

CSCI 599: Deep Learning and its Applications

Lecture 5

Spring 2019
Joseph J. Lim

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Office Hours

- Instructor OH @ RTH 402
 - Tuesday 2-3:30pm
 - This is NOT for homework related questions.
- TA OH @ SAL 125
 - Monday 5-6pm & Tuesday 12:30-1:30pm
 - Extra OH for assignment due & midterm weeks:
Monday 4-5pm & Tuesday 11:30-12:30pm

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

CSCI 599/699

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Communication

- Please use **Piazza** for any general communication including questions
<https://piazza.com/usc/spring2019/csci599/home>
- E-mails will be ignored.
- **Register TODAY. Look for your project team mates!**

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Assignment 1

- Assignment 1 is released on the website (<https://csci599-dl.github.io/>)
- DUE February 19th, 2019 (week 7)
- Collaboration is OK! Please list all collaborators' names.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

** Date changes **

- Midterm
- Project proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Important Dates

- Entrance exam: 1/15
- **Assignment 1: 2/19**
- **Midterm: 2/26**
- Project meeting with Instructor #1: 3/6 — 3/8
- Assignment 2: 3/26
- Project meeting with Instructor #2: 4/1 — 4/3
- Project meeting with TA: 2 times (arranged later)
- Final presentation: 4/23 5-9:00pm **4 hours**

Subject to change!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Course Project

Subject to change!

- Computational resource (**be considerate**)
\$150 Google Cloud credit per student
\$125 Amazon AWS credit per student
- Tentative Schedule for Project
 - **Week 5 (2/5): Course Project Team**
 - Week 8 (2/26): Course Project Proposal
 - Week 13 (4/2): Mid-report
 - Week 16 (4/23): Project Presentation (5-9pm) + Report

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Course Project

- Team-based project (4-5 students per team)
- Each team will have at least 1 dedicated TA
 - 2 Mandatory meetings with TA
 - 2 meetings with me
- Create your own problems (extra points)
 - **Talk and discuss** with your TAs and me!
 - In the worst case, we will give a project idea
 - Less fun, Less points!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Project Evaluation

- Creativity and difficulty of the problem setup
- Novelty of the approach
- Thoroughness of the experiments
- Quality of student's presentation, report, and meetings with TA/instructor

Extra credit for creating your own project (OK to discuss and get help from TAs)

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Project

- Report: a blog post (rather than a paper format)
- Lots of “what if we fail” questions

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Team Formation

- Due today!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Project meeting with TA #1

- Try to bring something concrete!
- We will post on Piazza your TA assignment
 - Temporary until the proposal due date
- Please sign up for your team on your TA's doodle link.
- Please be available — if not, we may not be able to schedule one with you.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Today's agenda

- Recap: CNNs
- Training Neural Networks
- Recurrent Neural Networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Today's agenda

- Recap: CNNs
- Training Neural Networks
- Recurrent Neural Networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Computational Graph

Function: $f(x, y, z) = x \cdot y + z$

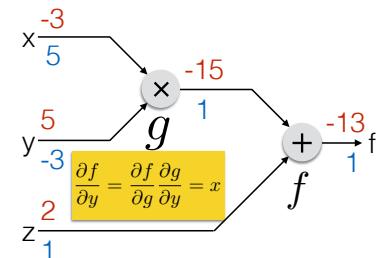
Backpropagation

$$\frac{\partial f}{\partial x} = 5$$

$$\text{We get: } \frac{\partial f}{\partial y} = -3$$

$$\frac{\partial f}{\partial z} = 1$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial y} = x$$

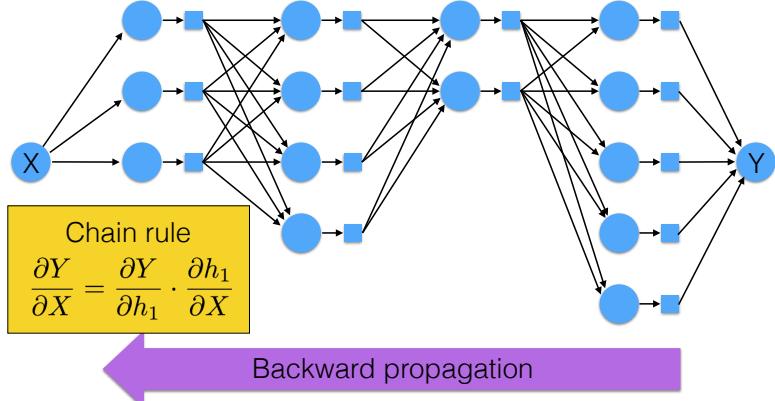


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Recap: Backward Propagation

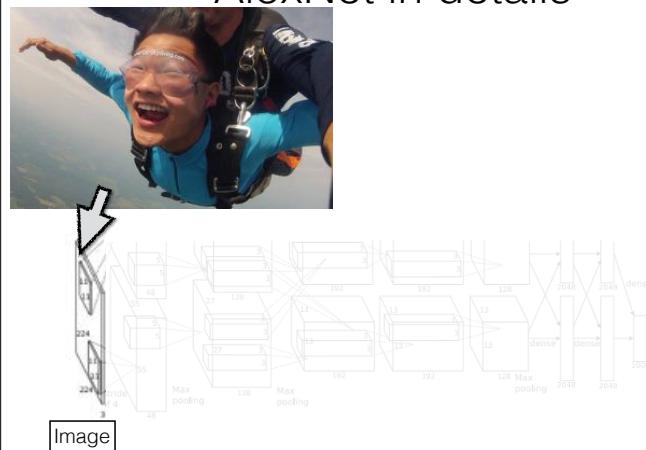


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details



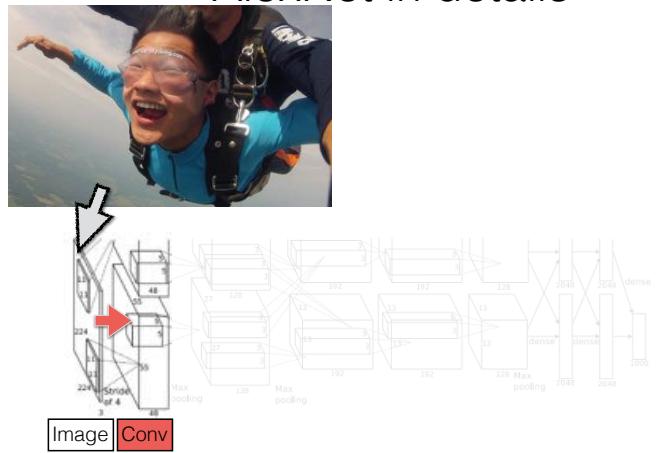
PC: Daniel Yang

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

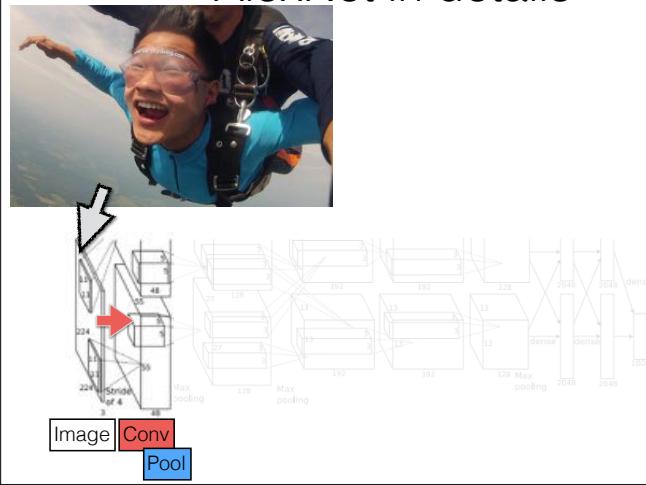


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

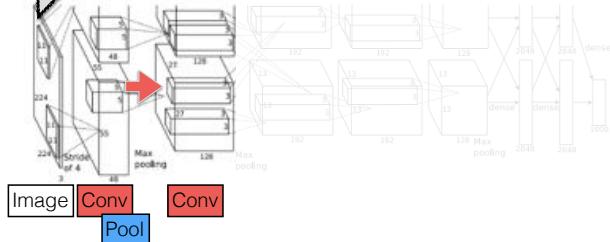


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

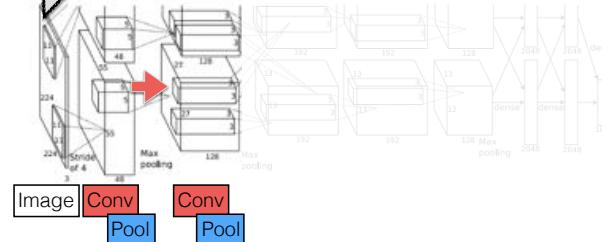


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

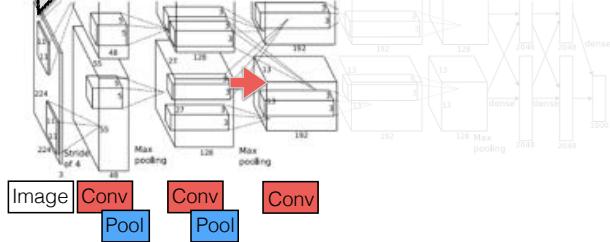


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

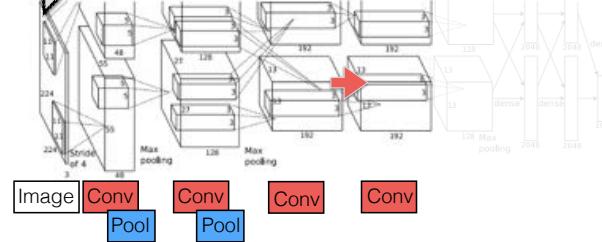


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

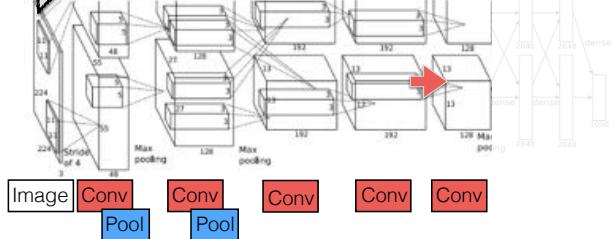


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

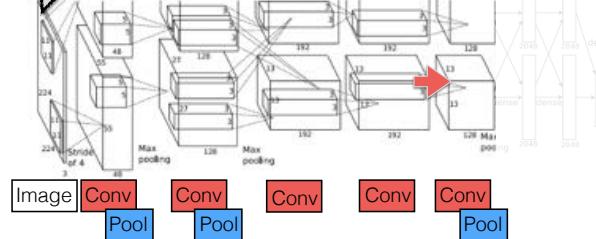


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

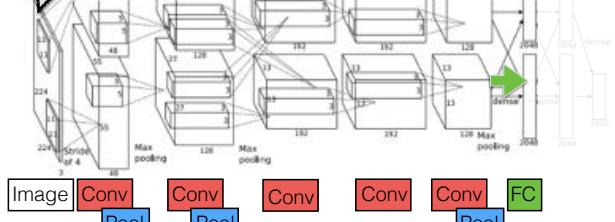


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

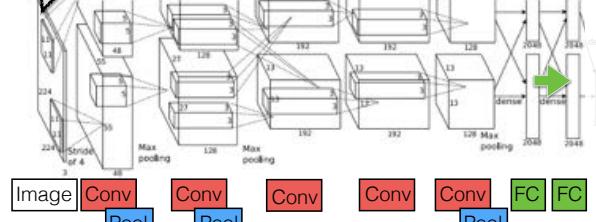


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

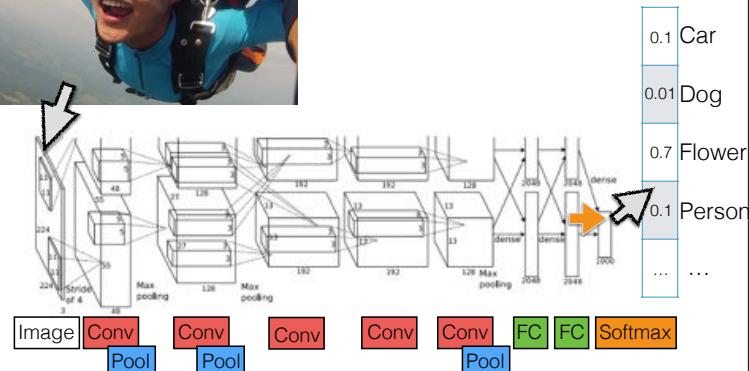


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details

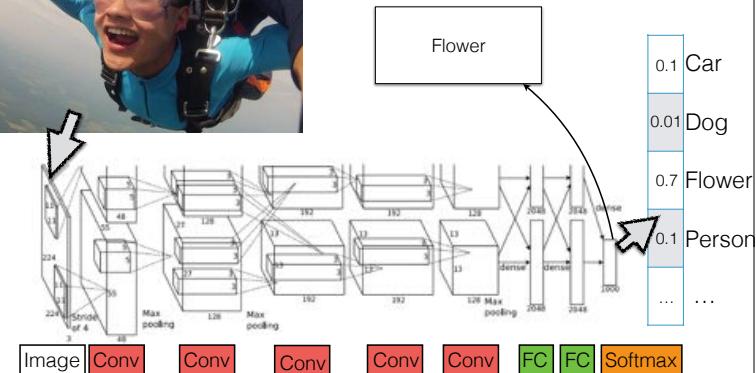


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

What are these layers?

