# Machine Learning

CS 539

Worcester Polytechnic Institute

Department of Computer Science

Instructor: Prof. Kyumin Lee

# Upcoming Schedule

- Project Proposal
  - https://canvas.wpi.edu/courses/58900/assignments/355606
  - Due date is June 25


- HW 3
  - https://canvas.wpi.edu/courses/58900/assignments/356655
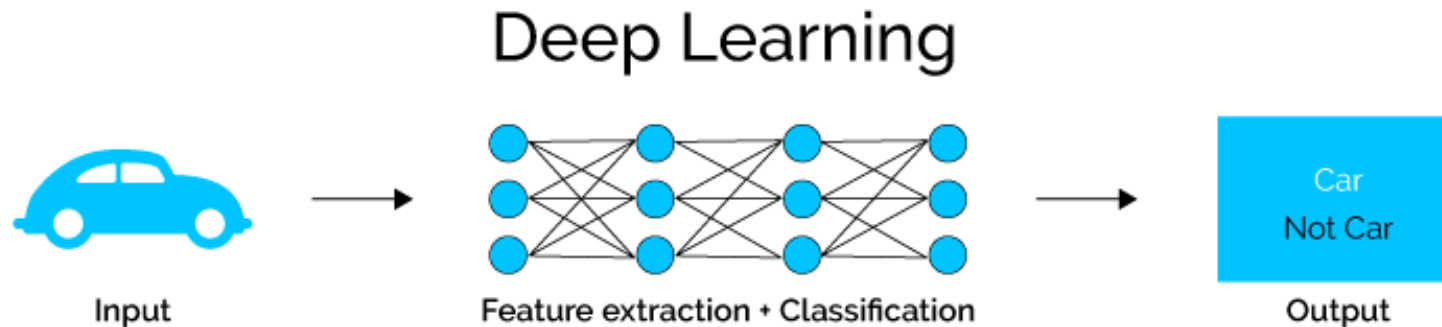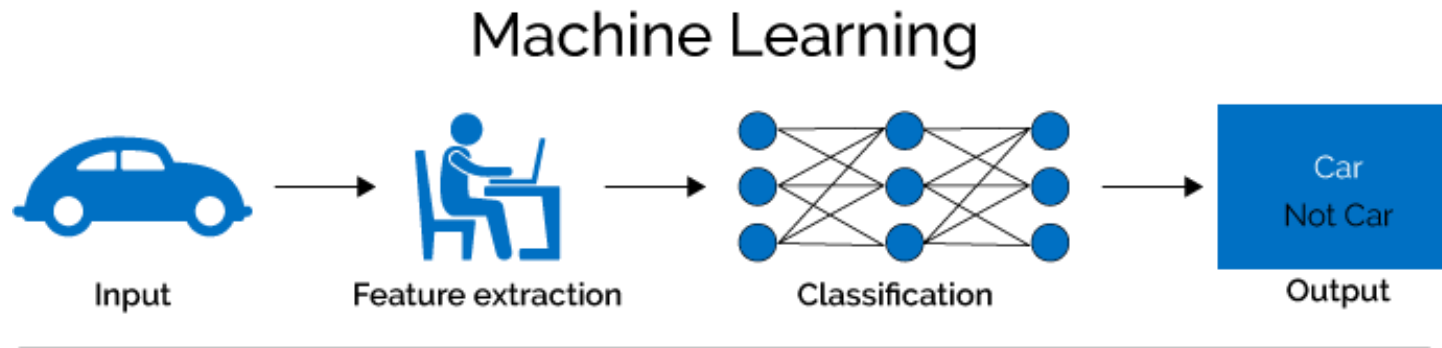  - Due date is July 2
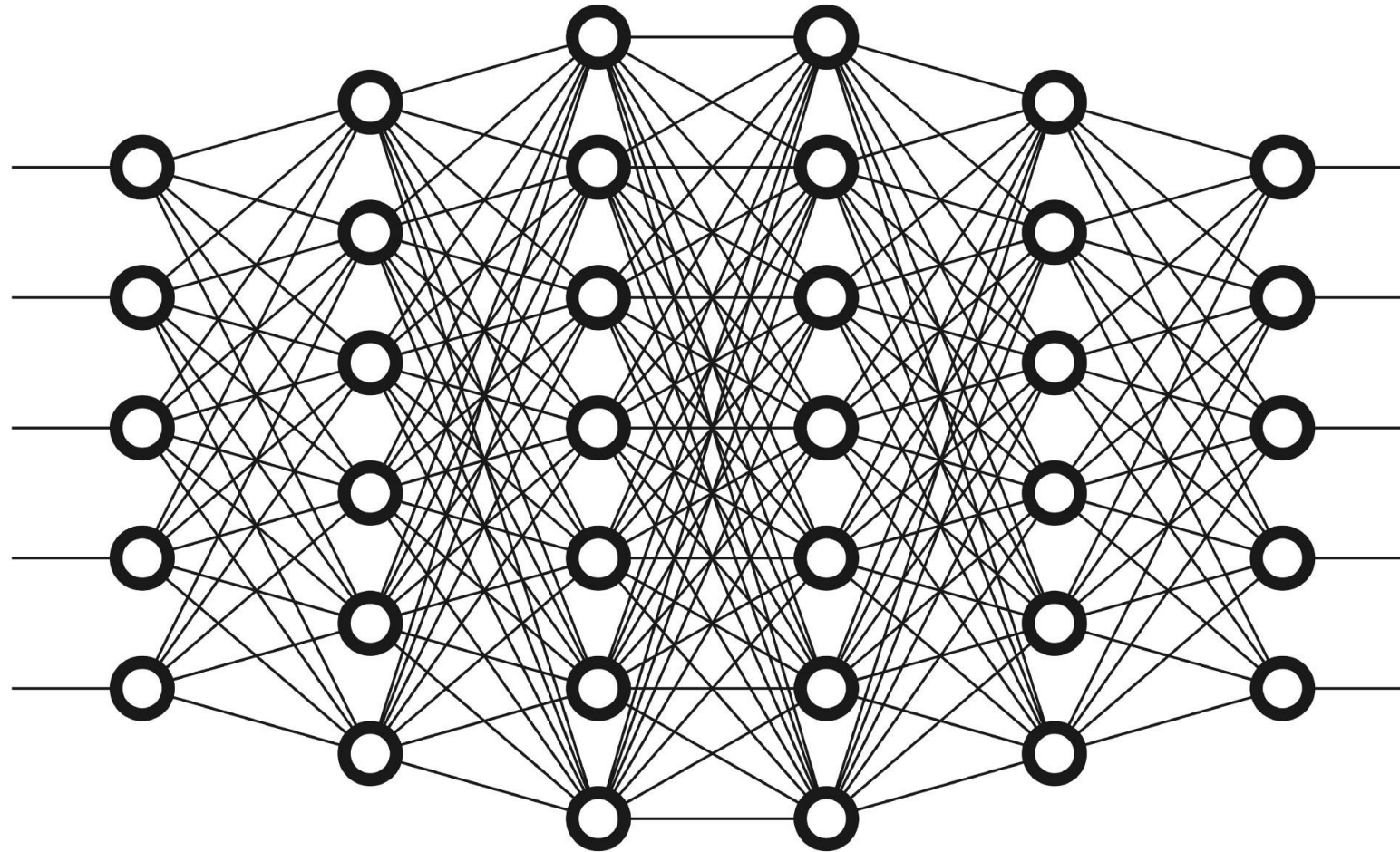
# Deep Learning

# What is Deep Learning (DL) ?

A machine learning subfield of learning **representations** of data. Exceptional effective at **learning patterns**.

