

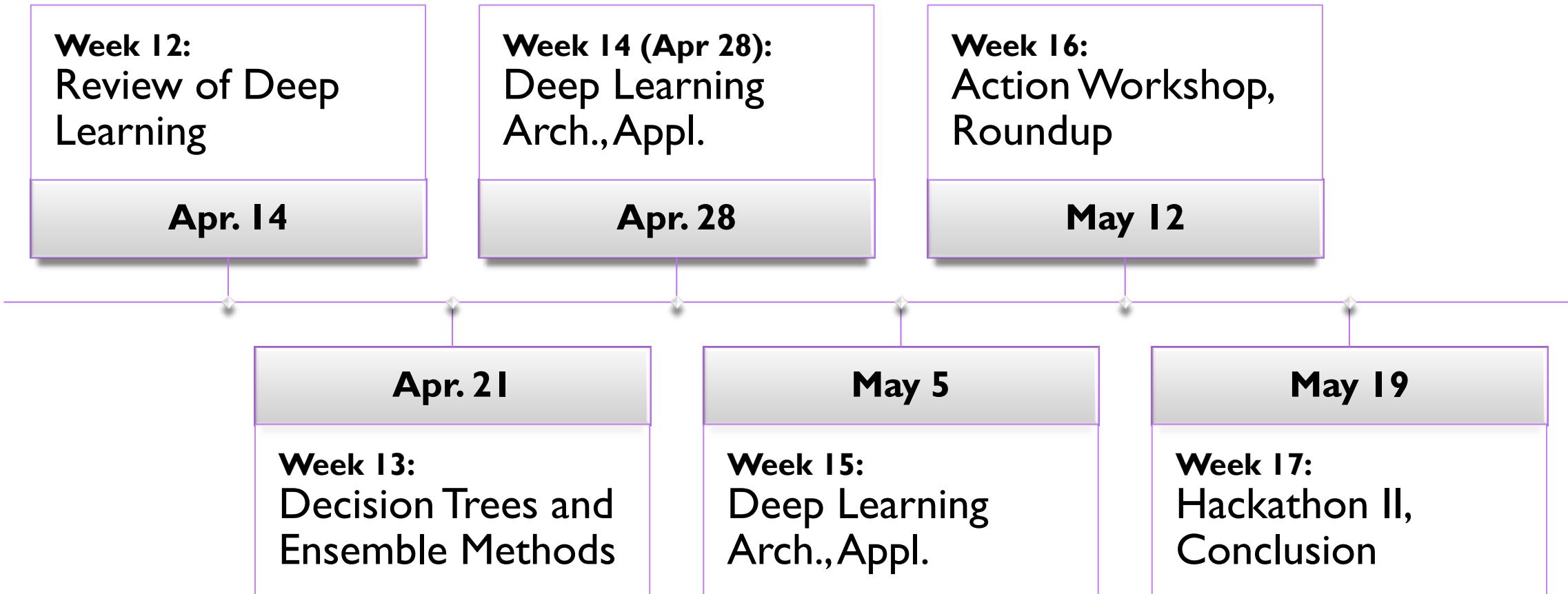
Administrivia: Road Ahead

Lectures, Mini Hackathons, Hackathon



Weeks Ahead

- No Extension of the course
- Resequence between Roundup and Lecture





This Week



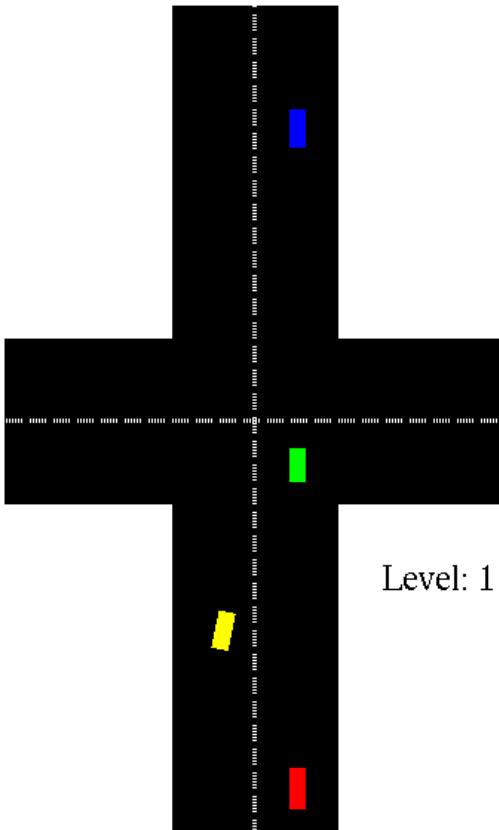
- Lecture
 - Review of Deep Neural Networks
- Lab
 - Catch-up with Experiments 9/10/11
 - Mini Hackathon: Drive a car



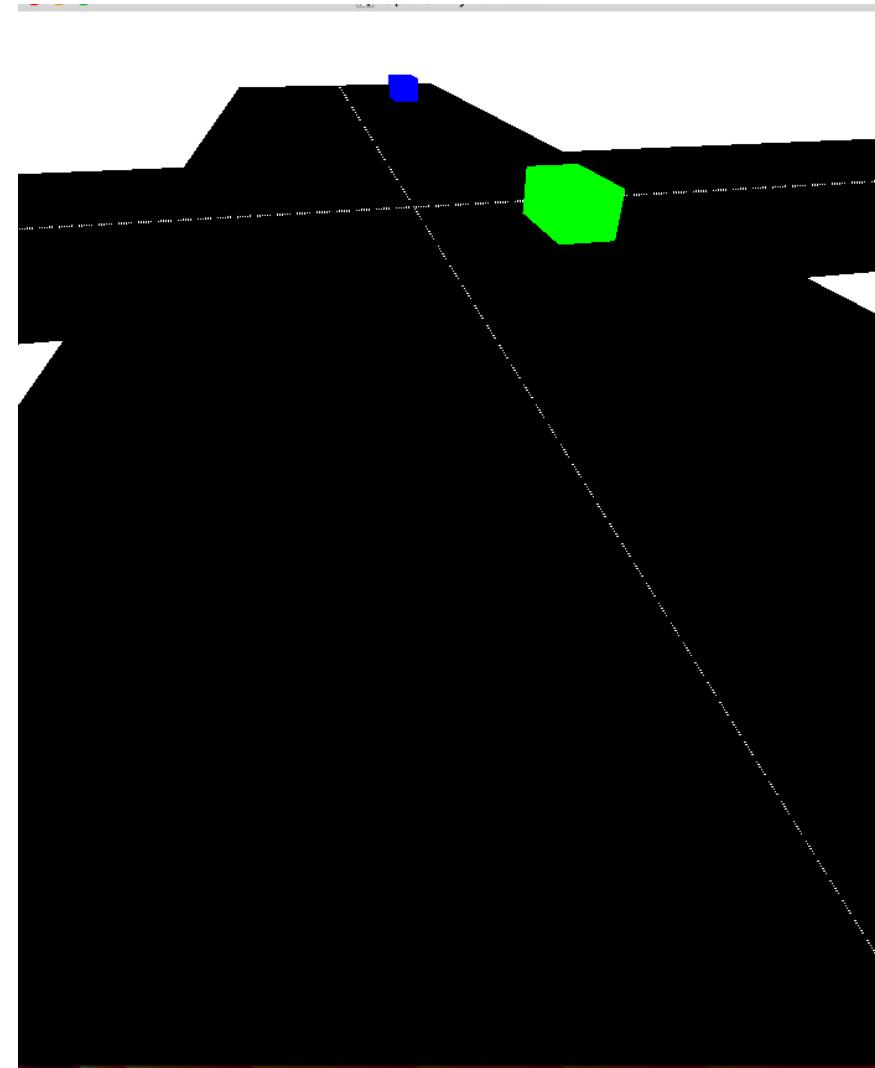
The Task



TAKE RIGHT TO WIN



Score: 0.657

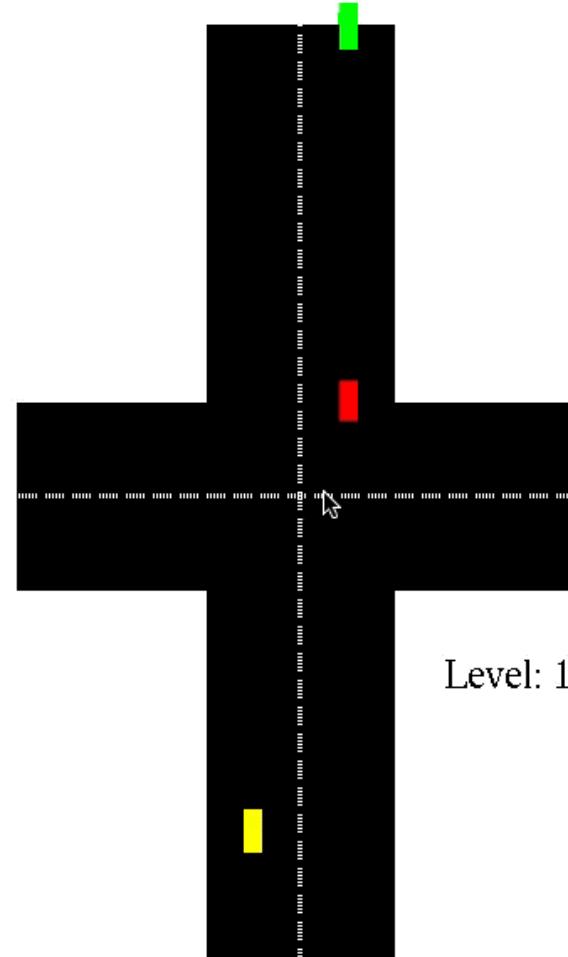




TAKE RIGHT TO WIN



The Task





Lab Plan



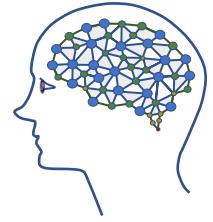
- First 30 minutes: Generate Training data and upload
- Work on experiments for 9/10/11
- Download combined training data
- Use it to train an MLP
 - How do you deal with noise in data?



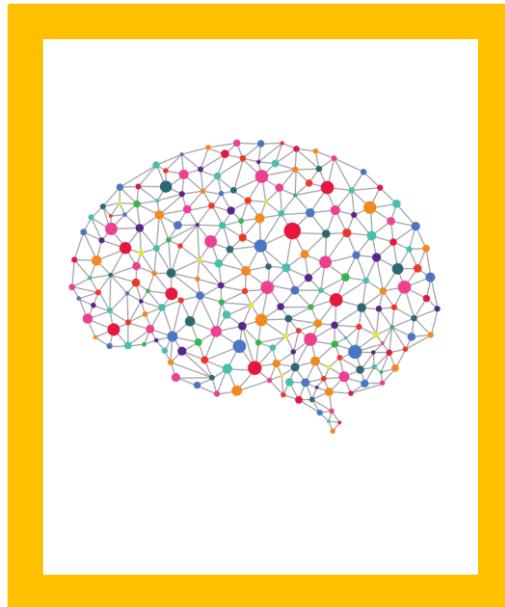
Mini Hackathons



- Mechanism:
 - Add a group exercise to the regular lab sessions
 - We will start in the labs, you can complete it later in the week
- Goal:
 - Enable you to do larger tasks
 - Enable peer learning
 - Form the groups early to facilitate final Hackathon
 - Build some background for the final Hackathon

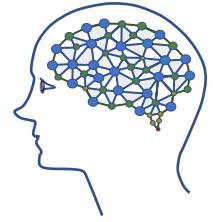


Questions?



Deep Neural Networks: A Review

CNNs, RNNs and Auto Encoders



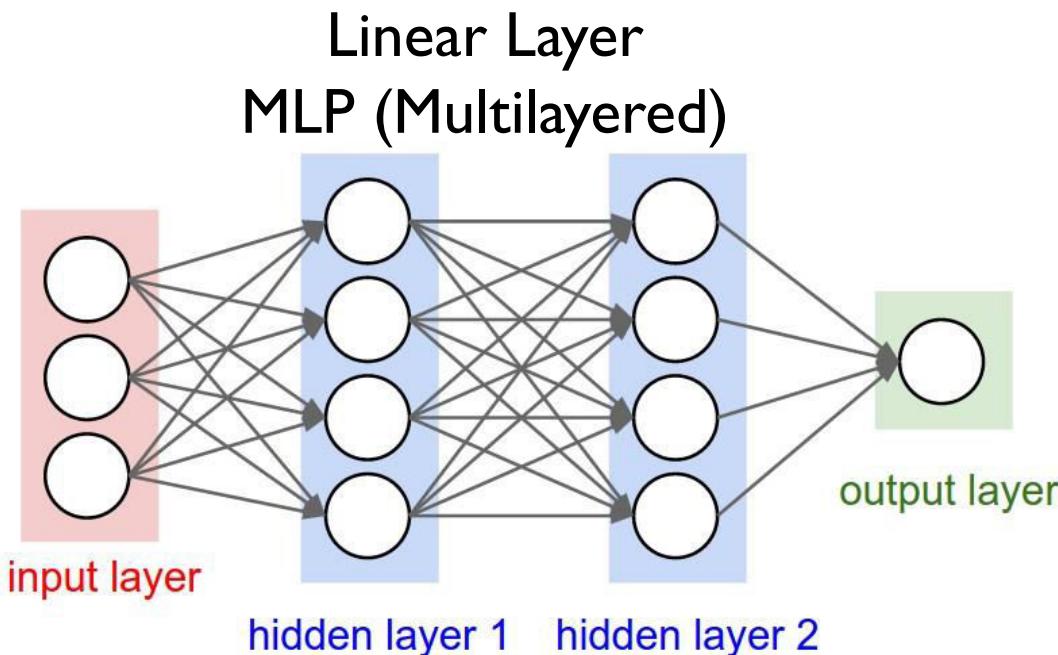
Convolutional Neural Networks

Inference and Training

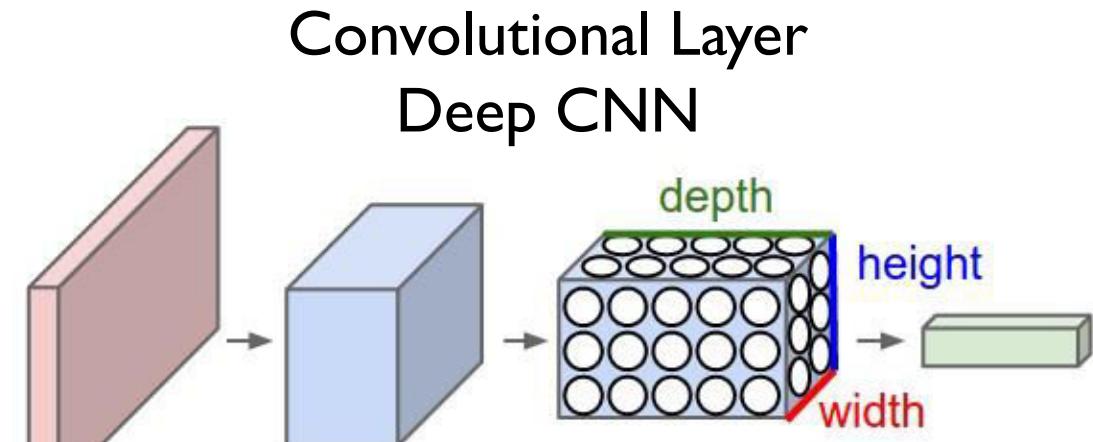


Convolutional Neural Networks

When Input is a vector
(fixed length).



When Input is a 1D/2D/3D sequence.
(variable length)

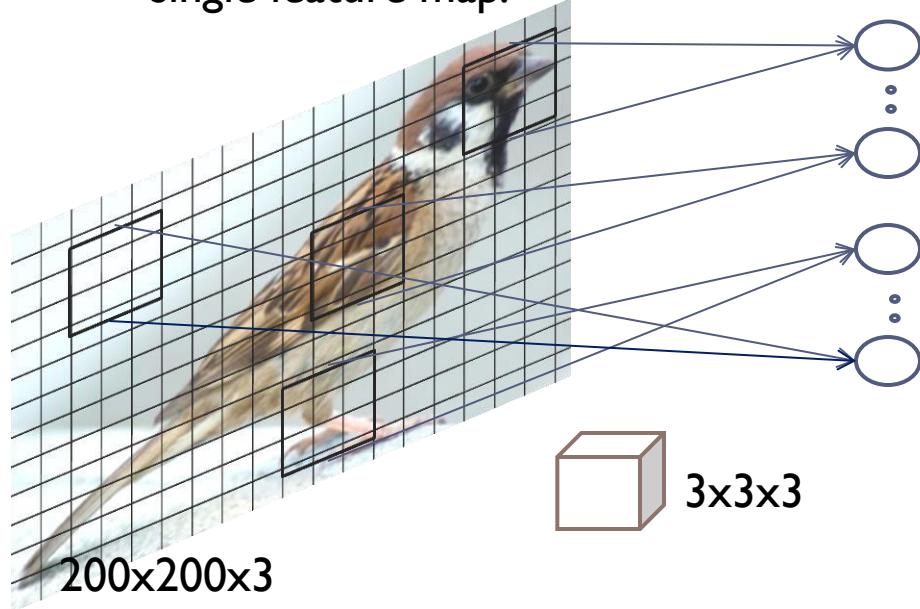




Convolutional Neural Network

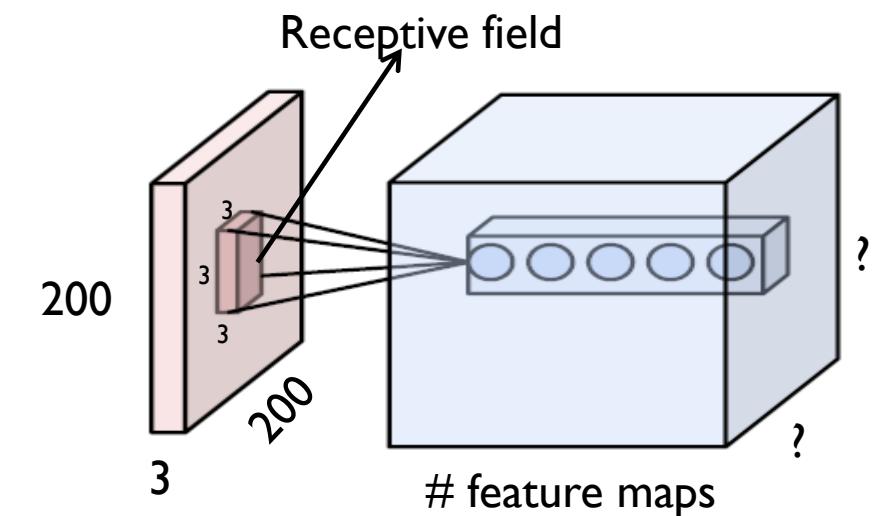
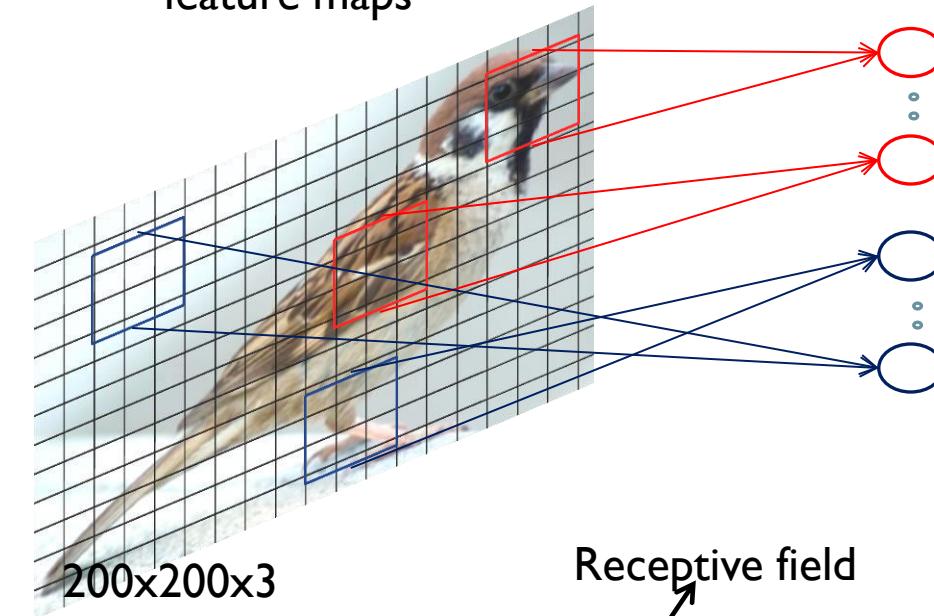


Convolutional layer with single feature map.



