



2019 TJMSC Tech. Courses

Computer Vision

Bingchen Zhao
zhaobc@tongji.edu.cn
Tongji Microsoft Student Club

May, 26th
Room 516, Ji Shi Building
Tongji Univ

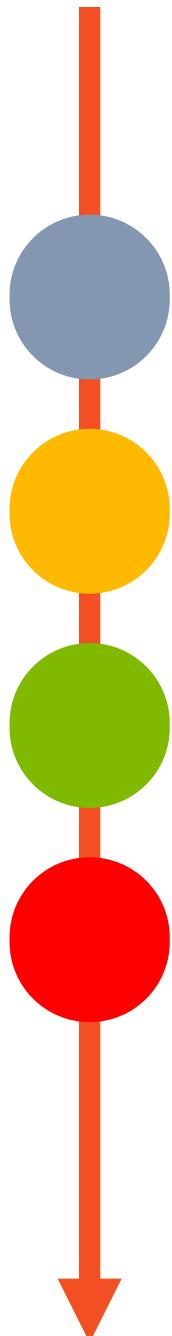


Image representation and classification

Start simple: Linear hypothesis

Non-linearity and feature composition

Interpretations of learned features

Classification: a core task in CV

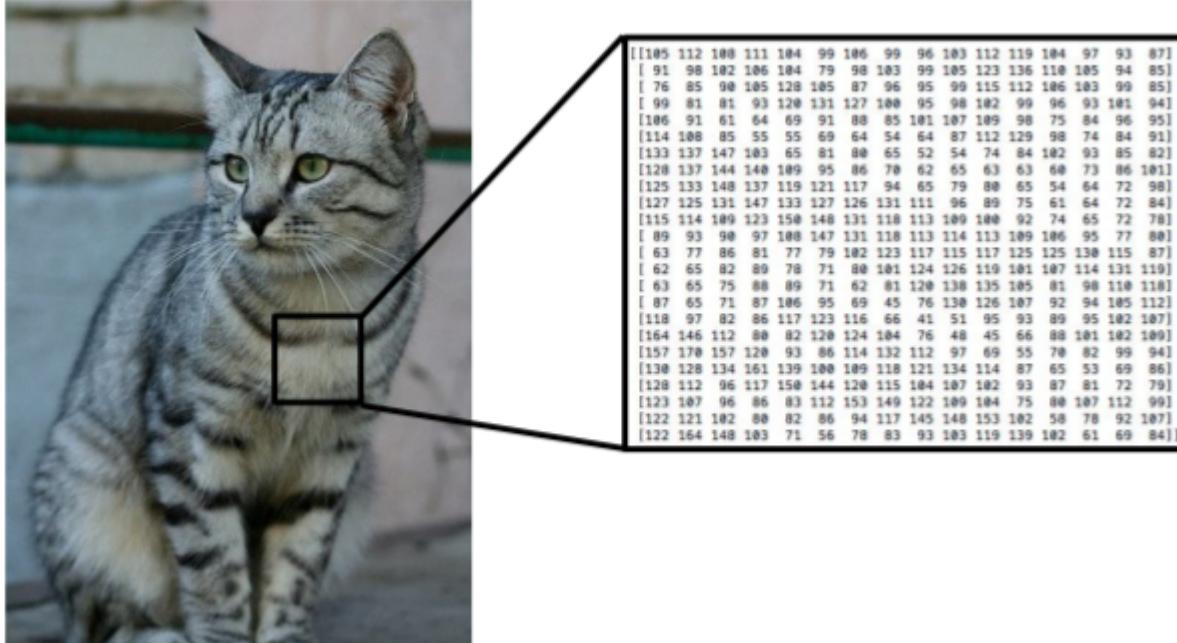


→Cat!

Maybe sounds easy for human,
but hard for computer.

```
def classify_image(image):  
    # Implement some good old dark magic here  
    return class_label
```

Challenges we are facing: semantic gap



We are seeing a cat, but what computer sees is just a big grid of numbers between [0,255]

Challenges : Background Clutter



[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

Challenges : Illumination



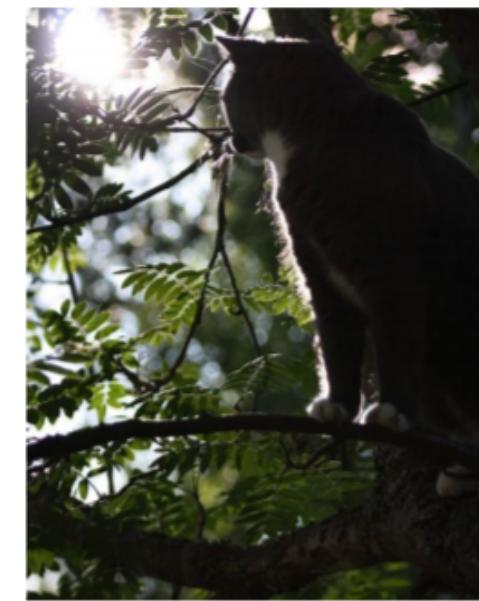
[This image is CC0 1.0 public domain](#)



[This image is CC0 1.0 public domain](#)

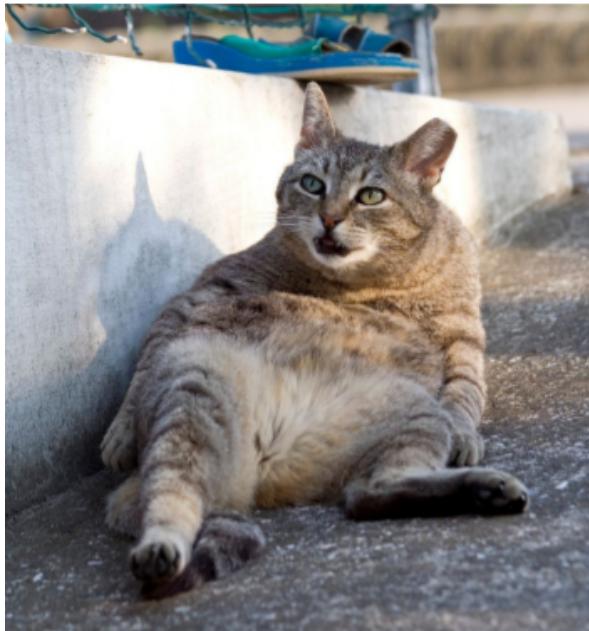


[This image is CC0 1.0 public domain](#)

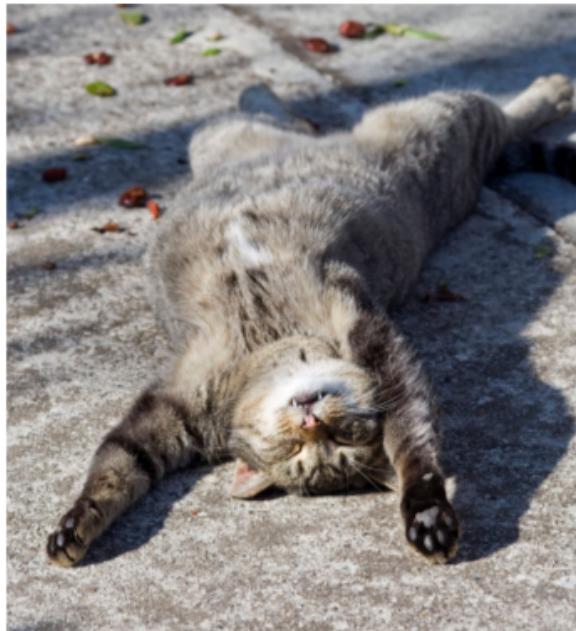


[This image is CC0 1.0 public domain](#)