Deep learning algorithms attempt to learn (multiple levels of) representation by using a **hierarchy of multiple layers**
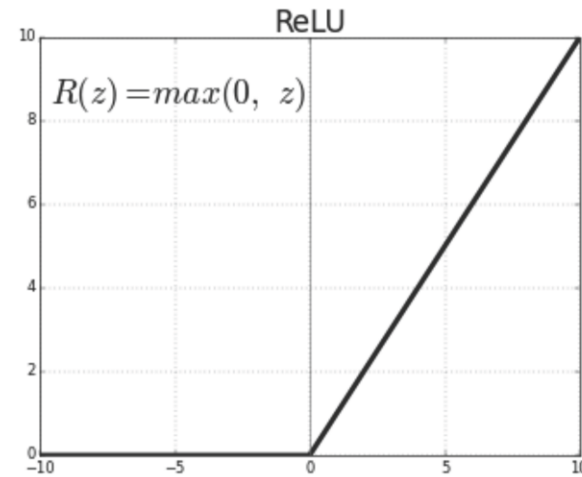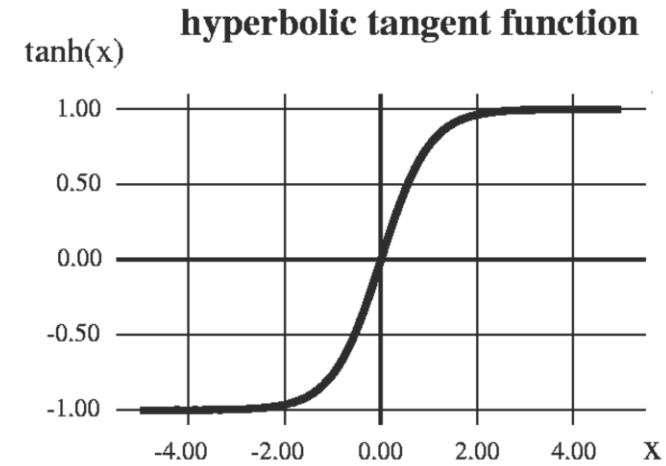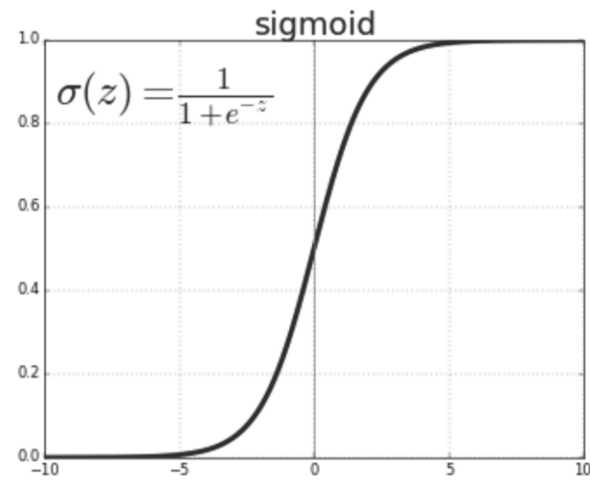
If you provide the system **tons of information**, it begins to understand it and respond in useful ways.

# Deep Neural Network

# Deep Net Activation Functions
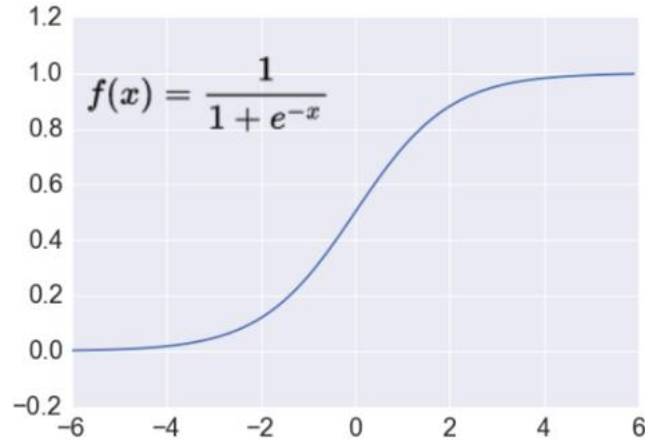


sigmoid

$$\sigma(z) = \frac{1}{1+e^{-z}}$$

hyperbolic tangent function

tanh(x)

ReLU

$$R(z) = max(0, \ z)$$

Rectified Linear Unit (ReLU)

# Activation: Sigmoid

$$f(x) = \frac{1}{1 + e^{-x}}$$

*http://adilmoujahid.com/images/activation.png*
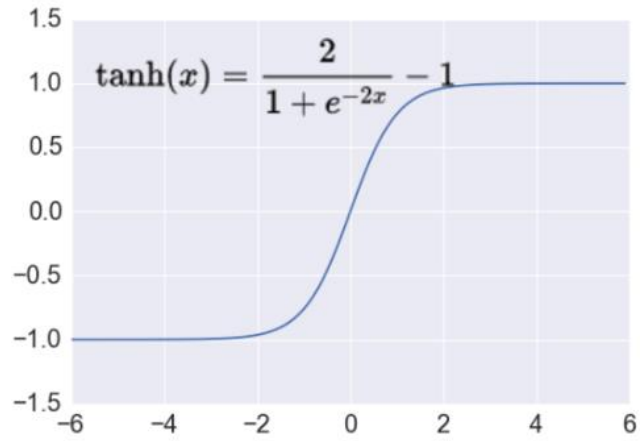
Takes a real-valued number and "squashes" it into range between 0 and 1.

$$R^n \rightarrow [0,1]$$

+ Nice interpretation as the **firing rate** of a neuron
  - 0 = not firing at all
  - 1 = fully firing

- Sigmoid neurons **saturate** and **kill gradients**, thus NN will barely learn
  - when the neuron's activation are 0 or 1 (saturate)
    - 😖 gradient at these regions almost zero
    - 😖 almost no signal will flow to its weights
    - 😖 if initial weights are too large then most neurons would saturate

# Activation: Tanh

$$\tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

http://adilmoujahid.com/images/activation.png
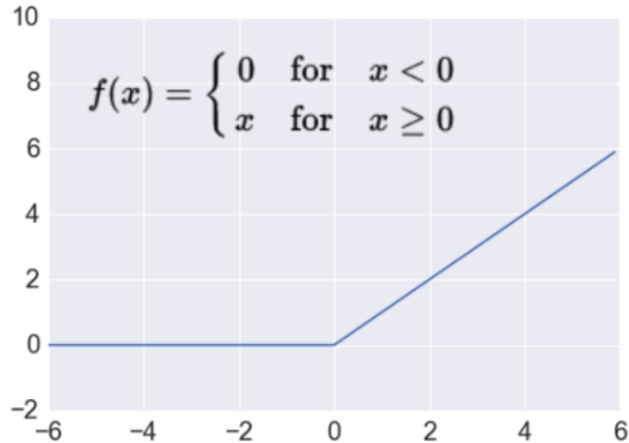
Takes a real-valued number and "squashes" it into range between -1 and 1.

$$R^n \rightarrow [-1,1]$$

- Like sigmoid, tanh neurons **saturate**
- Unlike sigmoid, output is **zero-centered**
- Tanh is a **scaled sigmoid**: $\tanh(x) = 2sigm(2x) - 1$

# Activation: ReLU

$$f(x) = \begin{cases} 0 & \text{for} \quad x < 0 \\ x & \text{for} \quad x \geq 0 \end{cases}$$

Takes a real-valued number and thresholds it at zero $\quad f(x) = \max(0, x)$

$$R^n \rightarrow R_+^n$$

Most Deep Networks use ReLU nowadays

🙂 Trains much **faster**
  - accelerates the convergence of SGD
  - due to linear, non-saturating form

🙂 Less expensive operations
  - compared to sigmoid/tanh (exponentials etc.)
  - implemented by simply thresholding a matrix at zero

🙂 Prevents the **gradient vanishing problem**

# Hardware Power

a single-core 2.5 GHz CPU has a theoretical performance of 10 billion FLOPS = 10 GFLOPS

