# CX4240 Spring 2026
# Recurrent Neural Networks
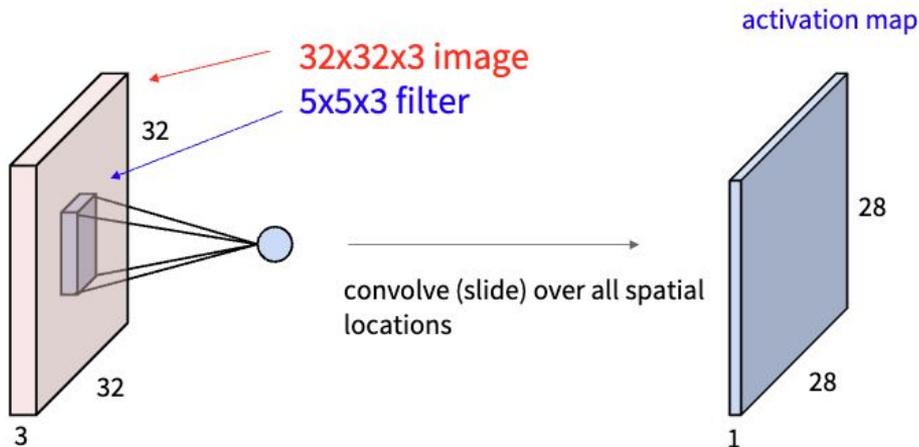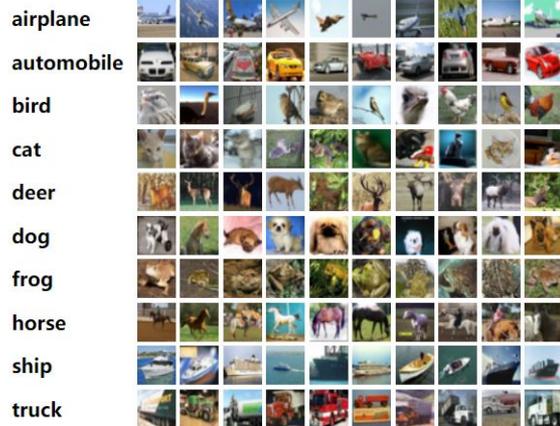
Bo Dai
School of CSE, Georgia Tech
bodai@cc.gatech.edu

Based on Fei-Fei Li & Ehsan Adeli's Slides

# Multiclass Logistic Regression Algorithms

Training Data
$$\{x^i, y^i\}_{i=1}^m$$

Learning Algorithm

Classifier
$$f : X \rightarrow Y$$

Multiclass Classification $\quad Y \in \{0, 1, \ldots, k\}$

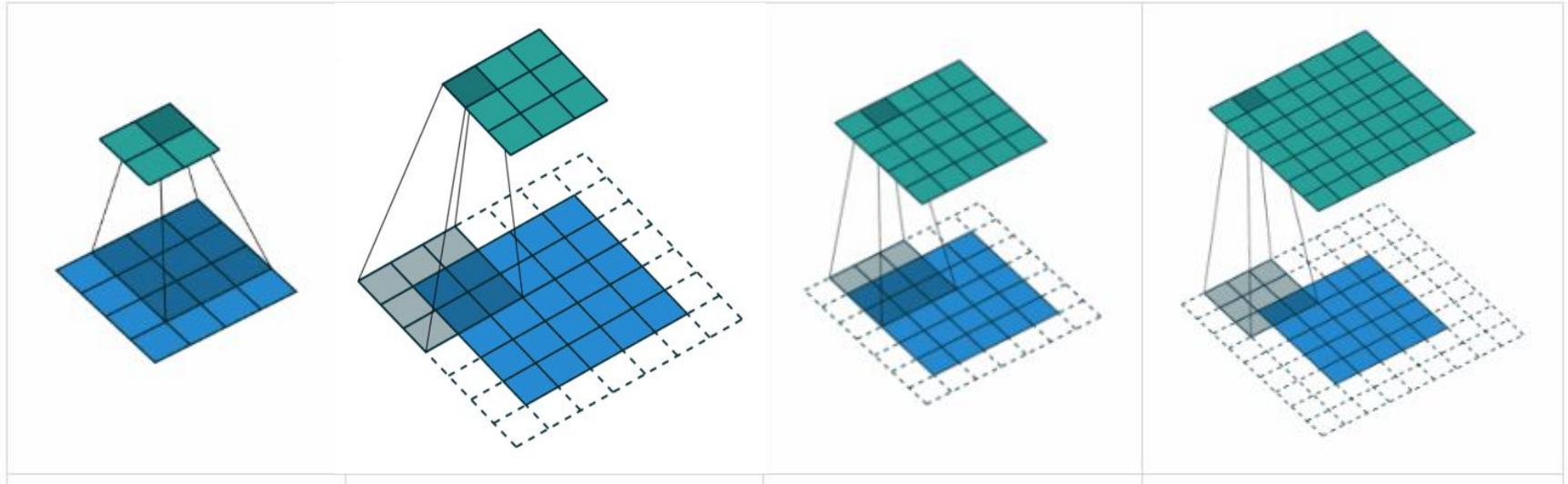## Multiclass Logistic Regression Pipeline

1. Build probabilistic models:
   Categorical Distribution + Conv NN
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

# Revisit Convolution Neural Networks



airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

activation map

28

28

1

# Convolution for 2D Images



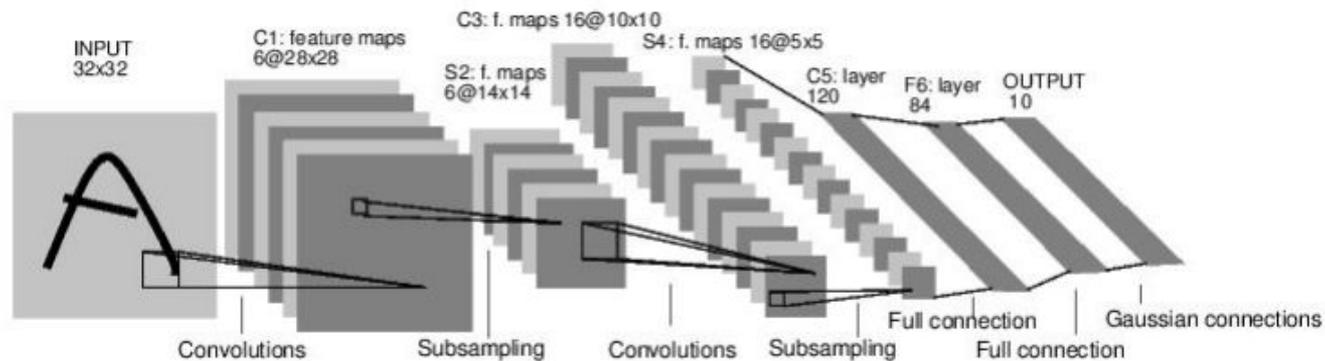padding = 0, stride = 1    padding = 1, stride = 2    padding = 1, stride = 1    padding = 2, stride = 1

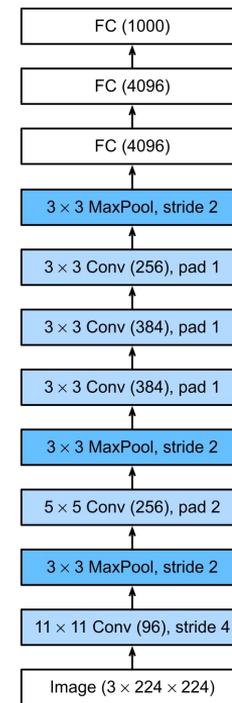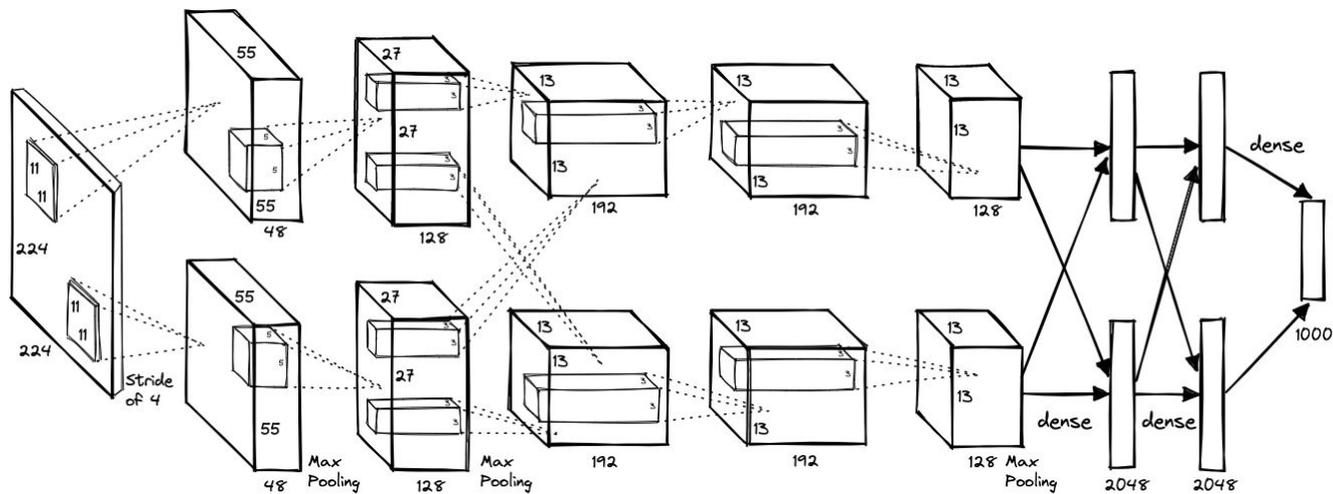$$W_{out} = \frac{W - F + 2P}{S} + 1$$
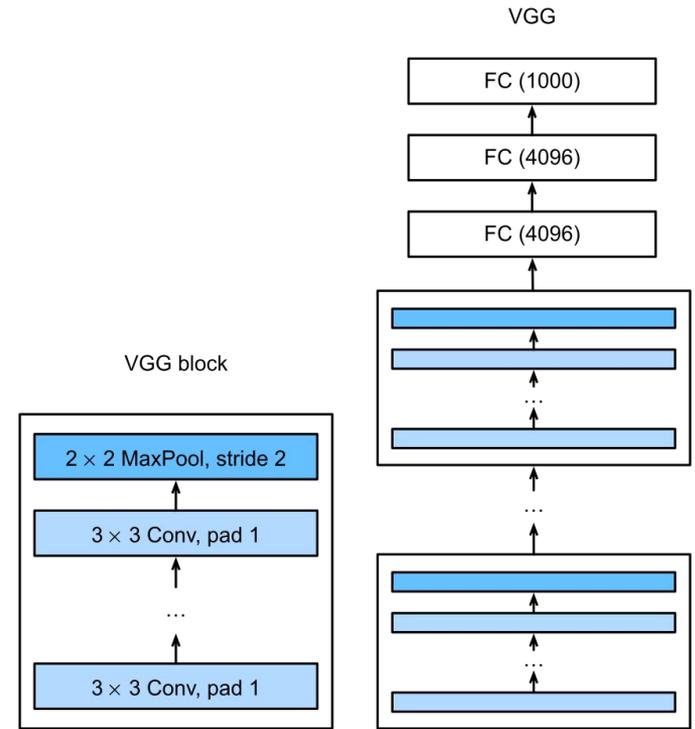
Animation from Hochschule der Medien

