

Overview of Modern ML Applications: Convolutional Neural Network (CNN)

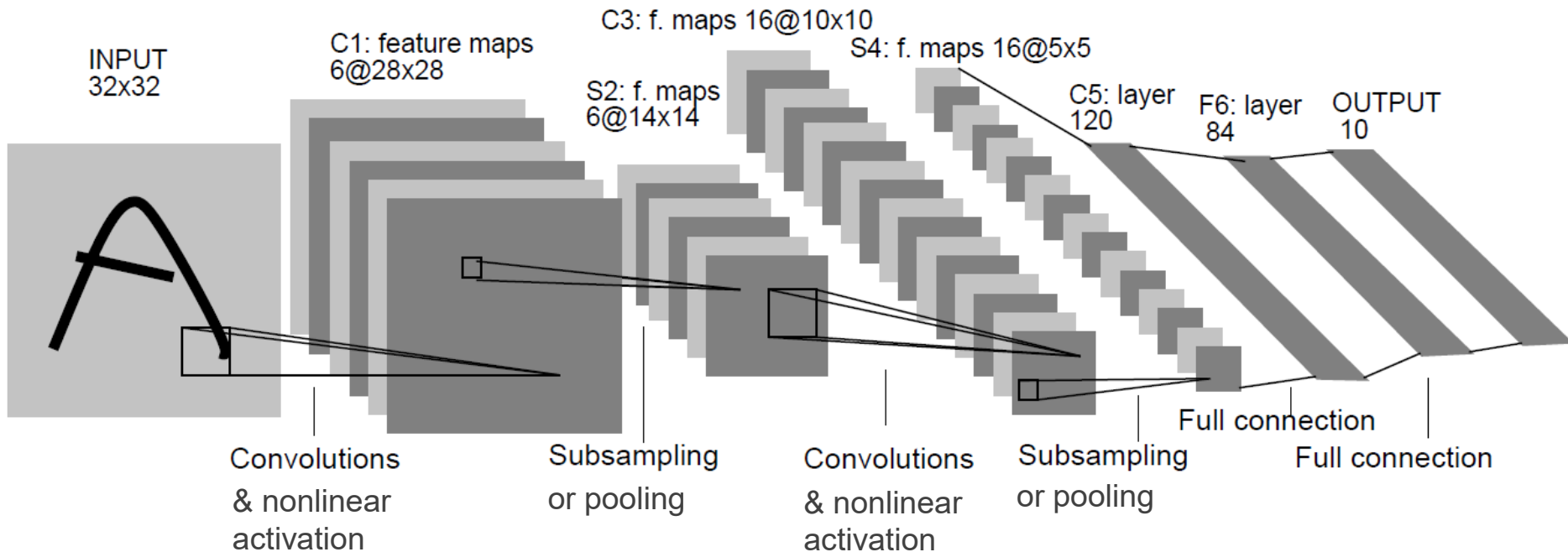
Learning objectives

- Describe the structure of CNN
- Build and train simple CNNs using a deep learning package

(Ref: Ch 9 of [Goodfellow et al. 2016](#))

Convolutional Neural Network (CNN)

The **single** most important technology that fueled the rapid development of **deep learning** and **big data** in the past decade.



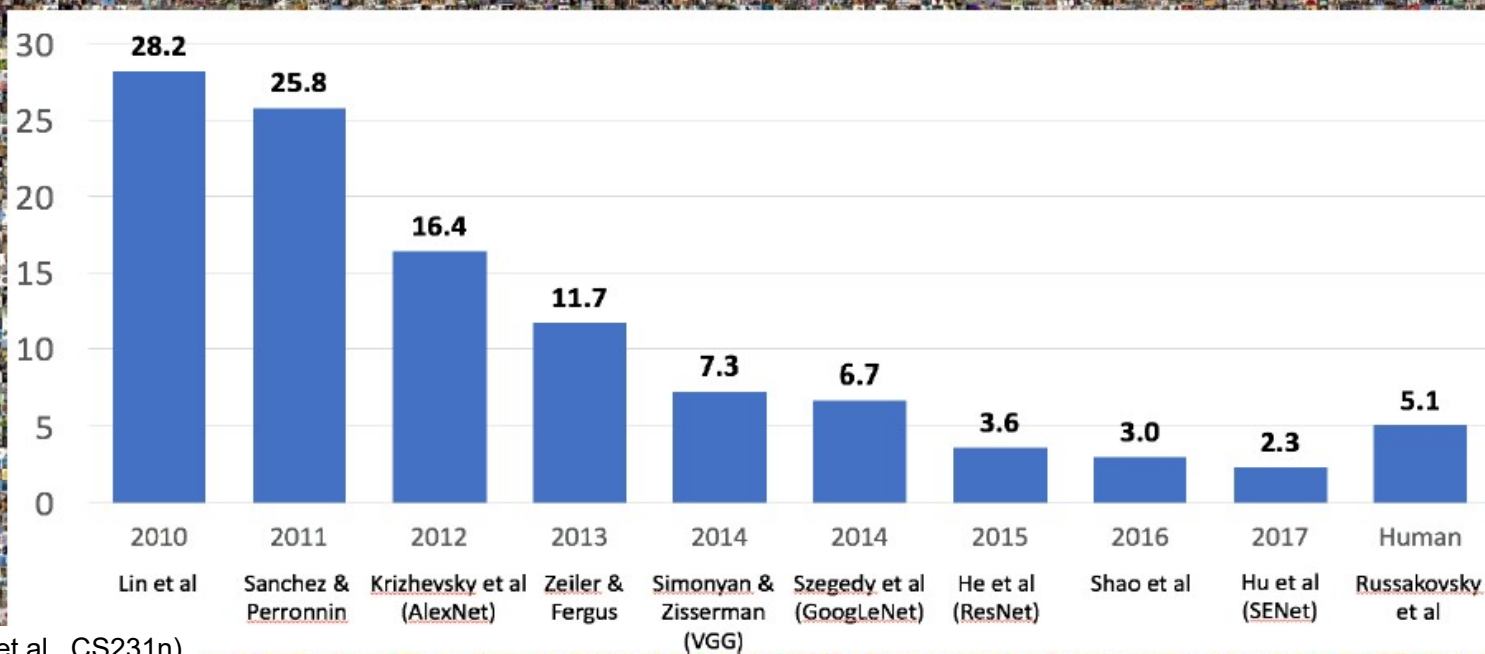
Why is Deep Learning so Successful?

- 1. Improved model:** convolutional layer, more layers (“deep”), simpler activation (i.e., ReLU), skip/residual connection (i.e., ResNet), attention (i.e., Transformer)
 - 2. Big data:** huge dataset, transfer learning
 - 3. Powerful computation:** graphical processing units (GPUs)
- ◆ Example of big data: ImageNet (22K categories, 15M images)



IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:
1,000 object classes
1,431,167 images



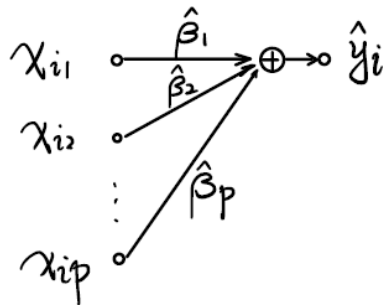
Linear Model to Neural Network

Recall linear model w/ multiple predictors / features / inputs.

$$\underbrace{y_i}_{\text{true output}} = \sum_{j=1}^p x_{ij} \beta_j + e_i = \underbrace{[\beta_1, \dots, \beta_p]}_{\text{true weights}} \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix} + e_i, \quad i=1, \dots, n.$$

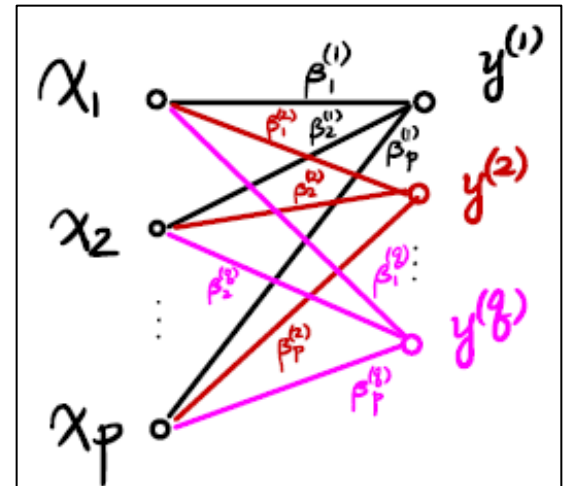
$$\underbrace{\hat{y}_i}_{\text{predicted output}} = \sum_{j=1}^p x_{ij} \hat{\beta}_j = \underbrace{[\hat{\beta}_1, \dots, \hat{\beta}_p]}_{\text{estimated weights}} \begin{bmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{bmatrix}, \quad i=n+1, \dots, n+m.$$

Graphically we have:



① Use multiple linear models

② Simplify the notations.



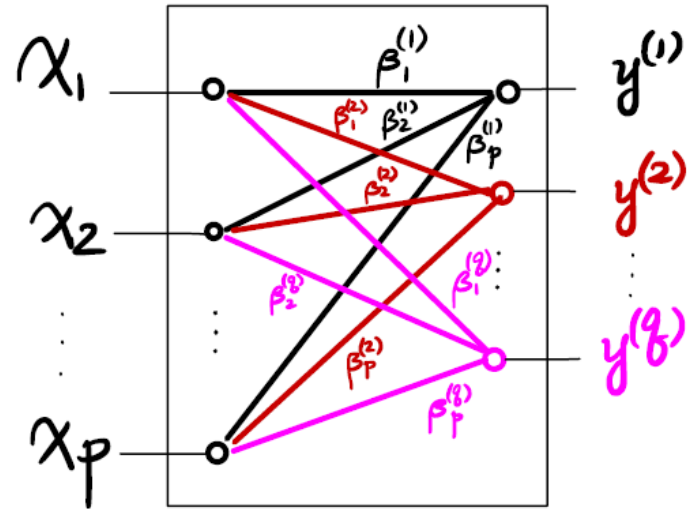
Fully-Connected Layer for 1D Signal

$$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(g)} \end{bmatrix} = \begin{bmatrix} \beta_1^{(1)} & \beta_2^{(1)} & \dots & \beta_p^{(1)} \\ \beta_1^{(2)} & \beta_2^{(2)} & \dots & \beta_p^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_1^{(g)} & \beta_2^{(g)} & \dots & \beta_p^{(g)} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$$

layer output, $y \in \mathbb{R}^g$

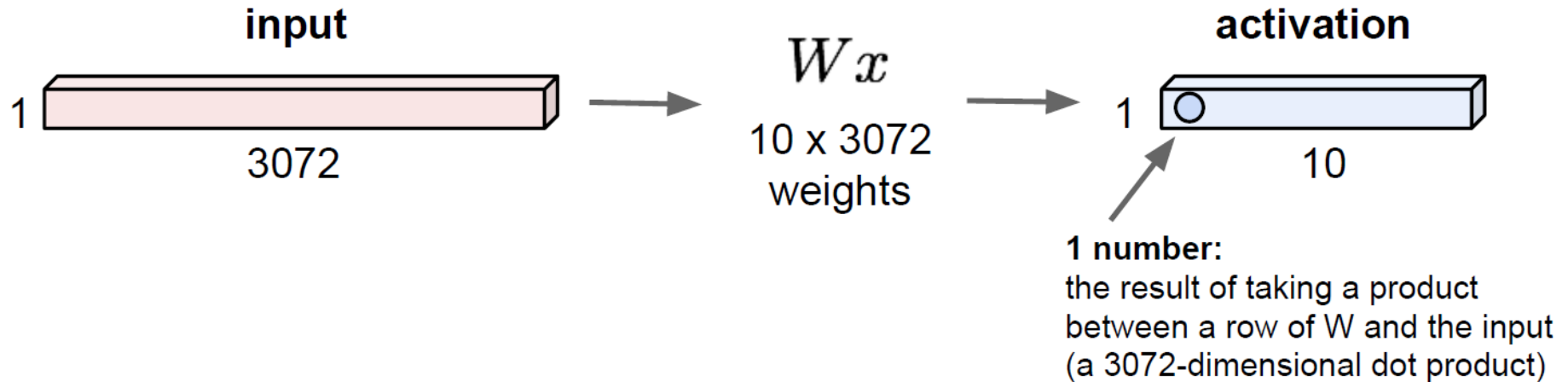
dense weight matrix
 $B \in \mathbb{R}^{g \times p}$

layer input, $x \in \mathbb{R}^p$



Fully-Connected Layer for RGB Image

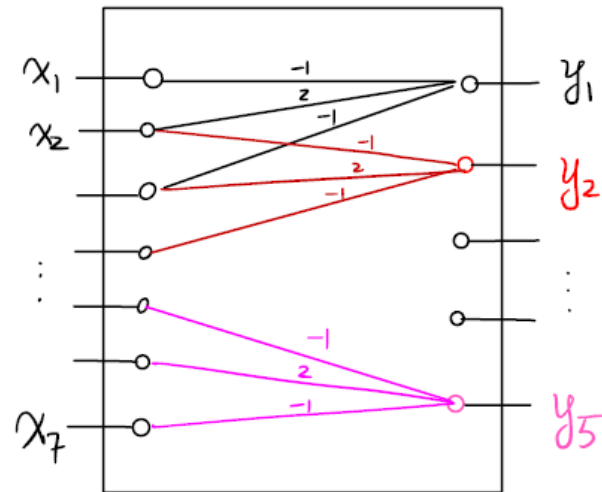
32x32x3 image -> stretch to 3072 x 1



Convolutional Layer for 1D Signal

$$\begin{bmatrix} y_1 \\ \vdots \\ y_5 \end{bmatrix} = \begin{bmatrix} -1 & 2 & -1 & & & & \\ & -1 & 2 & -1 & & & \\ & & \vdots & & -1 & 2 & -1 & \vdots \\ & & & & & -1 & 2 & -1 \\ & & & & & & -1 & 2 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_7 \end{bmatrix}$$

Sparse weight matrix



Input

x_1	x_2	x_3	x_4	x_5	x_6	x_7
-------	-------	-------	-------	-------	-------	-------

 length 7

Convolution/
filter mask

*		
-1	2	-1

 → length 3

Output

y_1	y_2	y_3	y_4	y_5
-------	-------	-------	-------	-------

 length $7 - (3 - 1) = 5$
(w/o boundary elements)