1 CPU core = 10 GFlops

Graphics processing unit

1 GPU

1 multi-core Hyper-threading CPU

= 8 Cores X 2 thread

= 16 CPUs

CPU
MULTIPLE CORES

GPU
THOUSANDS OF CORES

# What is Flops?

FLoating-point Operations Per Second

measure of computer performance, useful in fields of scientific computations that require floating-point calculations
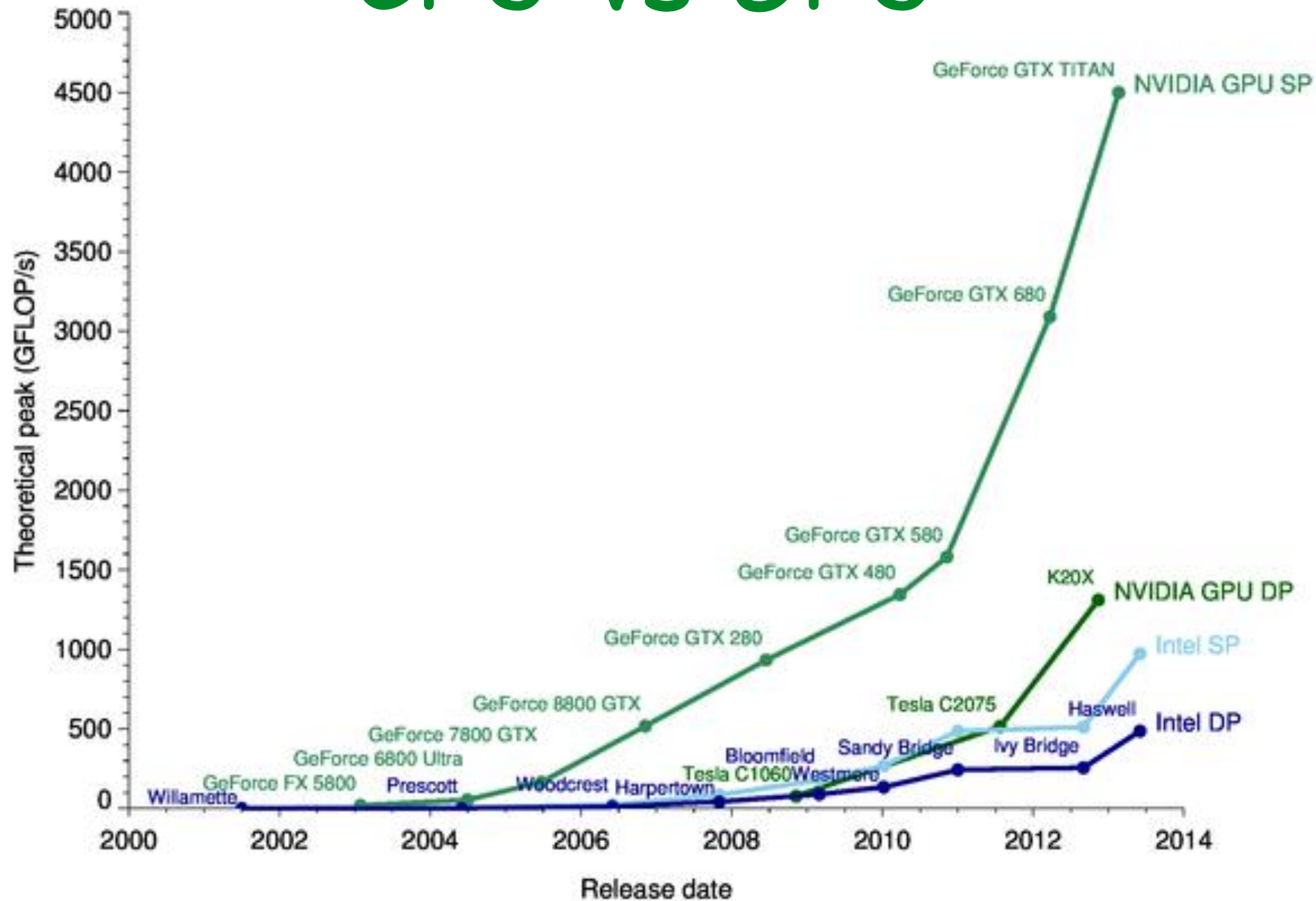
Loan: 325,000

Interest Rate(APR): 2.225%

Term: 1 year

Interest Paid: ?        7,312.5

# CPU vs GPU

# GeForce RTX 4090

## BEYOND FAST

The NVIDIA® GeForce RTX™ 4090 is the ultimate GeForce GPU. It brings an enormous leap in performance, efficiency, and AI-powered graphics. Experience ultra-high performance gaming, incredibly detailed virtual worlds, unprecedented productivity, and new ways to create. It's powered by the NVIDIA Ada Lovelace architecture and comes with 24 GB of G6X memory to deliver the ultimate experience for gamers and creators.
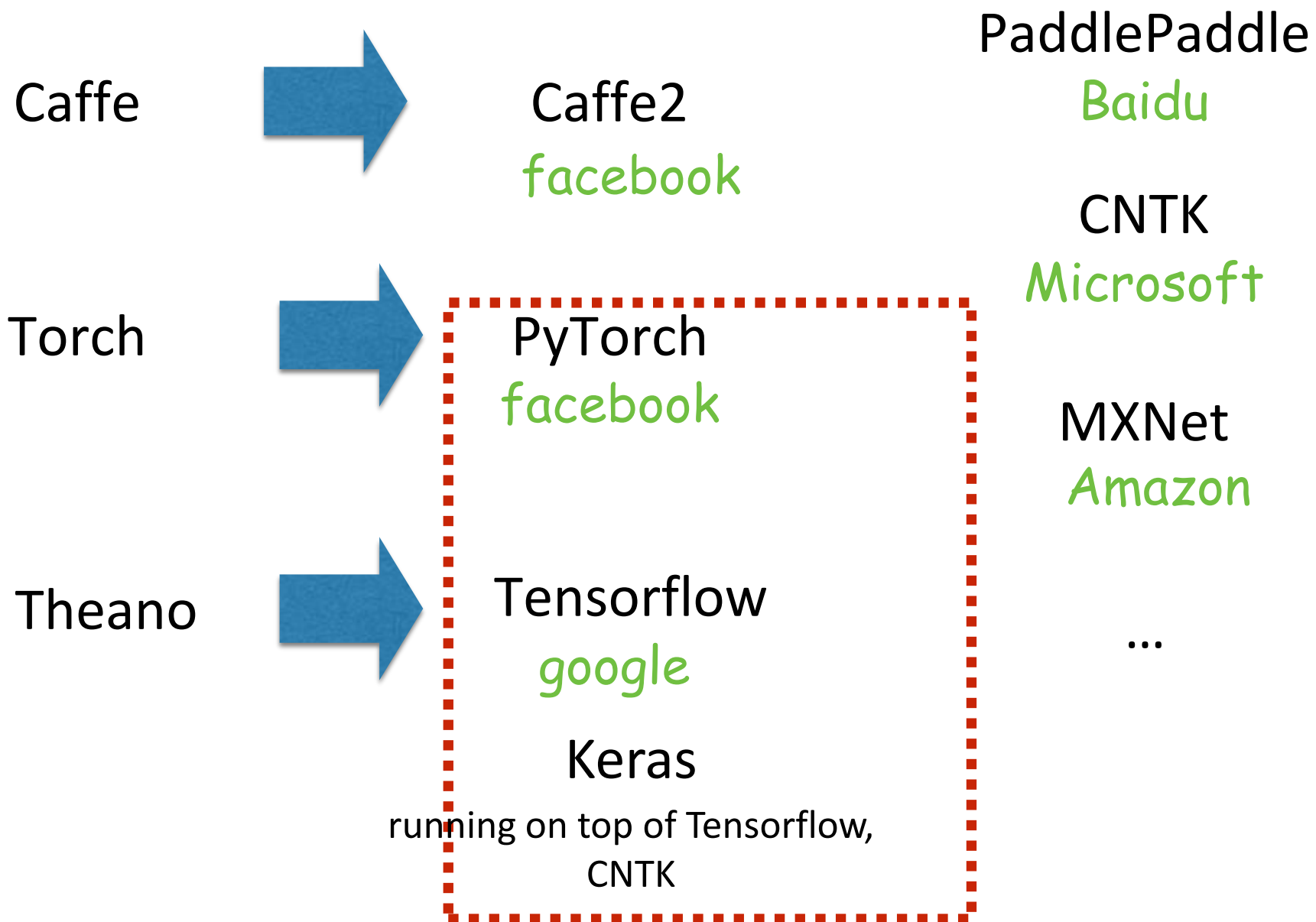
**Starting at $1599.00**

See All Buying Options

| | RTX 4090 | RTX 3090 Ti | RTX 3090 |
|---|---|---|---|
| **Architecture** | Ada Lovelace | Ampere | Ampere |
| **GPU** | AD102 | GA102 | GA102 |
| **Process node** | 5nm TSMC | 8nm Samsung | 8nm Samsung |
| **CUDA cores** | 16,384 | 10,752 | 10,496 |
| **RT cores** | 144 3rd-generation | 84 | 82 |
| **Tensor cores** | 576 4th-gen | 336 | 328 |
| **Base clock** | 2235MHz | 1560MHz | 1395Mhz |
| **Boost clock** | 2520MHz | 1860MHz | 1695MHz |
| **Memory** | 24GB GDDR6X | 24GB GDD6X | 24GB GDDR6X |
| **Memory speed** | 21Gbps | 21Gbps | 19.5Gbps |
| | 83 TFLOPS | 40 TFLOPS | 36 TFLOPS |