Challenges : Rare pose



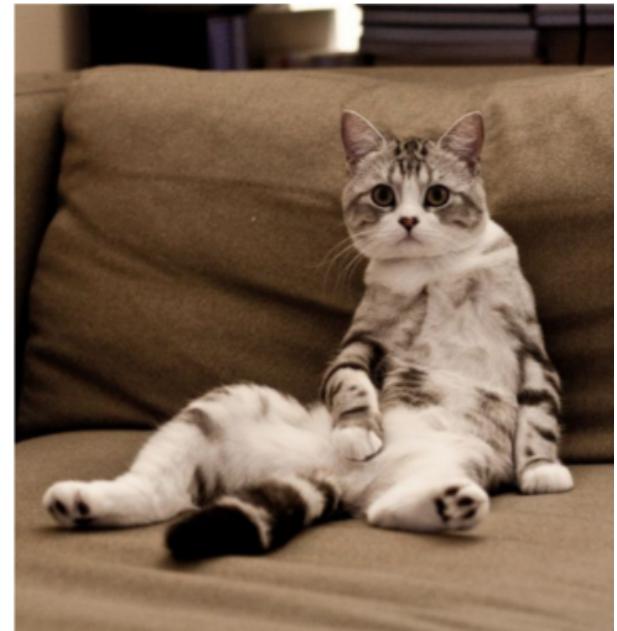
[This image by Umberto Salvagnin](#)
is licensed under [CC-BY 2.0](#)



[This image by Umberto Salvagnin](#)
is licensed under [CC-BY 2.0](#)



[This image by sare bear](#) is
licensed under [CC-BY 2.0](#)



[This image by Tom Thai](#) is
licensed under [CC-BY 2.0](#)

Challenges : Occlusion



[This image](#) is [CC0 1.0](#) public domain



[This image](#) is [CC0 1.0](#) public domain



[This image](#) by [jonsson](#) is licensed
under [CC-BY 2.0](#)

Start simple: CIFAR-10

airplane



automobile



bird



cat



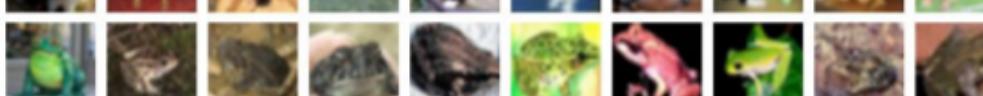
deer



dog



frog



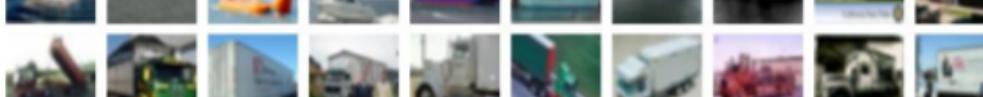
horse



ship



truck

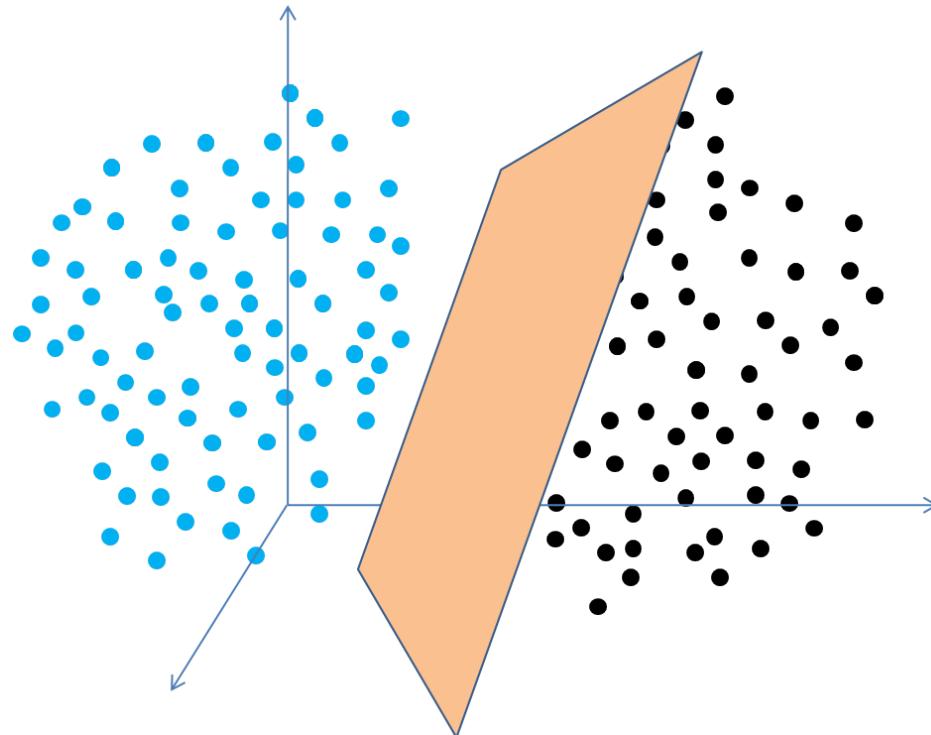


50,000 training images
each image is **32x32x3**

10,000 test images.

Start simple: Linear hypothesis

Linear hypothesis is that we suppose the data is linear separable



In three dimension, linear separable mean a plane



An image in CIFAR-10 is a array of 32x32x3 numbers(3072 dimension space)