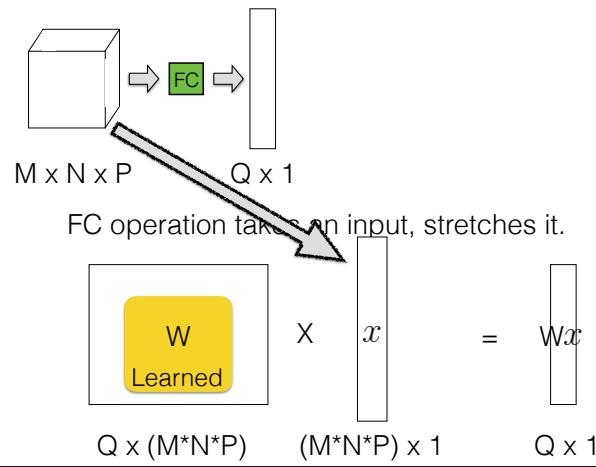
- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Fully connected layer (FC)



Joseph J. Lim

CSCI 599 @ USC

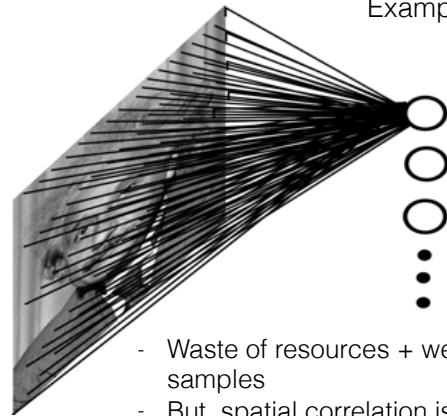
Lecture 5

Fully connected layer (FC)

Example: 200x200 image

40k hidden units

~2B parameters!!!



- Waste of resources + we have not enough training samples
- But, spatial correlation is local!

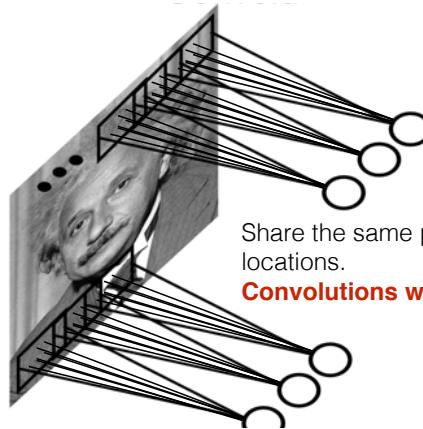
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Slide credit: Marc'Aurelio Ranzato

Convolutional layer



Share the same parameters across different locations.
Convolutions with learned kernels (weights)

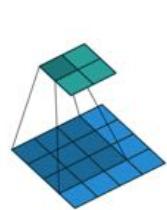
Slide credit: Marc'Aurelio Ranzato

Joseph J. Lim

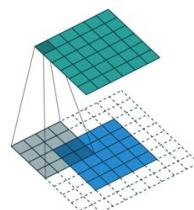
CSCI 599 @ USC

Lecture 5

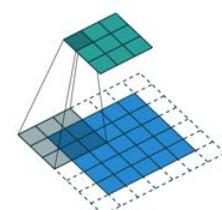
Convolutional layer



padding 0, stride 1



padding 2, stride 1



padding 1, stride 2

Note: it can have **depth** (3rd dimension).

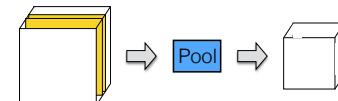
Image credit: vdumoulin (github): https://github.com/vdumoulin/conv_arithmetic

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Max pooling layer



M x N x P

m x n x P

1	3	2	7	2
5	1	3	0	5
4	8	5	0	6
3	2	6	0	8
0	7	4	6	9

max pool with
3x3 filters & stride 2

8	7
8	9

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Activation functions

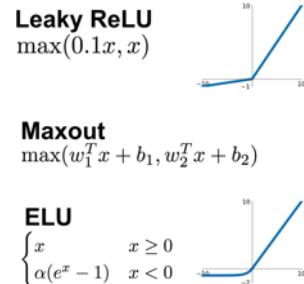
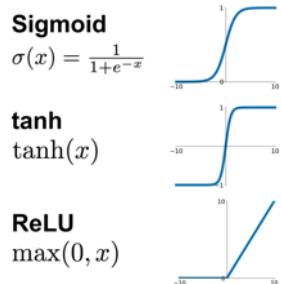


Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Summary

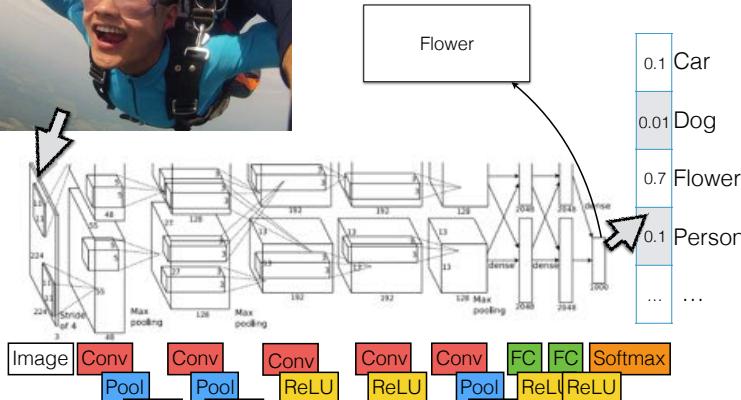
- Fully connected layer
- Convolutional layer
- Deconvolutional layer
- Pooling layers
- Activation layers (e.g. ReLU)

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

AlexNet in details



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Encoder

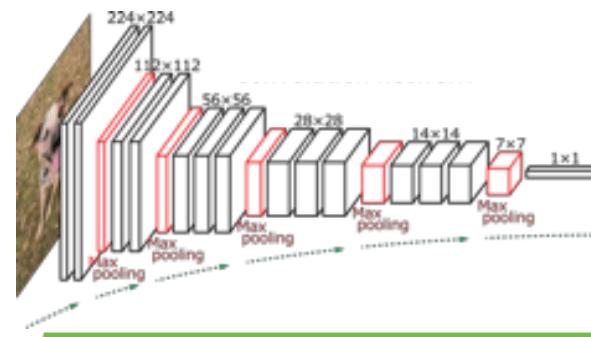


Image credit: Noh, et. al. Learning Deconvolution Network for Semantic Segmentation. ICCV 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Decoder

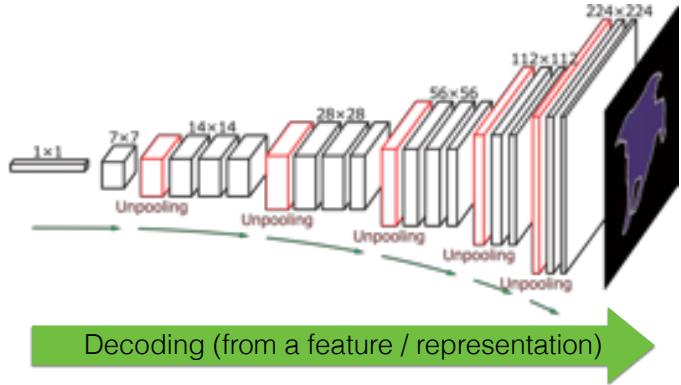


Image credit: Noh, et. al. Learning Deconvolution Network for Semantic Segmentation. ICCV 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Let's revisit CV tasks!

CNN applications

- Image classification
- Object detection
- Object segmentation
- Image resolution
- Human pose detection
- Action classification
- Image captioning

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image Classification



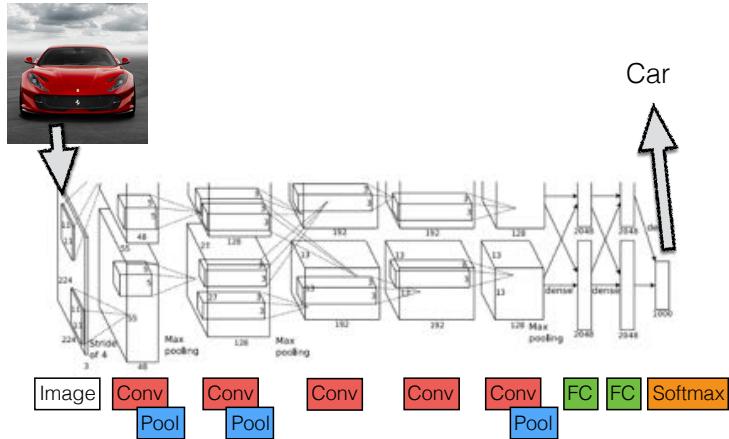
→ Car

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image Classification



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



1. Object proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



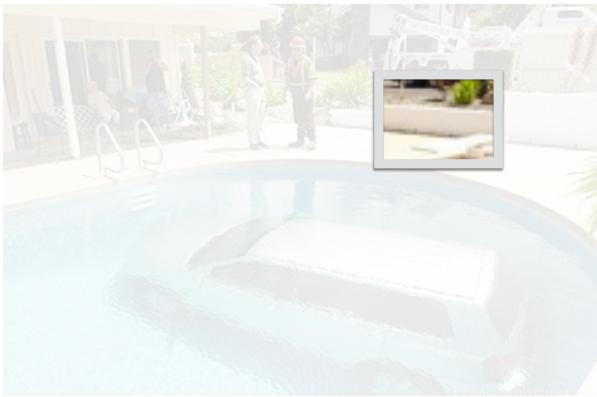
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



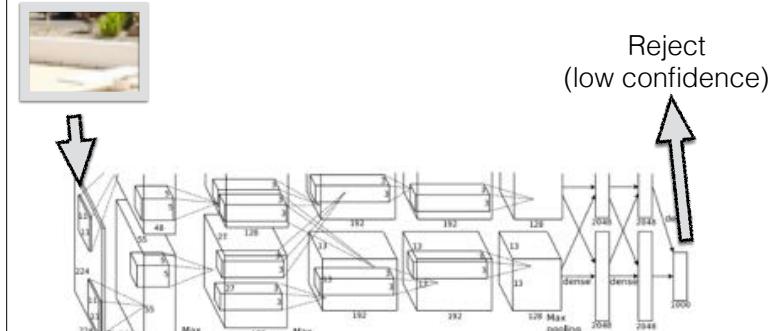
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



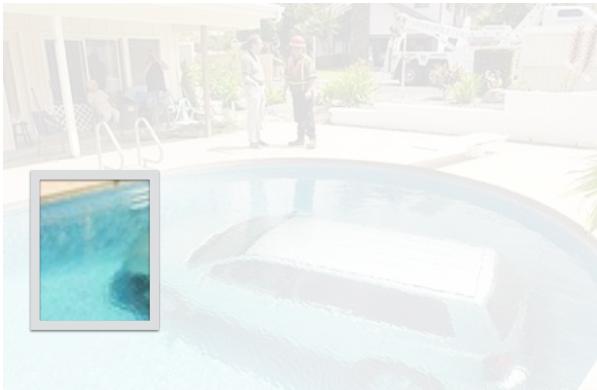
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



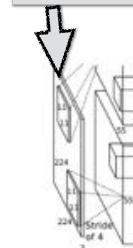
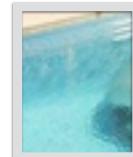
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



Reject
(low confidence)

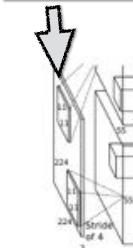
2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



Car

2. Evaluate each proposal

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Object Detection



For more:

R-CNN (<https://github.com/rbgirshick/rcnn>) and Fast R-CNN (<https://github.com/rbgirshick/fast-rcnn>)

