

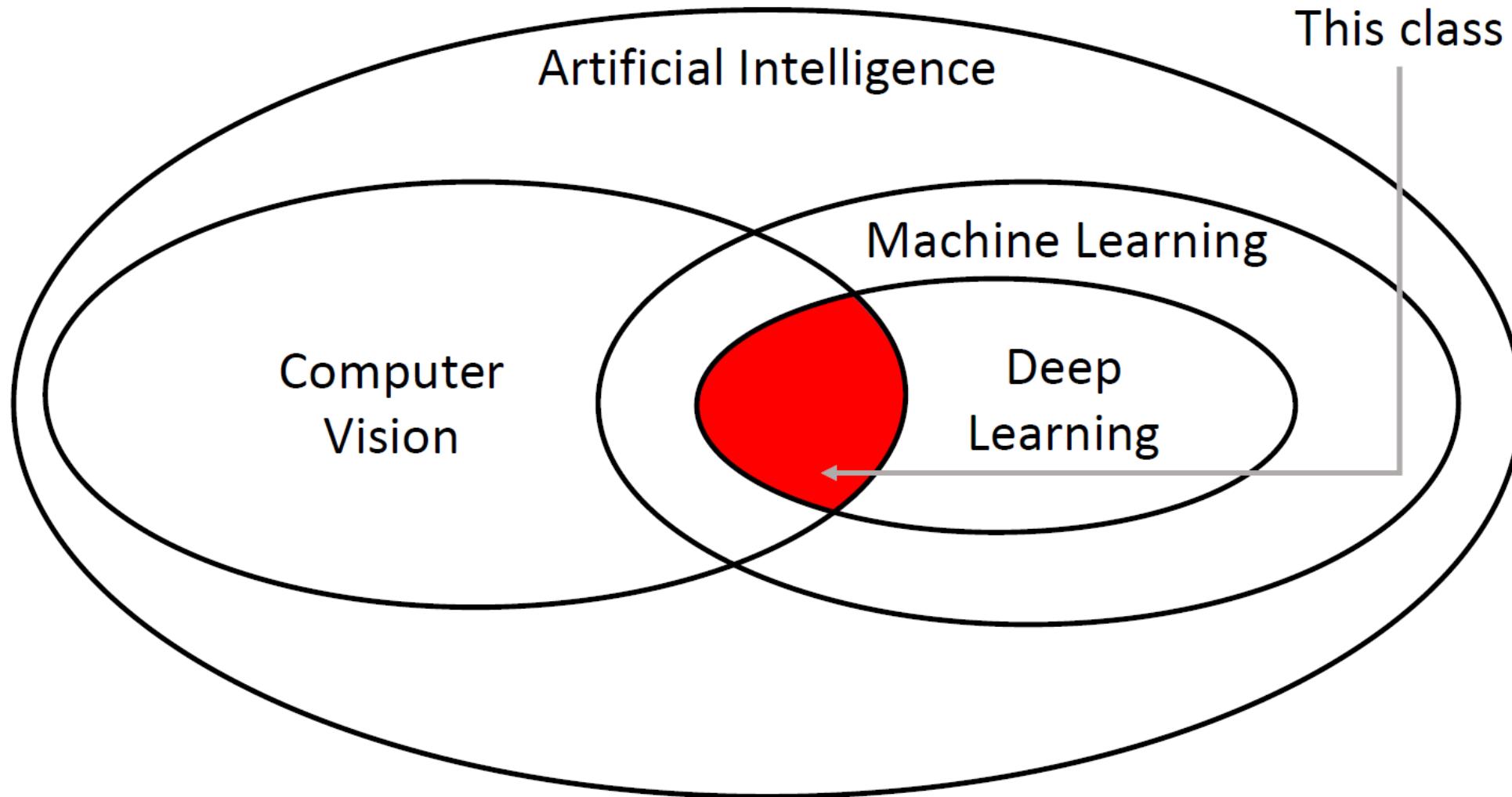


CS7GV1 Computer vision

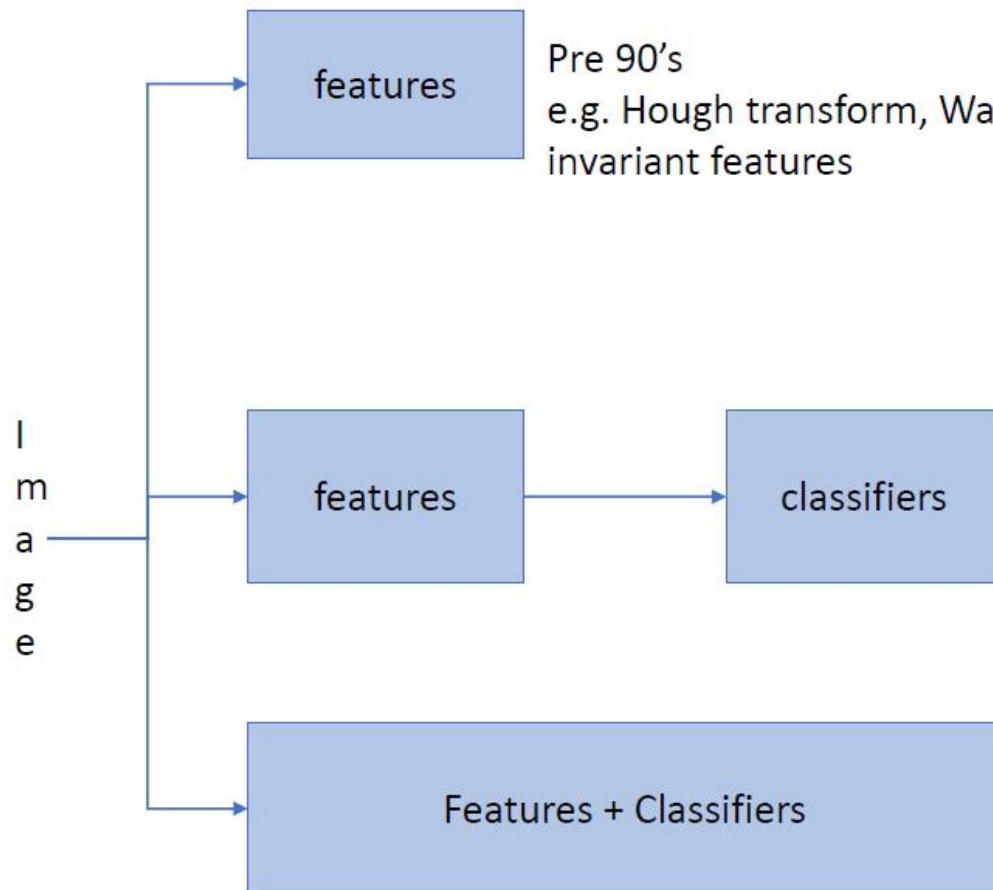
Convolutional Neural Networks

Dr. Martin Alain

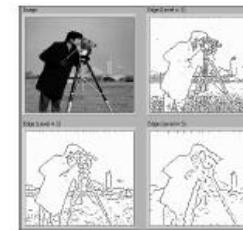
Introduction



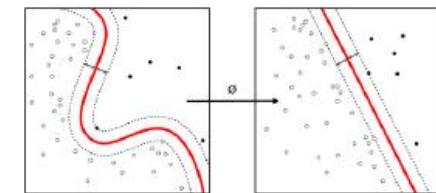
Computer vision and machine learning



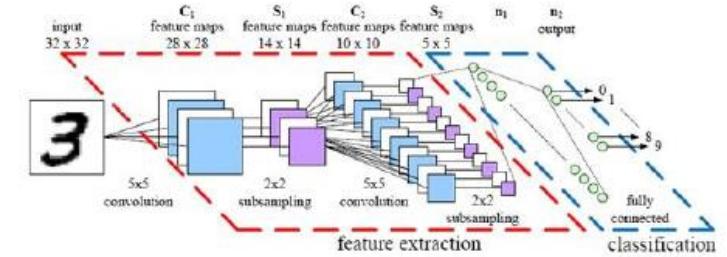
Pre 90's
e.g. Hough transform, Wavelets,
invariant features



~ 2000's e.g.
Adaboost (**Viola&Jones**),
Support Vector Machines (SVM)
Random Forest



~ 2010's e.g.
Deep Learning (CNN)



Machine learning

Background

A Recipe for Machine Learning

1. Given training data:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

2. Choose each of these:

– Decision function

$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

– Loss function

$$\ell(\hat{\mathbf{y}}, \mathbf{y}_i) \in \mathbb{R}$$

3. Define goal:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^N \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

4. Train with SGD:

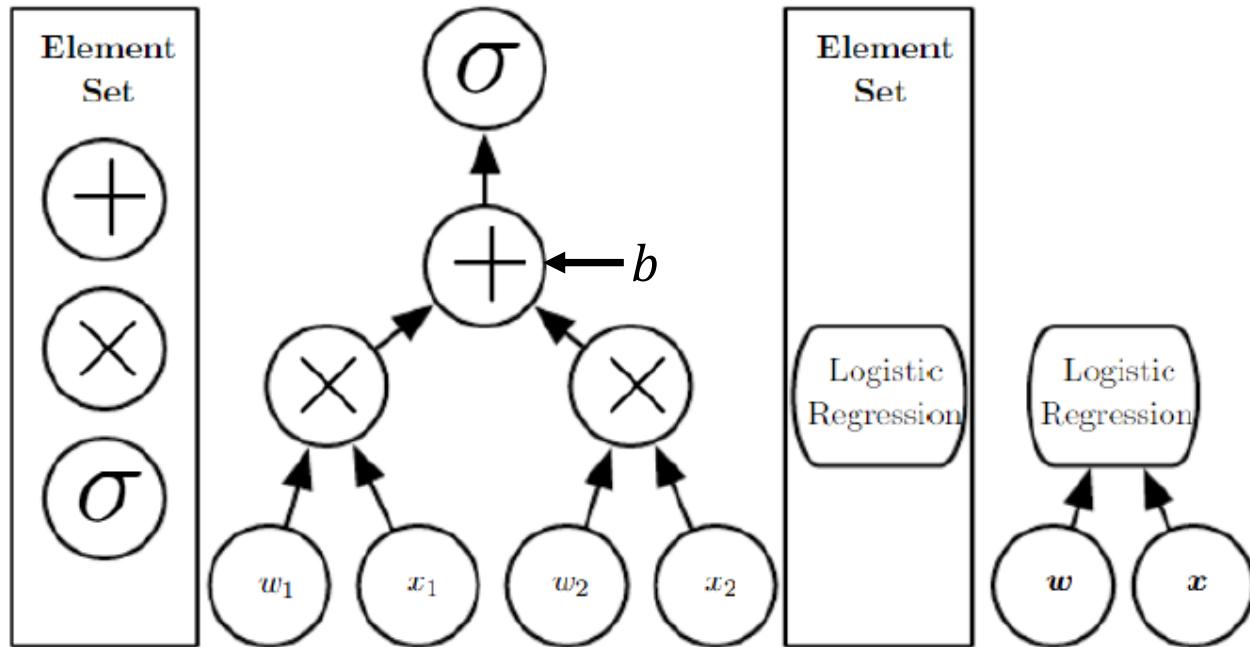
(take small steps
opposite the gradient)

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\mathbf{x}_i), \mathbf{y}_i)$$

Source:

[https://www.cs.cmu.edu/~mgormley/
courses/10601-s17/slides/lecture20-
backprop.pdf](https://www.cs.cmu.edu/~mgormley/courses/10601-s17/slides/lecture20-backprop.pdf)

Neuron



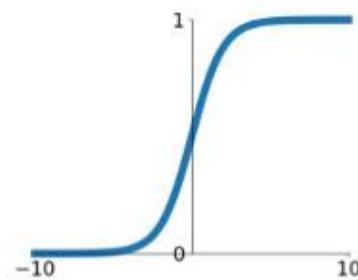
σ is an activation function chosen here as the logit function but other can be chosen as part of a network

Figure 1.3: Illustration of computational graphs mapping an input to an output where each node performs an operation. Depth is the length of the longest path from input to output but depends on the definition of what constitutes a possible computational step. The computation depicted in these graphs is the output of a logistic regression model, $\sigma(\mathbf{w}^T \mathbf{x})$, where σ is the logistic sigmoid function. If we use addition, multiplication and logistic sigmoids as the elements of our computer language, then this model has depth three. If we view logistic regression as an element itself, then this model has depth one.

Activation function

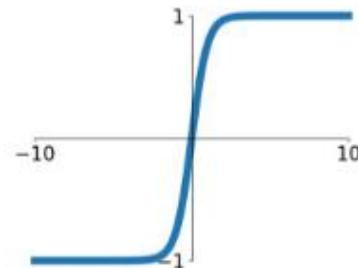
Sigmoid

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



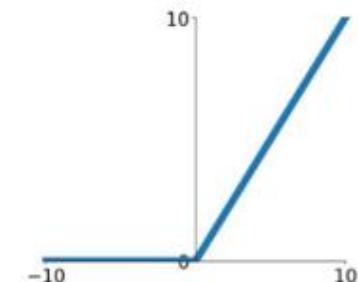
tanh

$$\tanh(x)$$



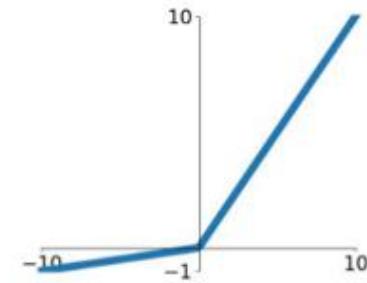
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

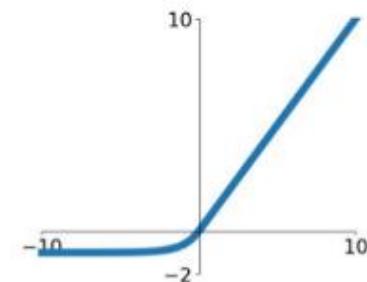


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

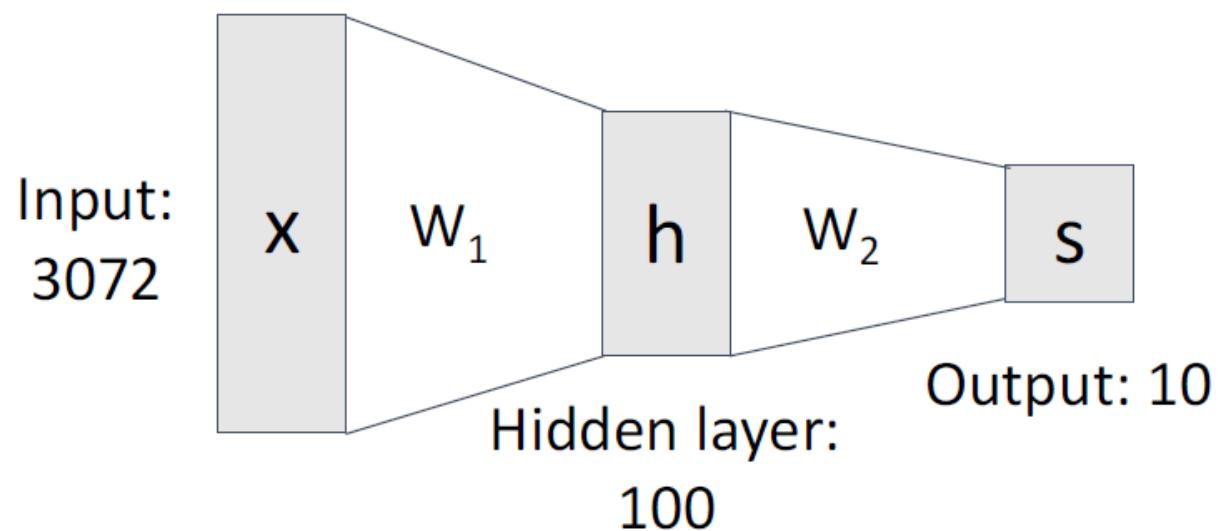
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



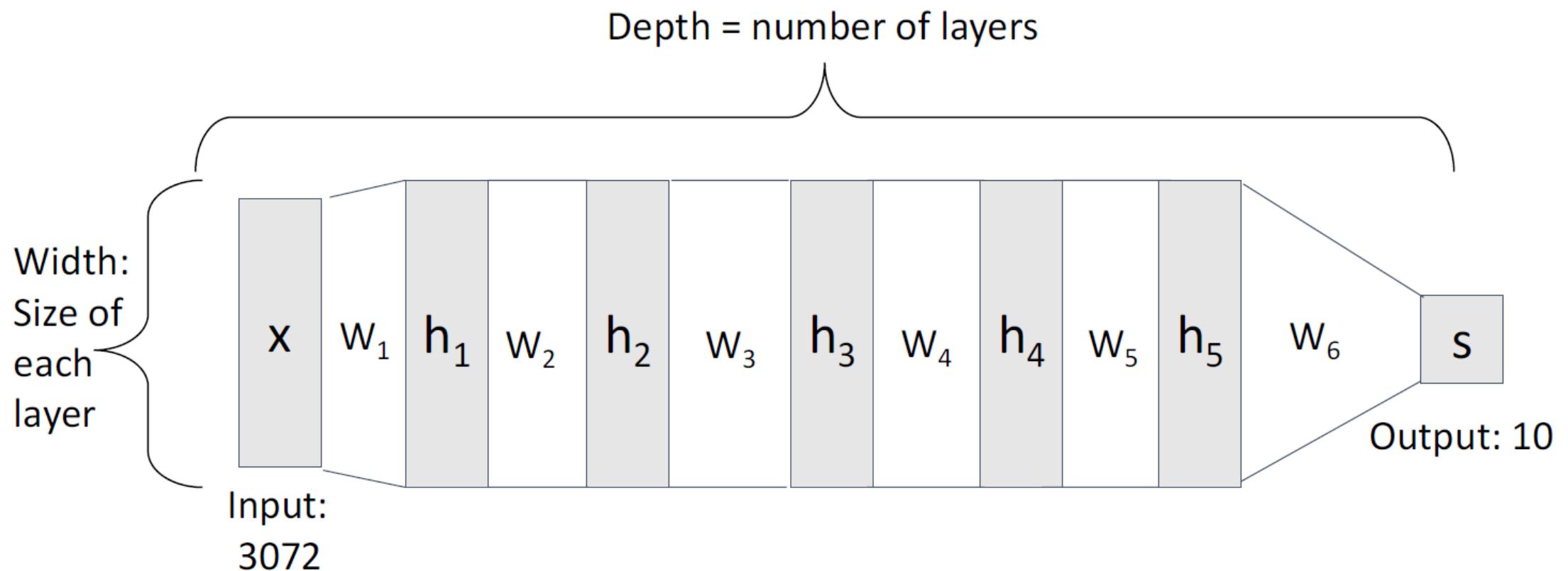
Neural Networks

- Fully connected neural network also called Multi-layer perceptron (MLP)