# Deep Learning Software

# Deep Learning Software
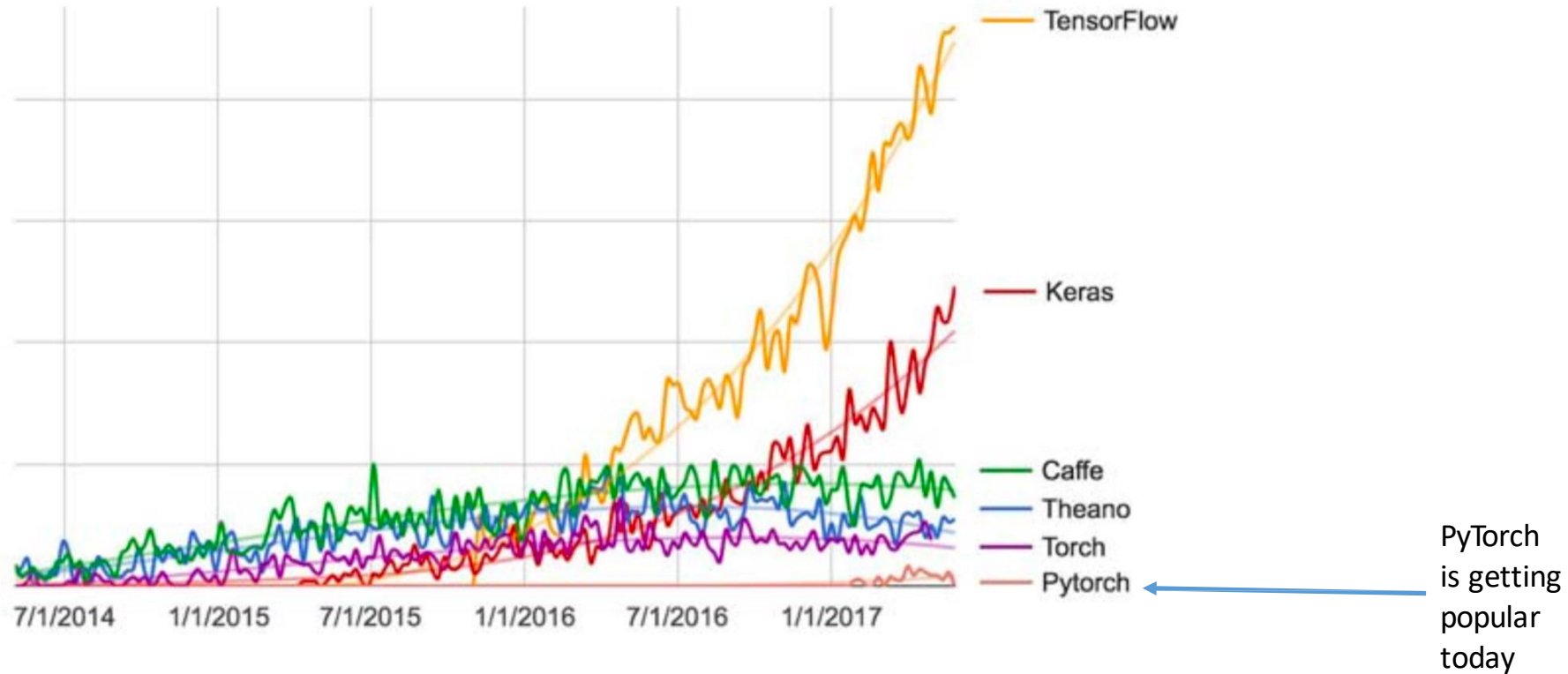


Figure 3.2 Google web search interest for different deep-learning frameworks over time
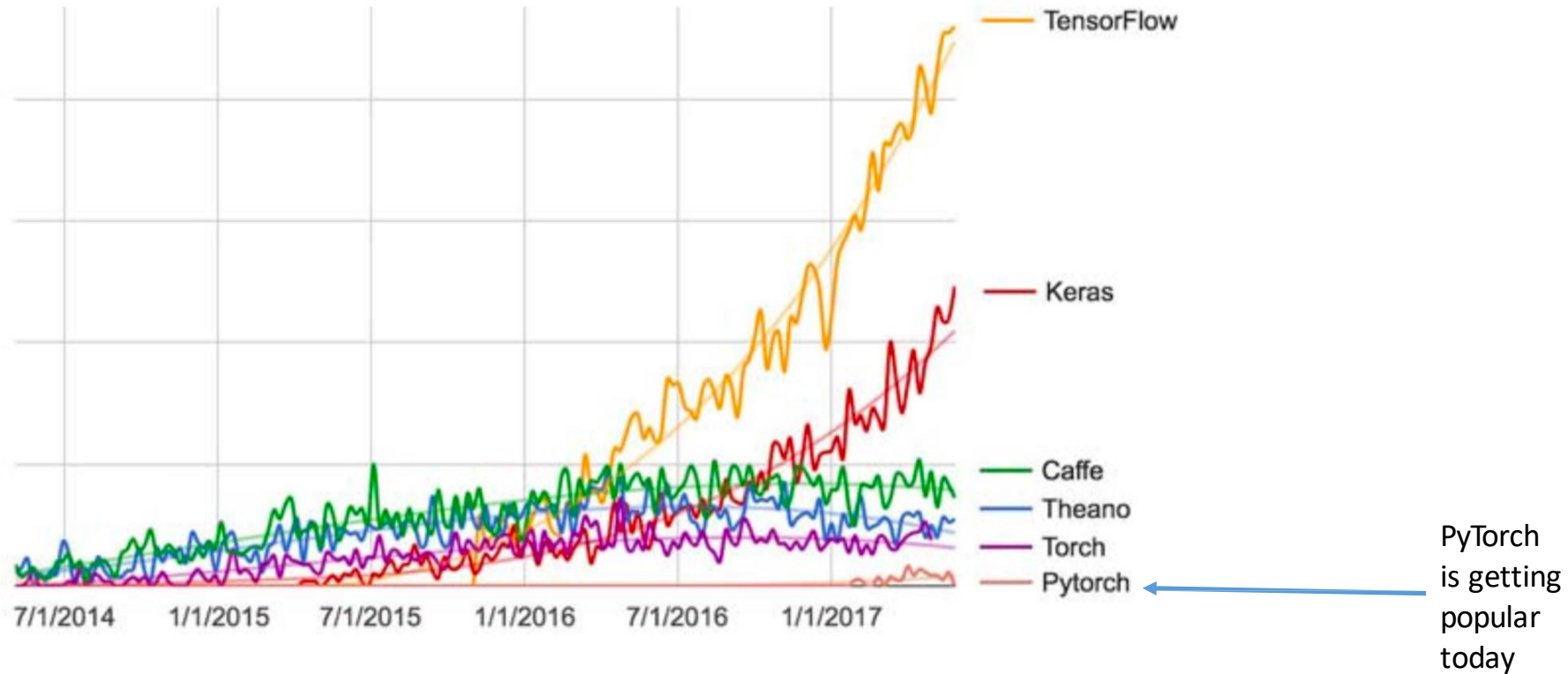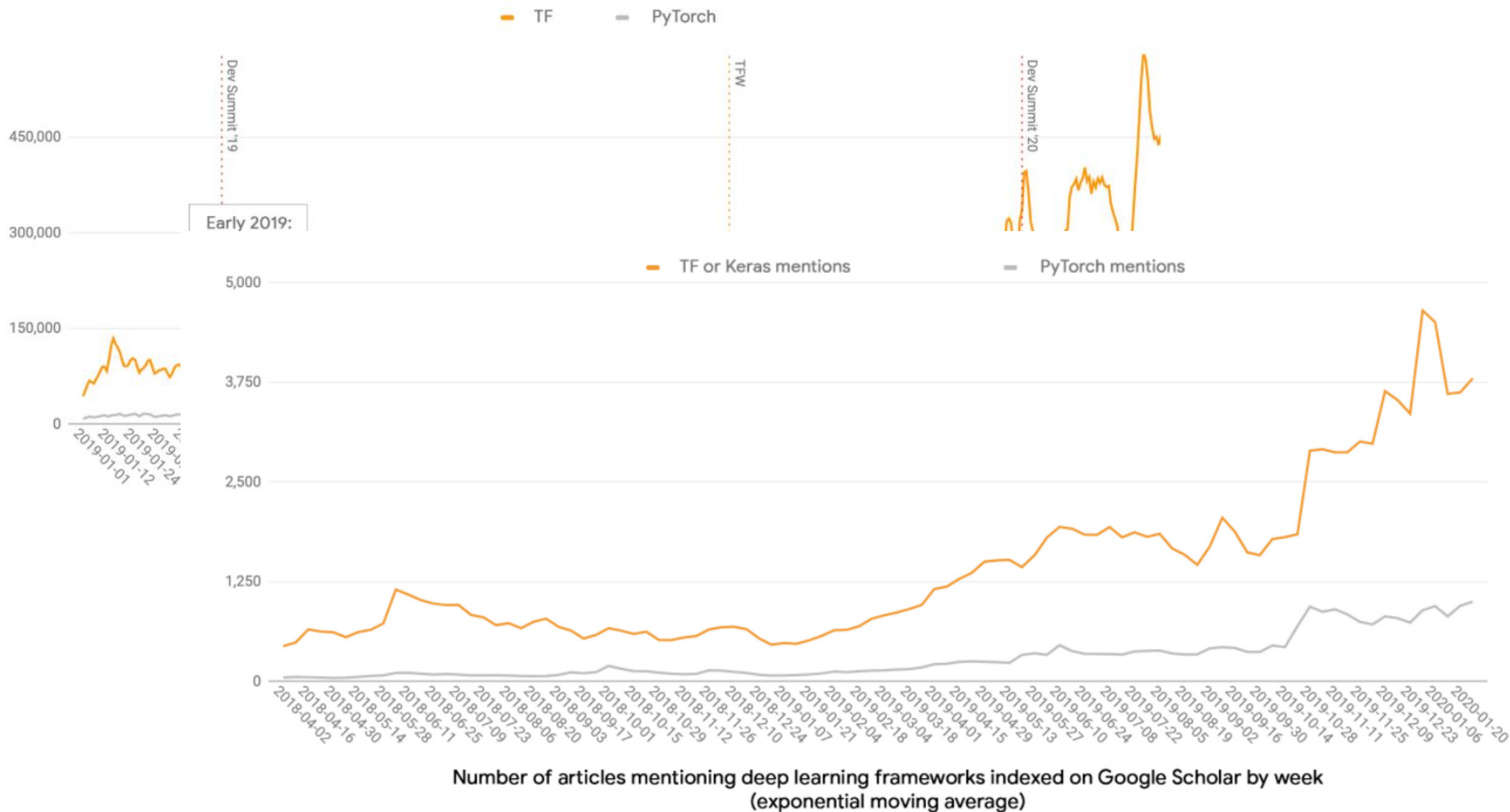
# Deep Learning Software