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

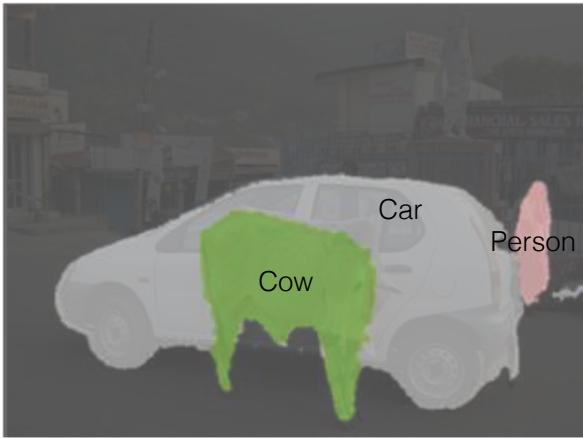


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

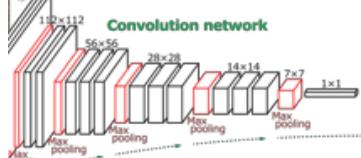


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

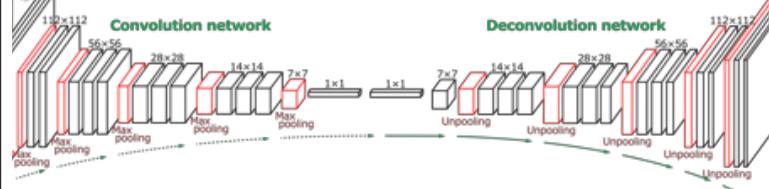


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

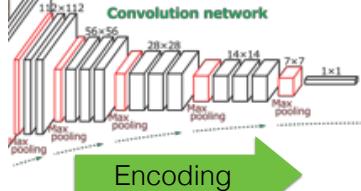


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

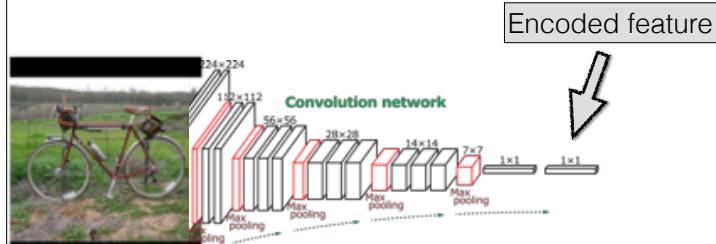


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

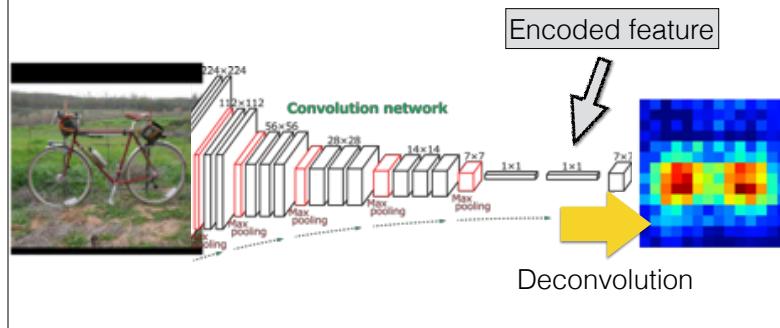


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

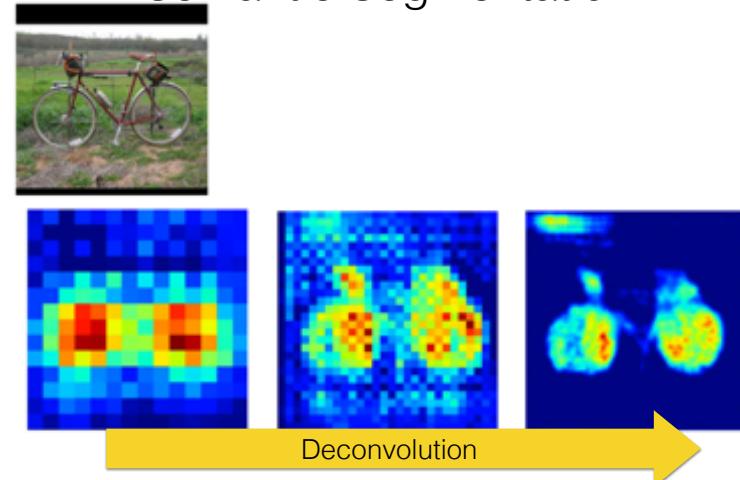


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation

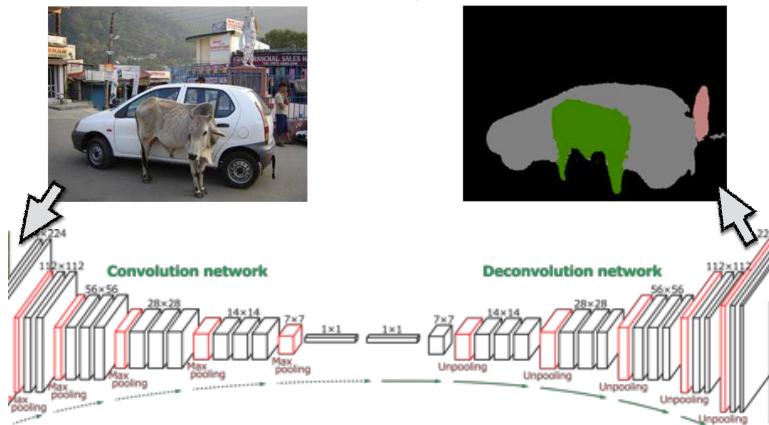


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Semantic Segmentation



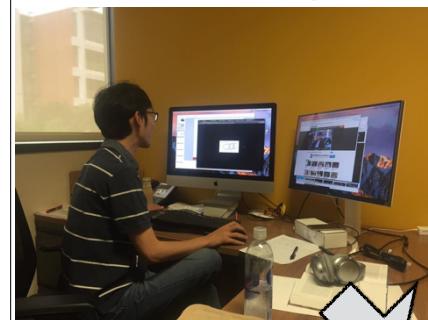
For more: Noh, et.al. Learning Deconvolution Network for Semantic Segmentation. ICCV 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Image captioning

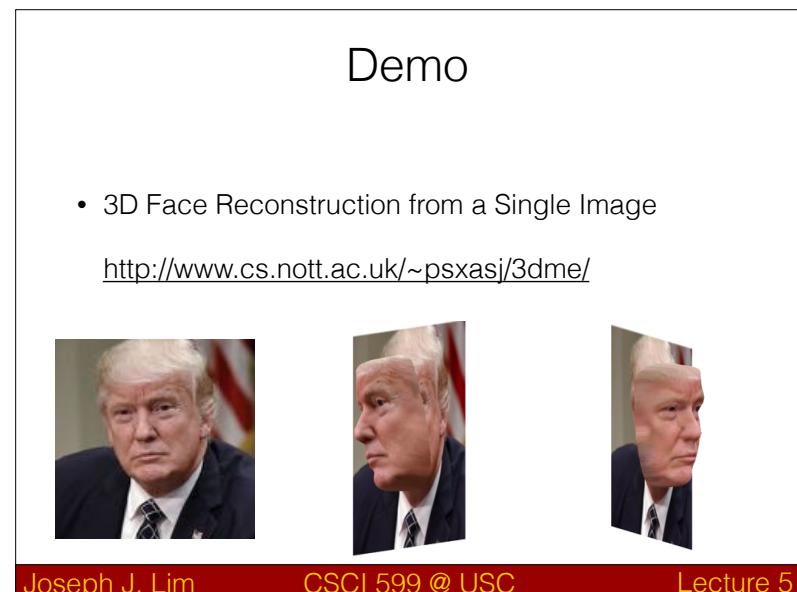
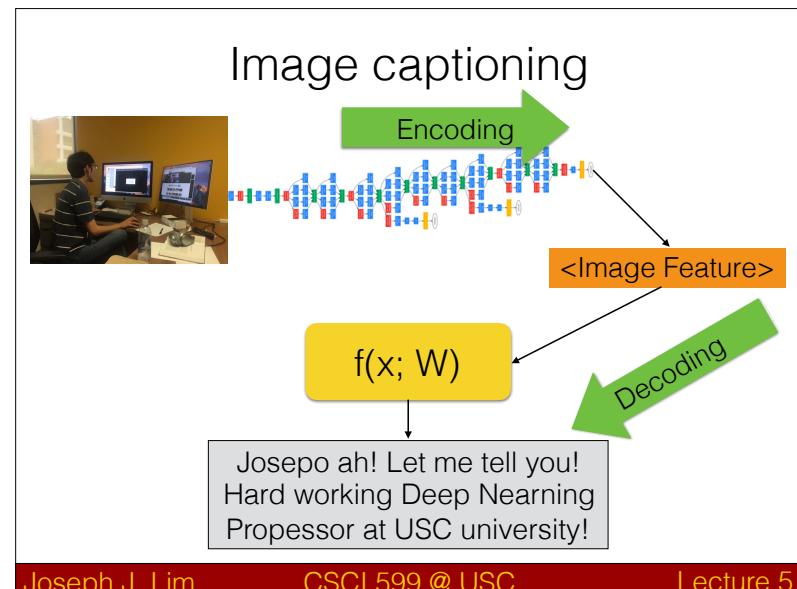
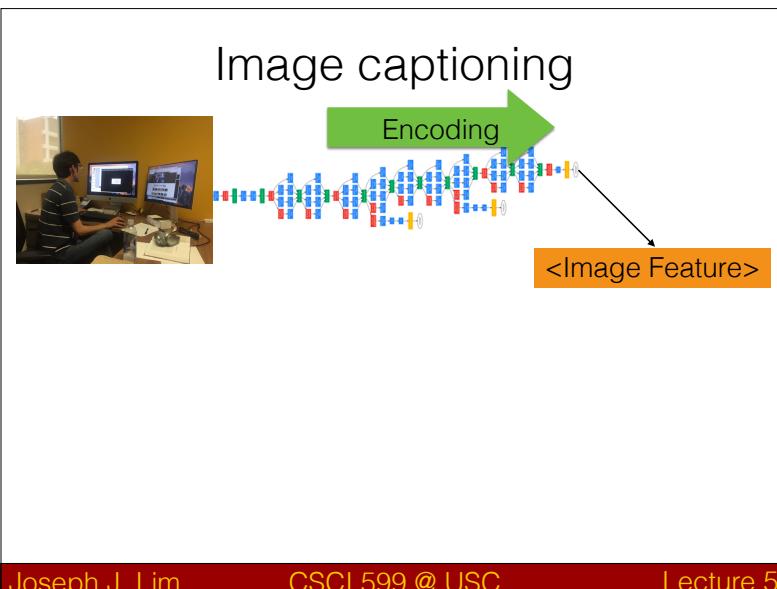


Hao Li (USC prof.) said
Josepo ah! Let me tell you!
Hard working Deep Nearning
Propessor at USC university!

Joseph J. Lim

CSCI 599 @ USC

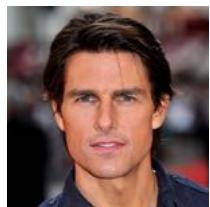
Lecture 5



Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Demo

- 3D Face Reconstruction from a Single Image

<http://www.cs.nott.ac.uk/~psxasj/3dme/>



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Today's agenda

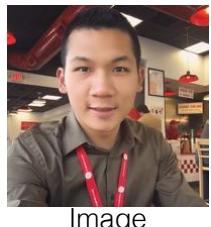
- Recap: CNNs
- Training Neural Networks
- Recurrent Neural Networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

All inputs are just numbers



a	b	c	d	e
f	g	h	i	j
k	l	m	n	o
p	q	r	s	t
u	v	w	x	y



a	b	c	d	e
---	---	---	---	---

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data preprocessing

- Zero-center data
- Normalize data
- PCA Whitening

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening

Joseph J. Lim

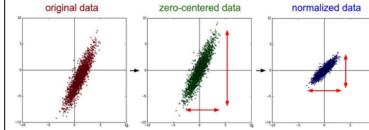
CSCI 599 @ USC

Lecture 5

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening



Activation functions (ReLU) works around **zero**.

If input data is biased toward positive or negative, randomly initialized layers produce biased output.

Figure from Stanford cs231n lecture slides

Joseph J. Lim

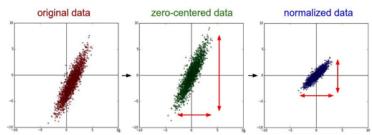
CSCI 599 @ USC

Lecture 5

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening



A feature with small scale has a negligible effect on backprop.

Figure from Stanford cs231n lecture slides

Joseph J. Lim

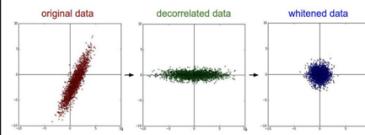
CSCI 599 @ USC

Lecture 5

Data preprocessing

Why do we need preprocessing?

- Zero-center data
- Normalize data
- PCA Whitening



Remove correlation between input features

Figure from Stanford cs231n lecture slides

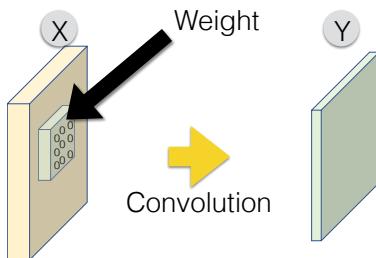
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Weight initialization

If all weights are 0?



Slide credit: CS 231N

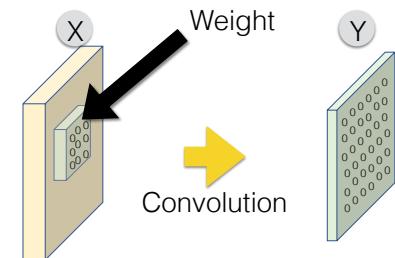
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Weight initialization

If all weights are 0?



Gradient vanishing problem

Output Y will be always 0 and then gradients will go to 0.

Slide credit: CS 231N

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Weight initialization

- Small random numbers
- Xavier
- Xavier / 2 (He et. al.)

