

Deep Learning

Mayank Vatsa

ADAM

For each Parameter w^j

(j subscript dropped for clarity)

$$\nu_t = \beta_1 * \nu_{t-1} + (1 - \beta_1) * g_t$$

$$s_t = \beta_2 * s_{t-1} + (1 - \beta_2) * g_t^2$$

$$\Delta\omega_t = -\eta \frac{\nu_t}{\sqrt{s_t + \epsilon}} * g_t$$

$$\omega_{t+1} = \omega_t + \Delta\omega_t$$

η : Initial Learning rate

g_t : Gradient at time t along ω^j

ν_t : Exponential Average of gradients along ω_j

s_t : Exponential Average of squares of gradients along ω_j

β_1, β_2 : Hyperparameters

ADAM

For each Parameter w^j

(j subscript dropped for clarity)

Exponential Weighted Averages for past gradients

$$\nu_t = \beta_1 * \nu_{t-1} + (1 - \beta_1) * g_t$$
$$s_t = \beta_2 * s_{t-1} + (1 - \beta_2) * g_t^2$$

$$\Delta\omega_t = -\eta \frac{\nu_t}{\sqrt{s_t + \epsilon}} * g_t$$

$$\omega_{t+1} = \omega_t + \Delta\omega_t$$

η : Initial Learning rate

g_t : Gradient at time t along ω^j

ν_t : Exponential Average of gradients along ω_j

s_t : Exponential Average of squares of gradients along ω_j

β_1, β_2 : Hyperparameters

ADAM

For each Parameter w^j

(j subscript dropped for clarity)

$$\nu_t = \beta_1 * \nu_{t-1} + (1 - \beta_1) * g_t$$

$$s_t = \beta_2 * s_{t-1} + (1 - \beta_2) * g_t^2$$

Exponential Weighted Averages for past squared gradients

- more weight to recent gradients

$$\Delta\omega_t = \frac{\nu_t}{\sqrt{s_t + \epsilon}} * g_t$$

$$\omega_{t+1} = \omega_t + \Delta\omega_t$$

η : Initial Learning rate

g_t : Gradient at time t along ω^j

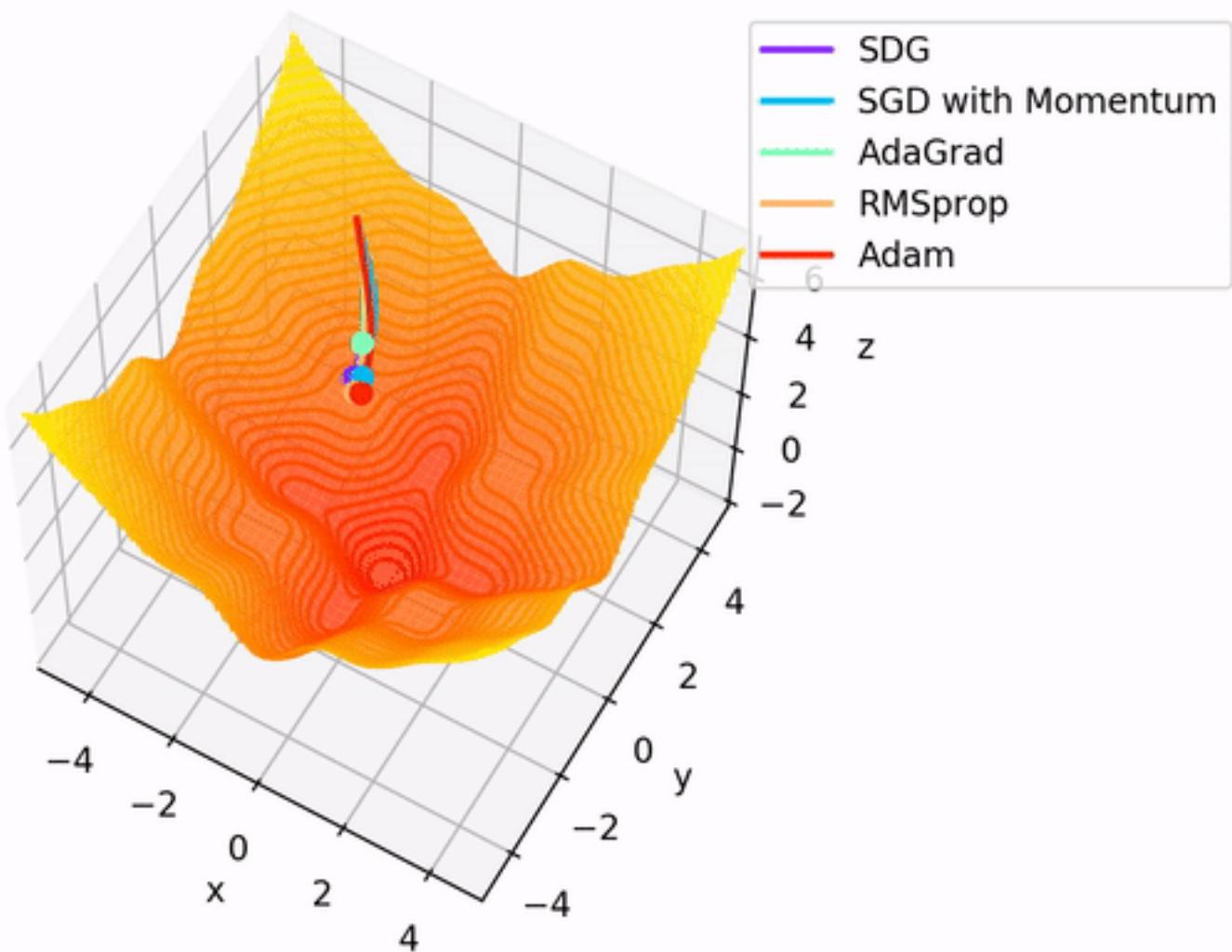
ν_t : Exponential Average of gradients along ω_j

s_t : Exponential Average of squares of gradients along ω_j

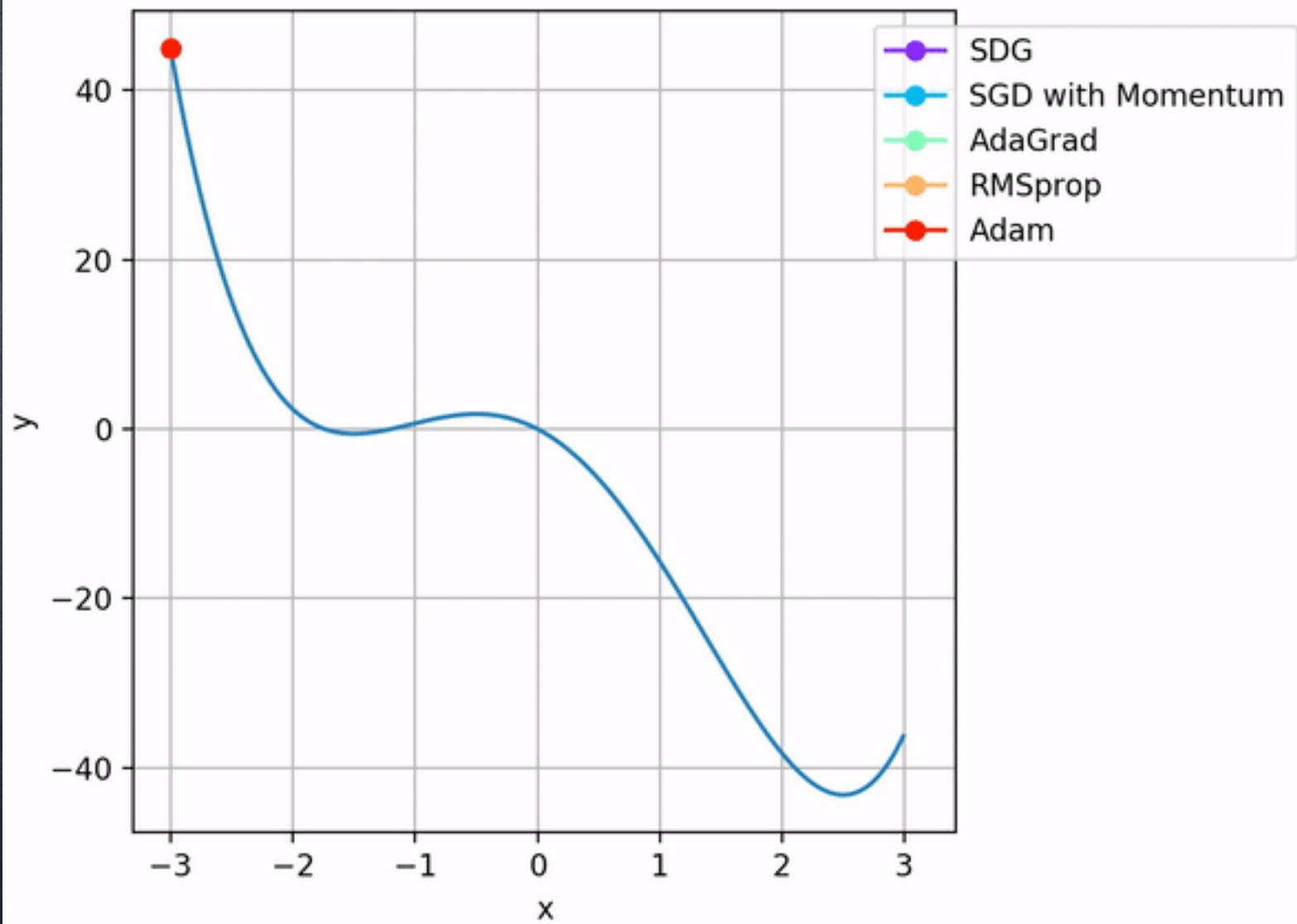
β_1, β_2 : Hyperparameters

ADAM

Optimizer Comparison



Optimizer Comparison

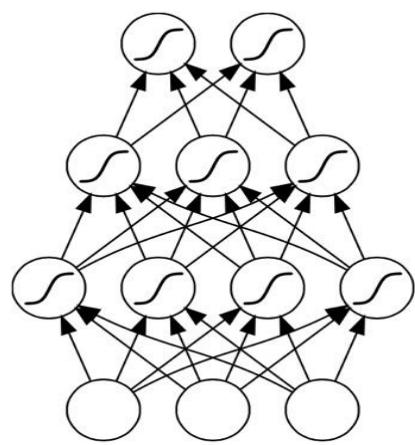


Designing CNN for a given problem

- Image/data size: for example 28X28
- Overall data size: n (50K)
- Design constraints
 - Number of layers (conv, FC, etc)
 - Loss constraints (problem specific)
- Regularization(s) and any other constraints

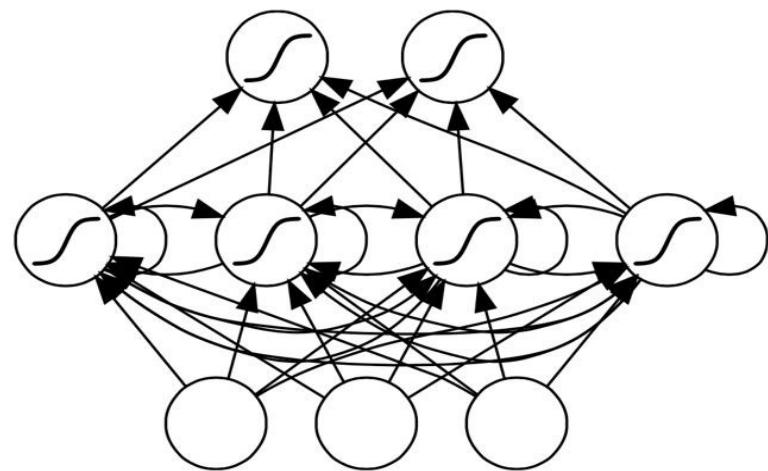
RNN: Another Class
of Neural Network

RNN



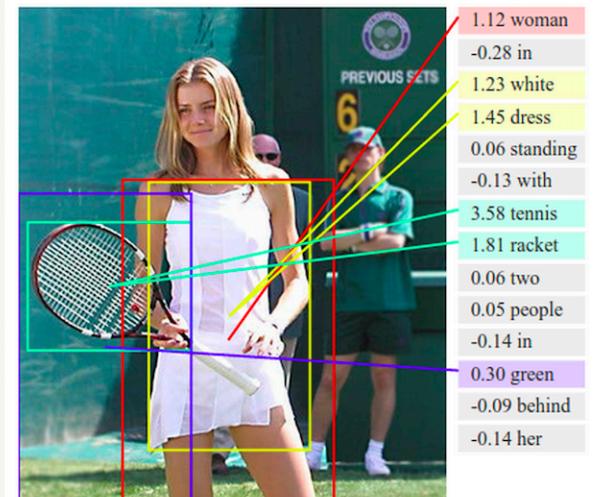
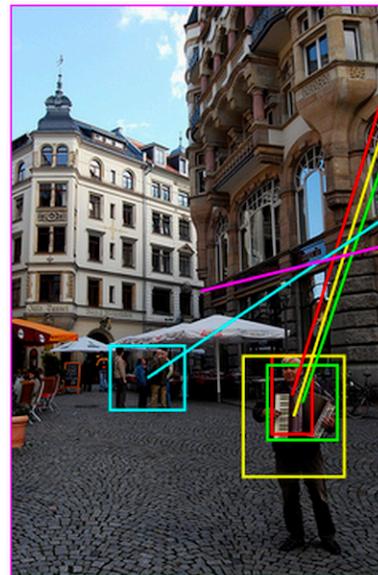
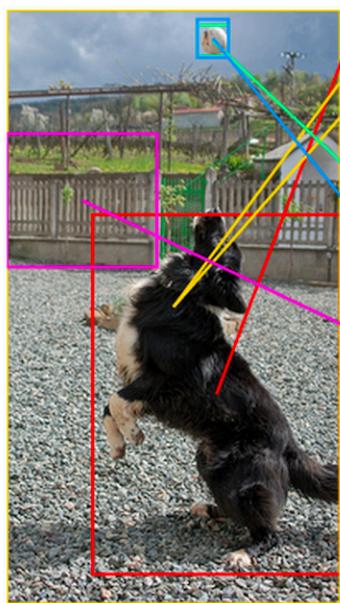
Multi-layer
Perceptron

