

ECE284 Fall 21 W2S2

Low-power VLSI Implementation for Machine Learning

Prof. Mingu Kang

UCSD Computer Engineering

HW1 Graded & HW2 Posted & Typo in Slides

HW1 graded

- Sign extension
- Unused wire
- Parallel if loop
- Reset use case

HW2 posted

- gedit, and using 4 spaces
- PDF download from jupyter notebook

Typo corrected

- W2S1 slides with purple box

Batch Size Choice

Software

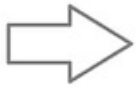
- Pros: fast run-time
- Cons: Large memory consumption -> cause memory fault error
(check with “nvidia-smi” command)

Hardware

- More data re-use opportunity
- Large latency

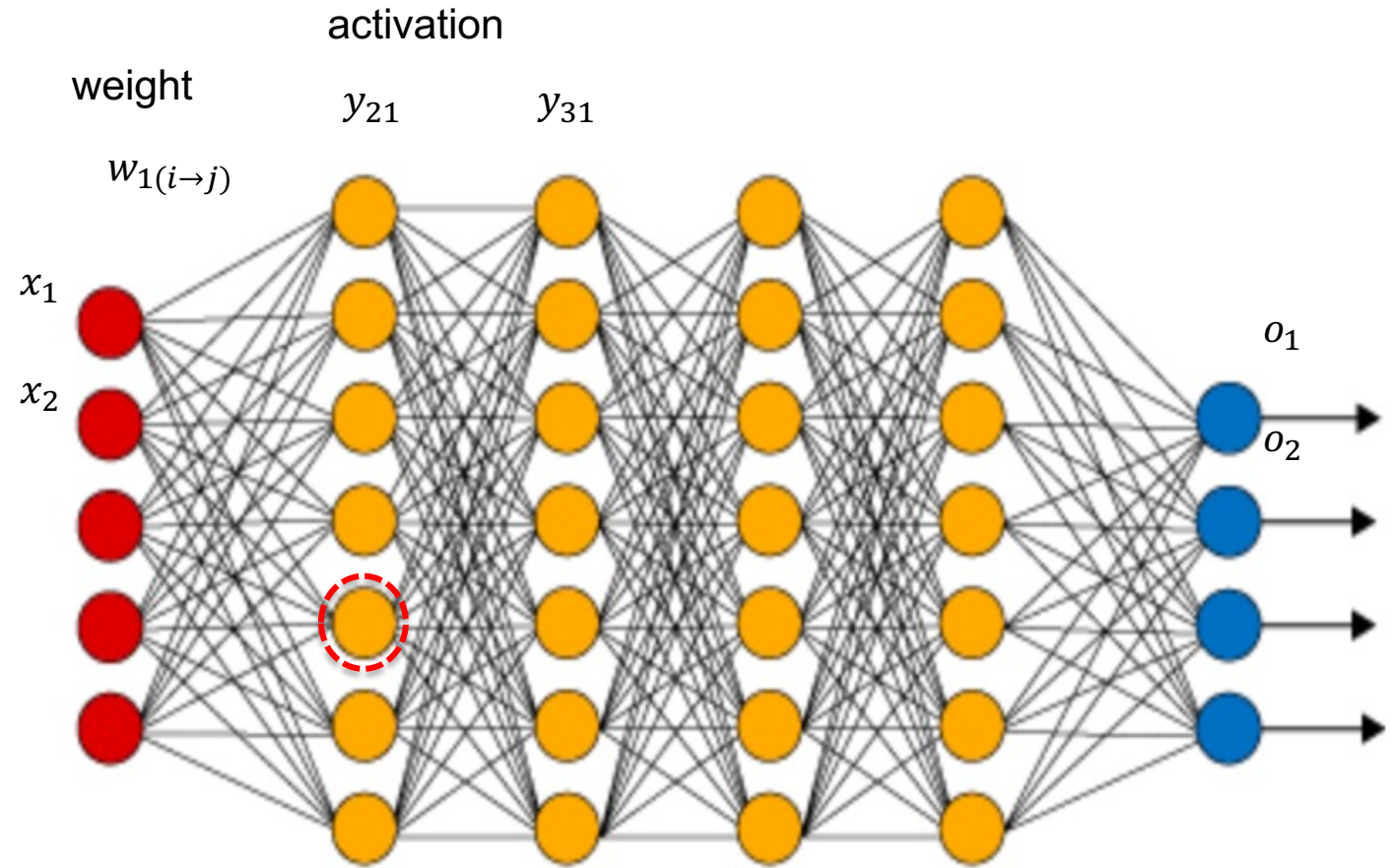