Joseph J. Lim

CSCI 599 @ USC

Slide credit: CS 231N

Lecture 5

Small Random Numbers

- Sample from a gaussian distribution with 0 mean and given standard deviation (such as 0.01)

$$W \sim \mathcal{N}(\mu, \sigma^2)$$

- Doesn't work for deeper networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Small Random Numbers

If the weights are too small,
the activation signals shrink quickly

If the weights are too large,
the activation signals explode

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Xavier Initialization

- Set the variance according to # of input neurons

$$W \sim \mathcal{N}(\mu, \sigma^2) \quad \leftarrow \quad Var(W) = \frac{1}{n_{in}}$$

- Consider # output neurons for back-prop

$$Var(W) = \frac{2}{n_{in} + n_{out}}$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

He et. al. Initialization

Random Numbers of dimension (D_{in}, D_{out})

$$\sqrt{D_{in}/2}$$

- Default choice of training a deep network

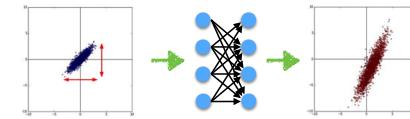
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Normalization

- After each layer, the distribution of activation signals changes (Internal Covariate Shift)



- As a network becomes deeper, distribution shifts more.
- Recall why we need data preprocessing

Figure from Stanford cs231n lecture slides

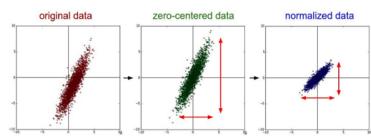
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Normalization

- Normalize the distribution of activations to have zero mean and unit variance



Same as data preprocessing

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

Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Why do we need normalization?
 - Internal Covariate Shift
 - To train deeper networks on larger data
 - Make each dimension Unit Gaussian

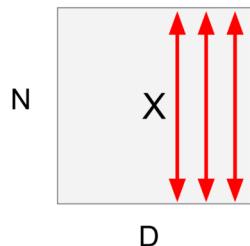
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Compute the empirical mean and variance independently for each dimension and normalize.



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Vanilla Normalization:

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

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Squashing the range:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

- Network can learn:

$$\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \text{E}[x^{(k)}]$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

[Ioffe and Szegedy, 2015]

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Why is it good?
 - Improves gradient flow
 - Allows higher learning rates
 - Reduces strong dependence on initialization
 - Regularization

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- At test time
 - The mean and variance are not computed based on the batch.
 - Empirical mean of activations during training is used.

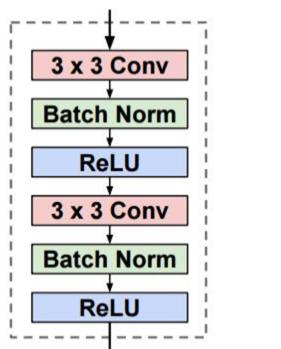
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch normalization

- Usually added after convolutional layers or fully connected layers
- before non-linearities



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Training details

- Learning rate
- Batch size
- Hyperparameter
- Loss curve

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

How **fast** or **slow** we update the model parameters

$$\theta_{t+1} = \theta_t - \eta_t \nabla f(\theta_t)$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

How **fast** or **slow** we update the model parameters

$$\theta_{t+1} = \theta_t - \eta_t \boxed{\nabla f(\theta_t)}$$

Gradients

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

Too large

- Explode

Too small

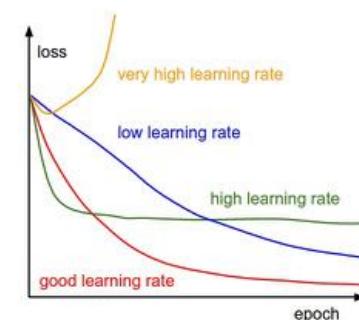
- Get stuck in local minimum

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate



ref: <http://cs231n.github.io/neural-networks-3/#anneal>

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

Too large

- Explode

Too small

- Get stuck in local minimum

Anneal the learning rate

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

- Constant
- Step decay
- Exponential decay
- Time-based decay
- Etc

Joseph J. Lim

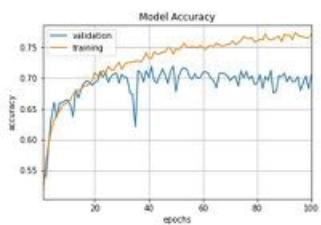
CSCI 599 @ USC

Lecture 5

Learning Rate

Constant

$$lr = lr_0$$



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

Joseph J. Lim

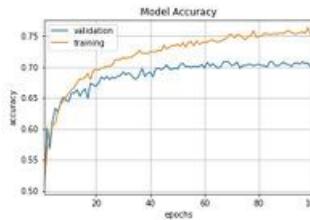
CSCI 599 @ USC

Lecture 5

Learning Rate

Time-based decay

$$lr = \frac{lr_0}{1 + k * iter}$$



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

Joseph J. Lim

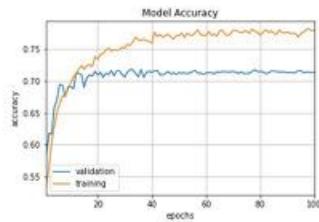
CSCI 599 @ USC

Lecture 5

Learning Rate

Step decay

$$lr = lr_0 * \text{floor}\left(\frac{\text{epoch}}{\text{epoch_drops}}\right)$$



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

Joseph J. Lim

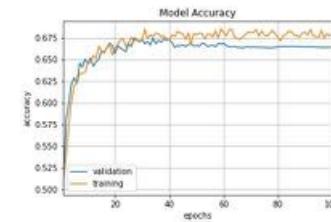
CSCI 599 @ USC

Lecture 5

Learning Rate

Exponential decay

$$lr = \frac{lr_0}{\exp^{k*iter}}$$



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

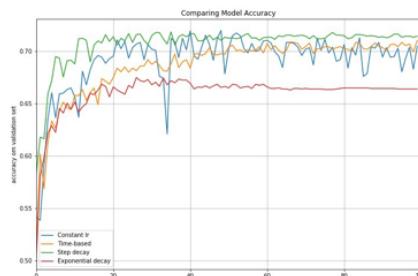
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

Comparisons



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

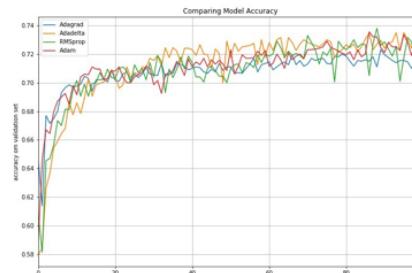
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Learning Rate

Adaptive learning rate methods



Dataset: CIFAR-10

ref: <https://goo.gl/xnrb3N>

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch Size

Number of samples in one forward/backward pass

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch Size

Too large

- Memory inefficient
- Computation inefficient
- Converge to sharp minimizers of the training function

Too small

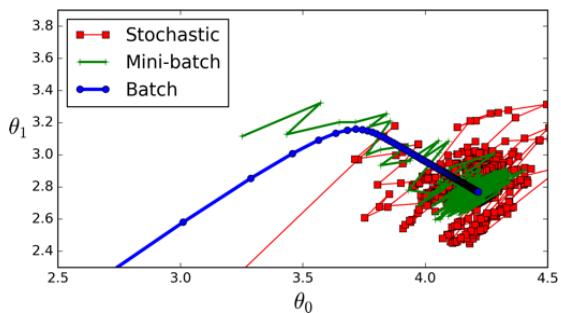
- Less accurate gradients at each iteration
- Converge to flat minimizers of the training function

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Batch Size



Stochastic: batch size is 1

ref: <https://stats.stackexchange.com/questions/153531/what-is-batch-size-in-neural-network>

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Optimizers and Algorithms

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Gradient Descent

- Gradient descent
- Stochastic gradient descent
- Mini-batch gradient descent
- Momentum
- Nesterov accelerated gradient
- Etc

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Gradient Descent

Computes the gradient of the cost function w.r.t. to the parameters θ for **the entire training dataset at each iteration**

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t)$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Stochastic Gradient Descent

Computes the gradient of the cost function w.r.t. to the parameters θ for **each data point at each iteration**

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t; x^{(i)})$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Mini-batch Gradient Descent

Computes the gradient of the cost function w.r.t. to the parameters θ for performs an update for **every mini-batch of n training examples at each iteration**

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} f(\theta_t; x^{(i:i+n)})$$

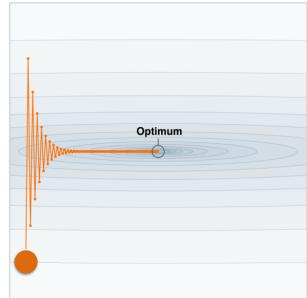
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

GD family updates parameters inefficiency



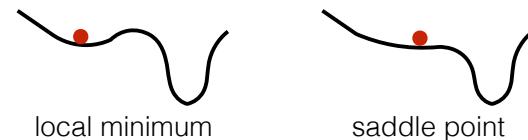
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

GD family gets stuck at local minimums and saddle points easily



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

Consider the history:

Adding a fraction γ of **the update vector of the past time step** to the current update vector

$$\nu_t = \gamma \nu_{t-1} + \eta \nabla_{\theta_t} f(\theta_t)$$

$$\theta_{t+1} = \theta_t - \nu_t$$

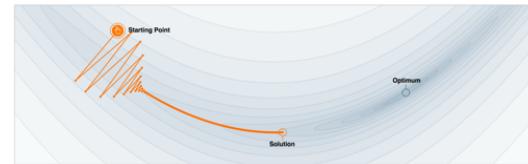
Joseph J. Lim

CSCI 599 @ USC

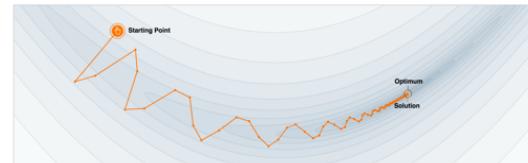
Lecture 5

Momentum

Without momentum



With momentum



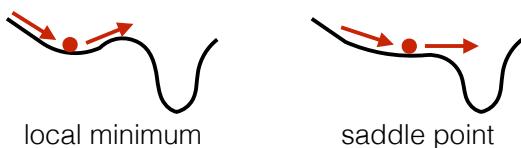
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Momentum

GD family gets stuck at local minimums and saddle points easily



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Nesterov accelerated gradient

Calculate the gradient not w.r.t. to our current parameters but w.r.t. **the approximate future position of our parameters**:

$$\nu_t = \gamma\nu_{t-1} + \eta\nabla_{\theta_t} f(\theta_t - \gamma\nu_{t-1})$$

$$\theta_{t+1} = \theta_t - \nu_t$$

Joseph J. Lim

CSCI 599 @ USC

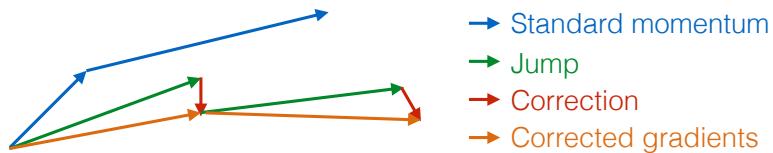
Lecture 5

Nesterov accelerated gradient

Calculate the gradient not w.r.t. to our current parameters but w.r.t. **the approximate future position of our parameters**:

$$\nu_t = \gamma\nu_{t-1} + \eta\nabla_{\theta_t} f(\theta_t - \gamma\nu_{t-1})$$

$$\theta_{t+1} = \theta_t - \nu_t$$



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Regularization

- Extra terms
- Dropout
- Data augmentation

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Regularization

$$L = L_{original} + \lambda R(\theta)$$

L2 norm regularizer $R(\theta) = \sum_i \sum_j \theta_{i,j}^2$

L1 norm regularizer $R(\theta) = \sum_i \sum_j |\theta_{i,j}|$

L1+L2 norm regularizer $R(\theta) = \sum_i \sum_j |\theta_{i,j}| + \alpha \theta_{i,j}^2$

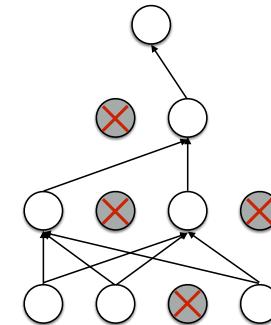
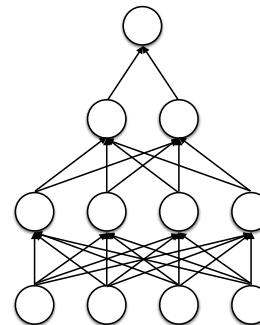
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Dropout

At each iteration, randomly shutdown some neurons with a certain probability (usually 0.5)



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Dropout

Intuition

- Forces the network to have a redundant representation
- Training a bunches of small models sharing parameters
- Each binary mask is one model

Joseph J. Lim

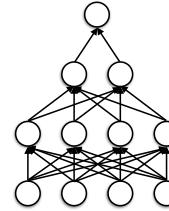
CSCI 599 @ USC

Lecture 5

Dropout

Testing phase

- All neurons are active always
- Scale the activations so that for each neuron
 - Testing output = expected output at training time