PyTorch is getting popular today

Figure 3.2    Google web search interest for different deep-learning frameworks over time

**TF** — **PyTorch**

Dev Summit '19

TFW

Dev Summit '20

Early 2019:

450,000

300,000

150,000

0

2019-01-01  2019-01-12  2019-01-24

**TF or Keras mentions** — **PyTorch mentions**

5,000

3,750

2,500

1,250

0

2018-04-02  2018-04-16  2018-04-30  2018-05-14  2018-05-28  2018-06-11  2018-06-25  2018-07-09  2018-07-23  2018-08-06  2018-08-20  2018-09-03  2018-09-17  2018-10-01  2018-10-15  2018-10-29  2018-11-12  2018-11-26  2018-12-10  2018-12-24  2019-01-07  2019-01-21  2019-02-04  2019-02-18  2019-03-04  2019-03-18  2019-04-01  2019-04-15  2019-04-29  2019-05-13  2019-05-27  2019-06-10  2019-06-24  2019-07-08  2019-07-22  2019-08-05  2019-08-19  2019-09-02  2019-09-16  2019-09-30  2019-10-14  2019-10-28  2019-11-11  2019-11-25  2019-12-09  2019-12-23  2020-01-06  2020-01-20

Number of articles mentioning deep learning frameworks indexed on Google Scholar by week
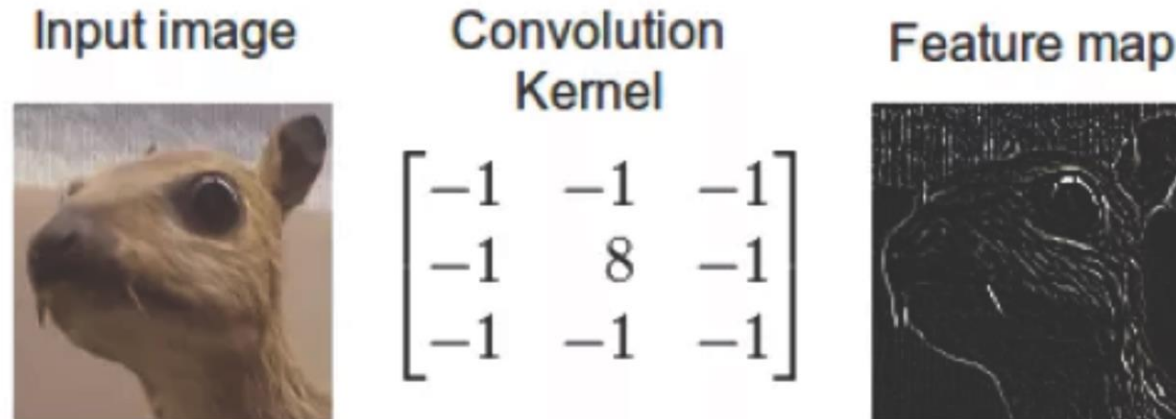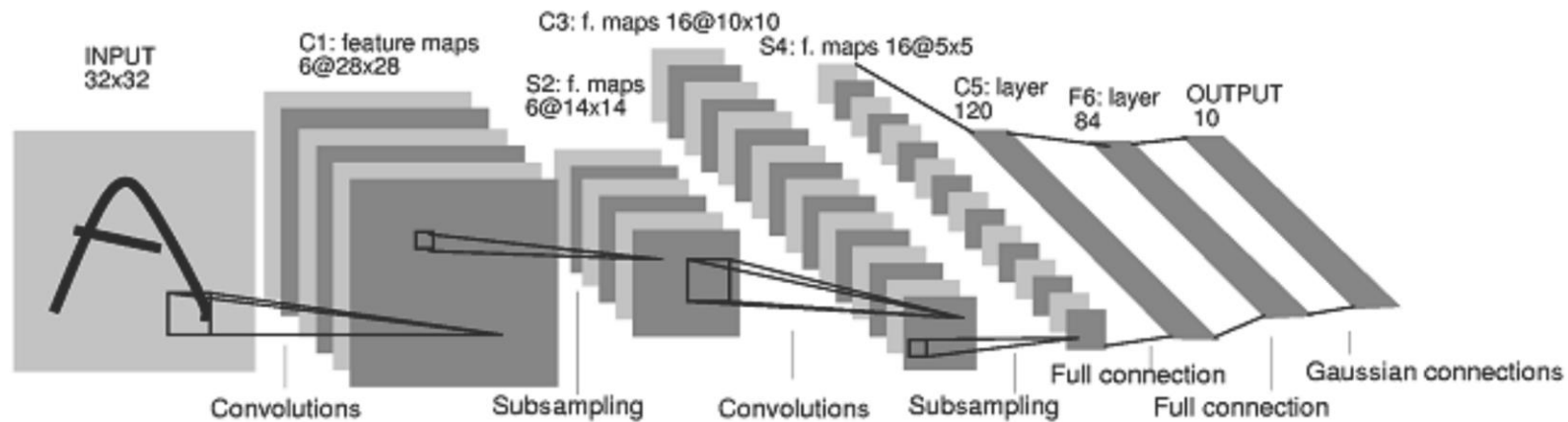(exponential moving average)

# Convolutional Neural Network (CNN)

**Convolution:**

Convolution is a mathematical operation on two functions to produce a third function that expresses how the shape of one is modified by the other. It is a term that the neural network community (Hinton, specifically) adopted from the signal and image processing communities to architect the "feature detection" layers of "deep" neural networks.
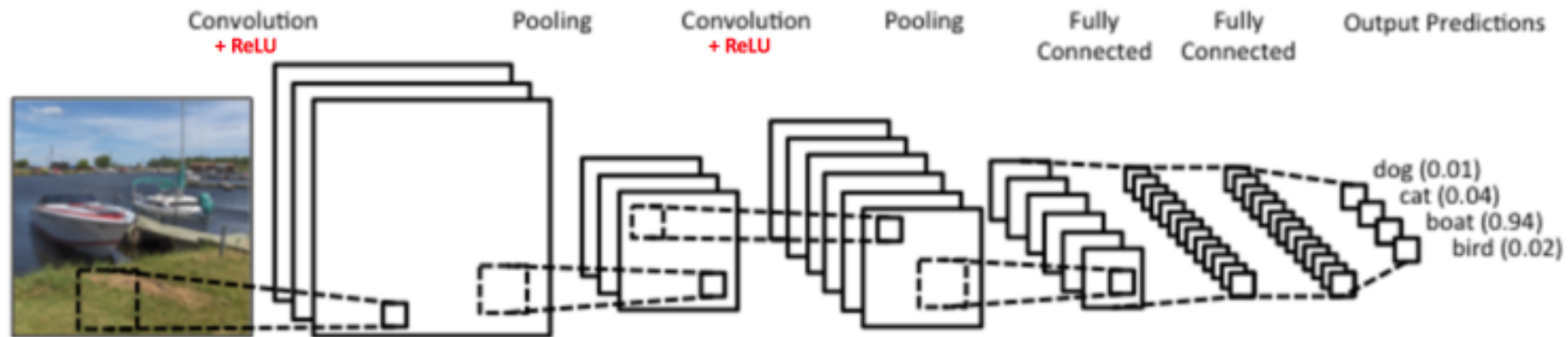
# Convolutional Neural Network



Input image    Convolution Kernel    Feature map

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

[Image credit: http://timdettmers.com/2015/03/26/convolution-deep-learning/]



INPUT 32x32

C1: feature maps 6@28x28

C3: f. maps 16@10x10

S2: f. maps 6@14x14

S4: f. maps 16@5x5

C5: layer 120

F6: layer 84

OUTPUT 10

Convolutions    Subsampling    Convolutions    Subsampling    Full connection    Gaussian connections