Start simple: Linear hypothesis

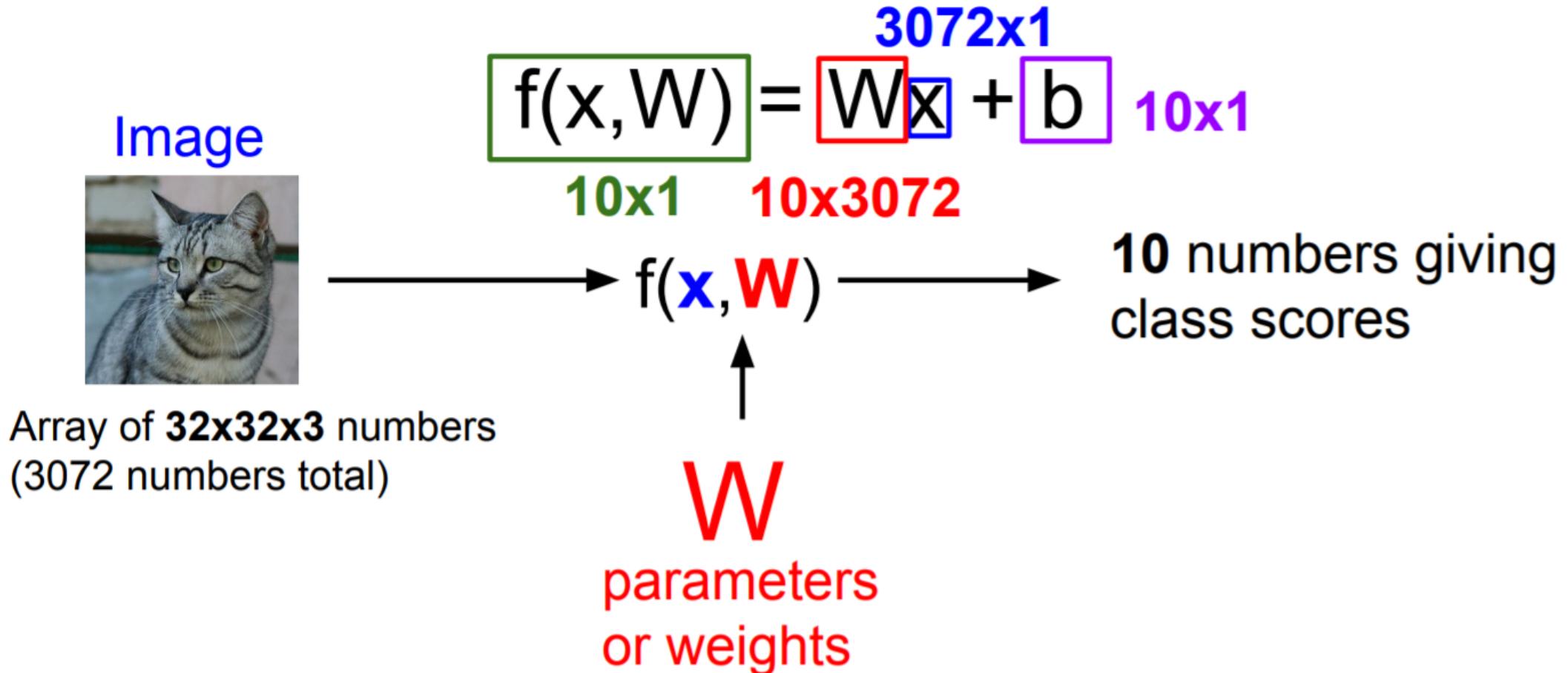
Suppose we already know the parameter of the hyper-plane(plane in high dimension), we can use it to define our magical function

$$f(x, W) = Wx + b$$

3072x1
10x1 **10x3072**

```
def classify_image(image, w, b):  
    # params:  
    #   image is a flattened 3072 array  
    #   w and b is our hyper plane parameters  
    class_label = image * w + b  
    # class_label is predicted class scores  
    return class_label
```

Start simple: Linear hypothesis



Start simple: Linear hypothesis

How do we find the parameters for our classifier?

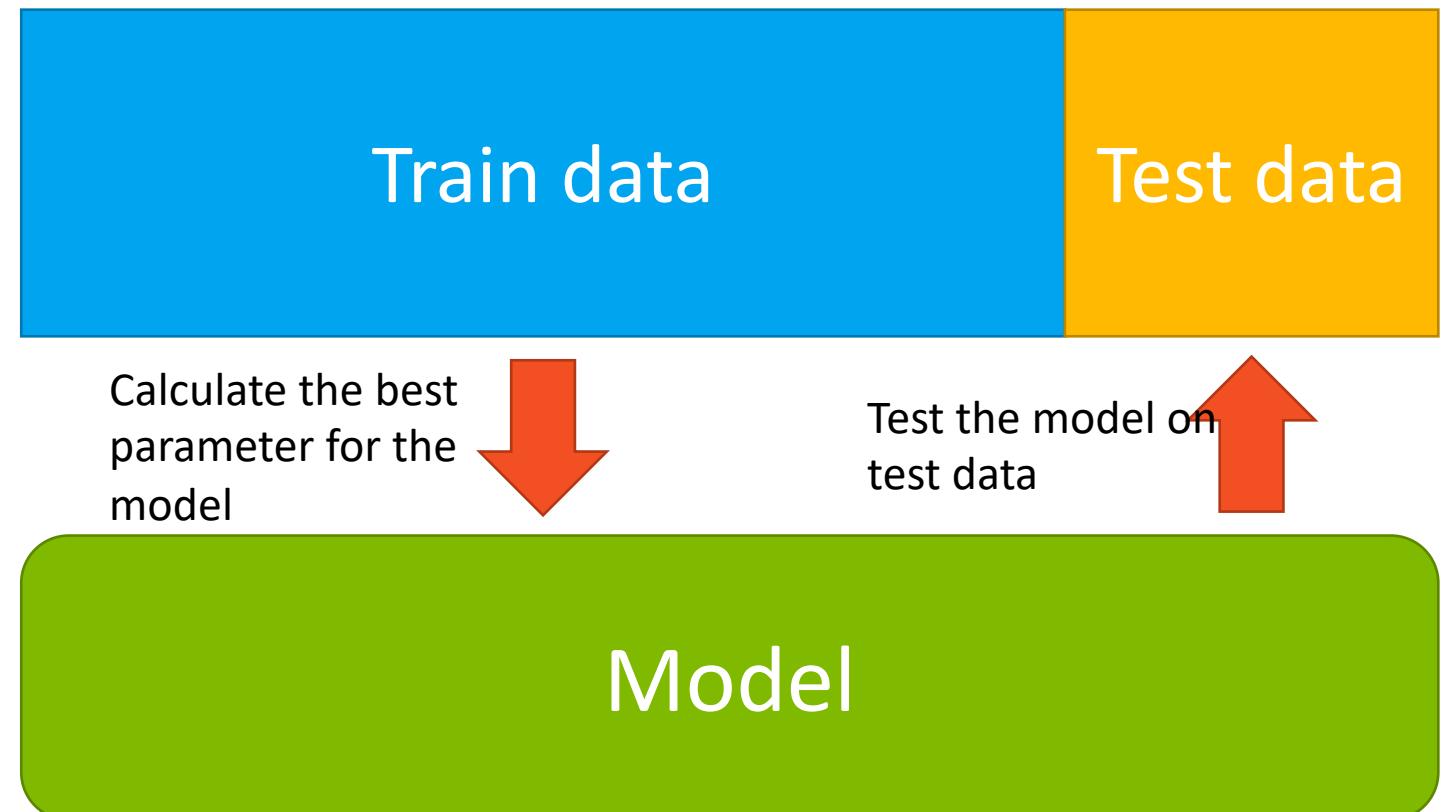
Objective is to find a W minimize the following equation

$$\text{argmin}_W \|f(X^{(train)}; W) - y^{(train)}\|_2$$

Least square method gives us

$$W = (X^{(train)T} X^{(train)})^{-1} X^{(train)T} y^{(train)}$$

But this is infeasible cause this need to solve a inverse problem, (sometime you can't)



Start simple: Linear hypothesis

Stochastic Gradient Descent(SGD)