Multi layer Perceptron

1	1	0
4	2	1
0	2	1

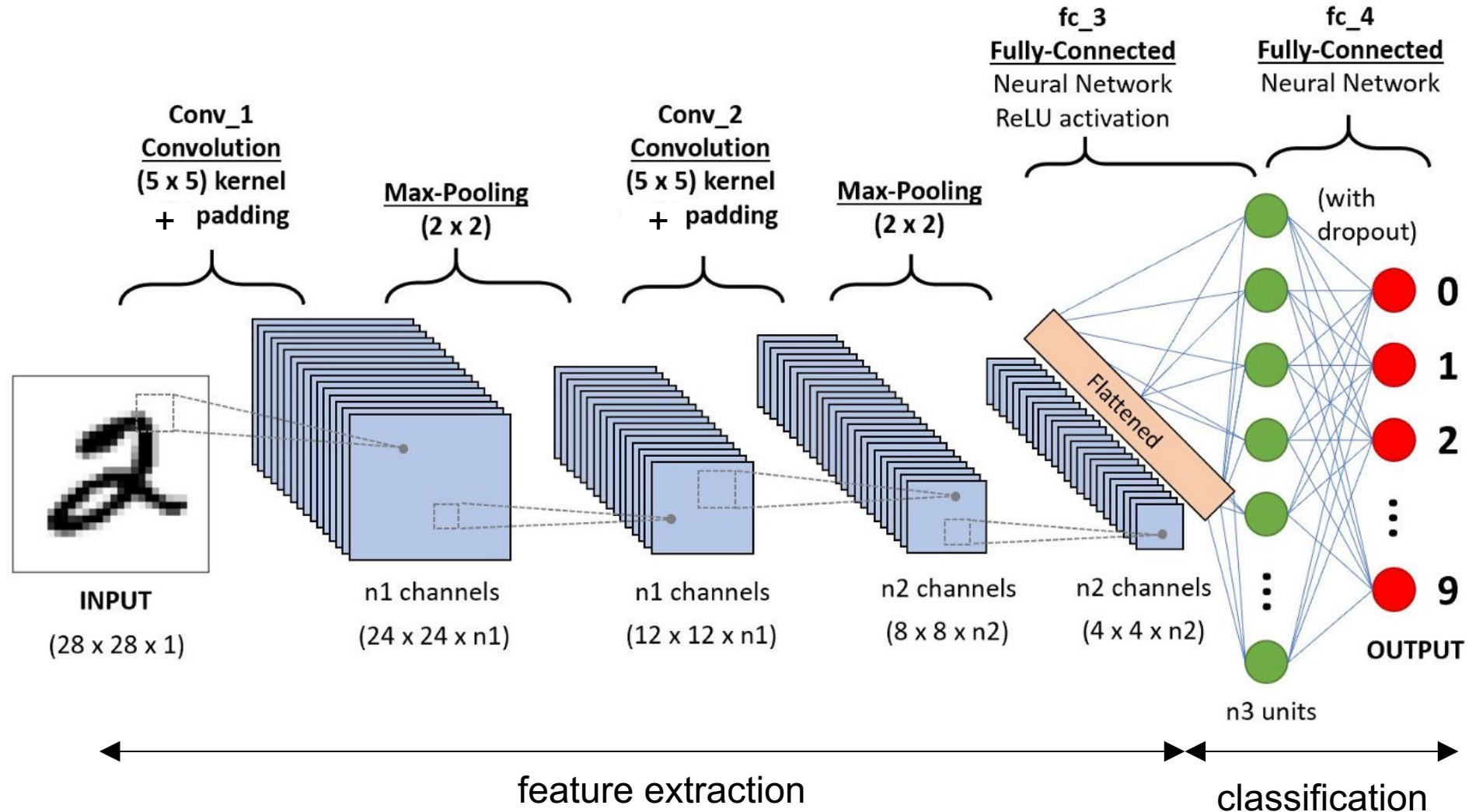


1
1
0
4
2
1
0
2
1

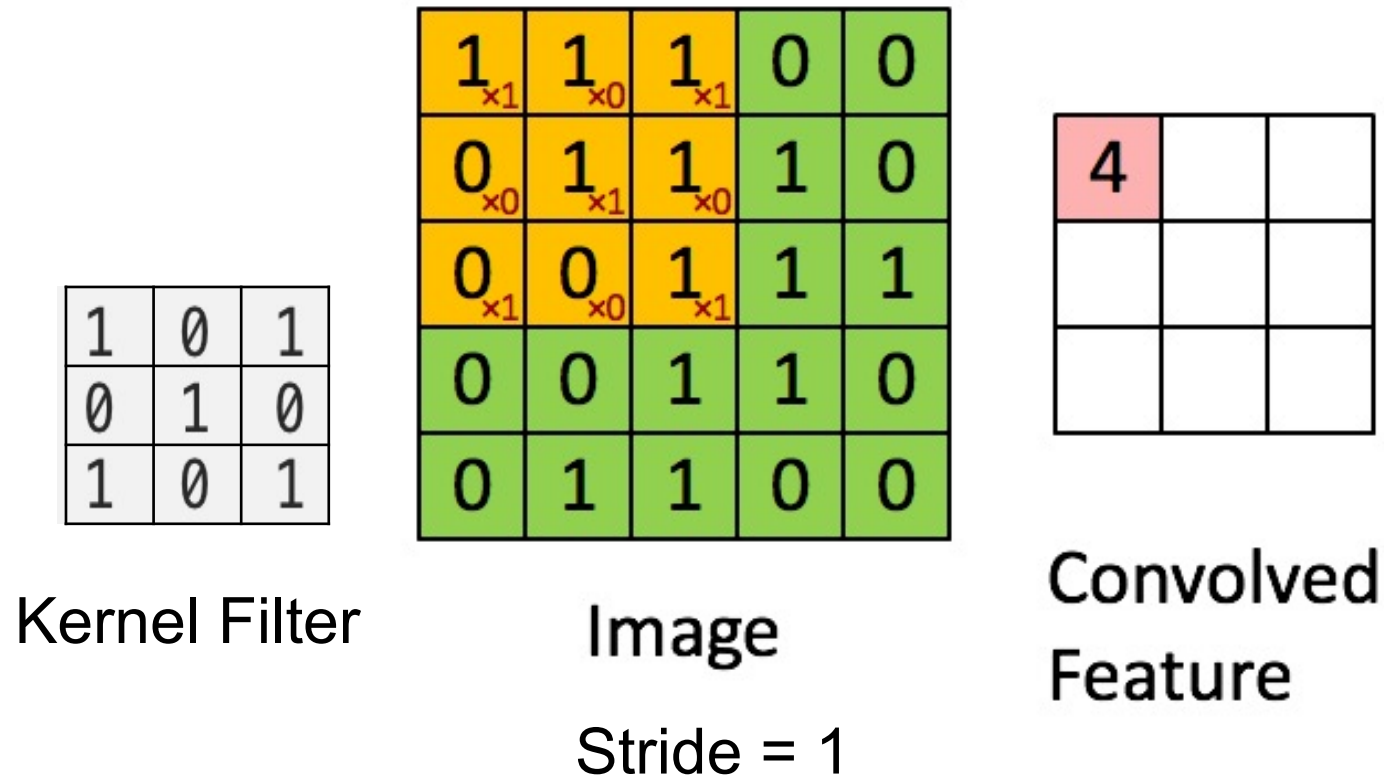
Input layer generation
from 2D image



Convolutional Neural Network



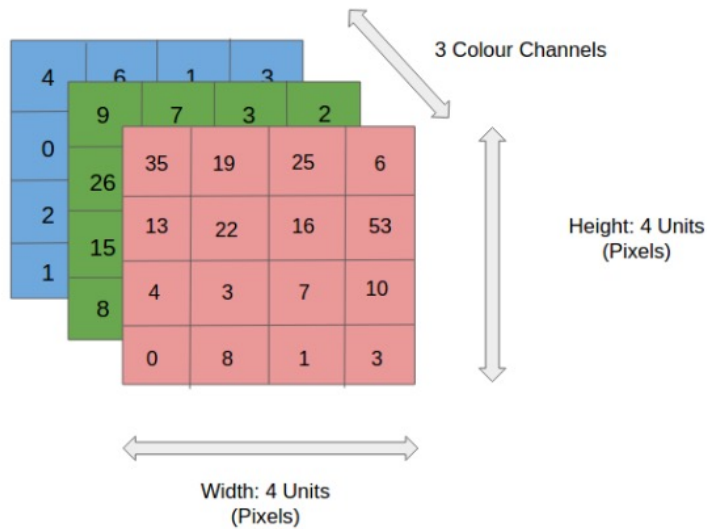
Convolution Layer within Single Channel



Why convolution ?