# Put Everything Together for Images



[LeNet-5, LeCun 1980]

# AlexNet (2012)



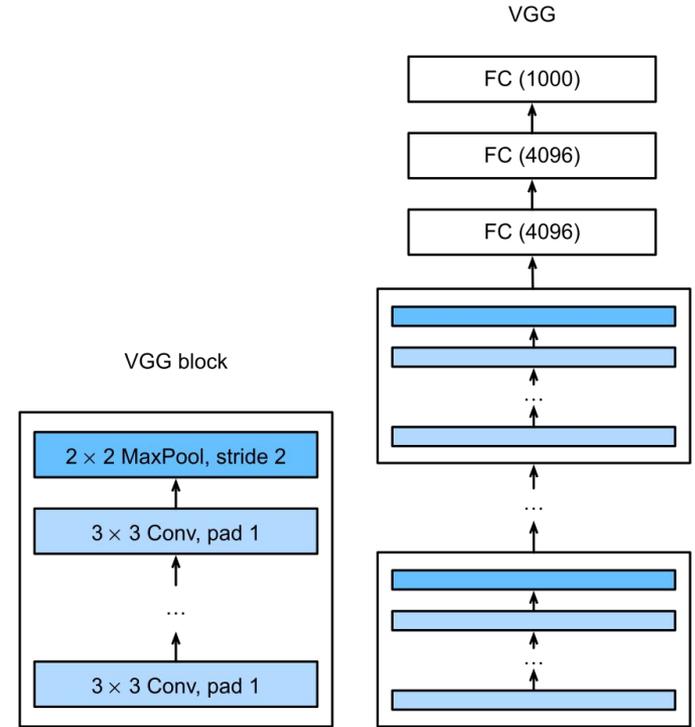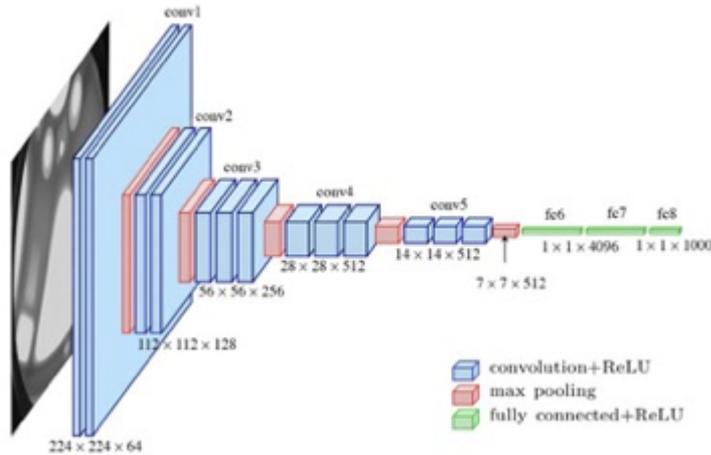| Layer |
|---|
| FC (1000) |
| FC (4096) |
| FC (4096) |
| 3 × 3 MaxPool, stride 2 |
| 3 × 3 Conv (256), pad 1 |
| 3 × 3 Conv (384), pad 1 |
| 3 × 3 Conv (384), pad 1 |
| 3 × 3 MaxPool, stride 2 |
| 5 × 5 Conv (256), pad 2 |
| 3 × 3 MaxPool, stride 2 |
| 11 × 11 Conv (96), stride 4 |
| Image (3 × 224 × 224) |

# VGGNet (2014)

- Very Deep CNN
- With only 3*3 conv filters
  - Fewer parameters, deeper nonlinear layers



VGG block



| 2 × 2 MaxPool, stride 2 |
| 3 × 3 Conv, pad 1 |
| … |
| 3 × 3 Conv, pad 1 |

VGG

FC (1000)

FC (4096)

FC (4096)

# VGGNet (2014)

- Very Deep CNN
- With only 3*3 conv filters
  - Fewer parameters, deeper nonlinear layers



VGG

FC (1000)

FC (4096)

FC (4096)

VGG block

2 × 2 MaxPool, stride 2

3 × 3 Conv, pad 1

...

3 × 3 Conv, pad 1
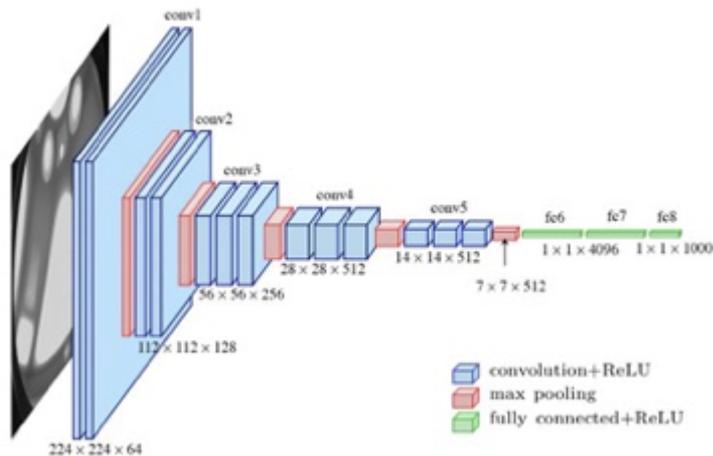
Keep increasing the depth?
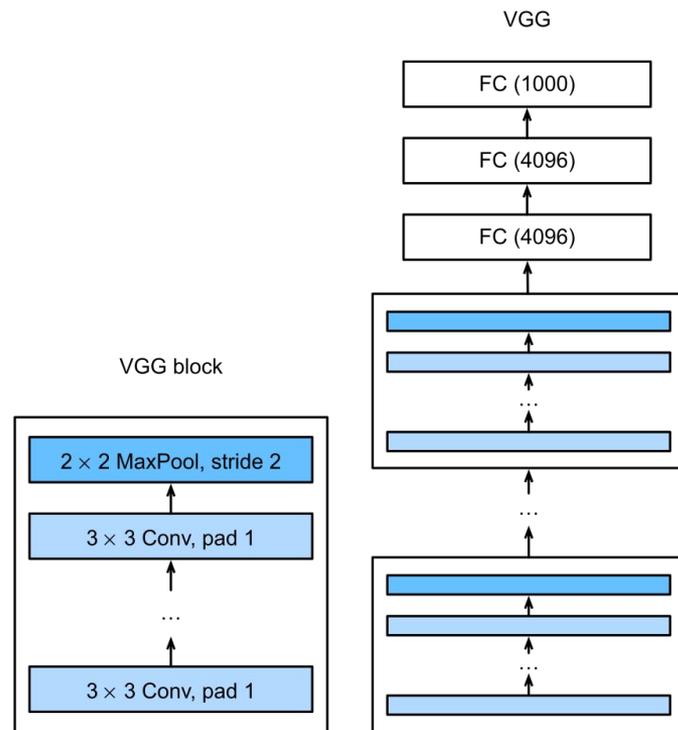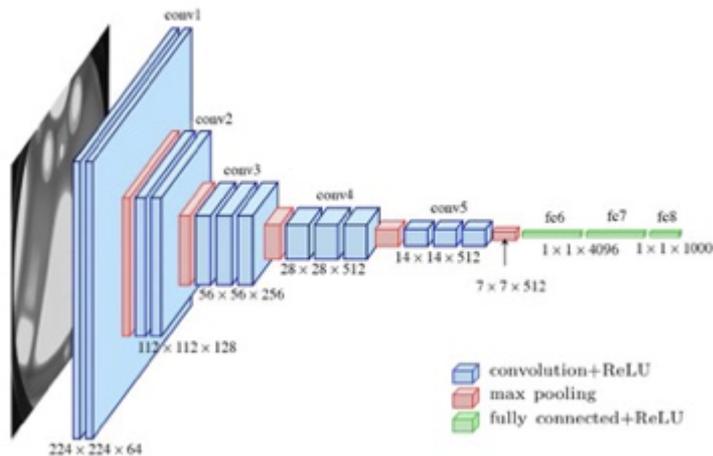
# VGGNet (2014)