Define the loss function as follows

$$L(W; X, y) = \|f(X; W) - y\|_2$$

Update the W according to the gradient

$$W_{n+1} = W_n - \epsilon \frac{dL}{dW}$$



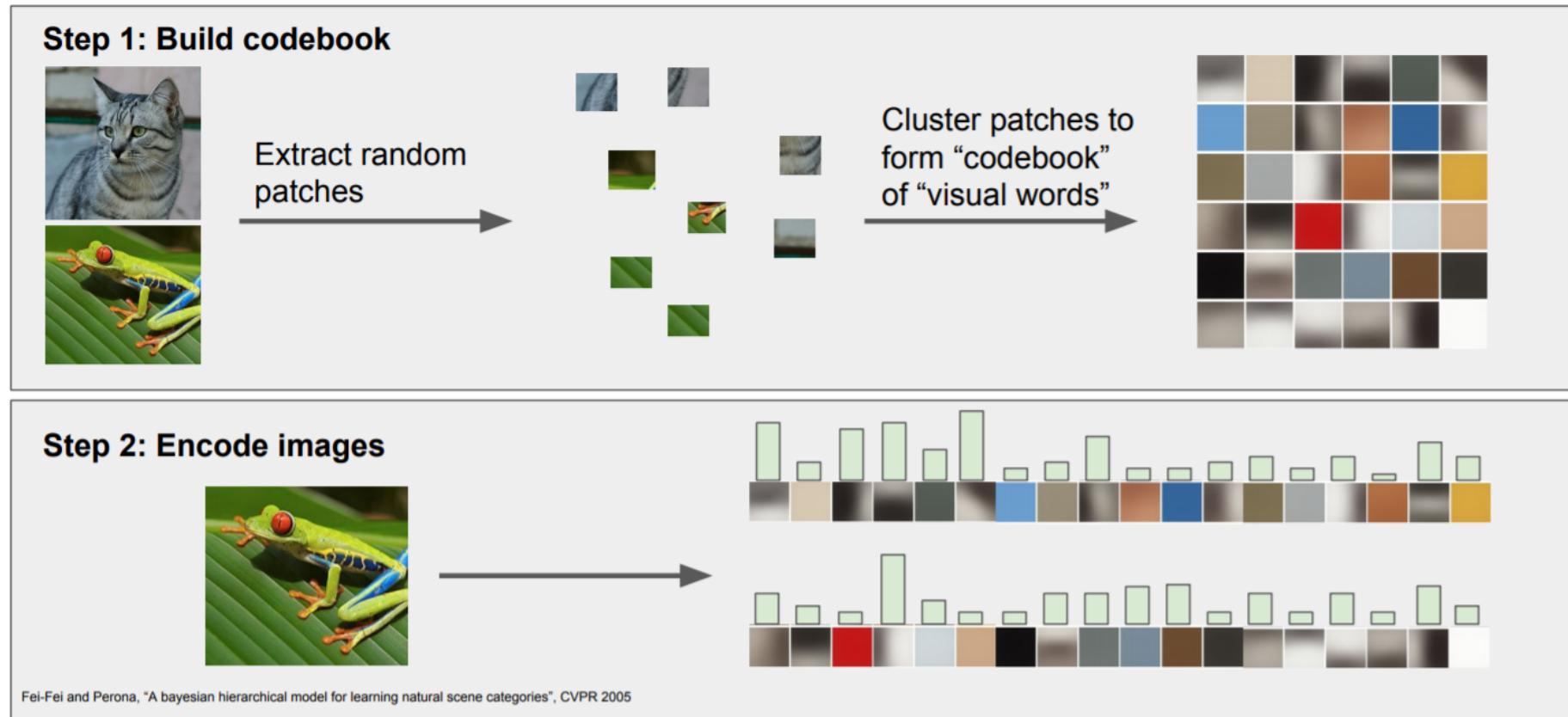
Start simple: Linear hypothesis

Time to put hands on

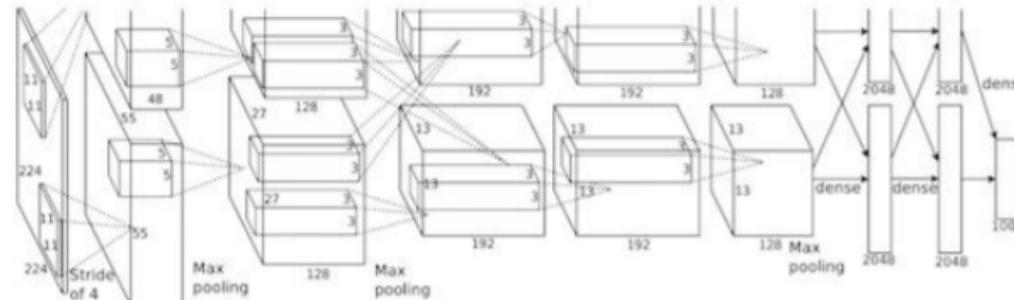
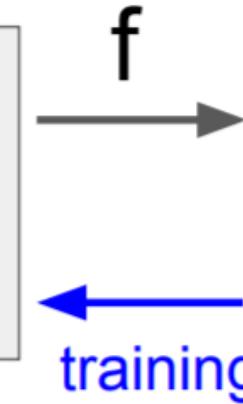
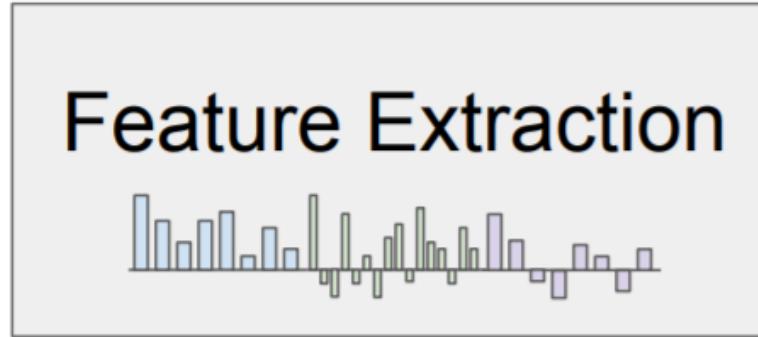
All code and notebook is on Github

Improve the model: feature composition

Previously we are treating every pixel as independent features, what about group nearby pixel together to get more complex features?



Improve the model: ConvNet



Krizhevsky, Sutskever, and Hinton, "Imagenet classification with deep convolutional neural networks", NIPS 2012.
Figure copyright Krizhevsky, Sutskever, and Hinton, 2012.
Reproduced with permission.