- capture the Spatial or temporal dependency well
- reduced the data volume in kernel

Convolution Layer across Channels (Color Image)



0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

→ Padding

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2

0	1	1
0	1	0
1	-1	1

Kernel Channel #3

308

+

-498

+

164

+

1 = -25

↑
Bias = 1

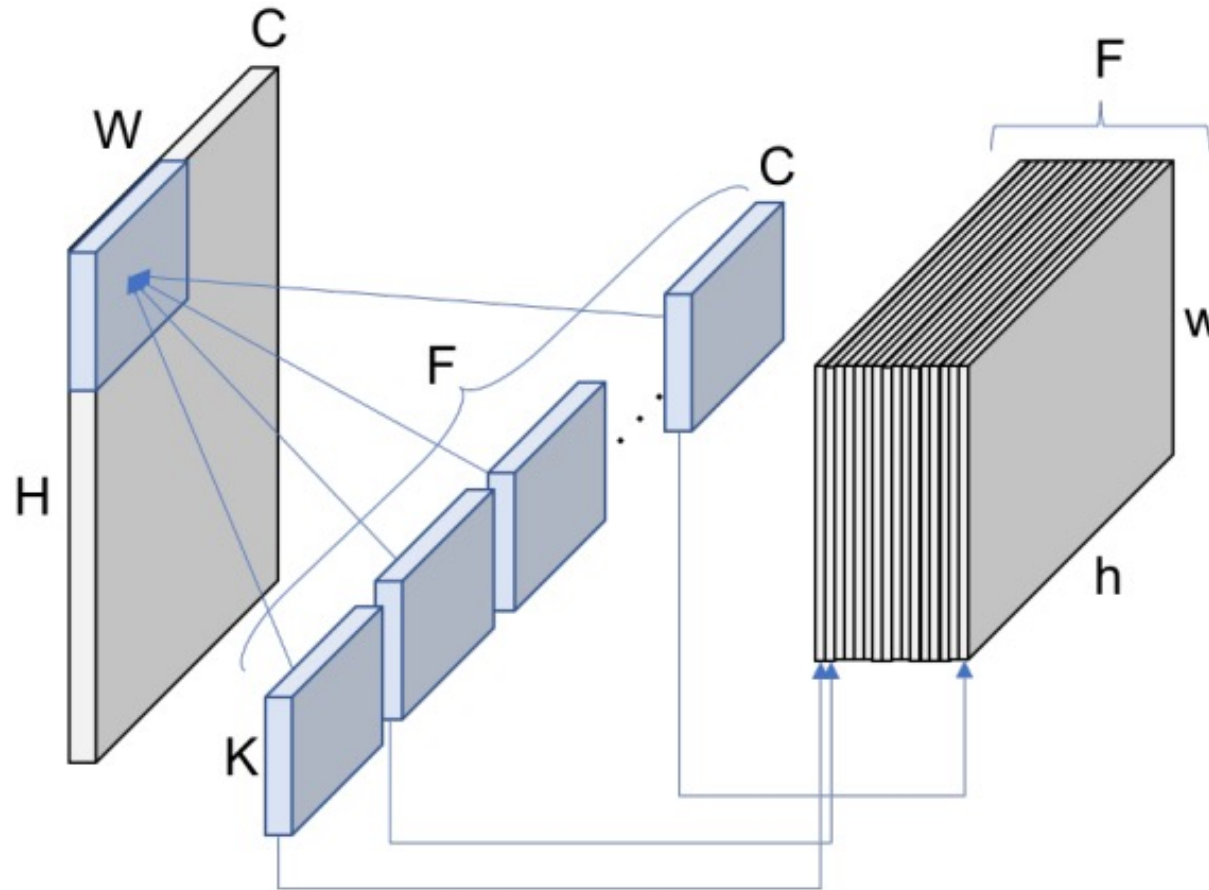
- Above filters are for only 1 output channel
- We need "num_out_ch" such sets of filters

Output

-25				...
				...
				...
				...
...