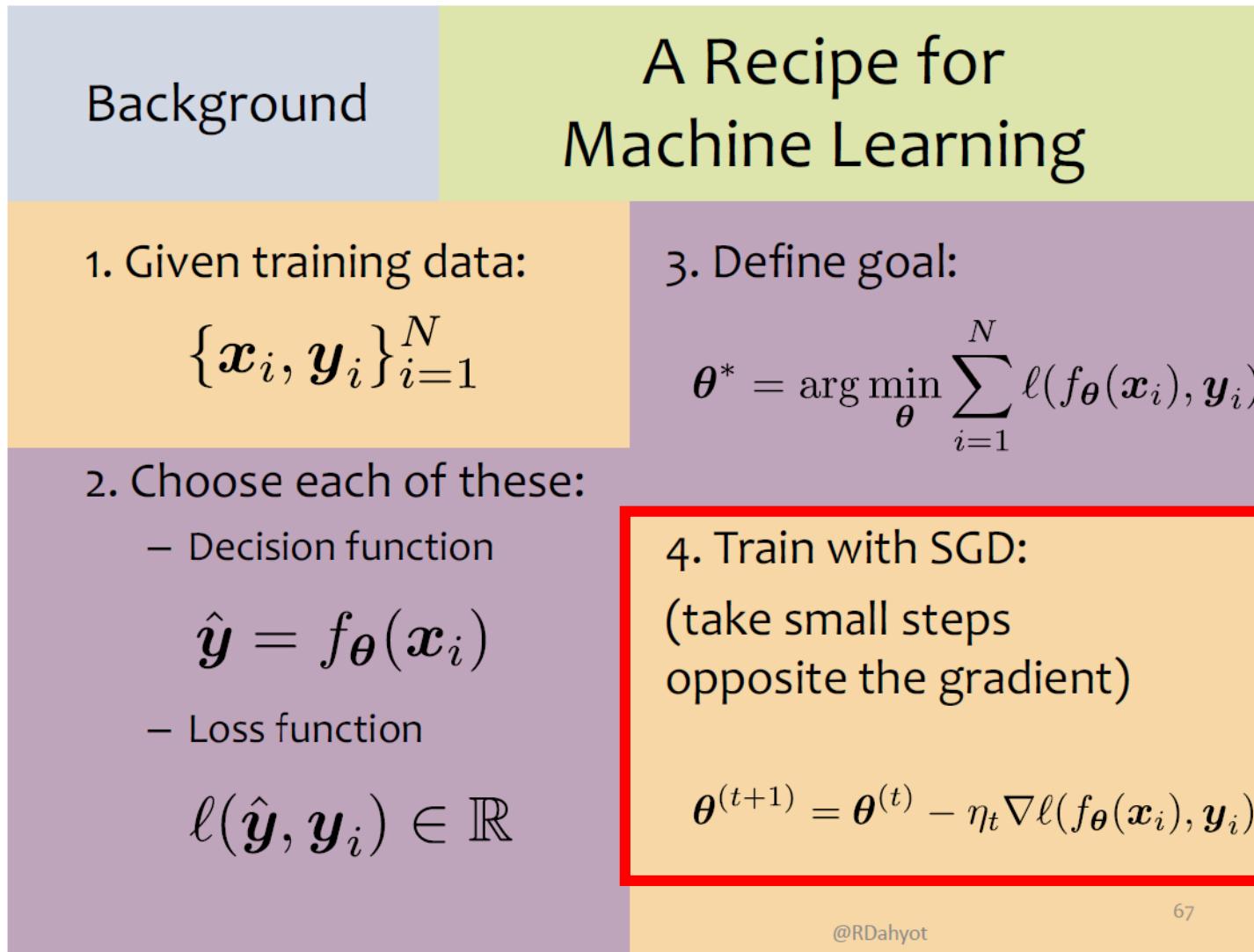


$$x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$$

Deep Neural Networks



Machine learning



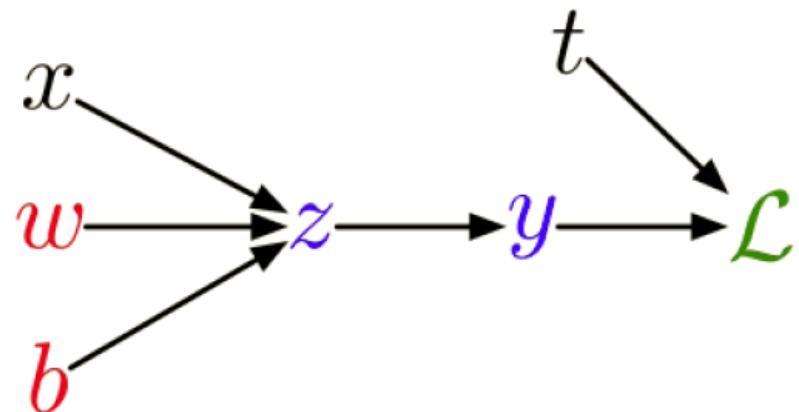
Backpropagation

- We've already been using the univariate Chain Rule.
- Recall: if $f(x)$ and $x(t)$ are univariate functions, then

$$\frac{d}{dt} f(x(t)) = \frac{df}{dx} \cdot \frac{dx}{dt}.$$

Backpropagation

Recall: Univariate logistic least squares model



$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

Let's compute the loss derivatives.

Backpropagation

How you would have done it in calculus class

$$\begin{aligned}\mathcal{L} &= \frac{1}{2}(\sigma(wx + b) - t)^2 \\ \frac{\partial \mathcal{L}}{\partial w} &= \frac{\partial}{\partial w} \left[\frac{1}{2}(\sigma(wx + b) - t)^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial w} (\sigma(wx + b) - t)^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial w} (\sigma(wx + b) - t) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial w} (wx + b) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b)x\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial b} &= \frac{\partial}{\partial b} \left[\frac{1}{2}(\sigma(wx + b) - t)^2 \right] \\ &= \frac{1}{2} \frac{\partial}{\partial b} (\sigma(wx + b) - t)^2 \\ &= (\sigma(wx + b) - t) \frac{\partial}{\partial b} (\sigma(wx + b) - t) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b) \frac{\partial}{\partial b} (wx + b) \\ &= (\sigma(wx + b) - t) \sigma'(wx + b)\end{aligned}$$

What are the disadvantages of this approach?

Backpropagation

A more structured way to do it

Computing the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\mathcal{L} = \frac{1}{2}(y - t)^2$$

(Forward pass)

Computing the derivatives:

$$\frac{d\mathcal{L}}{dy} = y - t$$

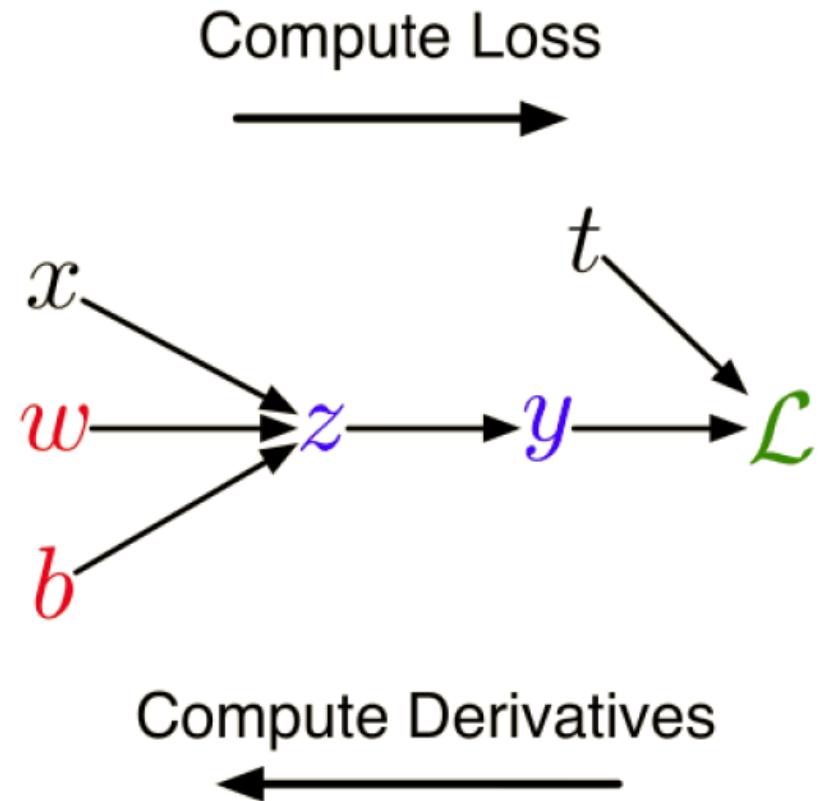
$$\frac{d\mathcal{L}}{dz} = \frac{d\mathcal{L}}{dy} \cdot \sigma'(z)$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{d\mathcal{L}}{dz} \cdot x$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{d\mathcal{L}}{dz}$$

(backpropagation)

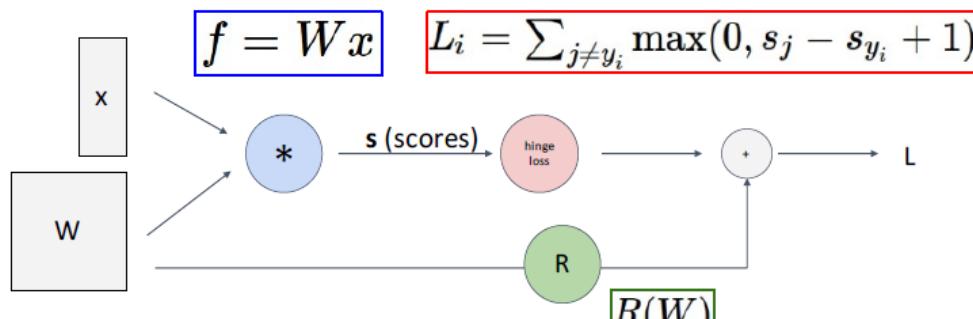
Remember, the goal isn't to obtain closed-form solutions, but to be able to write a program that efficiently computes the derivatives.



Backpropagation: summary

Last Time: Backpropagation

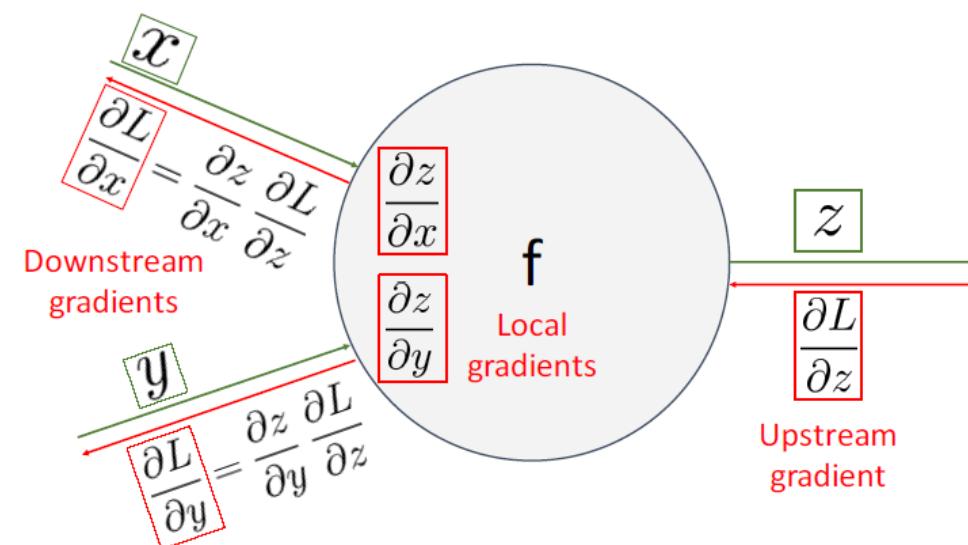
Represent complex expressions
as **computational graphs**



Forward pass computes outputs

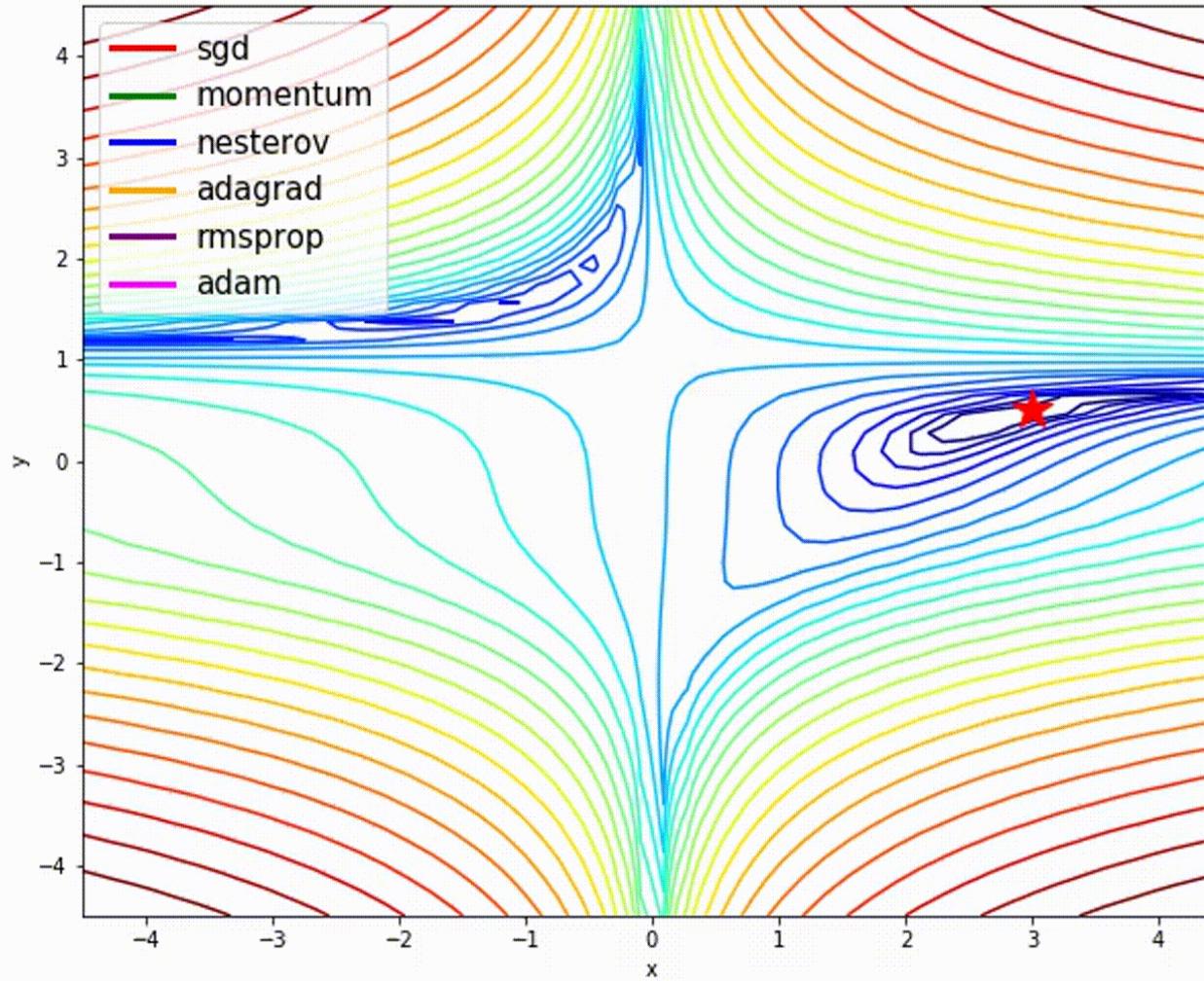
Backward pass computes gradients

During the backward pass, each node in the graph receives **upstream gradients** and multiplies them by **local gradients** to compute **downstream gradients**

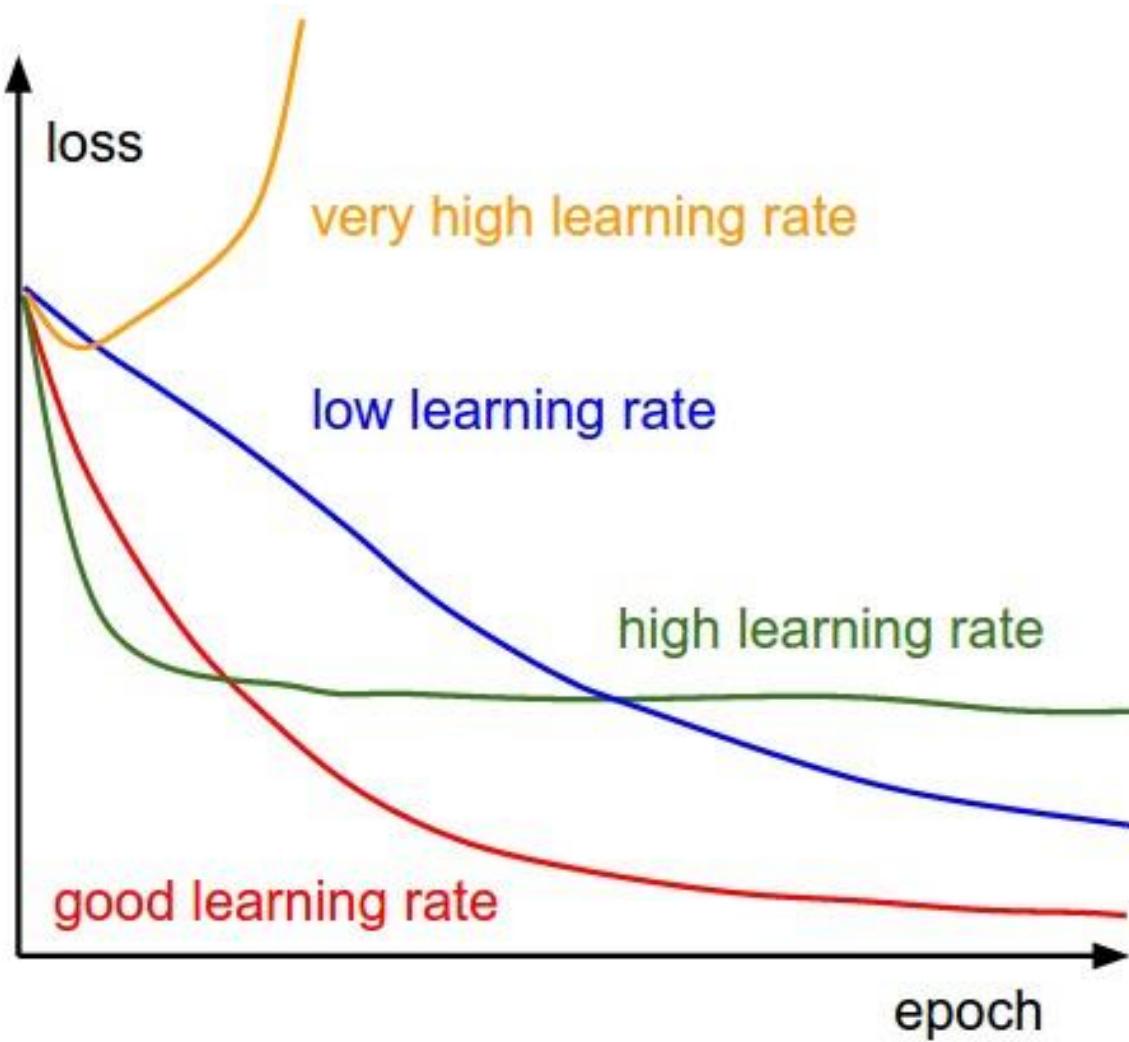


Source: [Justin Johnson](#)