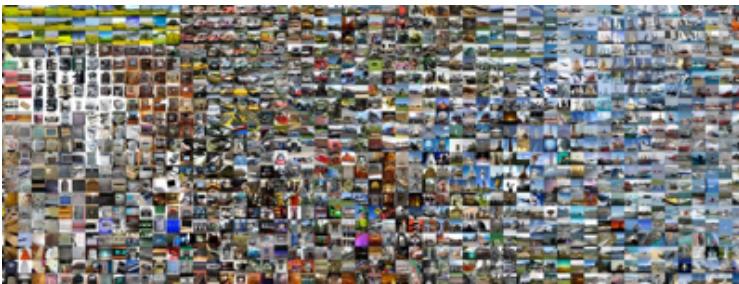
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Deep networks need to be trained on a huge number of training images



What if... we don't have enough amount of data?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Some tricks for image data augmentation

- Flip
- Scale
- Rotate
- Crop
- Color jitter
- etc

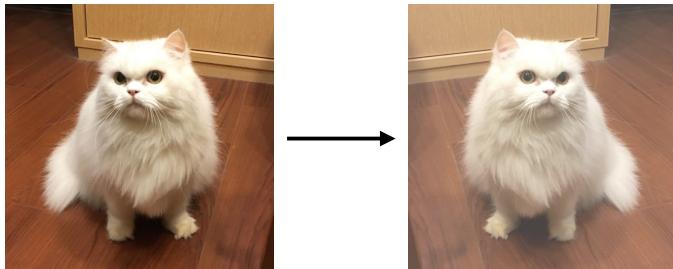
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Flip



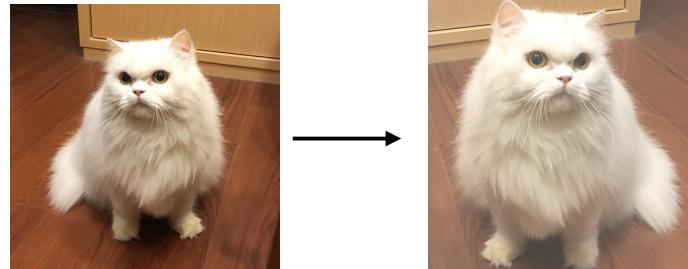
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Scale



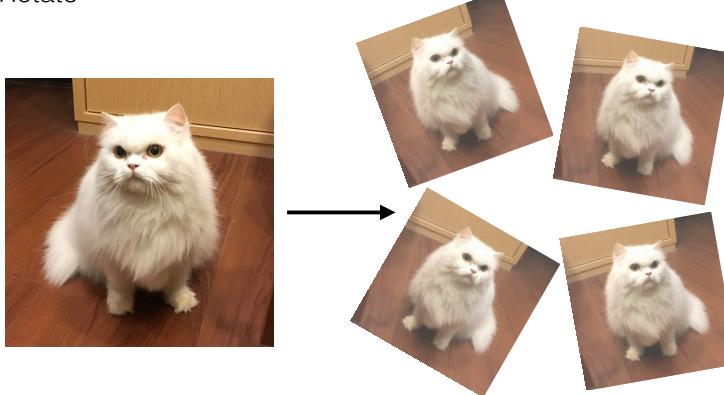
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Rotate



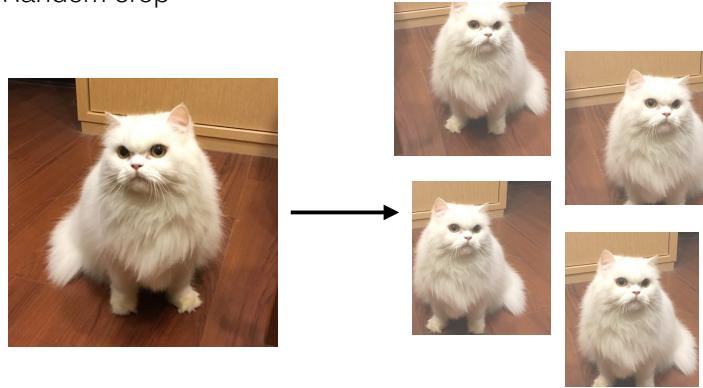
Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation

Random crop

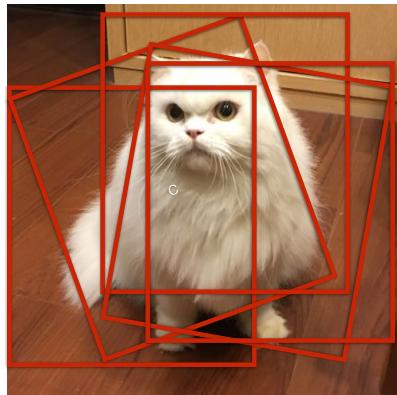


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Data Augmentation



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Summary

Deep learning is all about getting everything right...

- Data
- Model
 - Architecture
- Optimizer
- Hyperparameters
- Training tricks
 - Initialization, normalization, regularization, etc

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Questions?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Todo

- Find Teammates for your Project
- Assignment 1 will be out next week!

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Questions?

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Today's agenda

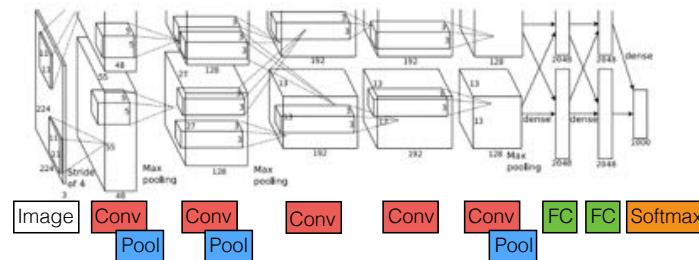
- Recap: CNNs
- Training Neural Networks
- Recurrent Neural Networks

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Recap: CNN (AlexNet)



A Krizhevsky, et. al. ImageNet Classification with Deep Convolutional Neural Networks. NIPS 2012.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recap: CNNs

- Convolutional layer
- Pool layer
- Fully connected layer
- Softmax layer
- Activation layer (e.g. ReLU)



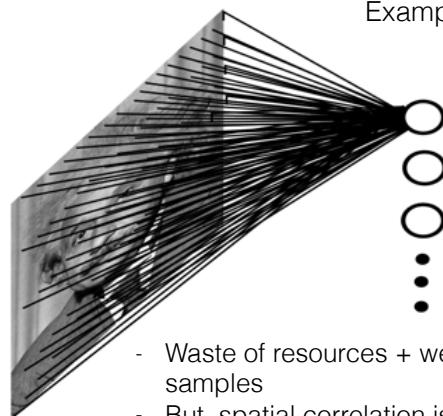
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recap: Fully connected layer (FC)

Example: 200x200 image
40k hidden units
~2B parameters!!!



- Waste of resources + we have not enough training samples
- But, spatial correlation is local!

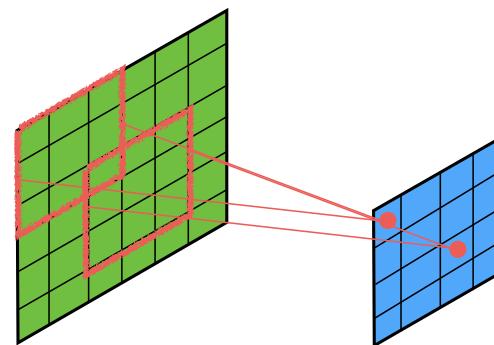
Slide credit: Marc'Aurelio Ranzato

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recap: Convolutional layer



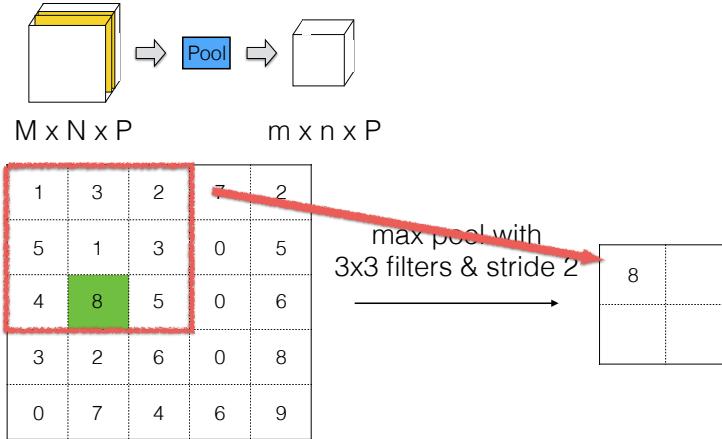
Inspired by Marc'Aurelio Ranzato's slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recap: Max pooling layer



Joseph J. Lim

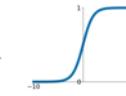
CSCI 599 @ USC

Lecture 8

Recap: Activation functions

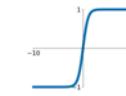
Sigmoid

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



tanh

$$\tanh(x)$$



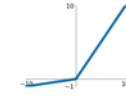
ReLU

$$\max(0, x)$$



Leaky ReLU

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



Maxout

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

ELU

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

Figure from Stanford cs231n lecture slides

Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Recap: VGGNet



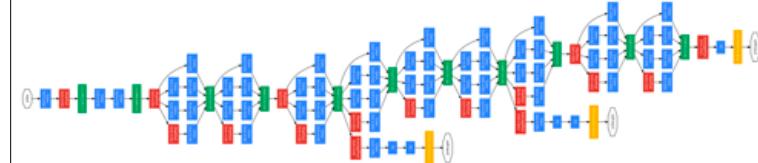
Simonyan and Zisserman. Very Deep Convolutional Networks for Large-scale Image Recognition. ICLR 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recap: GoogLeNet



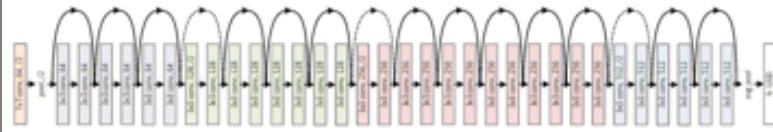
Szegedy, et. al. Going Deeper with Convolution. CVPR 2015.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recap: ResNet



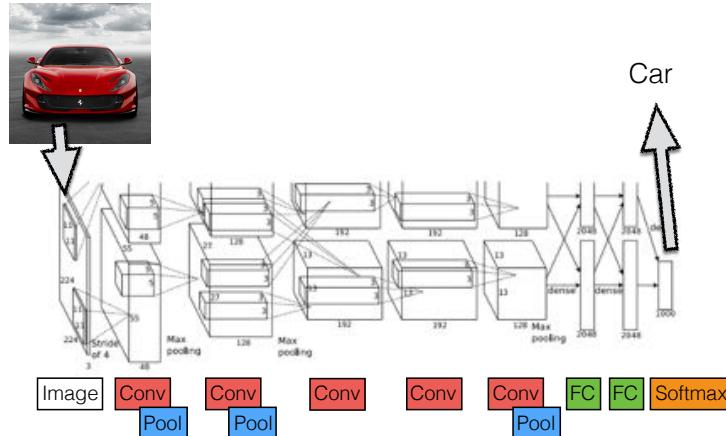
He, et. al. Deep **Residual** Learning for Image Recognition. CVPR 2016.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recap: Image Classification

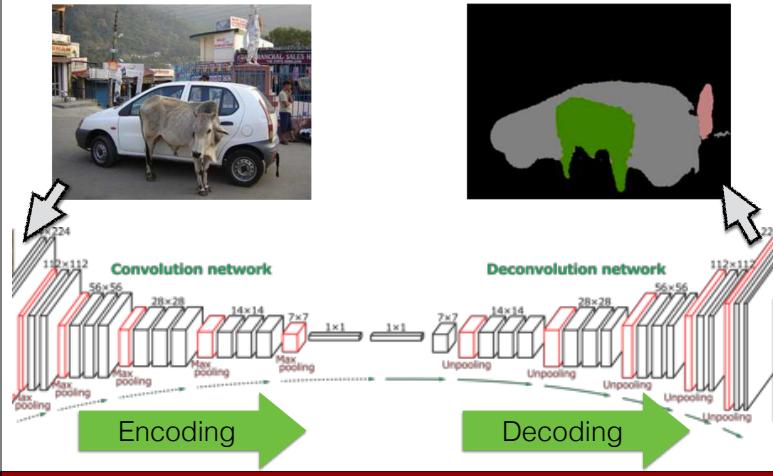


Joseph J. Lim

CSCI 599 @ USC

Lecture 5

Recap: Semantic Segmentation



Joseph J. Lim

CSCI 599 @ USC

Lecture 5

What about speech recognition?



Slide credit: Dhruv Batra
Image credit: Alex Graves and Kevin Gimpel

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

What about speech recognition?



Joseph J. Lim

CSCI 599 @ USC

Lecture 8

What about speech recognition?

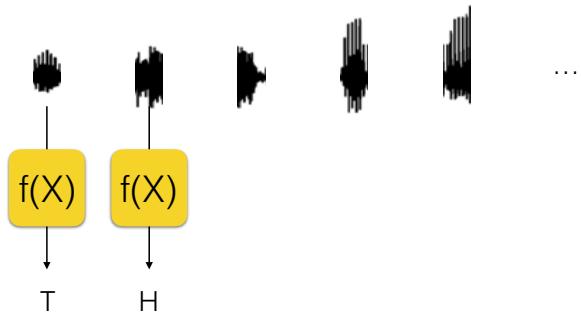


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

What about speech recognition?