- Very Deep CNN
- With only 3*3 conv filters
  - Fewer parameters, deeper nonlinear layers





vanishing or exploding

# VGGNet (2014)

- Very Deep CNN
- With only 3*3 conv filters
  - Fewer parameters, deeper nonlinear layers



VGG block

VGG

2 × 2 MaxPool, stride 2

3 × 3 Conv, pad 1

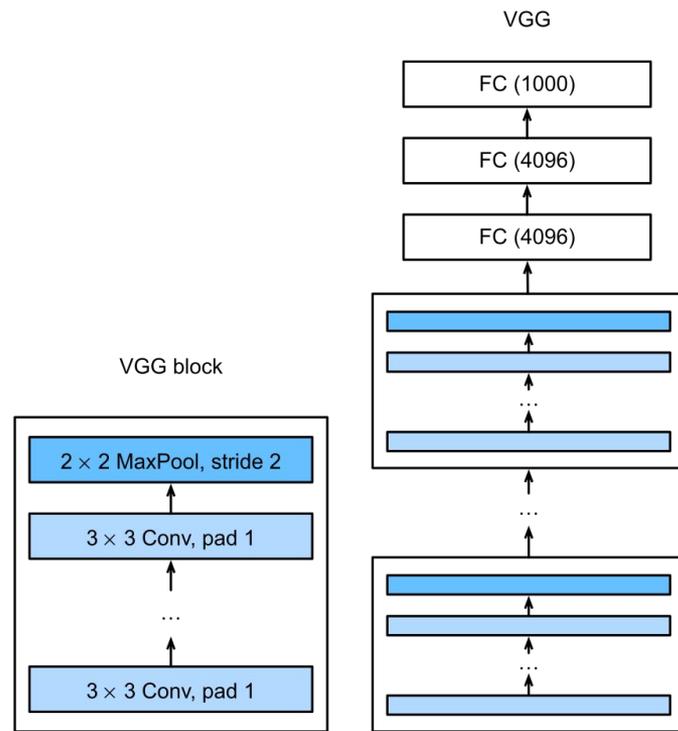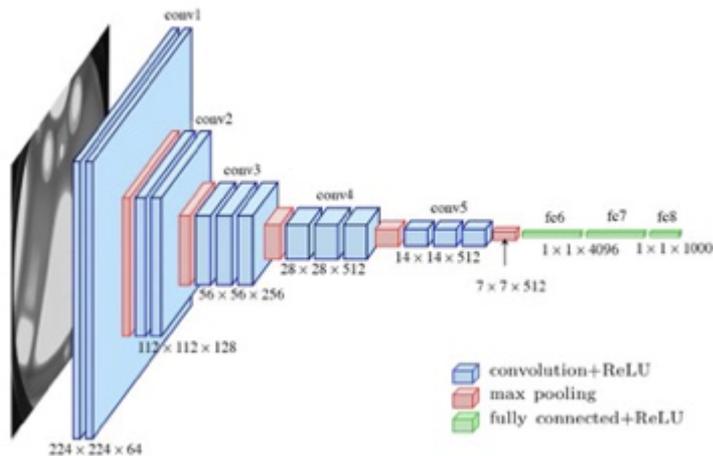3 × 3 Conv, pad 1

FC (1000)

FC (4096)

FC (4096)

vanishing or exploding

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a_n} \cdot \frac{\partial a_n}{\partial a_{n-1}} \cdot \frac{\partial a_{n-1}}{\partial a_{n-2}} \cdot \ldots \cdot \frac{\partial a_{i+1}}{\partial a_i} \cdot \frac{\partial a_i}{\partial w_i}$$

# VGGNet (2014)

- Very Deep CNN
- With only 3*3 conv filters
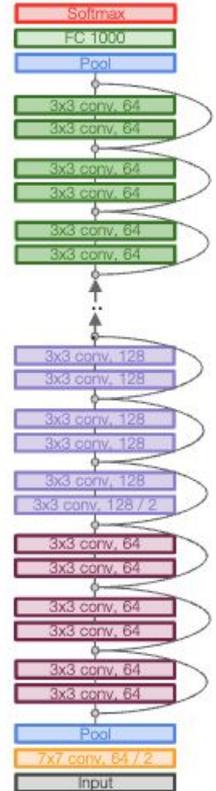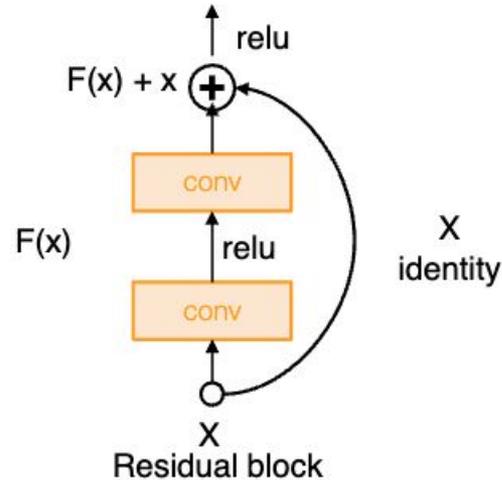  - Fewer parameters, deeper nonlinear layers



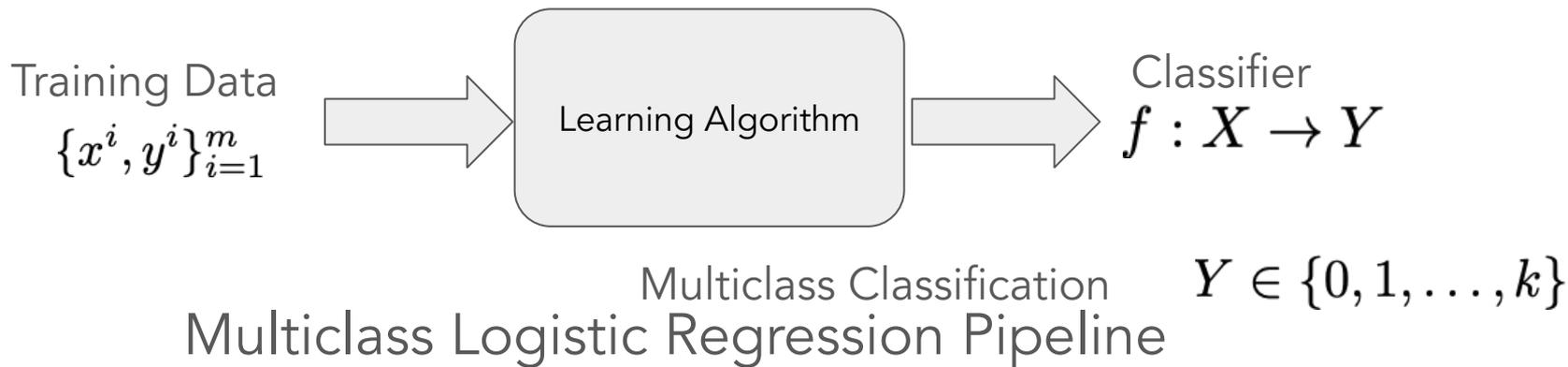VGG block

VGG

vanishing or exploding

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a_n} \cdot \frac{\partial a_n}{\partial a_{n-1}} \cdot \frac{\partial a_{n-1}}{\partial a_{n-2}} \cdot \ldots \cdot \frac{\partial a_{i+1}}{\partial a_i} \cdot \frac{\partial a_i}{\partial w_i} \qquad \frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a_n} \cdot \prod_{j=i}^{n} \frac{\partial a_j}{\partial a_{j-1}}$$

# ResNet (2015)

- Very Deep CNN with residual connections
  - 152-layer model for ImageNet
  - ILSVRC'15 classification winner (3.57% top 5 error)
  - Swept all classification and detection competitions in ILSVRC'15 and COCO'15!

# ResNet (2015)

- Very Deep CNN with residual connections
  - 152-layer model for ImageNet
  - ILSVRC'15 classification winner (3.57% top 5 error)
  - Swept all classification and detection competitions in ILSVRC'15 and COCO'15!
- Migrate gradient exploding and vanishing

$$y = x + F(x)$$

$$\frac{\delta E}{\delta x} = \frac{\delta E}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta E}{\delta y} * (1 + F'(x))$$

$$= \frac{\delta E}{\delta y} + \frac{\delta E}{\delta y} * F'(x)$$



F(x) + x

relu

F(x)

conv

relu

conv

X
identity

X
Residual block

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a_n} \cdot \prod_{j=i}^{n} \frac{\partial a_j}{\partial a_{j-1}}$$



Softmax
FC 1000
Pool
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128
3x3 conv, 128 / 2
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
3x3 conv, 64
Pool
7x7 conv, 64 / 2
Input

