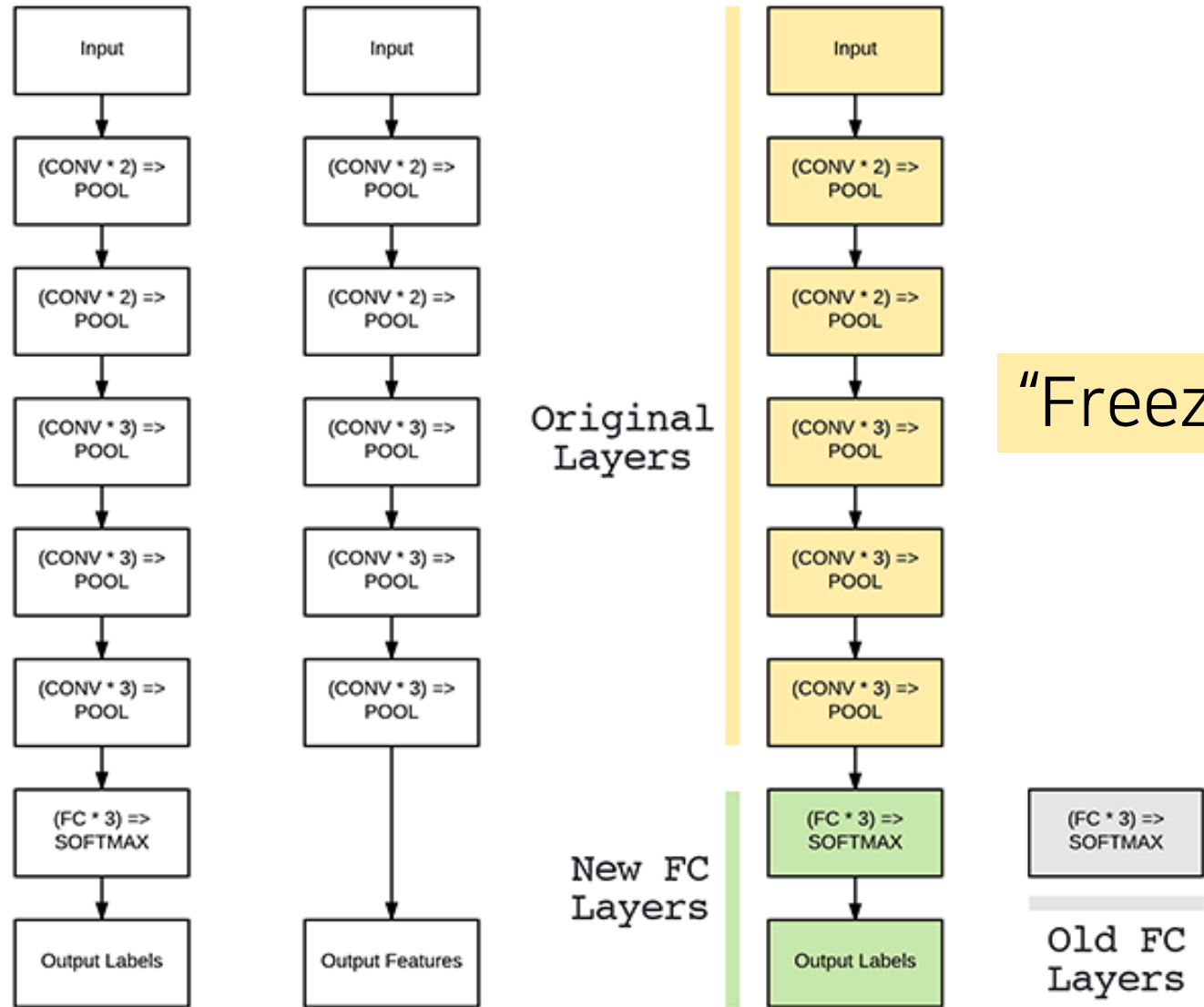
1. Via feature extraction
 - Pre-trained network = Feature extractor
2. Via fine-tuning