Optimizers

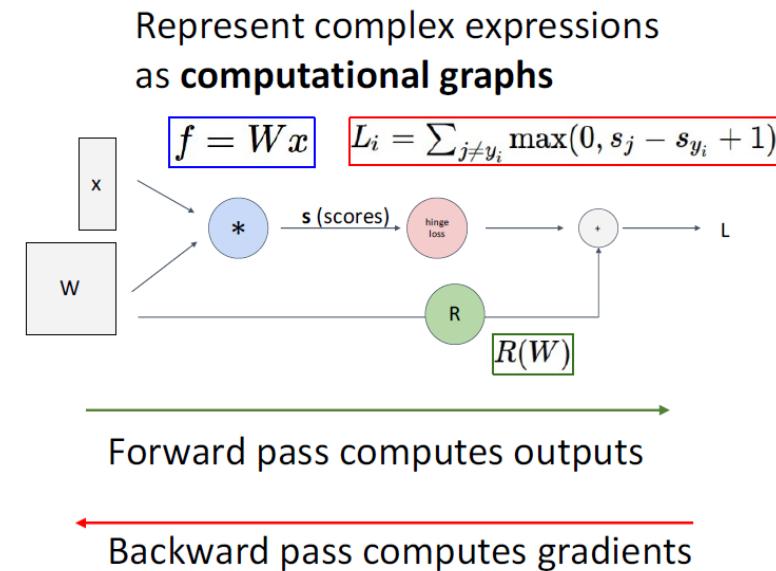


Learning rate



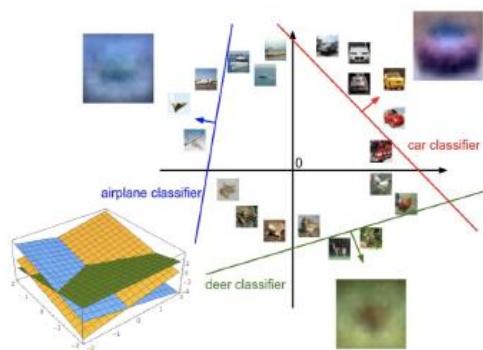
Epochs and (mini-)batches

- Epoch: Forward + backward passes for entire training dataset
- Batch: Subset of training dataset used for forward + backward passes
- Batch size: number of elements in a batch
- Iterations: number of batches needed to complete 1 epoch

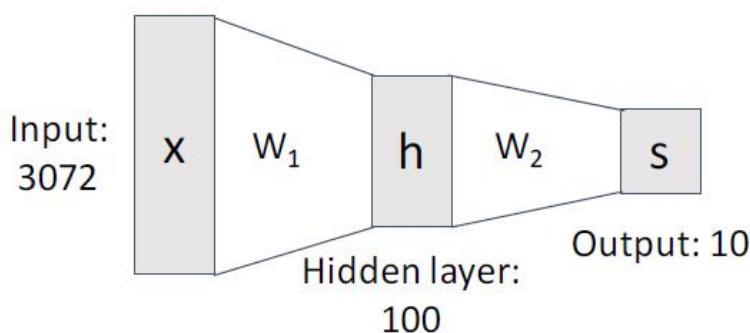


Why Convolutional Networks?

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



$$f = W_2 \max(0, W_1 x)$$



Stretch pixels into column



Input image
(2, 2)

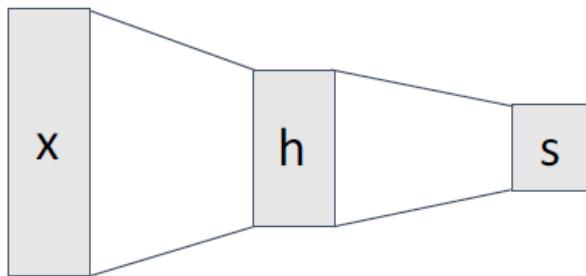
Problem: So far our classifiers don't respect the spatial structure of images!



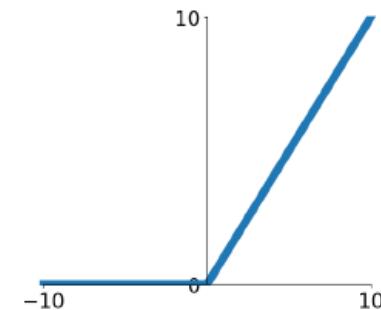
Solution: Define new computational nodes that operate on images!

Components of a Convolution Network

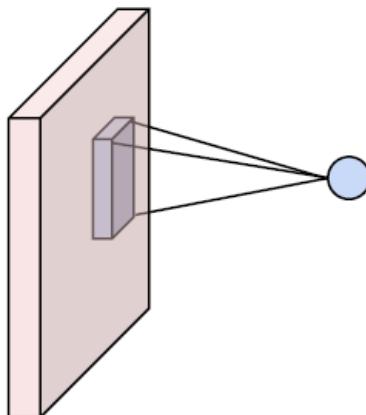
Fully-Connected Layers



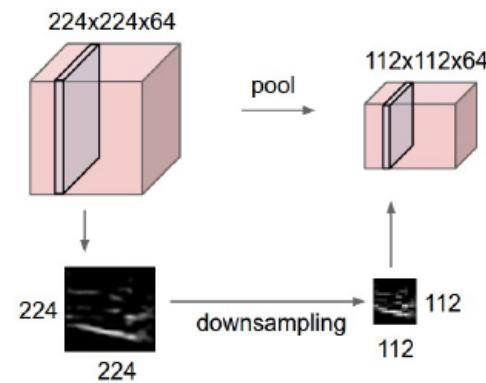
Activation Function



Convolution Layers



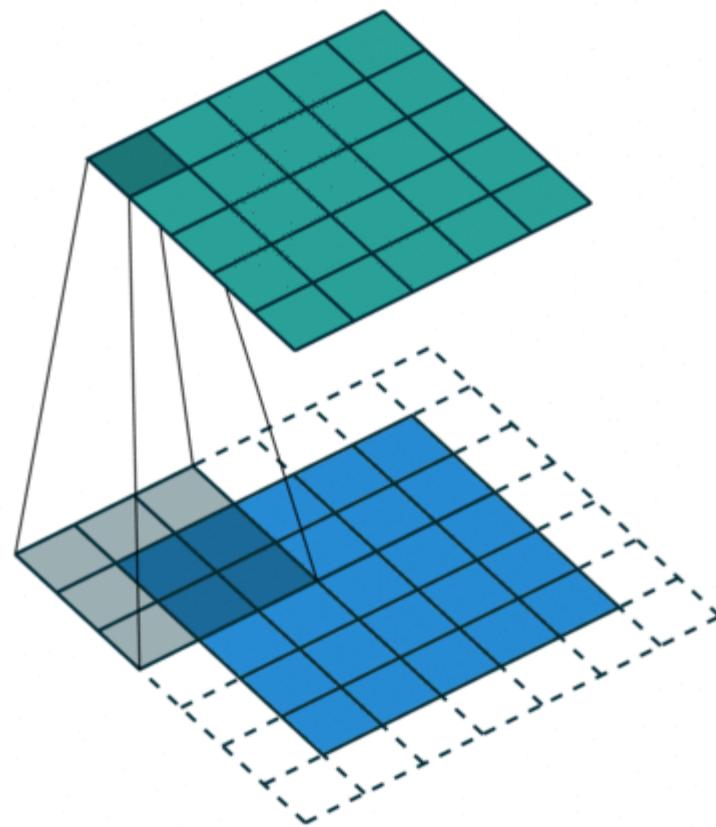
Pooling Layers



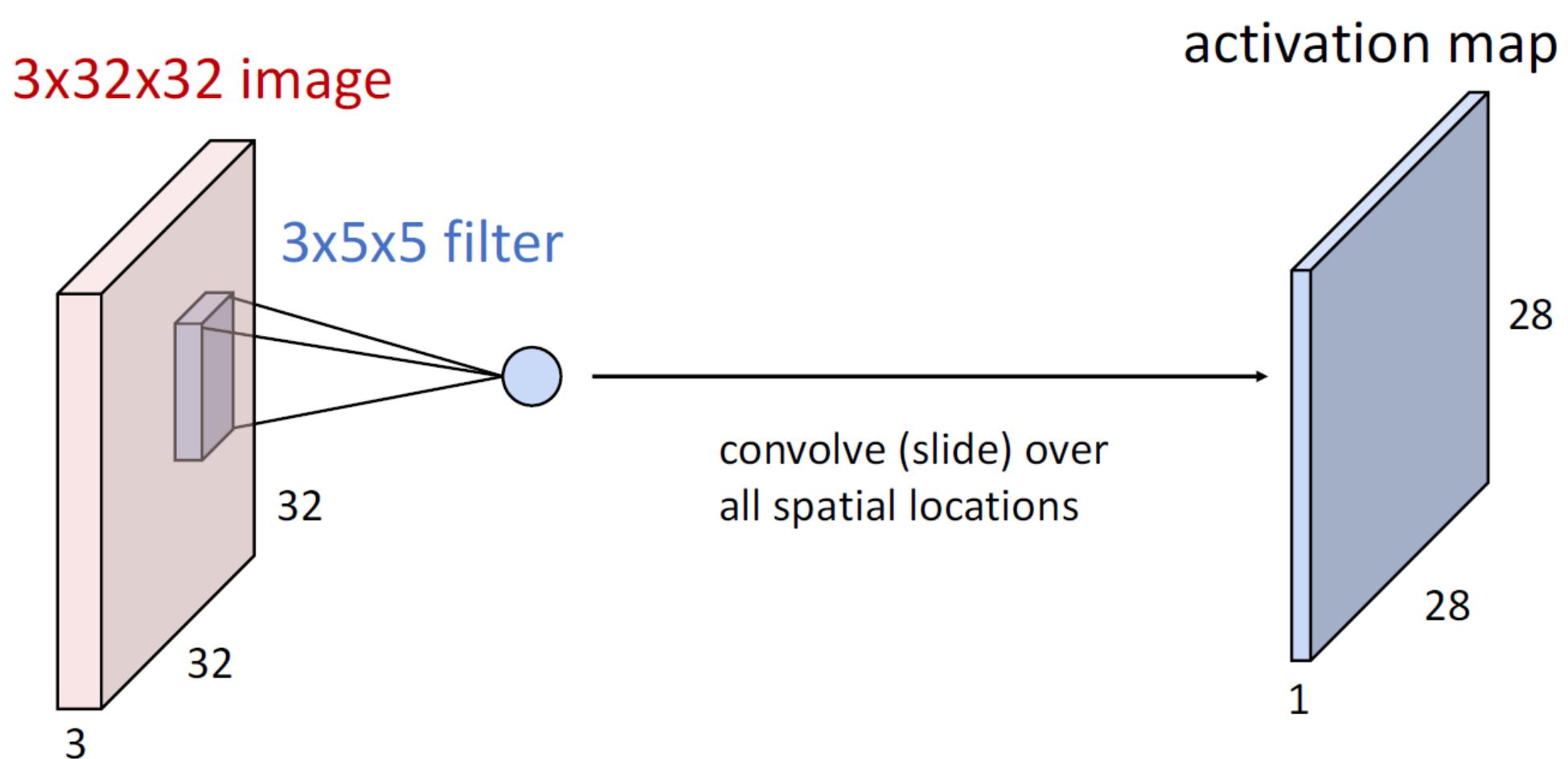
Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Convolution layers

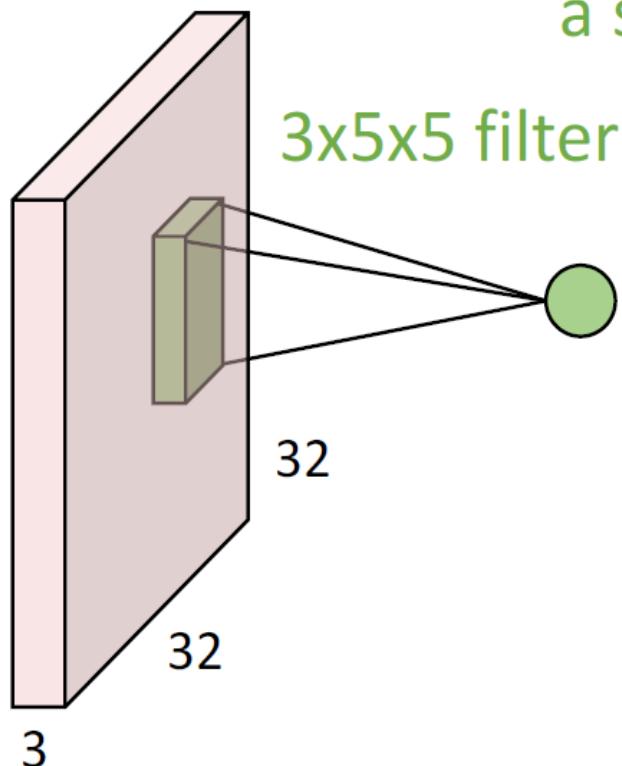


Convolution layers



Convolution layers

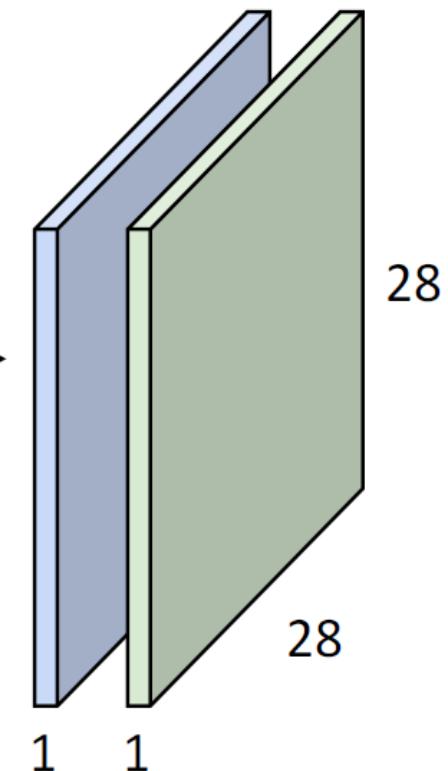
3x32x32 image



Consider repeating with
a second (green) filter:

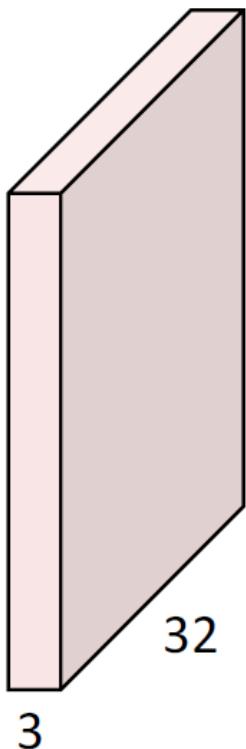
convolve (slide) over
all spatial locations

two 1x28x28
activation map

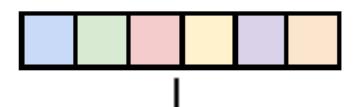


Convolution layers

3x32x32 image



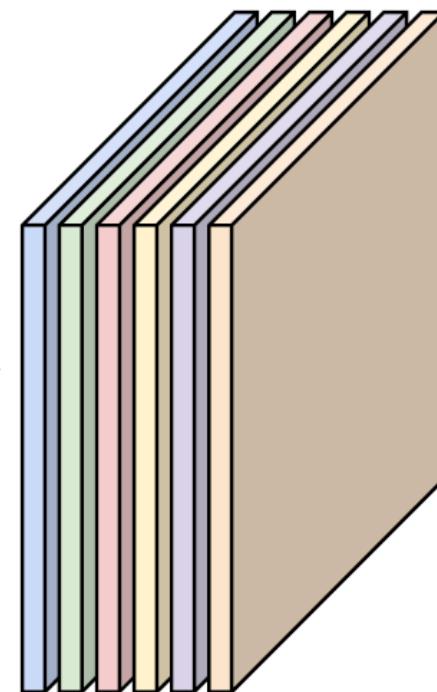
Also 6-dim bias vector:



6x3x5x5
filters

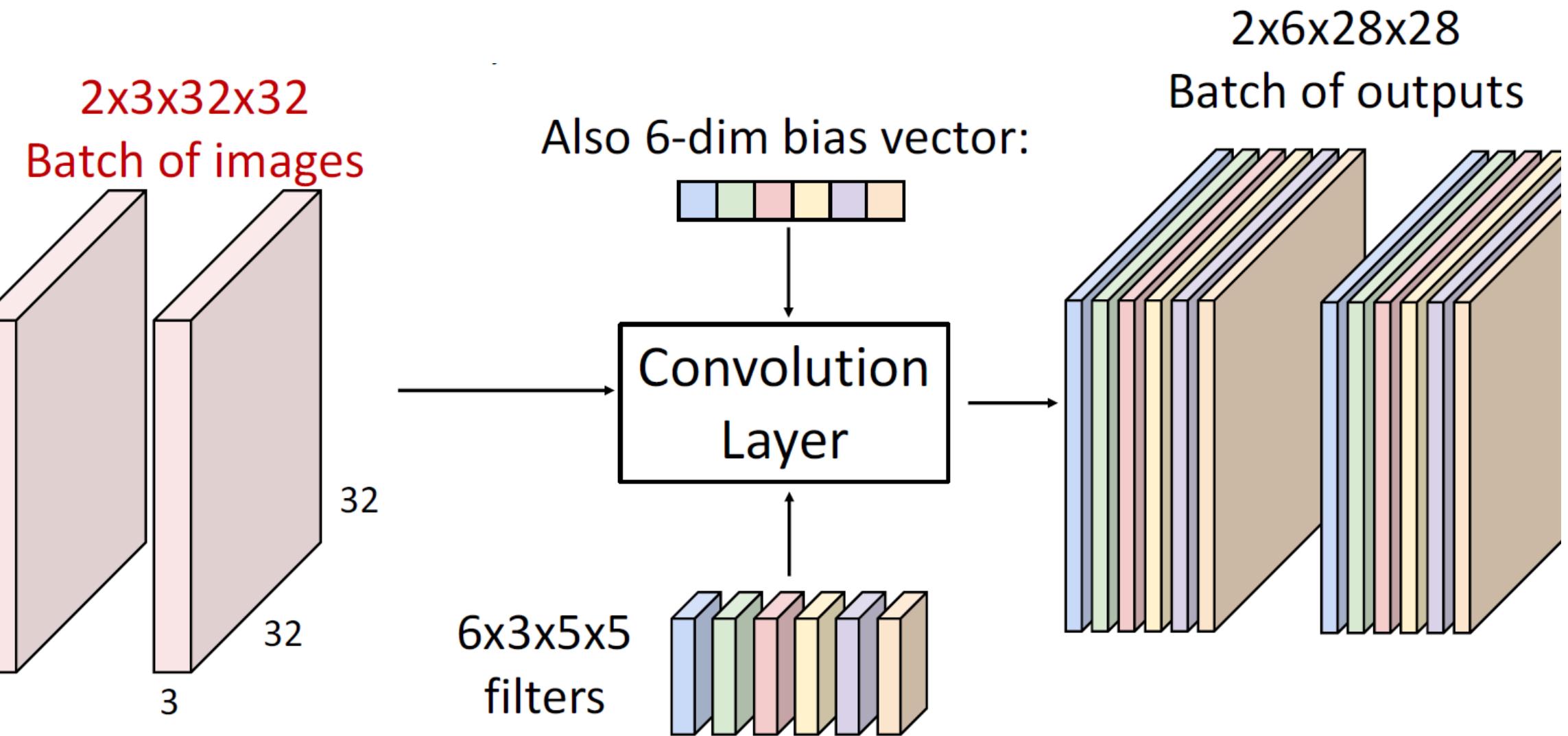


6 activation maps,
each 1x28x28

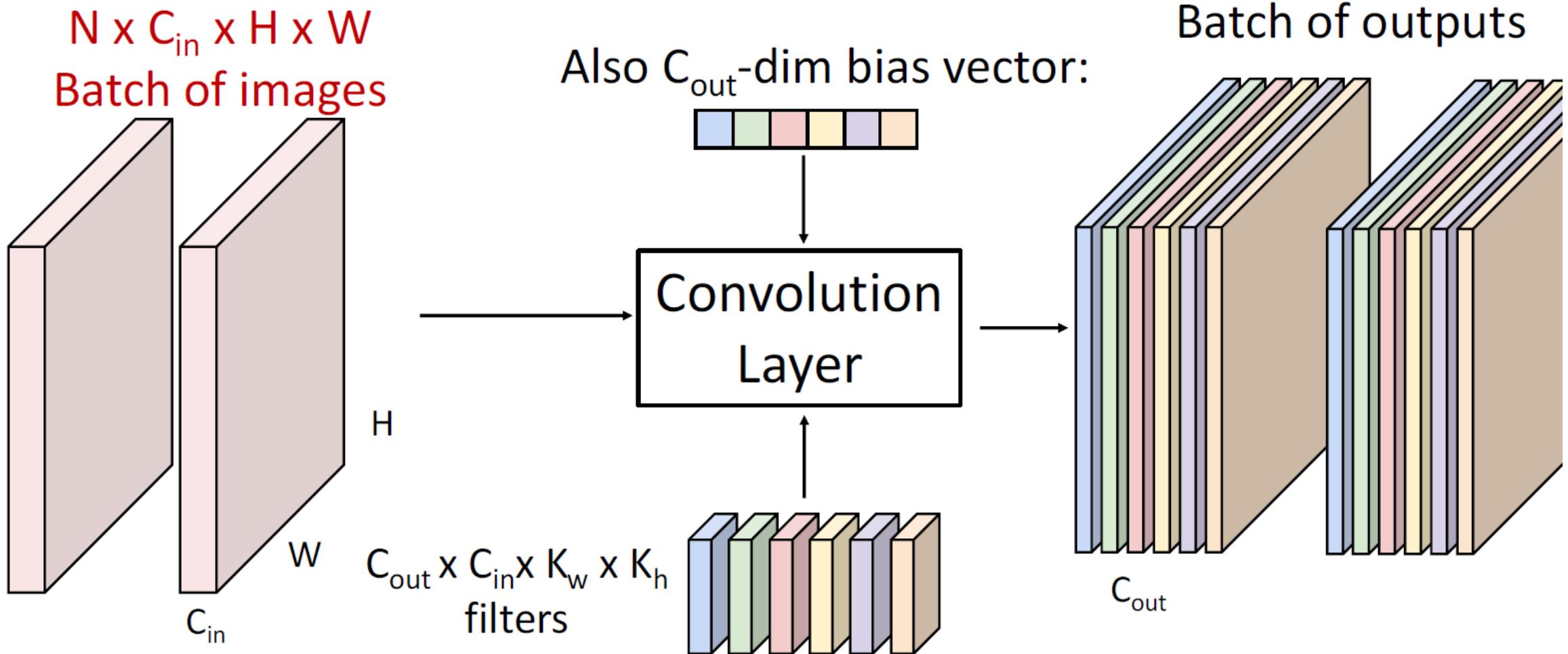


Stack activations to get a
6x28x28 output image!

Convolution layers



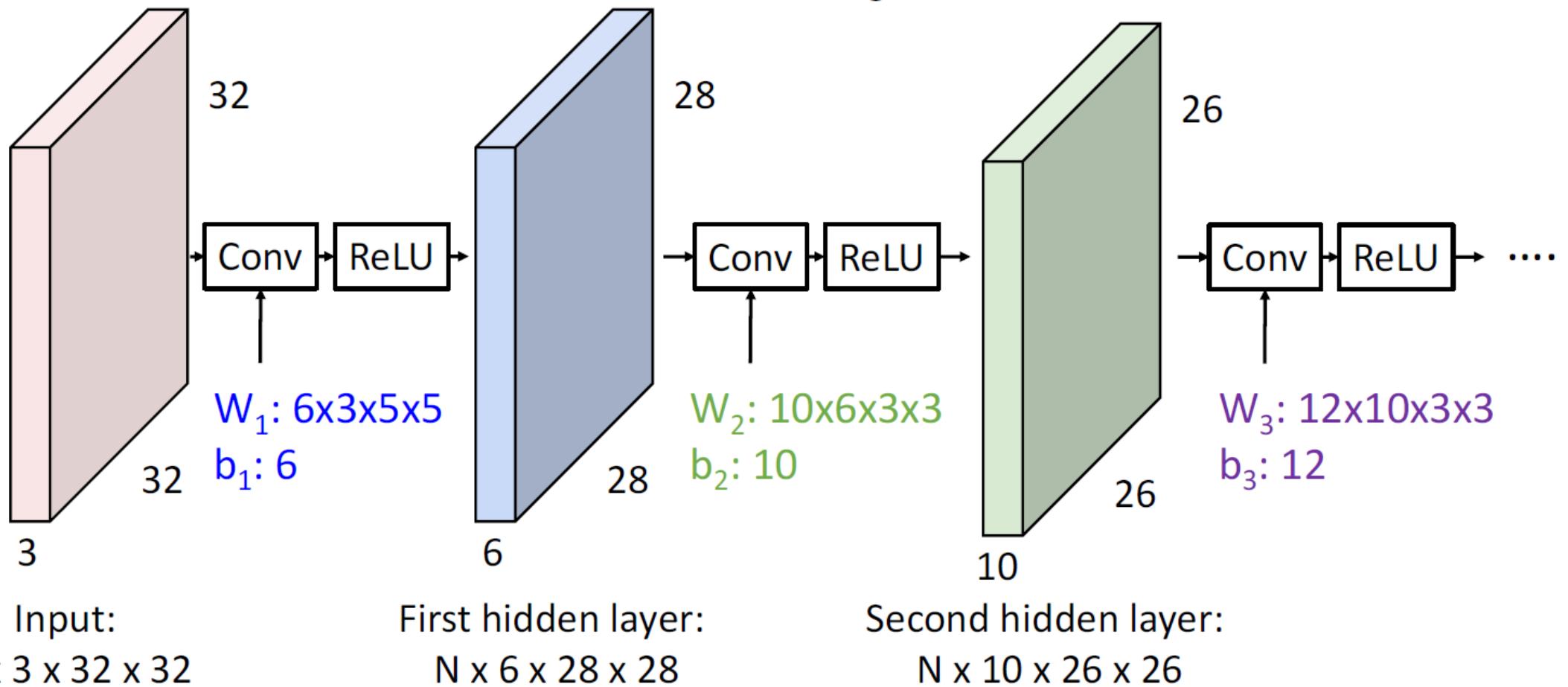
Convolution layers



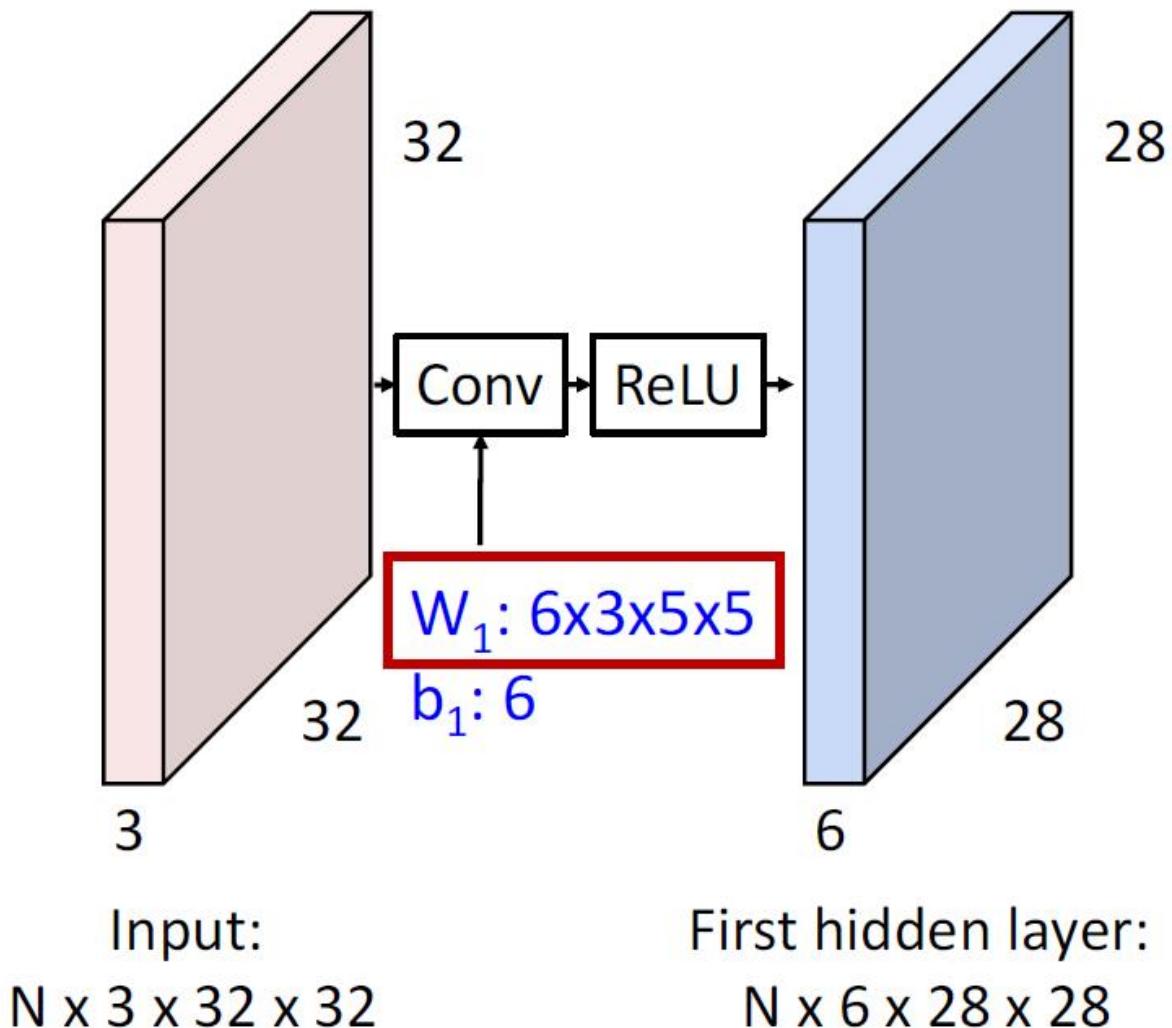
Stacking convolutions

Q: What happens if we stack two convolution layers? (Recall $y=W_2W_1x$ is a linear classifier)

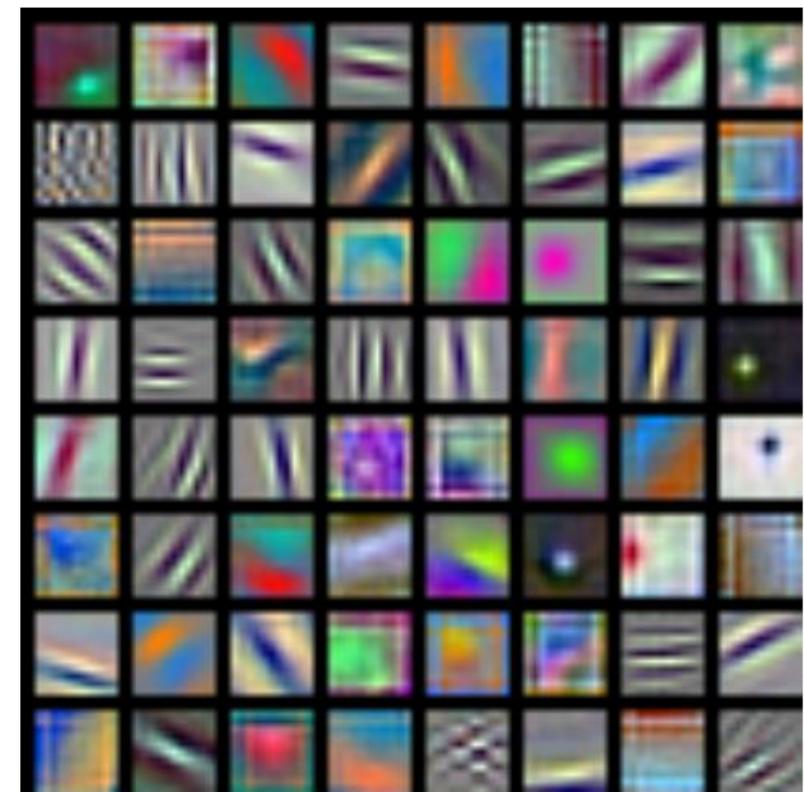
A: We get another convolution!



A closer look at convolutional filters



First-layer conv filters: local image templates
(Often learns oriented edges, opposing colors)



AlexNet: 64 filters, each $3 \times 11 \times 11$

Padding

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7

Filter: 3x3

Output: 5x5

In general:

Input: W

Filter: K

Padding: P

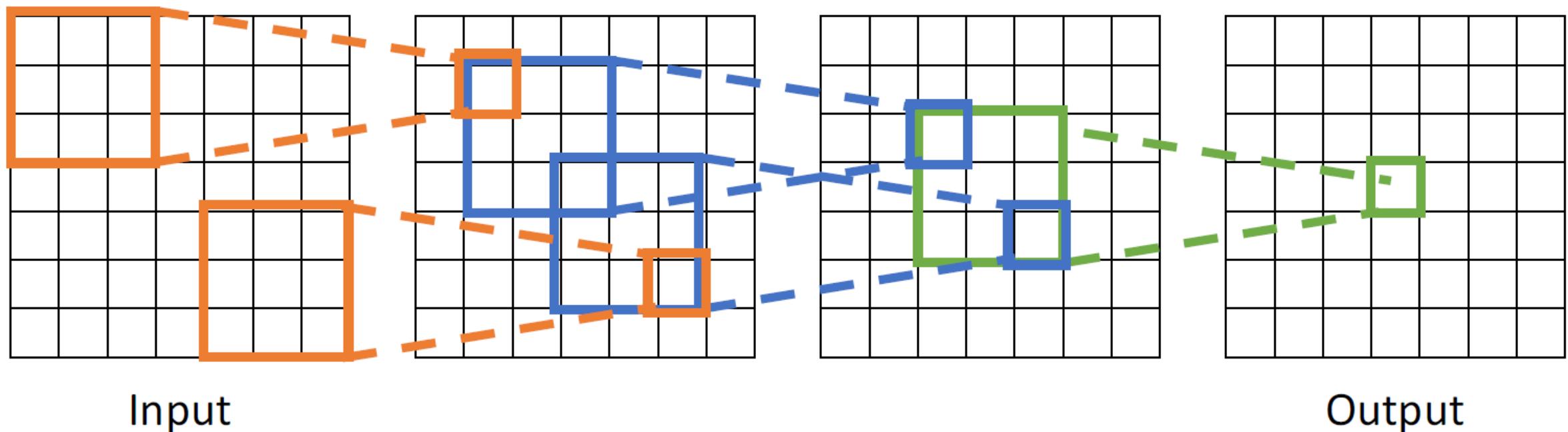
Output: $W - K + 1 + 2P$

Very common:

Set $P = (K - 1) / 2$ to
make output have
same size as input!