# Multiclass Logistic Regression Algorithms

Training Data
$\{x^i, y^i\}_{i=1}^m$

Learning Algorithm

Classifier
$f : X \rightarrow Y$

Multiclass Classification   $Y \in \{0, 1, \ldots, k\}$

## Multiclass Logistic Regression Pipeline

1. Build probabilistic models:
   Categorical Distribution + Conv NN
2. Derive loss function: MLE and MAP
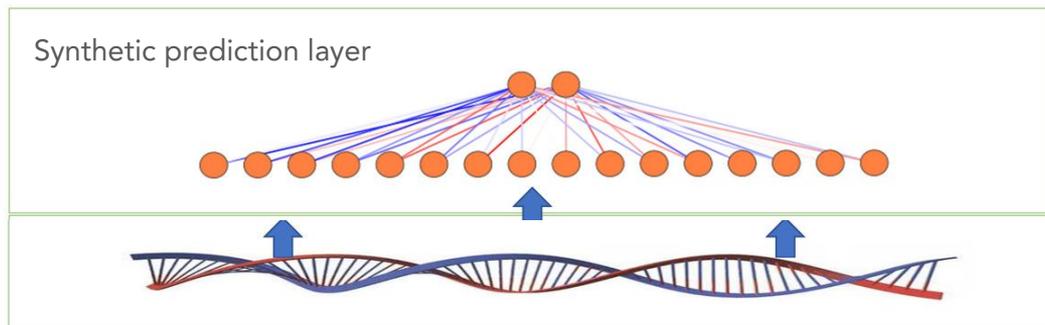3. Select optimizer: (Stochastic) Gradient Descent
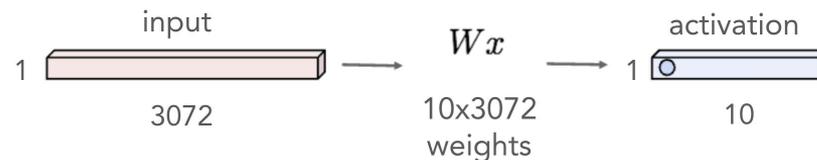
# Sequence Prediction



**texts[0]**

For a movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.
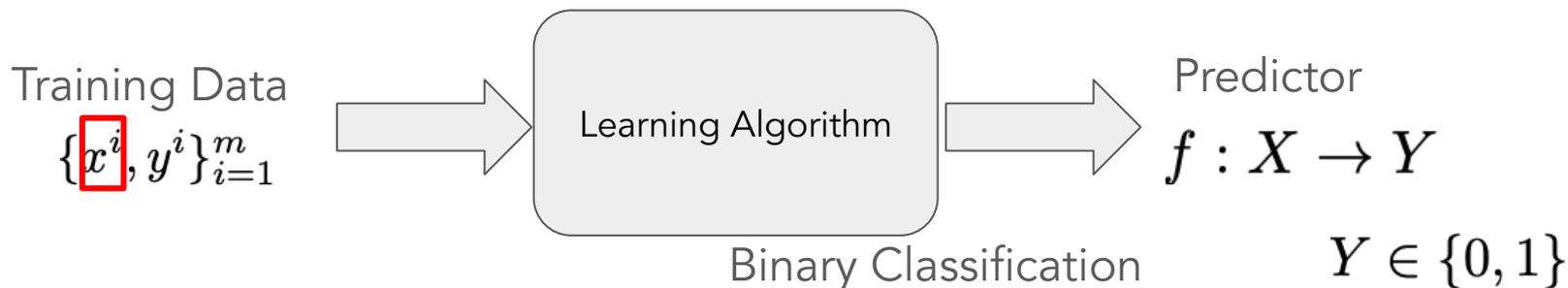
# Sequence Prediction

**texts[0]**

For a movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.

input

$Wx$

activation

1

3072

10x3072 weights

1

10

# Sequence Prediction

Synthetic prediction layer

I saw the movie with two grown children. Although it was not as clever as Shrek, I thought it was rather good. In a movie theatre surrounded by children who were on spring break, there was not a sound so I know the children all liked it. There parents also seemed engaged. The death and apparent death of characters brought about the appropriate gasps and comments. Hopefully people realize this movie was made for kids.

For a movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.

input

$Wx$

activation

1

3072

10x3072
weights

1

10

What if the length of sequences varies?

# Sequential Regression Algorithms

Training Data
$$\{x^i, y^i\}_{i=1}^m$$

Learning Algorithm

Predictor
$$f : X \rightarrow Y$$

$$Y \in \mathbb{R}$$

## Linear Regression Pipeline

1. Build probabilistic models:
   Gaussian Distribution + RNN
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent
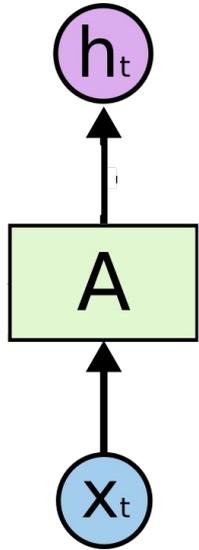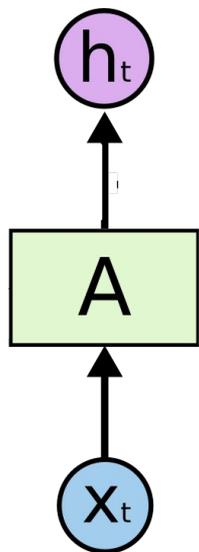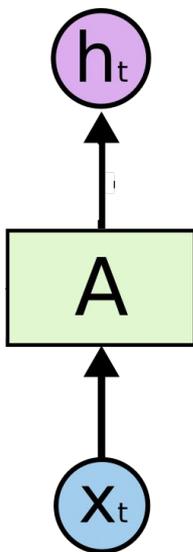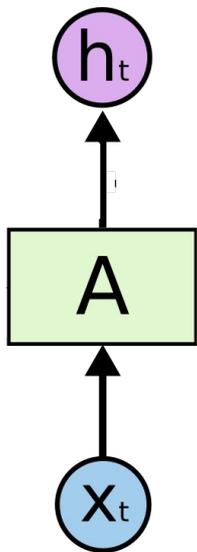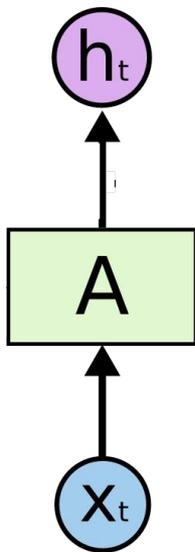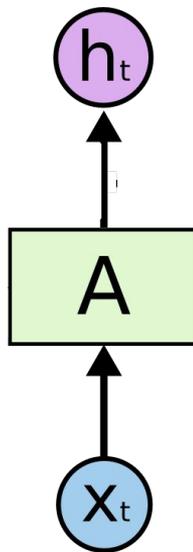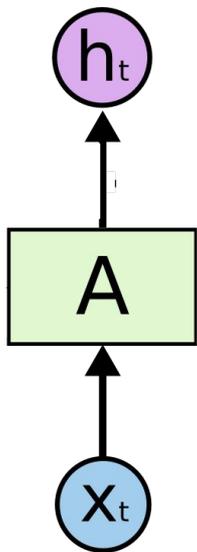
# Sequential Binary Classification Algorithms

Training Data
$$\{x^i, y^i\}_{i=1}^m$$

Learning Algorithm

Binary Classification

Predictor
$$f : X \to Y$$

$$Y \in \{0, 1\}$$

## Binary Logistic Regression Pipeline

1.  Build probabilistic models:
    Bernoulli Distribution + RNN
2.  Derive loss function: MLE and MAP
3.  Select optimizer: (Stochastic) Gradient Descent