Full connection

Reading: https://ujjwalkarn.me/2016/08/11/intuitive-explanation-convnets/

# Convolutional Neural Network
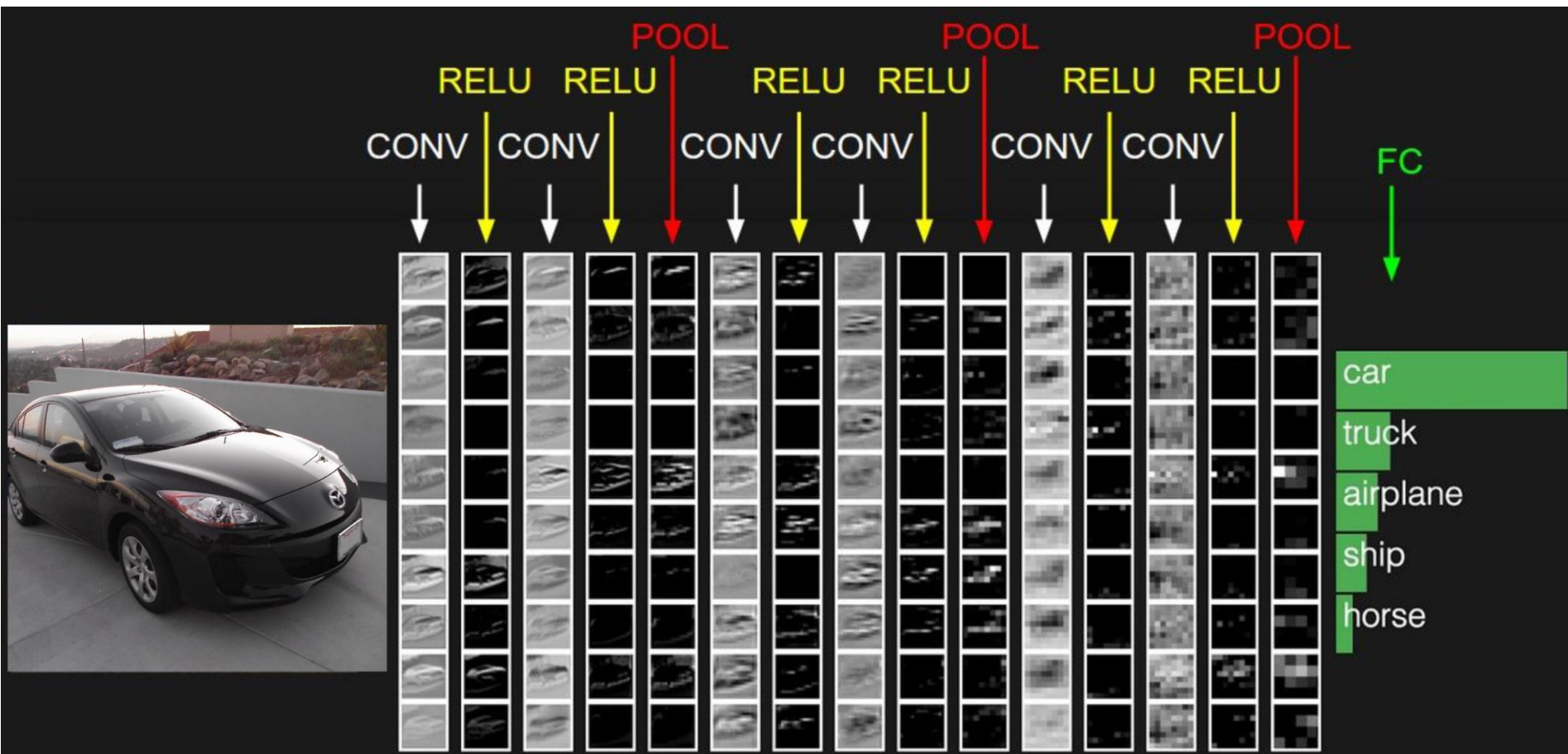


A simple ConvNet
http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/

# Convolutional Layer

- ## Filter  $\begin{pmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{pmatrix}$



Input Image



Convoluted Image

■ Inspired by the neurophysiological experiments conducted by Hubel and Wiesel 1962.

# Convolutional Layer



Input matrix

Convolutional 3x3 filter

Image

Convolved Feature

*http://deeplearning.stanford.edu/wiki/index.php/Feature_extraction_using_convolution*

Convolution operation captures the local dependencies in the ordinal image

# Pooling Layer

- Pooling reduces the dimensionality of each feature map
  - Max pooling: reports the maximum output within a rectangular neighborhood.
  - Average pooling: reports the average output of a rectangular neighborhood.

MaxPool with 2X2 filter with stride of 2

| 1 | 3 | 5 | 3 |
|---|---|---|---|
| 4 | 2 | 3 | 1 |
| 3 | 1 | 1 | 3 |
| 0 | 1 | 0 | 4 |