Convolutional Layer for 2D Matrix/Image

x_{11}			x_{15}
x_{21}			
x_{51}			x_{55}

Input image

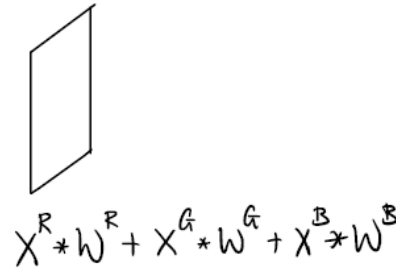
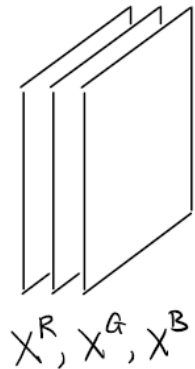
$$* \begin{array}{|c|c|} \hline \frac{1}{4} & \frac{1}{4} \\ \hline \frac{1}{4} & \frac{1}{4} \\ \hline \end{array} =$$

filter mask

y_{11}			y_{14}
y_{21}			
y_{41}			y_{44}

Activation map

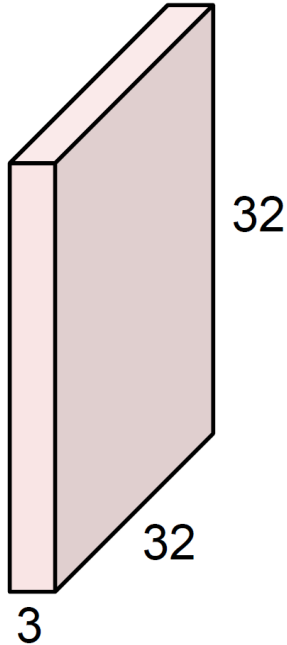
2D Convolution



Multiple color channels need multiple filter masks

Convolutional Layer for RGB Image

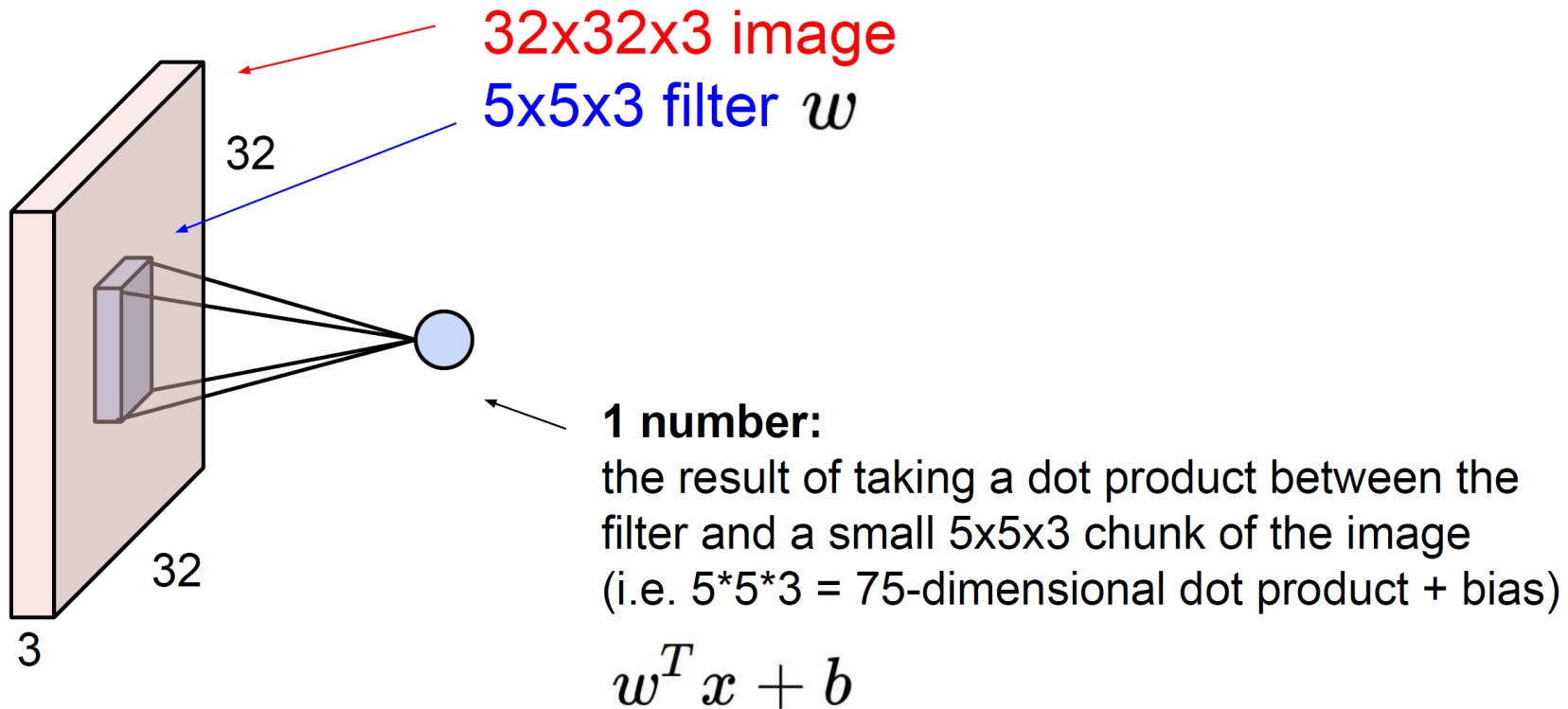
32x32x3 image



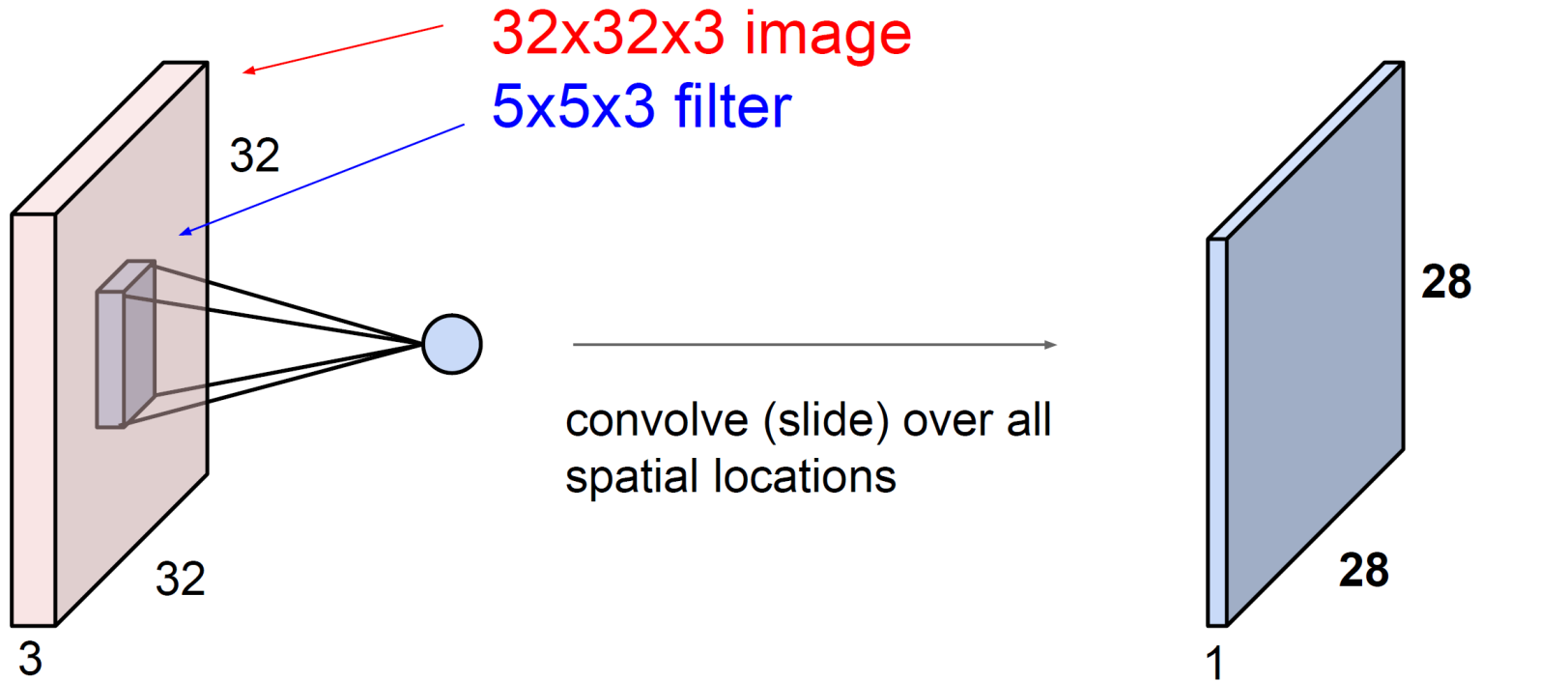
5x5x3 filter



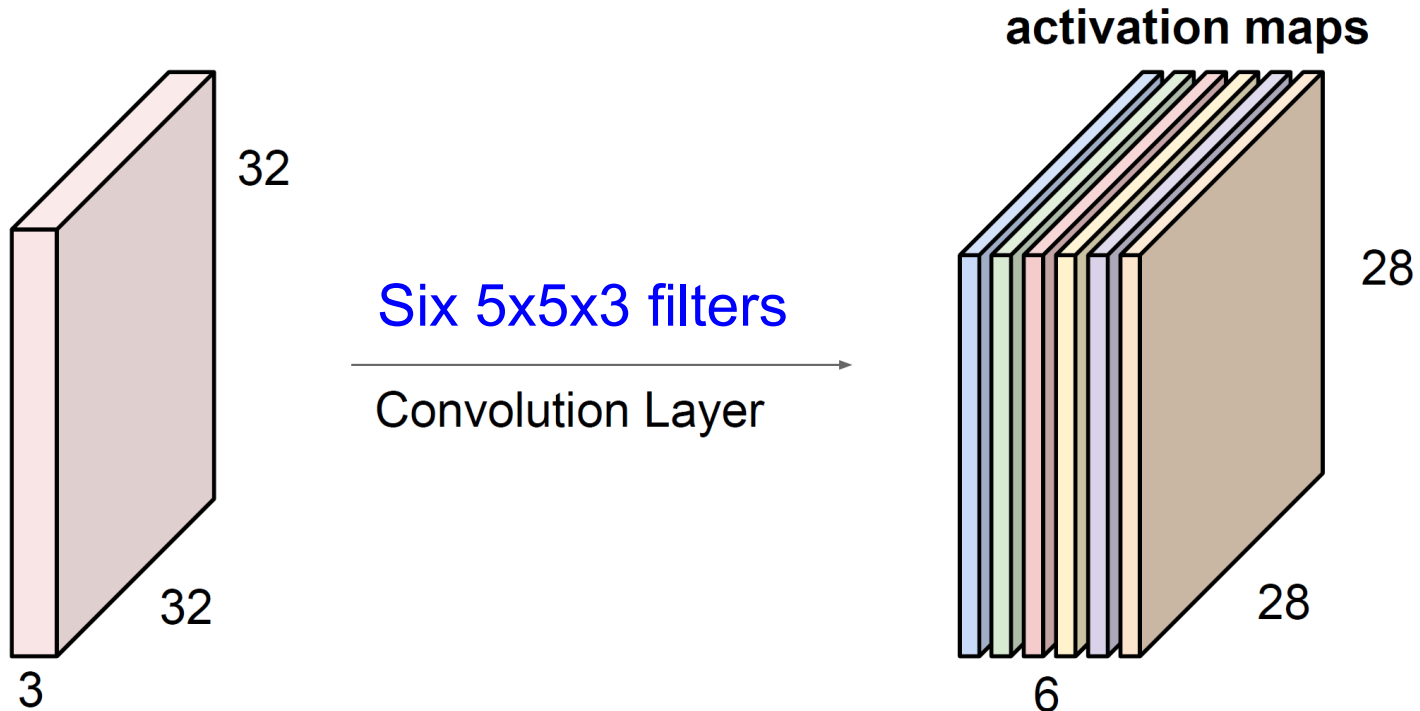
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”



A closer look at spatial dimensions:

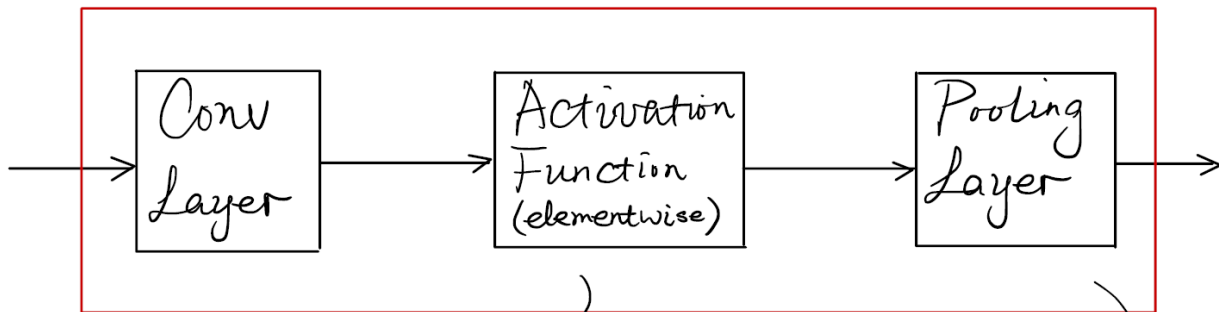


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

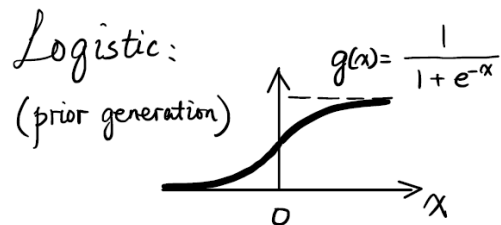
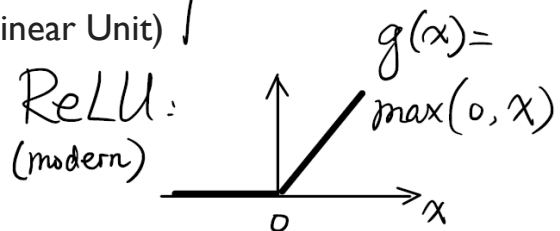


We stack these up to get a “new image” of size 28x28x6!