Q. 무조건 좋을까?

Feature extraction이 어떻게 되는 것에 따라 다르다. 데이터의 분포(특성)이 다를 경우

Transfer learning

VGG16

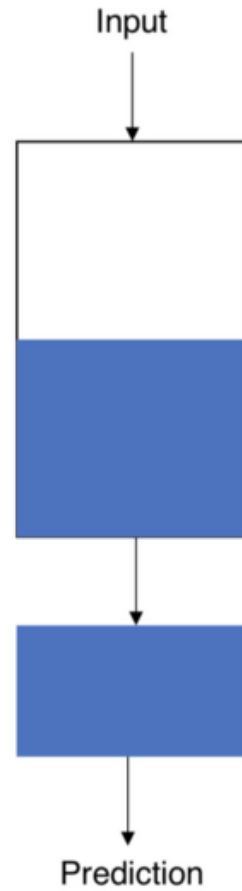


Transfer learning

Strategy 1
Train the
entire model



Strategy 2
Train some layers and
leave the others frozen



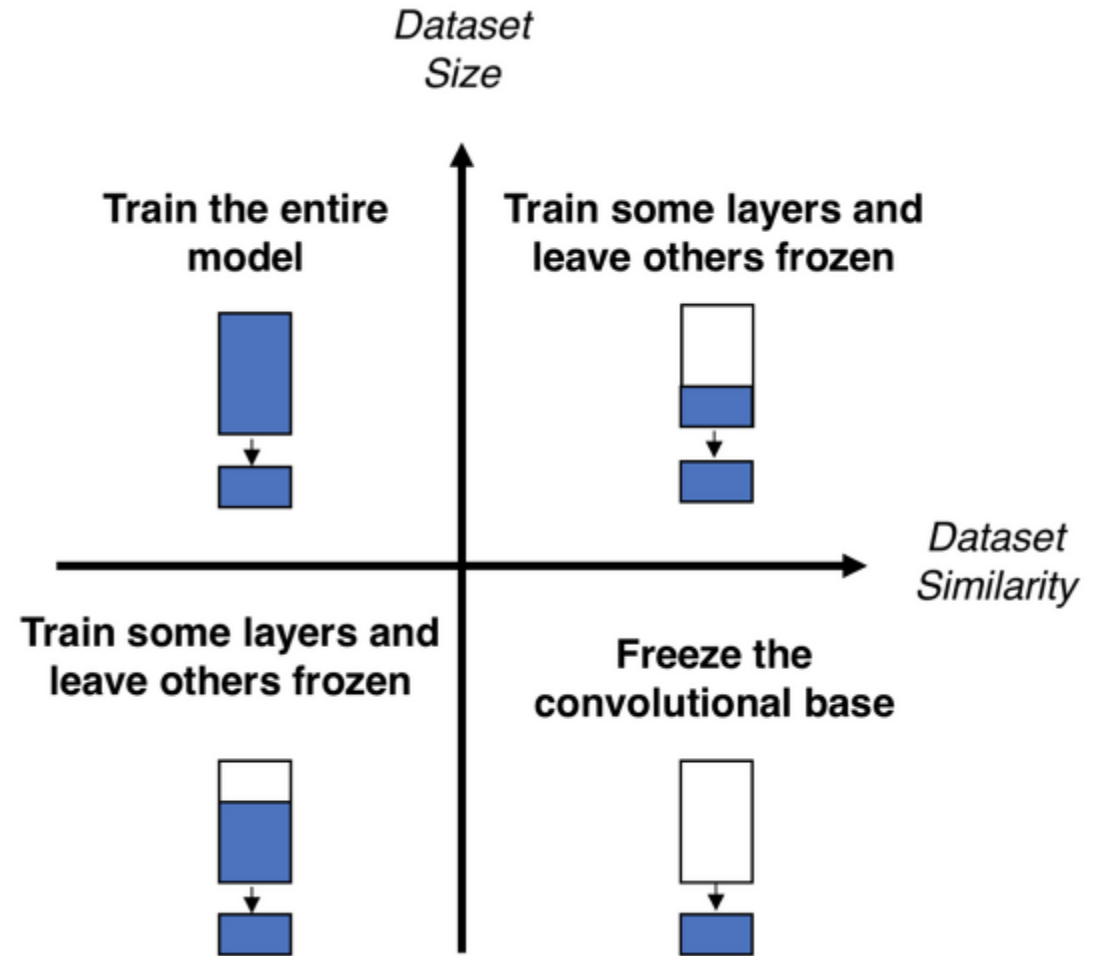
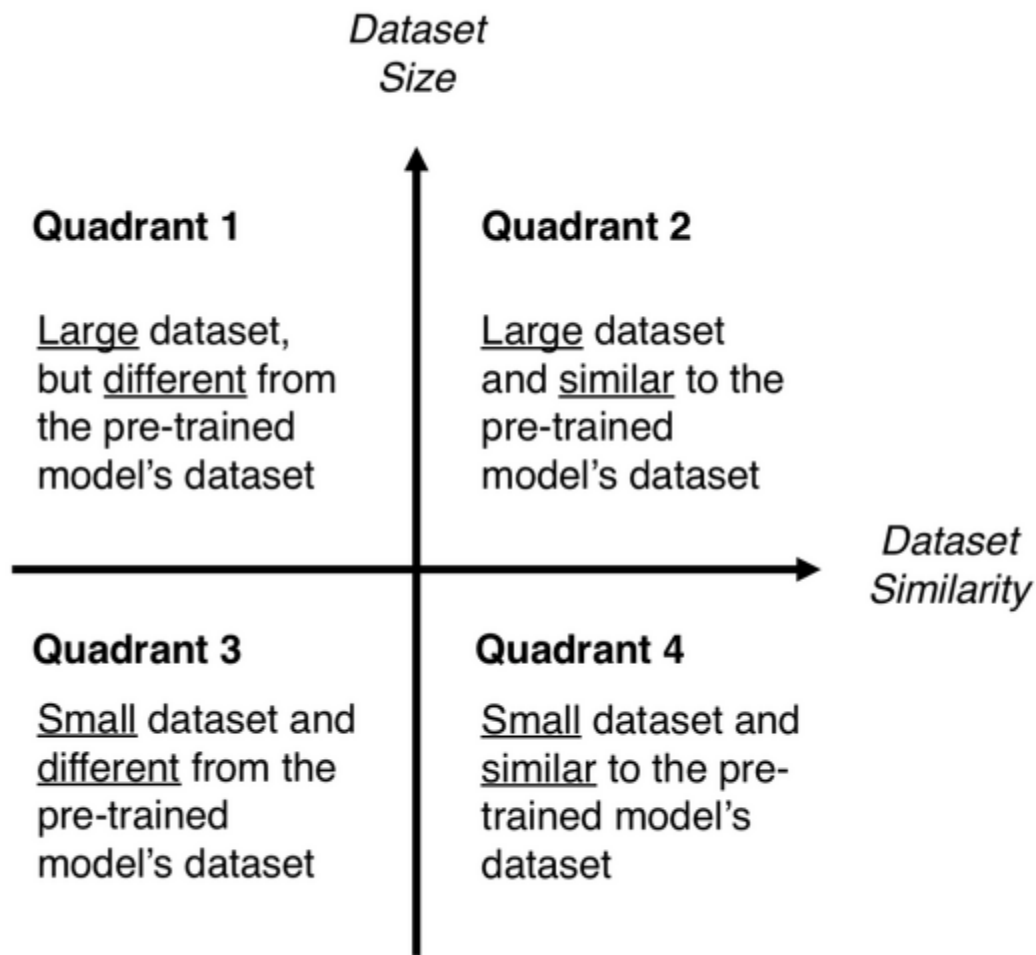
Strategy 3
Freeze the
convolutional base