Receptive fields

Each successive convolution adds $K - 1$ to the receptive field size
With L layers the receptive field size is $1 + L * (K - 1)$

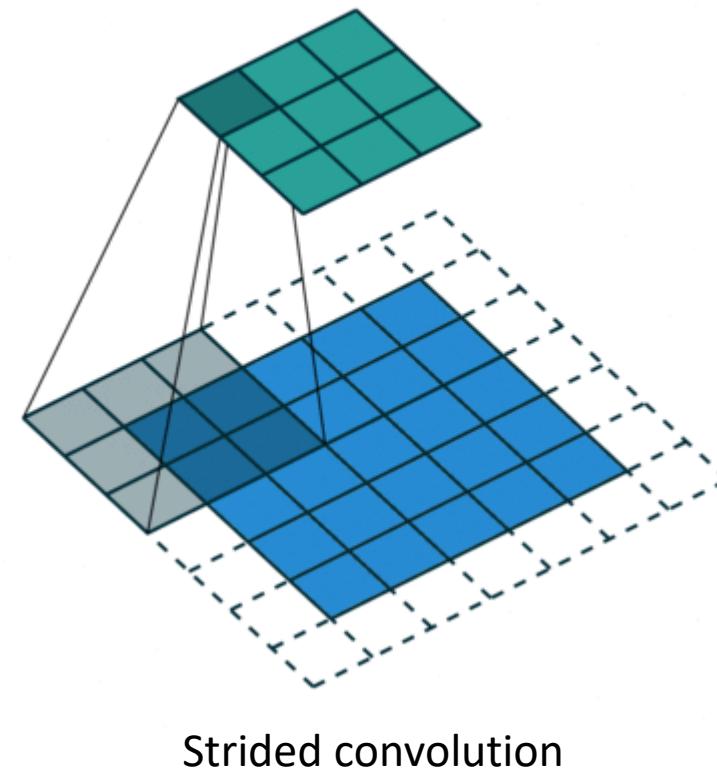
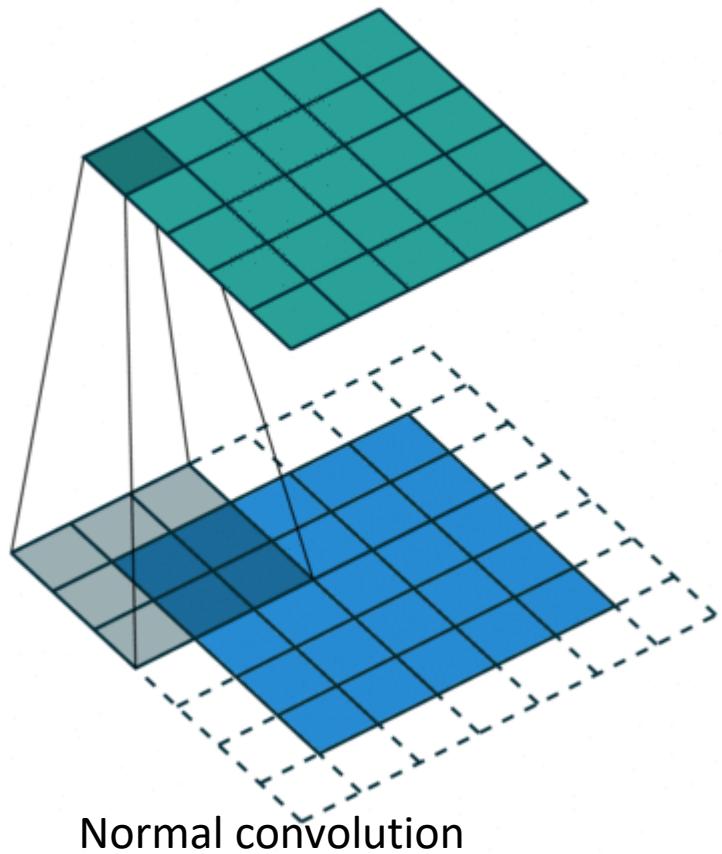


Input

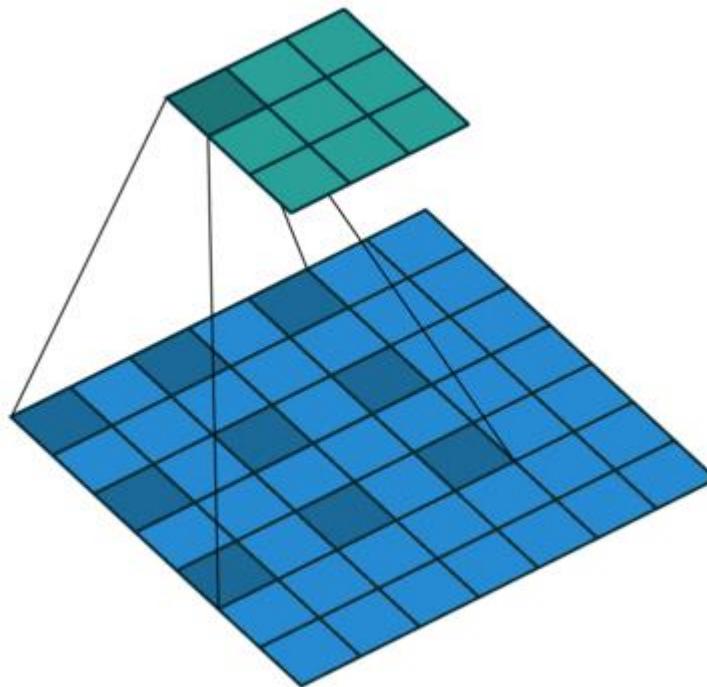
Output

Be careful – “receptive field in the input” vs “receptive field in the previous layer”
Hopefully clear from context!

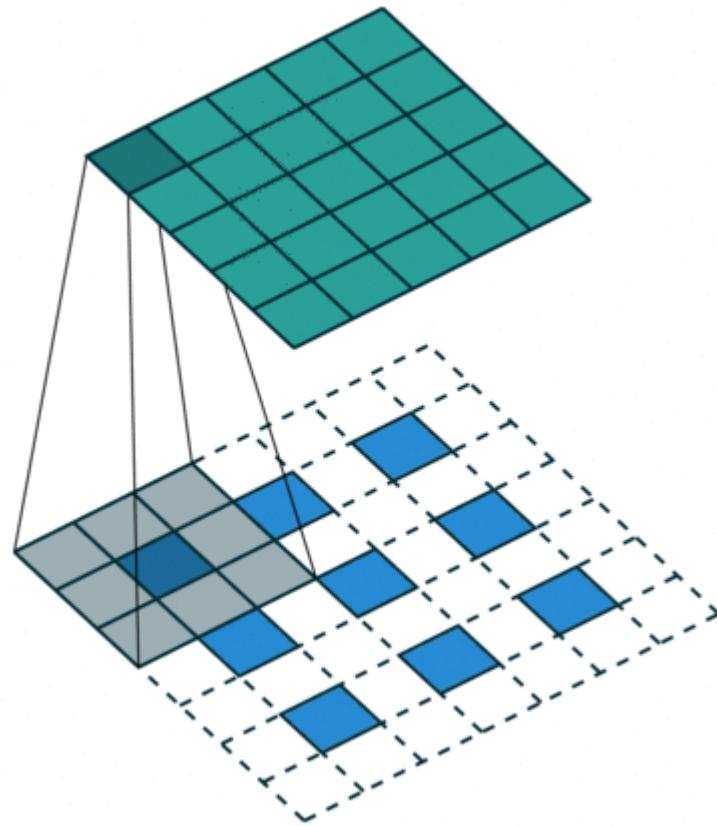
Strided convolution



Dilated convolution

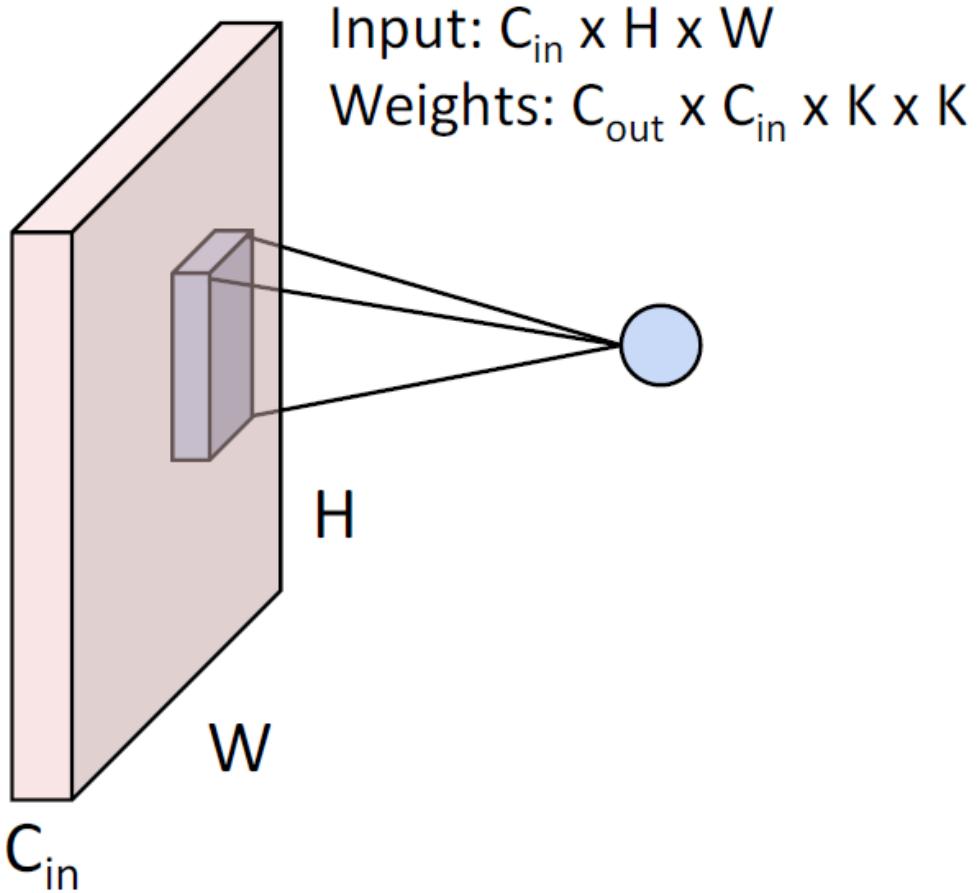


Transposed convolution

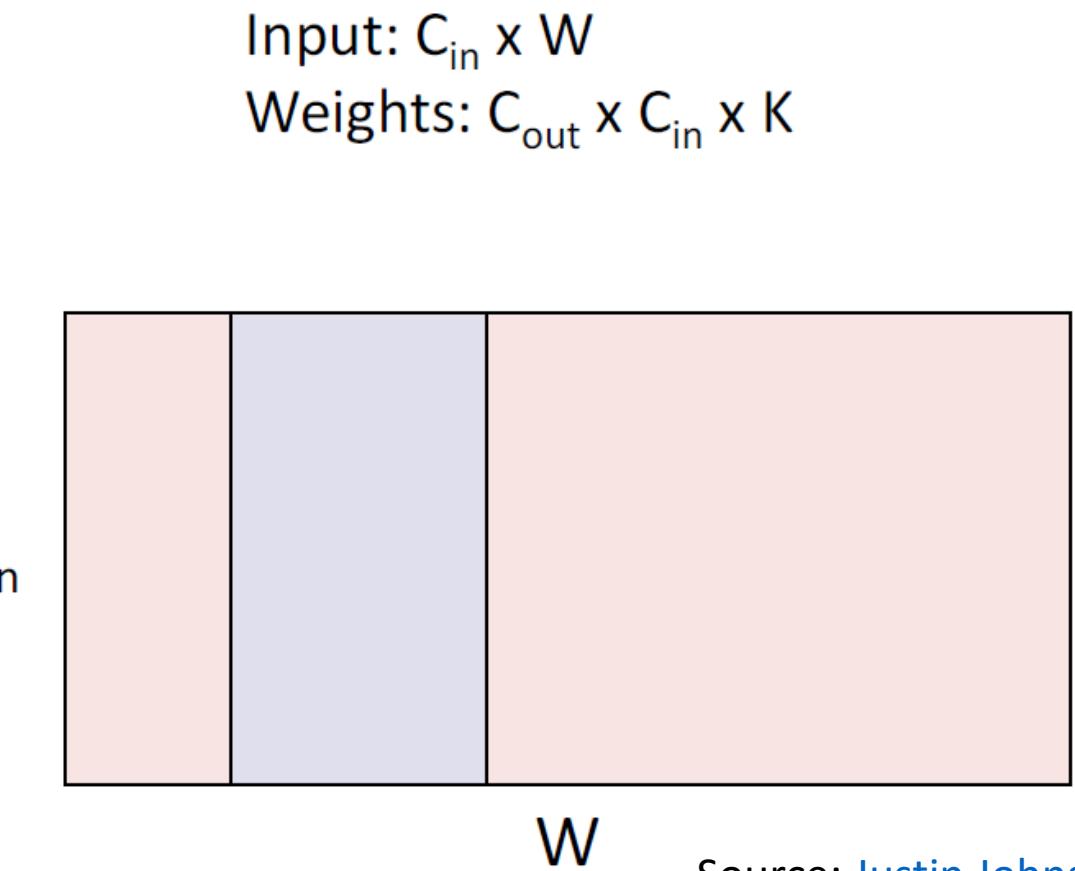


Other types of convolution

So far: 2D Convolution



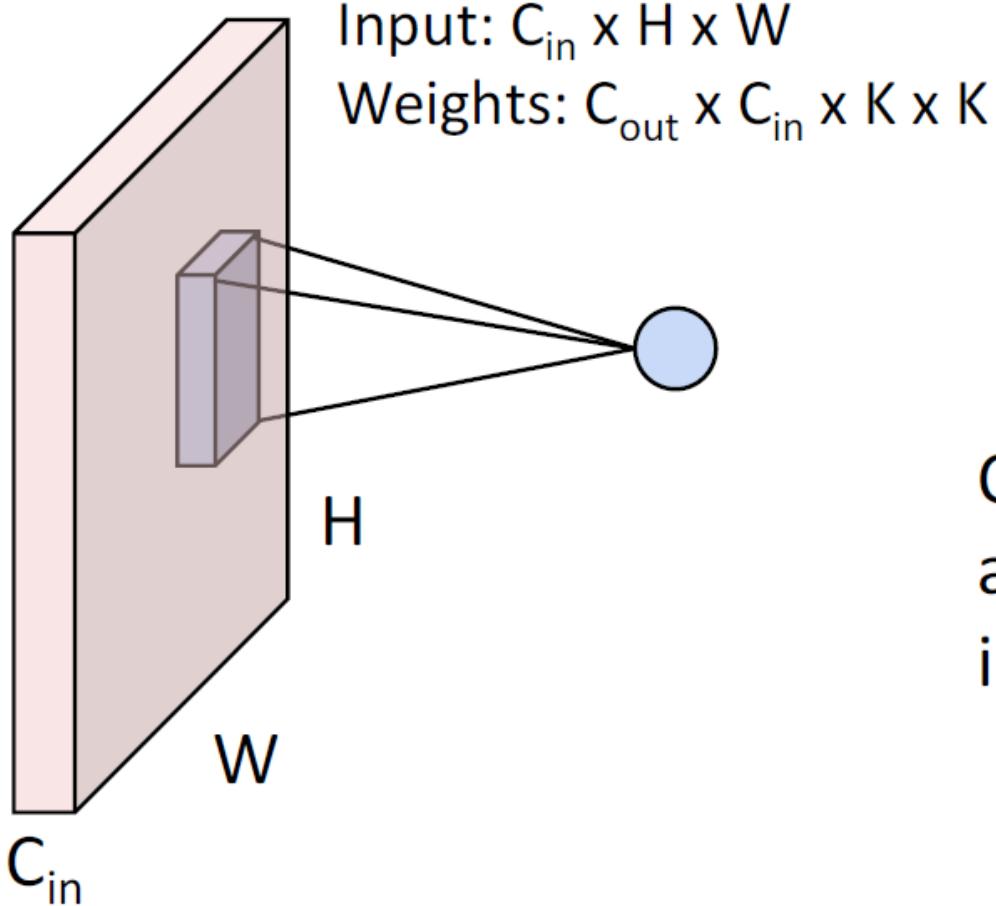
1D Convolution



Source: [Justin Johnson](#)

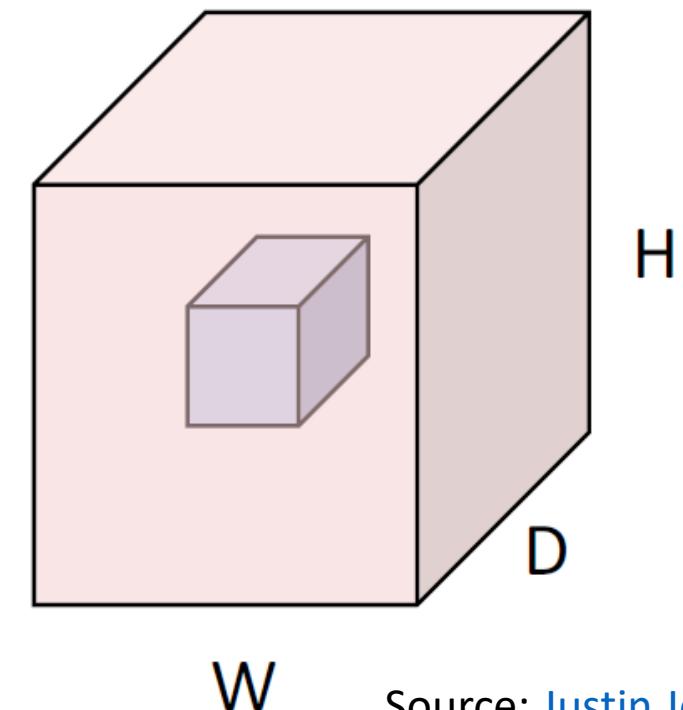
Other types of convolution

So far: 2D Convolution



3D Convolution

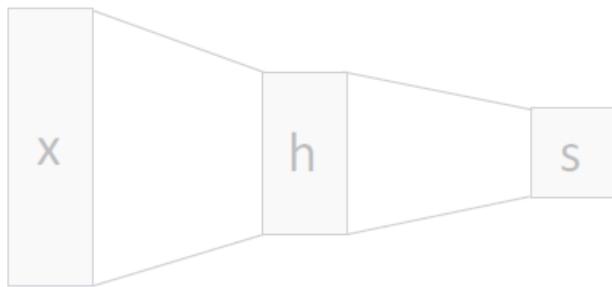
Input: $C_{in} \times H \times W \times D$
Weights: $C_{out} \times C_{in} \times K \times K \times K$



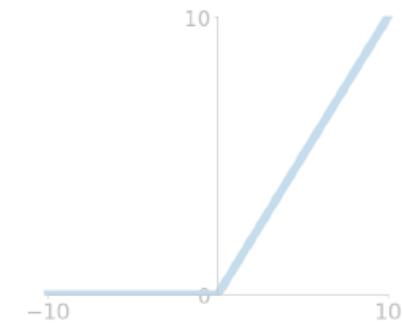
Source: [Justin Johnson](#)