- #Hidden Units: 120,000
- #Params: $27 \times \# \text{Feature Maps}$
- Sharing parameters
- Exploiting the stationarity property and preserves locality of pixel dependencies

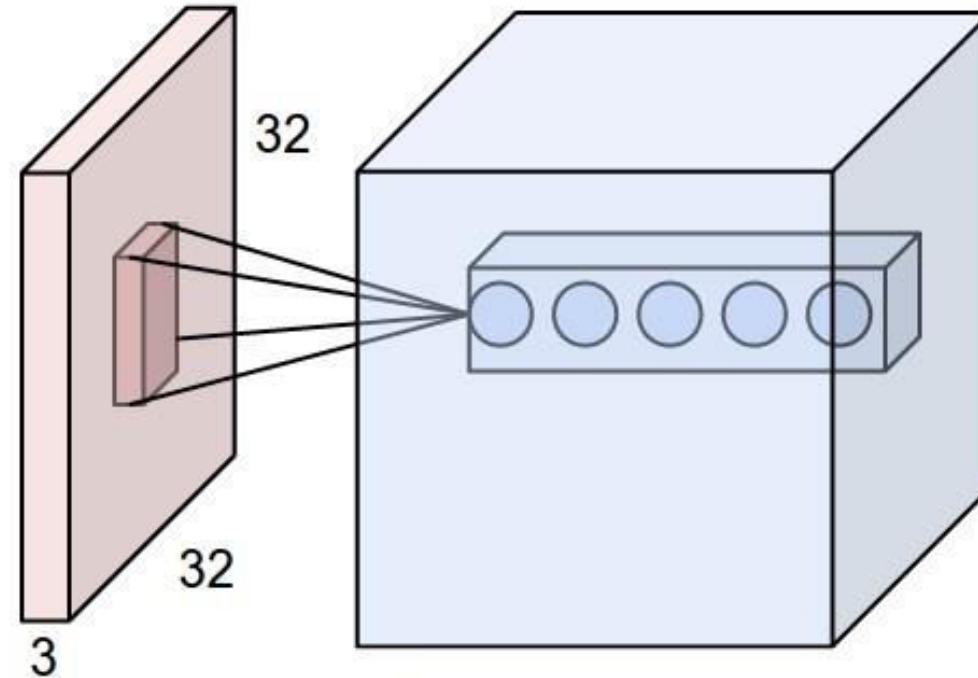
Convolutional layer with multiple feature maps





Input, Output Channels : Multiple Filters

A Closer Look at Convolution:

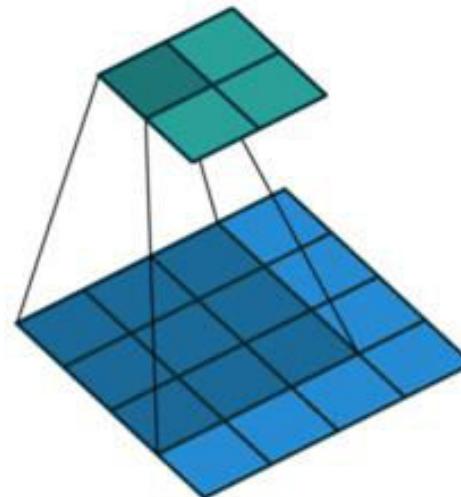




Convolution Layer



Window size



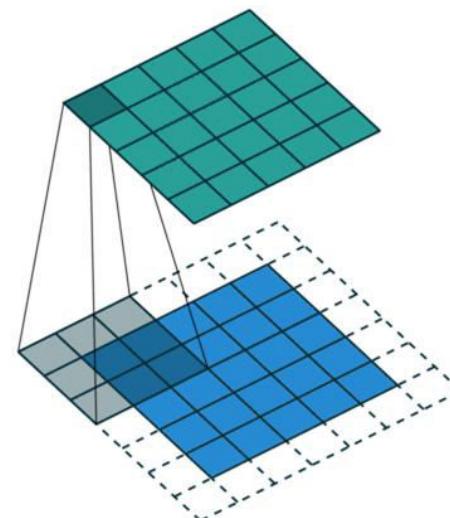
Window size: 3x3

Stride: 1

Padding: 0

Stride

Padding



Window size: 3x3

Stride: 1

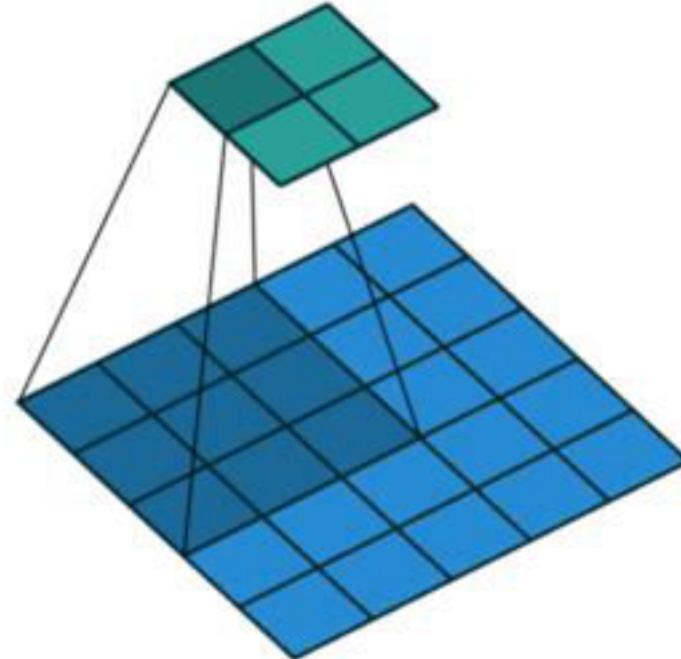
Padding: 1



Convolution Layer

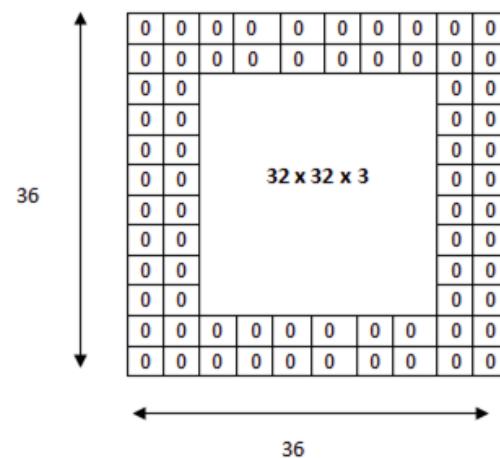


Strides reduces dimension



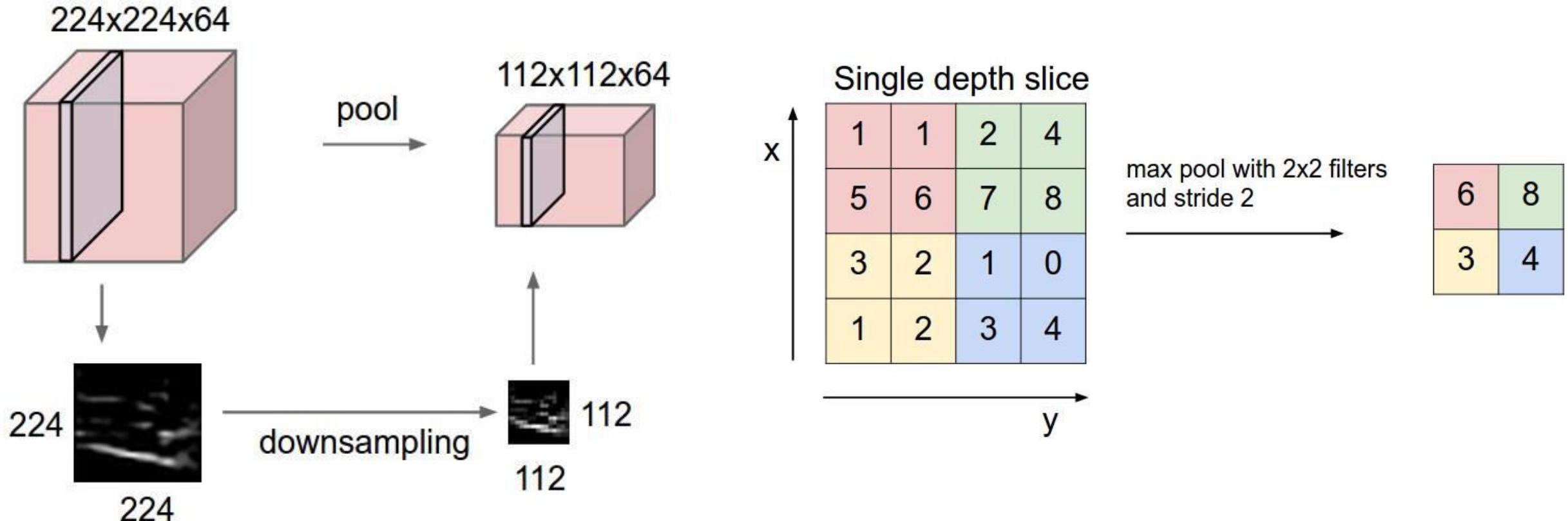
Window size: 3x3
Stride: 2
Padding: 0

$$O = \frac{(W - K + 2P)}{S} + 1$$





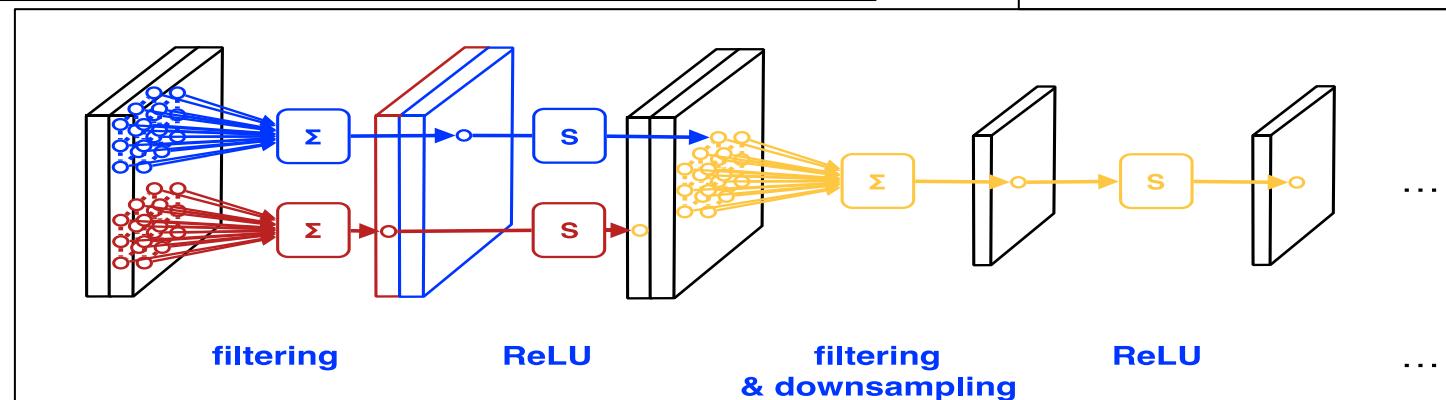
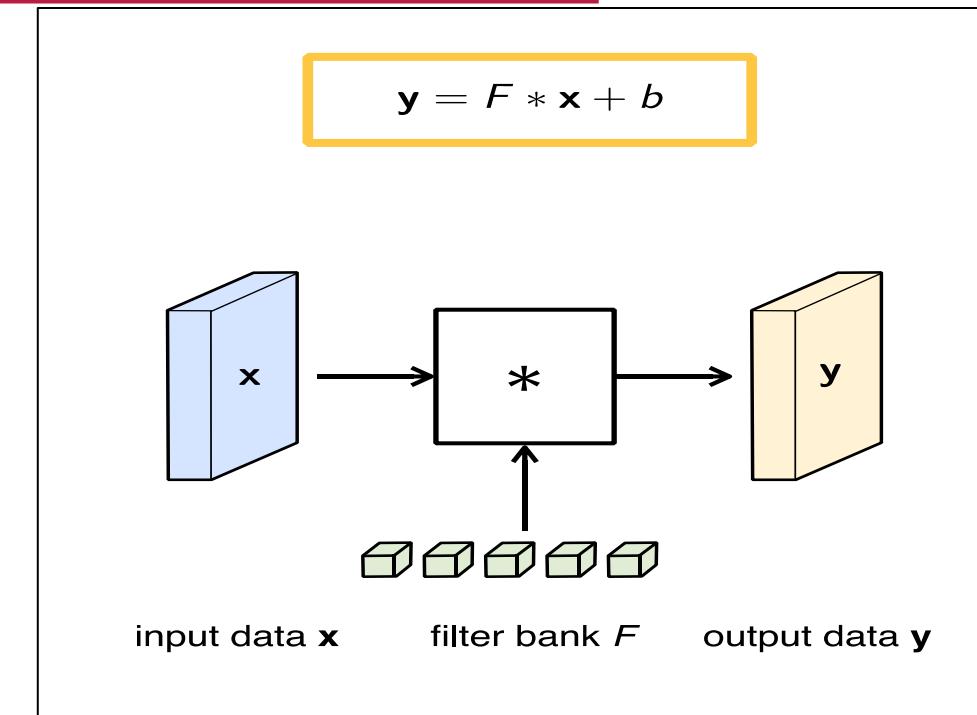
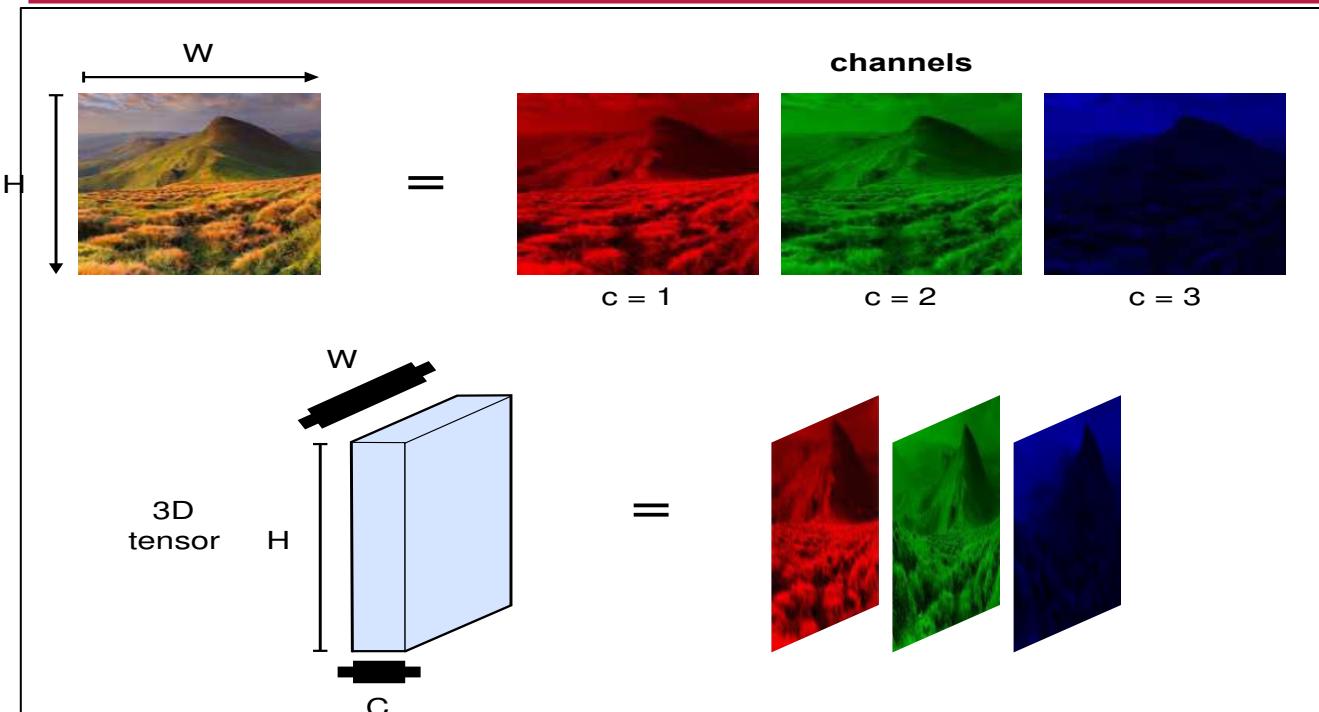
Pooling