# Sequential Multiclass Logistic Regression Algorithms

Training Data
$$\{x^i, y^i\}_{i=1}^m$$

Learning Algorithm

Classifier
$$f : X \rightarrow Y$$

Multiclass Classification $\quad Y \in \{0, 1, \ldots, k\}$

## Multiclass Logistic Regression Pipeline

1. Build probabilistic models:
   Categorical Distribution + RNN
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

# Inspiration from Convolution Layer



32

32

3

one computation cell is shared

# Recurrent Neural Network



$h_t$

A

$x_t$

For

For a movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.
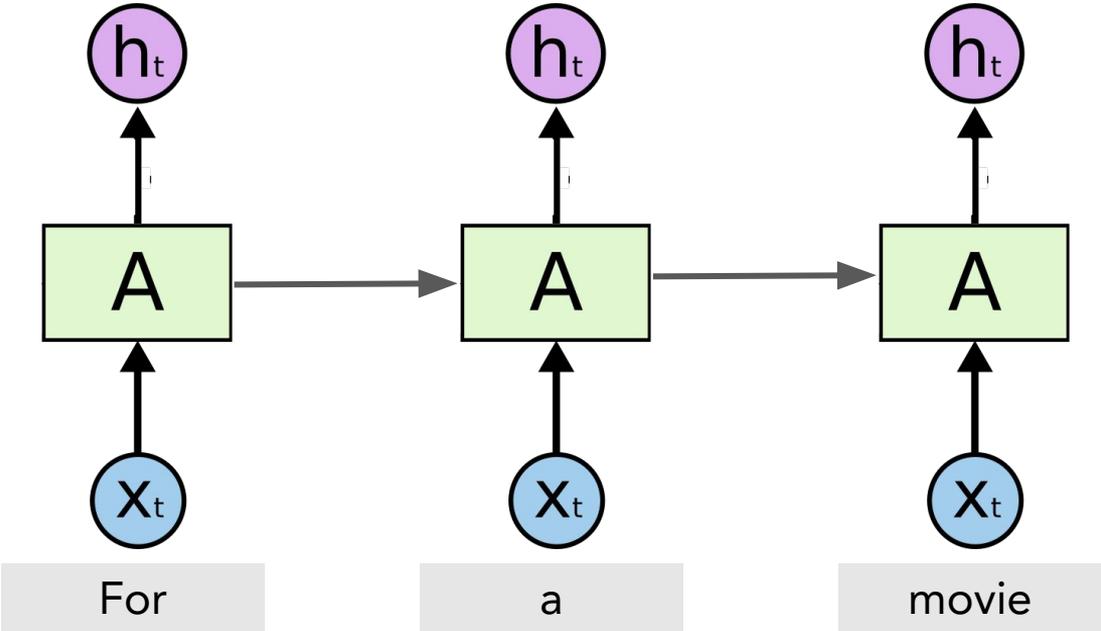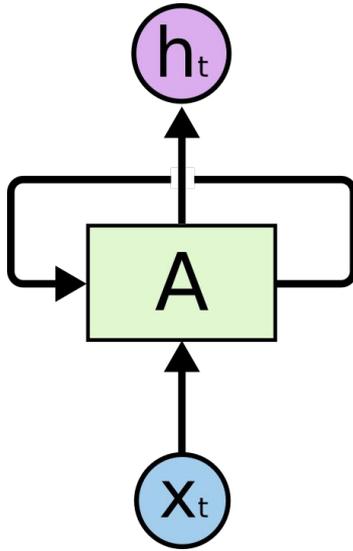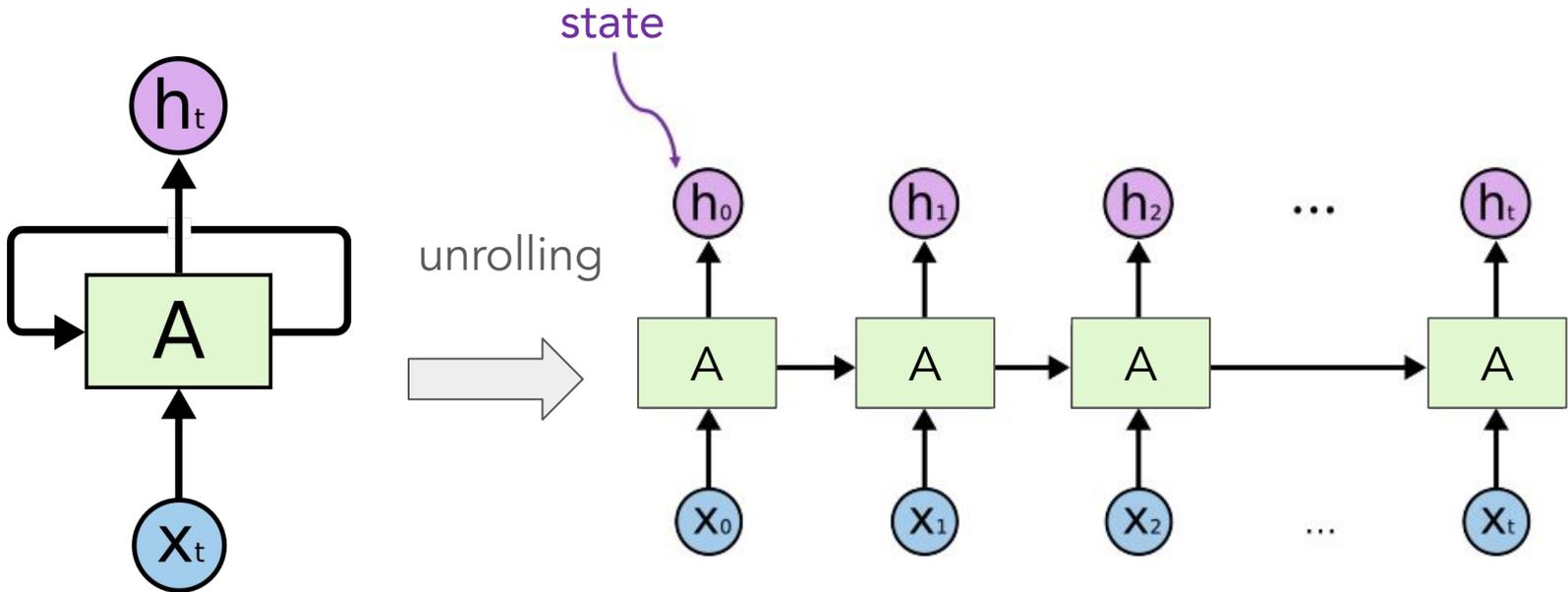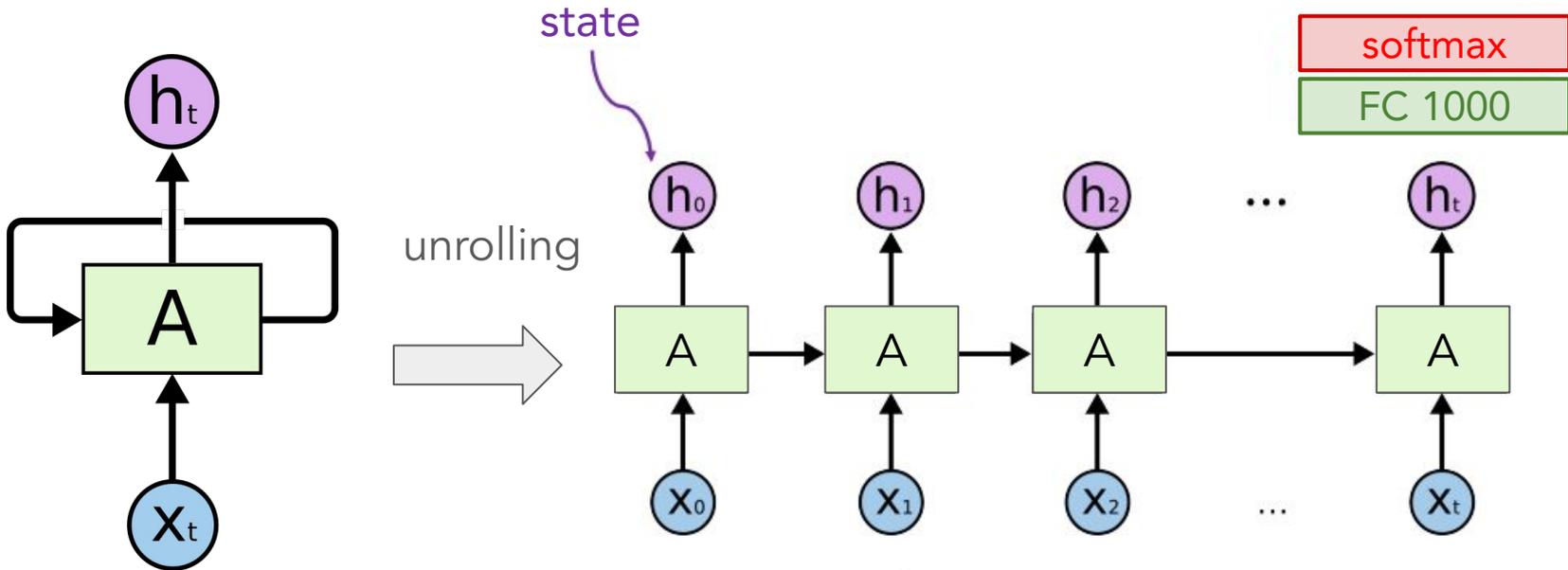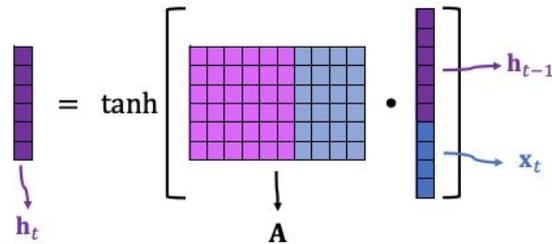
# Recurrent Neural Network



For $a$ movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.
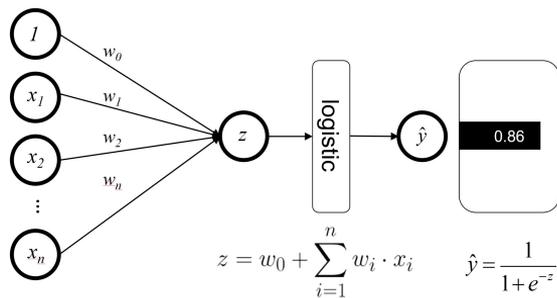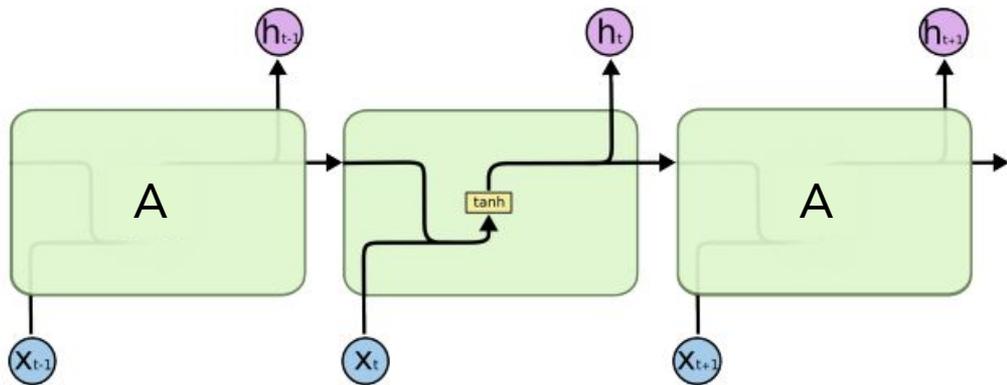
# Recurrent Neural Network



For a [movie] that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.

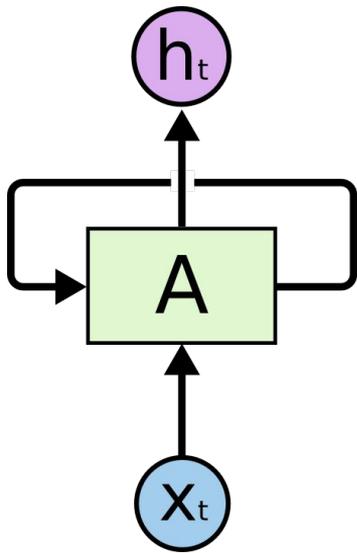# Recurrent Neural Network



Order Matters!