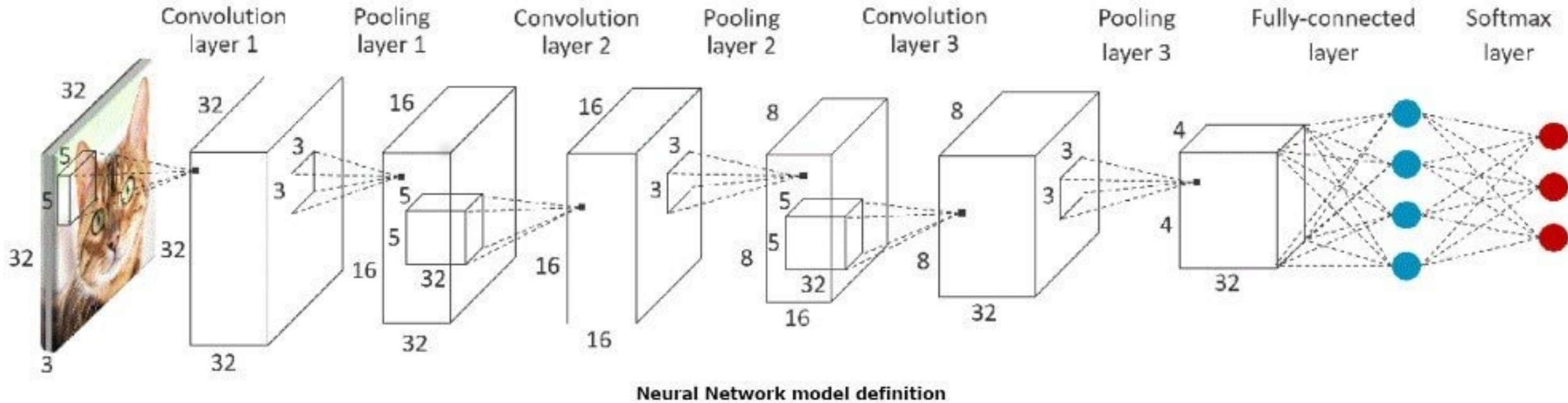
Convolution calculation across channel

3D Representation of Convolution



- C : number of input channels, F : number of output channels
- For each output channel, different kernel filters are required

Convolutional Neural Network for Color Image



Three Data Reuse Opportunities

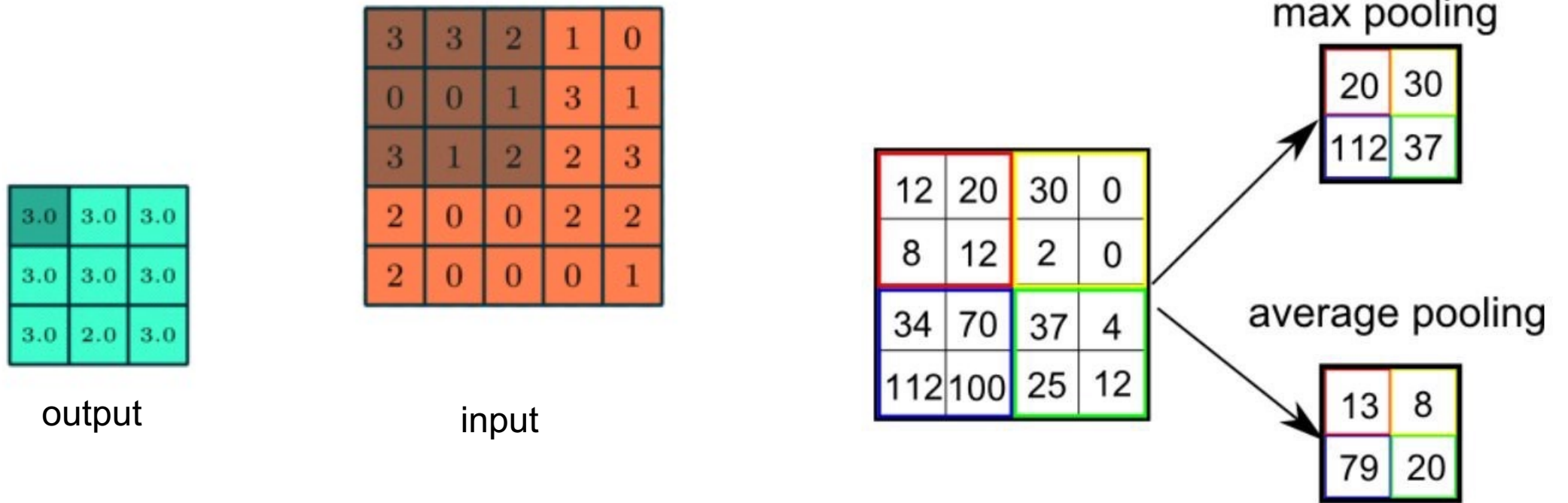
1. Filter (kernel) reuse across input feature map coordinate in convolution
2. Input feature map reuse across output channels
3. Filter (kernel) reuse across data points in the batch

Y.Chen, "Eyeriss: A Spatial Architecture for Energy-Efficient Dataflow for Convolutional Neural Networks", JSSC16

- Still has difficulty as:

- 1 .Convolution kernel is compute-intensive whereas
2. Fully-connected layer is memory-bounded