Legend:



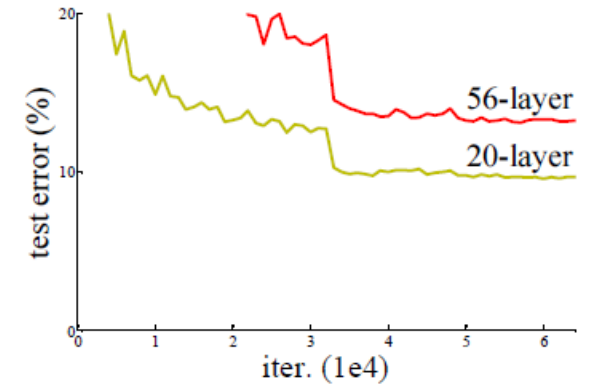
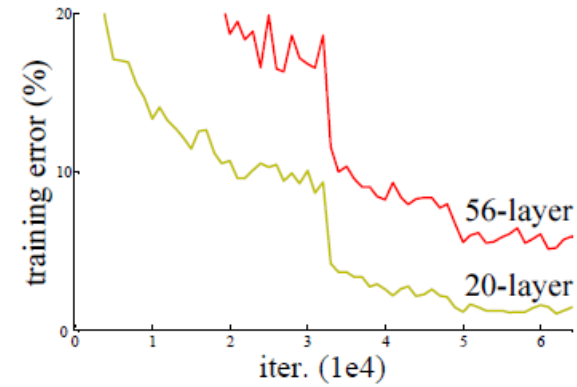
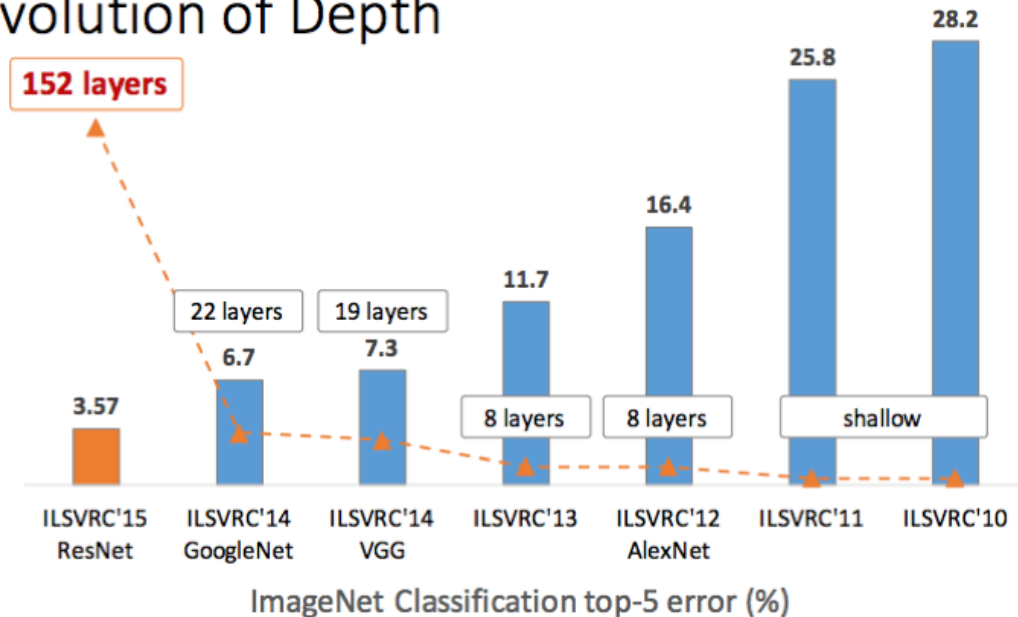
Transfer learning



ResNet-50

- 2015 ImageNet Large Scale Visual Recognition Challenge(ILSVRC) winner.
- Developed by Microsoft
- "Deep Residual Learning for Image Recognition" (Kaiming He et al.,)
- Gradient vanishing >> Residual block (skip connection)

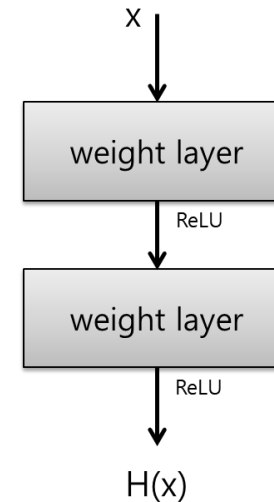
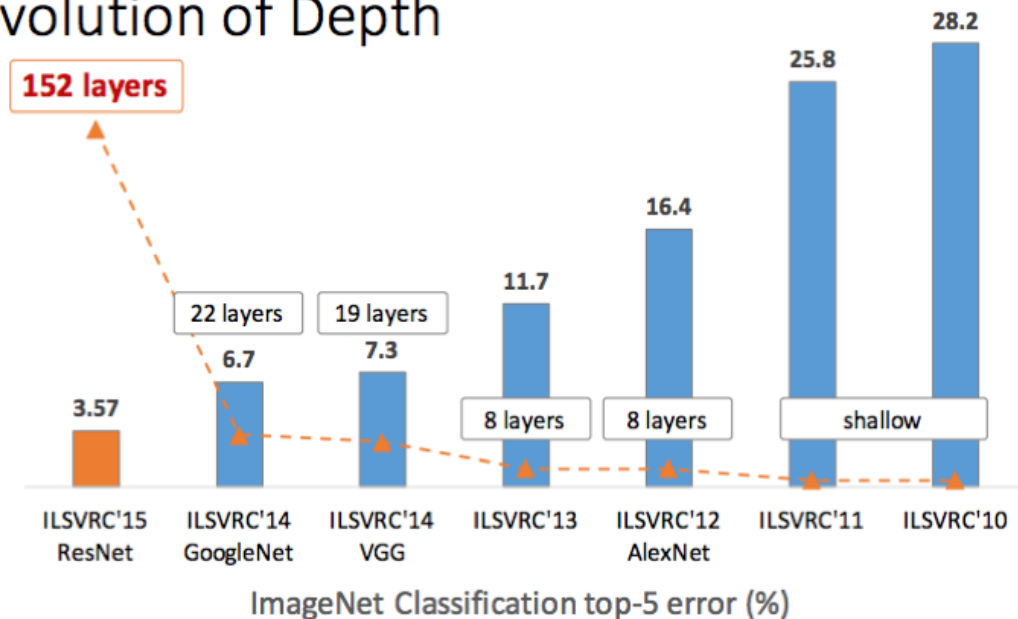
Revolution of Depth