Output Layer
Hidden Layers
Input Layer

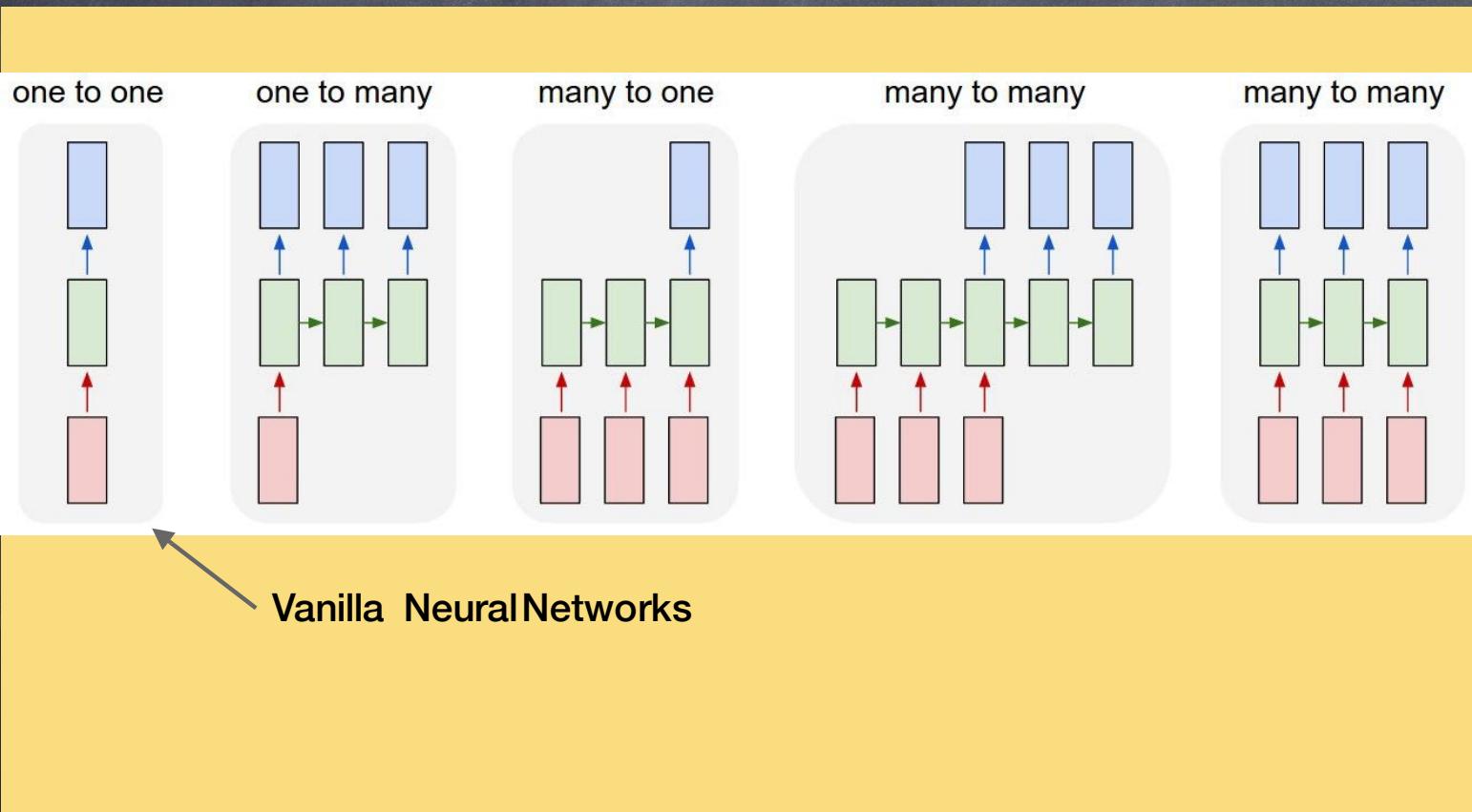


Recurrent Network

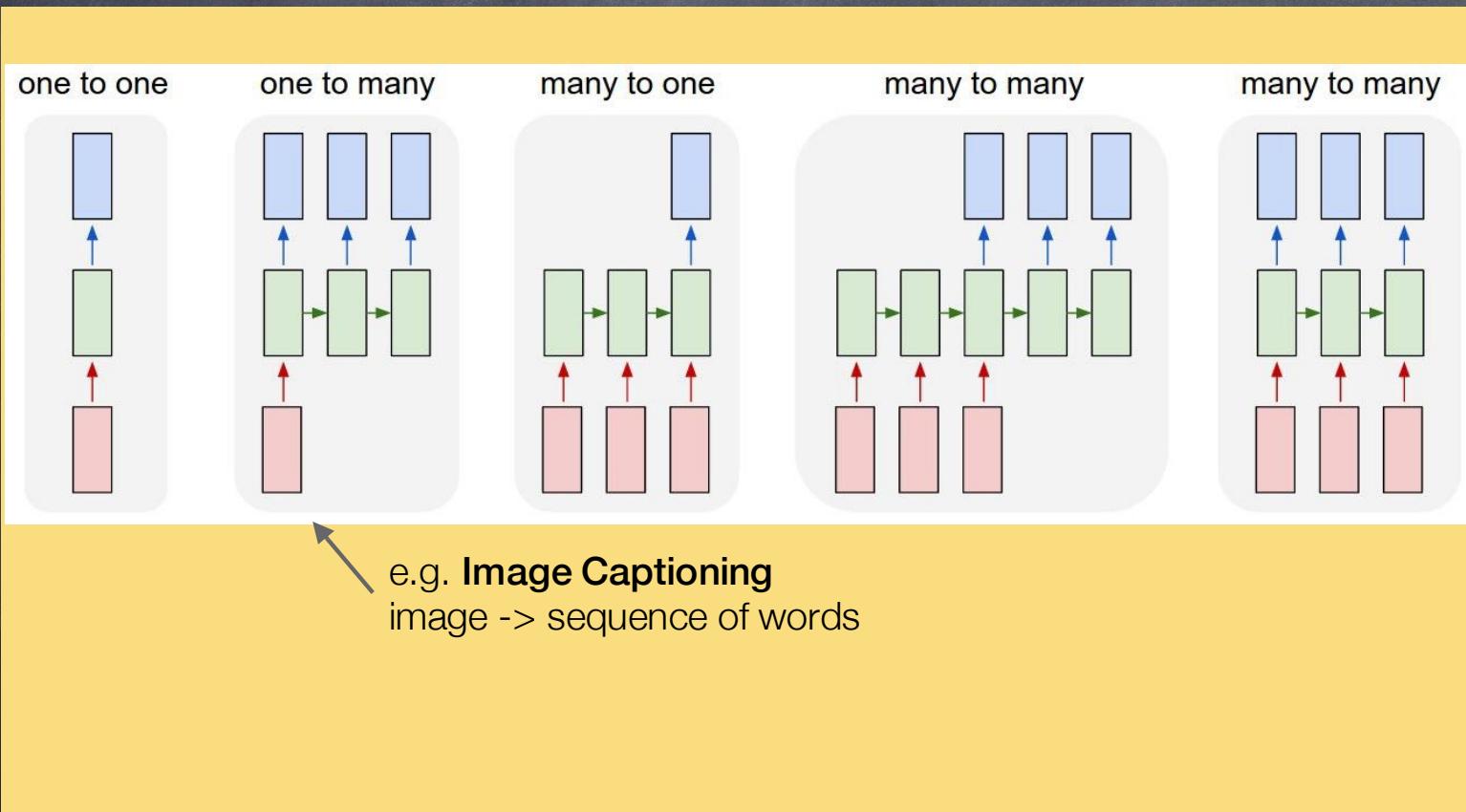
Sequence Modelling



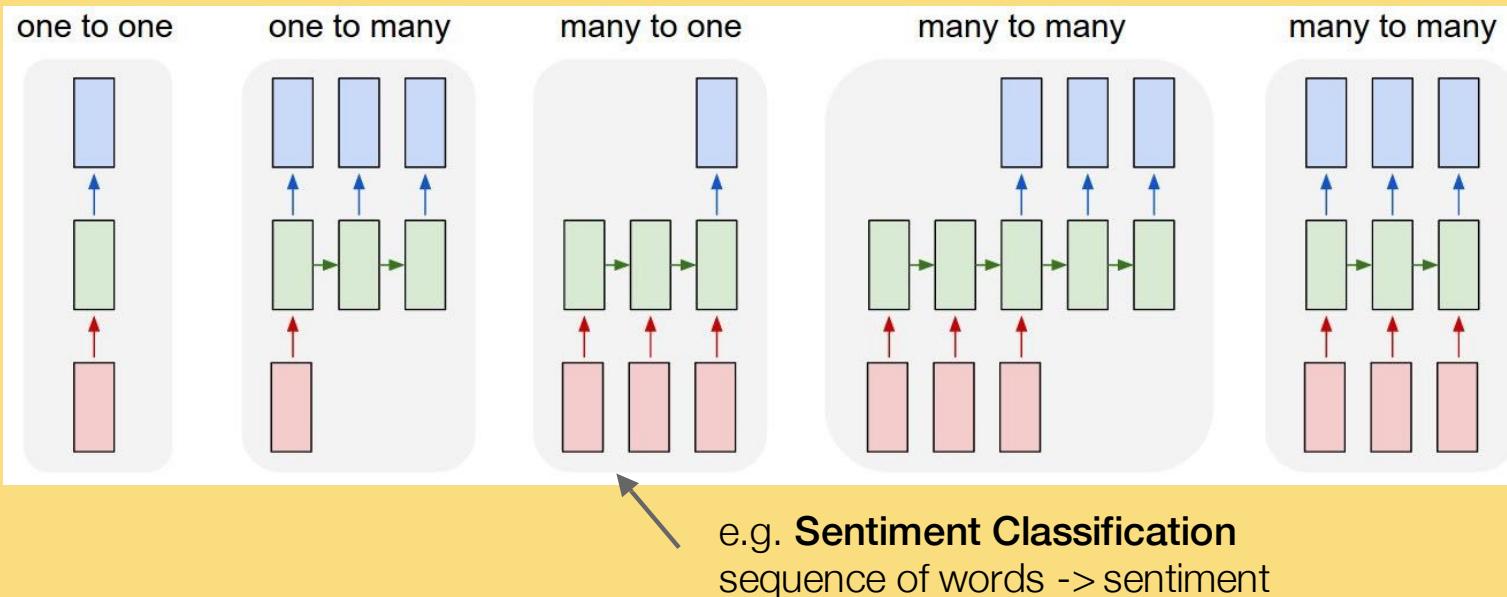
RNN



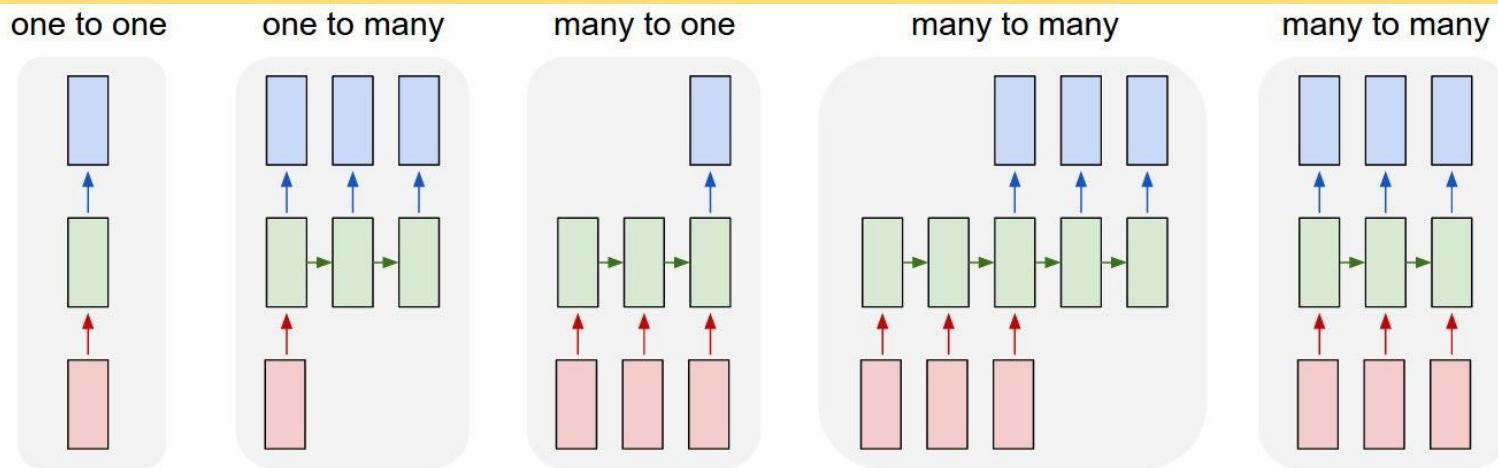
RNN



RNN

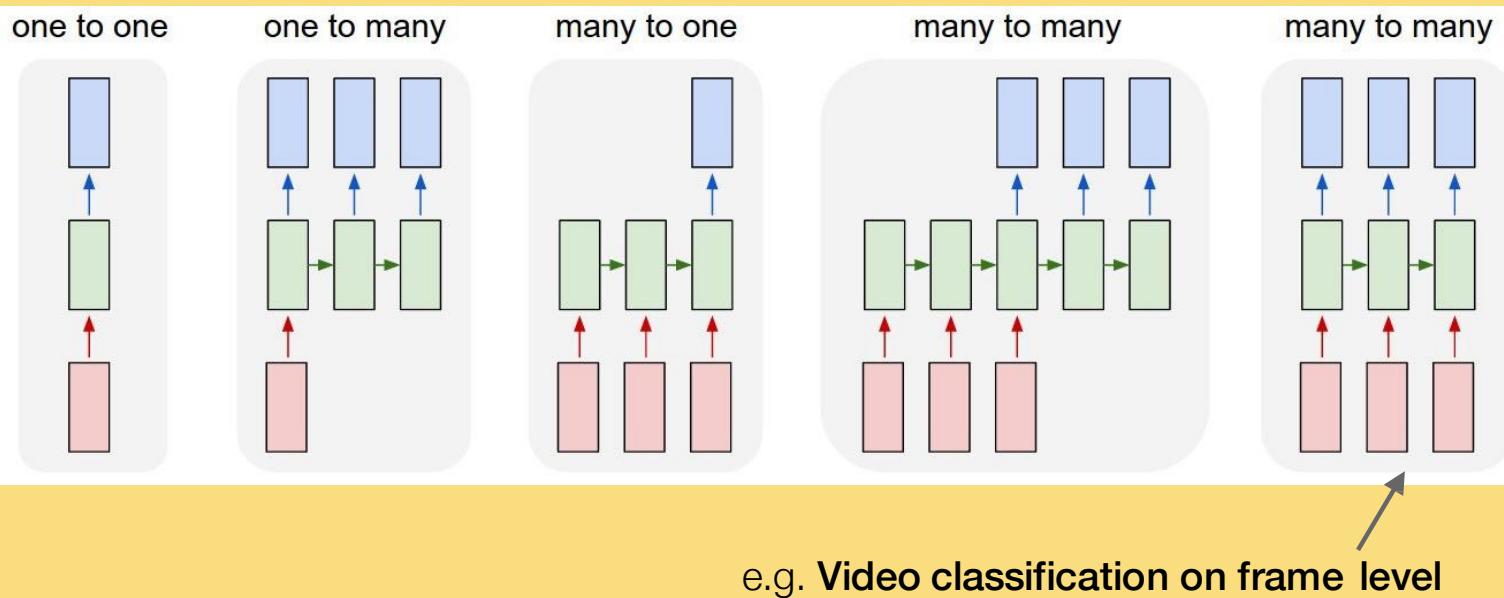


RNN

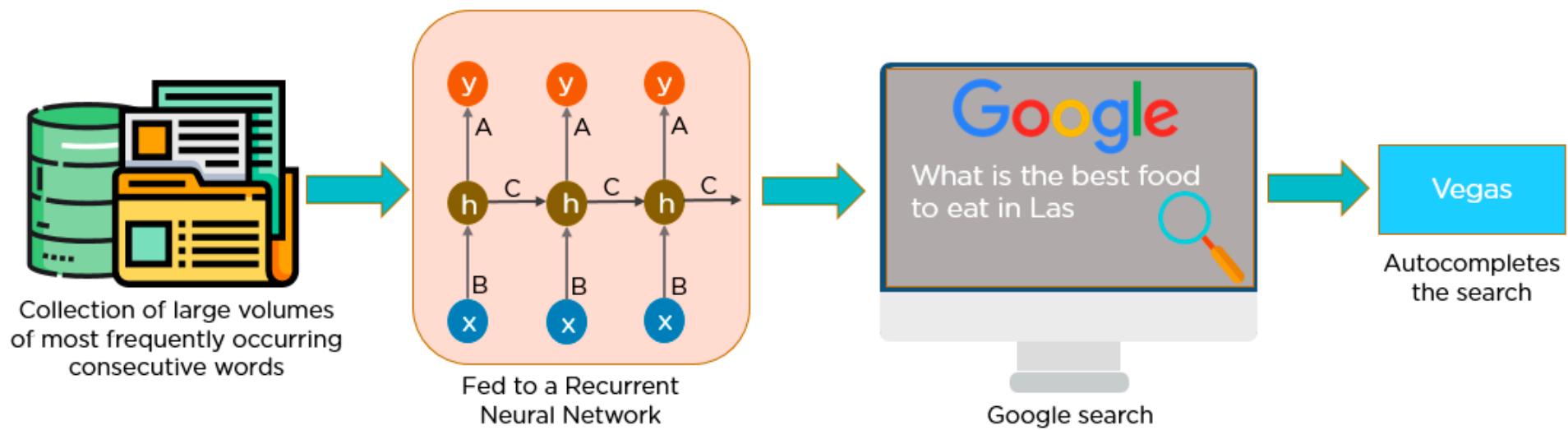


e.g. **Machine Translation**
seq of words -> seq of words

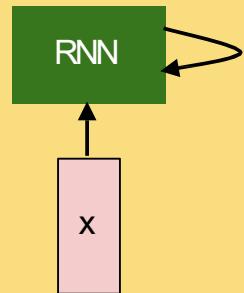
RNN



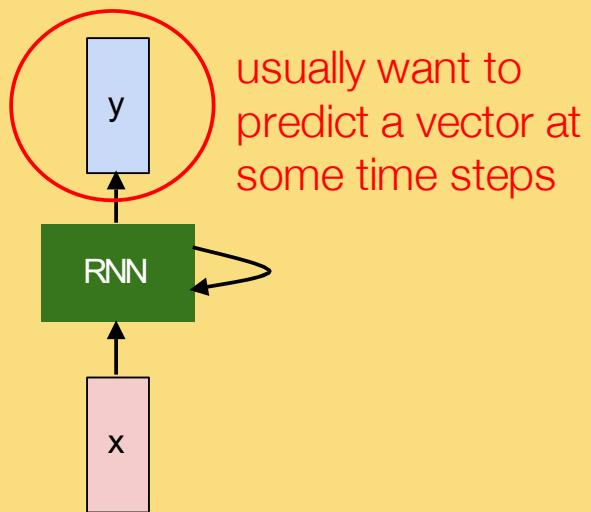
Example

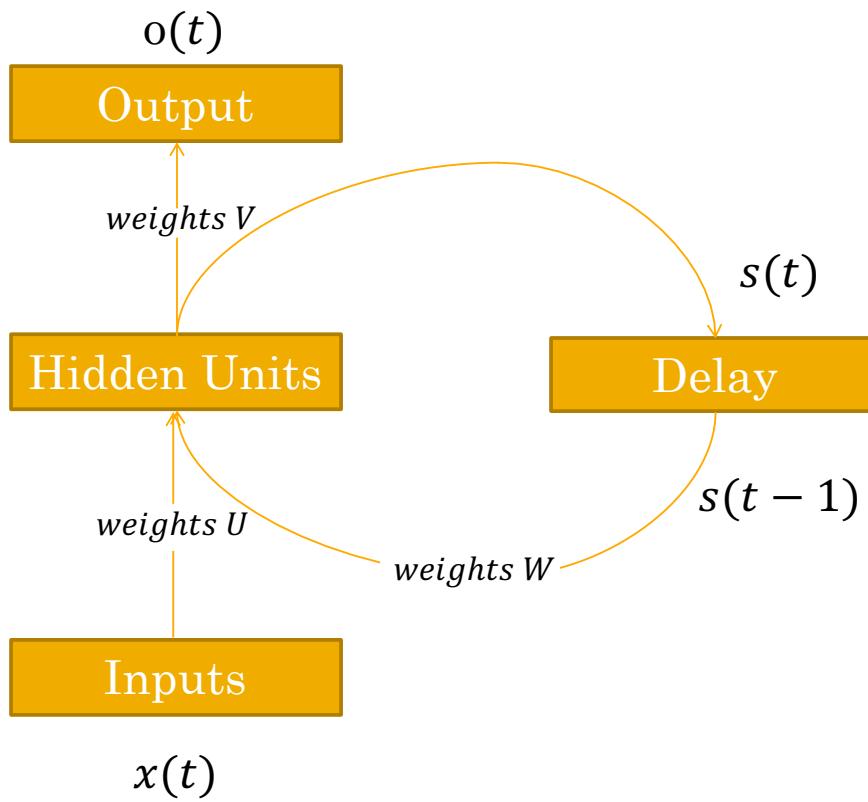


RNN: Basics



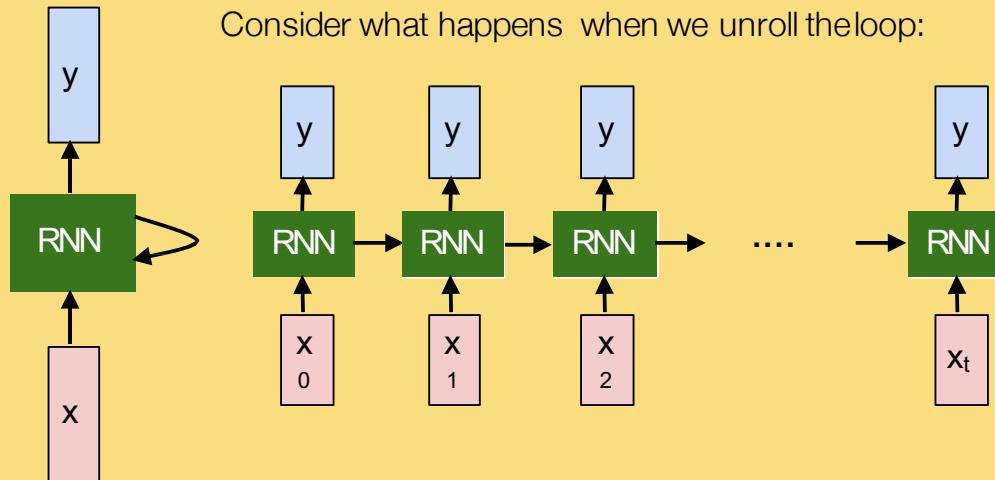
RNN: Basics





RNN: Basics

Consider what happens when we unroll the loop:



- The network input at time t:

$$a_h(t) = Ux(t) + Ws(t - 1)$$