Building Block for Modern CNN



(Rectified Linear Unit)



Ex: $\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix} : \{x_{ij}\}_{i,j=1}^2$

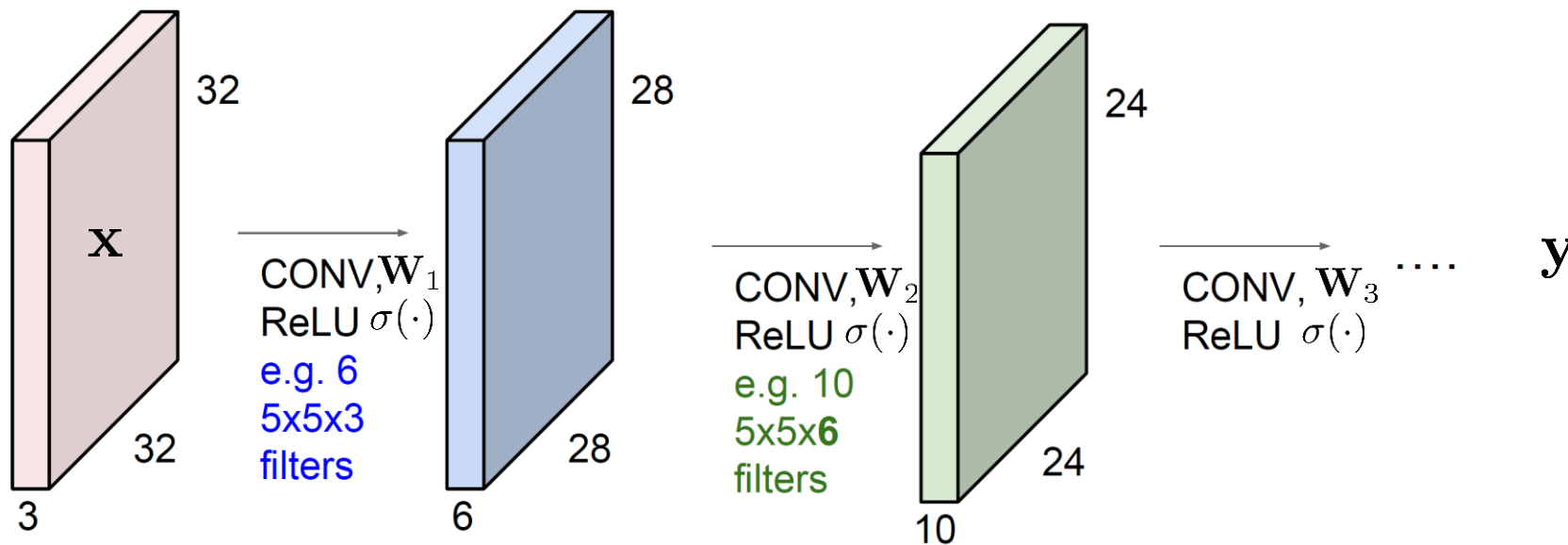
Max pooling:

$$g(\{x_{ij}\}) = \max(\{x_{ij}\})$$

Average pooling:

$$g(\{x_{ij}\}) = \frac{1}{|\{x_{ij}\}|} \sum_{i,j} x_{ij}$$

CNN is composed of a sequence of convolutional layers, interspersed with activation functions (ReLU, in most cases).

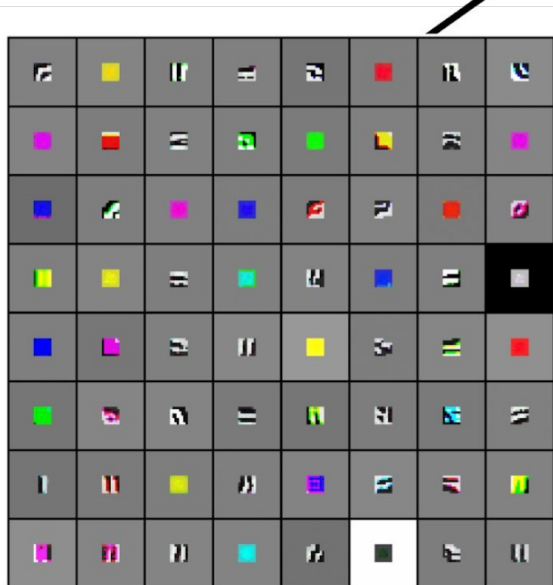


$$\mathbf{y} = \cdots \sigma \left(\mathbf{W}_3 \sigma \left(\mathbf{W}_2 \sigma \left(\mathbf{W}_1 \mathbf{x} \right) \right) \right) \cdots$$

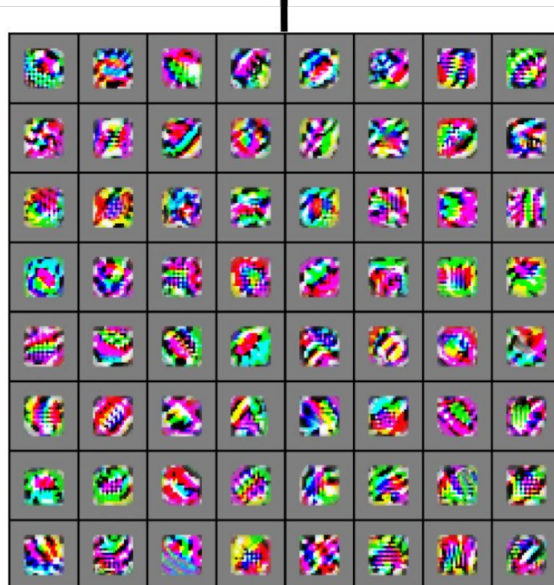
Source of nonlinearity, ReLU: $\sigma(x) := \max(0, x)$

[Zeiler and Fergus 2013]

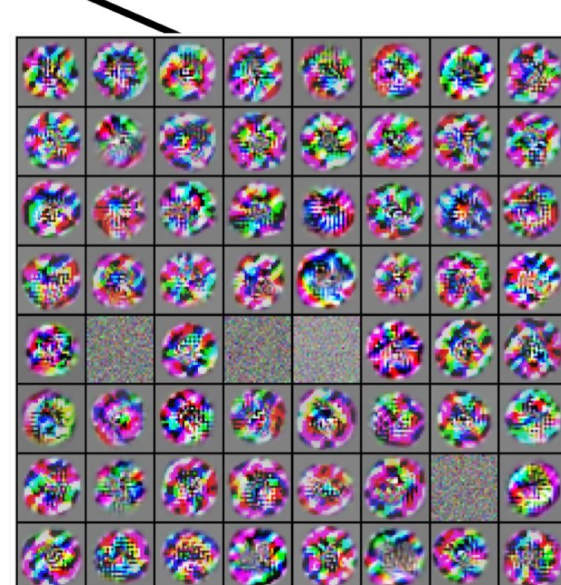
Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].



VGG-16 Conv1_1



VGG-16 Conv3_2



VGG-16 Conv5_3

IMAGENET Large Scale Visual Recognition Challenge

Year 2010

NEC-UIUC



Dense descriptor grid:
HOG, LBP

Coding: local coordinate,
super-vector

Pooling, SPM

Linear SVM

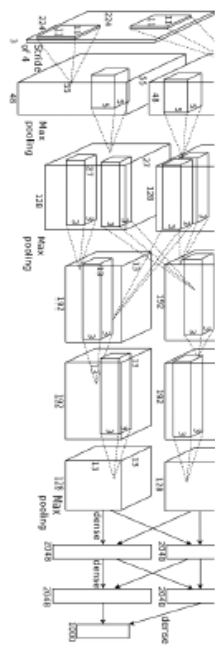
[Lin CVPR 2011]

Lion image by Swissfrog is licensed under [CC BY 3.0](https://creativecommons.org/licenses/by/3.0/)

AlexNet

Year 2012

SuperVision



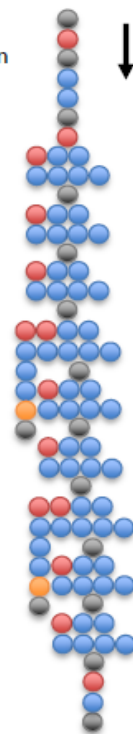
[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Year 2014

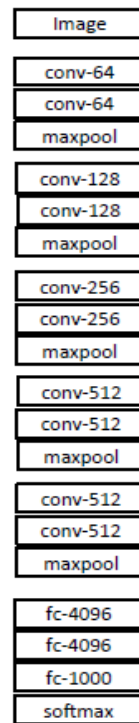
GoogLeNet

- Pooling
- Convolution
- Softmax
- Other



[Szegedy arxiv 2014]

VGG

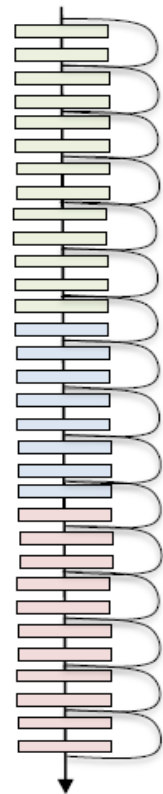


[Simonyan arxiv 2014]

ResNet

Year 2015

MSRA

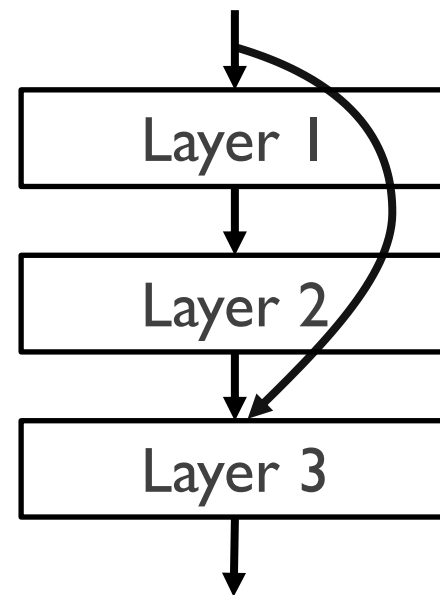


[He ICCV 2015]

(Fei-Fei Li et al., CS231n) 17

Residual Neural Network (ResNet) [Kaiming He et al., 2015]

- ◆ **Skip connections** or *shortcuts* are added.
- ◆ They can
 - ★ avoid “vanishing gradients”, and
 - ★ make optimization landscape flatter.
- ◆ From Taylor expansion perspective, the neural network only learns the higher-order error terms beyond the linear term \mathbf{x} .
- ◆ Has interpretations in PDE.
- ◆ Preferred modern NN structure.



$$y = \cdots \sigma \left(\mathbf{W}_3 \underbrace{[\mathbf{x} + \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))]}_{g(\mathbf{x})} \right) \cdots$$

When Output is Categorical / Qualitative

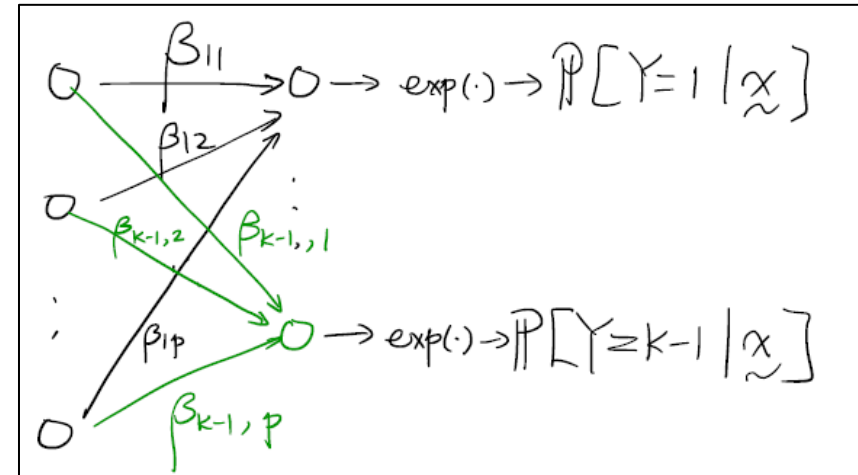
◆ A **softmax layer** is needed:

◆ Softmax function:

$$\sigma_i(\underline{z}) = \frac{e^{\beta z_i}}{\sum_{j=1}^K e^{\beta z_j}}$$

◆ Ex:

$$\begin{aligned} K=2 \quad \sigma_1 &= \frac{e^{\beta z_1}}{e^{\beta z_1} + e^{\beta z_2}} \\ &= \frac{1}{1 + e^{\beta(z_2 - z_1)}} \end{aligned}$$



When β very large,

$z_2 > z_1$ leads to $\begin{cases} \sigma_1 = 0 \\ \sigma_2 = 1 \end{cases}$

Winner takes all!

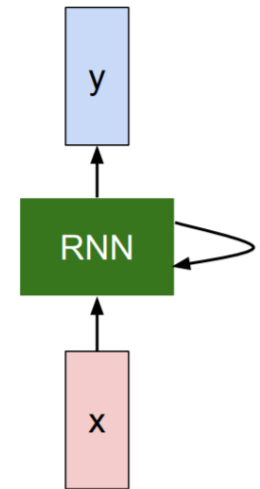
Other Essential Aspects of CNN

- ◆ Due to time constraints, this overview lecture covered only the structural elements of CNNs. Other essential aspects are:
 - ✦ Cost function, e.g., MSE, cross entropy. (will cover if time permits.)
 - ✦ How to train CNNs (estimate the weights), i.e., *backpropagation*. (will cover if time permits.)
 - ✦ Practical training considerations including
 - How to determine *number of hidden units/channels* to be used,
 - How to tune *learning rate* and *batch size*, and
 - When to stop training (number of *epochs*).
- ◆ For a more complete treatment on CNN, refer to the dedicate courses such as [CS231n CNNs for Visual Recognition](#).