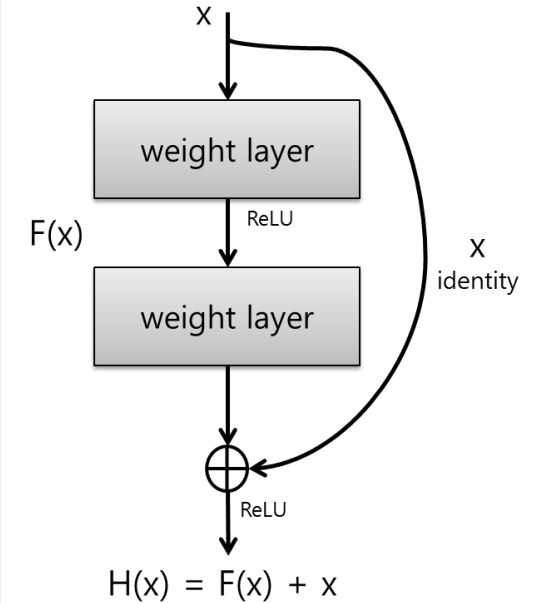
ResNet-50

- 2015 ImageNet Large Scale Visual Recognition Challenge(ILSVRC) winner.
- Developed by Microsoft
- "Deep Residual Learning for Image Recognition" (Kaiming He et al.,)
- Gradient vanishing >> Residual block (skip connection)

Revolution of Depth

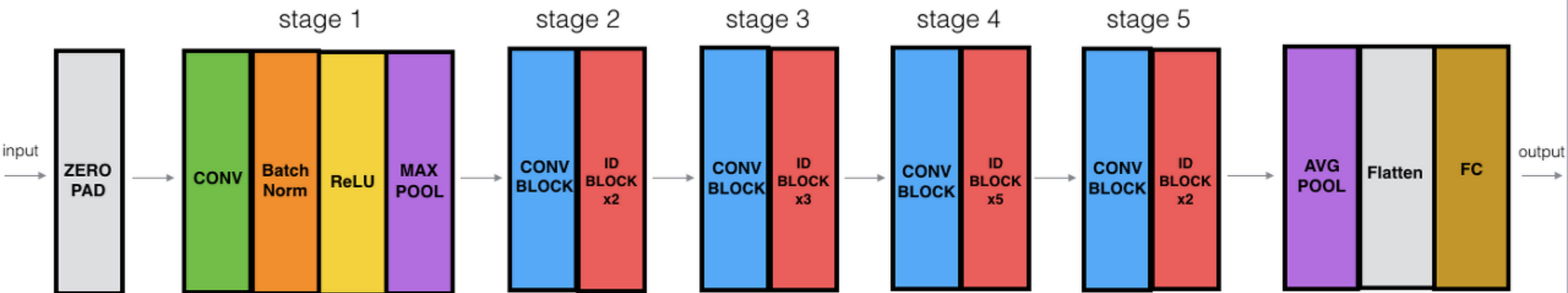


기존 방식



Residual block

ResNet-50



layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112	7×7, 64, stride 2				
conv2_x	56×56	3×3 max pool, stride 2				
		$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
FLOPs		1.8×10^9	3.6×10^9	3.8×10^9	7.6×10^9	11.3×10^9