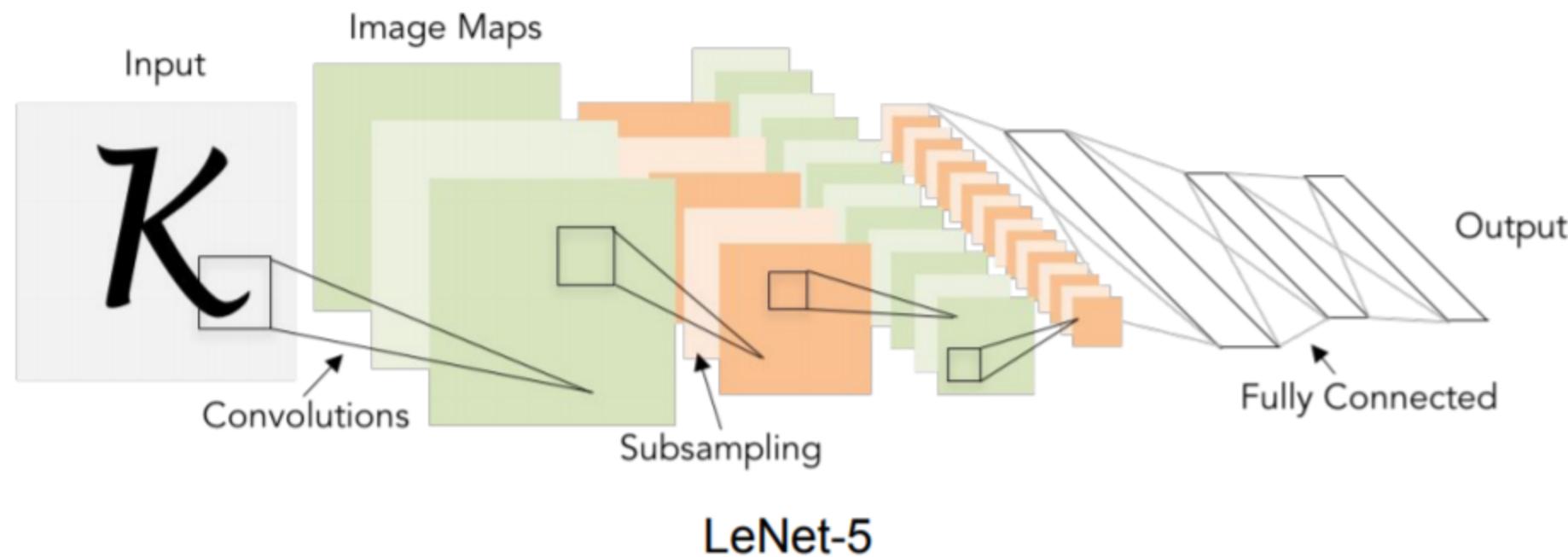


10 numbers giving scores for classes



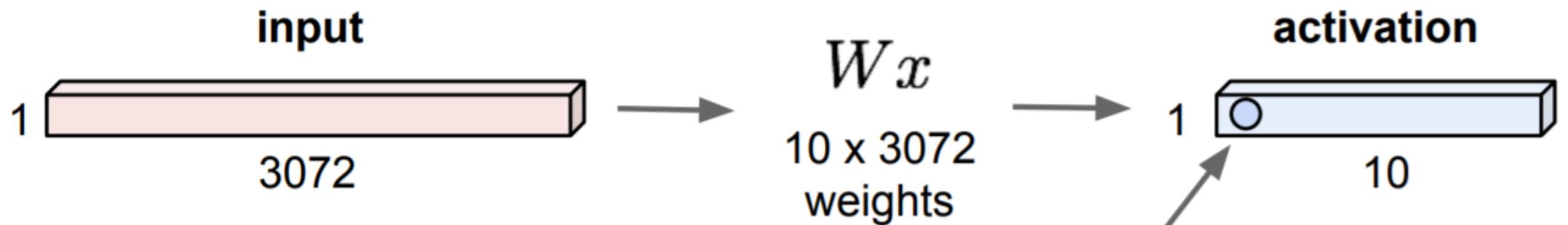
ConvNet

Using Convolutional layer to let the network learn the feature composition, ReLU layer can be used to add non-linearity.



ConvNet

32x32x3 image -> stretch to 3072 x 1

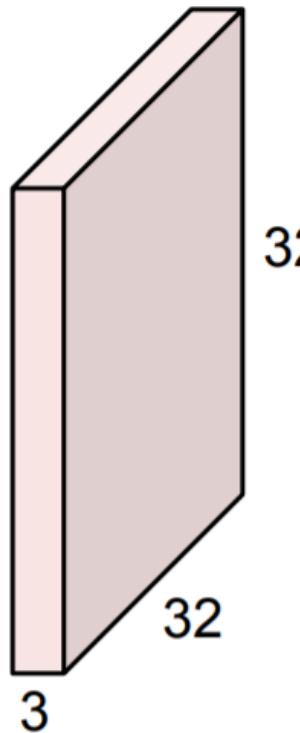


1 number:

the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

ConvNet

32x32x3 image

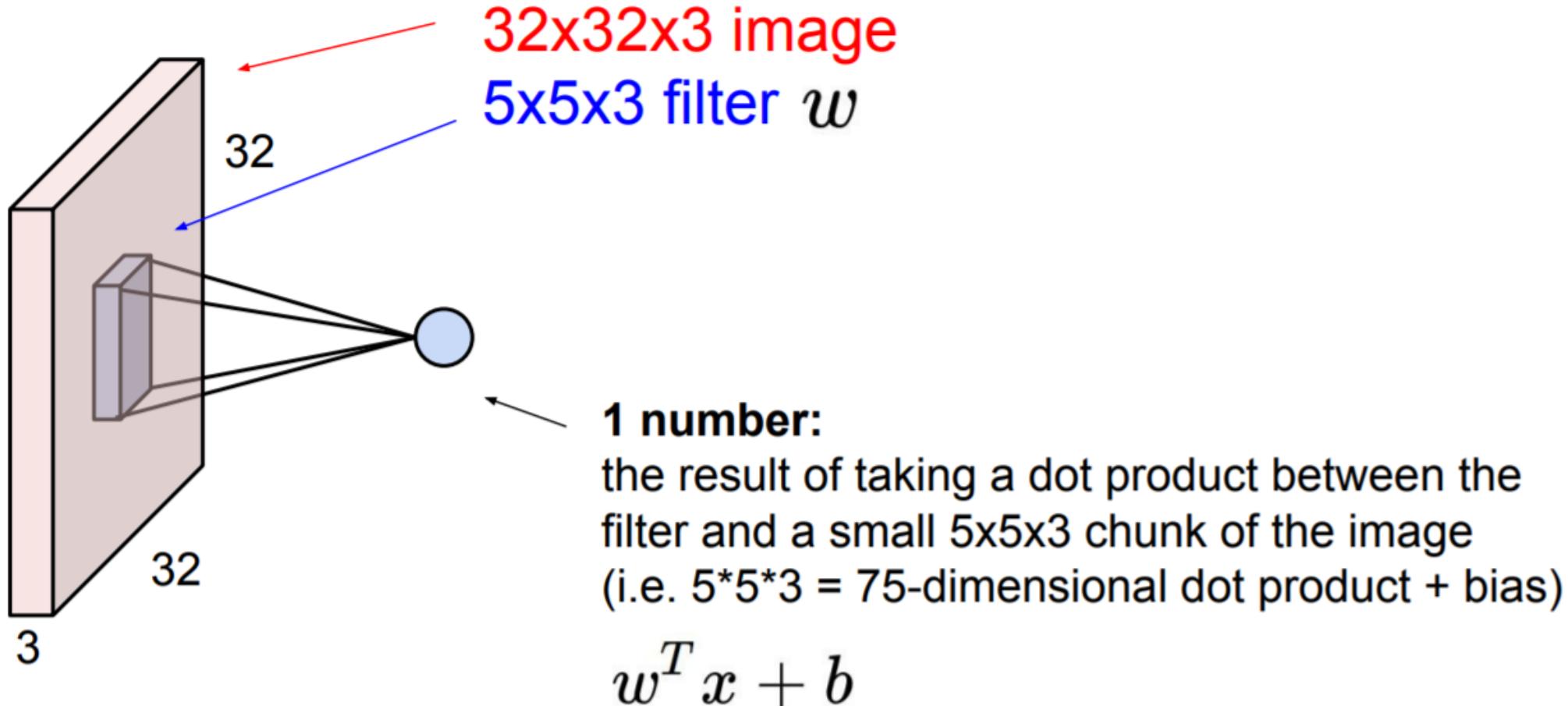


5x5x3 filter

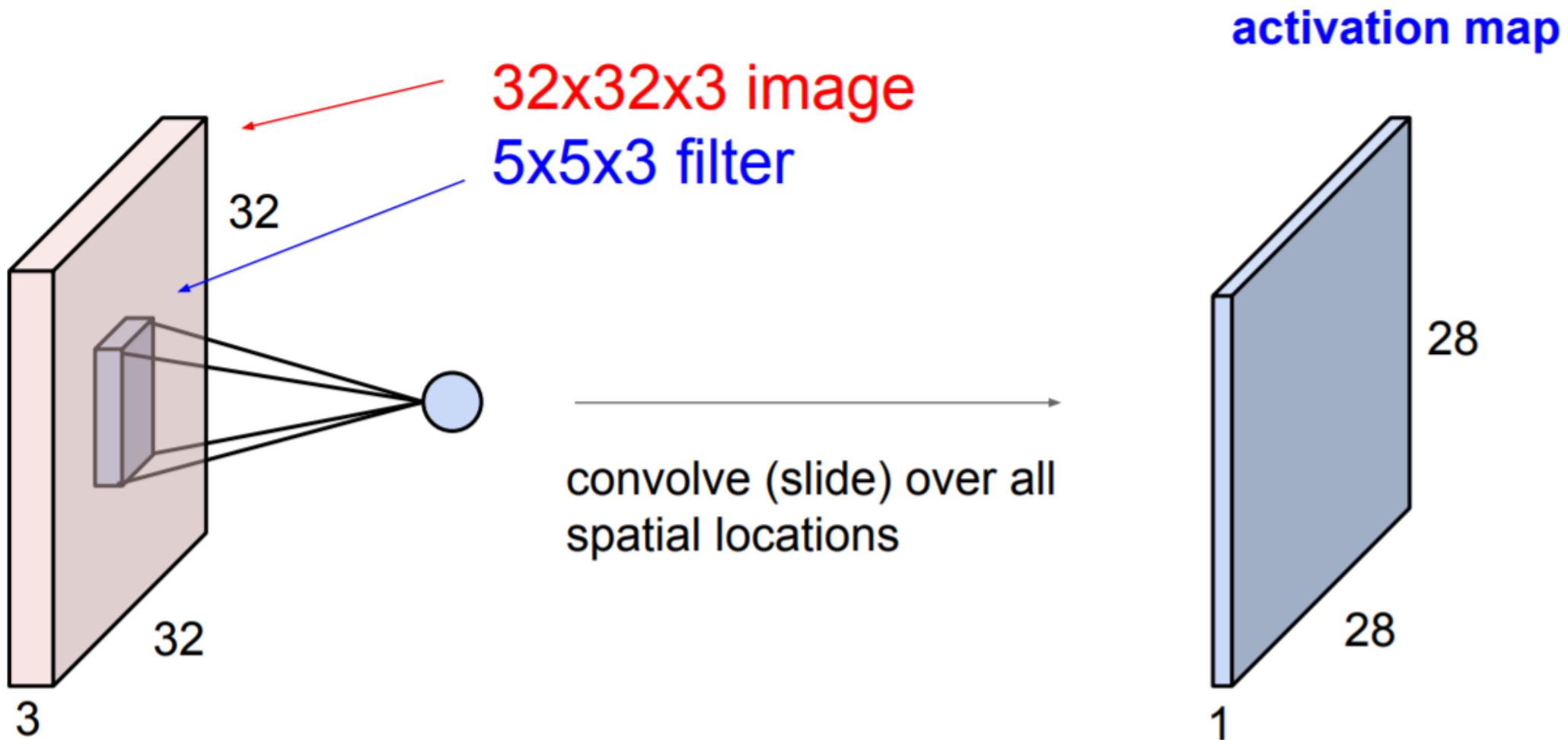


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