Machine Learning (ML) and Data Science (DS)

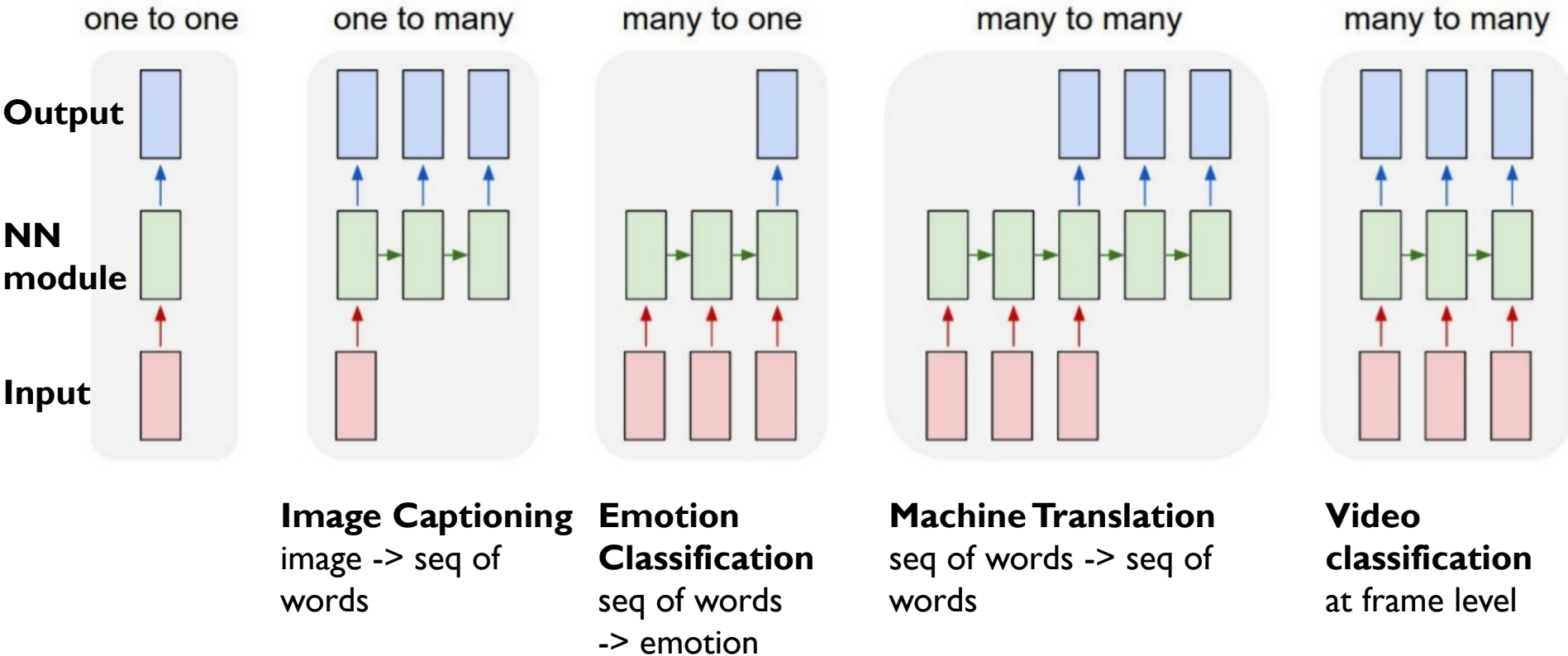
- ◆ Follow-up machine learning / data science courses:
 - ECE 542 Neural Nets and Intro to Deep Learning
 - ECE 592-6I Data Science
 - ECE 759 Pattern Recognition and Machine Learning
 - ECE 763 Computer Vision
 - ECE 792-4I Statistical Foundations for Signal Processing & Machine Learning
 - ★ Any courses/videos on YouTube, Coursera, etc.
- ◆ State-of-the-art theory & applications: ICML, NeurIPS, ICLR, AAAI
- ◆ Data science competitions: [kaggle.com](https://www.kaggle.com)
- ◆ Programming languages for ML/DS: Python, R, Matlab

Overview of Modern ML Applications: Recurrent Neural Network (RNN) and LSTM



A Sequence of **Identical** Neural Network Modules

Force the neural nets (in green) to be the same to lower the complexity!



Recurrent Neural Network (RNN): Definition

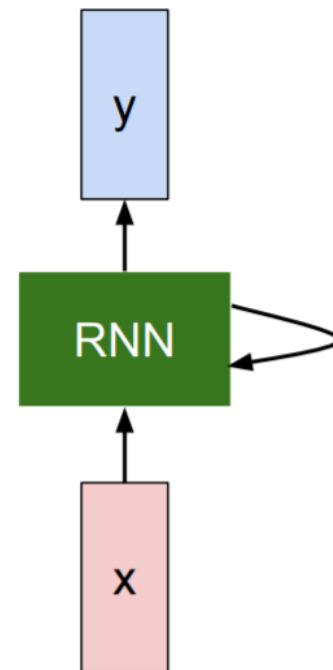
We can process a sequence of vectors \mathbf{x} by applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

h_t : state, x_t : input

f_W : neural network

$$y_t = W_{hy}h_t$$



Recurrent Neural Network (RNN): Implementation

$$h_t = f_W(h_{t-1}, x_t)$$

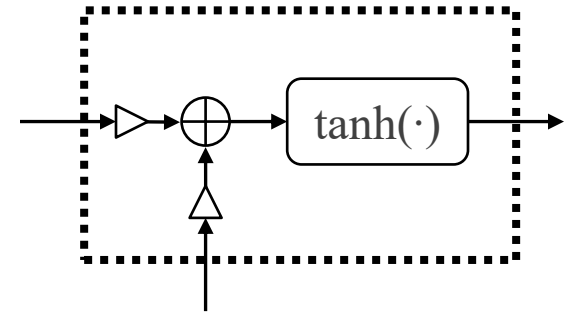


$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

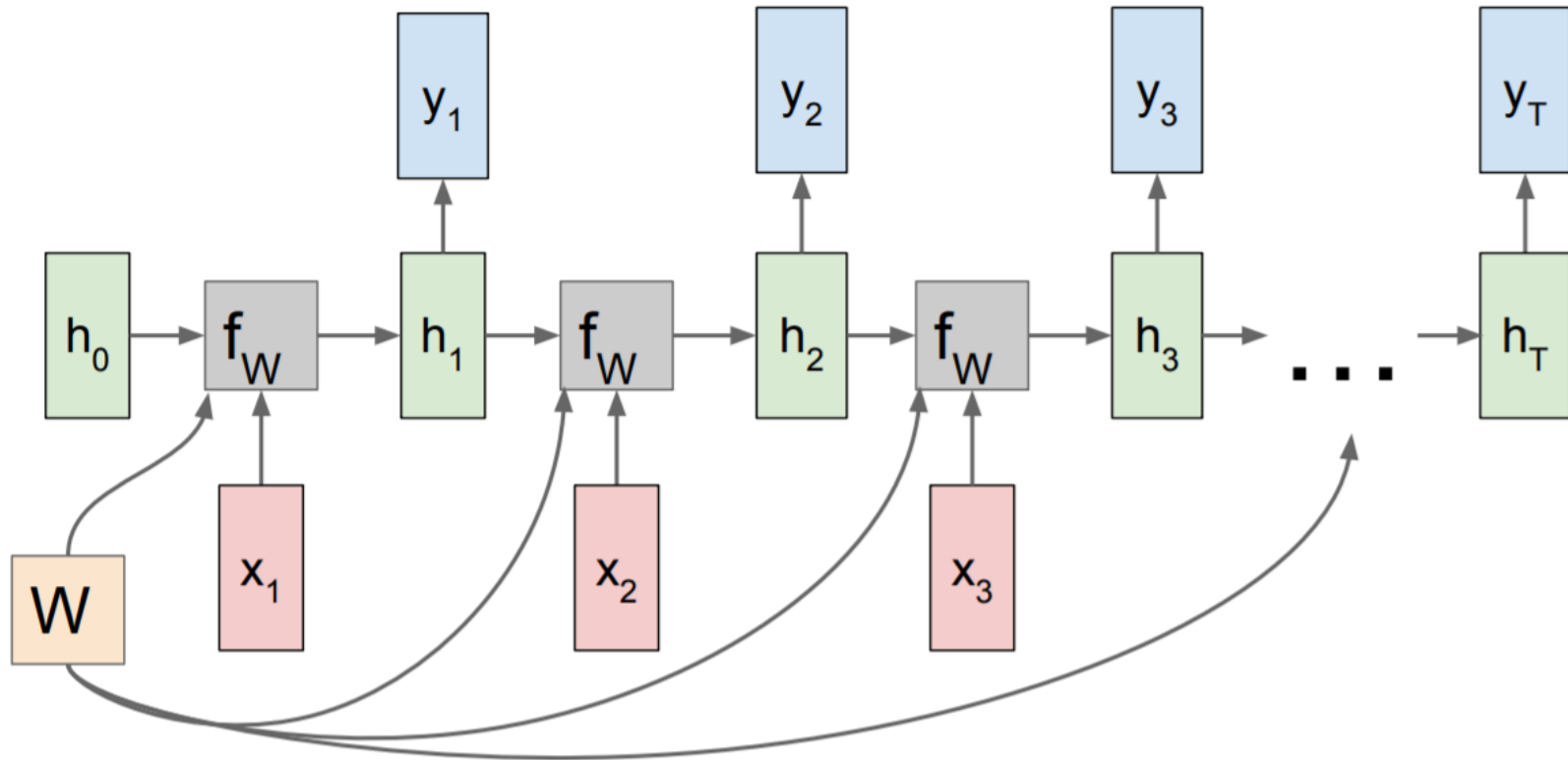
f_W is implemented via

- linear transforms W_{hh} and W_{xh} and
- elementwise **nonlinear** function $\tanh(\cdot)$

Diagram for f_W .
Can you label the details?



Recurrent Neural Network (RNN): Unrolled

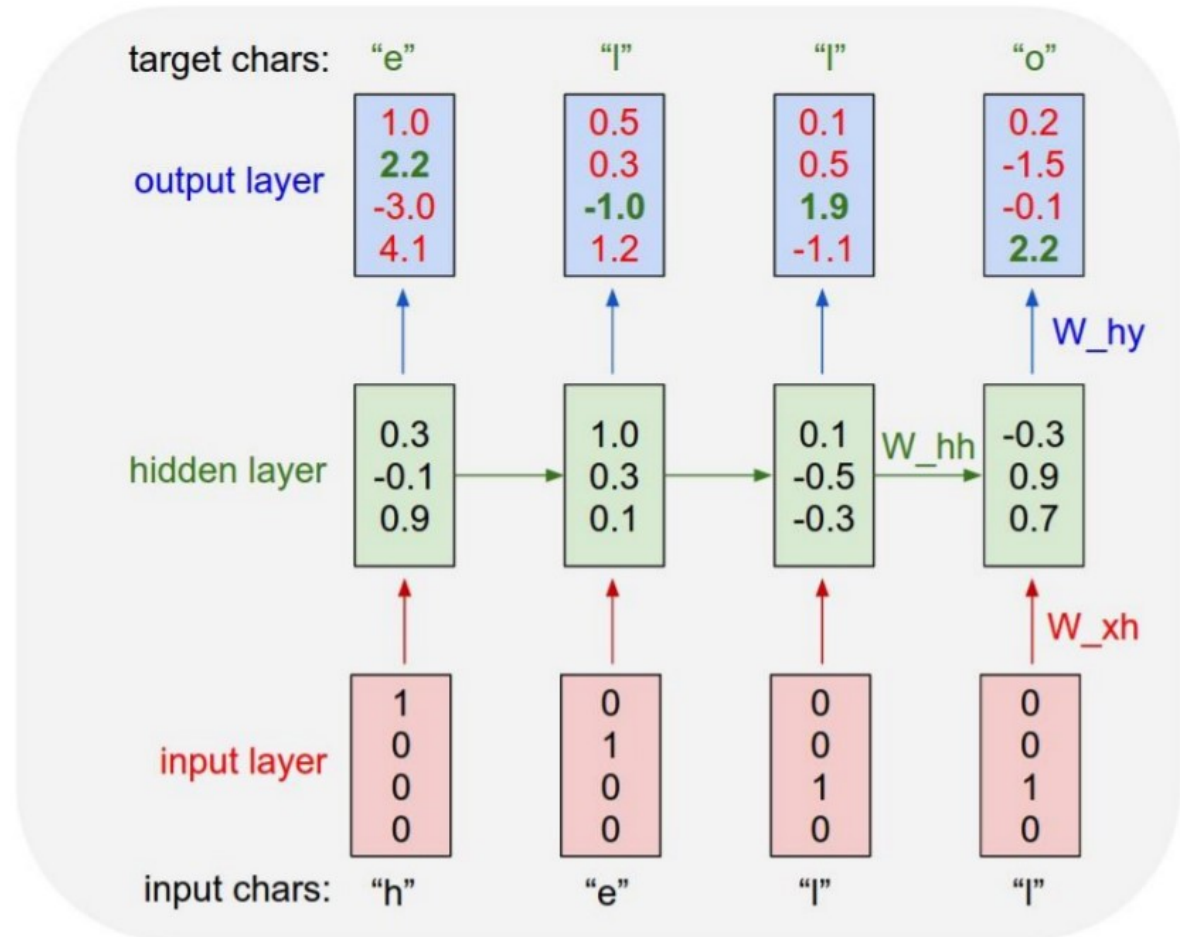


Example: Character-Level Language Model

Vocabulary:
[h, e, l, o]

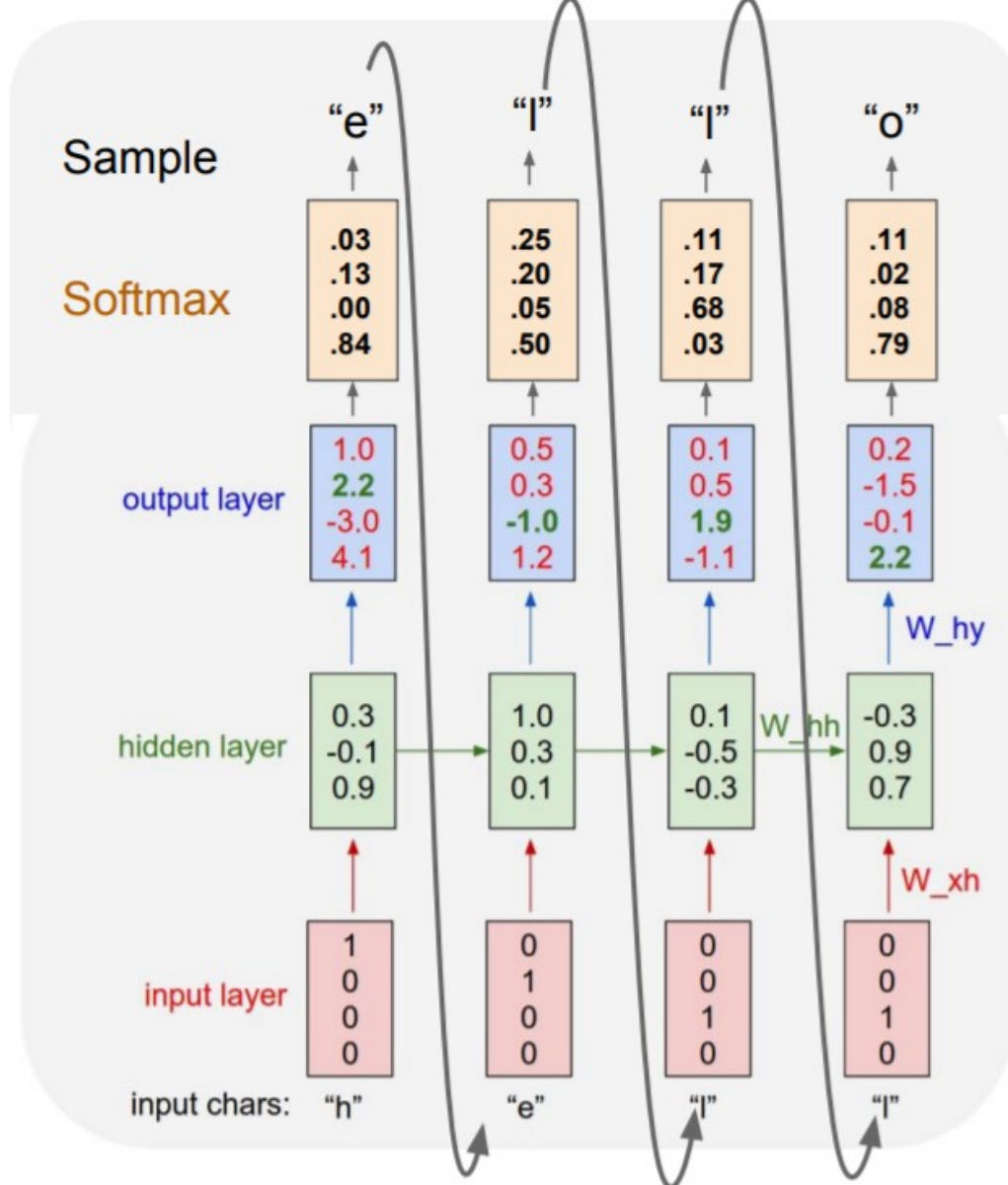
Embedding:

Example training
sequence:
"hello"



Vocabulary:
[h, e, l, o]

At test time, sample characters one at a time, feed back to model



Long Short-Term Memory (LSTM) Network

- ◆ RNN has the “vanishing gradient” problem!
- ◆ Resolved by long short-term memory (LSTM) units.

RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

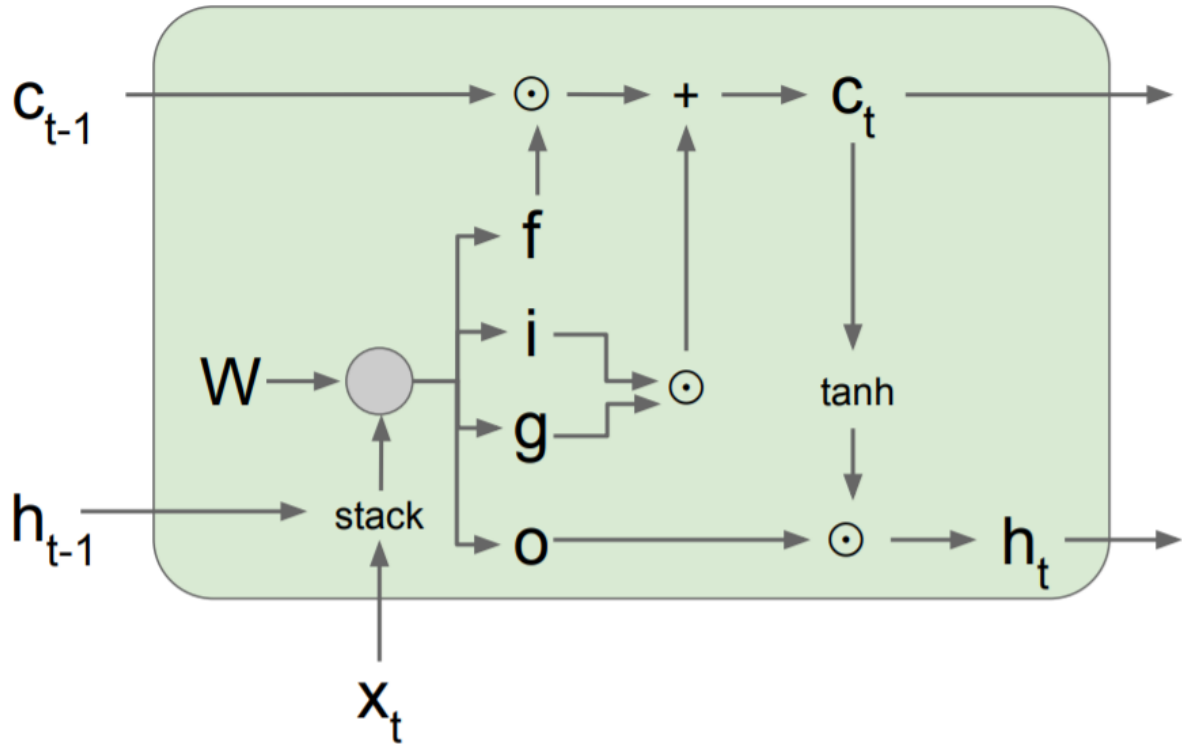
LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

An LSTM Unit [Hochreiter et al., 1997]

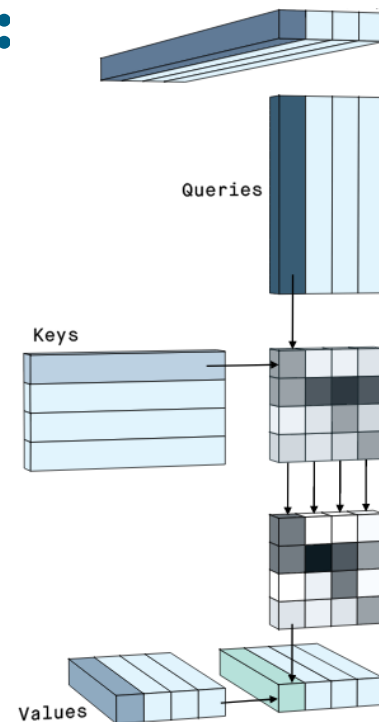


$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Overview of Modern ML Applications: Transformers and BERT



Acknowledgment: Some graphics and slides were adapted from <https://peltarion.com/blog/data-science/self-attention-video>

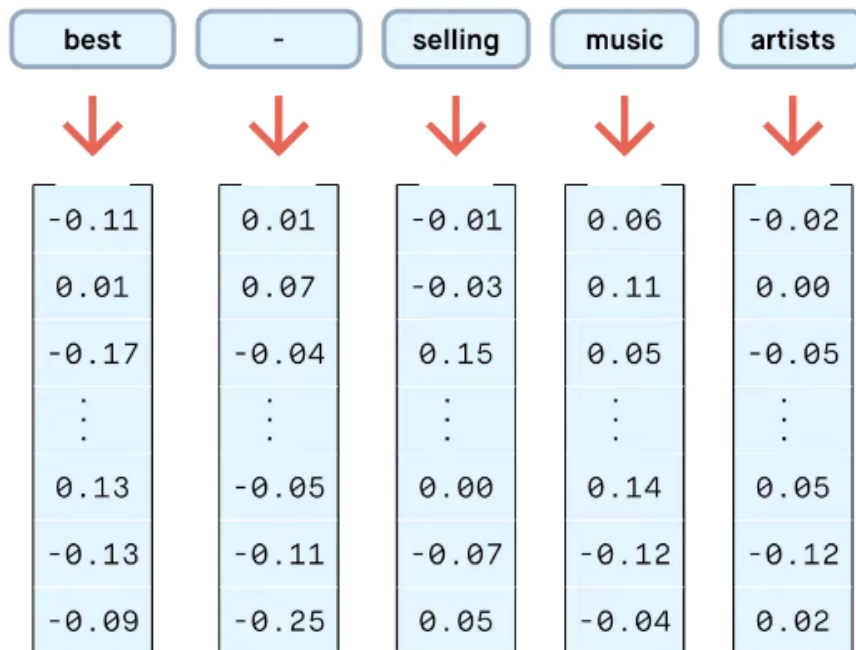
How to make good sense of language?

- ◆ Reading comprehension: If you were Google, what result(s) should you return for “brazil traveler to usa need a visa”?
 1. A webpage on U.S. citizens traveling to Brazil
 2. A webpage of the U.S. embassy/consulate in Brazil

- ◆ Contextualization is the key!
 - ✦ A nice walk by the river **bank**.
 - ✦ Walk to the **bank** and get cash.

Word Embeddings in Natural Language Processing (NLP)

WordPiece tokens:

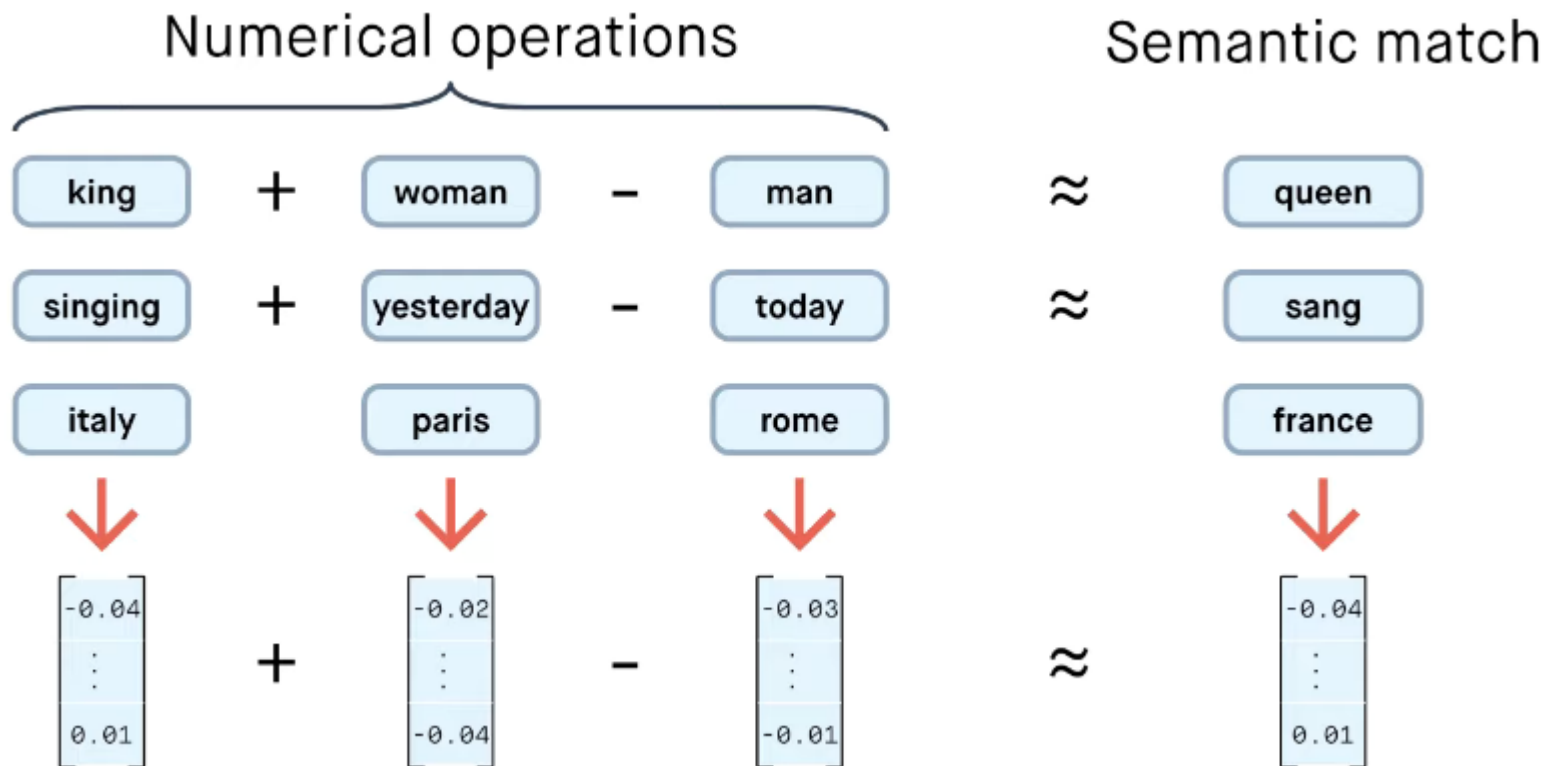


Embedding vectors:

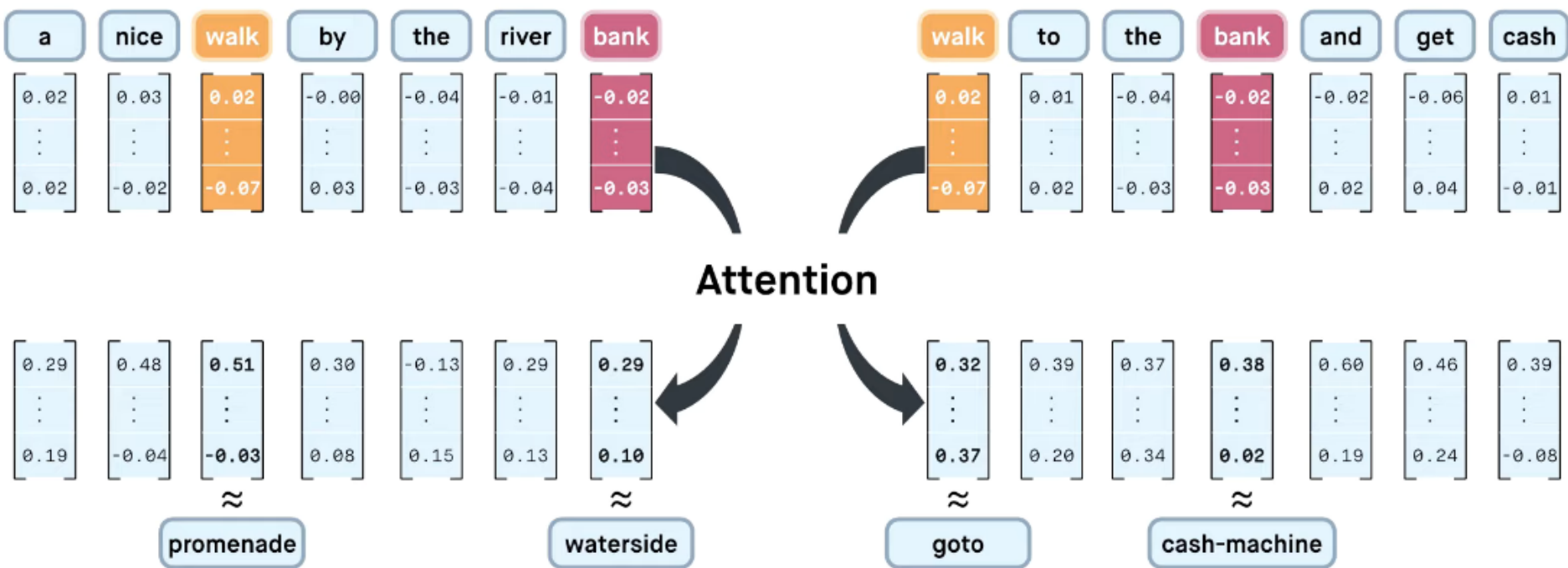
Values are *pretrained*

Embedding examples: Bag of Words (BoW), Word2Vec, ...