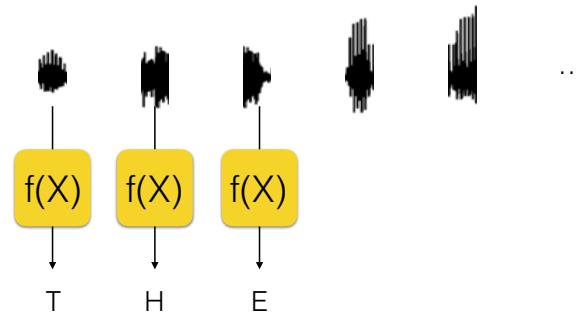


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

What about speech recognition?

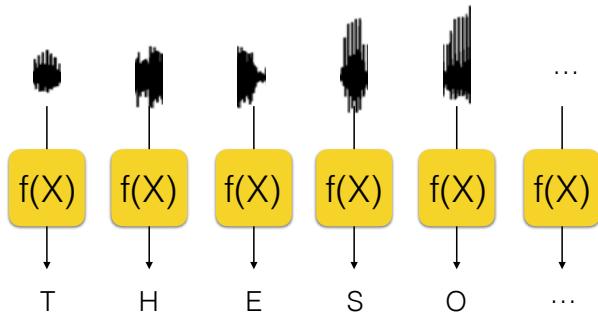


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

What about speech recognition?

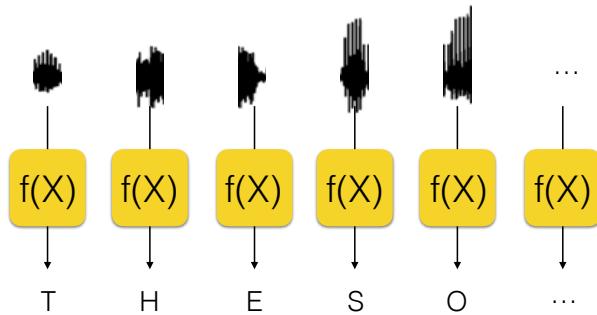


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

What about speech recognition?



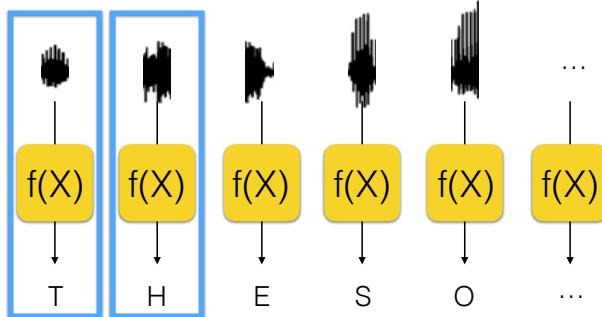
Any limitation on this approach?

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

A limitation of “normal” CNNs

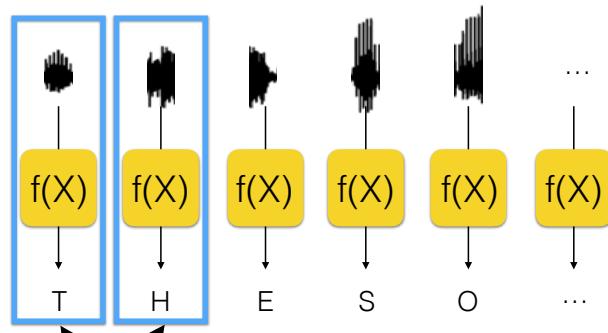


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

A limitation of “normal” CNNs



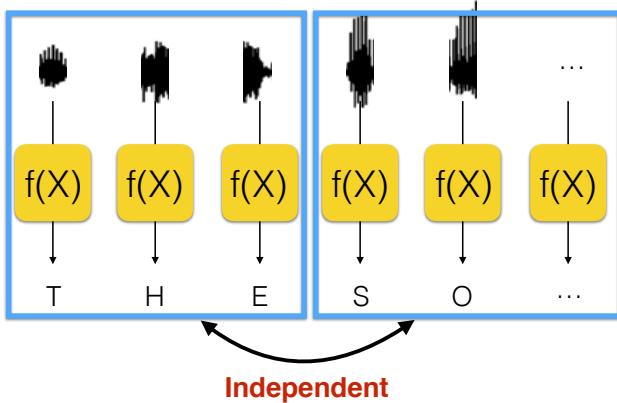
Independent

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

A limitation of “normal” CNNs



Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Remember Hidden Markov Model?

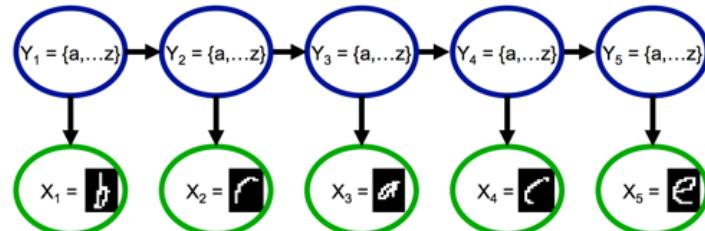


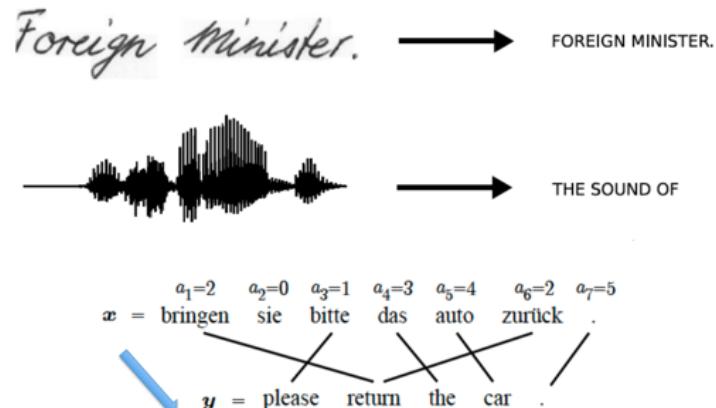
Figure credit: Carlos Guestrin

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Sequential data is everywhere



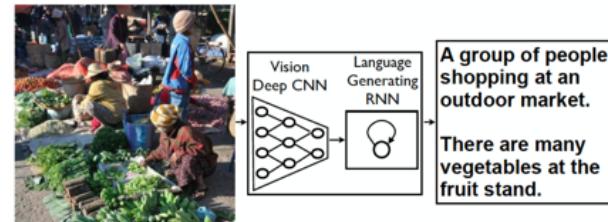
Slide credit: Dhruv Batra
Image credit: Alex Graves and Kevin Gimpel

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Sequential data is everywhere



John has a dog . →

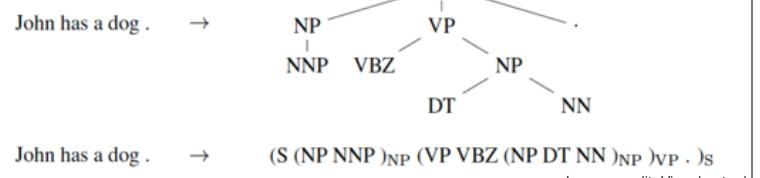


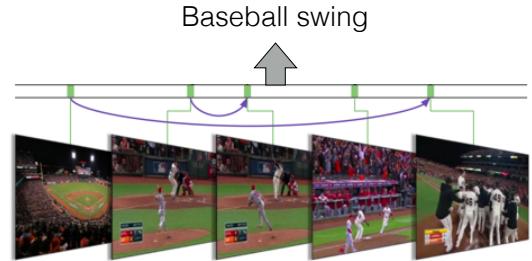
Image credit: Vinyals et. al.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Sequential data is everywhere



Yeung, et. al. End-to-end Learning of Action Detection
from Frame Glimpses in Videos. CVPR 2017.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

How do we model sequences?



one-to-one
("normal" CNNs)

Inspired from Stanford CS231n

Joseph J. Lim

CSCI 599 @ USC

Lecture 8



one-to-one
("normal" CNNs)

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



one-to-one
("normal" CNNs)

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



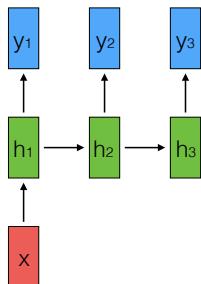
one-to-one
("normal" CNNs)

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



one-to-many
(e.g. image to sentence)

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



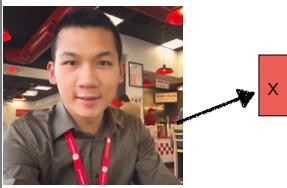
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

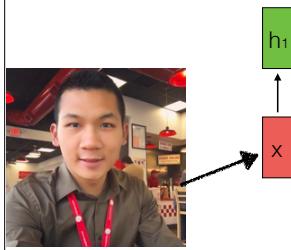


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

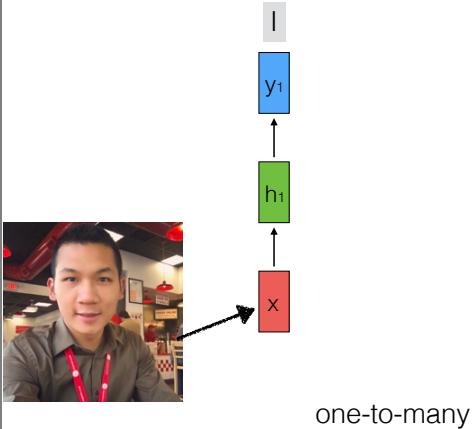


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

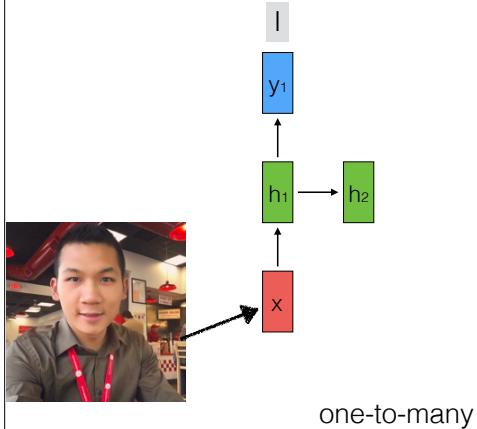


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

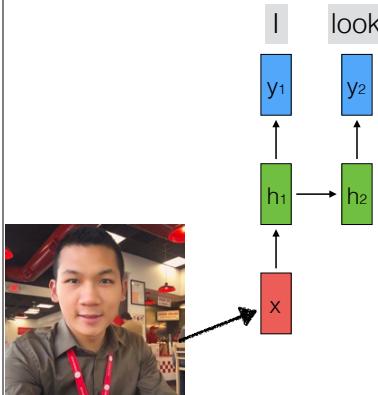


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



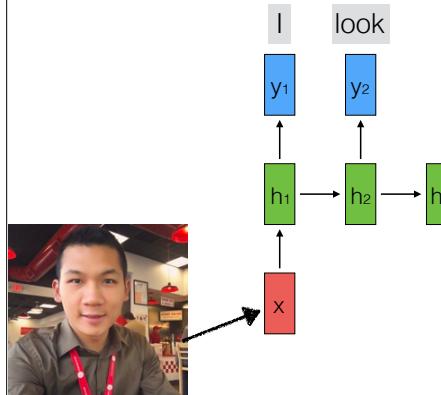
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



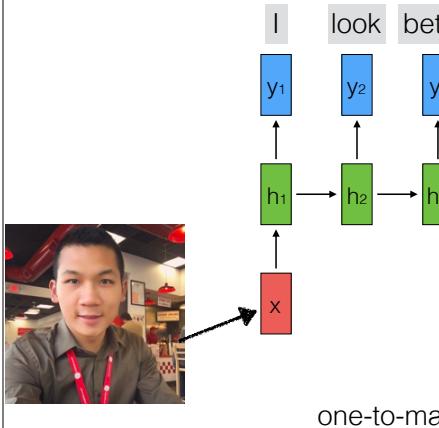
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



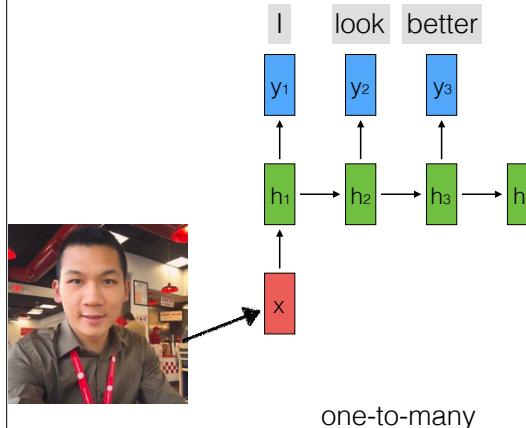
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



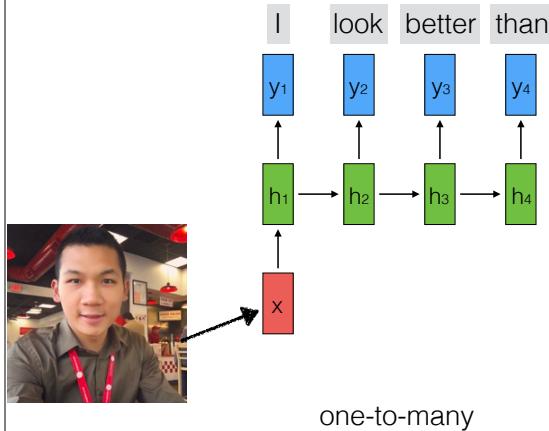
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