- The activation of the input at time t:

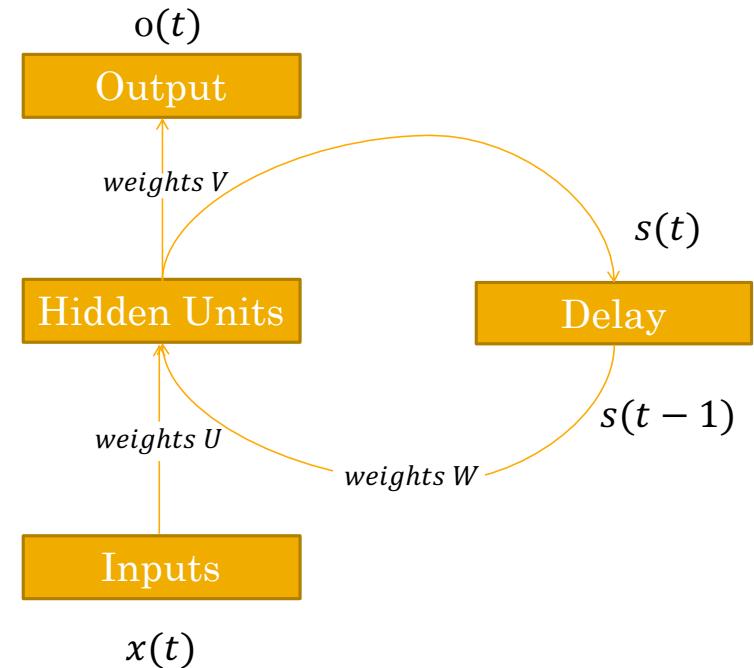
$$s(t) = f_h(a_h(t))$$

- The network input to the output unit at time t:

$$a_o(t) = Vs(t)$$

- The output of the network at time t is:

$$o(t) = f_o(a_o(t))$$



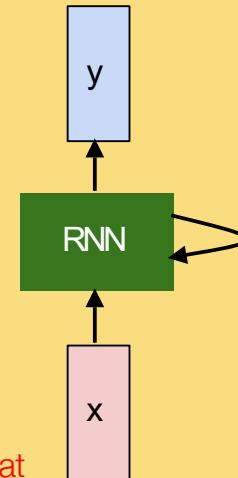
RNN: Basics

We can process a sequence of vectors x by applying a recurrence formula at every time step:

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

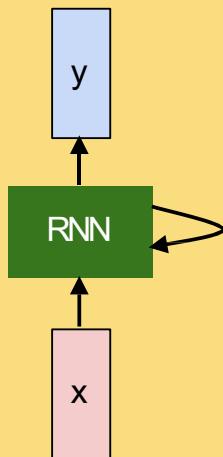
new state old state input vector at
some function with parameters W some time step

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



RNN: Basics

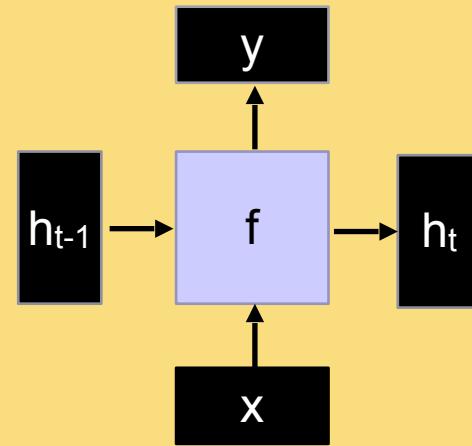
The state consists of a single “hidden” vector \mathbf{h} :