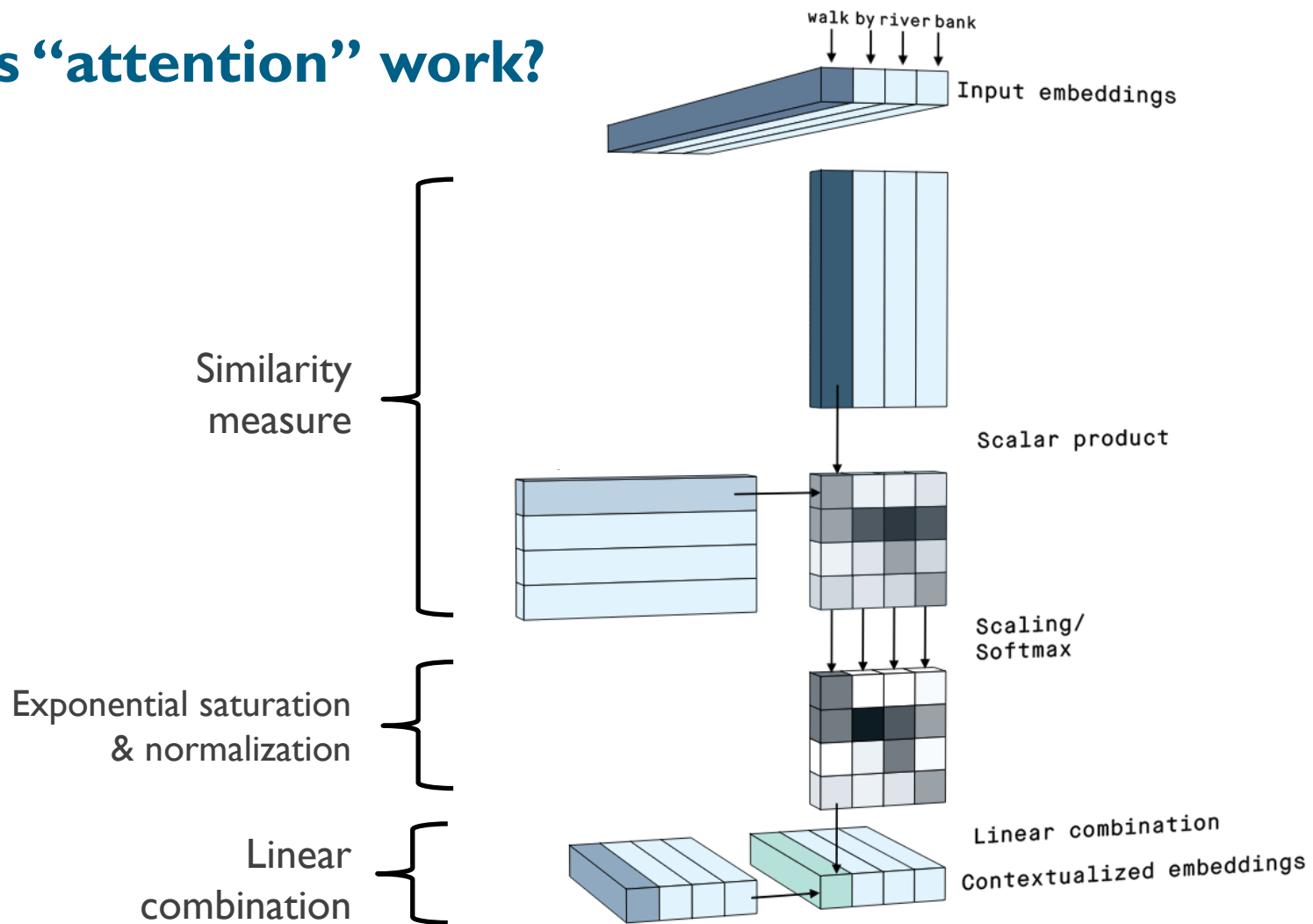
Word Embeddings are Meaningful Under “+” and “-”



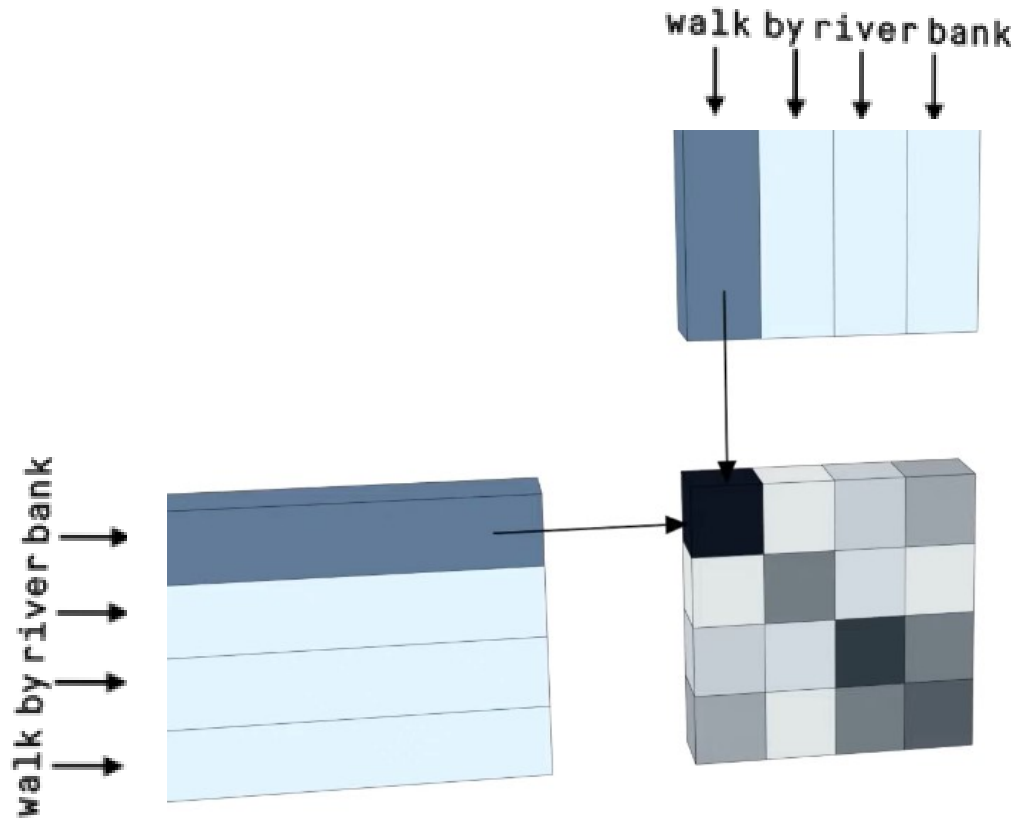
Contextualization by “Attention”



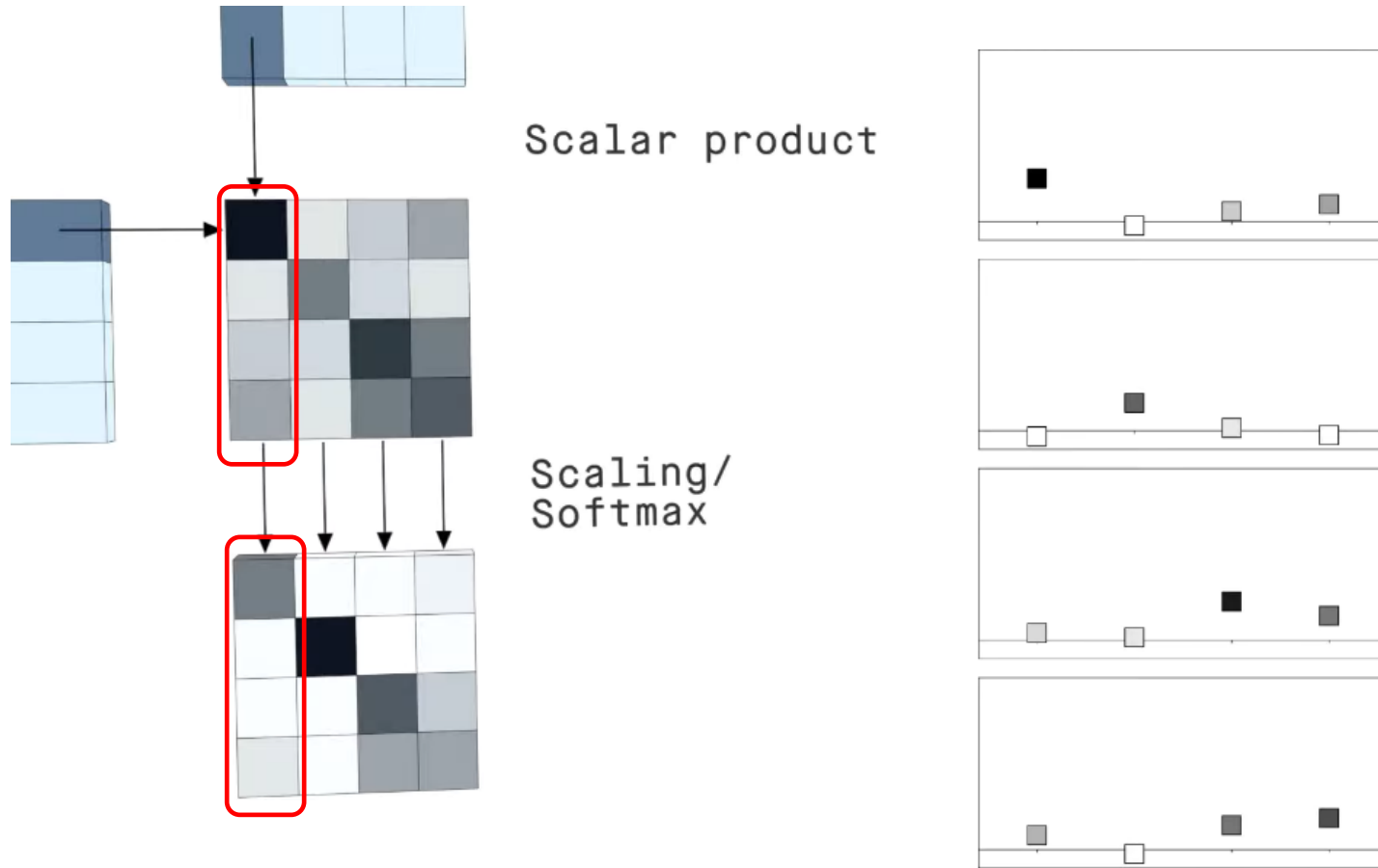
How does “attention” work?