No Learnable Parameters



Summary of CNNs





Learn the full pipeline

VISION

pixels → edge → texton → motif → part → object

SPEECH

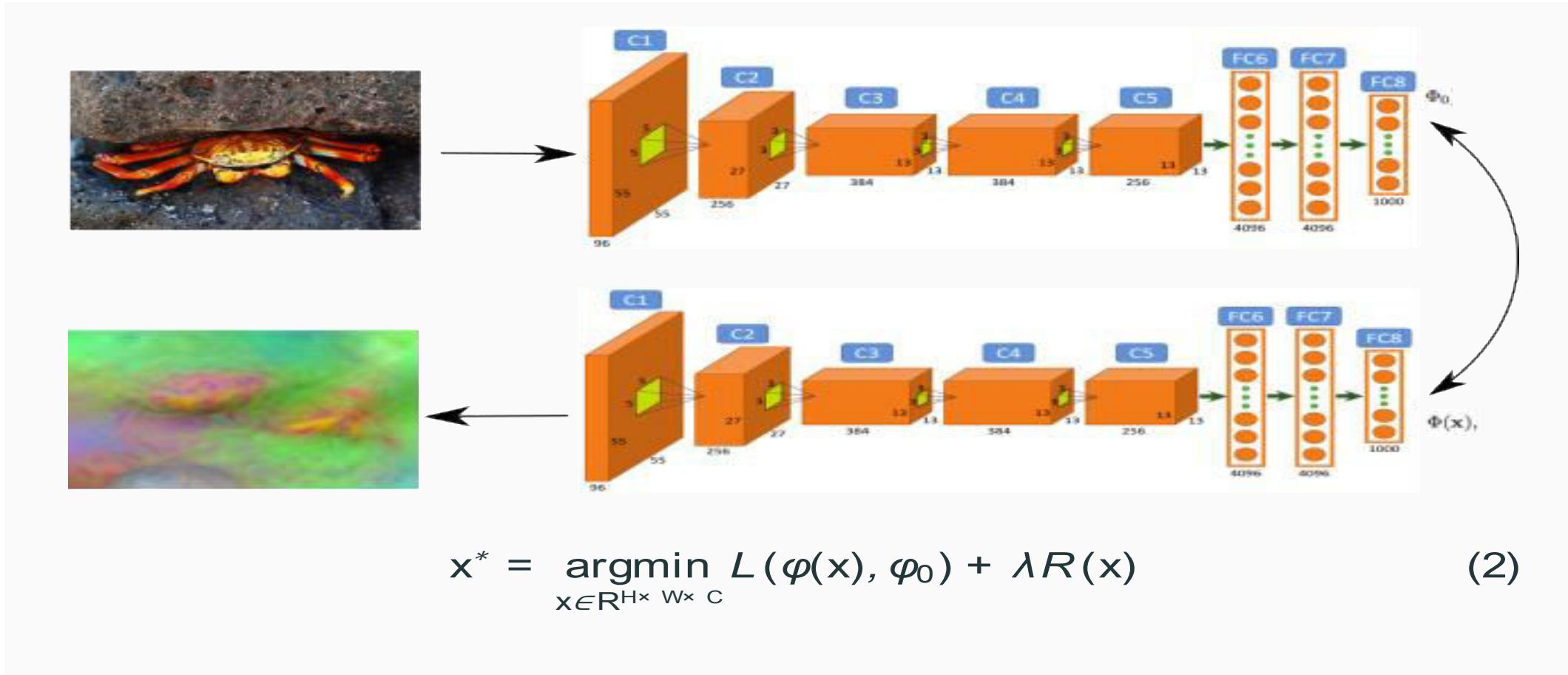
sample → spectral band → formant → motif → phone → word

NLP

character → word → NP/VP/.. → clause → sentence → story

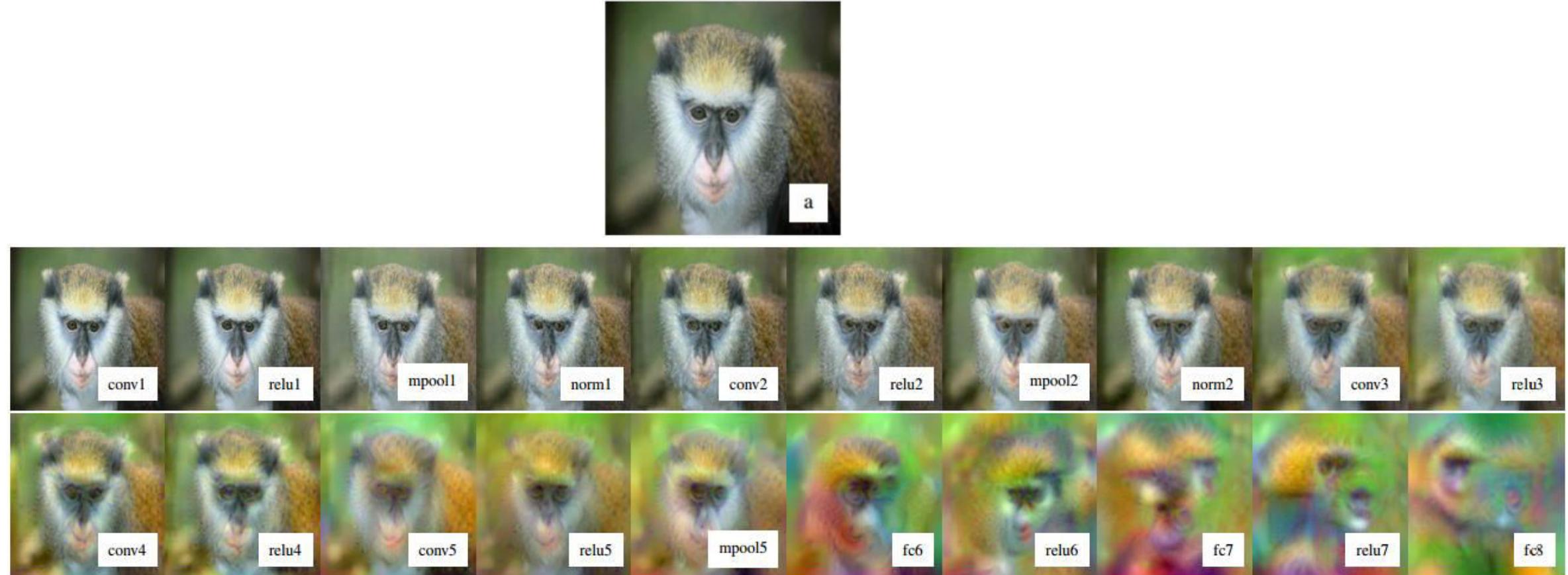


Inverting Specific Representation

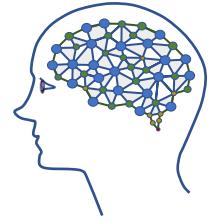
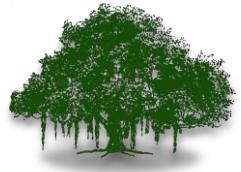




Inverting at Different Stages



Reconstructions from intermediate layers



Questions?



Auto Encoder

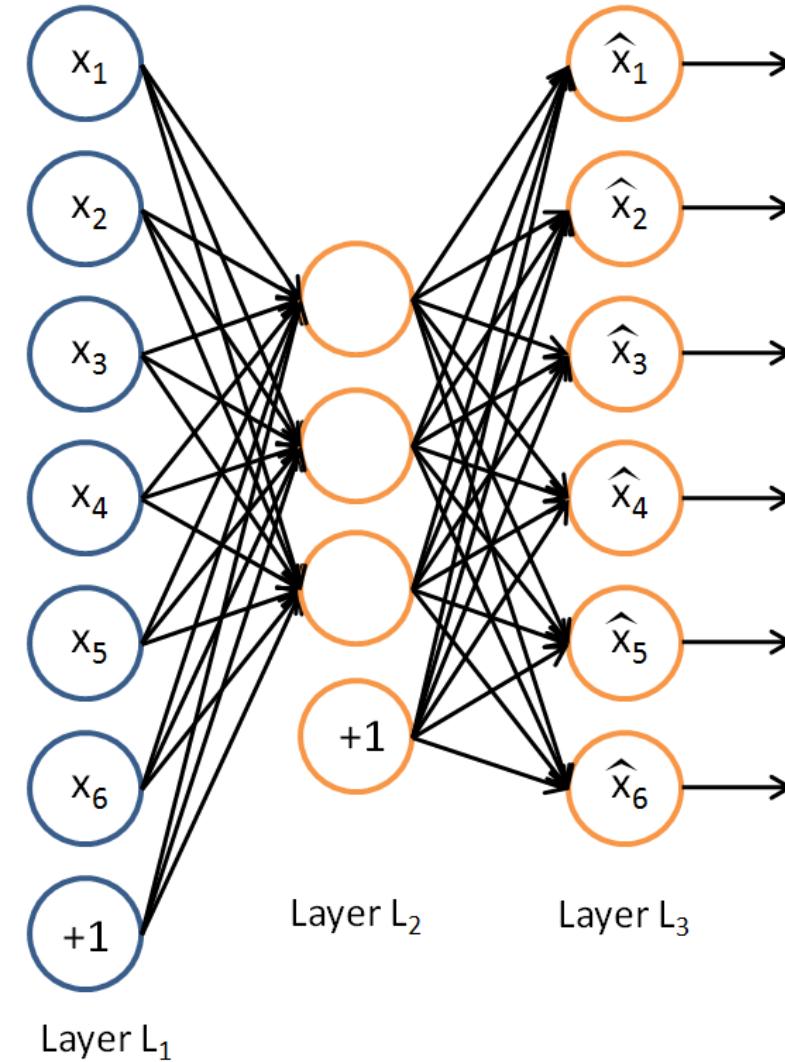
What if we do not have labels?



Autoencoder



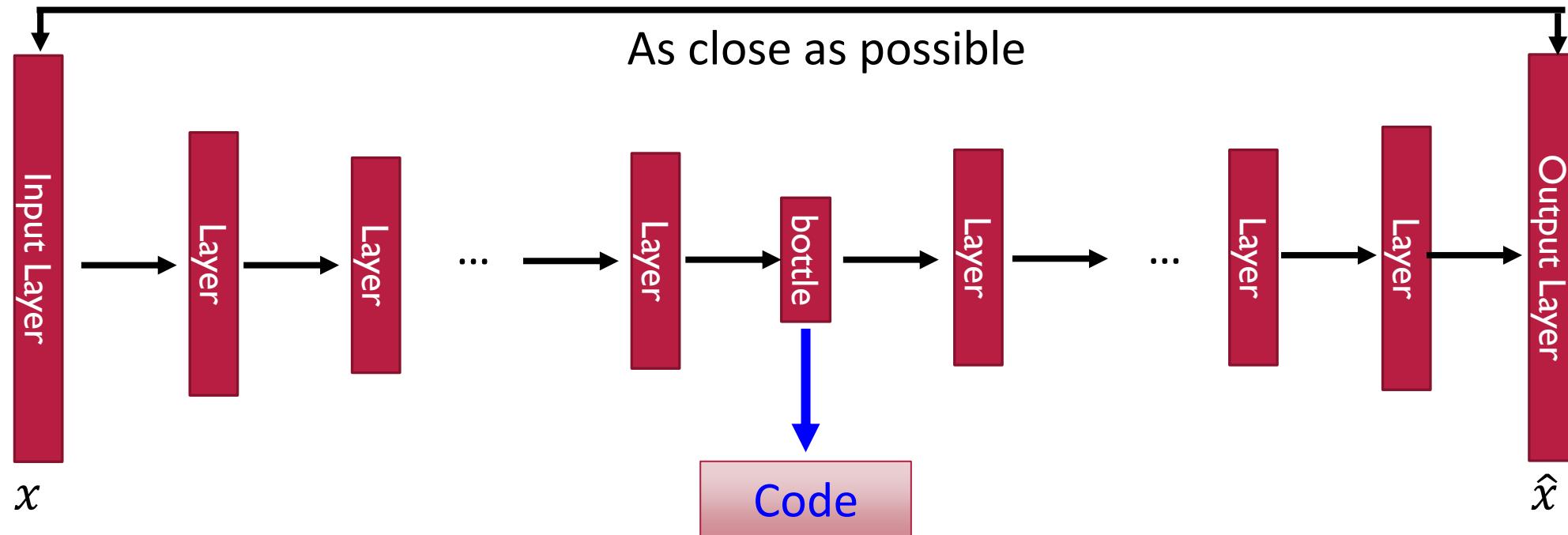
- Similar to MLP
- Input is same as output
- Network learns to reconstruct.
- “Bottleneck” layer learns a compact representation.





Deep Auto-encoder

- Of course, the auto-encoder can be deep Symmetric is not necessary.



Reference: Hinton, Geoffrey E., and Ruslan R. Salakhutdinov. "Reducing the dimensionality of data with neural networks." *Science* 313.5786 (2006): 504-507



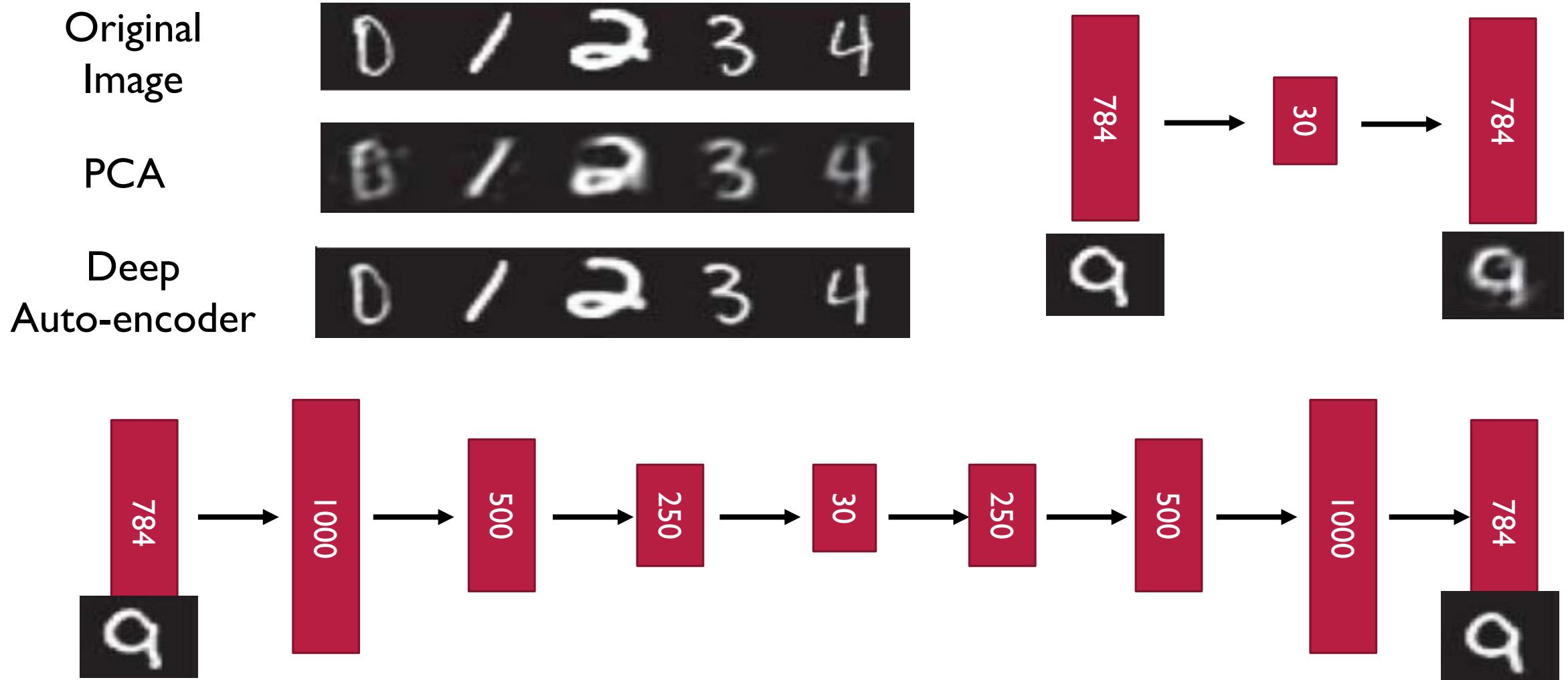
View Points

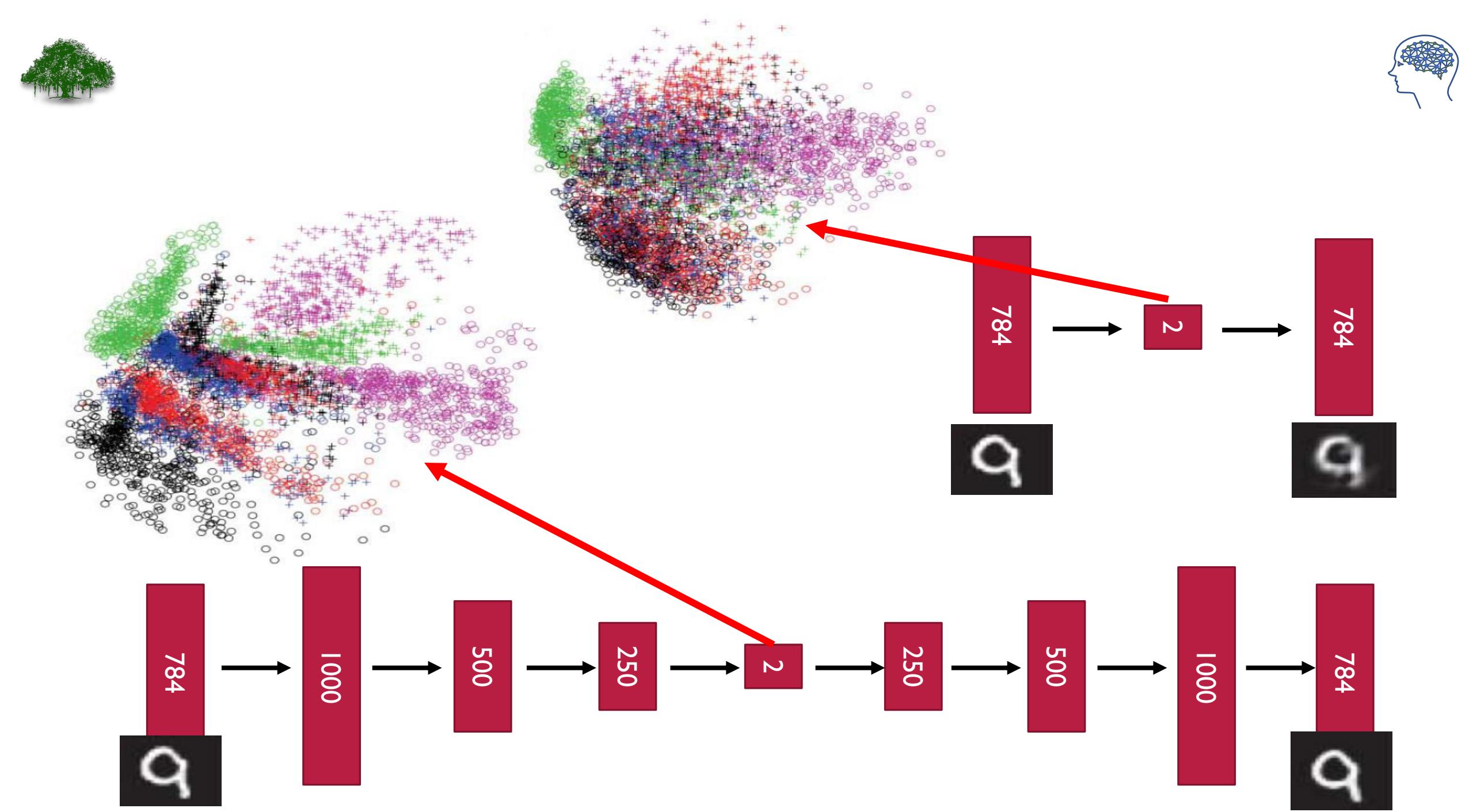


- Nonlinear PCA (Dimensionality Reduction)
- Unsupervised Learning
- Data compression