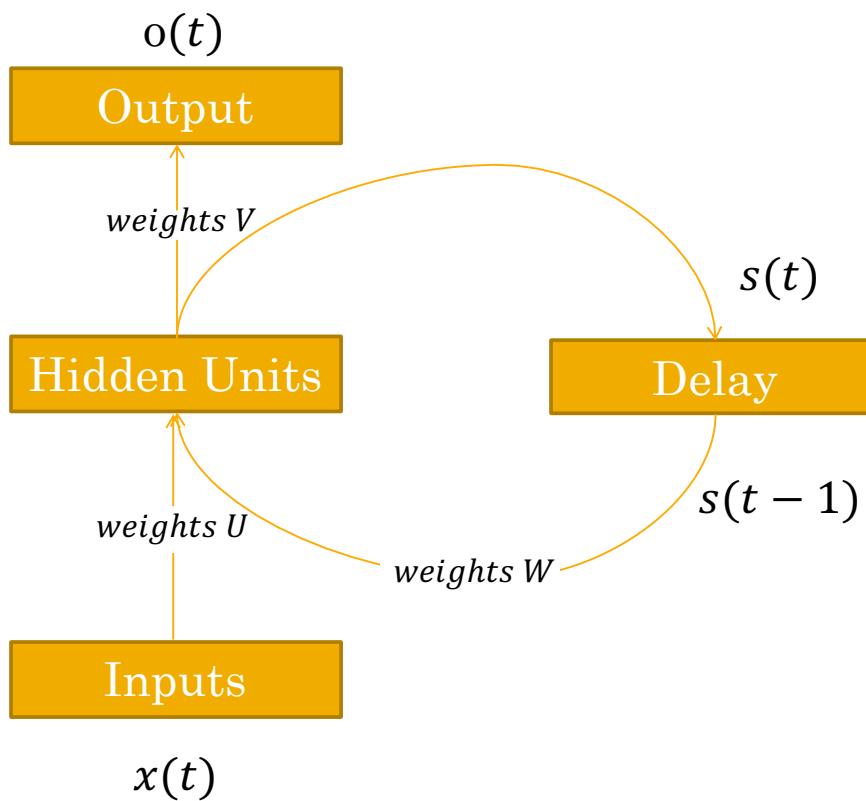
$$h_t = f_W(h_{t-1}, x_t)$$

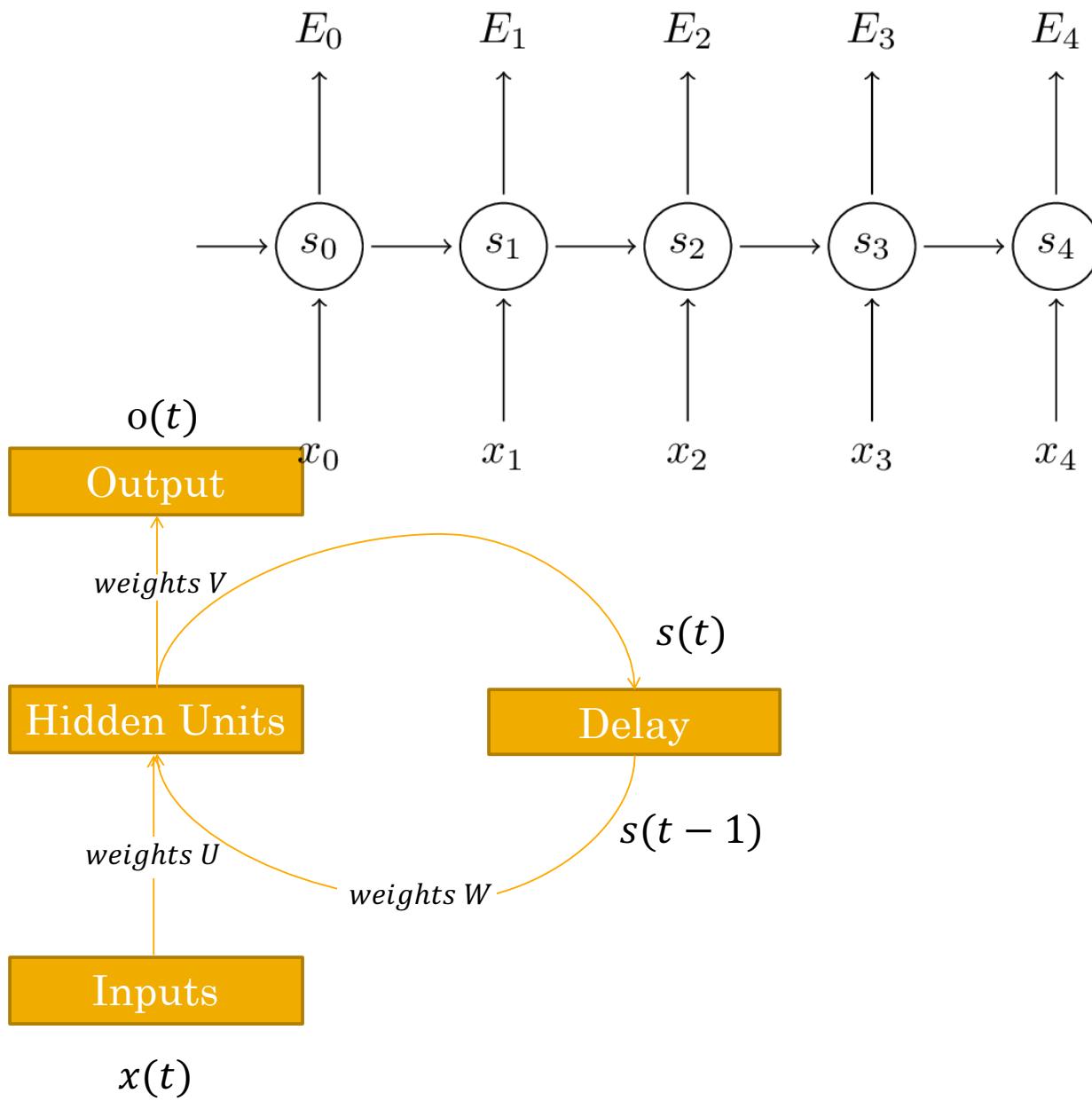
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

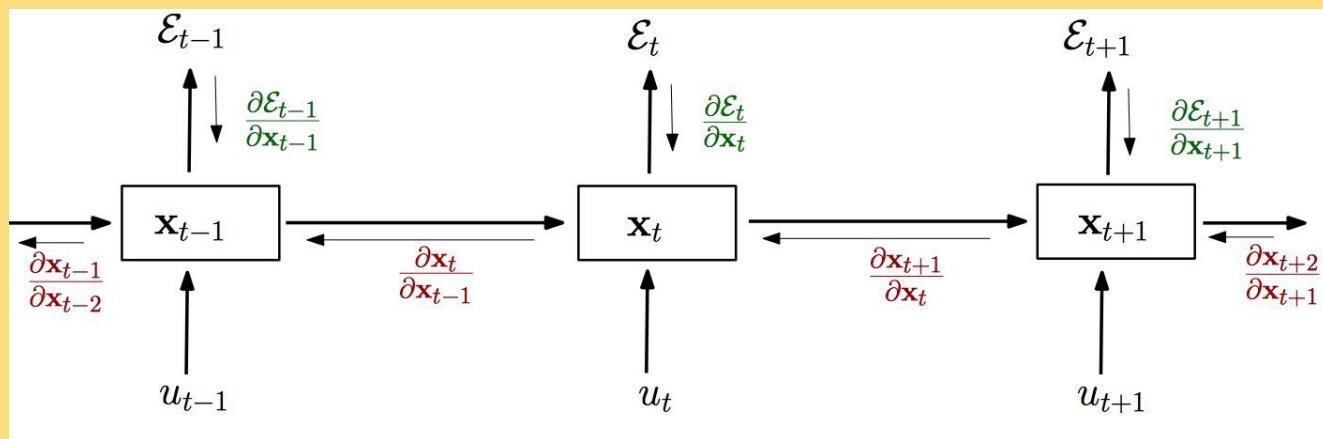
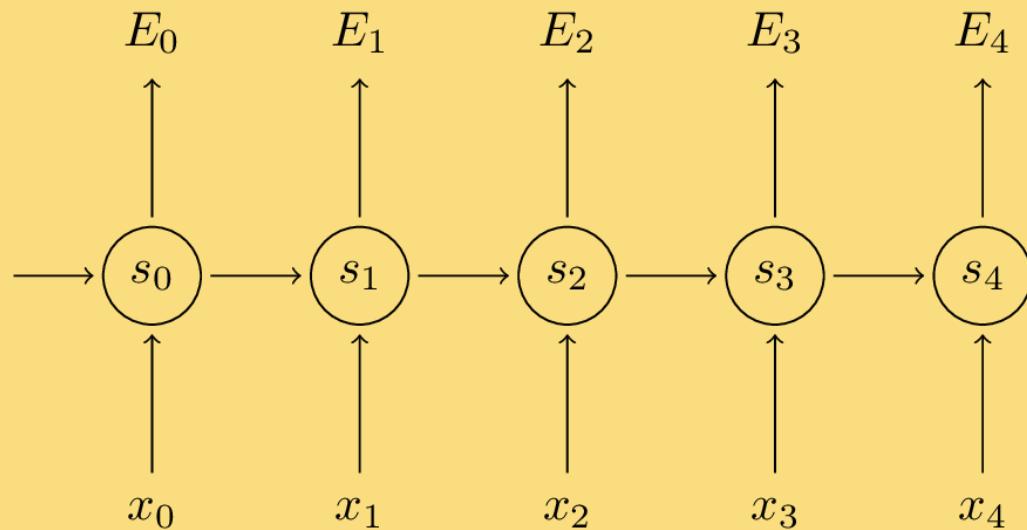
$$y_t = W_{hy}h_t$$



What are the sources of error backpropagation



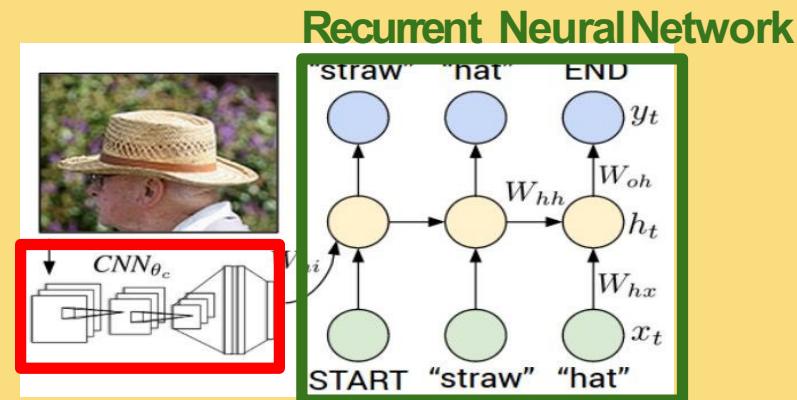




Backpropagation through time (BPTT)