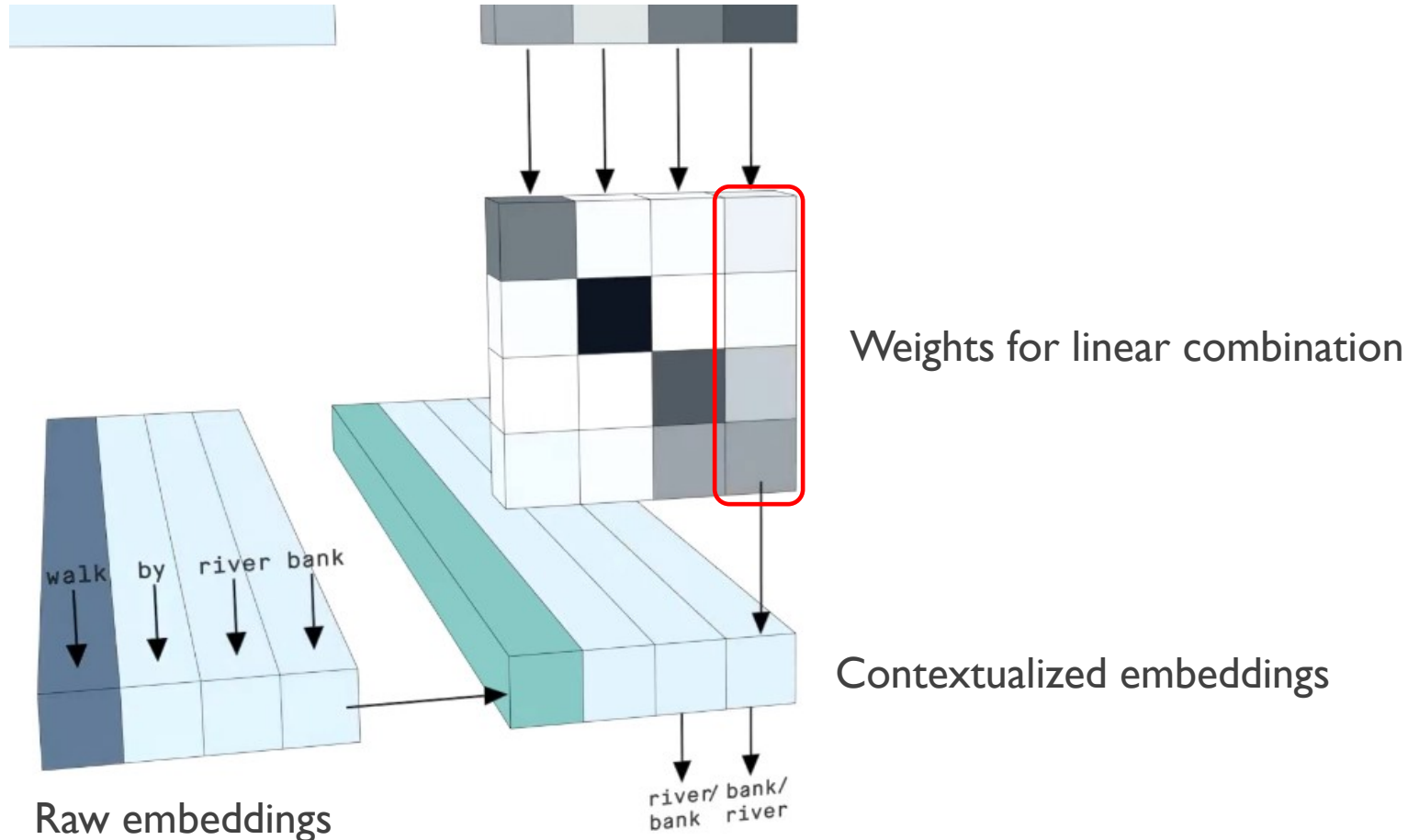
Similarity Measure via Inner Product



Exponential Saturation & Normalization via Softmax

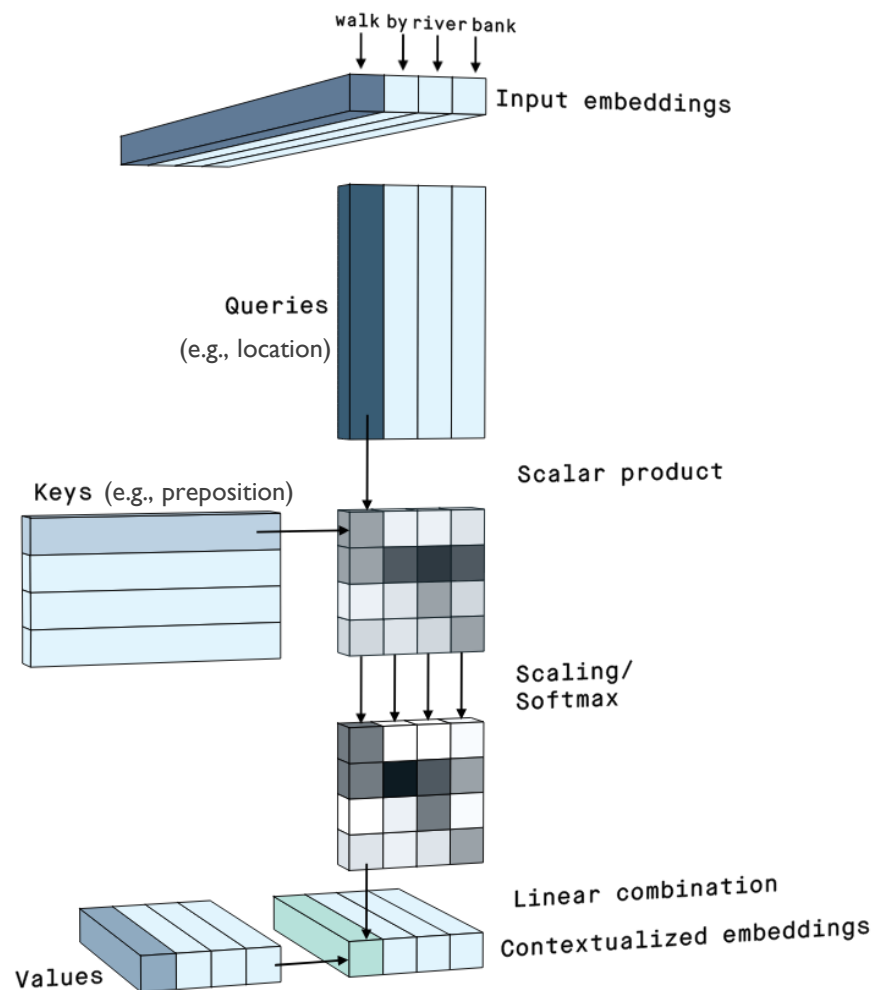


Contextualization via Linear Combination

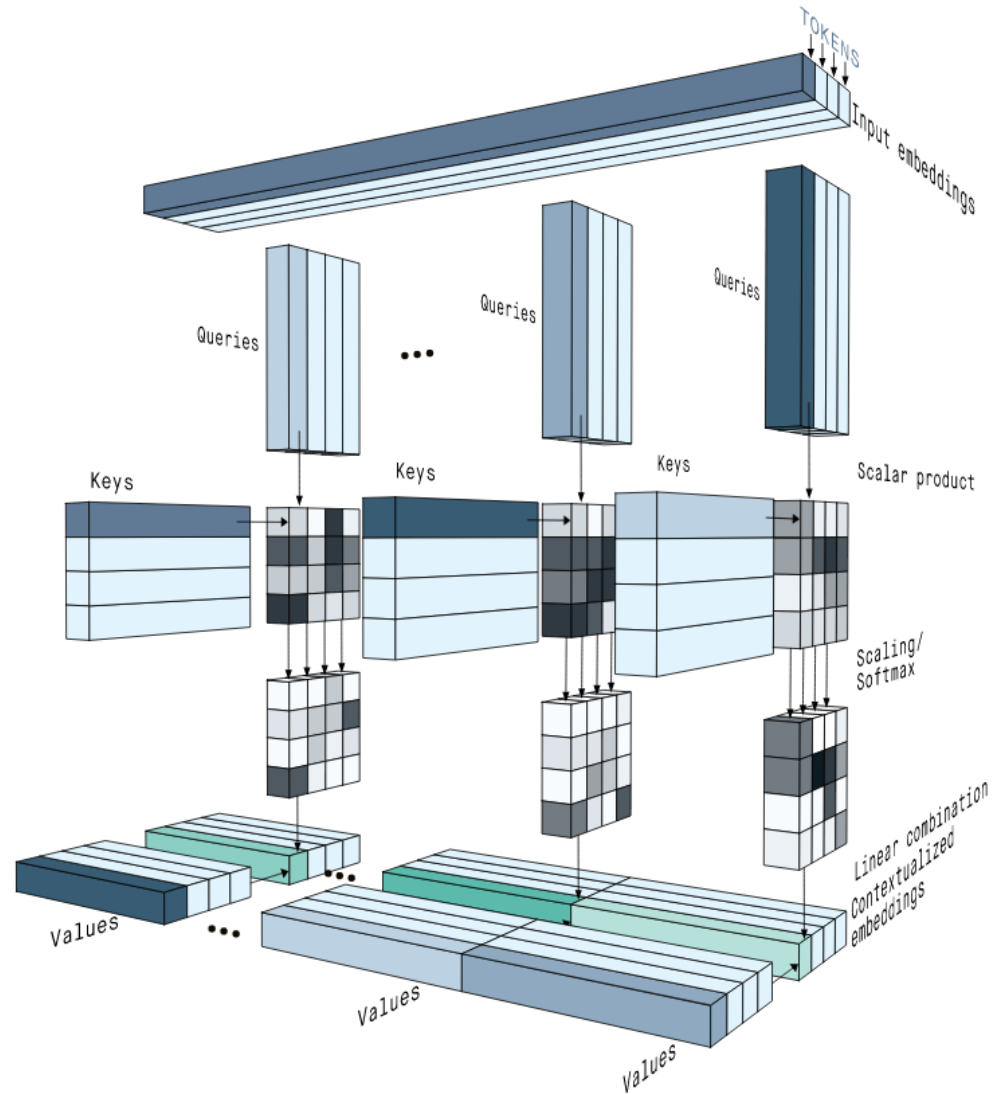


Key, Value, and Query

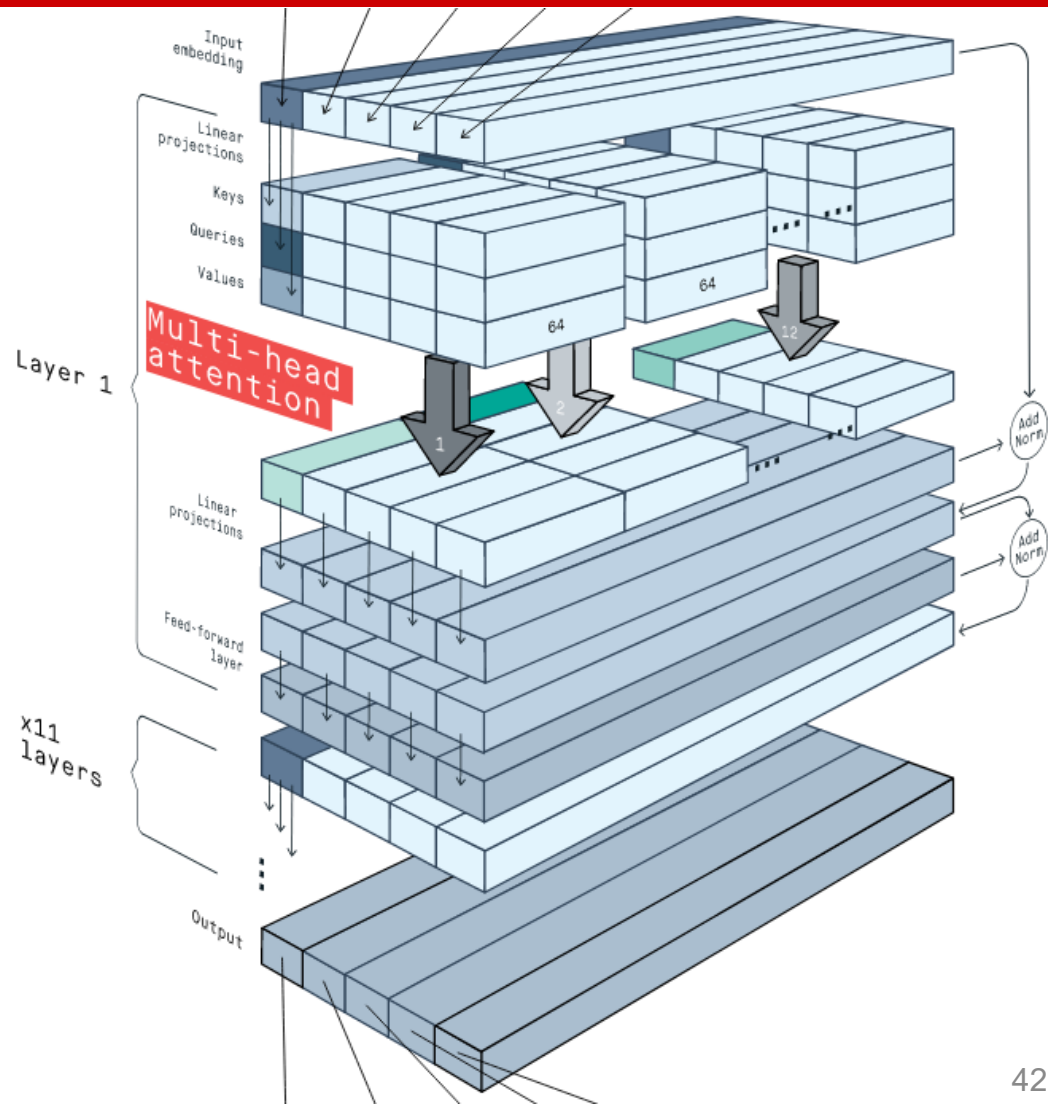
- ◆ “Key”, “value”, and “query” are three projections of an input embedding to three vector subspaces.
- ◆ Each subspace represents a unique semantic aspect.
- ◆ The projection operators / matrices provide **trainable** parameters for Transformer neural networks.



Multi-Head Attention



Bidirectional Encoder Representations from Transformers (BERT)



Neural Network Training: Backpropagation

Acknowledgment: Some graphics and slides were adapted from Profs. Jain (MSU), Min Wu (UMD), Fei-Fei Li (Stanford)

Some figures are from Duda-Hart-Stork textbook, Fei-Fei Li's slides

3-Layer Neural Network Structure

- ◆ A single “bias unit” is connected to each unit in addition to the input units

- ◆ Net activation:
$$net_j = \sum_{i=1}^d x_i w_{ji} + w_{j0} = \sum_{i=0}^d x_i w_{ji} \equiv \mathbf{w}_j^t \cdot \mathbf{x},$$

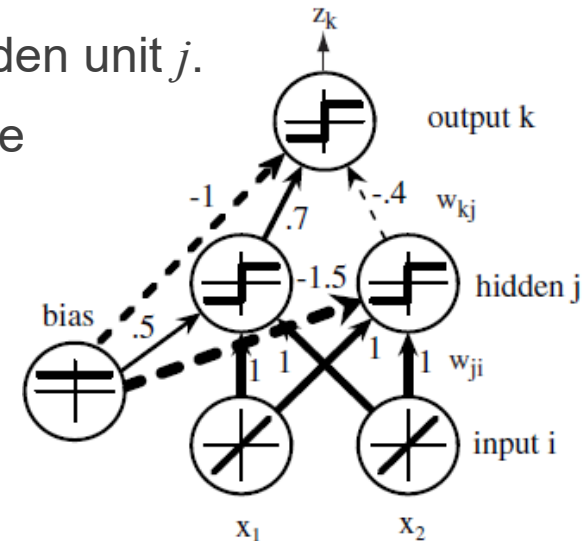
where the subscript i indexes units in the input layer, j indexes units in the hidden layer;

w_{ji} denotes the input-to-hidden layer weights at hidden unit j .

- ◆ In neurobiology, such weights or connections are called “synapses”

- ◆ Each hidden unit emits an output that is a nonlinear function of its activation

$$y_j = f(net_j)$$



Training Neural Networks / Estimating Weights

- Notations: $t_k \sim$ the k th target (or desired) output,
 $z_k \sim$ the k th estimated/computed output with $k = 1, \dots, c$.
 $w_{ij} \sim$ weight of the network

- Squared cost func:

$$J(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 = \frac{1}{2} \|\mathbf{t} - \mathbf{z}\|^2$$

- Learning based on gradient descent by iteratively updating the weights:

$$\mathbf{w}(m+1) = \mathbf{w}(m) + \Delta\mathbf{w}(m),$$

$$\Delta\mathbf{w} = -\eta \frac{\partial J}{\partial \mathbf{w}}$$

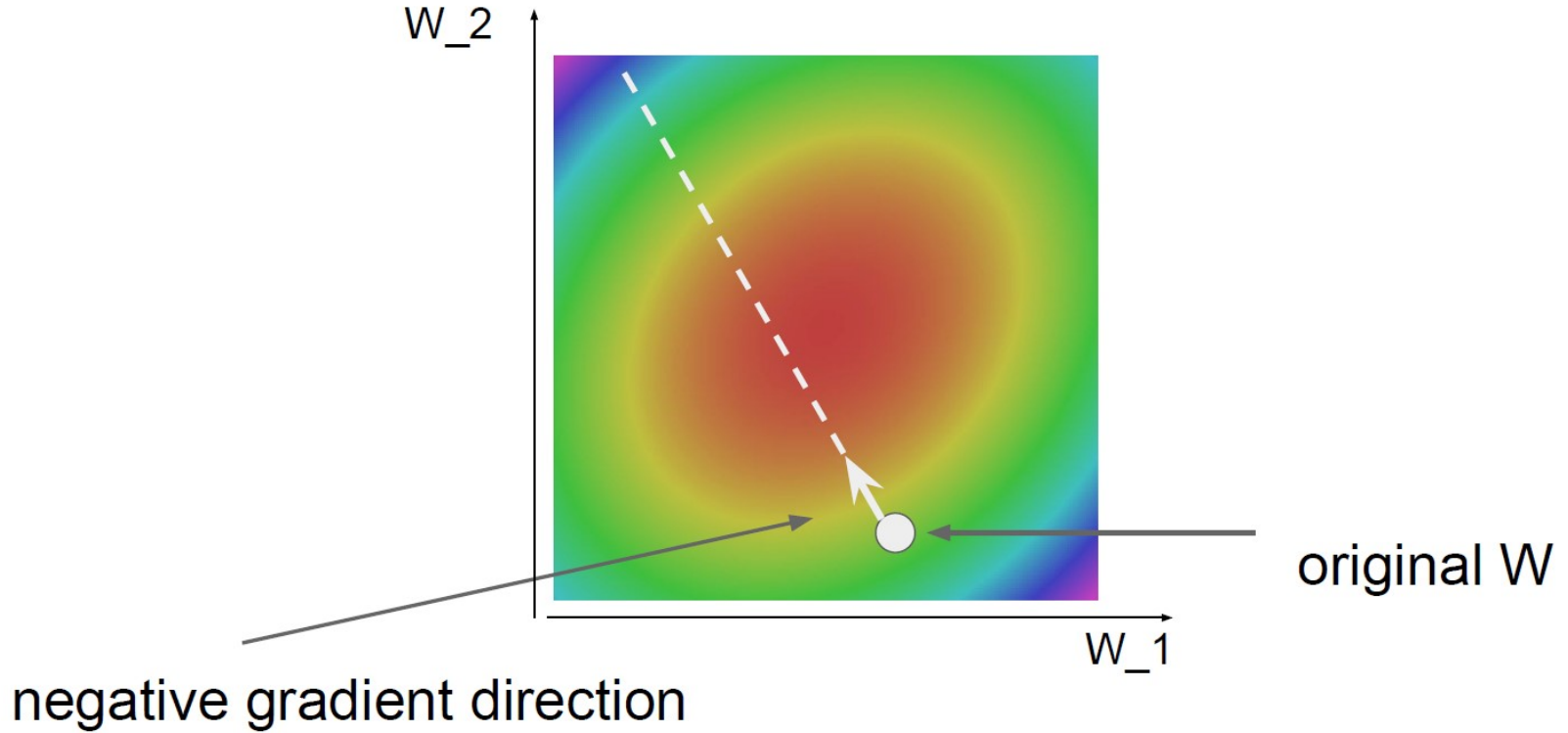
- The weights are initialized with random values, and updated in a direction to reduce the error.
- **Learning rate**, η , controls the step size of the update in weights.