Components of a Convolution Network

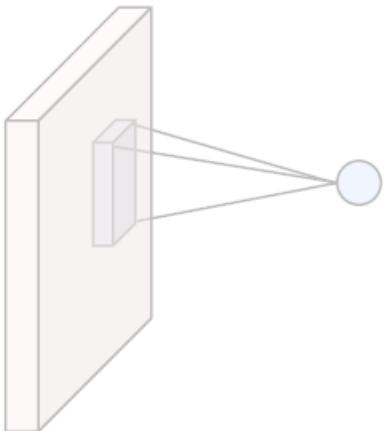
Fully-Connected Layers



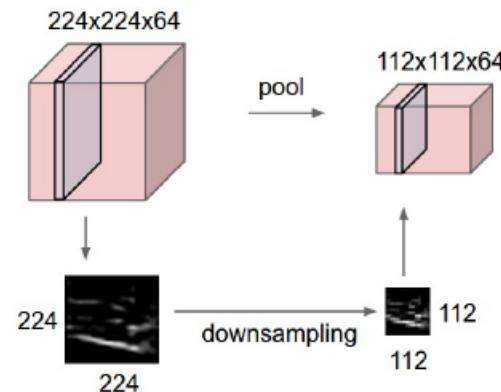
Activation Function



Convolution Layers



Pooling Layers

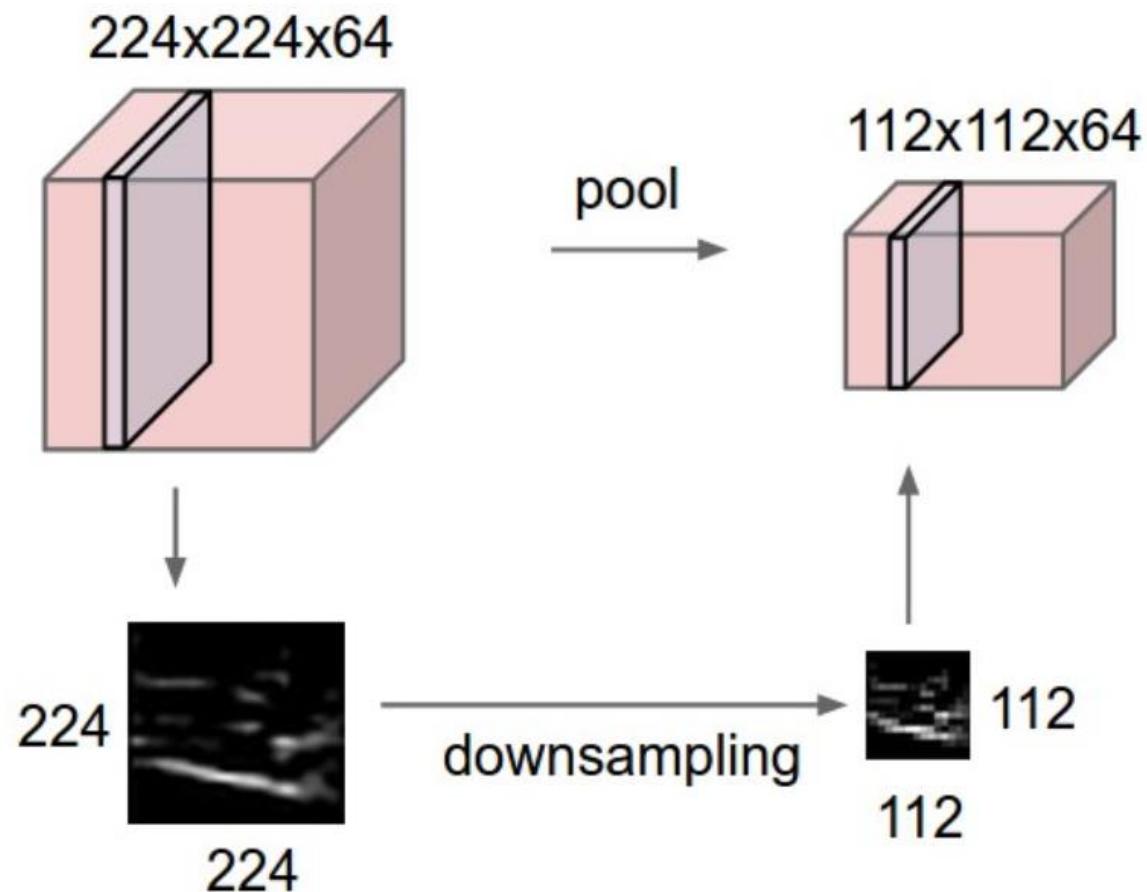


Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

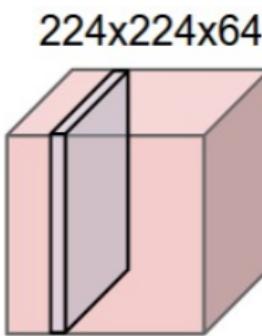
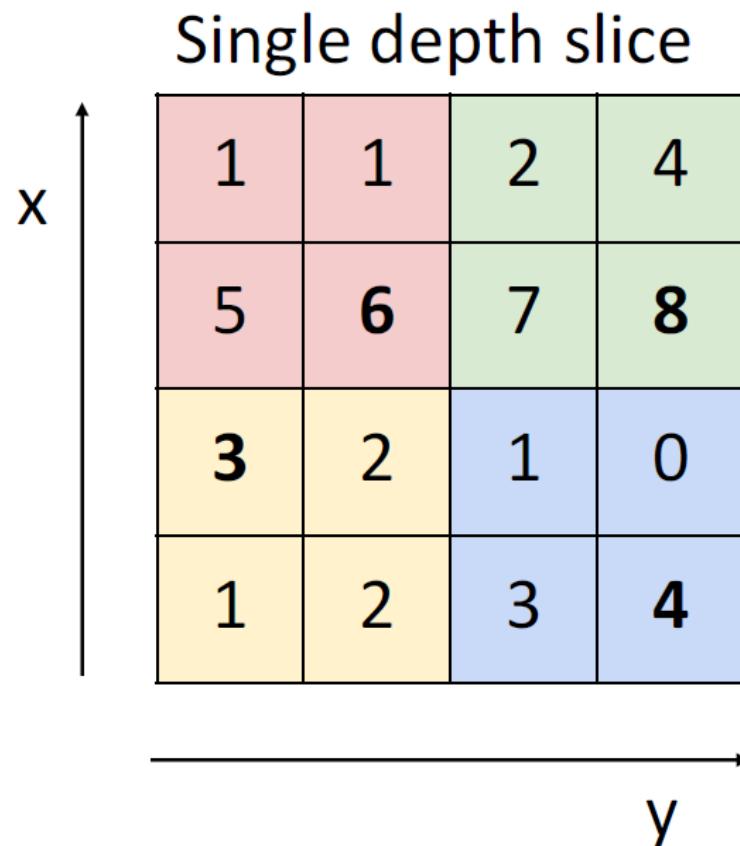
Source: [Justin Johnson](#)

Pooling layers



Hyperparameters:
Kernel Size
Stride
Pooling function

Max pooling



Max pooling with 2x2 kernel size and stride 2

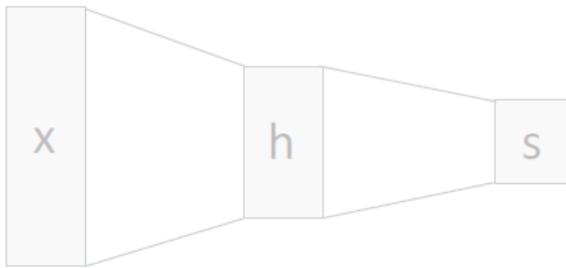
6	8
3	4

Introduces **invariance** to
small spatial shifts
No learnable parameters!

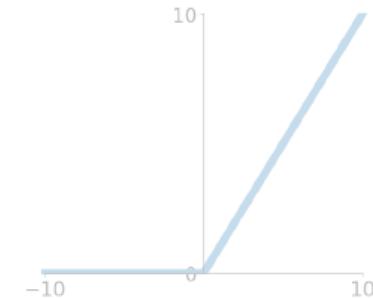
Geoffrey Hinton: “The pooling operation used in convolutional neural networks is a big mistake and the fact that it works so well is a disaster.”

Components of a Convolution Network

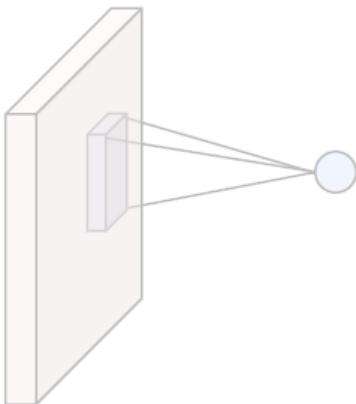
Fully-Connected Layers



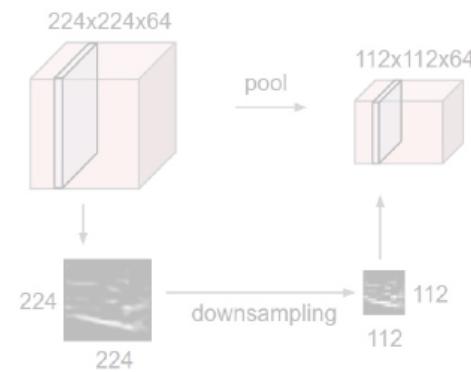
Activation Function



Convolution Layers



Pooling Layers



Normalization

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Batch normalization

Idea: “Normalize” the outputs of a layer so they have zero mean and unit variance

Why? Helps reduce “internal covariate shift”, improves optimization

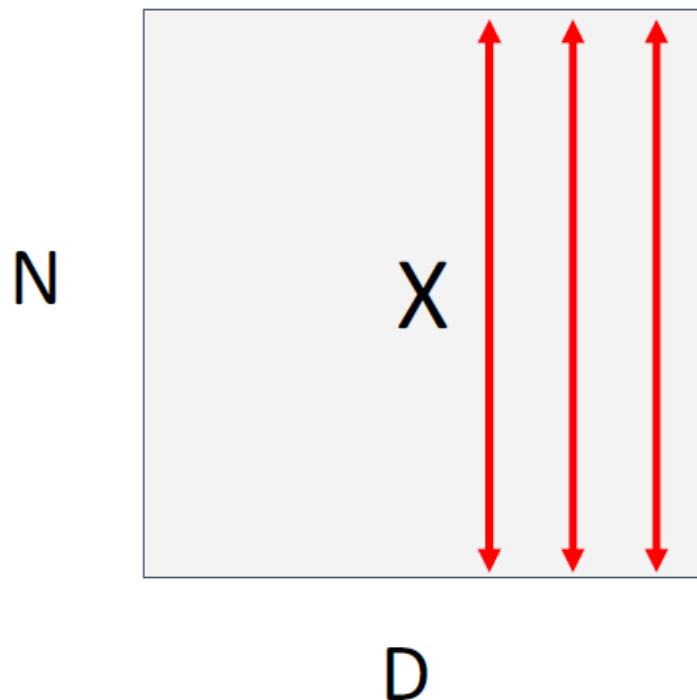
We can normalize a batch of activations like this:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

This is a **differentiable function**, so we can use it as an operator in our networks and backprop through it!

Batch normalization

Input: $x : N \times D$



$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$

Per-channel mean, shape is D

$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$

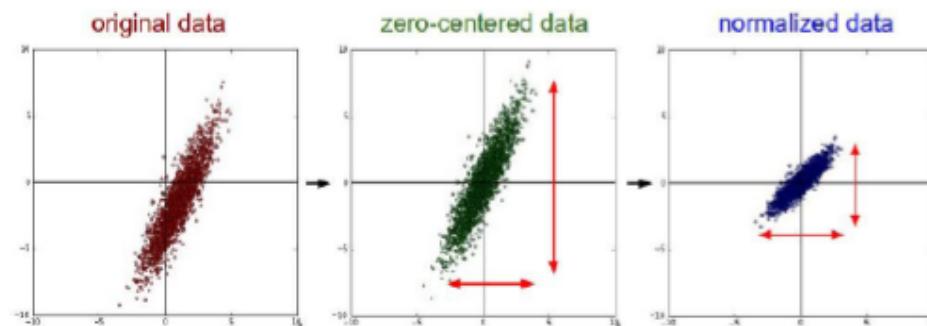
Per-channel std, shape is D

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Normalized x,
Shape is $N \times D$

Decorrelated Batch Normalization

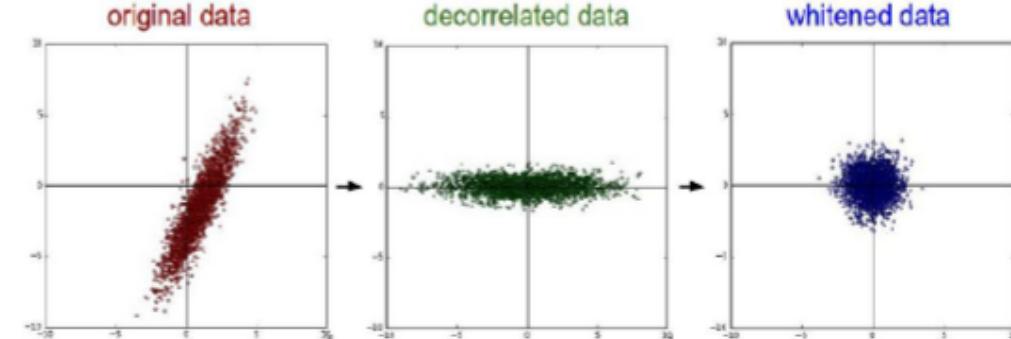
Batch Normalization



$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

BatchNorm normalizes the data, but cannot correct for correlations among the input features

Decorrelated Batch Normalization



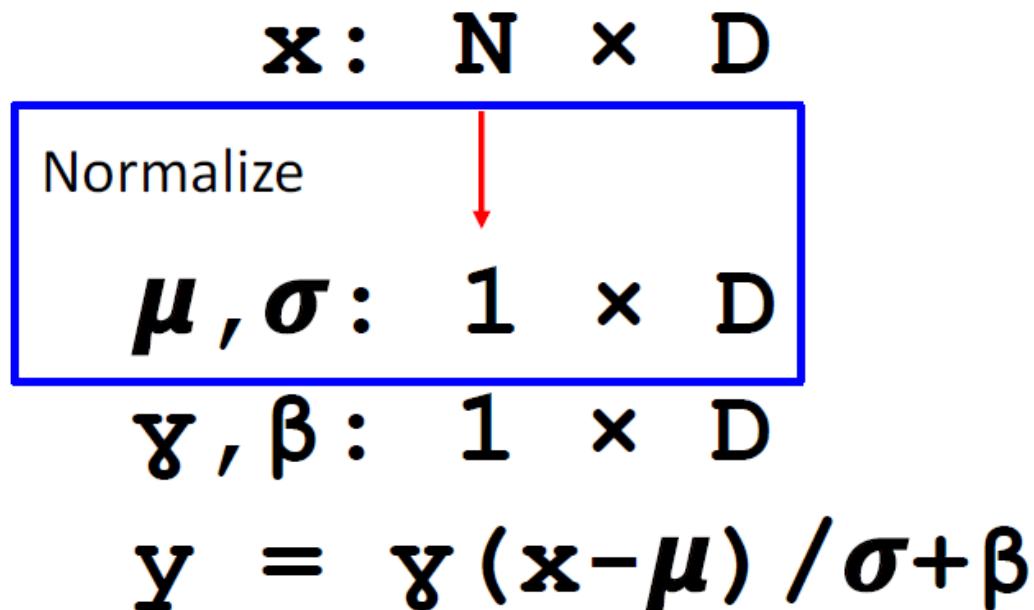
$$\hat{x}_i = \Sigma^{-\frac{1}{2}}(x_i - \mu)$$

DBN whitens the data using the full covariance matrix of the minibatch; this corrects for correlations

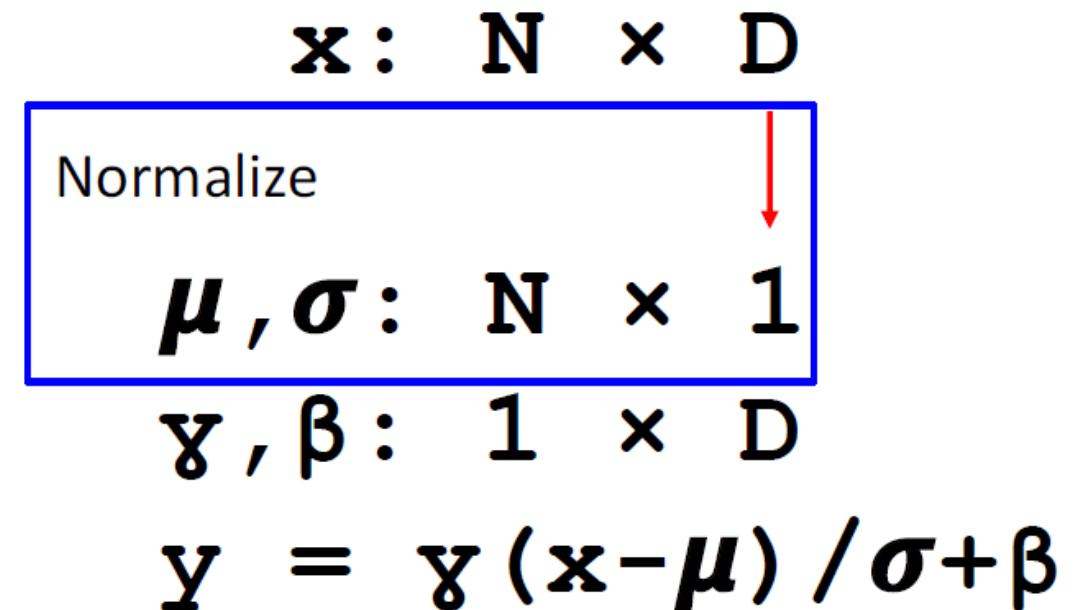
Huang et al, "Decorrelated Batch Normalization", arXiv 2018 (Appeared 4/23/2018)

Layer normalization

Batch Normalization for
fully-connected networks



Layer Normalization for fully-connected networks
Same behavior at train and test!
Used in RNNs, Transformers



Instance normalization

Batch Normalization for convolutional networks

$$\mathbf{x} : N \times C \times H \times W$$

Normalize



$$\mu, \sigma : 1 \times C \times 1 \times 1$$

$$\gamma, \beta : 1 \times C \times 1 \times 1$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

Instance Normalization for convolutional networks
Same behavior at train / test!