RNN: Example

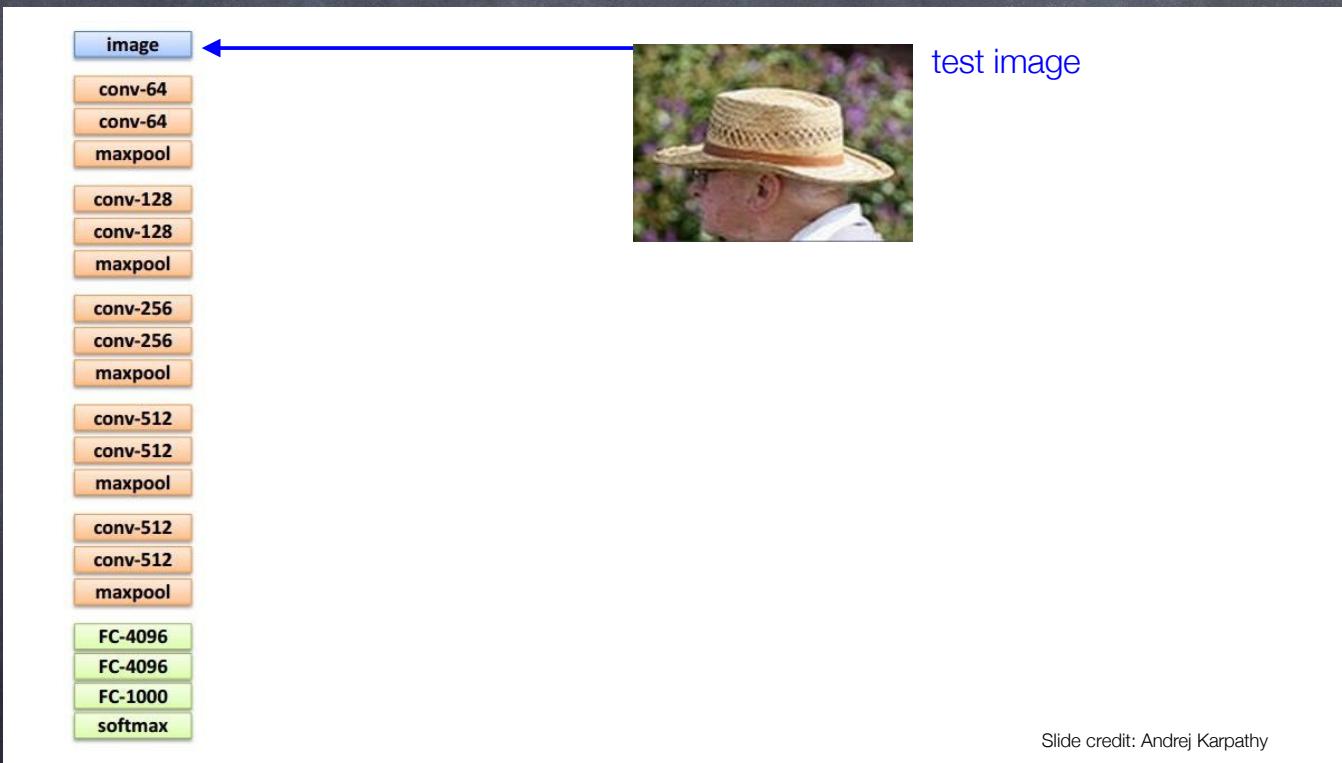
Image Captioning



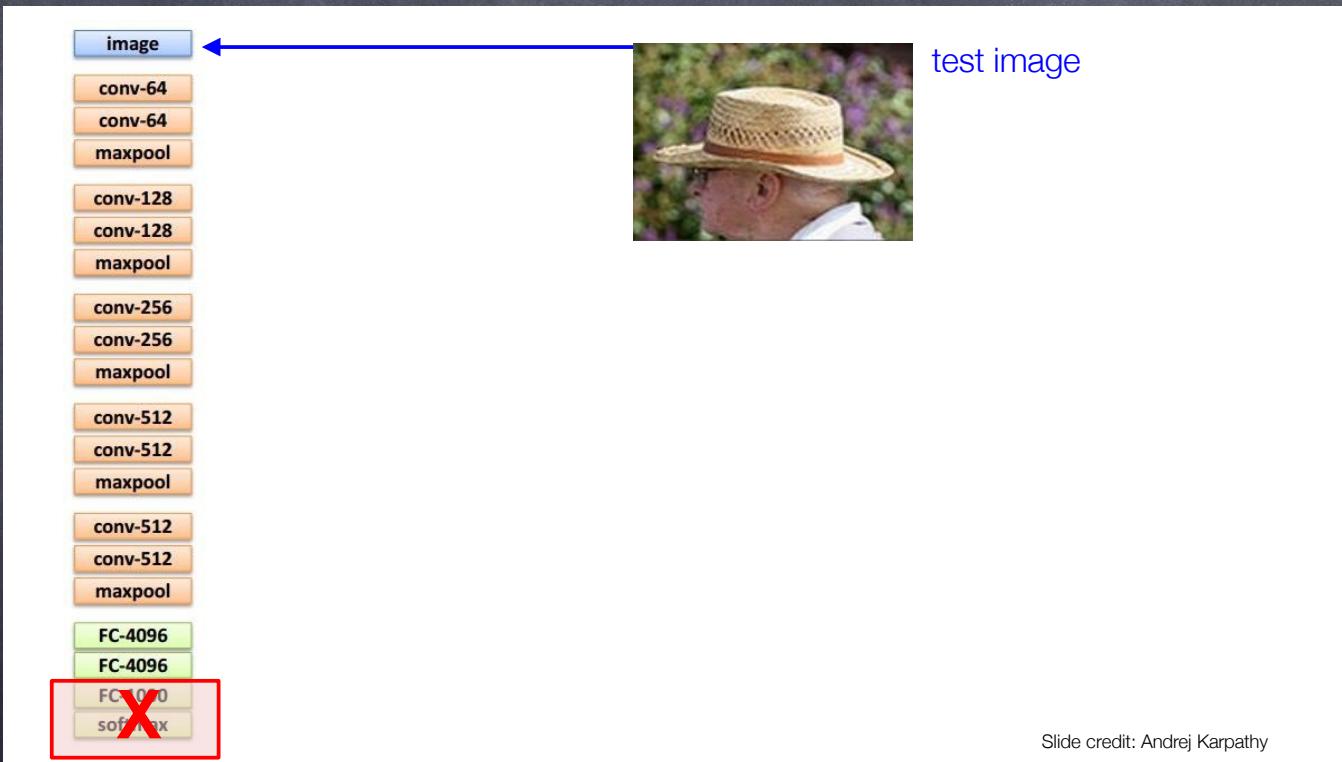


test image

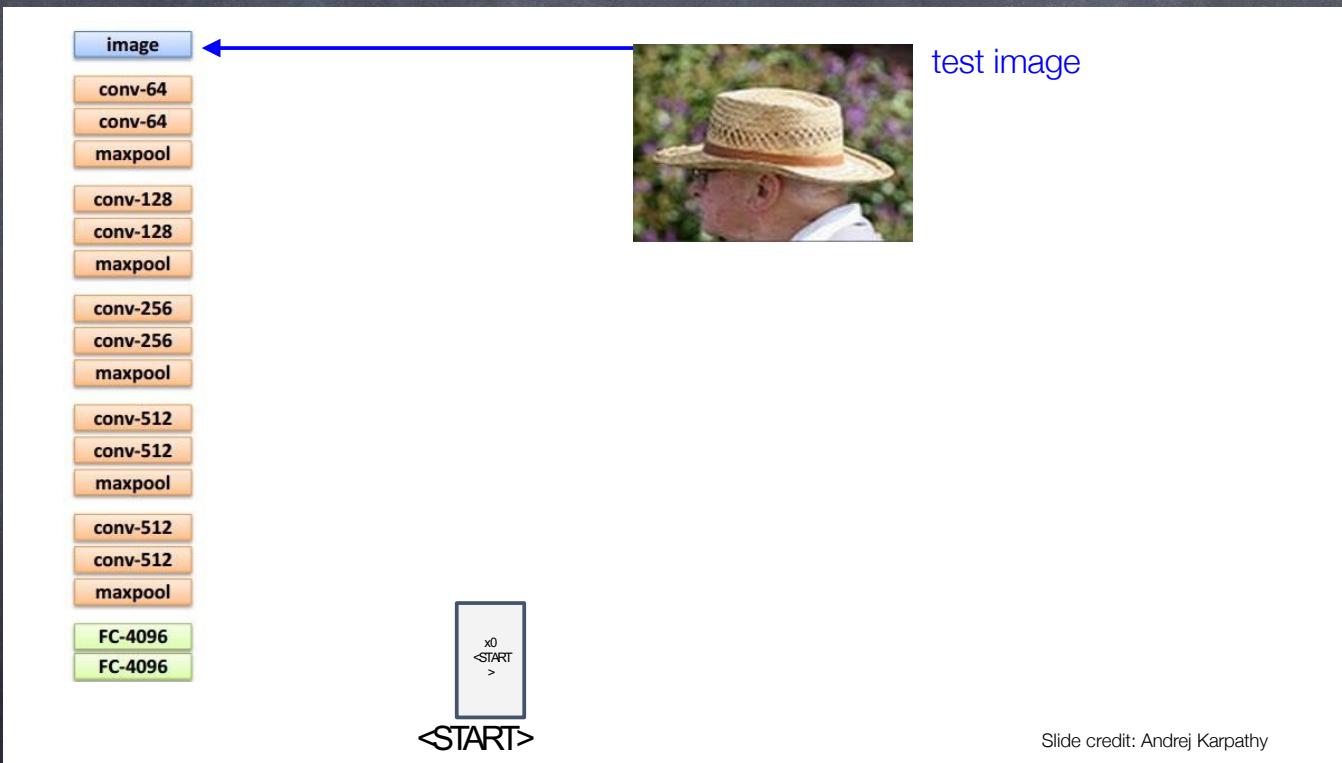
Slide credit: Andrej Karpathy



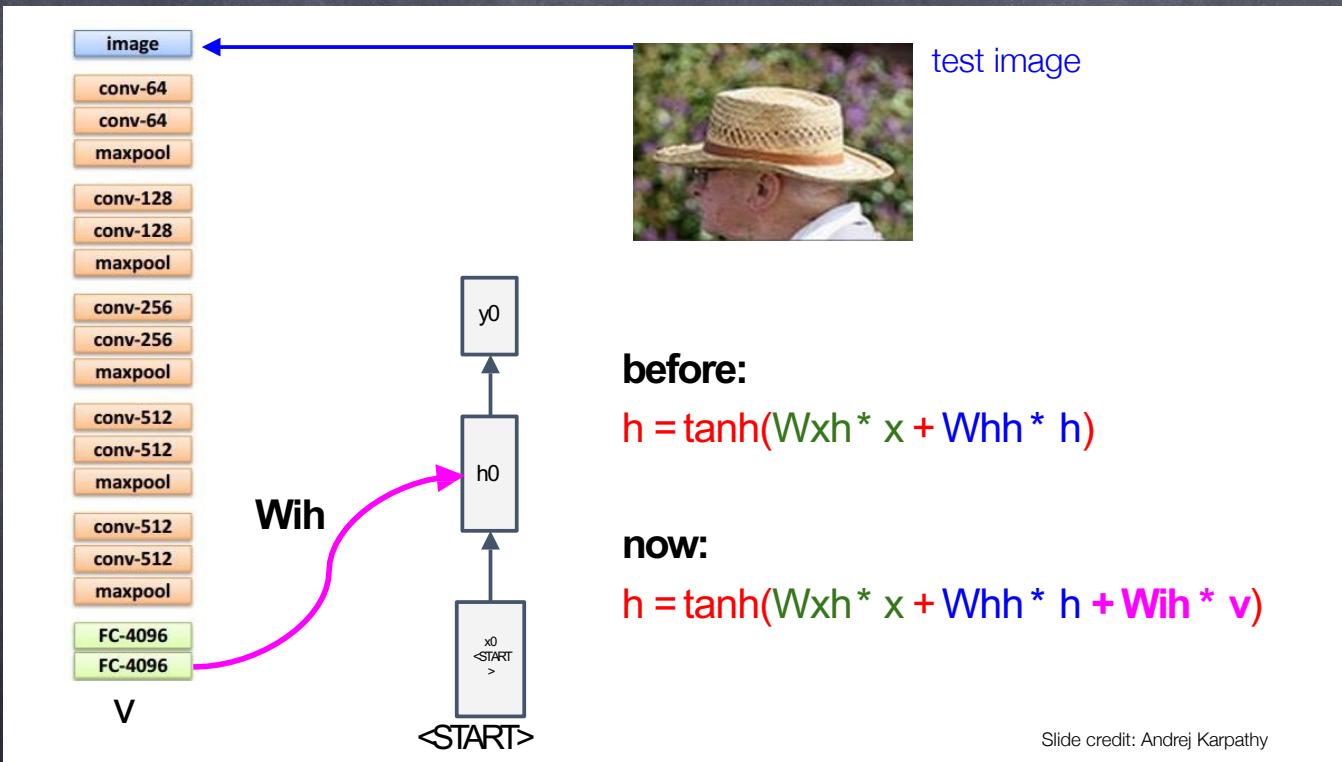
Slide credit: Andrej Karpathy

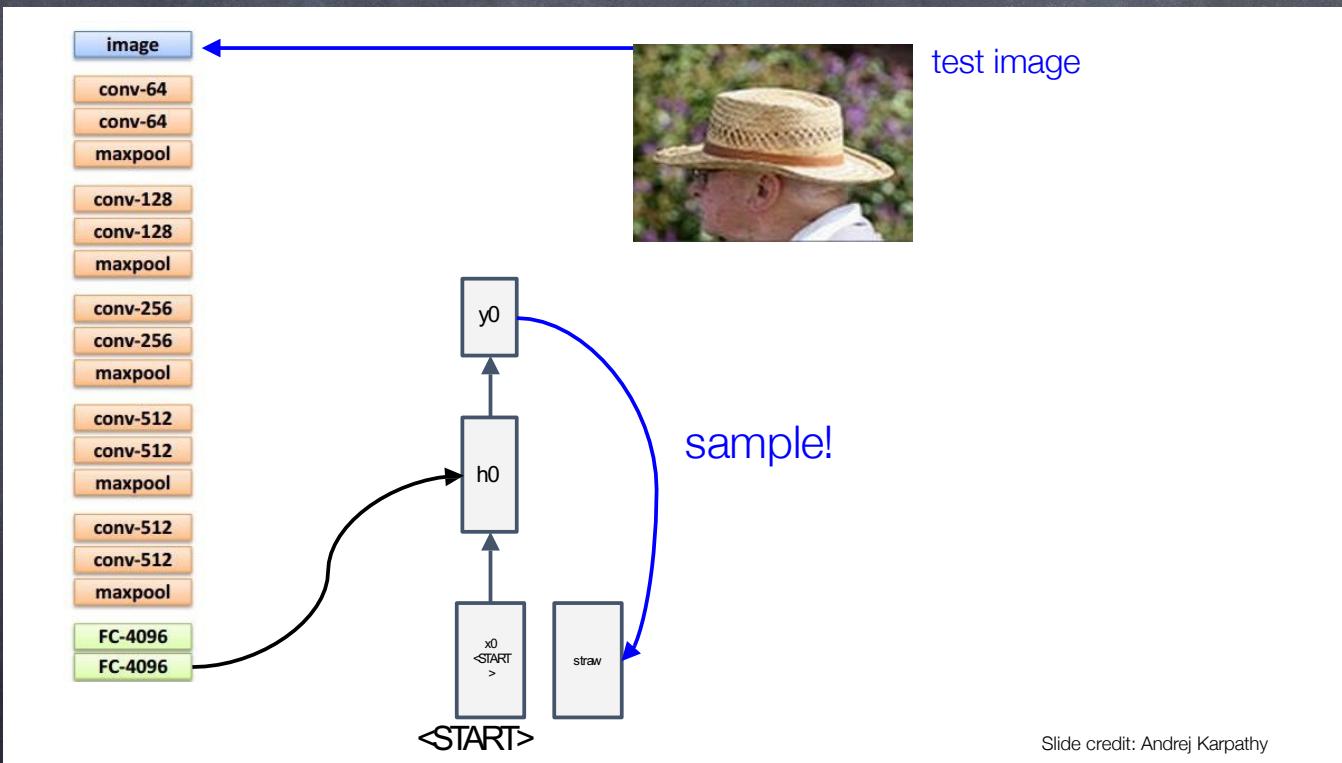


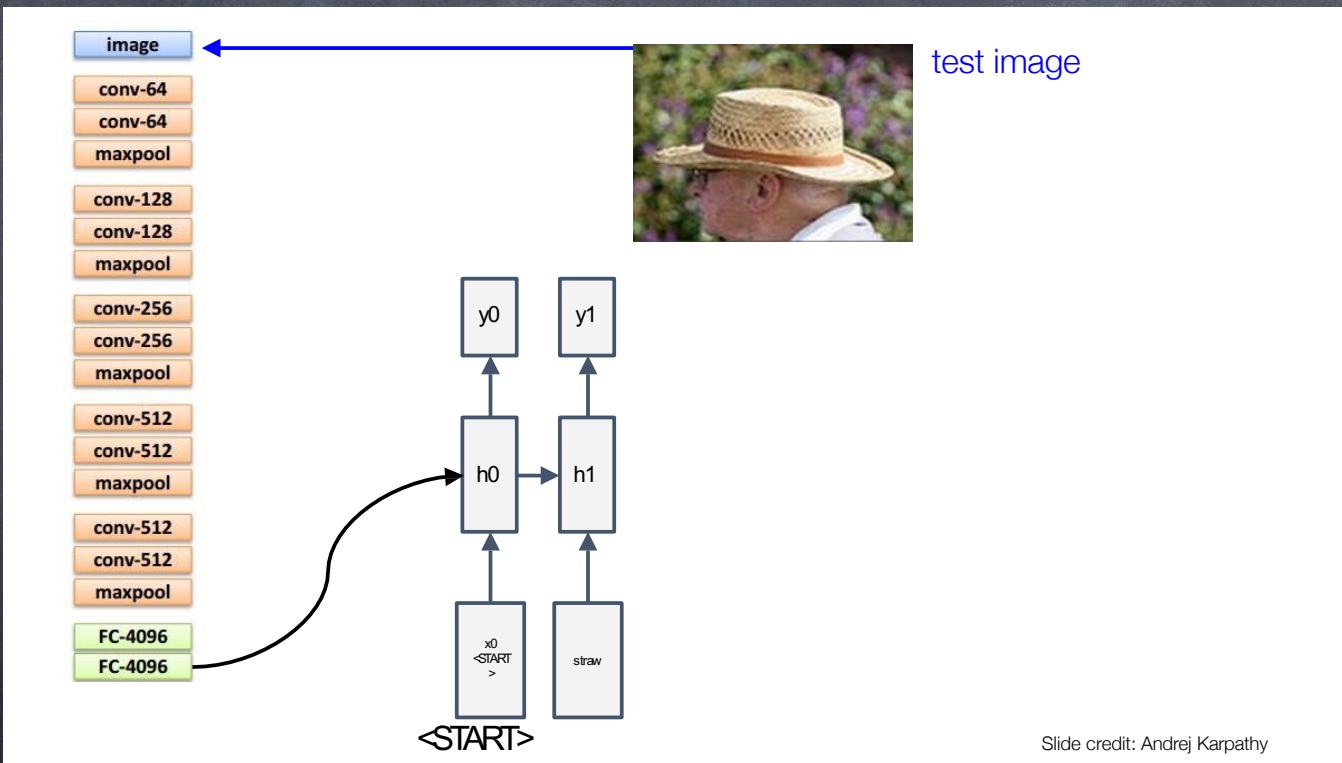
Slide credit: Andrej Karpathy

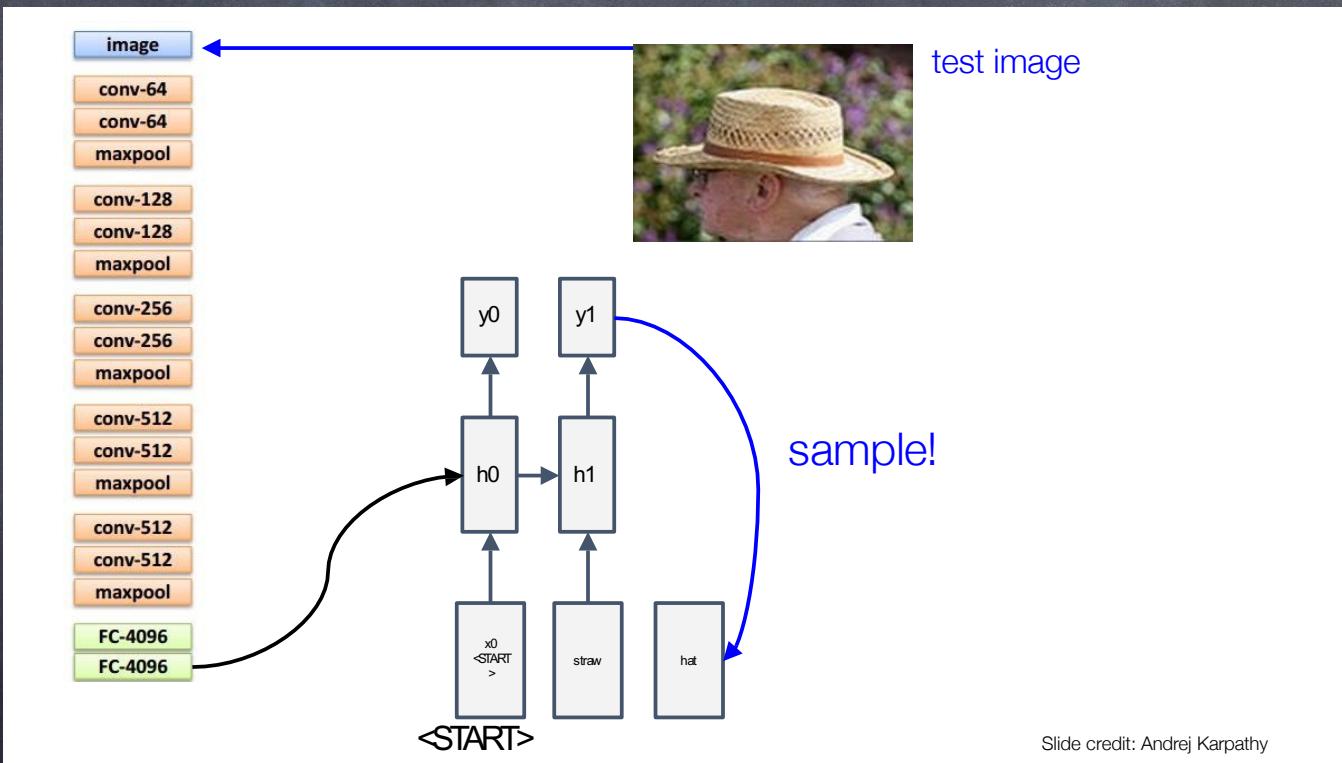


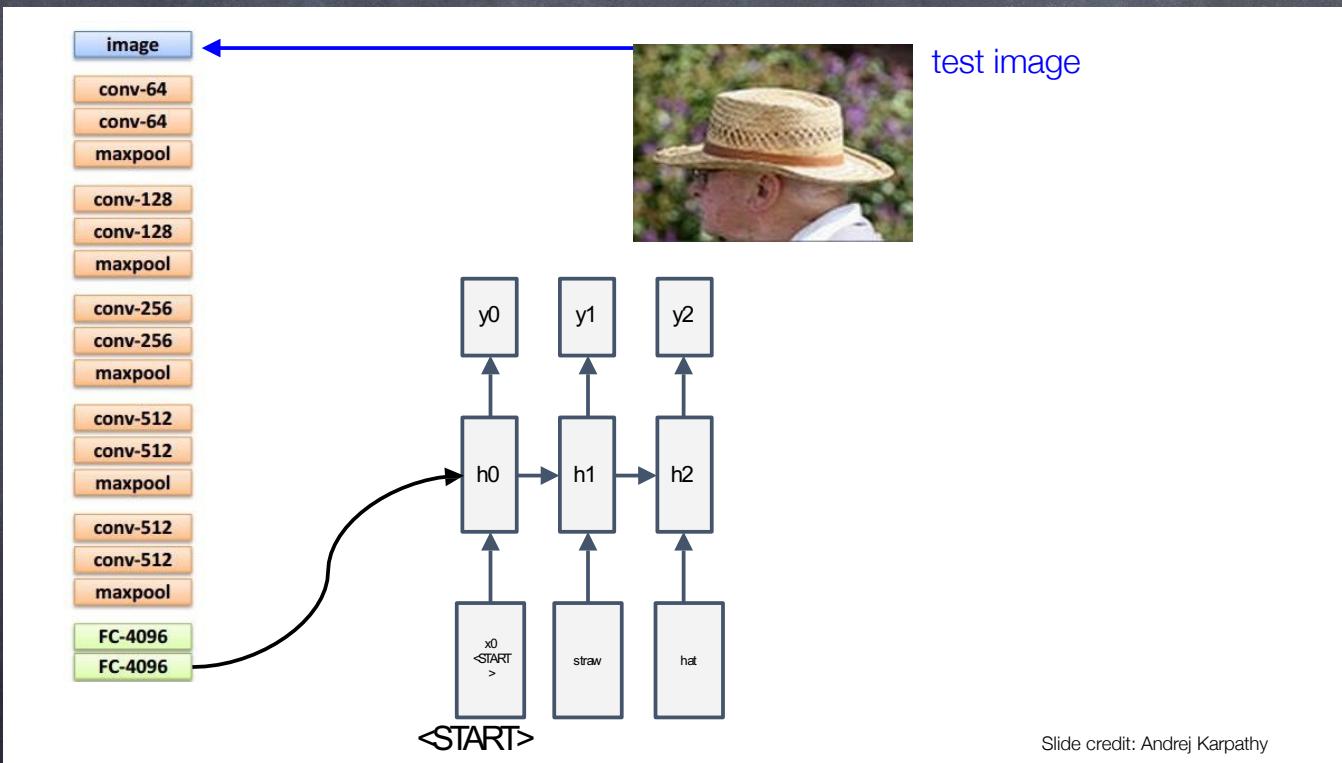
Slide credit: Andrej Karpathy

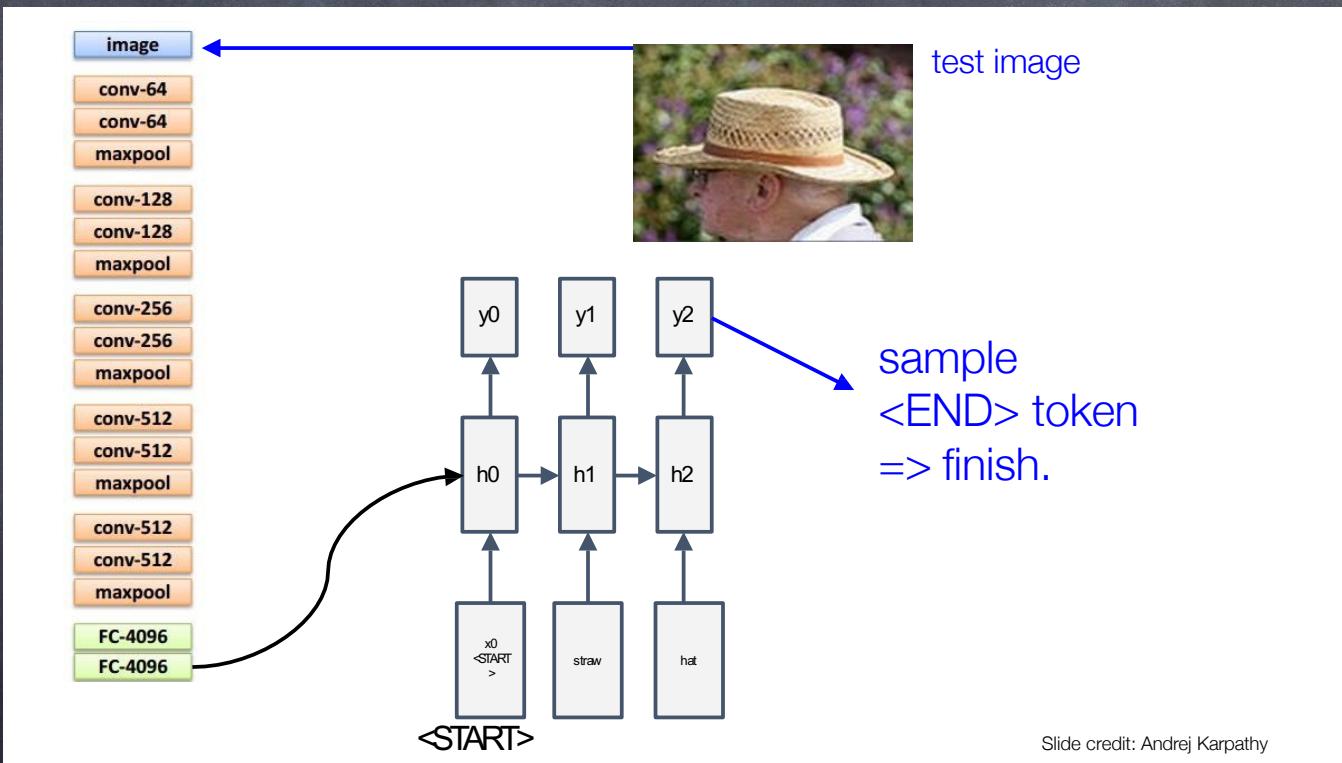












Example

- Ground truth - What are the places to go in Las Vegas?
- Prediction - What are the places to go in Las Vegas

Pros and Cons

- Will predict: "the clouds are in the sky"
- May not predict: "I grew up in India ... I speak fluent Hindi."