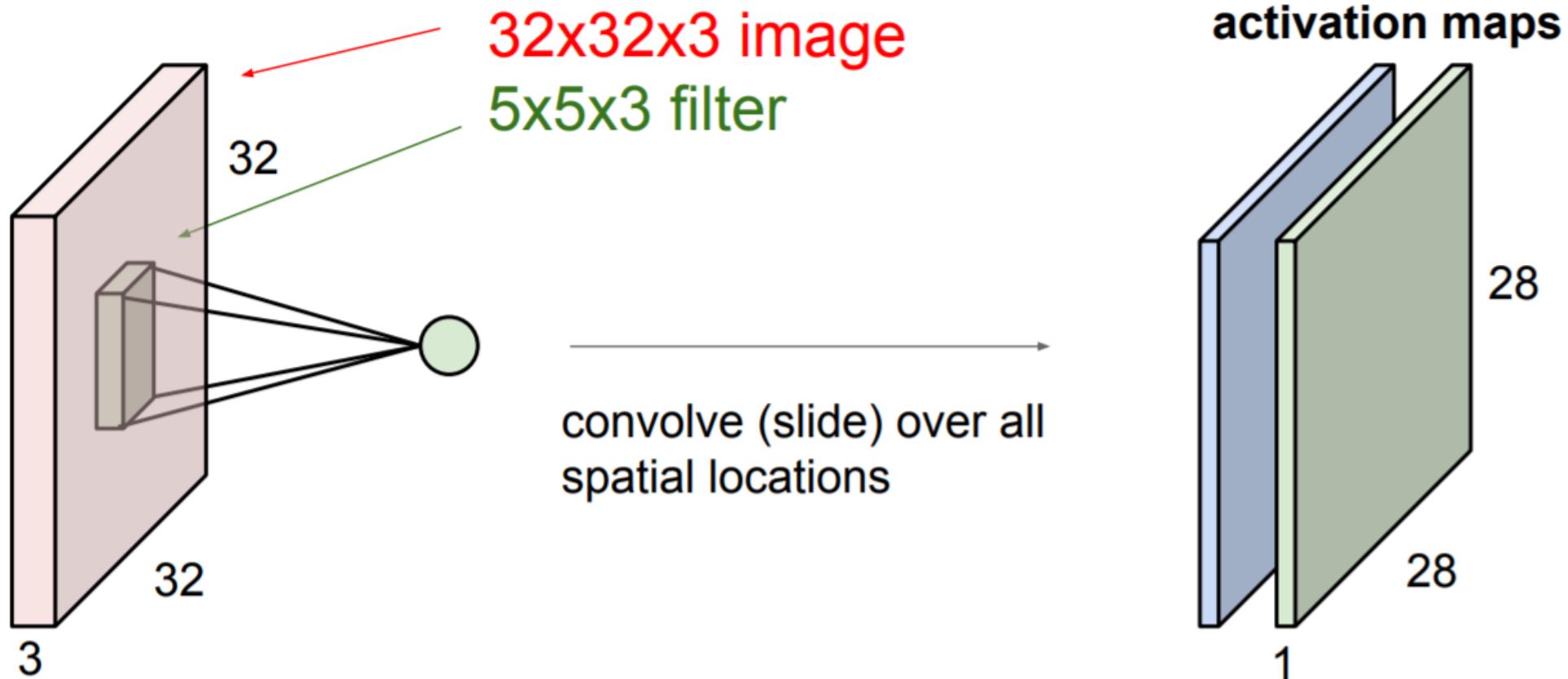
ConvNet



ConvNet

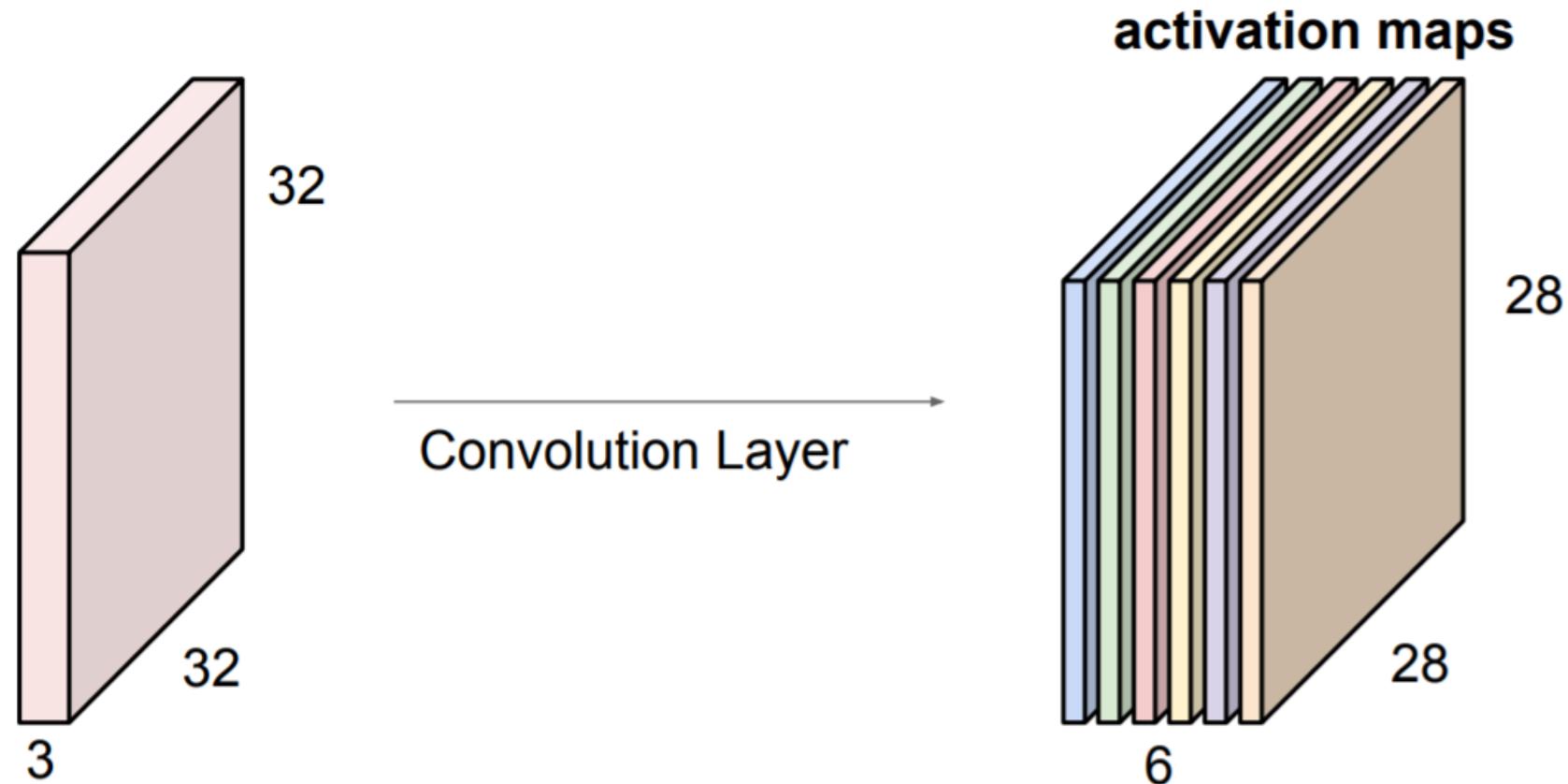


ConvNet



ConvNet

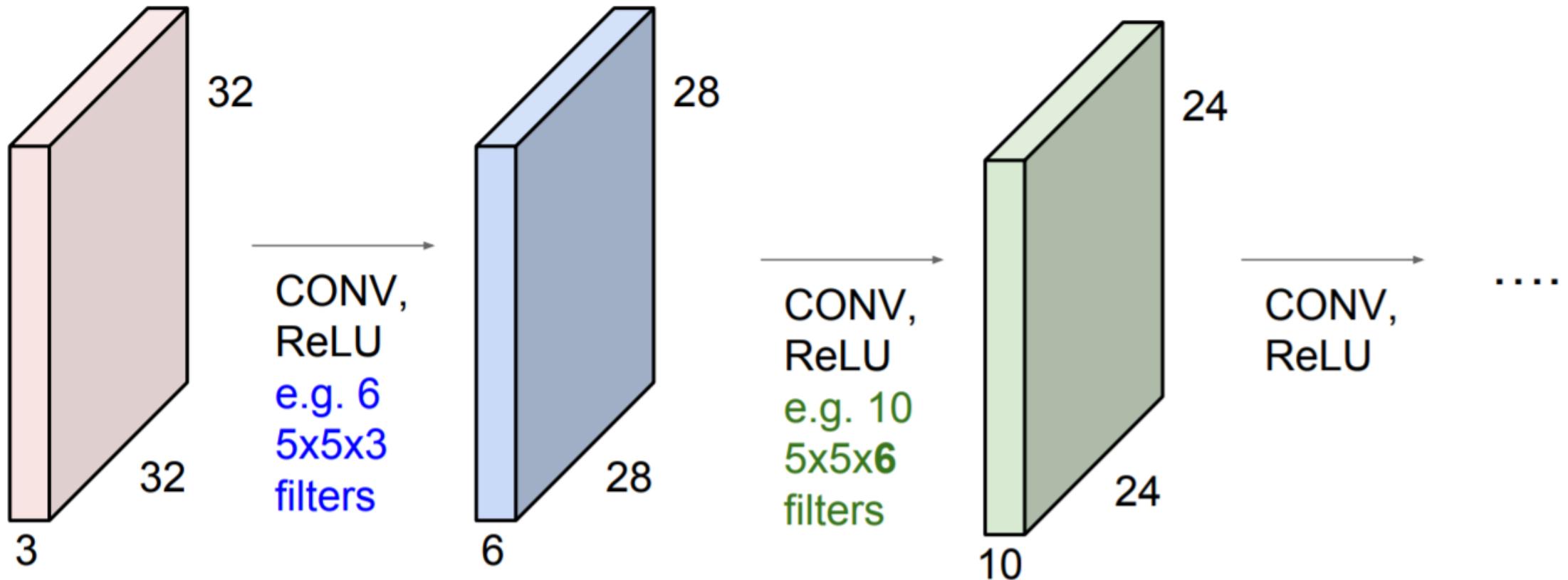
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