Deep Auto-encoder

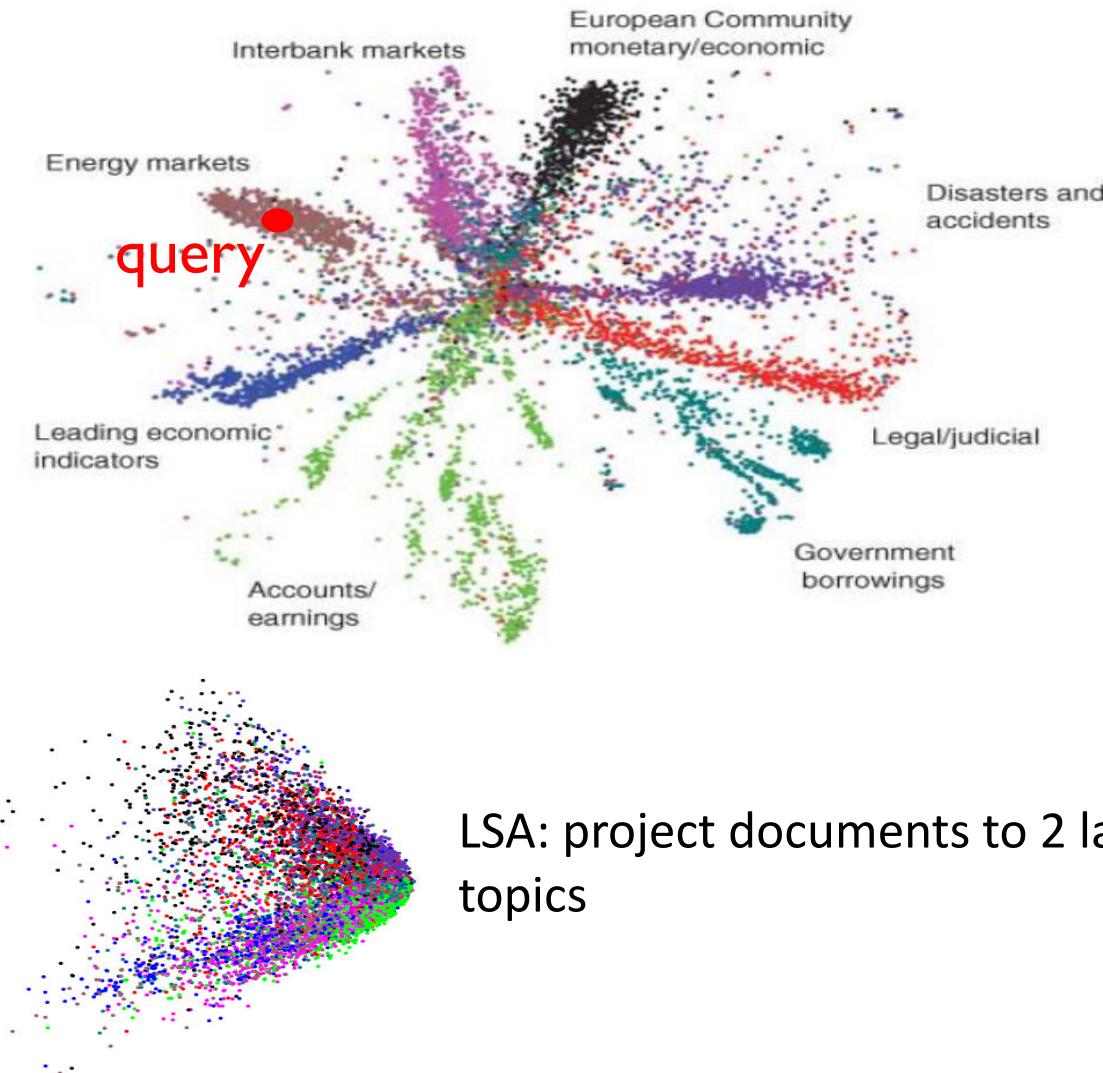
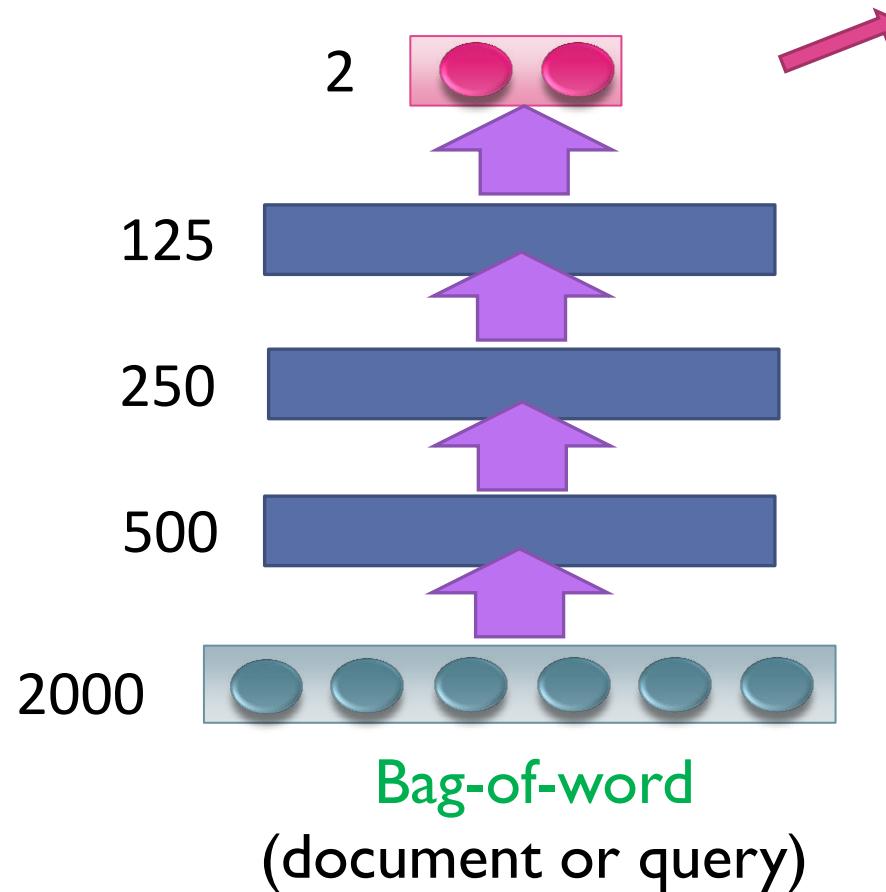






Auto-encoder – Text Retrieval

The documents talking about the same thing will have close code.





Variants

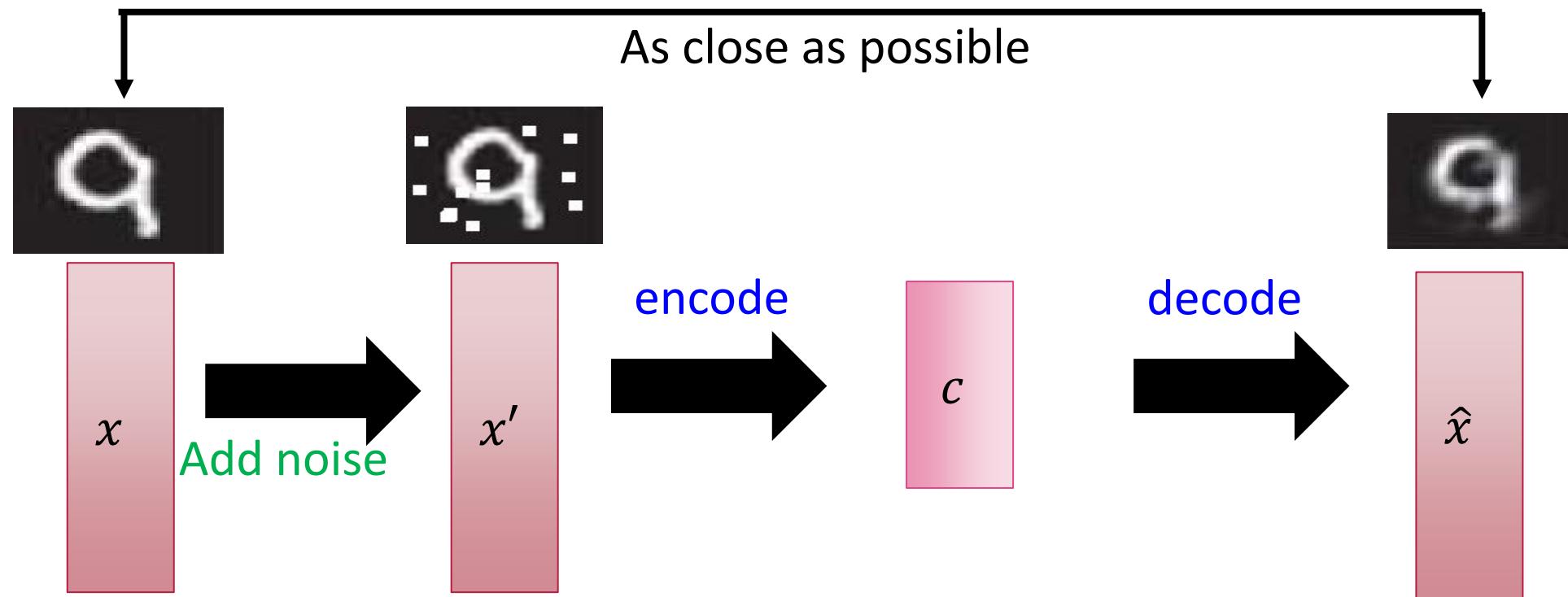


- Denoising Auto encoder
- Sparse Auto encoder

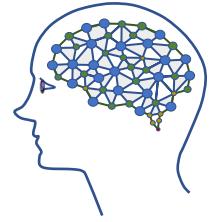


Auto-encoder

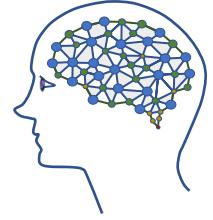
- De-noising auto-encoder



Vincent, Pascal, et al. "Extracting and composing robust features with denoising autoencoders." *ICML*, 2008.



Questions?



Siamese Networks

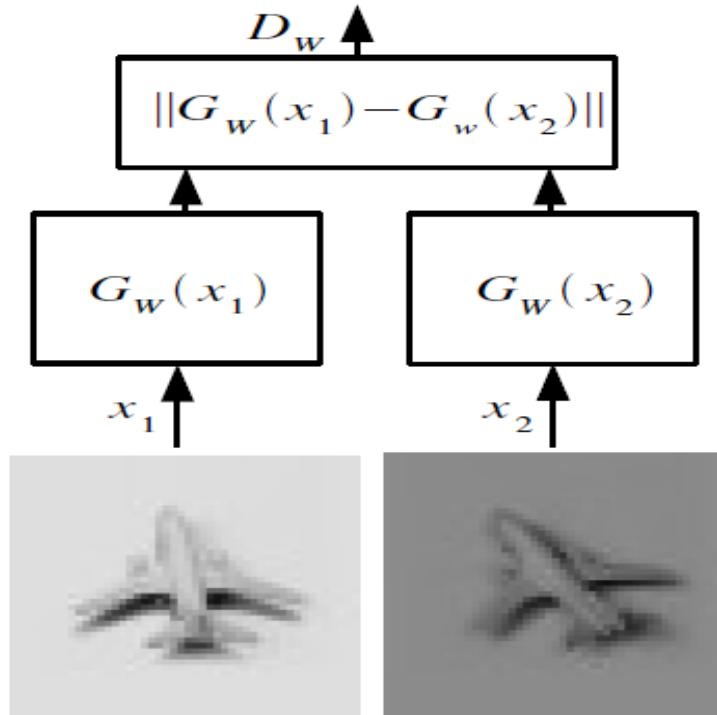
Learning from Pairs (Triplets)



Siamese: Loss

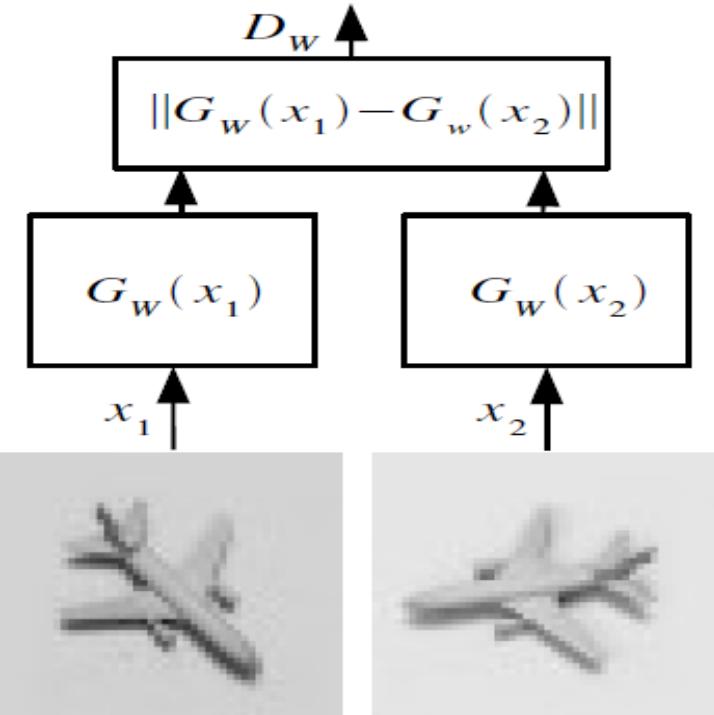


Make this small



Similar images (neighbors
in the neighborhood graph)

Make this large



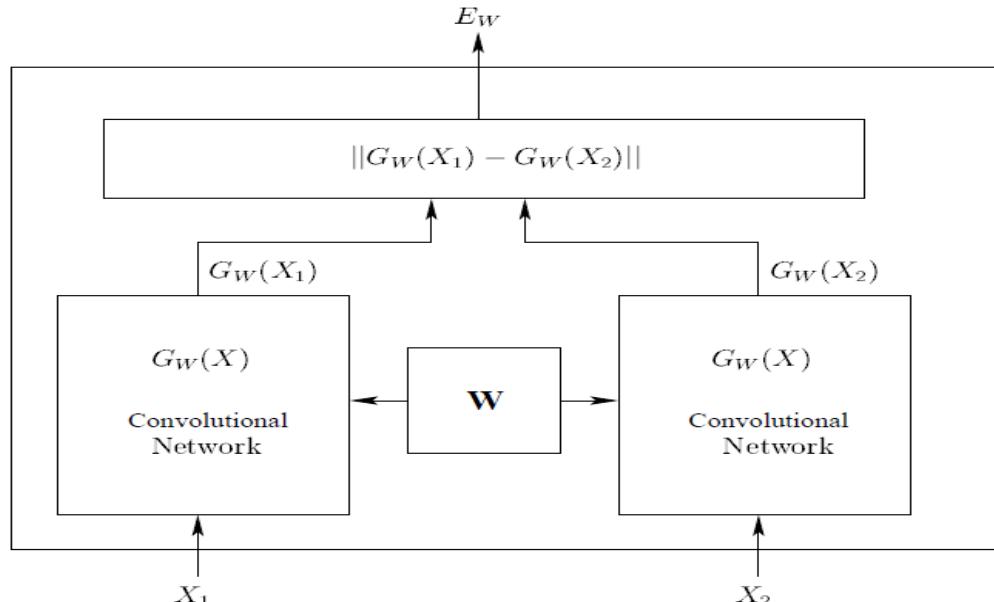
Dissimilar images
(non-neighbors in the
neighborhood graph)