Project(s)

- There will be two kinds of projects:
 - 3-0-0 part
 - 0-0-2/3 part (lab)

Project Team Size

- 3 students each project group
- BSB students should create their groups - their projects will be tailored towards BSB problems

Project 1

- For 3-0-0 part
 - Select a problem of your choice (BSB - from your core department)
 - The project should be about total 60 hours of worth of work (20 hours per student)
 - Kaggle competition, implementation of existing papers, comparisons of multiple algorithms on two or more databases

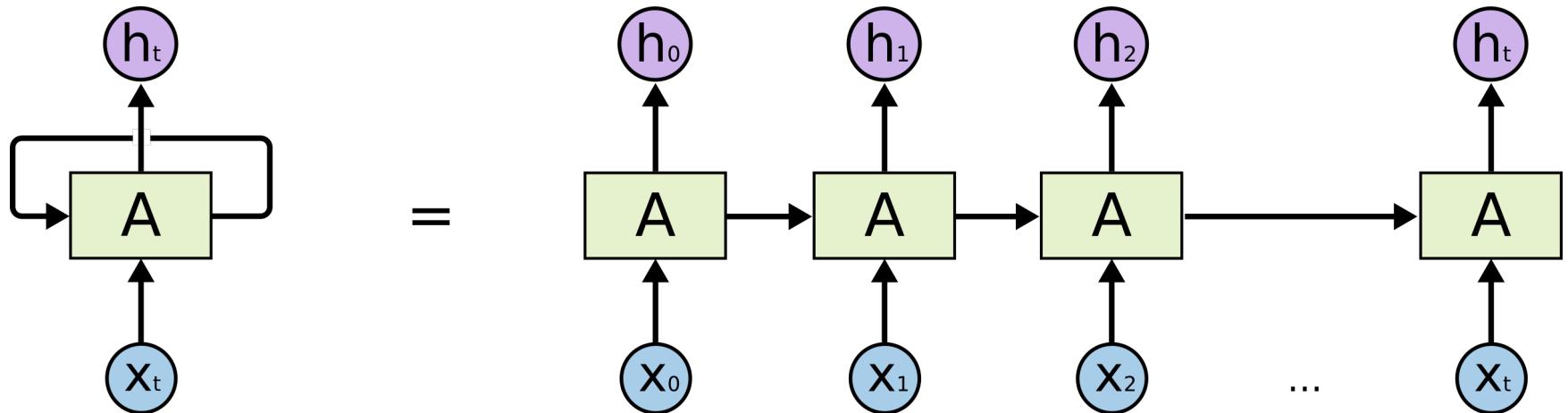
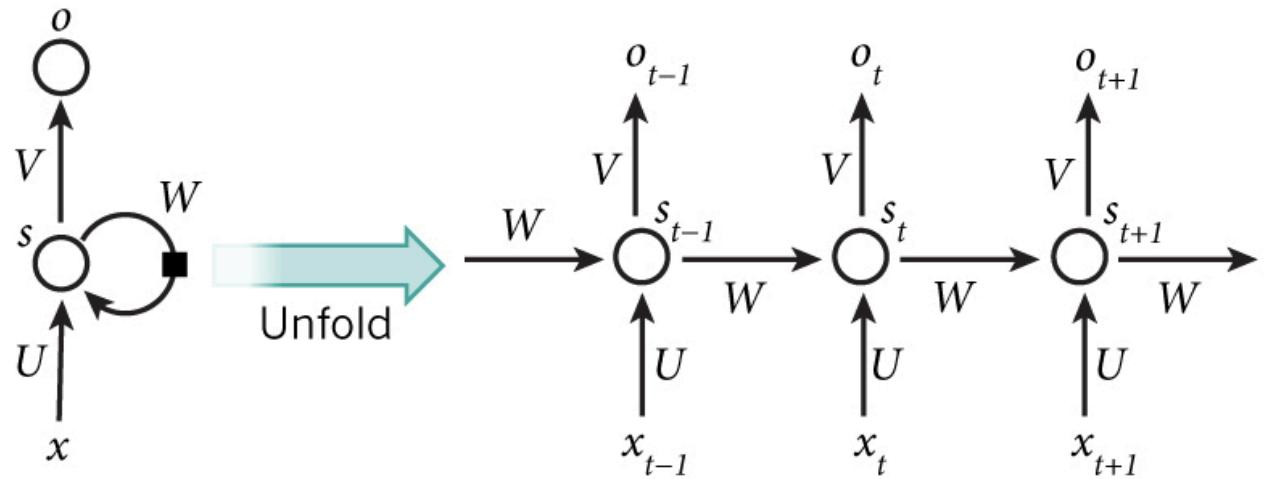
Project 2

- x-0-2/3 part
- The project should be real world application oriented implementation using learnings in the lab
- Distinct from 3-0-0 project
- Same team may do this project