Input Matrix

| 4 | 5 |
|---|---|
| 3 | 4 |

Output Matrix

# Fully Connected Layers



A simple ConvNet
http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/

# Convolution Layer

32x32x3 image -> preserve spatial structure

32 height
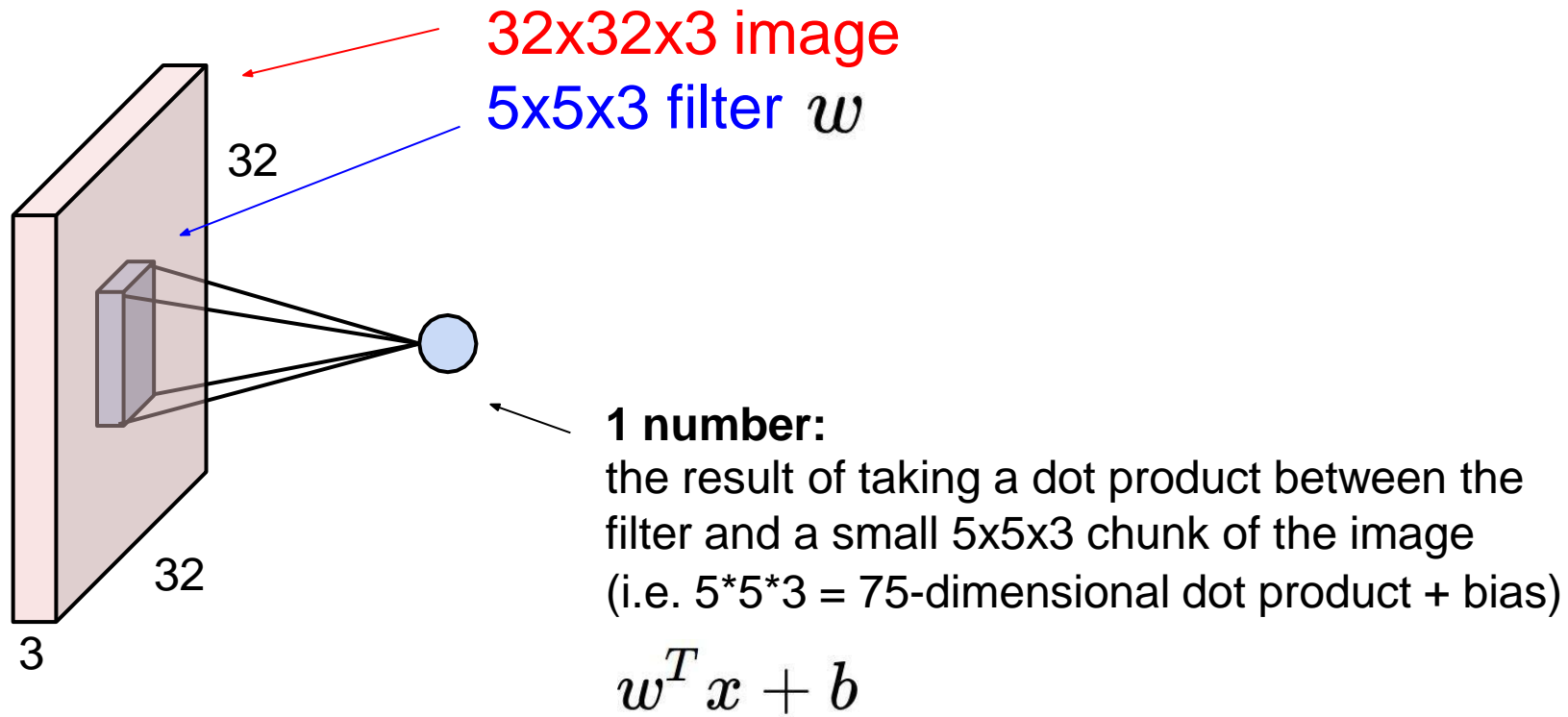
32 width

3 depth

# Convolution Layer

32x32x3 image



32

32

3

5x5x3 filter



**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"
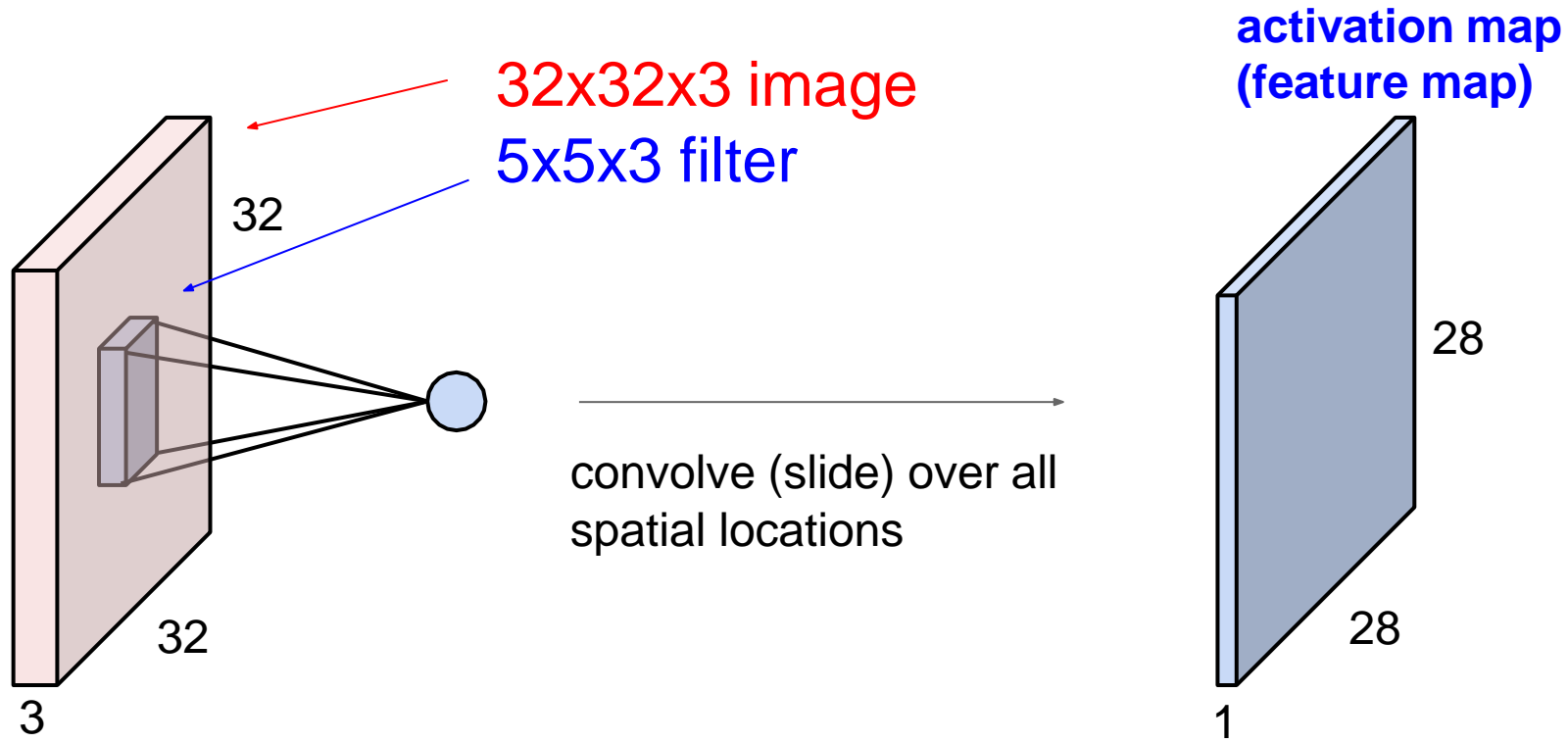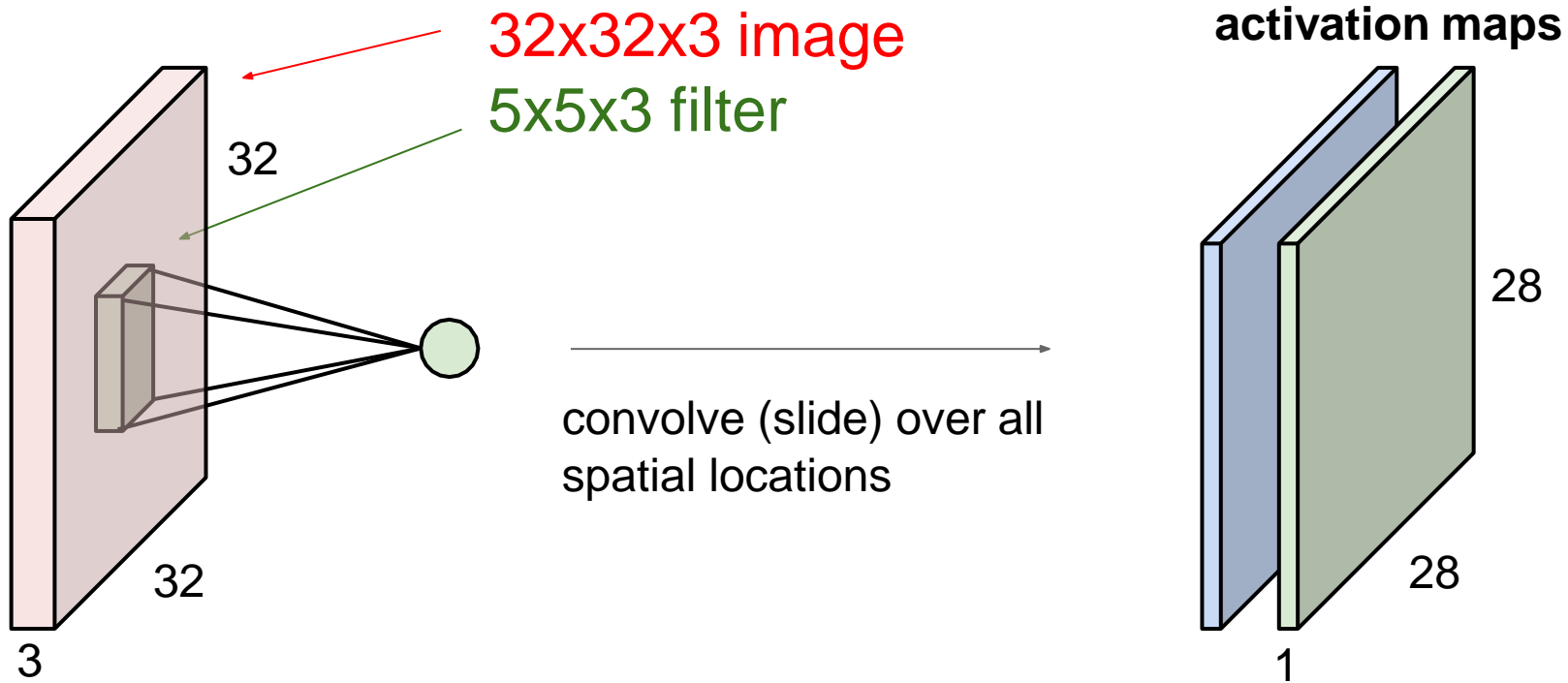
# Convolution Layer

32x32x**3** image

Filters always extend the full depth of the input volume

5x5x**3** filter

32

32

3

**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"

# Convolution Layer



**32x32x3 image**

**5x5x3 filter** $w$

**1 number:**
the result of taking a dot product between the
filter and a small 5x5x3 chunk of the image
(i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution Layer



32x32x3 image

5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation map
(feature map)**

28

28

1

# Convolution Layer

consider a second, green filter

32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

**activation maps**

28

28

1

# Convolution Layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

**activation maps**

32

32

3

Convolution Layer

28

28

6

We stack these up to get a "new image" of size 28x28x6!

# Pooling Layer

- makes the representations smaller and more manageable
- operates over each activation map independently:



Convolution Layer

Pooling Layer

Assume pool with 2x2 filters and stride 2

# Convolutional Neural Network



A simple ConvNet
http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/

# Why CNN for Image?

[Zeiler, M. D., *ECCV 2014*]
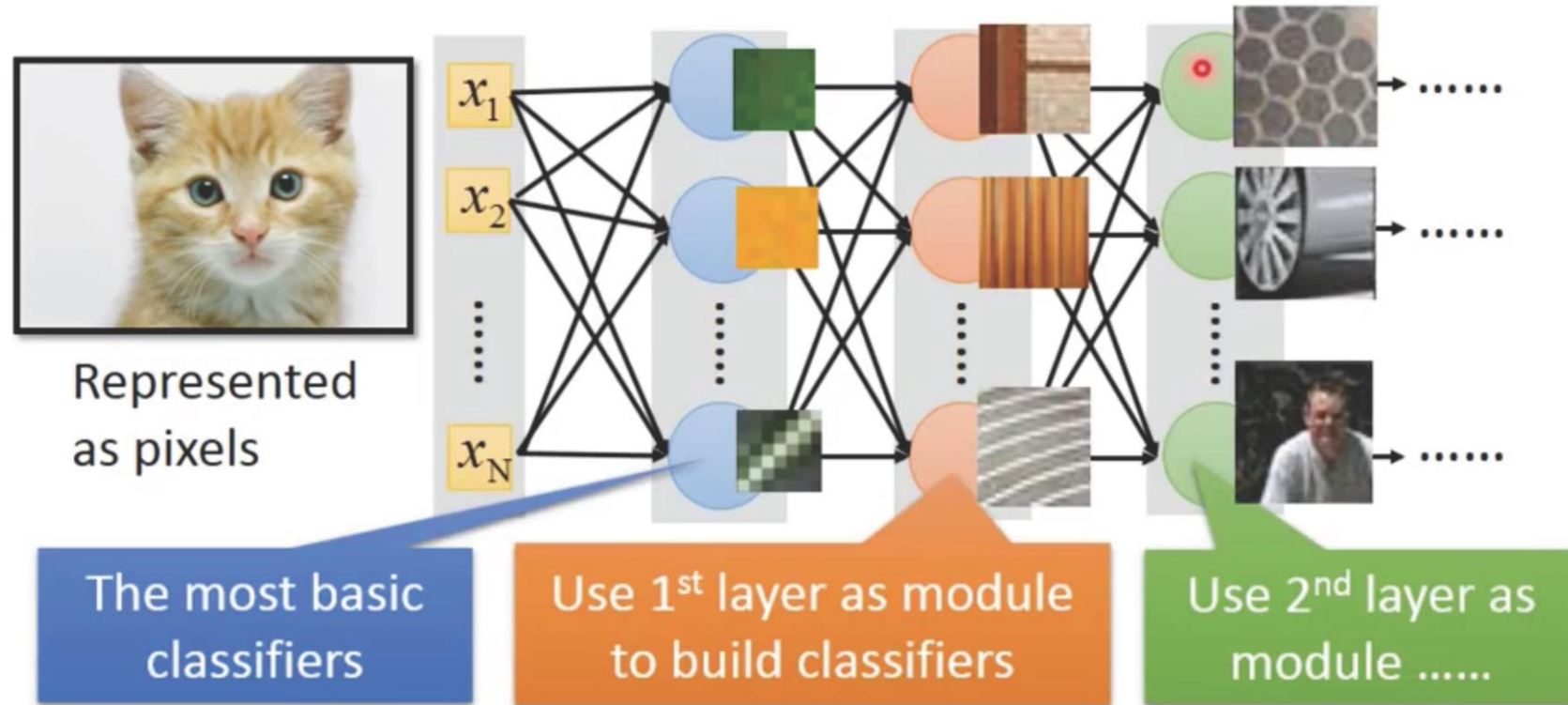
Represented as pixels

# Why CNN for Image?

[Zeiler, M. D., *ECCV 2014*]



Represented as pixels

$x_1$

$x_2$

$x_N$

The most basic classifiers

# Why CNN for Image?

Represented as pixels

The most basic classifiers

Use 1st layer as module to build classifiers

# Why CNN for Image?

[Zeiler, M. D., *ECCV 2014*]

Represented as pixels

$x_1$
$x_2$
$x_N$

The most basic classifiers

Use 1st layer as module to build classifiers

Use 2nd layer as module ......

Too many weights in a dense network!

Filter 1

6 x 6 image

6 x 6 image

Filter 1

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
⋮

3

Filter 1

6 x 6 image
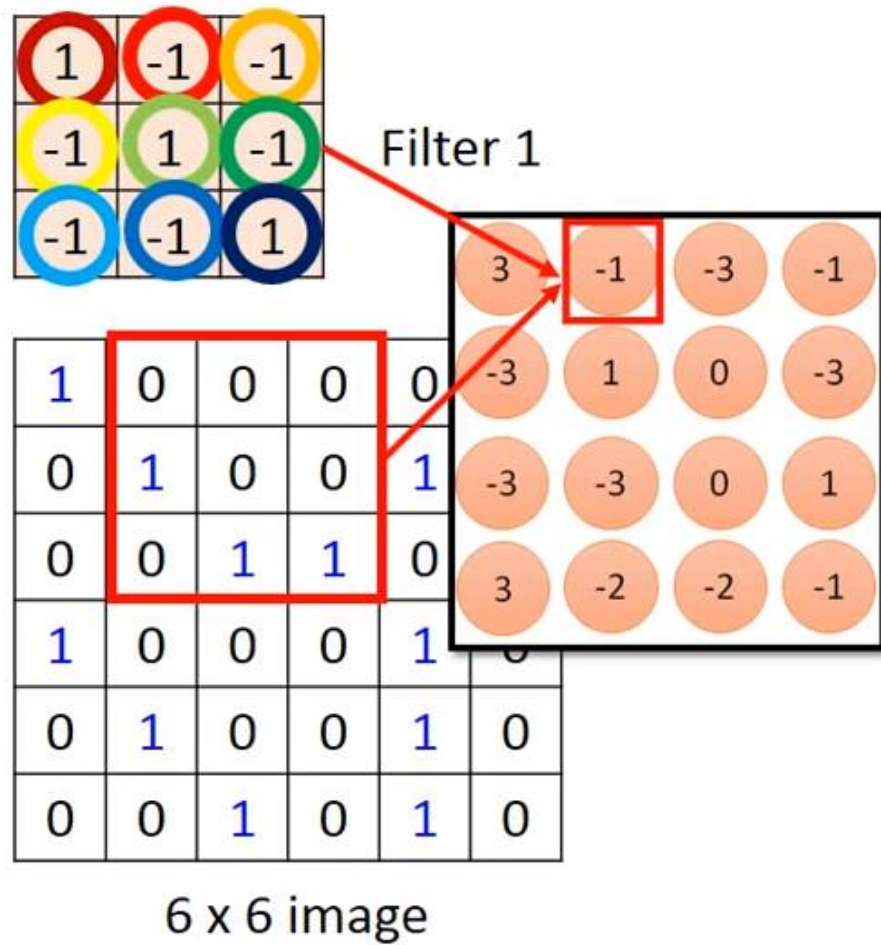
1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
14: 0
15: 1
16: 1
⋮

3

Filter 1

6 x 6 image

Only connect to 9 input, not fully connected

Filter 1

6 x 6 image

**Less parameters!**

1: 1
2: 0
3: 0
4: 0
⋮
7: 0
8: 1
9: 0
10: 0
⋮
13: 0
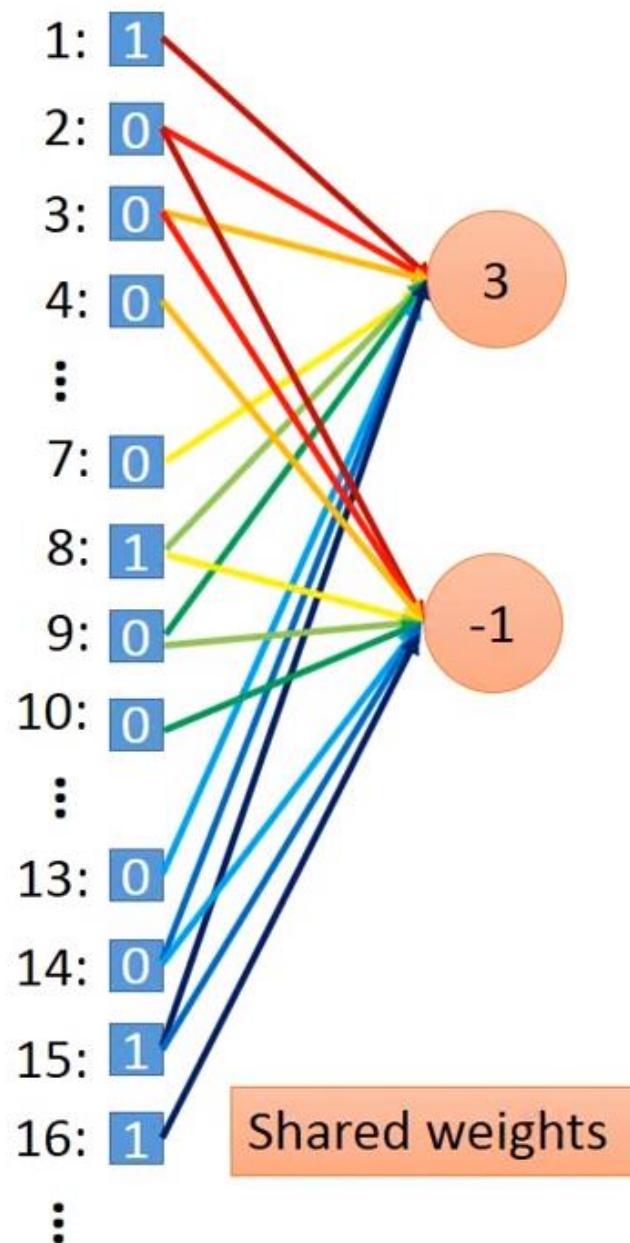14: 0
15: 1
16: 1
⋮

3

-1

Filter 1

6 x 6 image

Less parameters!

Even less parameters!

Shared weights

# The whole CNN