For a `movie` that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.
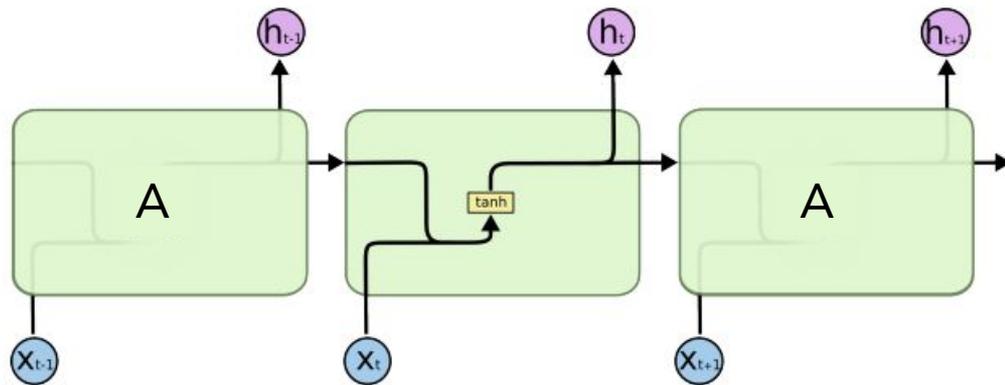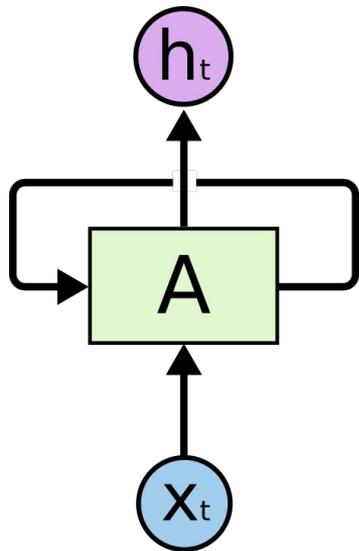
# Recurrent Neural Network



For a [movie] that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.
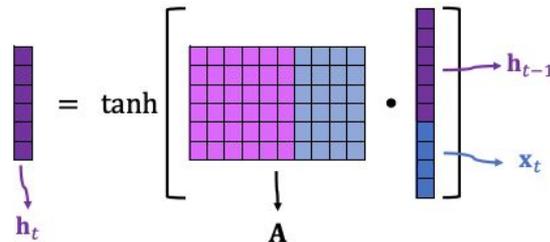
# Recurrent Neural Network



For a movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.
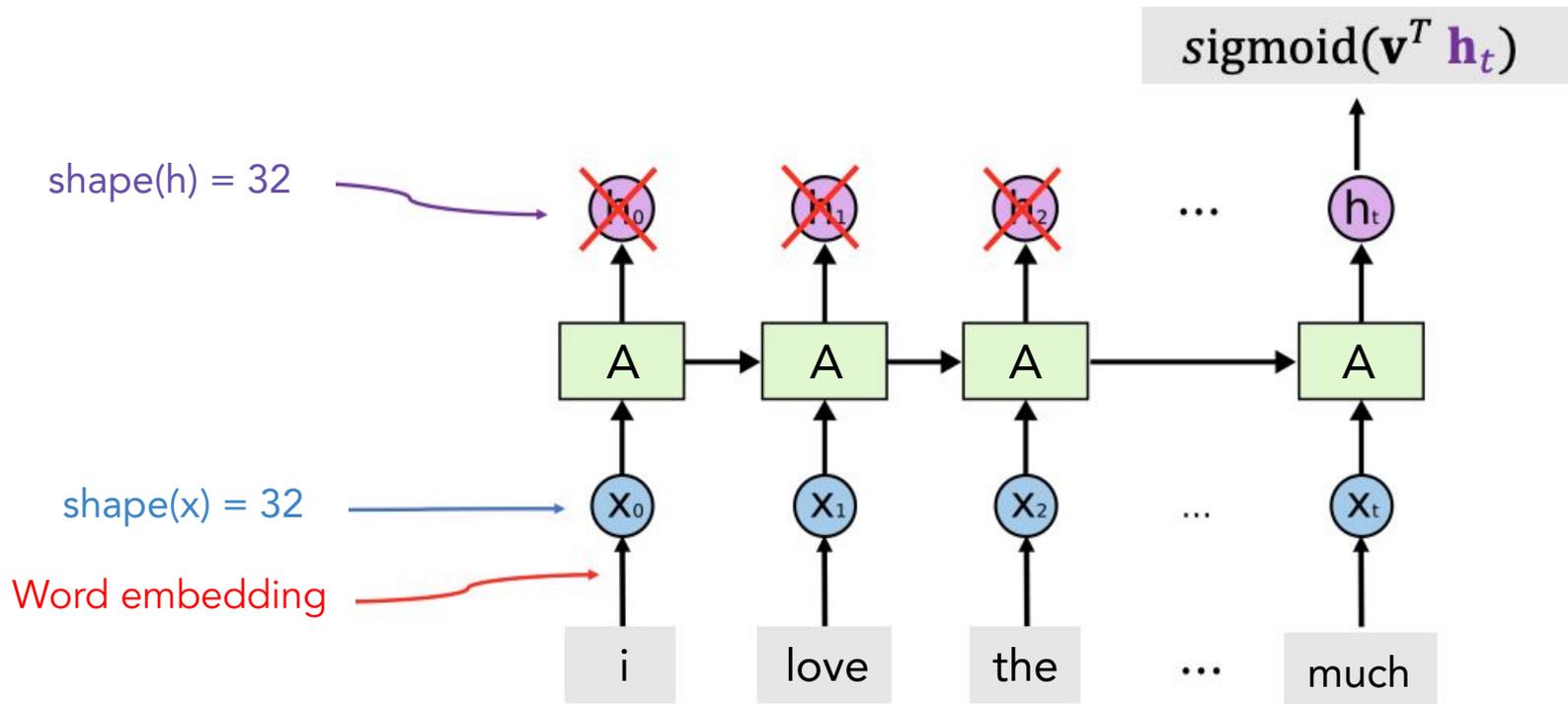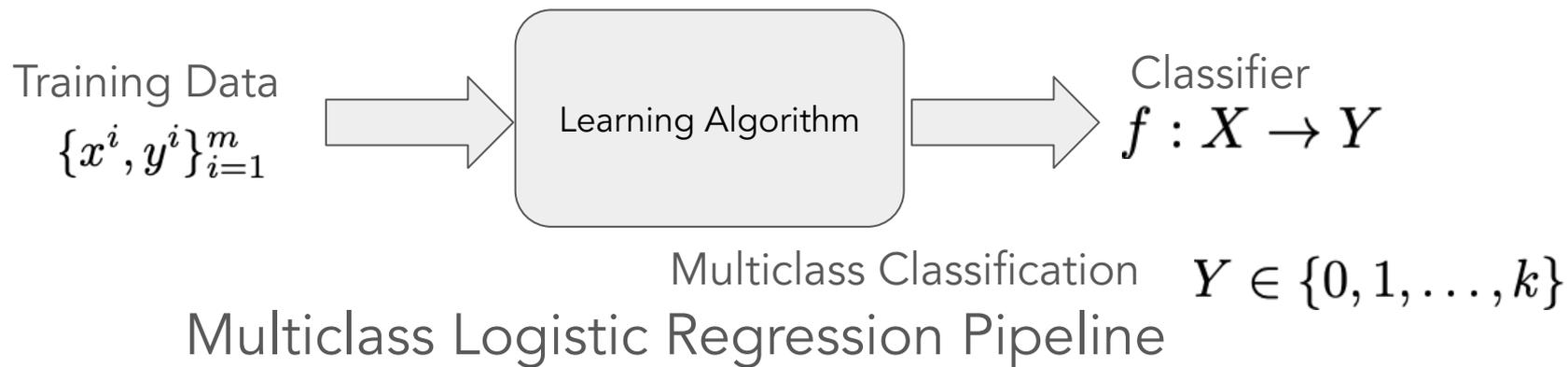
# Recurrent Neural Network



state

unrolling

For a movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.

# Recurrent Neural Network



state

softmax

FC 1000

unrolling

For a movie that gets no respect there sure are a lot of memorable quotes listed for this gem. Imagine a movie where Joe Piscopo is actually funny! Maureen Stapleton is a scene stealer. The Moroni character is an absolute scream. Watch for Alan "The Skipper" Hale jr. as a police Sgt.