Pooling Layer

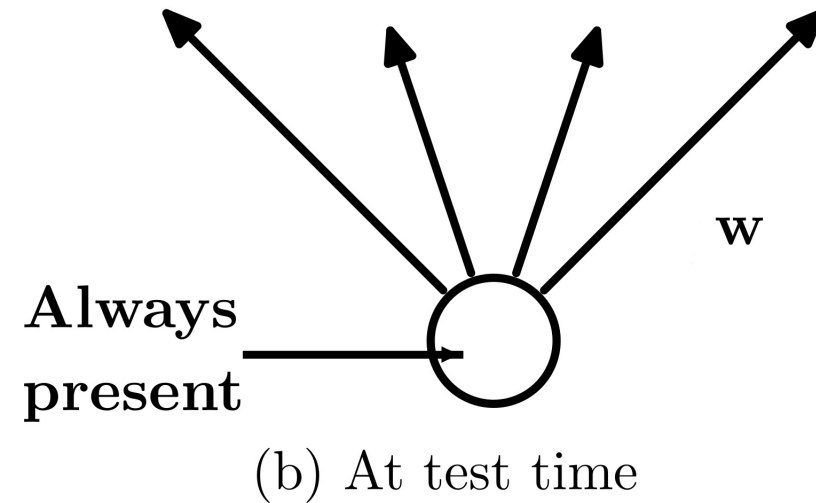
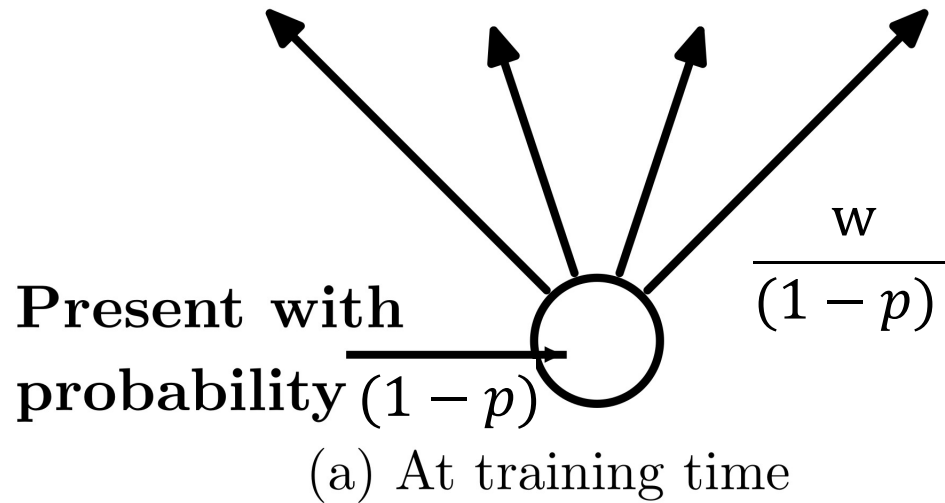


Max pooling operation

Why max or average pooling ?

- extract dominant feature
- reduce the computation power by reducing dimension

(Optional Layers) Dropout Layer



Dropout layer

- makes a certain node zero with a probability of p during training
- helps the “*overfitting*” and co-adaptation problems
- only during training