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

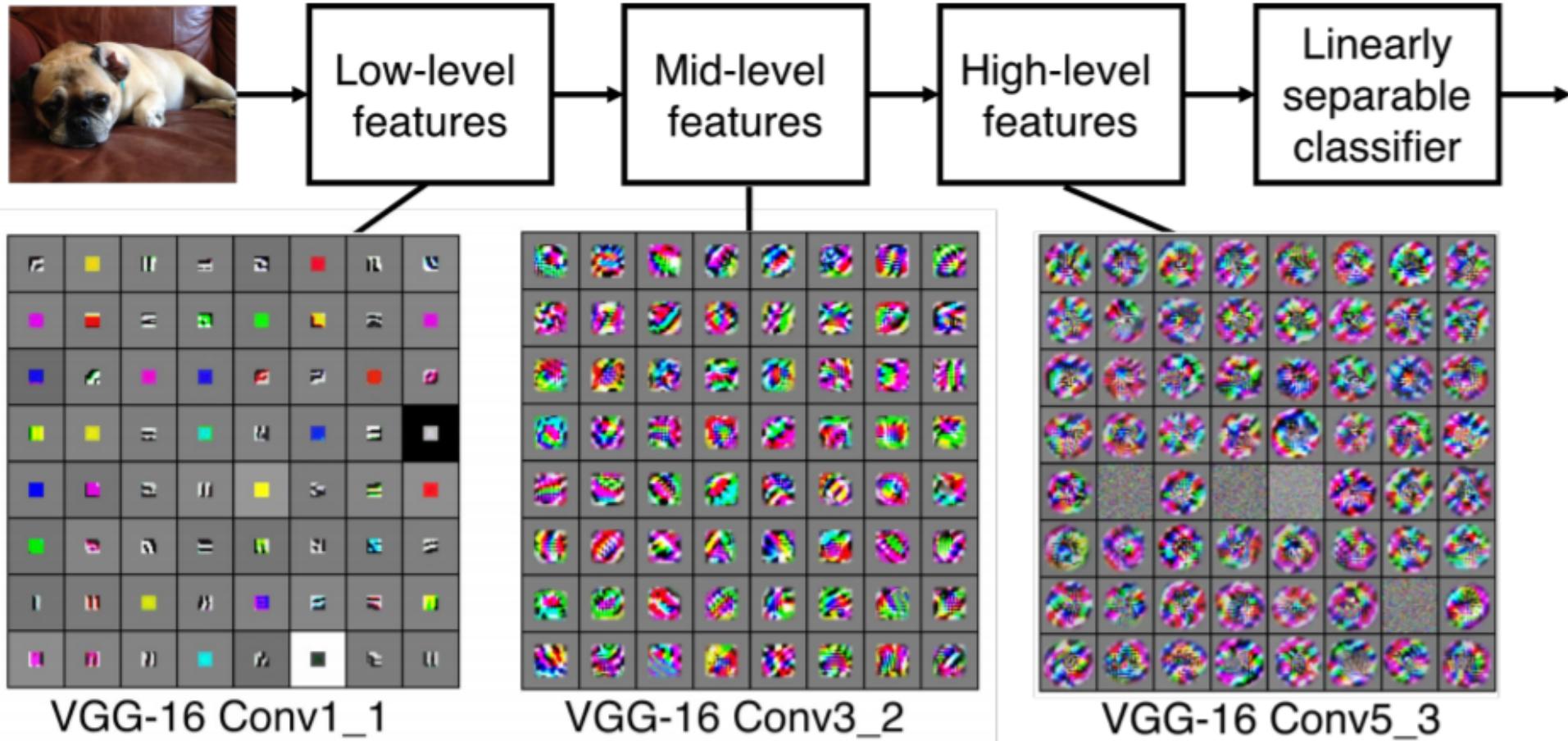
ConvNet

ConvNet is a sequence of Convolution Layers, interspersed with activation functions(ReLU which set value < 0 to 0)



ConvNet

Transfer learning



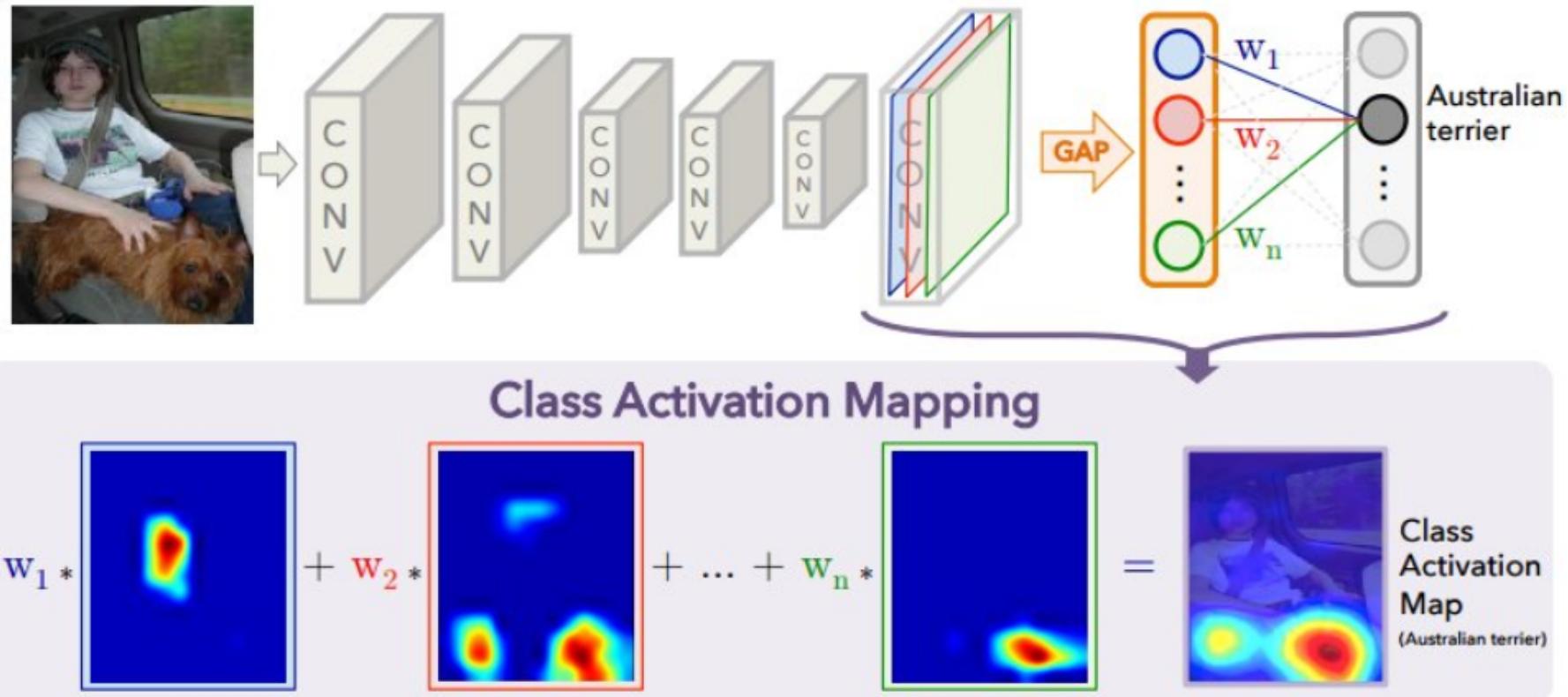
ConvNet : CIFAR10 Code walkthrough

```
Net(  
    (conv1): Conv2d(3, 6, kernel_size=(5, 5), stride=(2, 2))  
    (conv2): Conv2d(6, 16, kernel_size=(5, 5), stride=(2, 2))  
    (fc1): Linear(in_features=400, out_features=120, bias=True)  
    (fc2): Linear(in_features=120, out_features=84, bias=True)  
    (fc3): Linear(in_features=84, out_features=10, bias=True)  
)
```

ConvNet : dog-vs-cat Code walkthrough

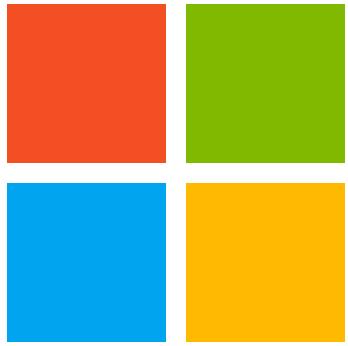


ConvNet : Interpretations



ConvNet : Interpretations





Microsoft



加入TJMSC QQ群



关注TJMSC微信公众号