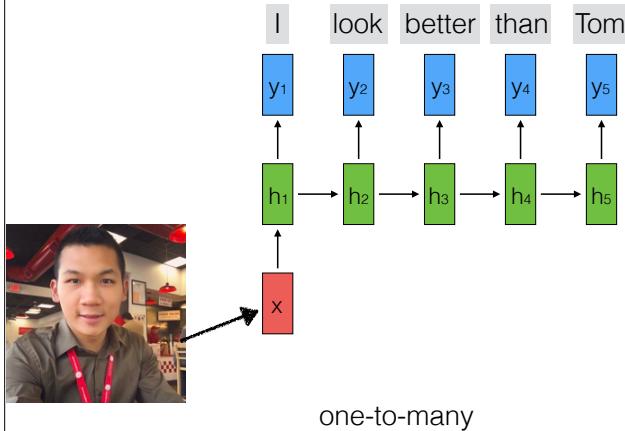
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



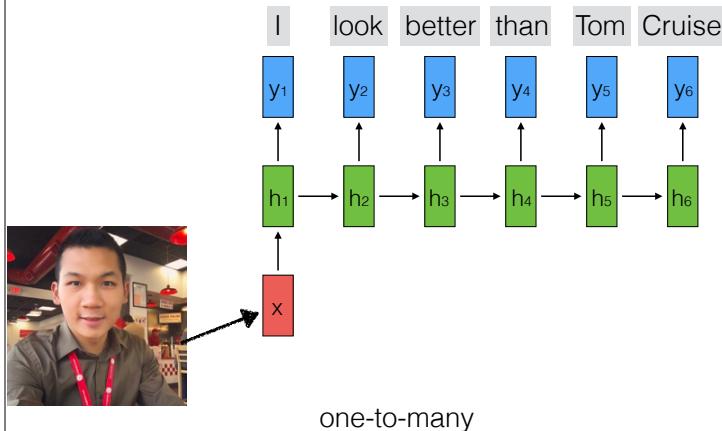
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



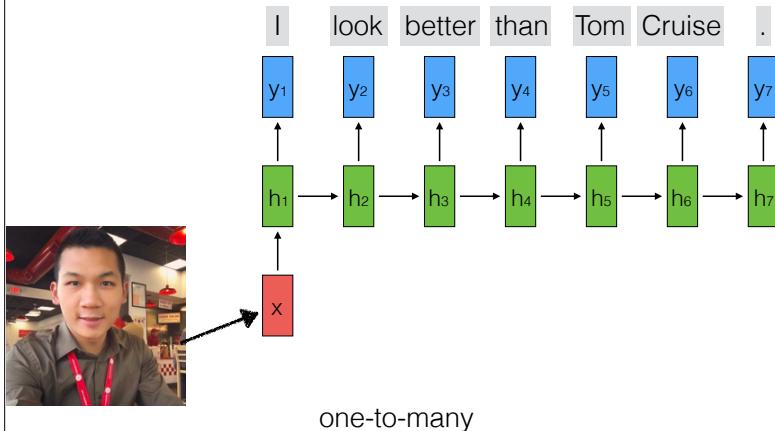
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

I look better than Tom Cruise .



This is why a training data is important!

The model failed...



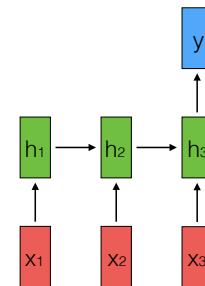
one-to-many

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



many-to-one

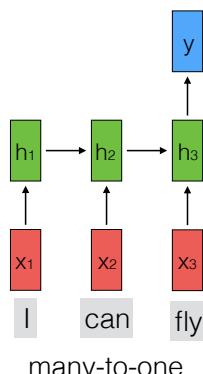
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

Wrong!



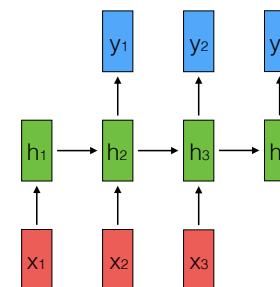
many-to-one

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



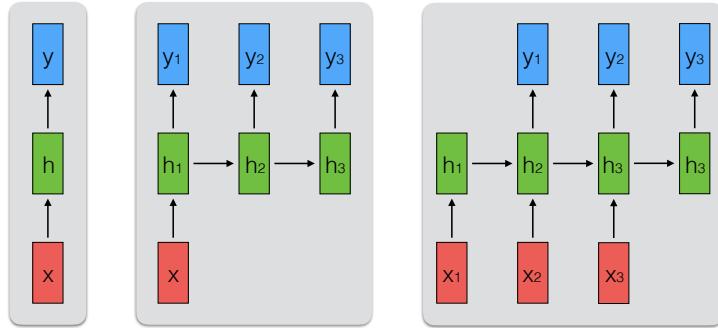
many-to-many
(e.g. machine translation)

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Different types of RNNs



Many more variations!

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

It can get tricky...

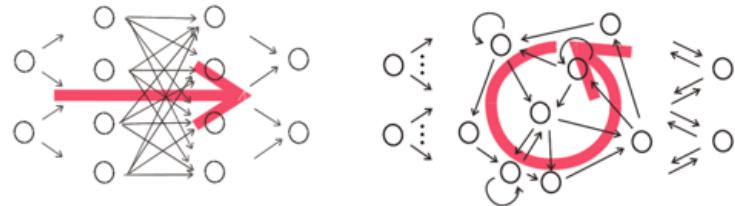


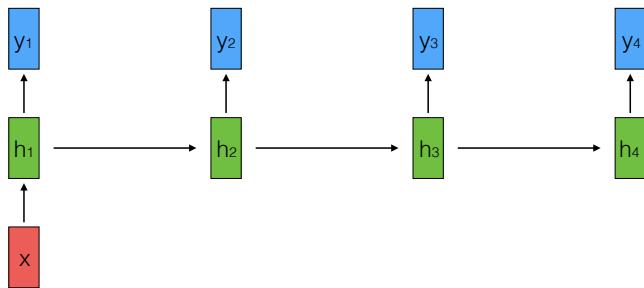
Image credit: Herbert Jaeger

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Let's zoom in!

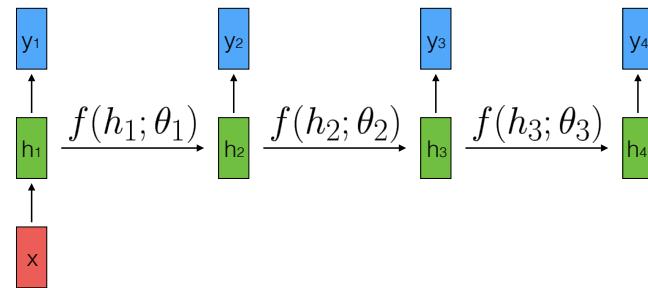


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Let's zoom in!

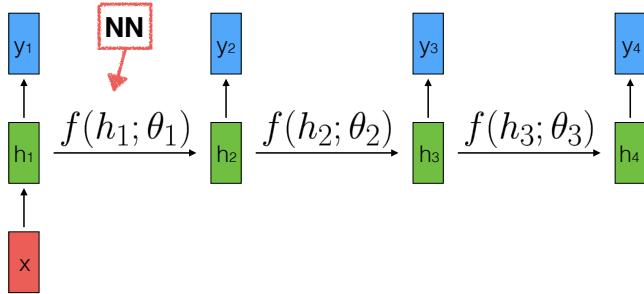


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Let's zoom in!

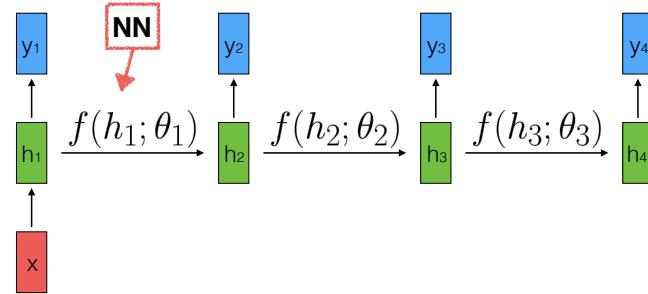


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Let's zoom in!



Any issue with this?

How do we model sequences?

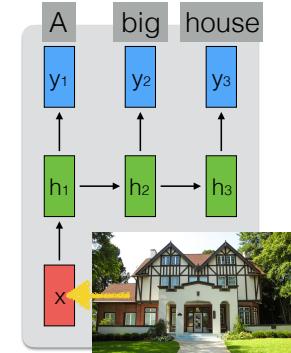


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

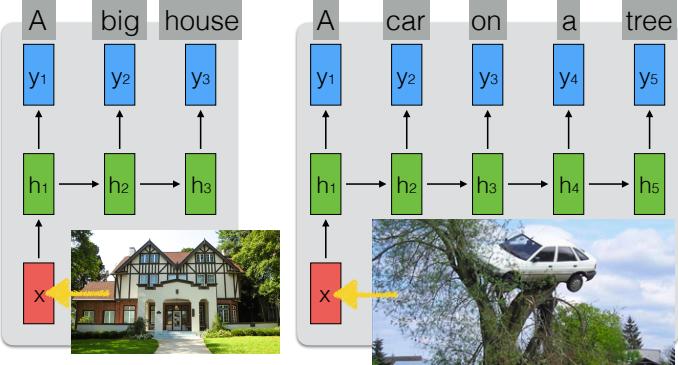


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?

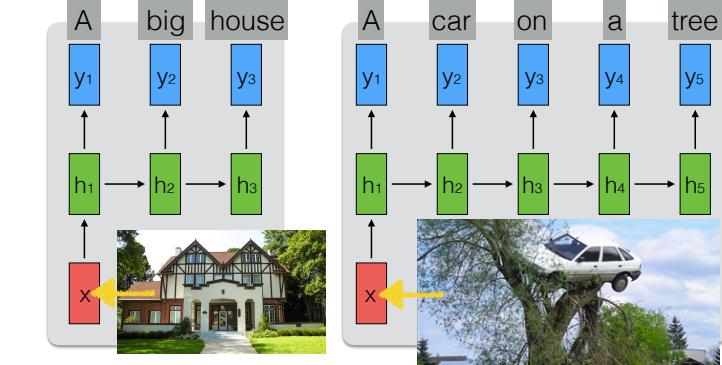


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

How do we model sequences?



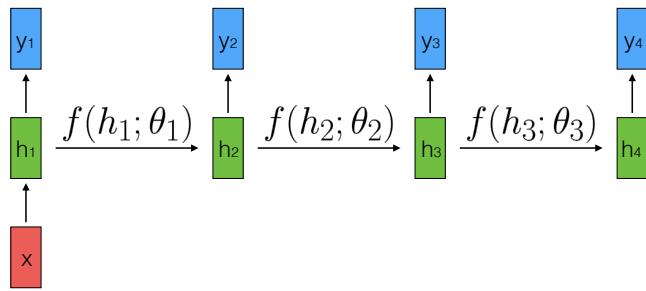
How to cover a variable input / output length?