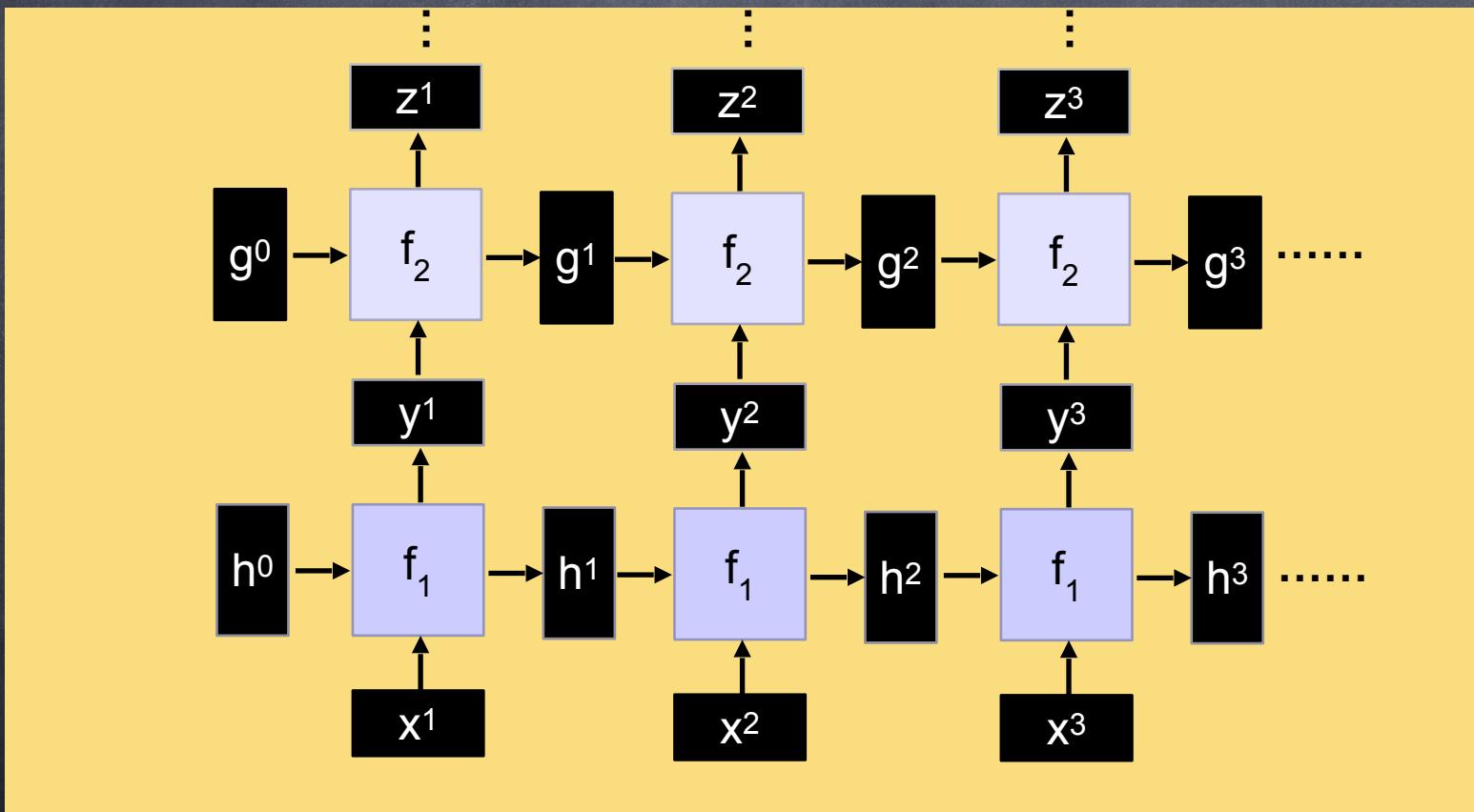
Deadlines

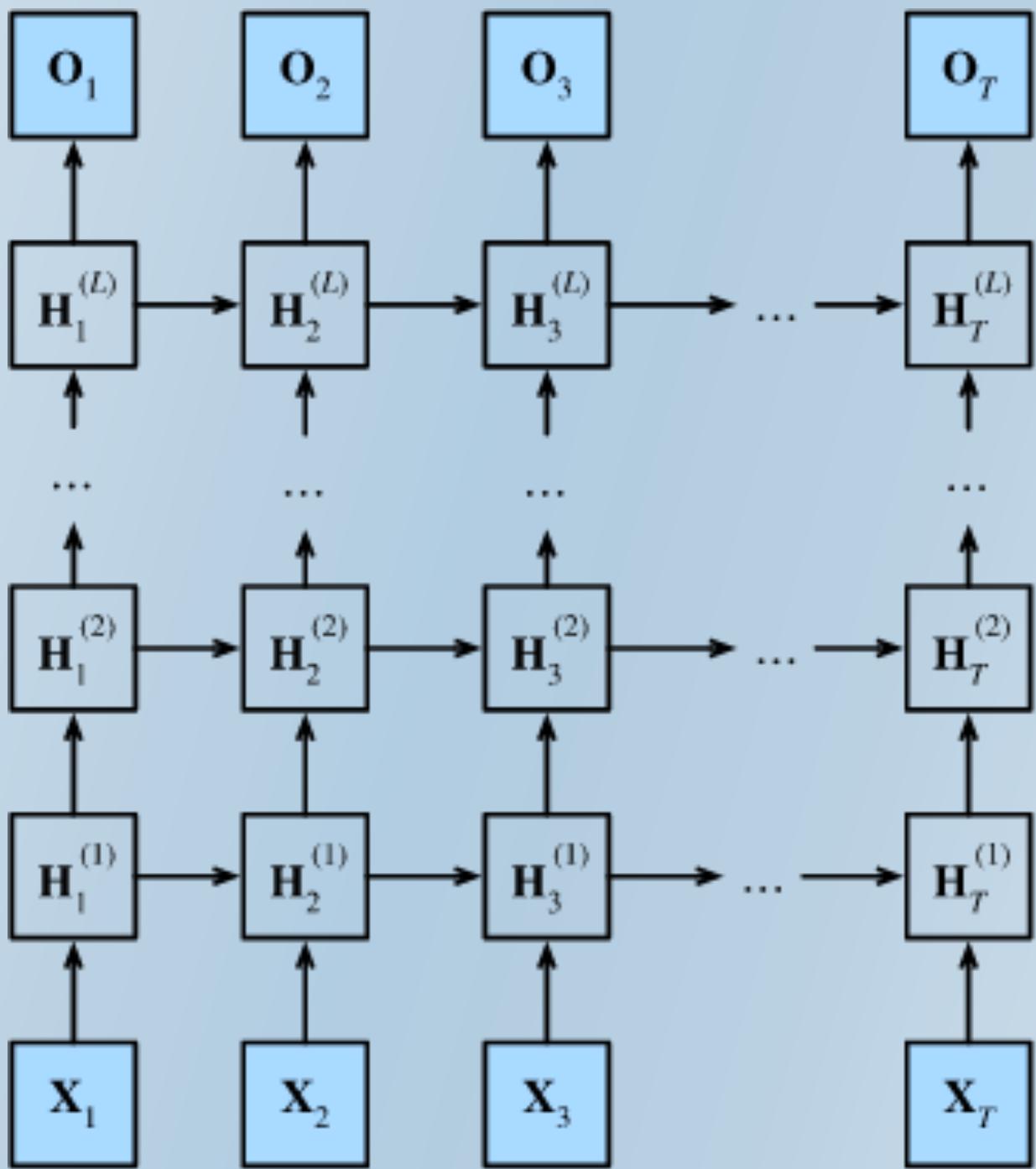
- Team formation: March 6
- Proposals from each team: March 15 (max 2 pages for each project)
 - Problem statement - 1 para
 - Lit review around problem statement - 1-2 para
 - Plan including databases and any other details
- Completing Project 2: April 15
- Completing Project 1: Just before exam - (April end/May Start)

Recurrent Neural Network



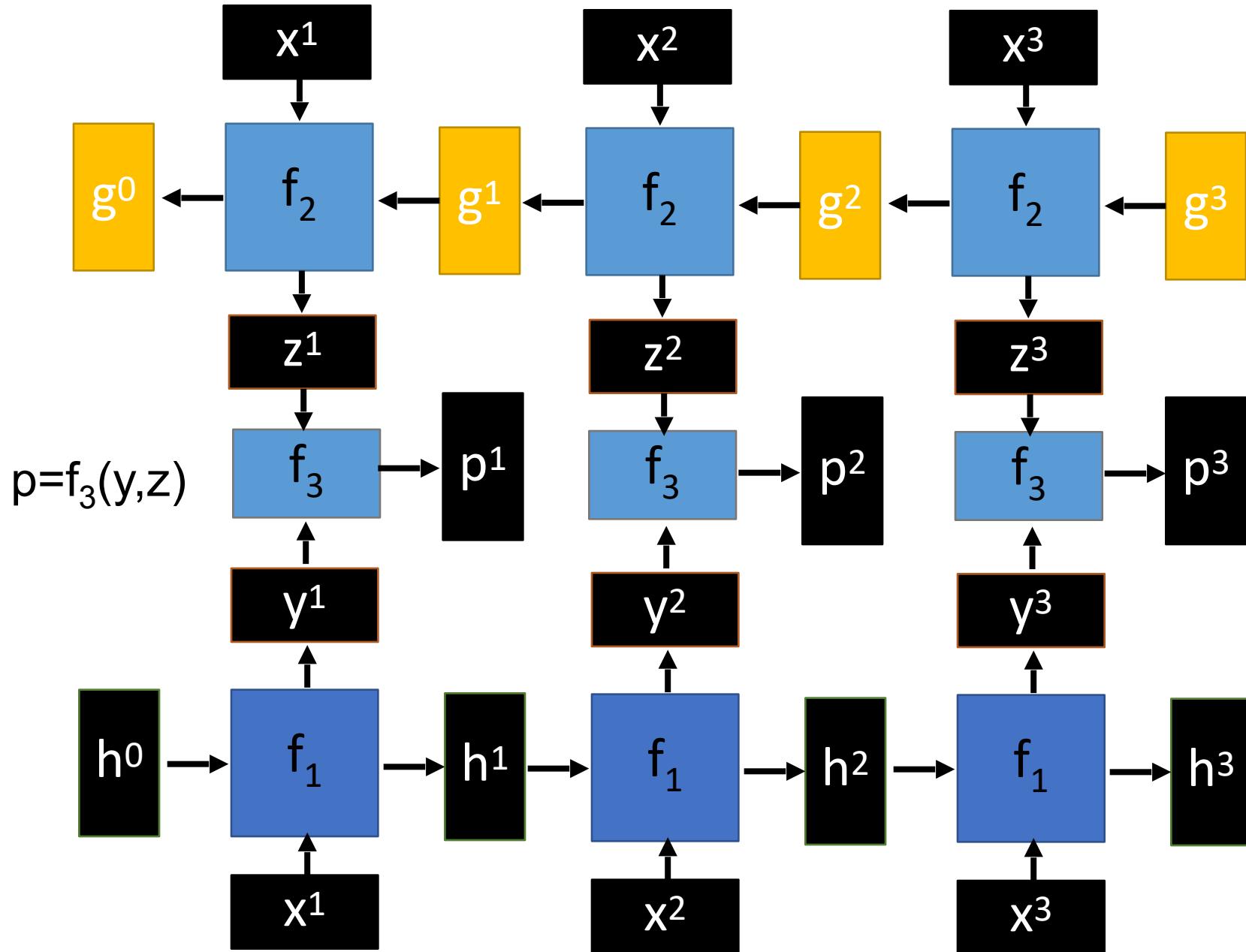
Deep RNN





Bidirectional RNN

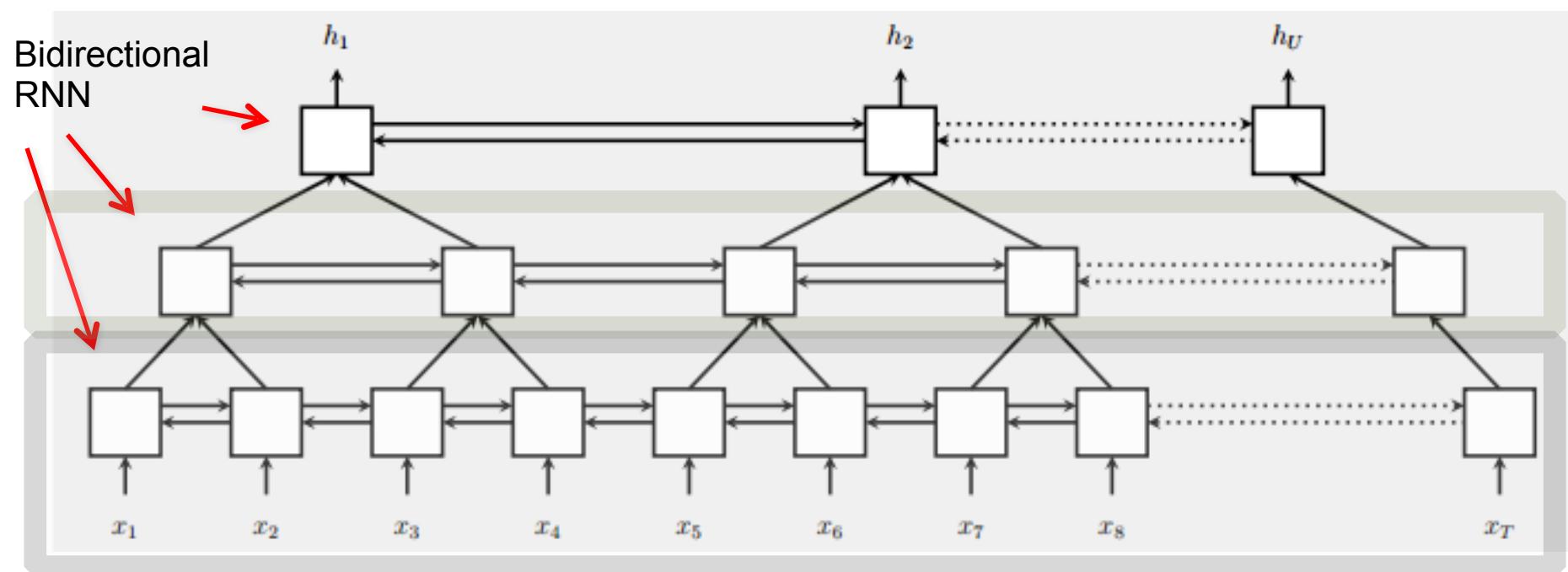
$$y, h = f_1(x, h) \quad z, g = f_2(g, x)$$



Pyramid RNN

Significantly speed up training

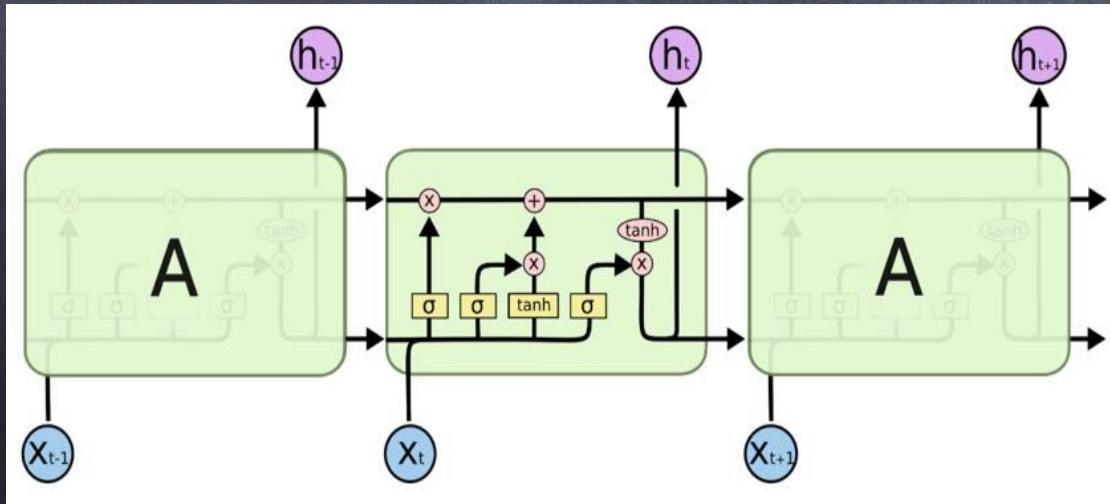
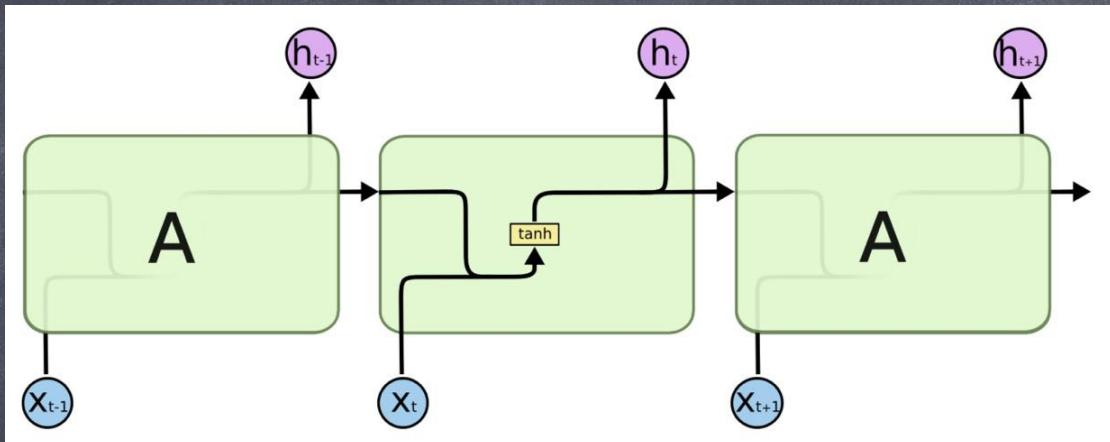
- Reducing the number of time steps