# Simple RNN Cell



$$z = w_0 + \sum_{i=1}^{n} w_i \cdot x_i \qquad \hat{y} = \frac{1}{1+e^{-z}}$$

# Simple RNN Cell



| | | |
|---|---|---|
| Binary step | | $\begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$ |
| Logistic, sigmoid, or soft step | | $\sigma(x) \doteq \dfrac{1}{1+e^{-x}}$ |
| Hyperbolic tangent (tanh) | | $\tanh(x) \doteq \dfrac{e^{x} - e^{-x}}{e^{x} + e^{-x}}$ |
| Rectified linear unit (ReLU)[13] | | $(x)^{+} \doteq \begin{cases} 0 & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ $= \max(0, x) = x\mathbf{1}_{x>0}$ |
| Gaussian Error Linear Unit (GELU)[5] | | $\dfrac{1}{2}x\left(1 + \operatorname{erf}\left(\dfrac{x}{\sqrt{2}}\right)\right)$ where erf is the $= x\Phi(x)$ gaussian error function. |
| Softplus[14] | | $\ln(1 + e^{x})$ |
| Exponential linear unit (ELU)[15] | | $\begin{cases} \alpha\left(e^{x} - 1\right) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$ with parameter $\alpha$ |

# Simple RNN for IMDB Review

$$\text{sigmoid}(\mathbf{v}^T \mathbf{h}_t)$$

shape(h) = 32

shape(x) = 32

Word embedding

| | | | |
|---|---|---|---|
| $h_0$ ✗ | $h_1$ ✗ | $h_2$ ✗ | $h_t$ |
| A | A | A | A |
| $x_0$ | $x_1$ | $x_2$ | $x_t$ |
| i | love | the | much |

# Sequential Multiclass Logistic Regression Algorithms

Training Data
$$\{x^i, y^i\}_{i=1}^m$$

Learning Algorithm

Classifier
$$f : X \to Y$$

Multiclass Classification

$$Y \in \{0, 1, \dots, k\}$$

## Multiclass Logistic Regression Pipeline

1. Build probabilistic models:
   Categorical Distribution + RNN
2. Derive loss function: MLE and MAP
3. Select optimizer: (Stochastic) Gradient Descent

# Backpropagation SGD

$$\ell \left[ \text{sigmoid}(\mathbf{v}^T \, \mathbf{h}_t) \quad y \right]$$

shape(h) = 32

shape(x) = 32

Word embedding

i    love    the    ...    much

# Backpropagation SGD

Layer $\Longleftrightarrow$ Time-step



shape(h) = 32

shape(x) = 32

Word embedding

# Backpropagation SGD

Backpropagation through Time
(BPTT)



$$\ell \left[ y \atop \text{sigmoid}(\mathbf{v}^T \mathbf{h}_t) \right]$$

much

the

love

i

shape(h) = 32

shape(x) = 32

Word embedding

# Backpropagation SGD



$$\ell \left[ \text{sigmoid}(\mathbf{v}^T \mathbf{h}_t) \quad y \right]$$

shape(h) = 32

shape(x) = 32

Word embedding

i    love    the    ...    much

# More Usages of RNNs

many to one



IMDB text review classification

# More Usages of RNNs

many to one

one to many

Image Captioning
image -> sequence of words

# More Usages of RNNs

one to many  many to one  many to many  many to many



Translation

# Training of RNNs

one to many          many to one          many to many          many to many

Backpropagation Through Unrolling Steps

# Simple RNN is not good at long-term dependence



$\mathbf{h}_{100}$ is almost irrelevant to $\mathbf{x}_1$: $\dfrac{\partial \mathbf{h}_{100}}{\partial \mathbf{x}_1}$ is near zero or exploding.

Gradient Vanishing Again!

# Simple RNN is not good at long-term dependence



$\mathbf{h}_{100}$ is almost irrelevant to $\mathbf{x}_1$: $\frac{\partial \mathbf{h}_{100}}{\partial \mathbf{x}_1}$ is near zero or exploding.

Gradient Vanishing Again! -> ShortCut connection

# Long Short Term Memory (LSTM)



Simple RNN

LSTM

Figures from Christopher Olah's blog

# LSTM Cell: Conveyor Belt

The past information directly flows to the future.



$$C_t = C_{t-1} \circ f_t + \hat{C}_t \circ i_t$$

# LSTM Cell: Conveyor Belt

The past information directly flows to the future.  (ShortCut connection)



$$C_t = C_{t-1} \circ f_t + \hat{C}_t \circ i_t$$

# LSTM Cell: Forget Gate

A value of *zero* means "let *nothing* through."

A value of *one* means "let *everything* through!"

$$f_t = \sigma(W_f x_t + U_f h_{t-1})$$



$$C_t = C_{t-1} \circ f_t + \hat{C}_t \circ i_t$$

# LSTM Cell: Input Gate

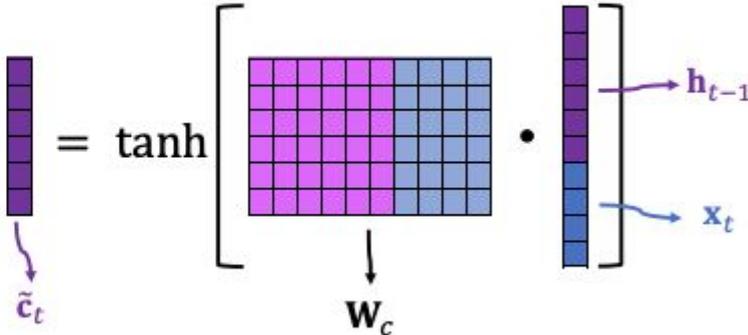How much information current context provided.

$$i_t = \sigma(W_i x_t + U_i h_{t-1})$$



$$C_t = C_{t-1} \circ f_t + \hat{C}_t \circ i_t$$

# LSTM Cell: New Value

"local" context, only up to immediately preceding state

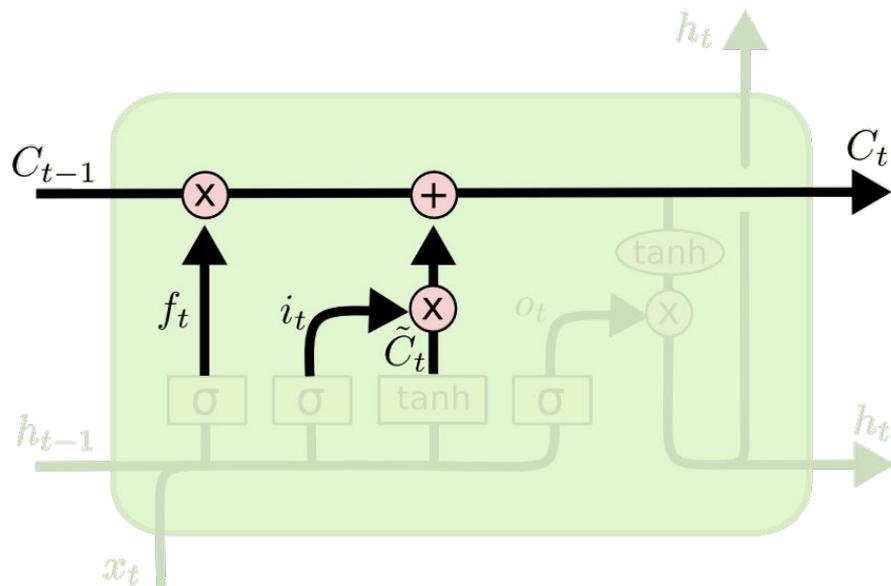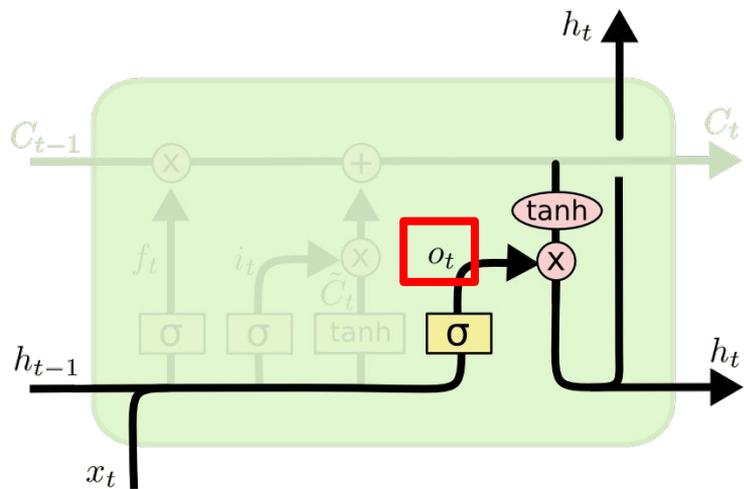$$\hat{C}_t = \sigma(W_c x_t + U_c h_{t-1})$$