$$\mathbf{x} : N \times C \times H \times W$$

Normalize

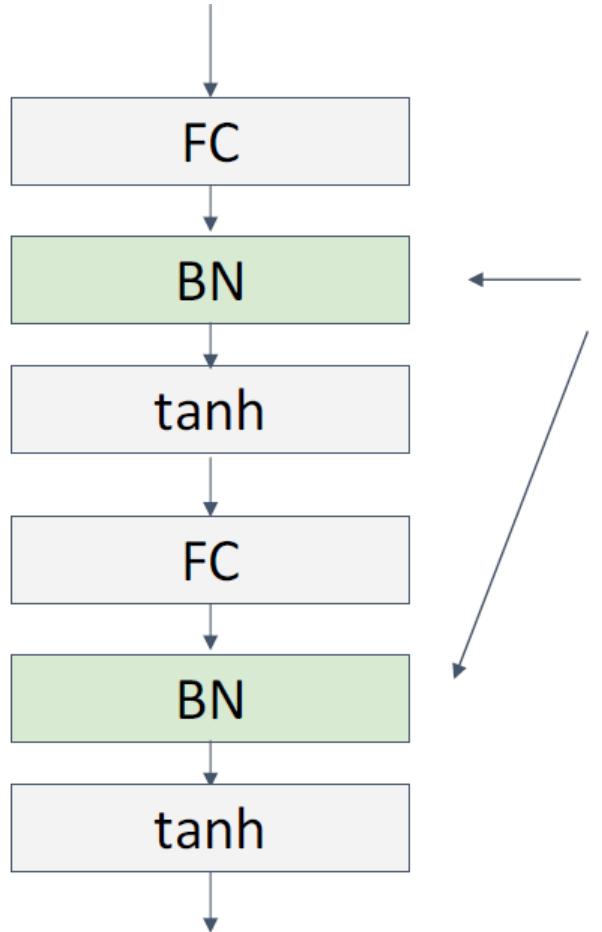


$$\mu, \sigma : N \times C \times 1 \times 1$$

$$\gamma, \beta : 1 \times C \times 1 \times 1$$

$$y = \gamma(x - \mu) / \sigma + \beta$$

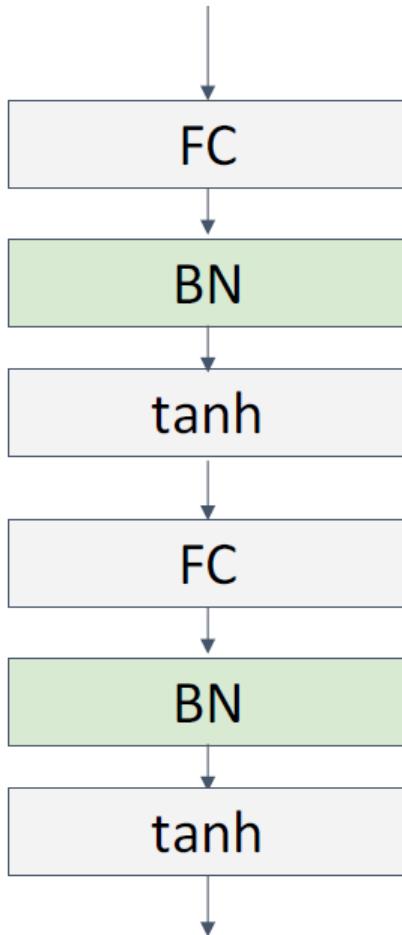
Batch normalization



Usually inserted after Fully Connected or Convolutional layers, and before nonlinearity.

$$\hat{x}^{(k)} = \frac{x^{(k)} - \text{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

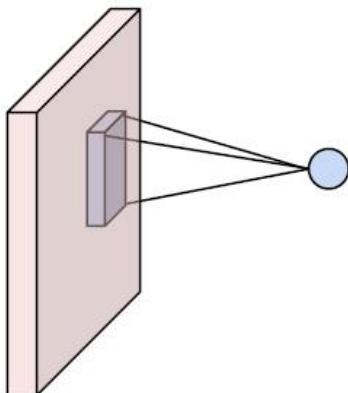
Batch normalization



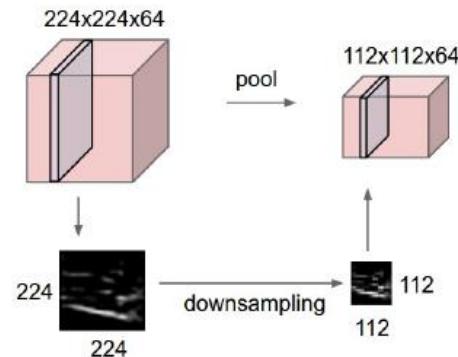
- Makes deep networks **much** easier to train!
- Allows higher learning rates, faster convergence
- Networks become more robust to initialization
- Acts as regularization during training
- Zero overhead at test-time: can be fused with conv!
- Not well-understood theoretically (yet)
- Behaves differently during training and testing: this is a very common source of bugs!

Summary: Components of a Convolution Network

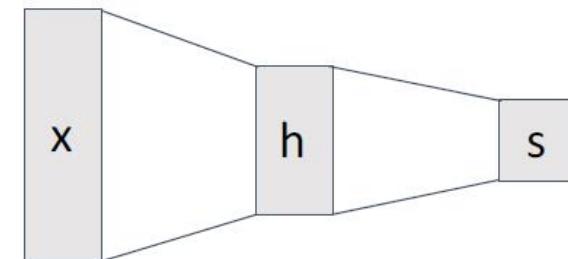
Convolution Layers



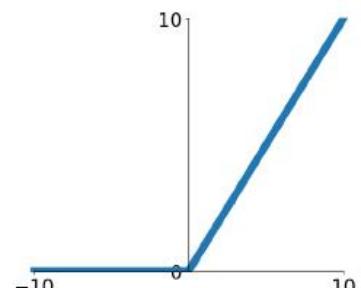
Pooling Layers



Fully-Connected Layers



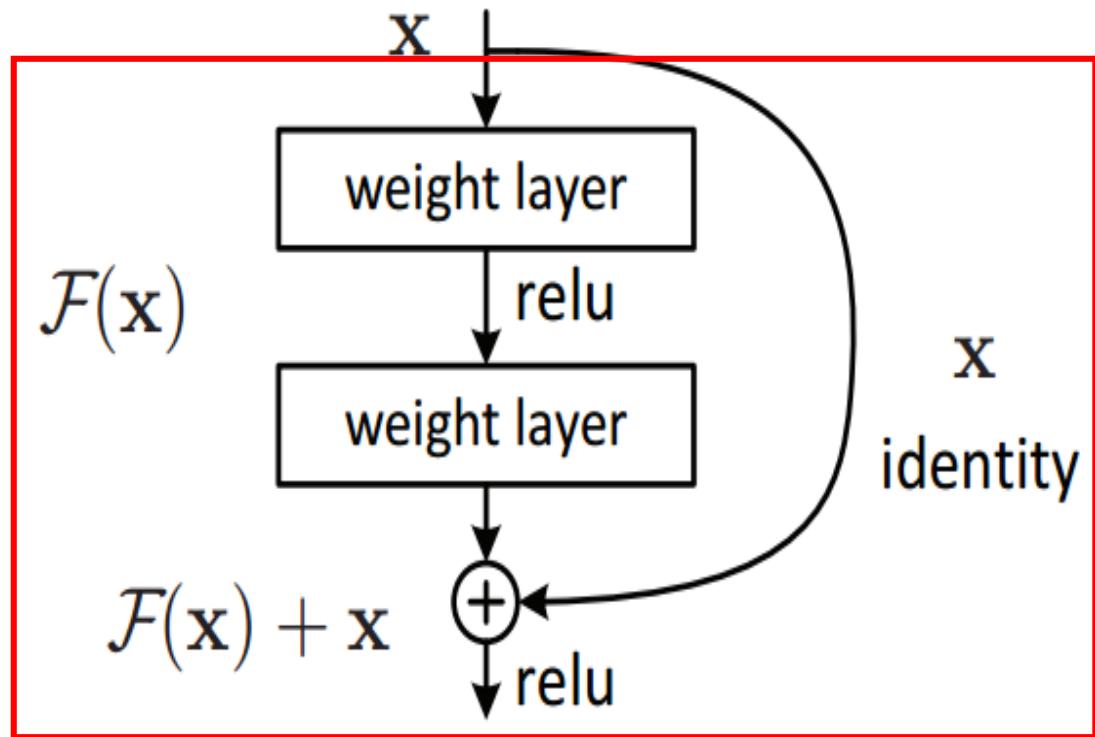
Activation Function



Normalization

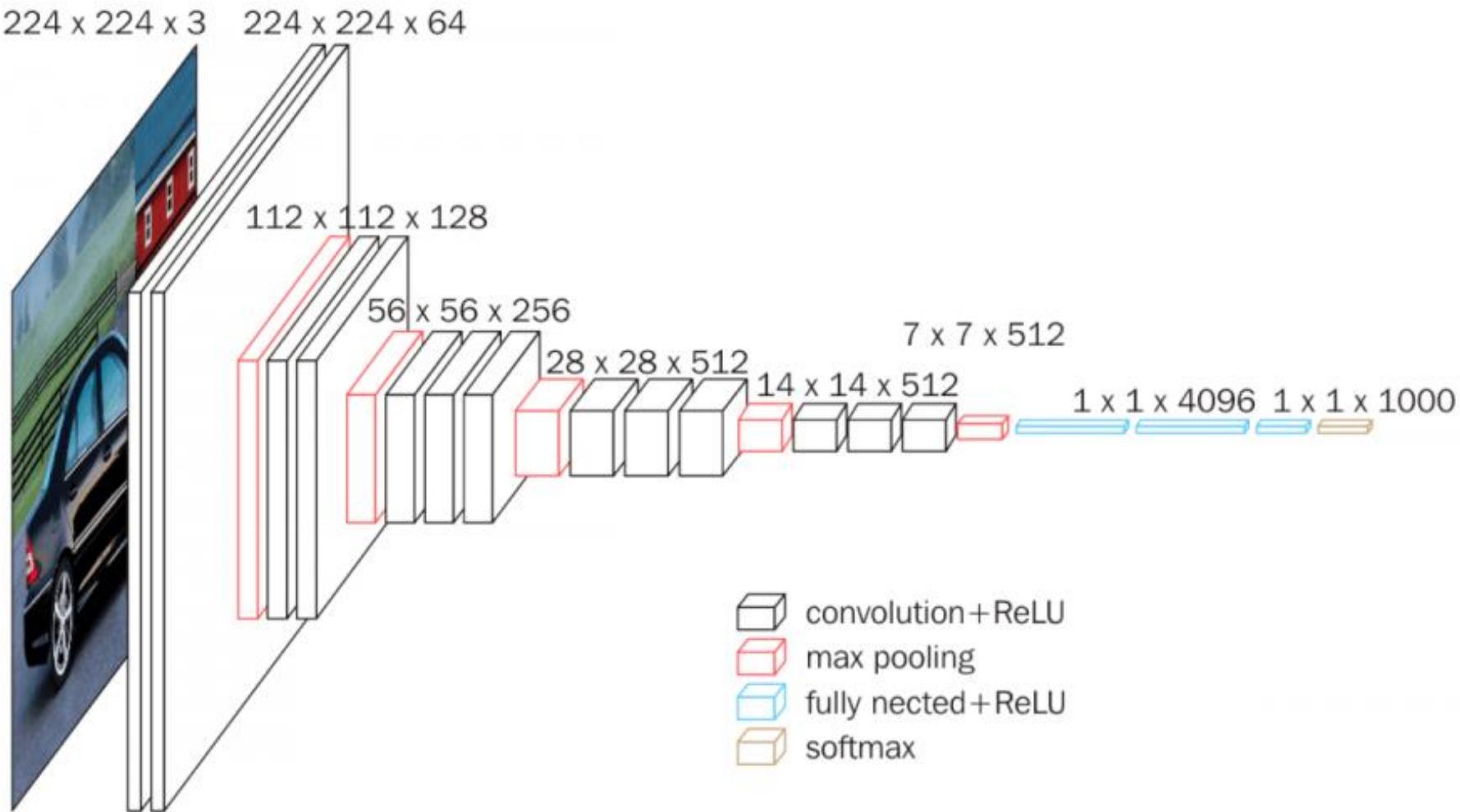
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