Siamese Architecture

- Given a family of functions $G_W(X)$ parameterized by W , find W such that the similarity metric $D_W(X_1, X_2)$ is small for similar pairs and large for dissimilar pairs:-

$$D_W(X_1, X_2) = ||G_W(X_1) - G_W(X_2)||$$



Loss function

$$\mathcal{L}(W) = \sum_{i=1}^P L(W, (Y, \vec{X}_1, \vec{X}_2)^i)$$

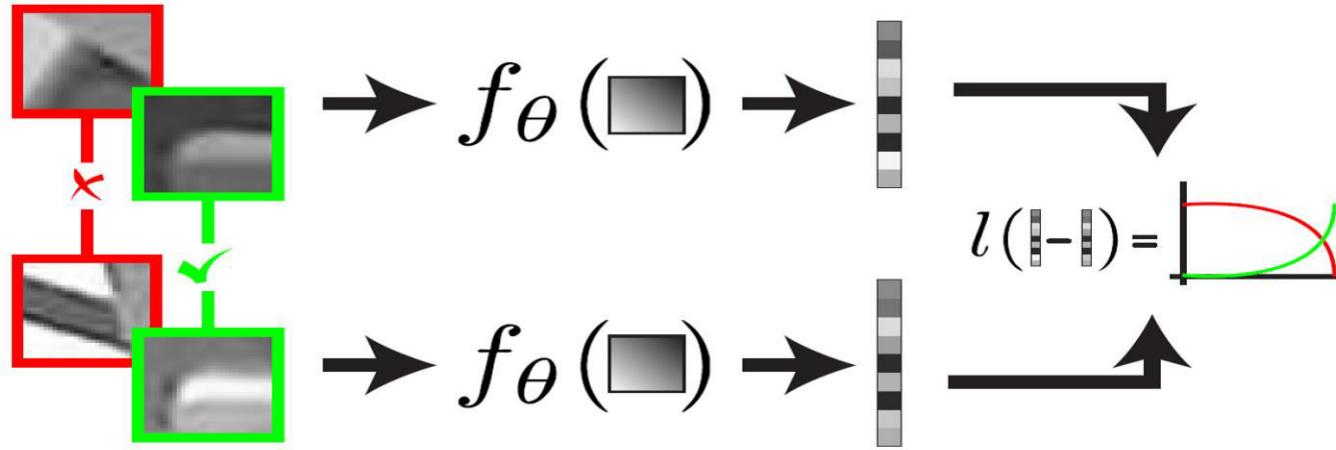
Loss function for similar pairs

$$L(W, (Y, \vec{X}_1, \vec{X}_2)^i) = (1 - Y)L_S(D_W^i) + YL_D(D_W^i)$$

Loss function for dissimilar pairs

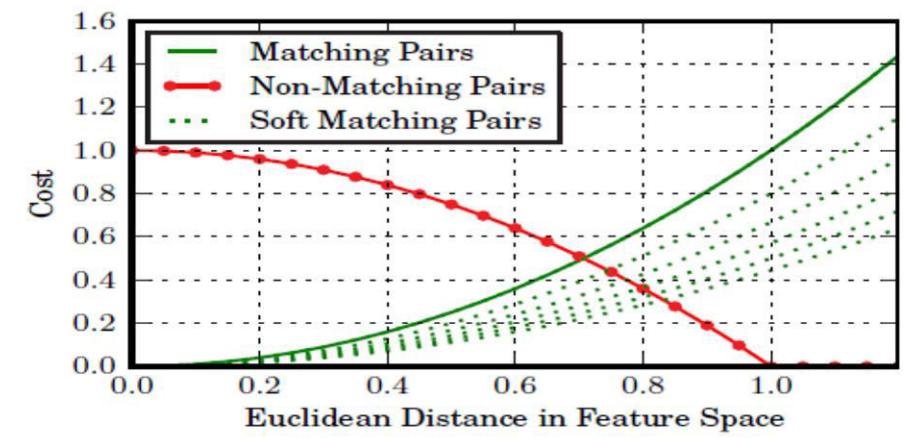


Learning to Match



Using the contrastive cost function

$$l_{\theta} (\mathbf{y}_i, \mathbf{y}_j) = \begin{cases} s_{ij} d_{ij}^2, & \text{if matching} \\ \max(1.0 - d_{ij}^2, 0), & \text{if non-matching} \end{cases}$$





Learning to Match

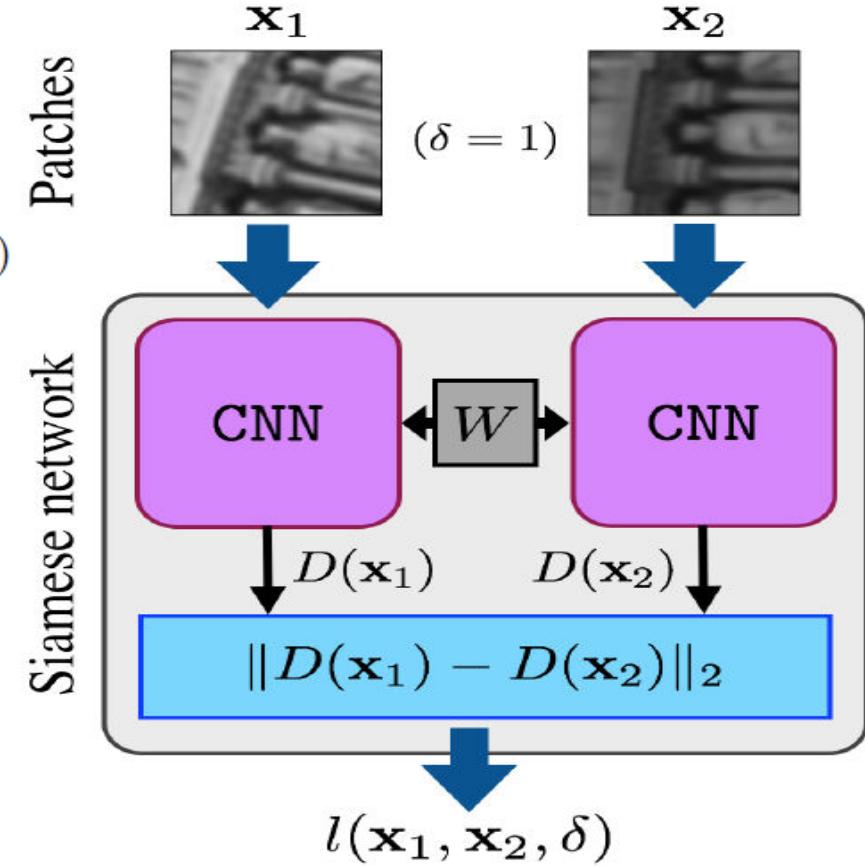


$$d_D(\mathbf{x}_1, \mathbf{x}_2) = \|D(\mathbf{x}_1) - D(\mathbf{x}_2)\|_2$$

$$l(\mathbf{x}_1, \mathbf{x}_2, \delta) = \delta \cdot l_P(d_D(\mathbf{x}_1, \mathbf{x}_2)) + (1 - \delta) \cdot l_N(d_D(\mathbf{x}_1, \mathbf{x}_2))$$

$$l_P(d_D(\mathbf{x}_1, \mathbf{x}_2)) = d_D(\mathbf{x}_1, \mathbf{x}_2)$$

$$l_N(d_D(\mathbf{x}_1, \mathbf{x}_2)) = \max(0, m - d_D(\mathbf{x}_1, \mathbf{x}_2))$$



Stochastic sampling and aggressive mining



Triplet Loss

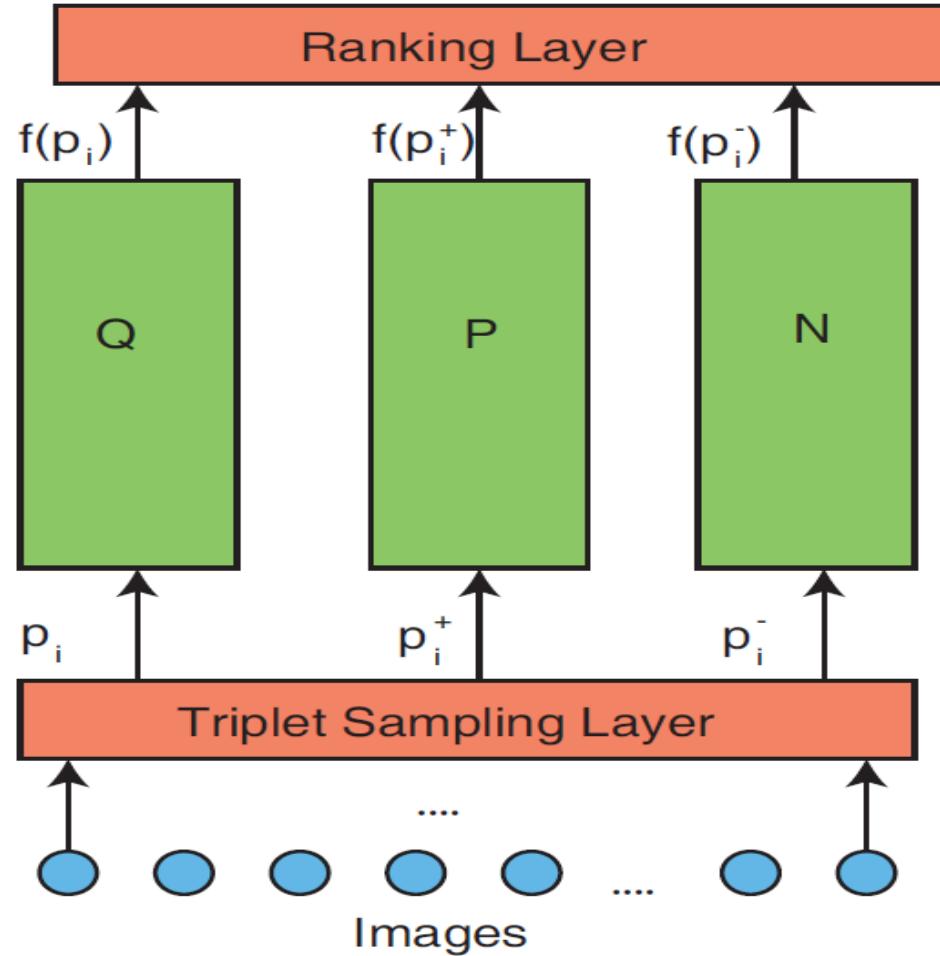
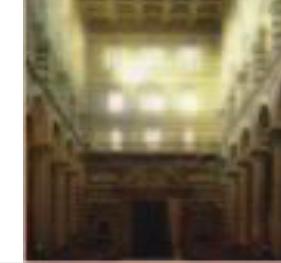
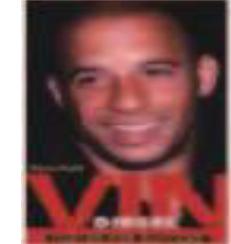
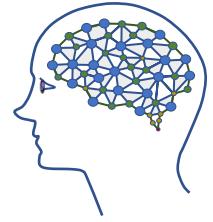
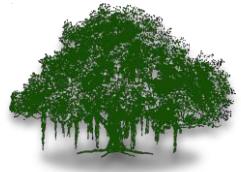




Image Retrieval (Ranking)

Query					
Positive					
Negative					



Questions?

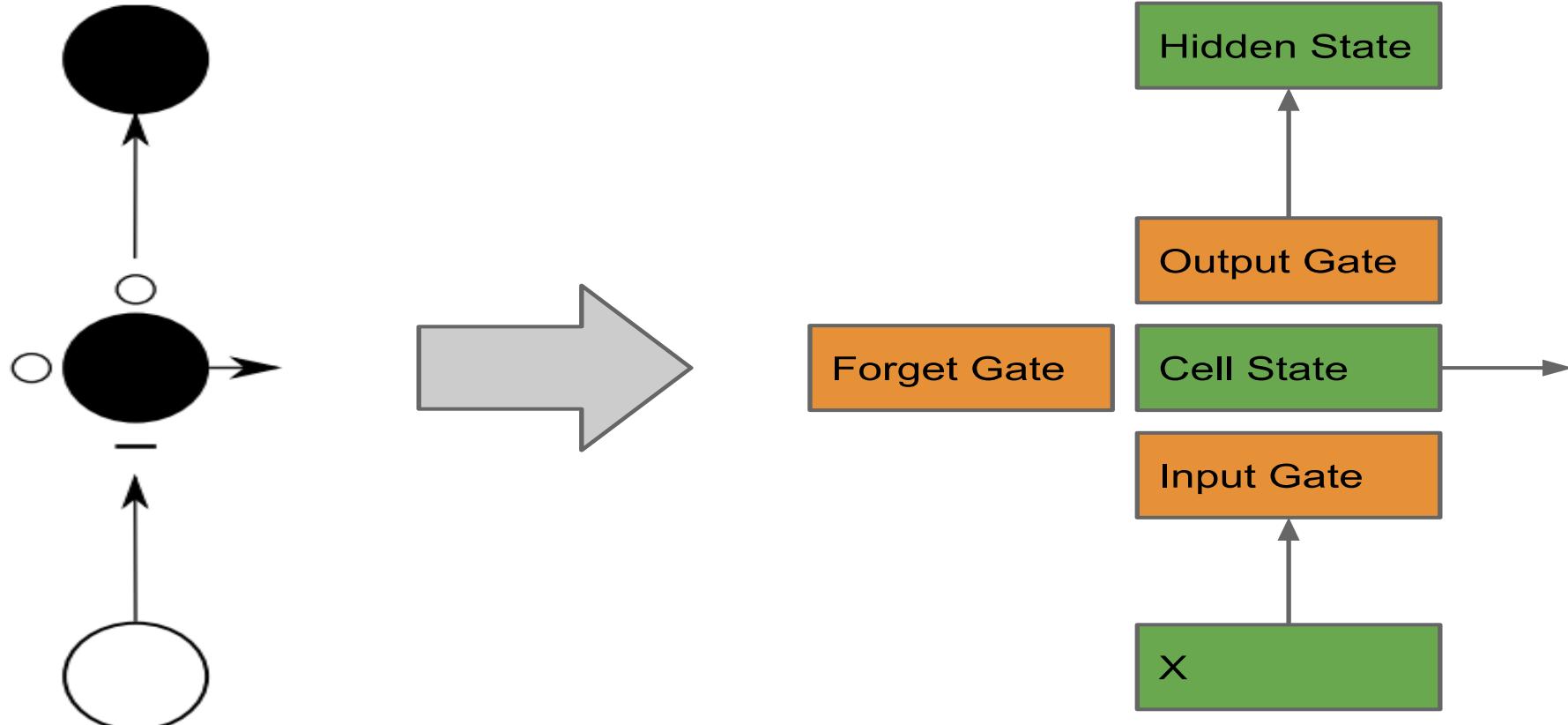


Recurrent Neural Networks

Architecture and Training

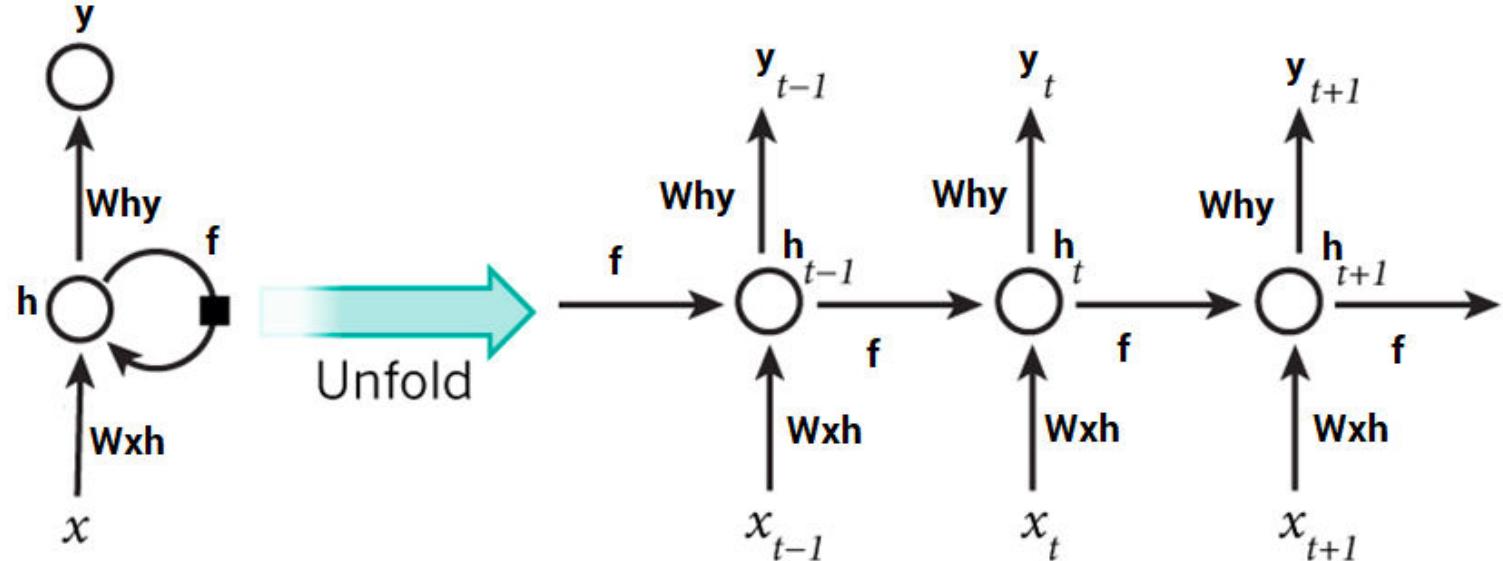
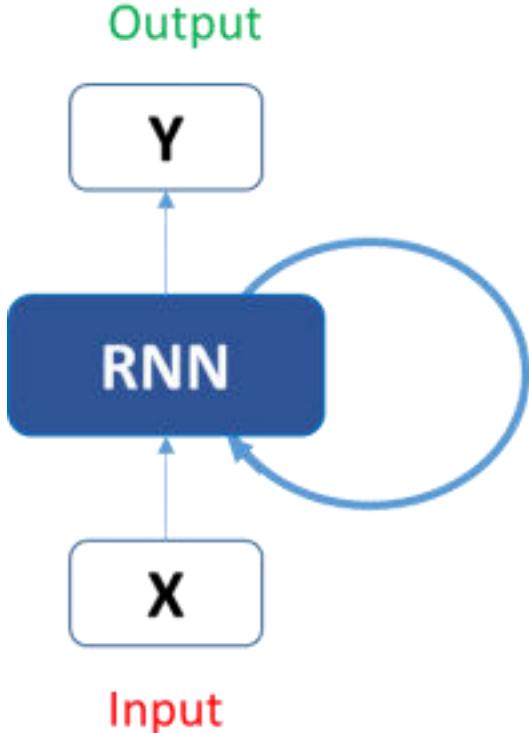


Cell with Memory





RNNs: Overview and Summary: Q?