$$C_t = C_{t-1} \circ f_t + \hat{C}_t \circ i_t$$

# LSTM Cell: Update Conveyor Belt
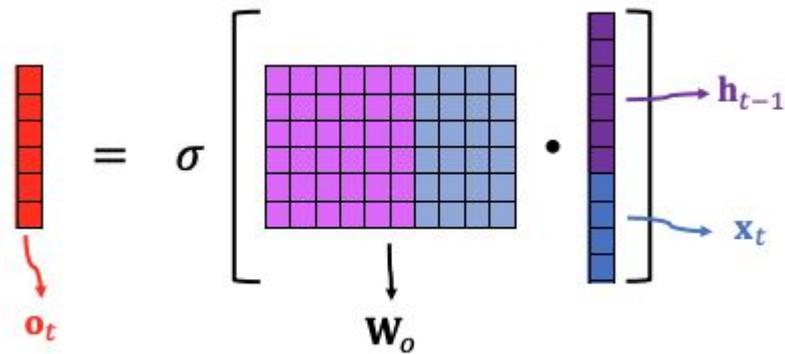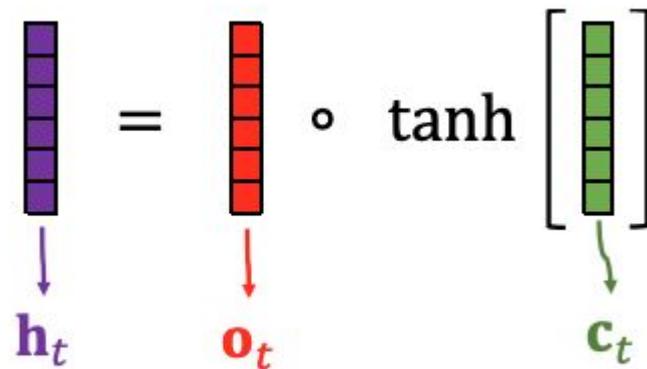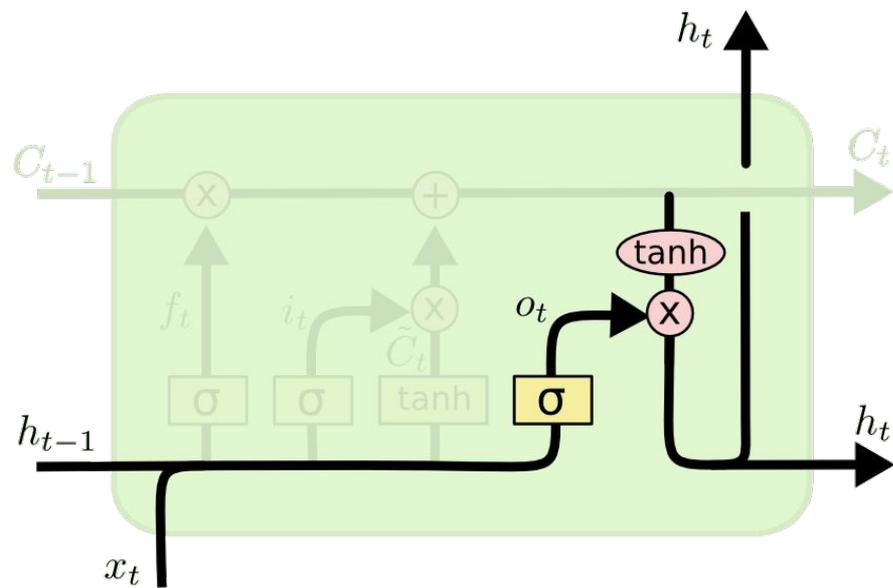


$$C_t = C_{t-1} \circ f_t + \hat{C}_t \circ i_t$$

# LSTM Cell: Output Gate
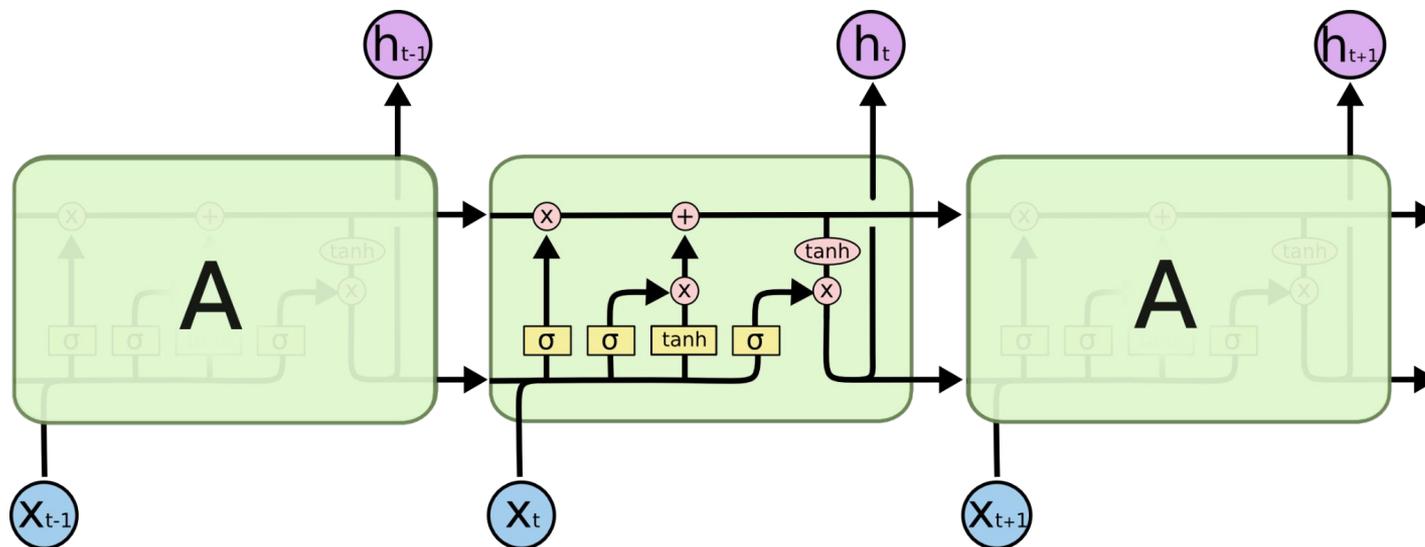
$$h_t = o_t * \tanh\left(C_t\right)$$

# LSTM Cell: Update State



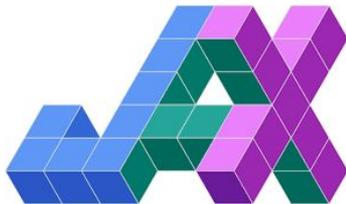$$h_t = o_t * \tanh\left(C_t\right)$$

# LSTM Cell
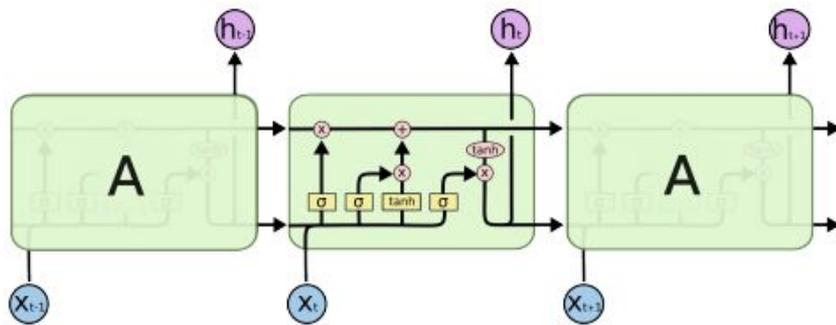
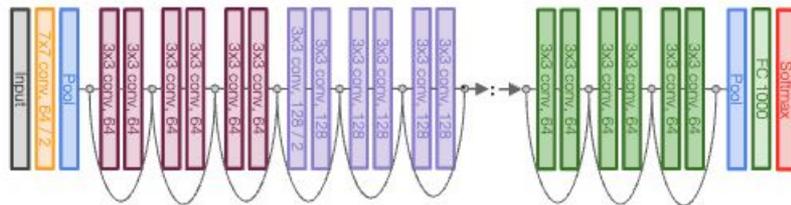# Auto-differentiation Packages

PyTorch

JAX

Tensorflow

# LSTM vs. ResNet



$$C_t = C_{t-1} \circ f_t + \hat{C}_t \circ i_t$$

LSTM

Similar to ResNet

# Other Variants of RNN

Gated Recurrent Unit



$$z_t = \sigma \left( W_z \cdot [h_{t-1}, x_t] \right)$$

$$r_t = \sigma \left( W_r \cdot [h_{t-1}, x_t] \right)$$

$$\tilde{h}_t = \tanh \left( W \cdot [r_t * h_{t-1}, x_t] \right)$$
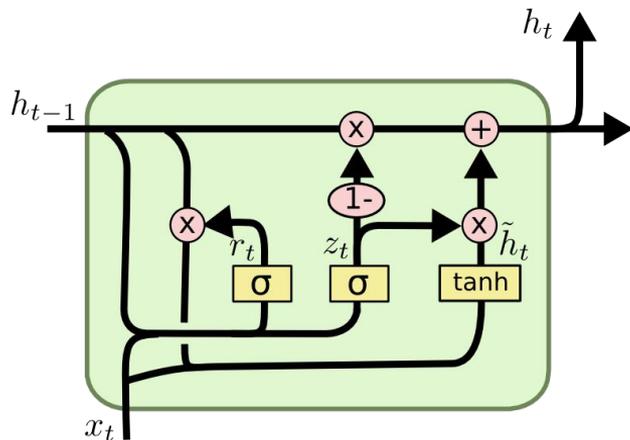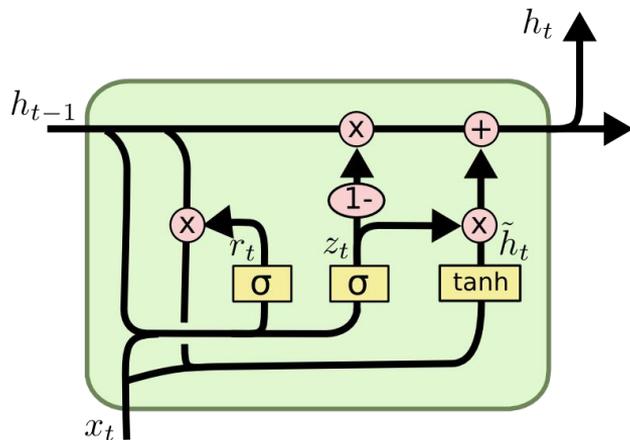
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Other Variants of RNN

Gated Recurrent Unit



$$z_t = \sigma\left(W_z \cdot [h_{t-1}, x_t]\right)$$

$$r_t = \sigma\left(W_r \cdot [h_{t-1}, x_t]\right)$$

$$\tilde{h}_t = \tanh\left(W \cdot [r_t * h_{t-1}, x_t]\right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Summary

- RNNs allow a lot of flexibility in architecture design

# Summary

- RNNs allow a lot of flexibility in architecture design

- Vanilla RNNs are simple but don't work very well

# Summary

- RNNs allow a lot of flexibility in architecture design

- Vanilla RNNs are simple but don't work very well

- Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)

PyTorch session will be online next Monday

# Q&A