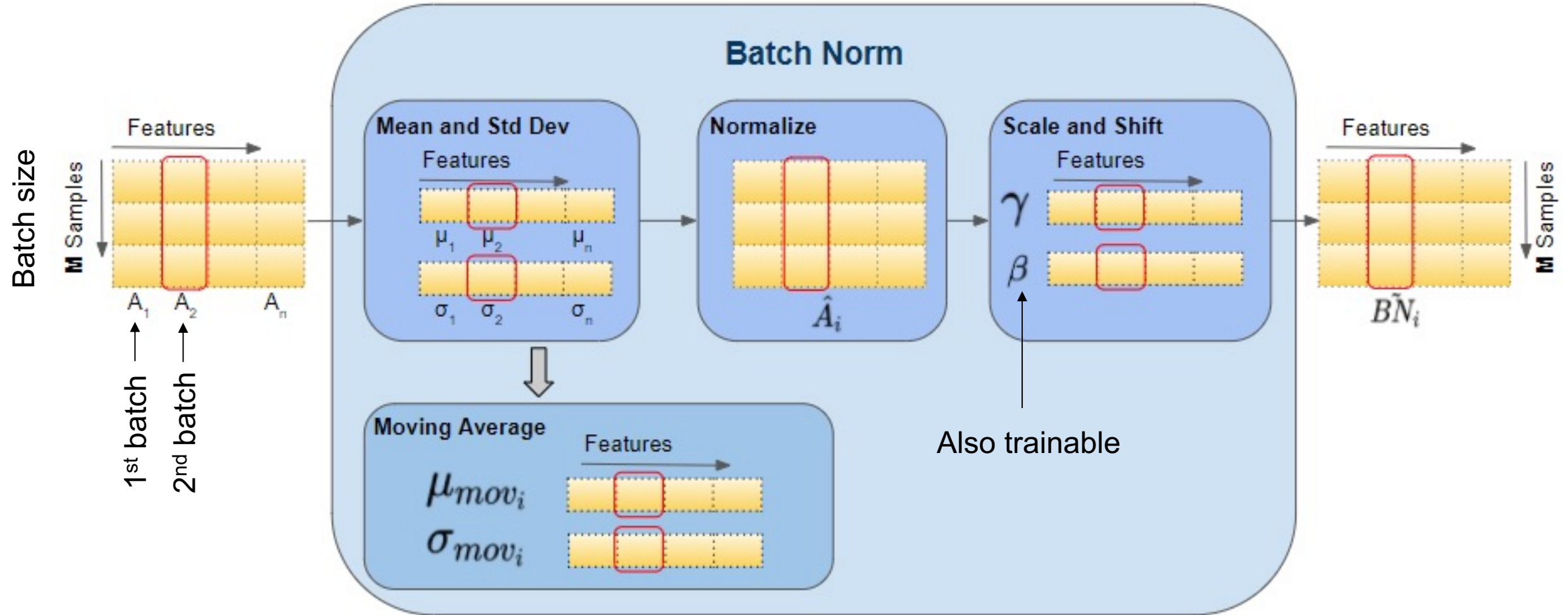
(Optional Layers) Batch Normalization

$$y = \frac{x - \mathbf{E}[x]}{\sqrt{\text{Var}[x] + \epsilon}} * \gamma + \beta$$

Batchnorm layer

- The mean and standard-deviation are calculated per-channel over the mini-batches during training
- γ (default: 1) and β (default: 0) are learnable parameter
- does not update during inference, but just calculate with fixed mean and var

Batchnorm Visualization



- The last value during the training is used for inference

[CODE] Batch-normalization Demo (Example 1)

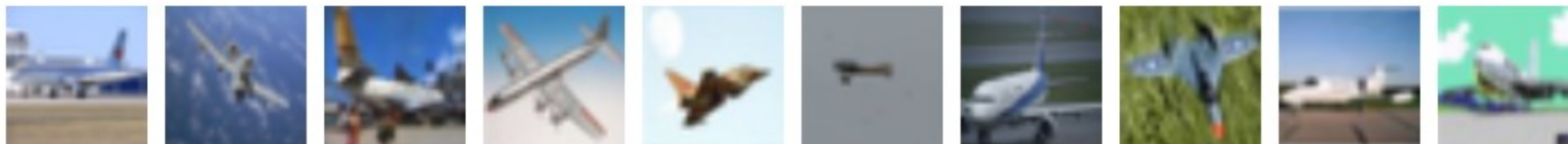
- Input is 2D data
- Check the mean with manual calculation
- Difference during training vs. inference

[CODE] CNN for MNIST (Example 2)

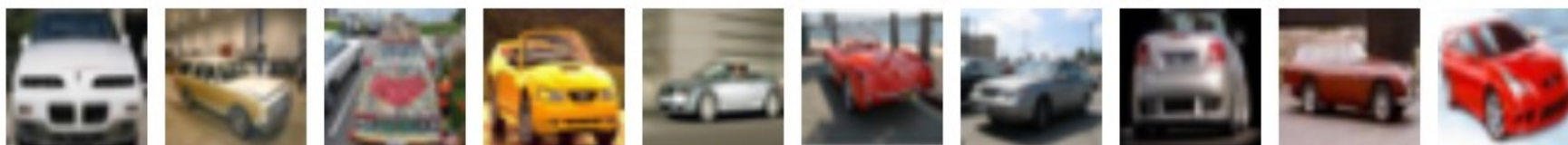
- Introducing GPU use-case
- padding option in conv
- dropout
- Check point save and load
- Network size calculation
 - 1st conv output size = $28 - (3 - 1) = 26$
 - 2nd conv output size = $26 - (3 - 1) = 24 \rightarrow \text{max pool} \rightarrow 12$
- fc1 input = $12^2 * \text{input channel (64)} = 9216$ (Analyze page 5 as well in the same way)
- Pre-hook use-case