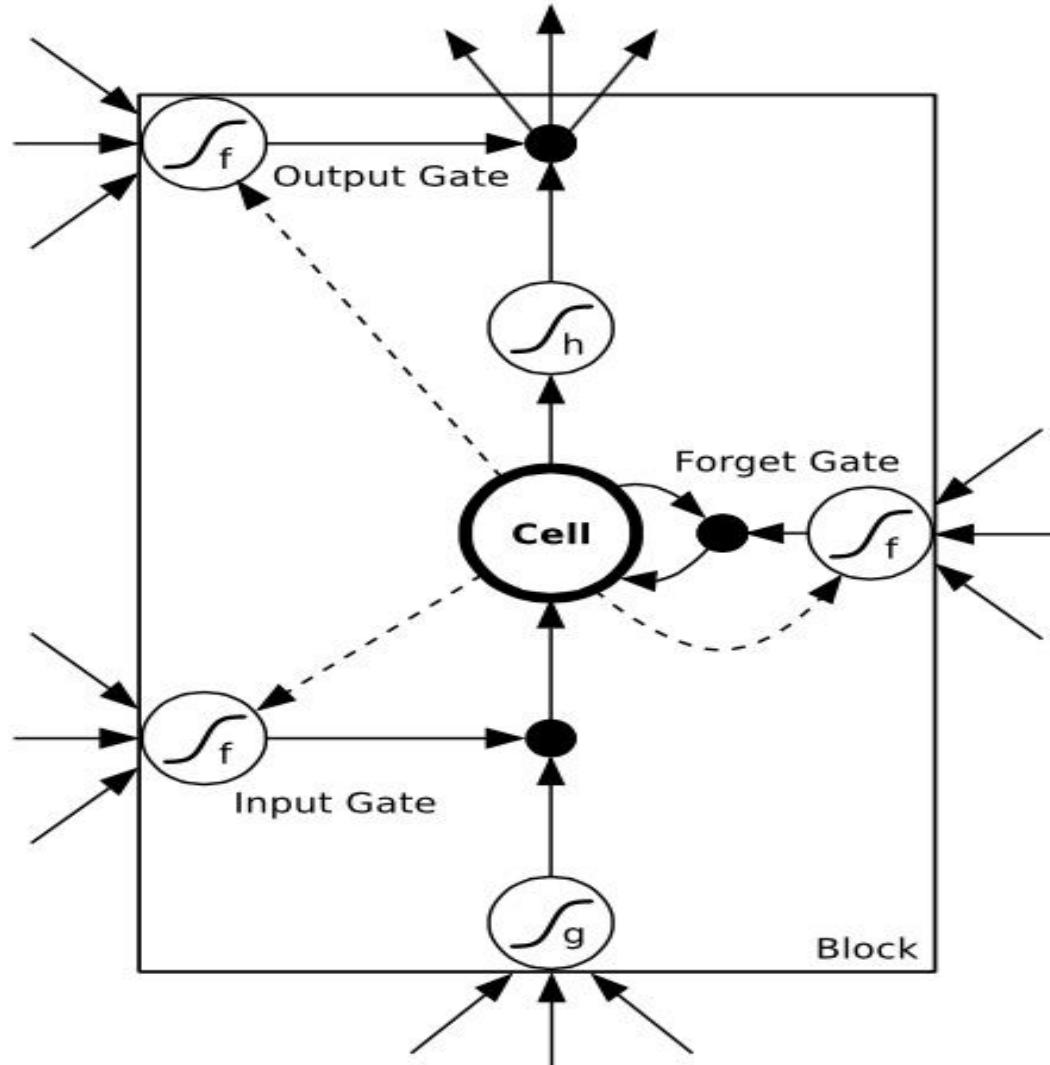


$$h_t = \tanh (W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$



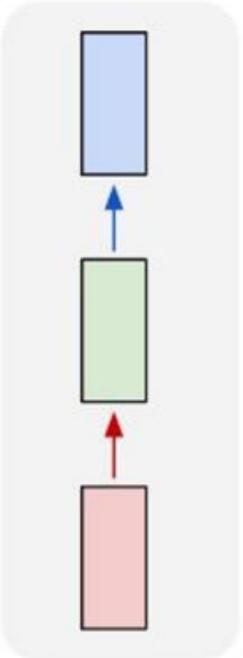
LSTM: Gates





Recurrent Neural Networks (RNNs)