Walking man image is CC0 1.0 public domain

Figure source: Stanford CS231n by Fei-Fei Li

Gradient Descent



Efficient Gradient Calculation: Backpropagation

- ◆ Computes $\partial J / \partial w_{ji}$ for a single input-output pair.
- ◆ Exploit the chain rule for differentiation, e.g.,

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y_j} \cdot \frac{\partial y_j}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}}$$

- ◆ Computed by **forward** and **backward** sweeps over the network, keeping track only of quantities local to each unit.
- ◆ Iterate backward one unit at a time from last layer. Backpropagation avoids redundant calculations.

Backpropagation (BP): An Example

Backpropagation: a simple example

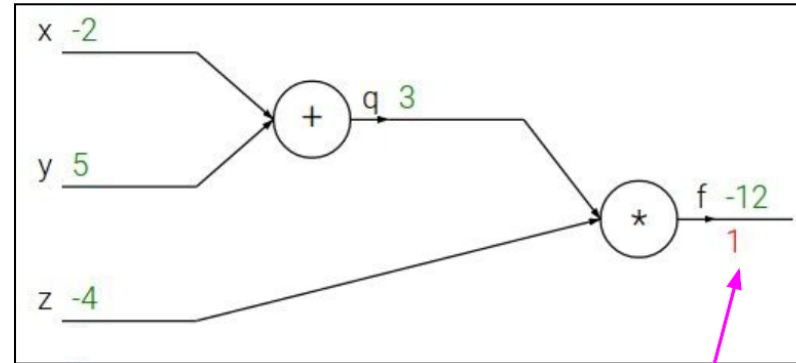
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial f}$$

Network Learning by BP (cont'd)

- ★ Error on hidden-to-output weight:

$$\frac{\partial J}{\partial w_{kj}} = \frac{\partial J}{\partial net_k} \cdot \frac{\partial net_k}{\partial w_{kj}} = -\delta_k \frac{\partial net_k}{\partial w_{kj}}$$

- ★ δ_k , the sensitivity of unit k :

describes how the overall error changes with the activation of the unit's net activation

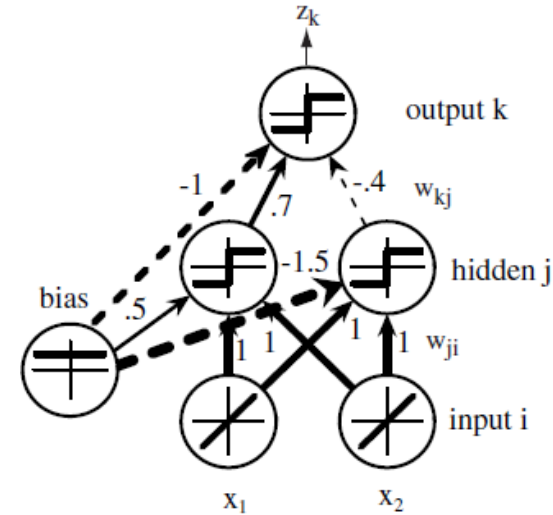
$$\delta_k \equiv -\frac{\partial J}{\partial net_k}$$

$$\delta_k = -\frac{\partial J}{\partial net_k} = -\frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial net_k} = (t_k - z_k) f'(net_k)$$

- ★ Since $net_k = w_k^t y$, we have $\frac{\partial net_k}{\partial w_{kj}} = y_j$

- ★ **Summary I:** weight update (or learning rule) for the hidden-to-output weight is:

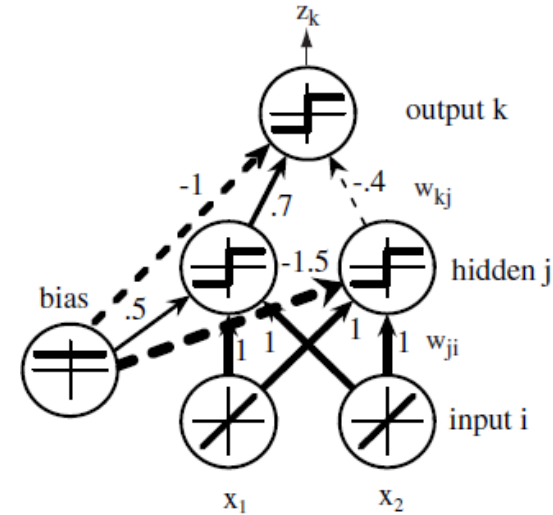
$$\Delta w_{kj} = \eta (t_k - z_k) f'(net_k) y_j = \eta \delta_k y_j$$



Network Learning by BP (cont'd)

★ Error on input-to-hidden weight:

- chain rule:
$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial J}{\partial y_j} \cdot \frac{\partial y_j}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}}$$
- $$\frac{\partial J}{\partial y_j} = \frac{\partial}{\partial y_j} \left[\frac{1}{2} \sum_{k=1}^c (t_k - z_k)^2 \right] = - \sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial y_j}$$
- $$= - \sum_{k=1}^c (t_k - z_k) \frac{\partial z_k}{\partial net_k} \cdot \frac{\partial net_k}{\partial y_j} = - \sum_{k=1}^c (t_k - z_k) f'(net_k) w_{kj}$$



★ Sensitivity of a hidden unit:
(Similarly defined as earlier)

$$\delta_j \equiv \frac{\partial J}{\partial net_j} = f'(net_j) \sum_{k=1}^c w_{kj} \delta_k$$

★ **Summary 2:** Learning rule for the input-to-hidden weight is:

$$\Delta w_{ji} = \eta \underbrace{[f'(net_j) \sum w_{kj} \delta_k]}_{\delta_j} x_i = \eta \delta_j x_i$$

Sensitivity at Hidden Node

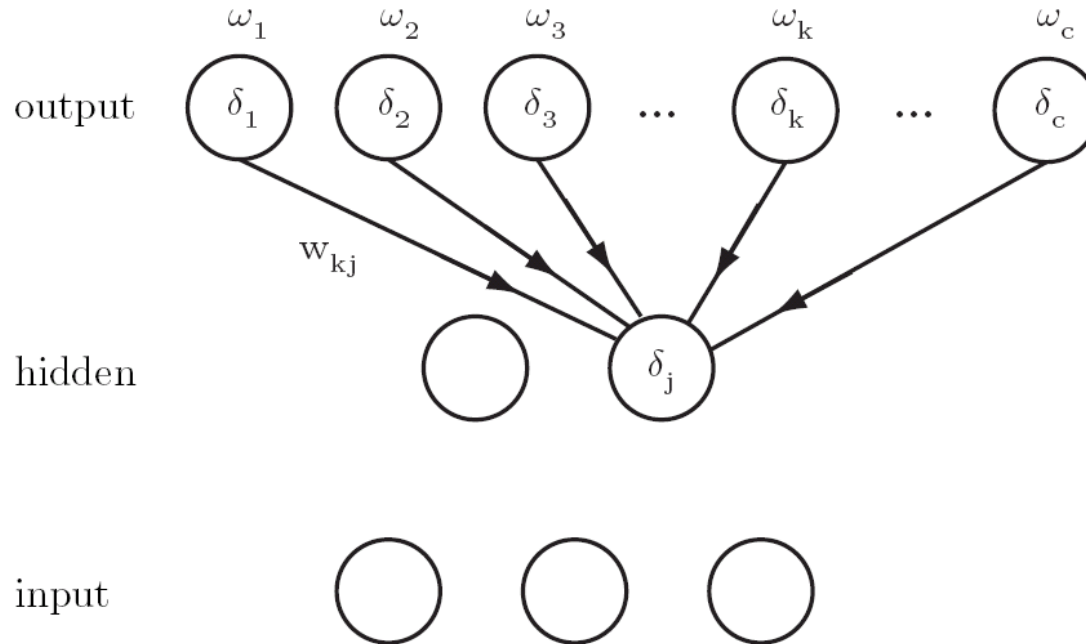


Figure 6.5: The sensitivity at a hidden unit is proportional to the weighted sum of the sensitivities at the output units: $\delta_j = f'(net_j) \sum_{k=1}^c w_{kj} \delta_k$. The output unit sensitivities are thus propagated “back” to the hidden units.

BP Algorithm: Training Protocols

◆ Training protocols:

- ★ **Batch:** Present all patterns before updating weights
- ★ **Stochastic:** patterns/input are chosen randomly from training set; network weights are updated for each pattern
- ★ **Online:** present each pattern once & only once (no memory for storing patterns)

◆ Stochastic backpropagation algorithm:

```

Begin  initialize  $n_H$ ;  $w$ , criterion thres,  $\eta$ ,  $m \leftarrow 0$ 
  do  $m \leftarrow m + 1$ 
     $x^m \leftarrow$  randomly chosen pattern
     $w_{ji} \leftarrow w_{ji} + \eta \delta_j x_i$ ;  $w_{kj} \leftarrow w_{kj} + \eta \delta_k y_j$ 
  until  $||\nabla J(w)|| <$  thres
  return  $w$ 

End

```

BP Algorithm: Stopping Criterion

- ◆ Algorithm terminates when the change in criterion function $J(w)$ is smaller than some preset thres
 - ✦ Also exist other stopping criteria with better performance
- ◆ A weight update may reduce the error on the single pattern being presented, but can increase the error on the full training set
 - ✦ In stochastic backpropagation and batch propagation, we should make several passes (epoches) through the training data

Learning Curves

- ★ Before training starts, the error on the training set is high; as the learning proceeds, error becomes smaller
- ★ Error per pattern depends on the amount of training data and expressive power (e.g. # of weights) in the network
- ★ Average error on an independent test set is always higher than on the training set, and it can decrease or increase
- ★ A validation set is used in order to decide when to stop training:
 - Avoid overfitting the network and decrease the power of the classifier's generalization

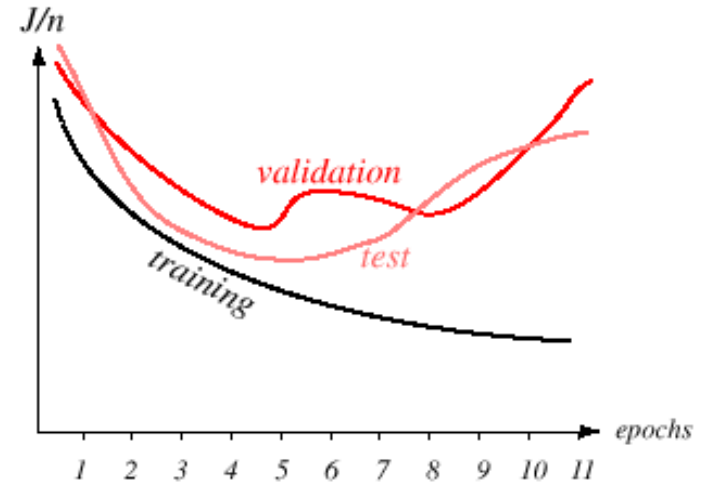


FIGURE 6.6. A learning curve shows the criterion function as a function of the amount of training, typically indicated by the number of epochs or presentations of the full training set. We plot the average error per pattern, that is, $1/n \sum_{p=1}^n J_p$. The validation error and the test or generalization error per pattern are virtually always higher than the training error. In some protocols, training is stopped at the first minimum of the validation set. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

“Stop training when the error on the validation set is minimum”

Practical Considerations: Learning Rate

◆ Learning Rate

- ★ Small learning rate: slow convergence
- ★ Large learning rate: high oscillation and slow convergence

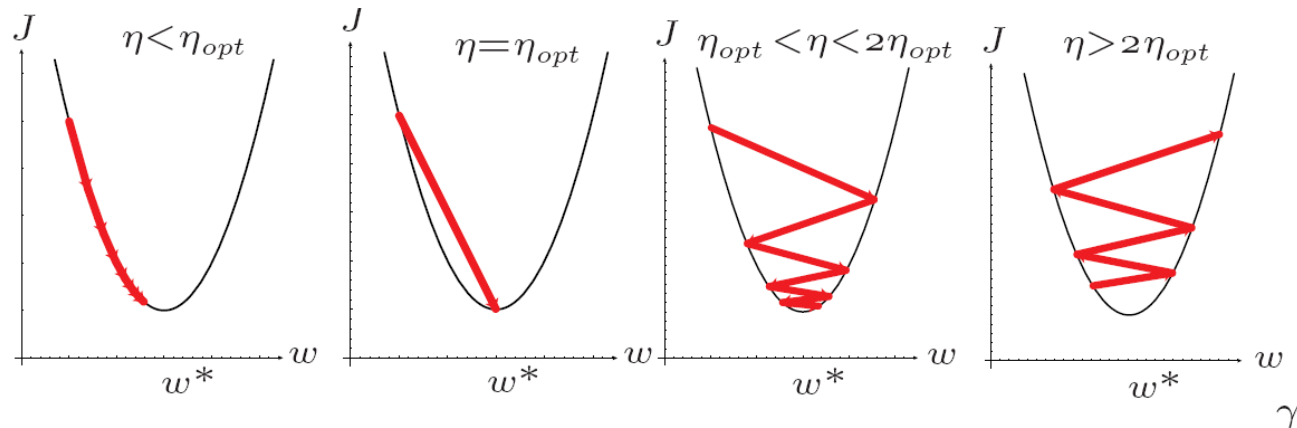


Figure 6.18: Gradient descent in a one-dimensional quadratic criterion with different learning rates. If $\eta < \eta_{opt}$, convergence is assured, but training can be needlessly slow. If $\eta = \eta_{opt}$, a single learning step suffices to find the error minimum. If $\eta_{opt} < \eta < 2\eta_{opt}$, the system will oscillate but nevertheless converge, but training is needlessly slow. If $\eta > 2\eta_{opt}$, the system diverges.