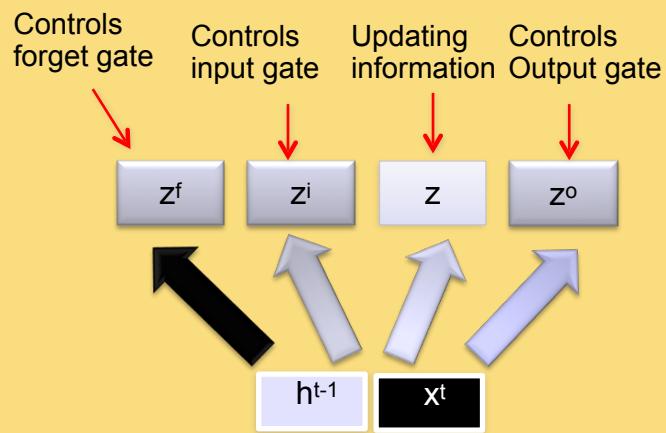
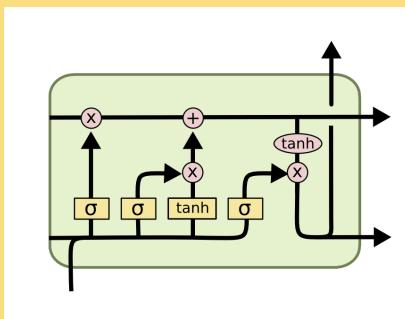
W. Chan, N. Jaitly, Q. Le and O. Vinyals, "Listen, attend and spell: A neural network for large vocabulary conversational speech recognition," ICASSP, 2016

- Long-term temporal dependencies

LSTM: Memory Extension of RNN

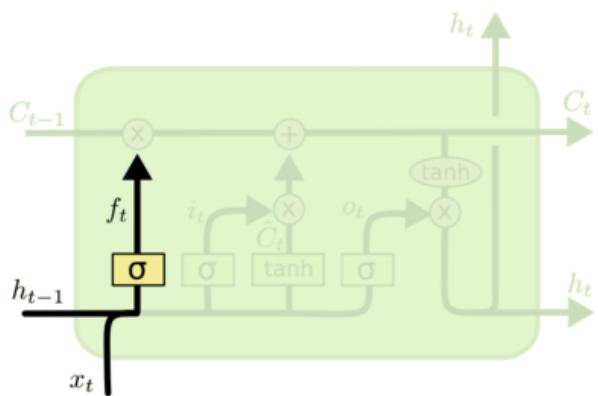


LSTM



LSTM

Walkthrough

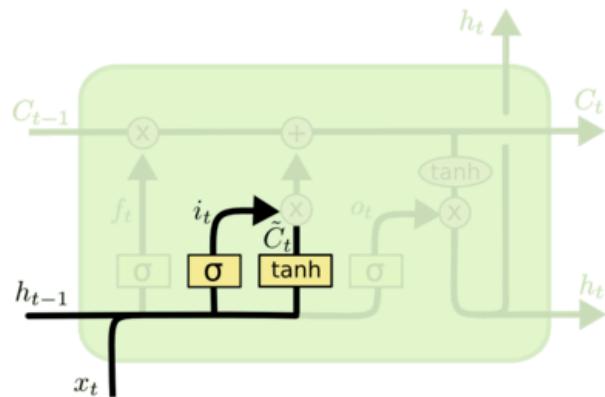


What part of memory
to “forget” – zero
means forget this bit

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

LSTM

Walkthrough



What bits to insert
into the next states

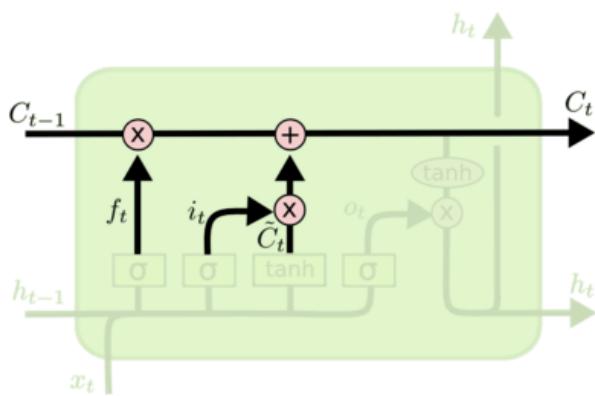
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

What content to store
into the next state

LSTM

Walkthrough

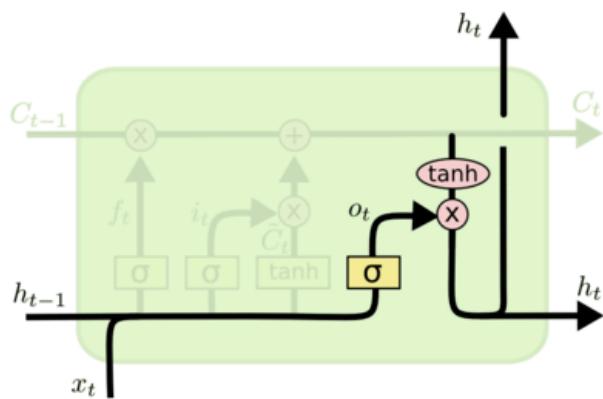


Next memory cell content – mixture of not-forgotten part of previous cell and insertion

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

LSTM

Walkthrough



What part of cell to output

$$o_t = \sigma (W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh (C_t)$$

tanh maps bits to [-1,+1] range



Neural Network
Layer



Pointwise
Operation



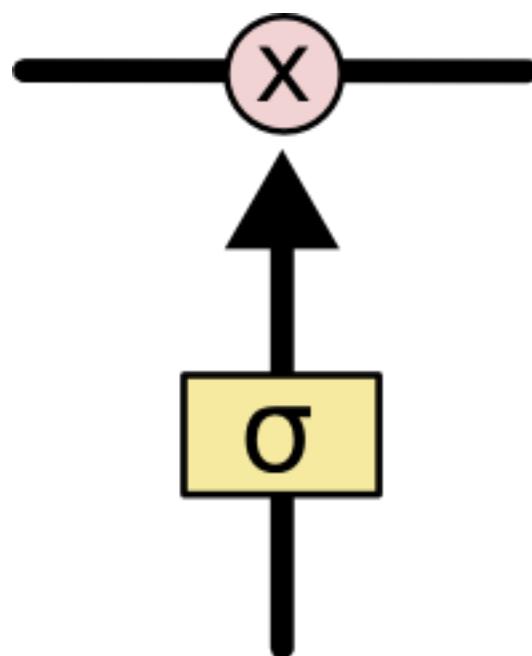
Vector
Transfer



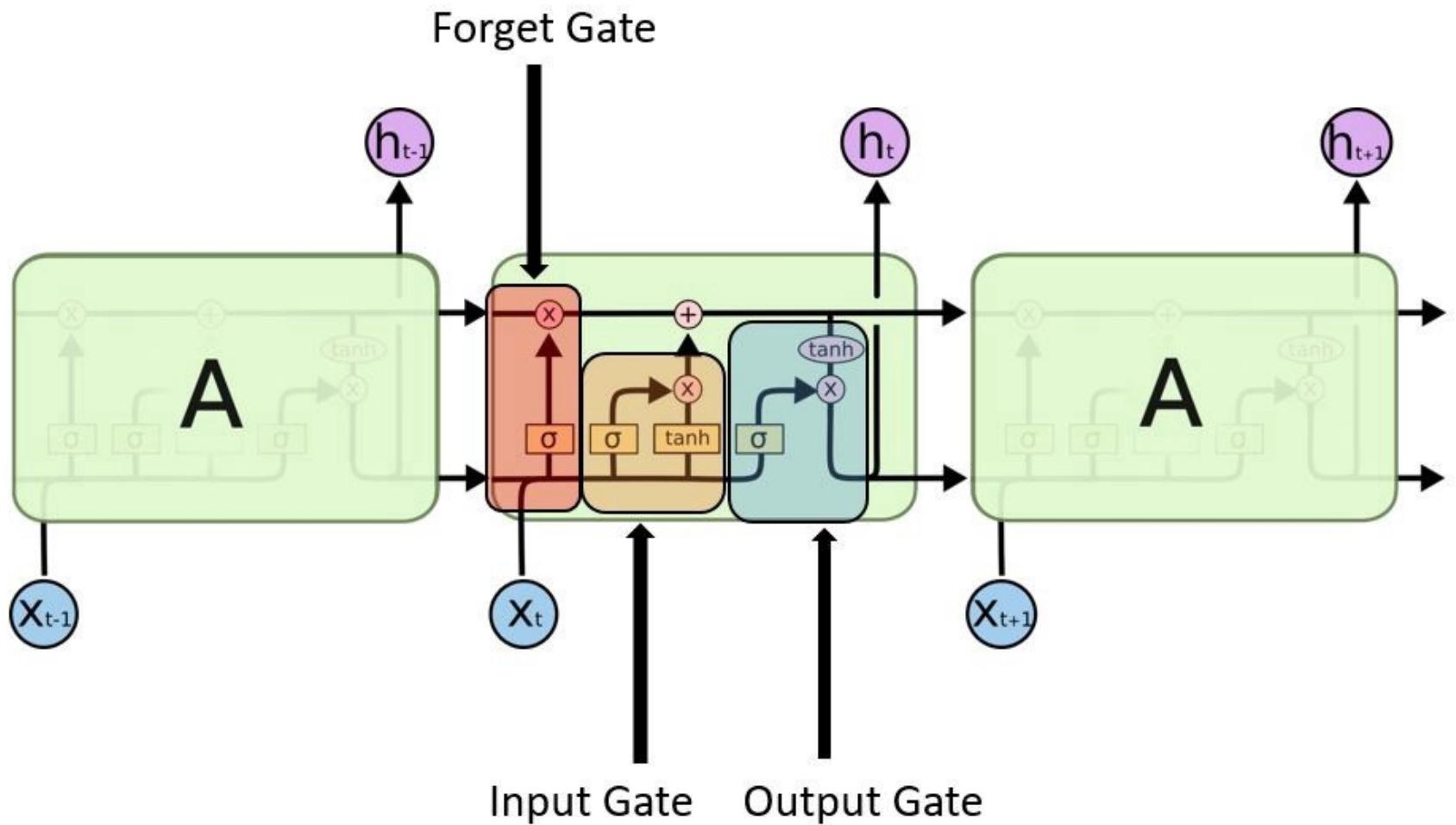
Concatenate



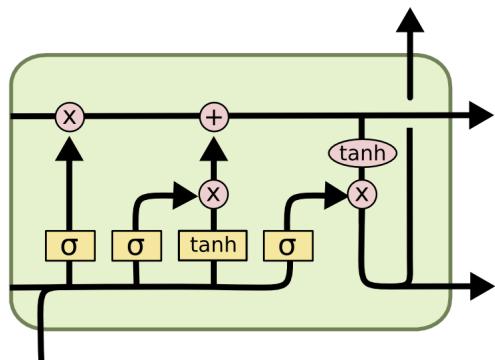
Copy



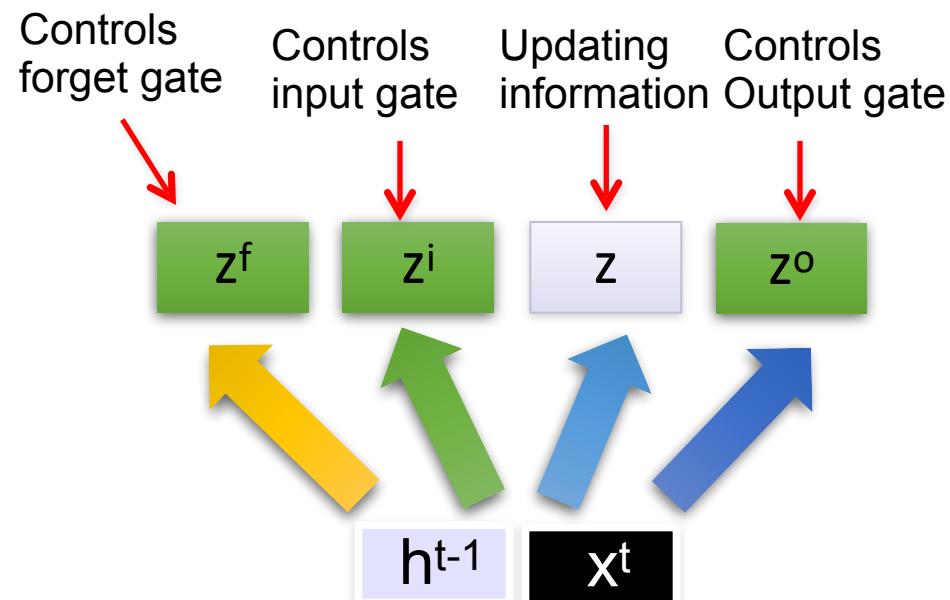
The sigmoid layer outputs numbers between 0-1 determine how much each component should be let through. Pink X gate is point-wise multiplication.



These 4 matrix computation should be done concurrently.



c^{t-1}



$$z = \tanh(W h^{t-1} + x^t)$$

$$z^i = \sigma(W_i h^{t-1} + x^t)$$

$$z^f = \sigma(W_f h^{t-1} + x^t)$$

$$z^o = \sigma(W_o h^{t-1} + x^t)$$

Information flow of LSTM