CIFAR10 Dataset

airplane



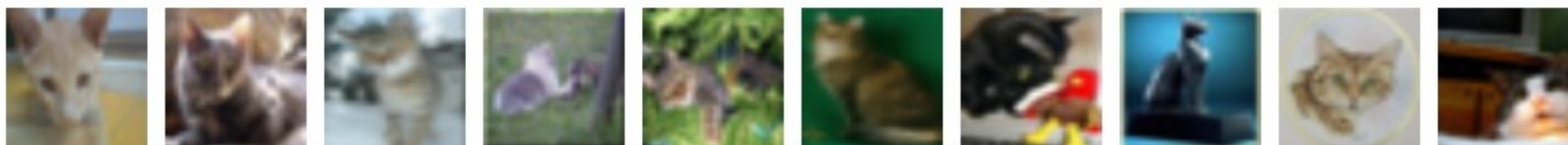
automobile



bird



cat



deer



dog



ImageNet Dataset

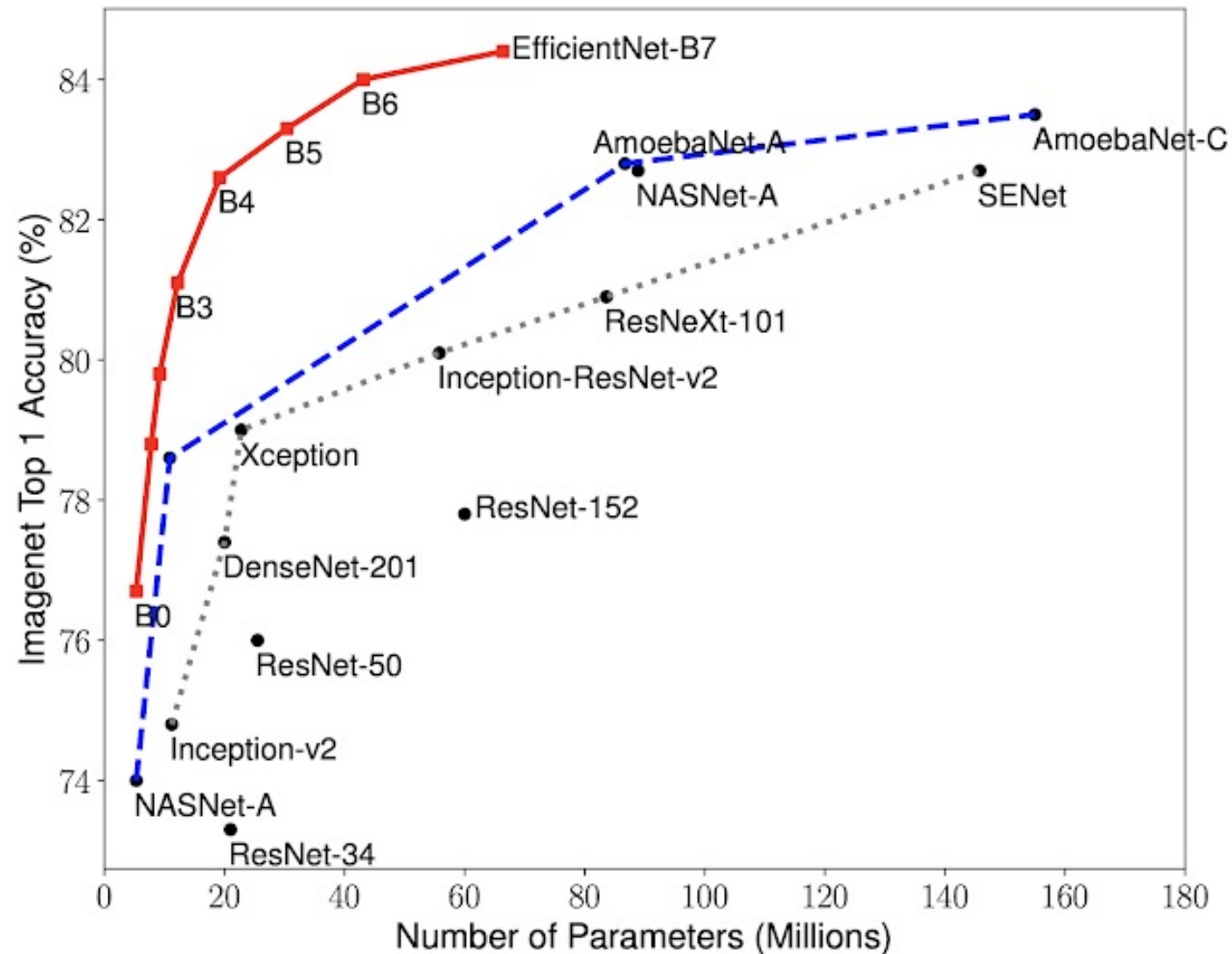


mammal → placental → carnivore → canine → dog → working dog → husky



vehicle → craft → watercraft → sailing vessel → sailboat → trimaran

ImageNet Accuracy Trend



- Total 1000 classes:

Full list of classes [Link](#)

[CODE] CNN Training for CIFAR10 (Example3)