one to one

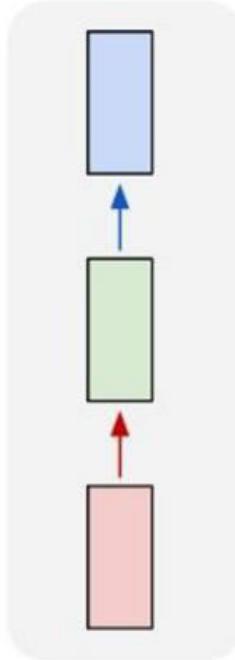


← **Vanilla Neural Networks**

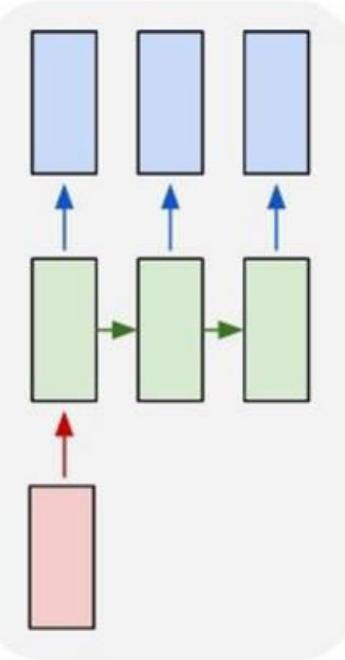


RNNs : Process Sequence

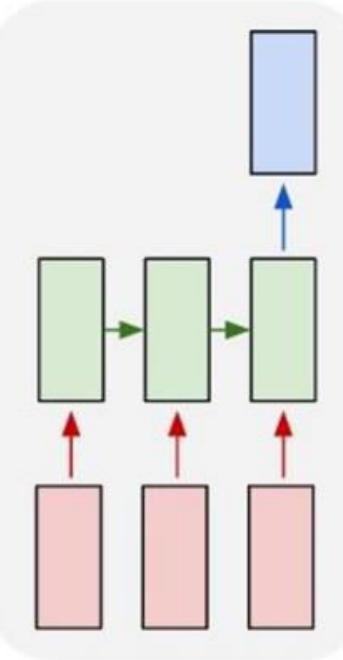
one to one



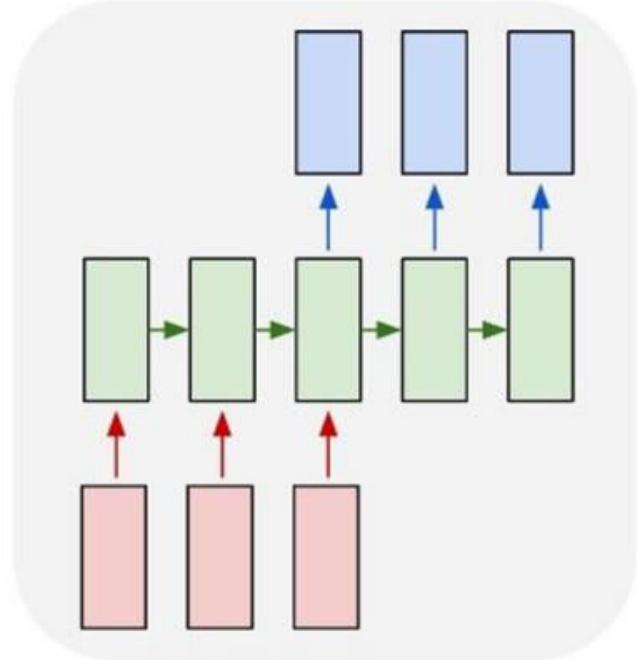
one to many



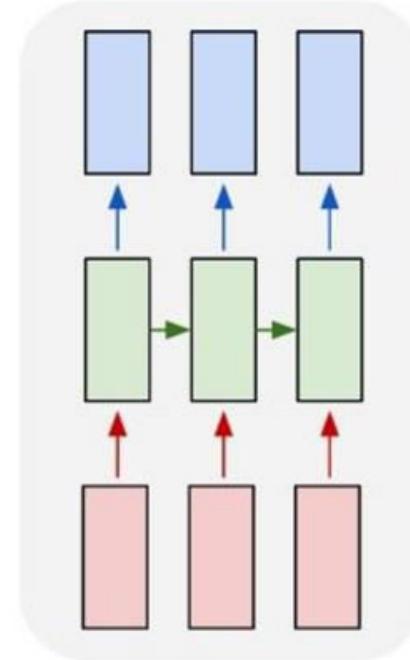
many to one



many to many



many to many

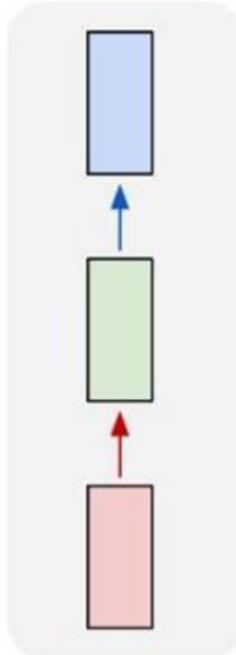


e.g. **Image Captioning**
image -> sequence of words

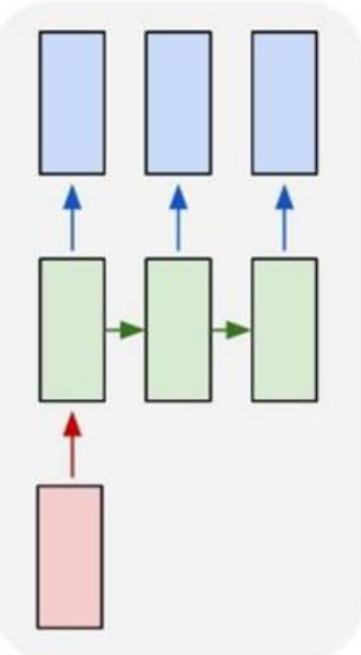


RNNs : Process Sequence

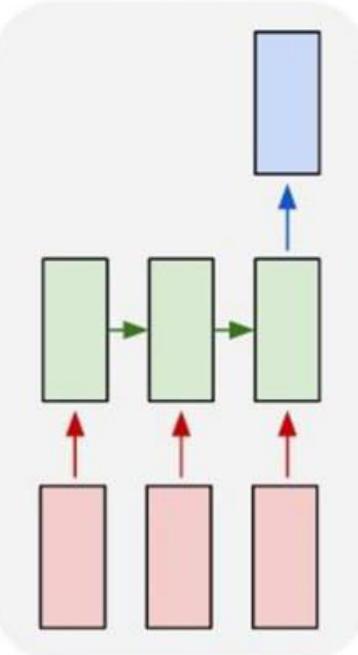
one to one



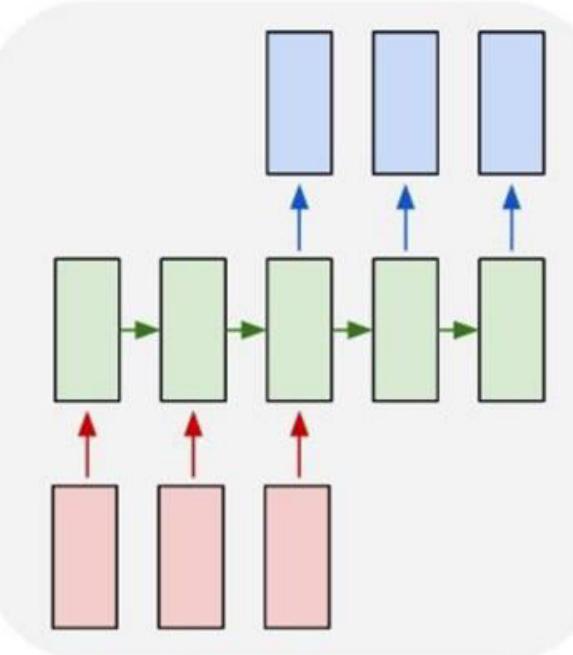
one to many



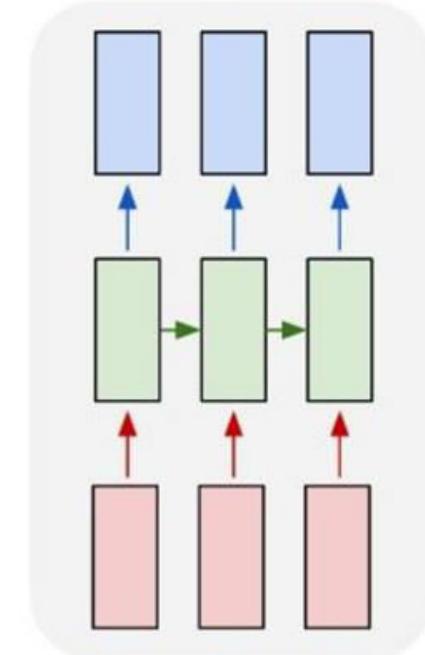
many to one



many to many



many to many

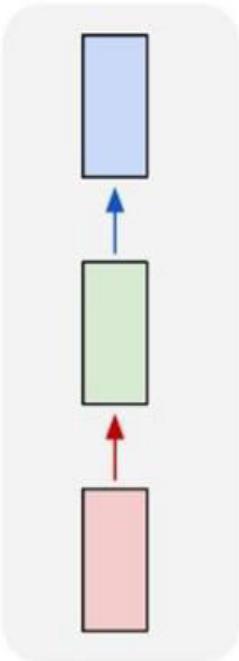


e.g. **Sentiment Classification**
sequence of words -> sentiment

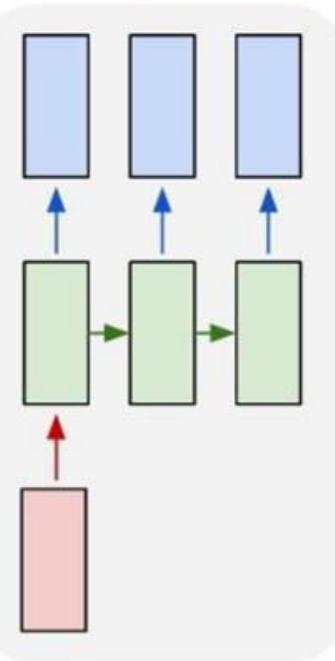


RNNs : Process Sequence

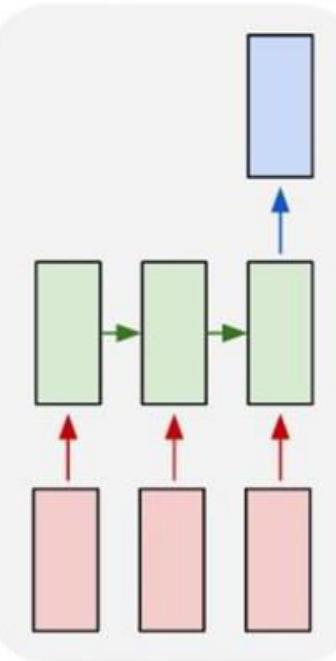
one to one



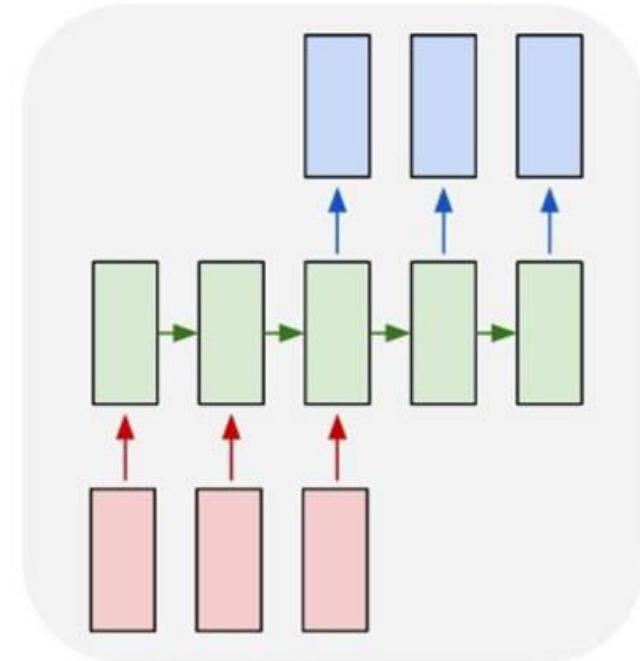
one to many



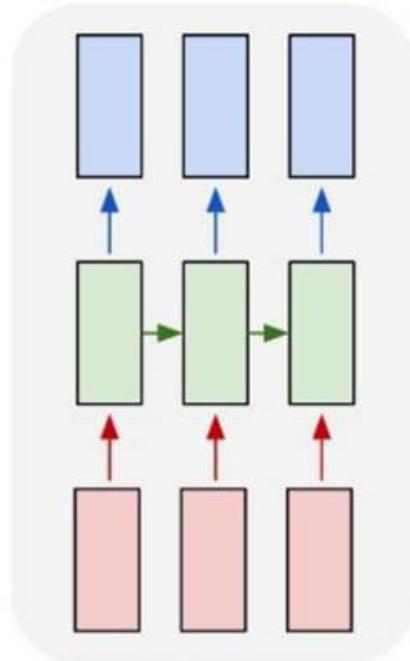
many to one



many to many



many to many

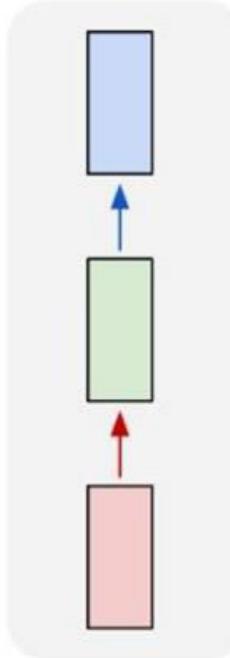


↑
e.g. **Machine Translation**
seq of words -> seq of words

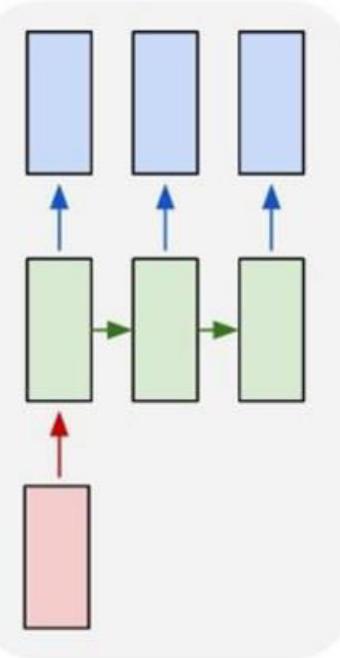


RNNs : Process Sequence

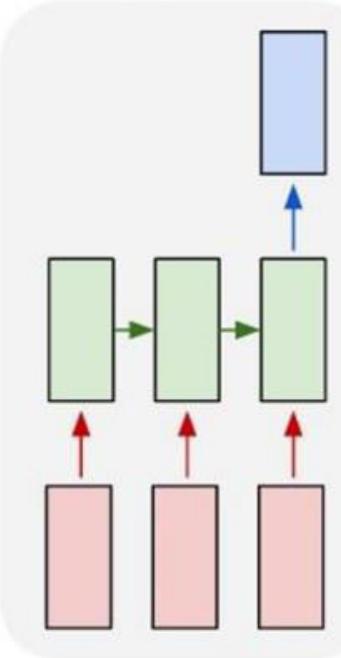
one to one



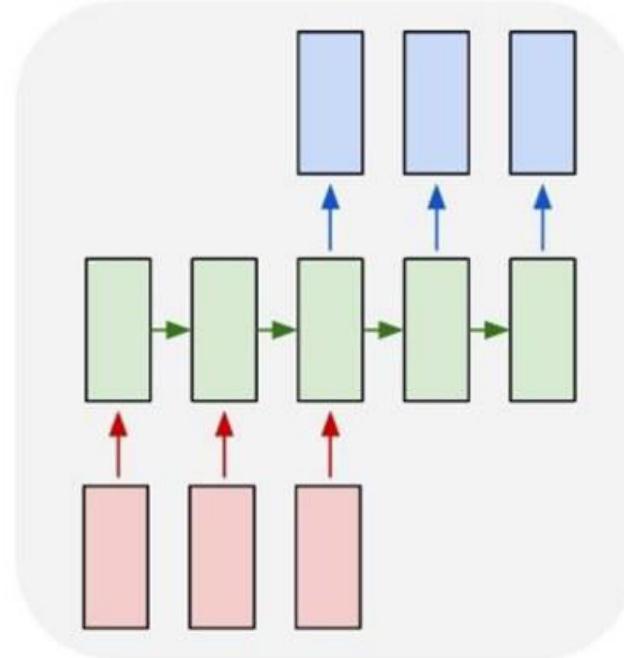
one to many



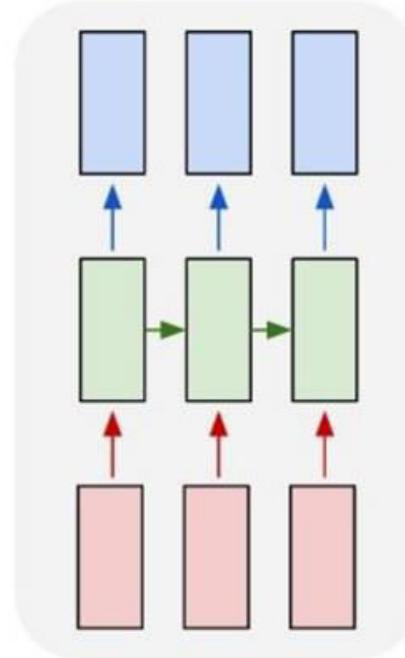
many to one



many to many



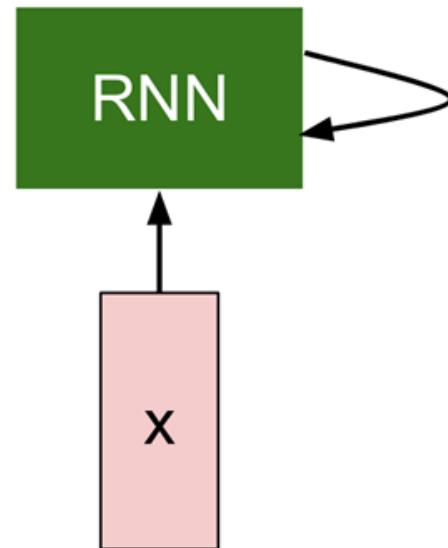
many to many



e.g. Video classification on frame level

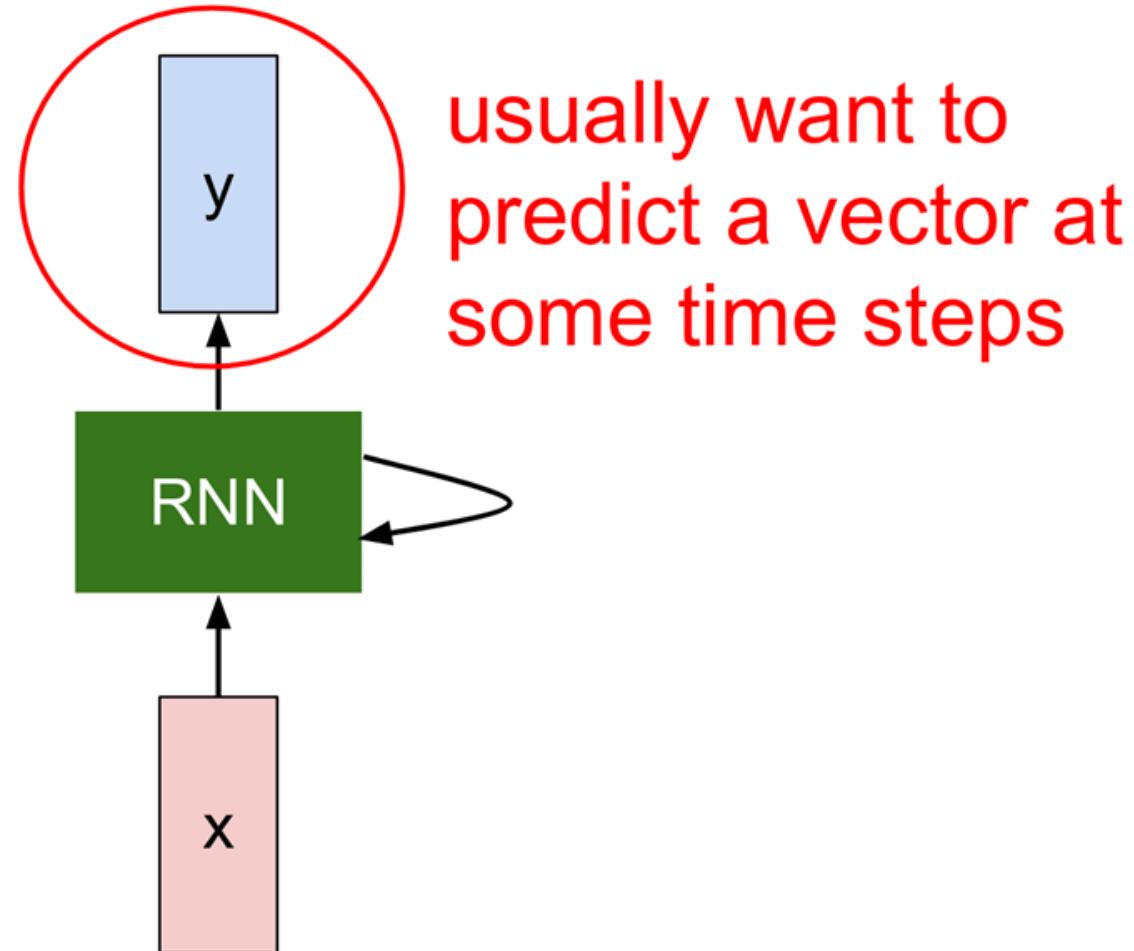


RNNs





RNNs

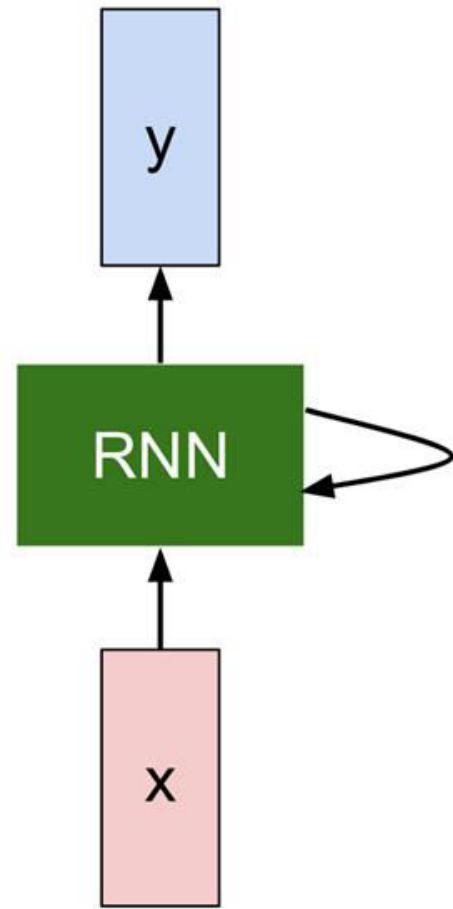




RNNs



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





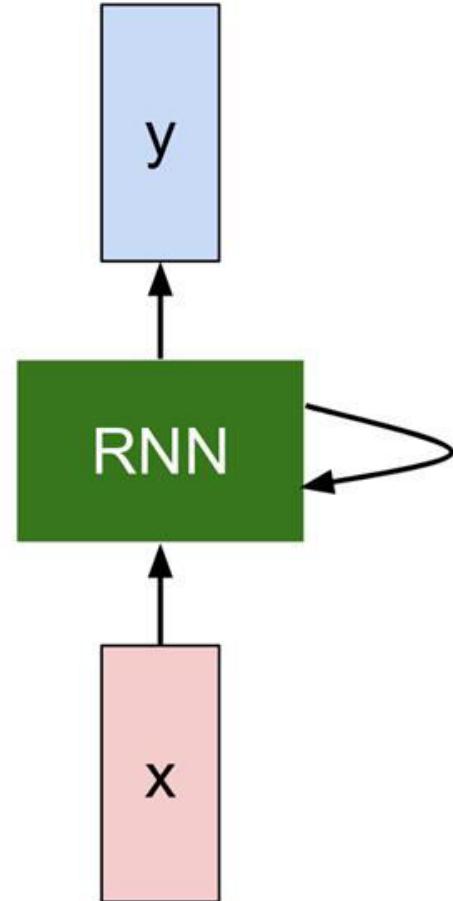
RNNs



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

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

Notice: the same function and the same set of parameters are used at every time step.

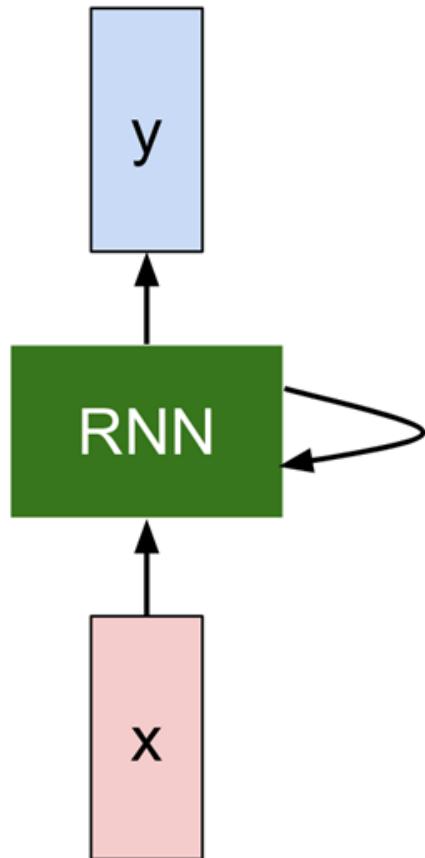




RNNs



The state consists of a single “*hidden*” vector \mathbf{h} :