Residual connections

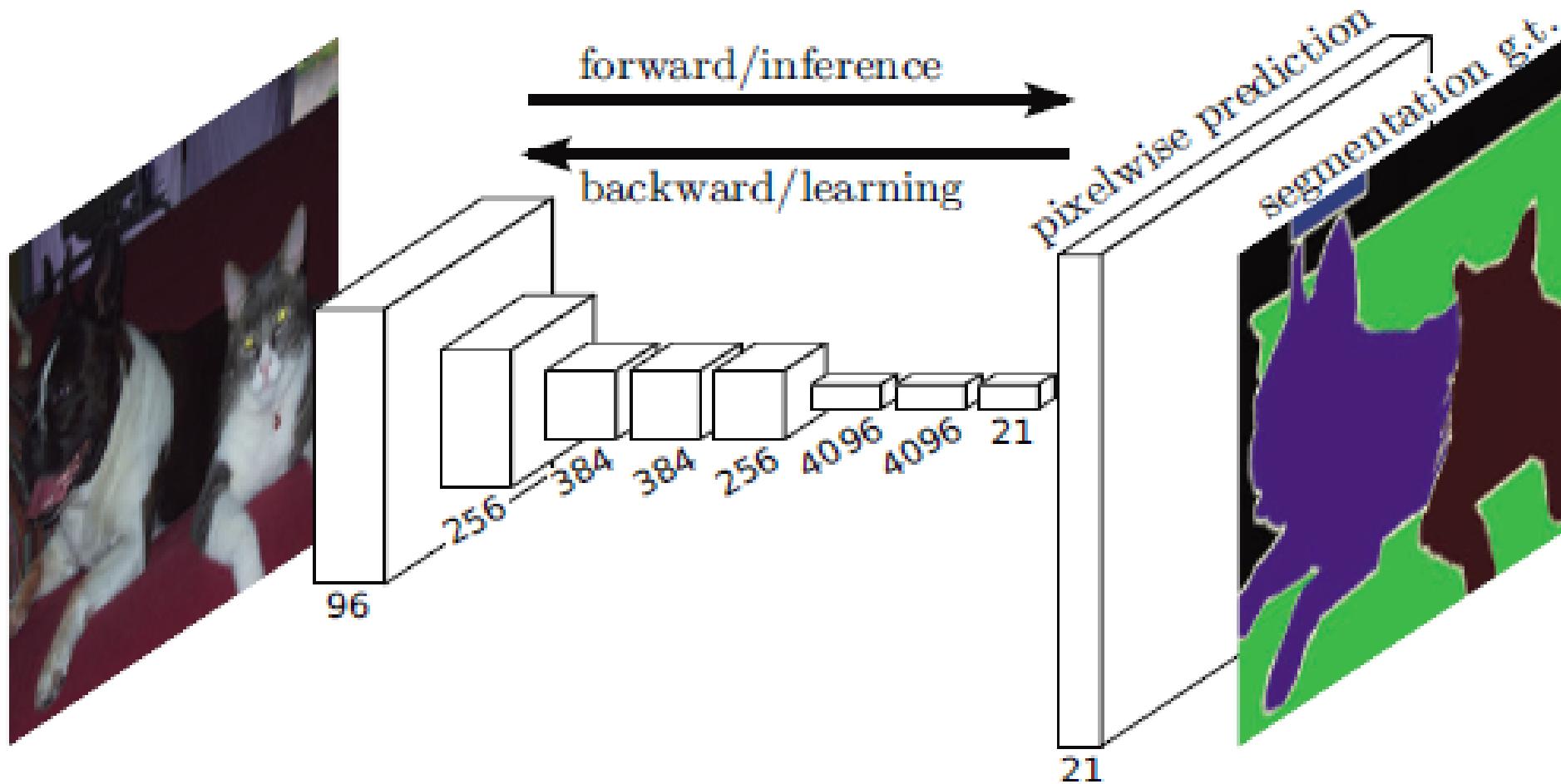


Residual block

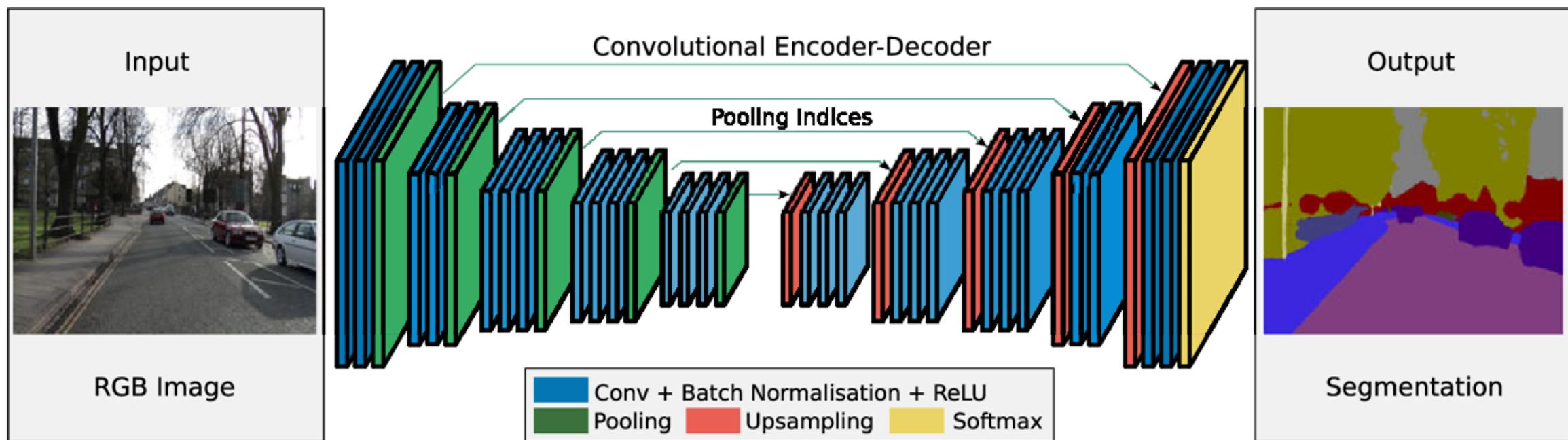
Architecture: Classification networks



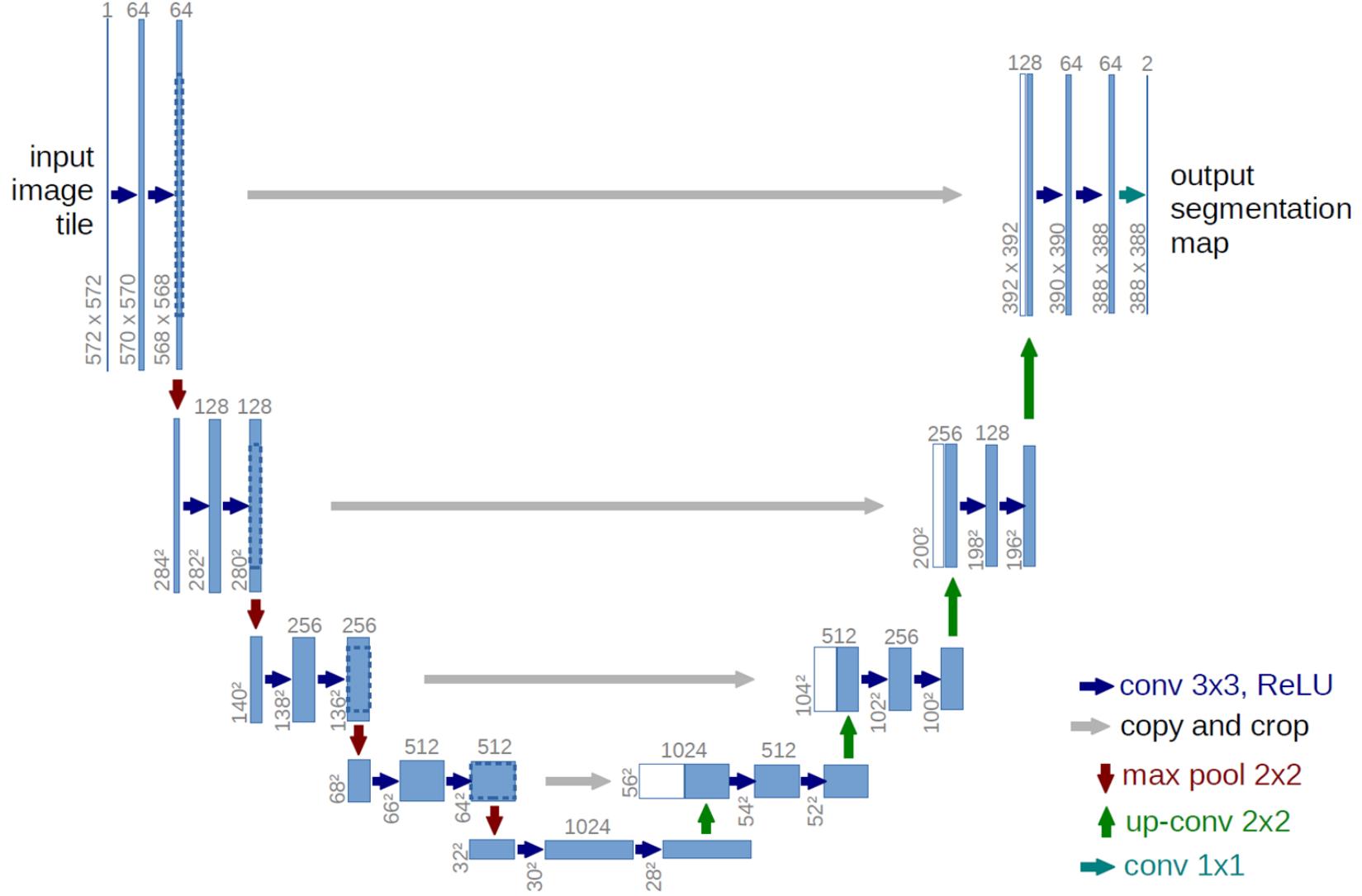
Architecture: Fully convolutional



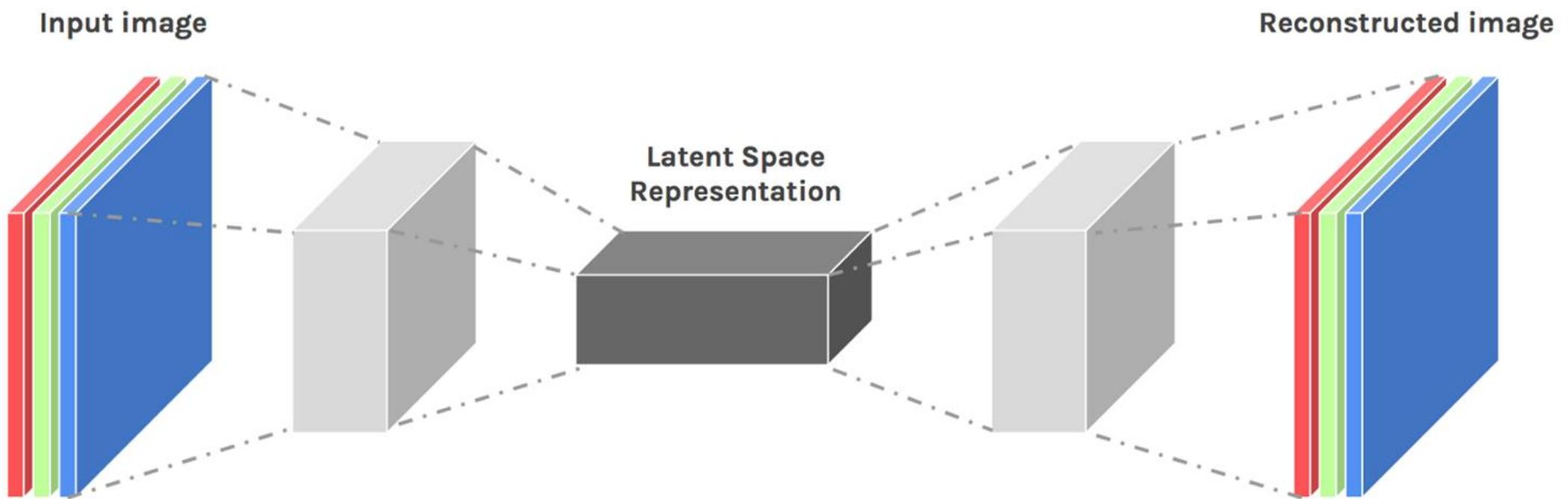
Architecture: Encoder/Decoder



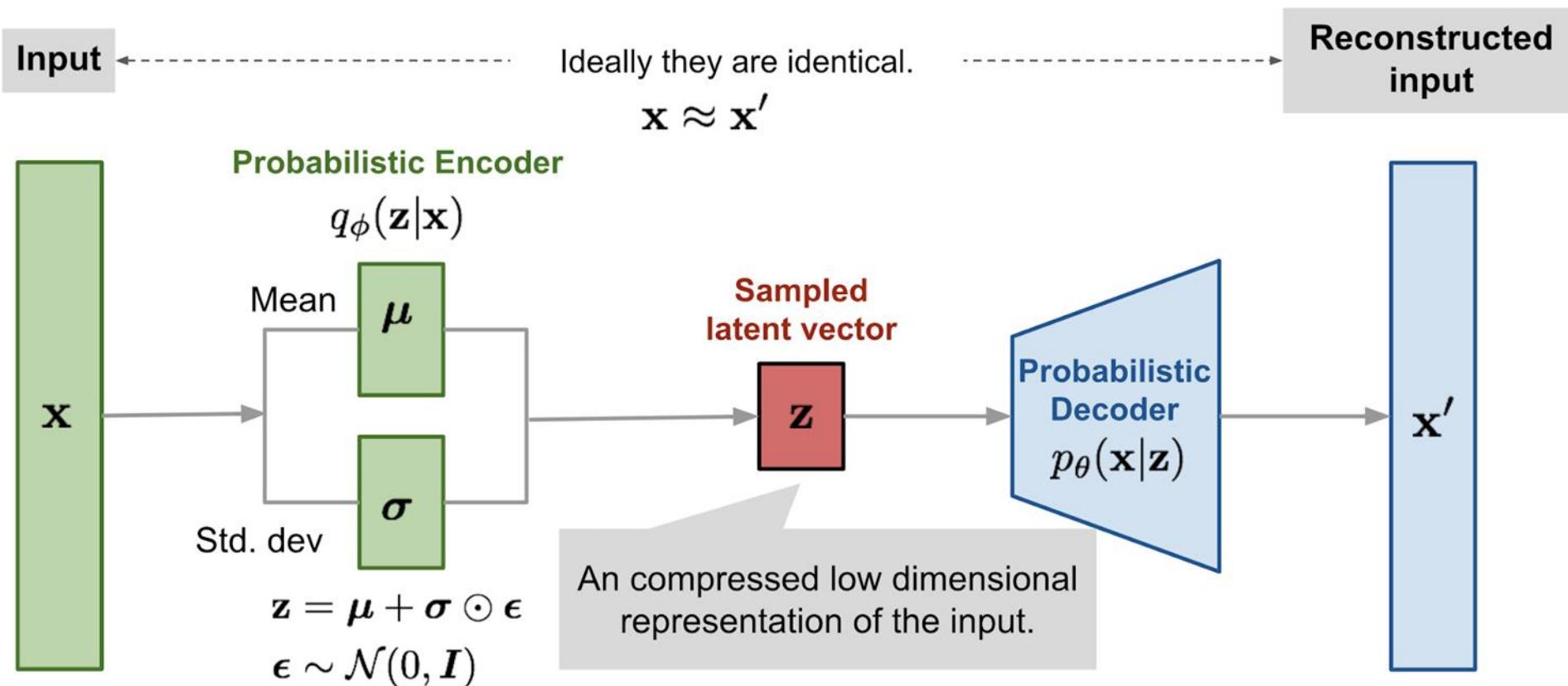
Architecture: U-Net



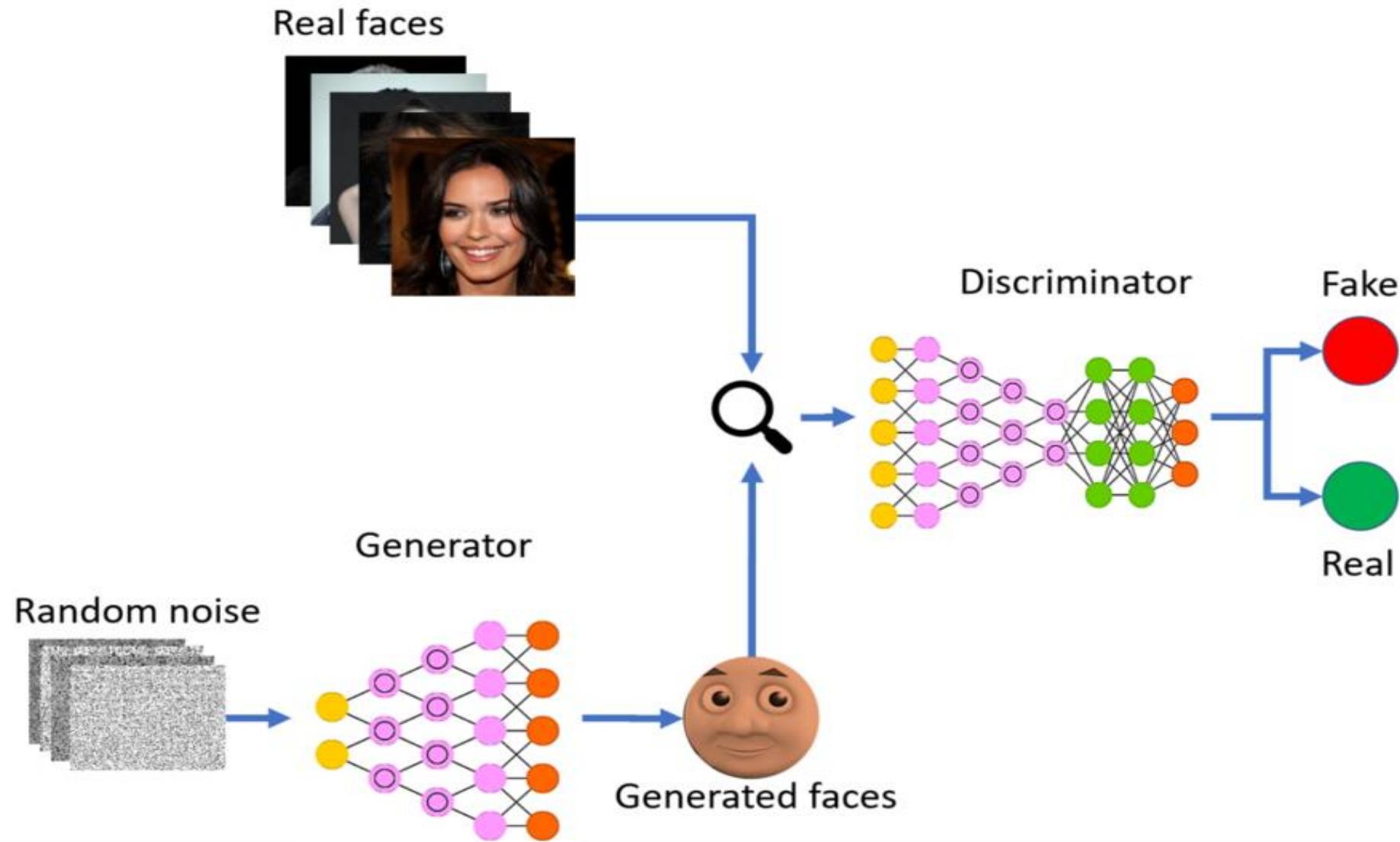
Architecture: Autoencoder



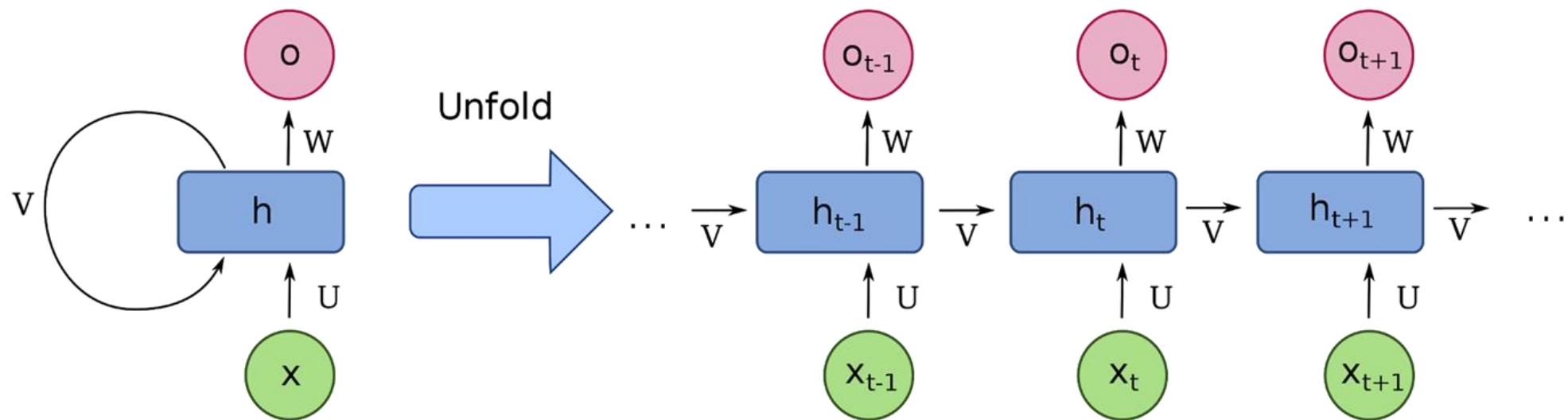
Architecture: Variational Autoencoder



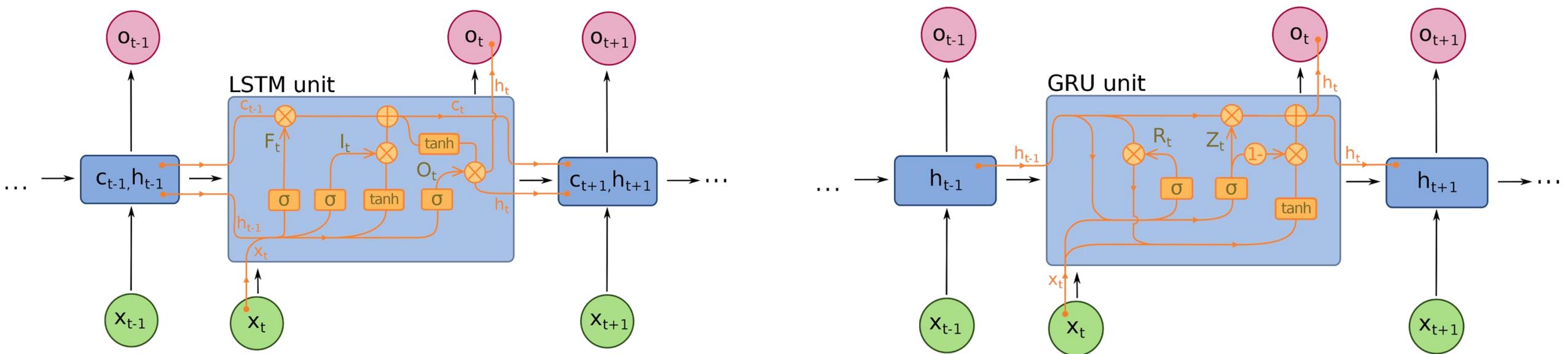
Architecture: Generative adversarial network



Architecture: Recurrent neural network



Architecture: LSTM/GRU



Architecture: Transformers

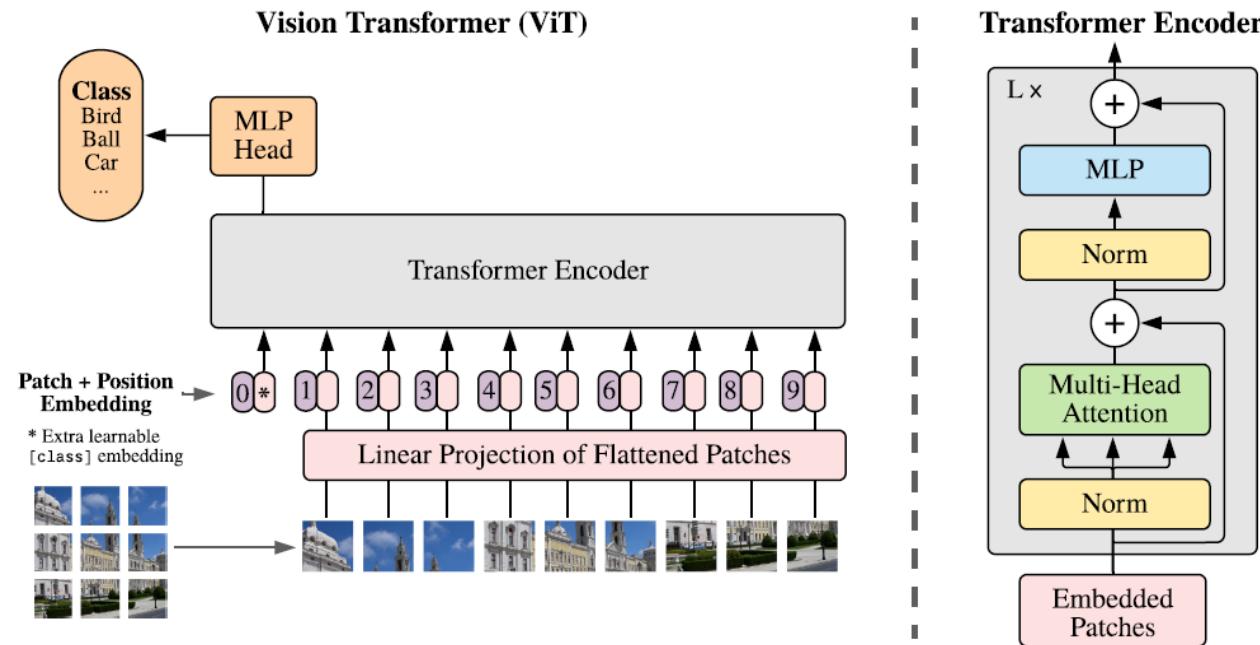


Figure 1: Model overview. We split an image into fixed-size patches, linearly embed each of them, add position embeddings, and feed the resulting sequence of vectors to a standard Transformer encoder. In order to perform classification, we use the standard approach of adding an extra learnable “classification token” to the sequence. The illustration of the Transformer encoder was inspired by Vaswani et al. (2017).