$$\mathbf{h}_t = f_W(\mathbf{h}_{t-1}, \mathbf{x}_t)$$

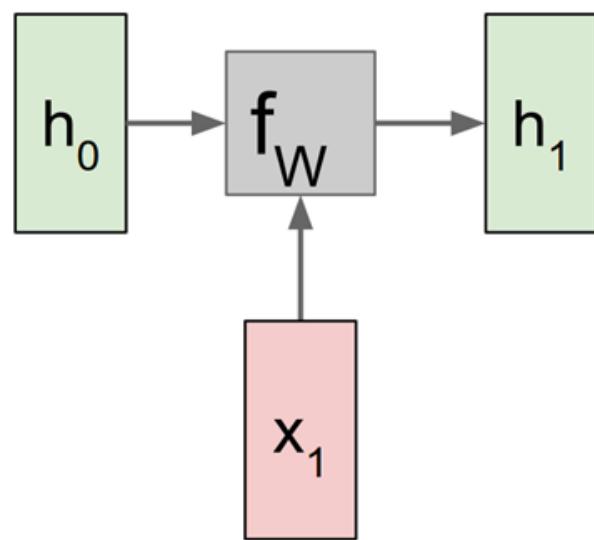


$$\mathbf{h}_t = \tanh(W_{hh}\mathbf{h}_{t-1} + W_{xh}\mathbf{x}_t)$$

$$y_t = W_{hy}\mathbf{h}_t$$

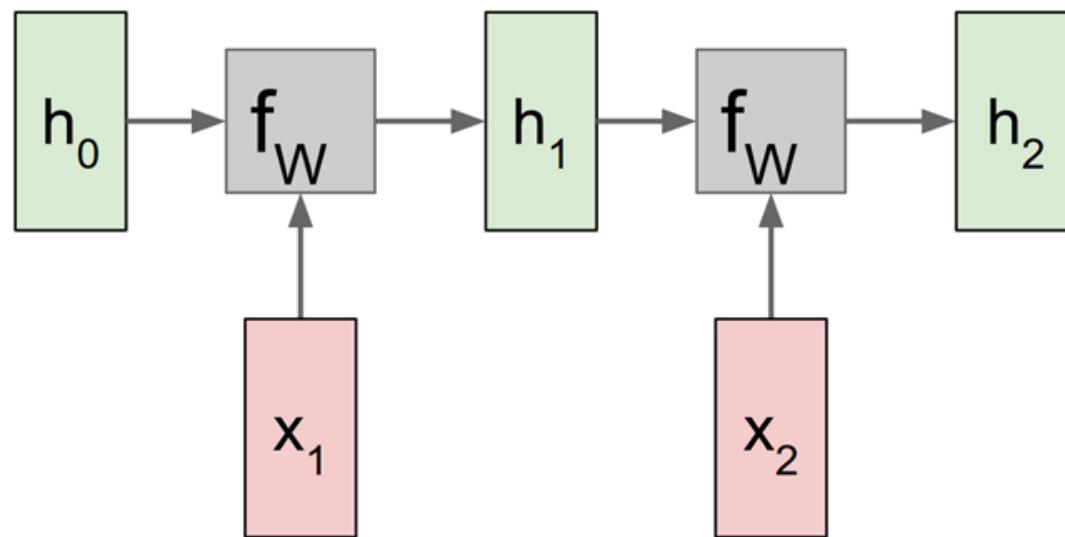


RNN : Computation Graph



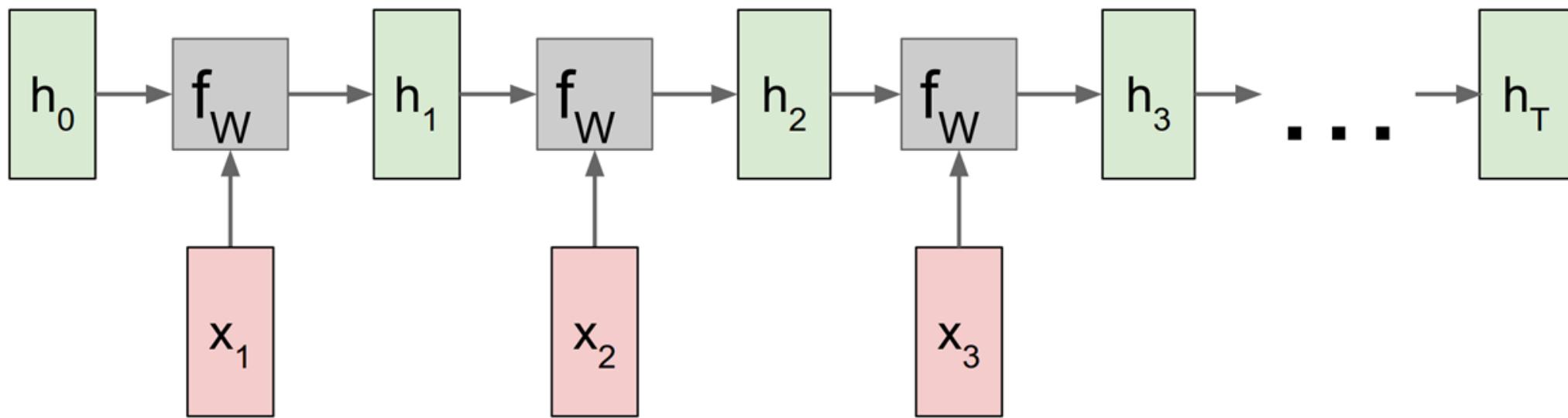


RNN : Computation Graph





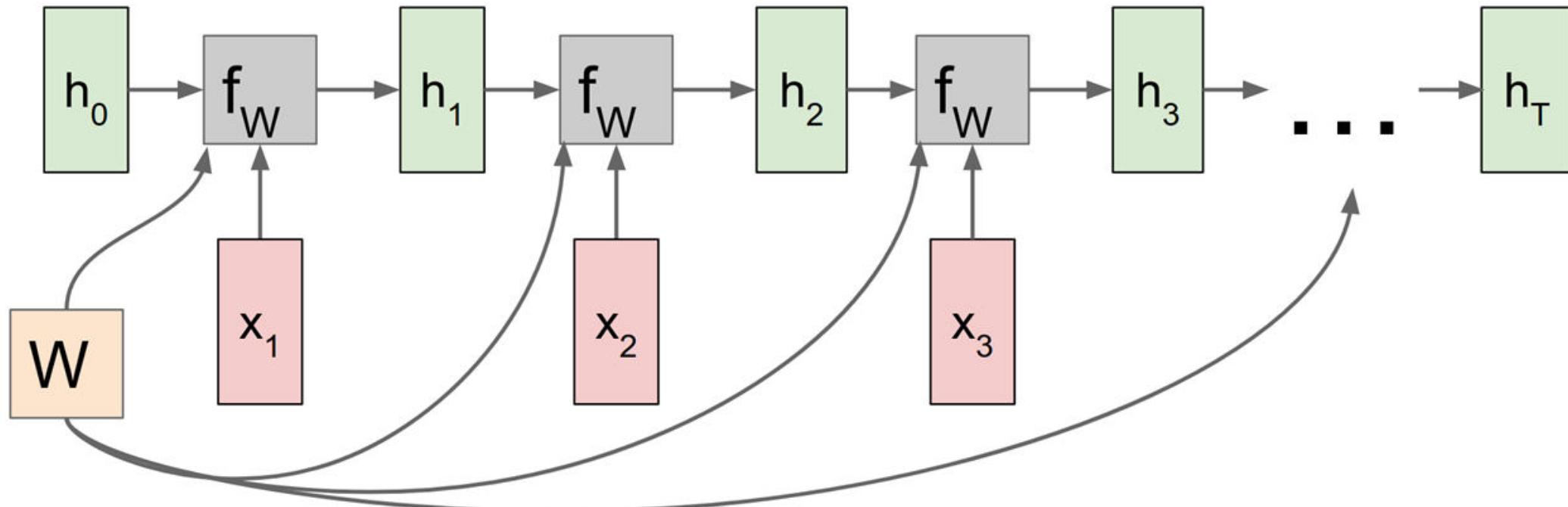
RNN : Computation Graph





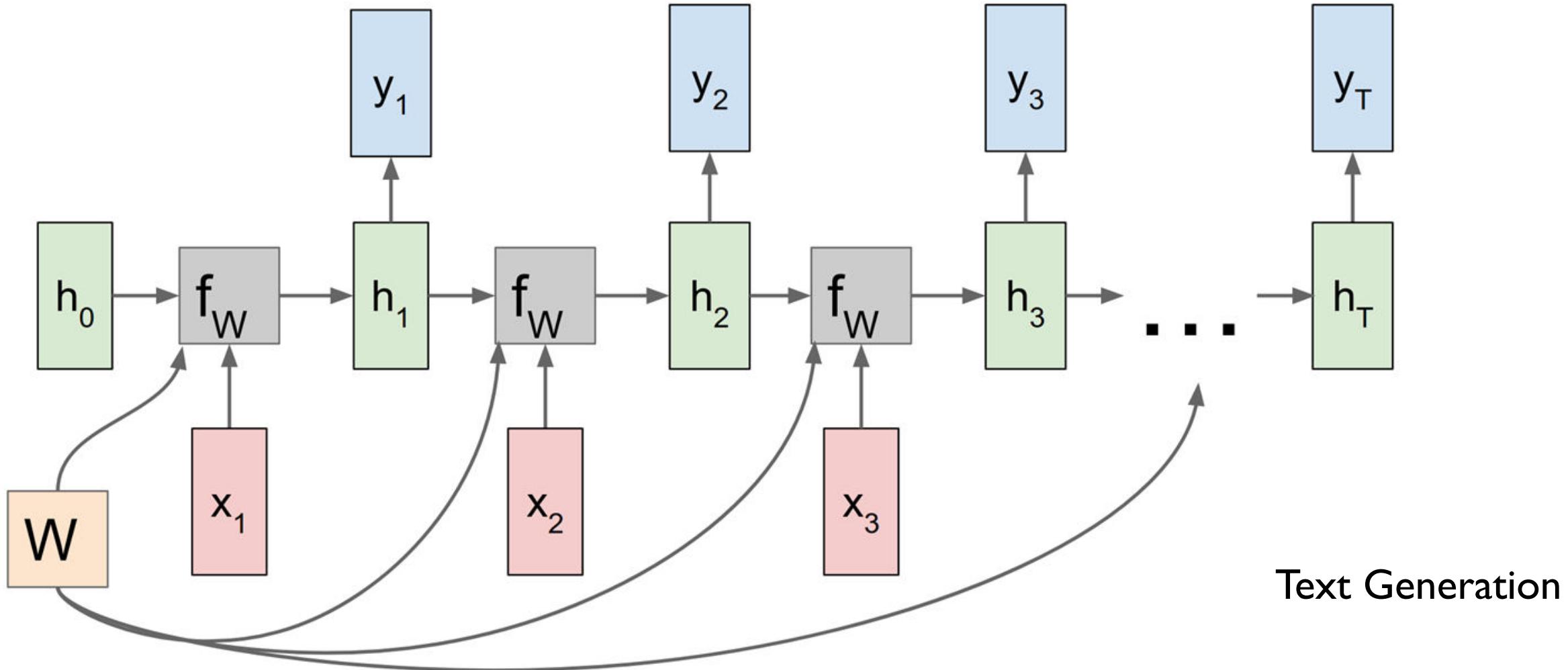
RNN : Computation Graph

Re-use the same weight matrix at every time-step



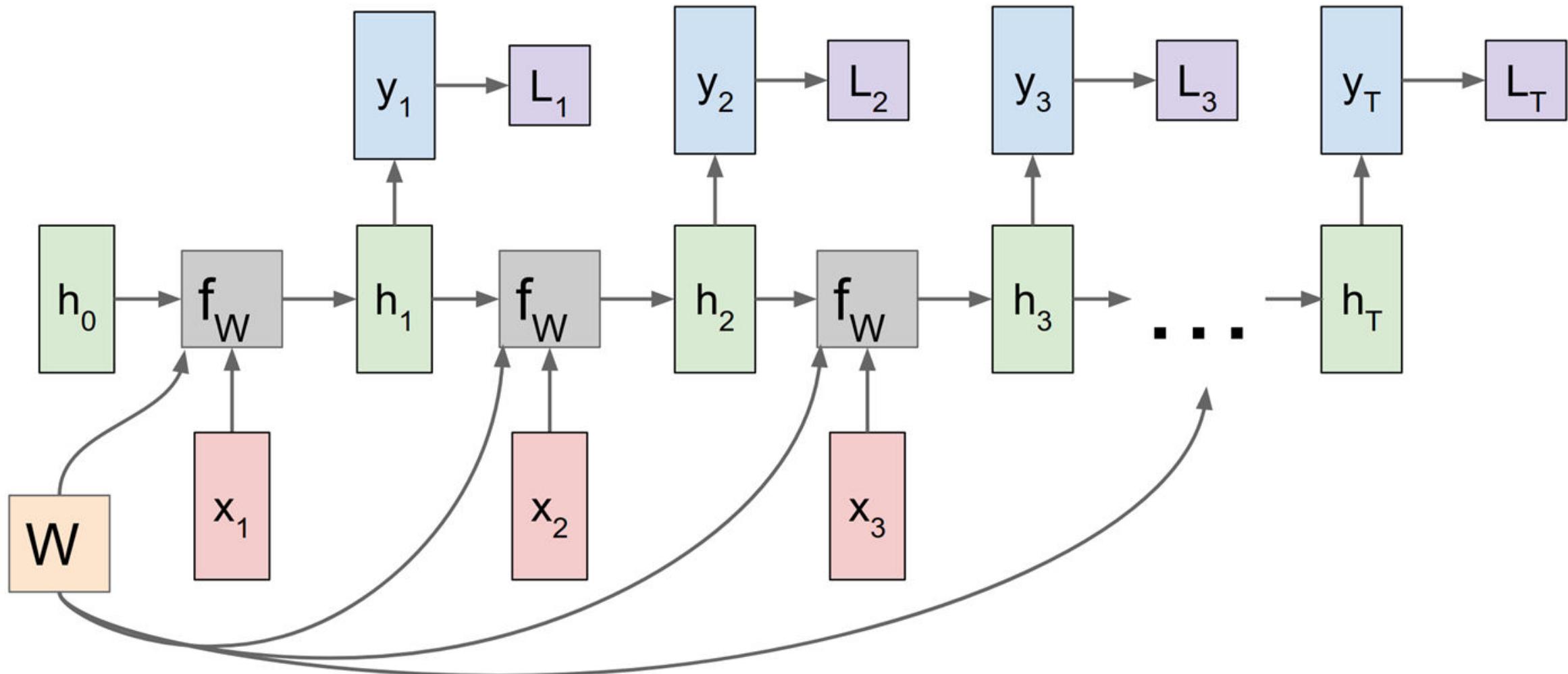


RNN: Many to Many





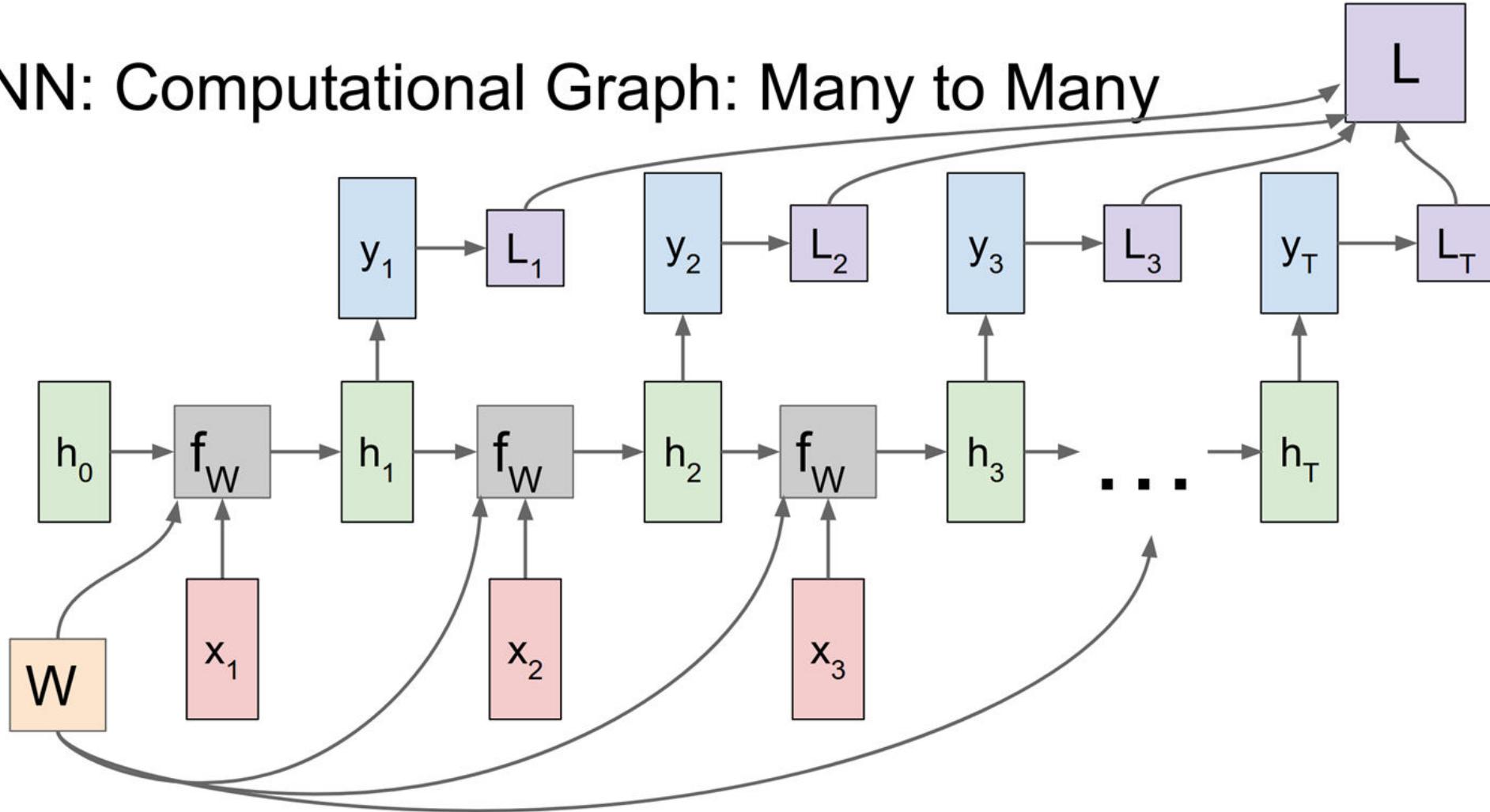
Loss Function for Many to Many





Loss Function for Many to Many

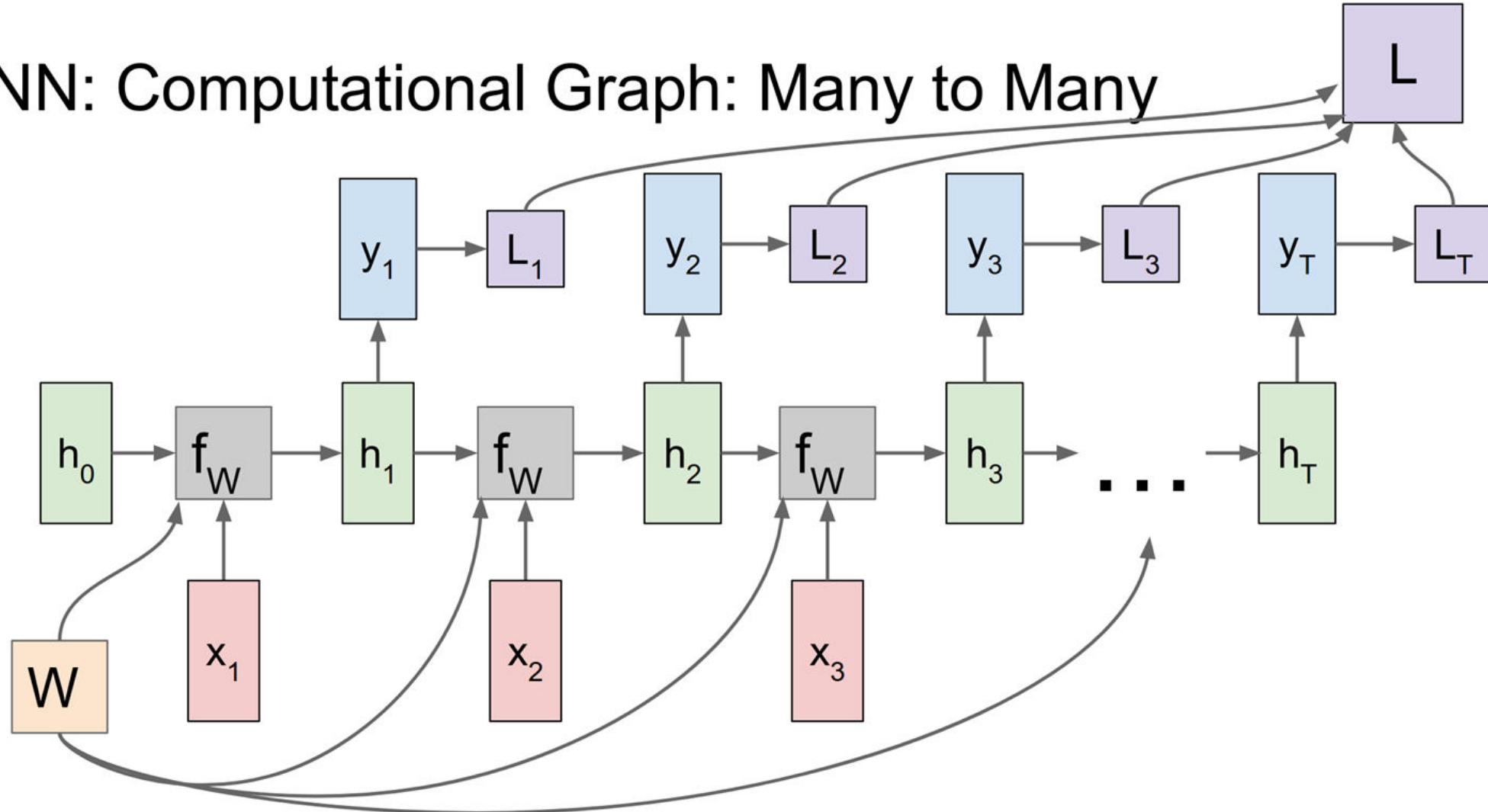
RNN: Computational Graph: Many to Many





Loss Function for Many to Many

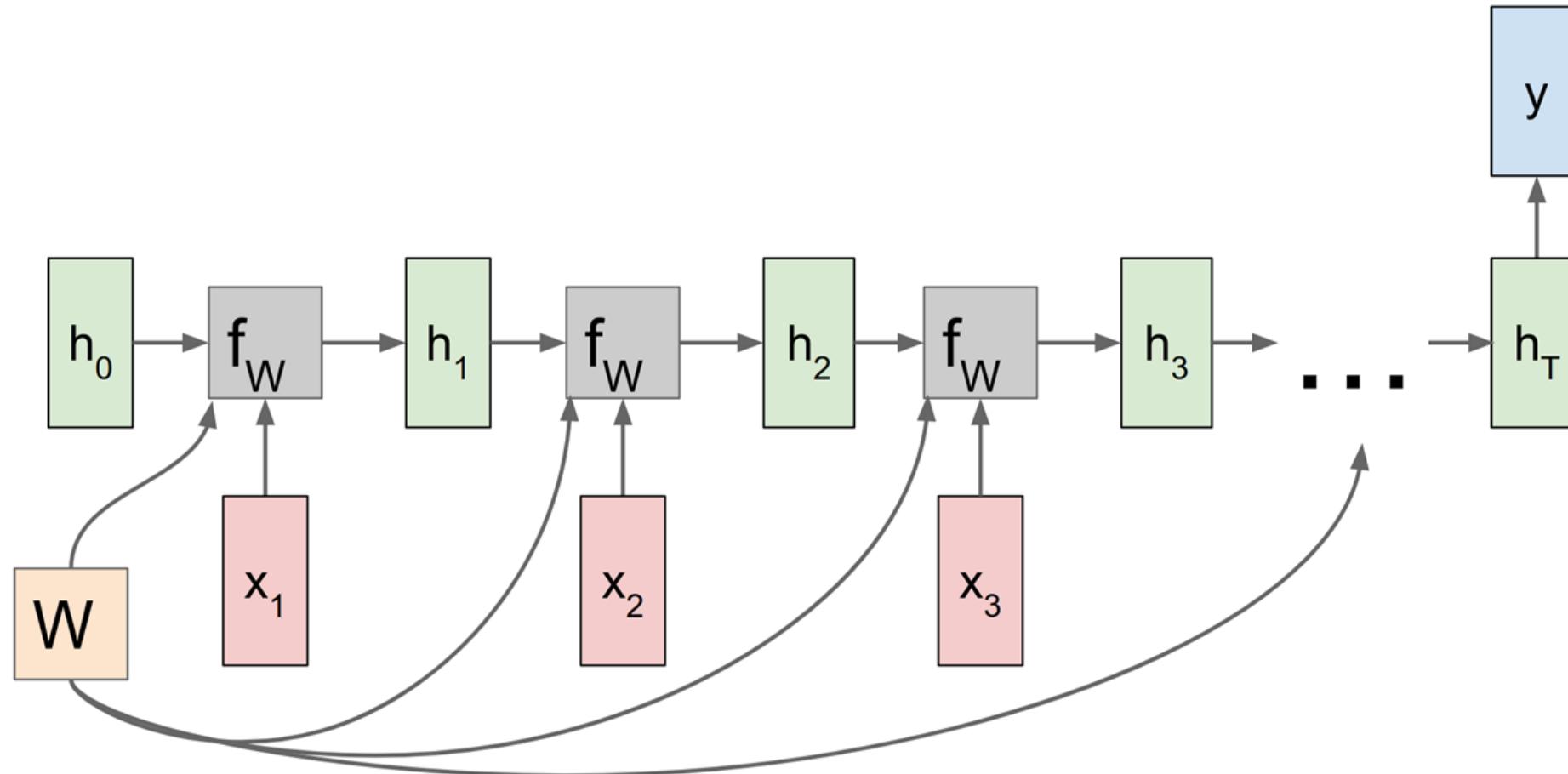
RNN: Computational Graph: Many to Many





Many to One (Sentiment Analysis)

RNN: Computational Graph: Many to One





One to Many (Image Captioning)

RNN: Computational Graph: One to Many