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recurrent Neural Network

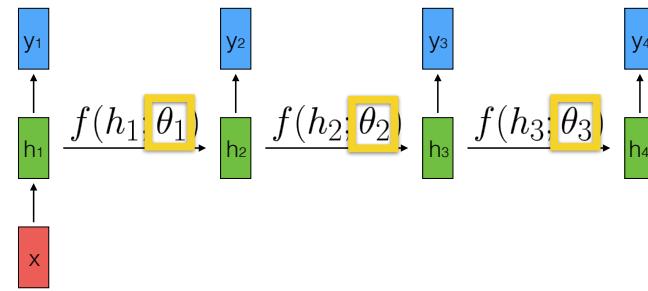


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recurrent Neural Network

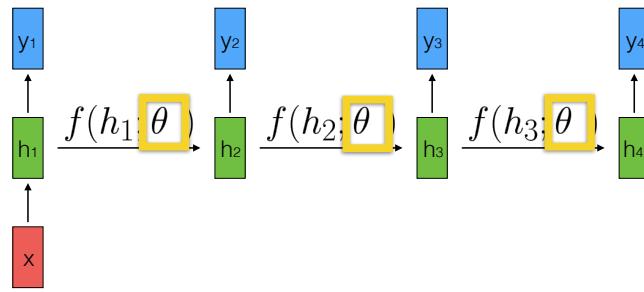


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recurrent Neural Network

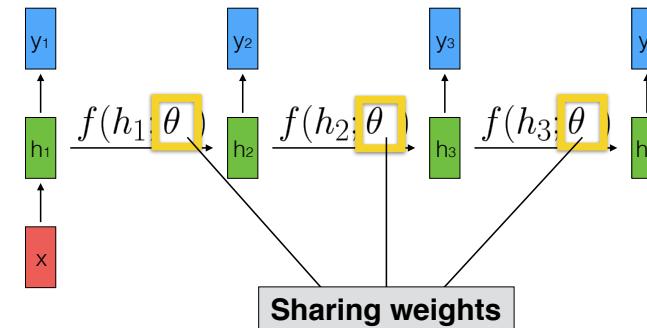


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recurrent Neural Network

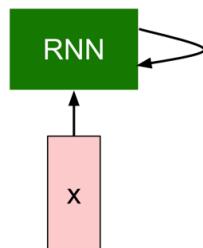


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recurrent Neural Network

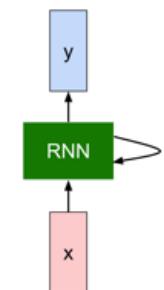


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recurrent Neural Network



Slide credit: Stanford CS231n

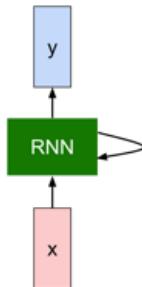
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Recurrent Neural Network

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



Slide credit: Stanford CS231n

Joseph J. Lim

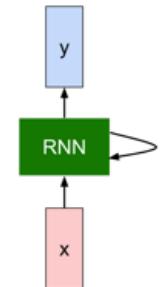
CSCI 599 @ USC

Lecture 8

Recurrent Neural Network

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

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



Slide credit: Stanford CS231n

Joseph J. Lim

CSCI 599 @ USC

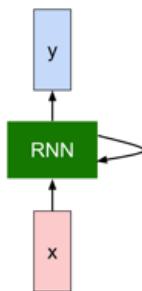
Lecture 8

Recurrent Neural Network

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

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

old state |
input vector at some time step



Slide credit: Stanford CS231n

Joseph J. Lim

CSCI 599 @ USC

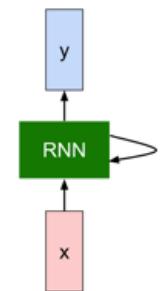
Lecture 8

Recurrent Neural Network

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

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

old state |
input vector at some time step



Slide credit: Stanford CS231n

Joseph J. Lim

CSCI 599 @ USC

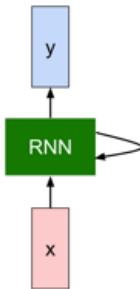
Lecture 8

Recurrent Neural Network

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

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

some function
with parameters W



Joseph J. Lim

CSCI 599 @ USC

Lecture 8

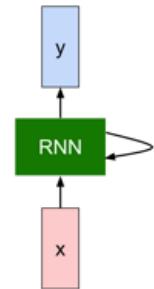
Slide credit: Stanford CS231n

Recurrent Neural Network

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

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

new state
old state
some function
with parameters W



Slide credit: Stanford CS231n

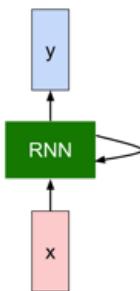
Recurrent Neural Network

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

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

new state
old state
some function
with parameters W

Note: the same function with the same parameters



Joseph J. Lim

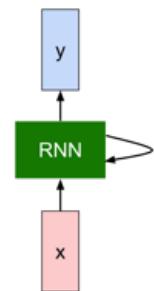
CSCI 599 @ USC

Lecture 8

Slide credit: Stanford CS231n

Recurrent Neural Network

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

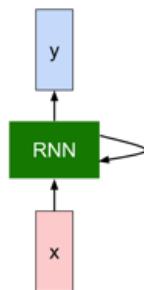


Slide credit: Stanford CS231n

Joseph J. Lim CSCI 599 @ USC Lecture 8

(Vanilla) Recurrent Neural Network

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



Slide credit: Stanford CS231n

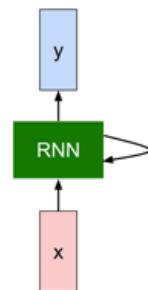
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

(Vanilla) Recurrent Neural Network

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



The state has a single “hidden” vector **h**

Slide credit: Stanford CS231n

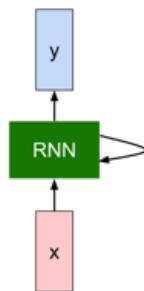
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

(Vanilla) Recurrent Neural Network

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



The state has a single “hidden” vector **h**

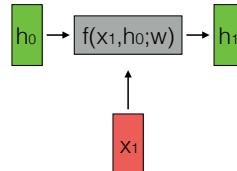
Slide credit: Stanford CS231n

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Computational Graph for RNN

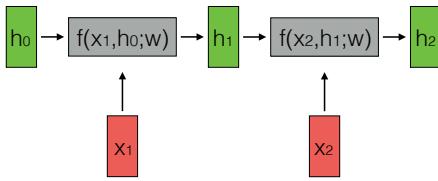


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Computational Graph for RNN

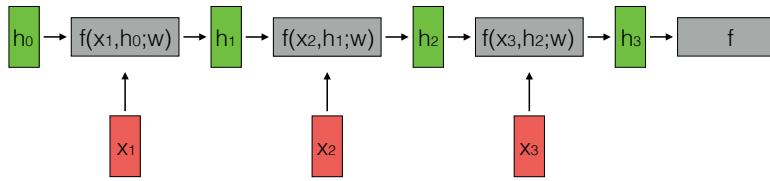


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Computational Graph for RNN

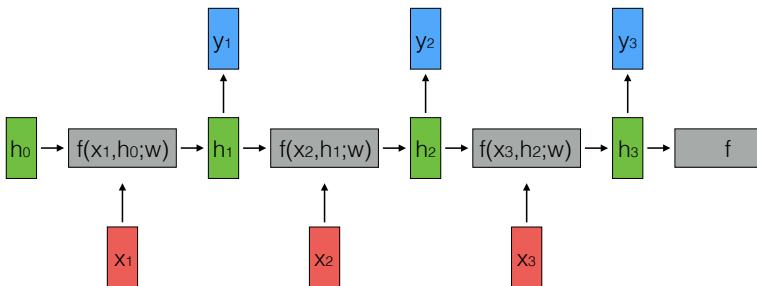


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Computational Graph for RNN

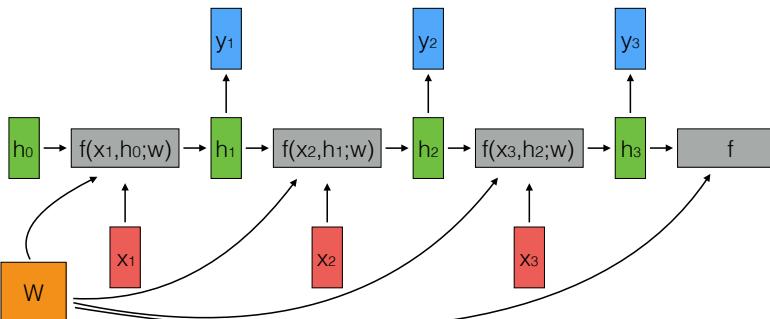


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Computational Graph for RNN

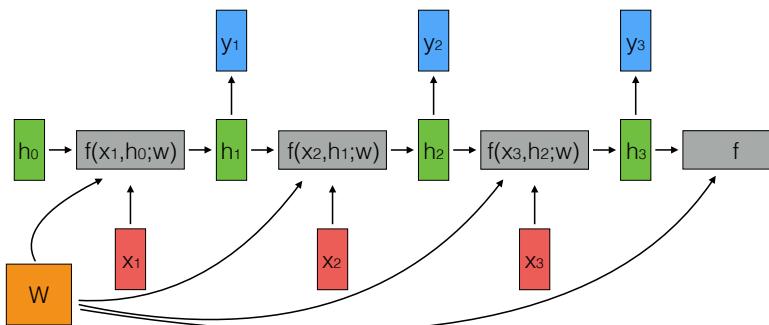


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Computational Graph for RNN

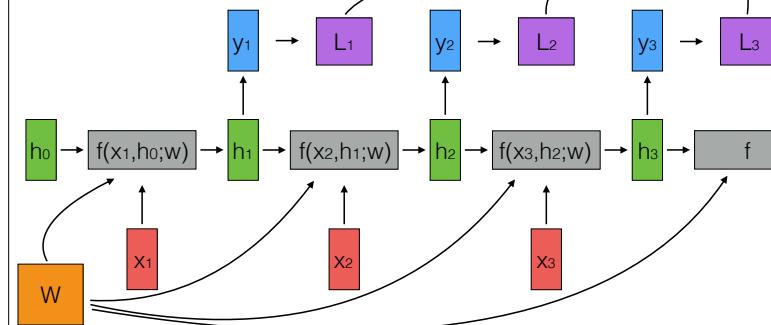


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Computational Graph for RNN



Joseph J. Lim

CSCI 599 @ USC

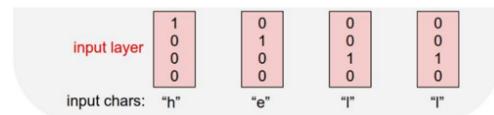
Lecture 8

An Example

Character-level Language Model

Vocabulary:
[h, e, l, o]

Sample Training
Sequence: "hello"



Credits: Stanford CS231n

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

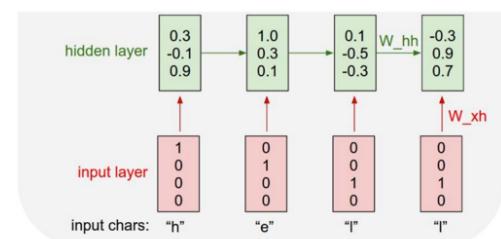
An Example

Character-level Language Model

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

Vocabulary:
[h, e, l, o]

Sample Training
Sequence: "hello"



Credits: Stanford CS231n

Joseph J. Lim

CSCI 599 @ USC

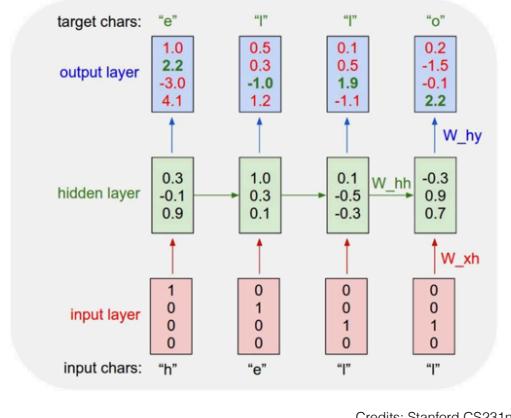
Lecture 8

An Example

Character-level Language Model

Vocabulary:
[h, e, l, o]

Sample Training
Sequence: "hello"



Joseph J. Lim

CSCI 599 @ USC

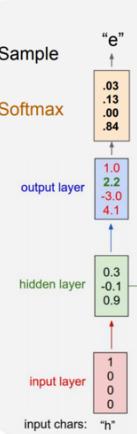
Lecture 8

An Example

Character-level Language Model

Vocabulary:
[h, e, l, o]

At **test-time** sample
characters one at a
time, and feed it
back to the model



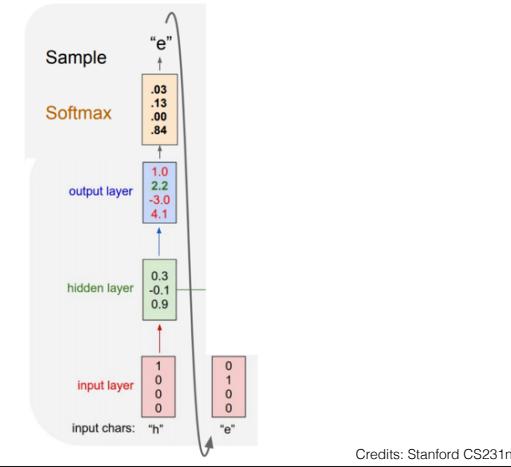
Joseph J. Lim CSCI 599 @ USC Lecture 8

An Example

Character-level Language Model

Vocabulary:
[h, e, l, o]

At **test-time** sample
characters one at a
time, and feed it
back to the model



Joseph J. Lim

CSCI 599 @ USC

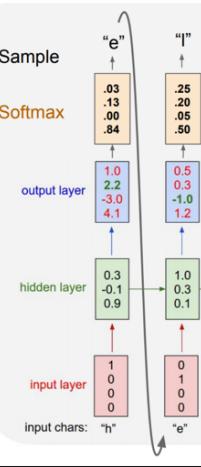
Lecture 8

An Example

Character-level Language Model

Vocabulary:
[h, e, l, o]

At **test-time** sample
characters one at a
time, and feed it
back to the model



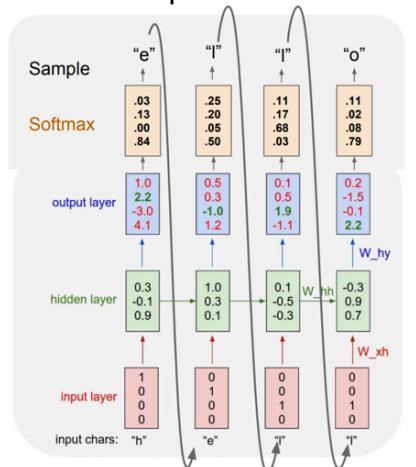
Joseph J. Lim CSCI 599 @ USC Lecture 8

An Example

Character-level Language Model

Vocabulary:
[h, e, l, o]

At **test-time** sample characters one at a time, and feed it back to the model

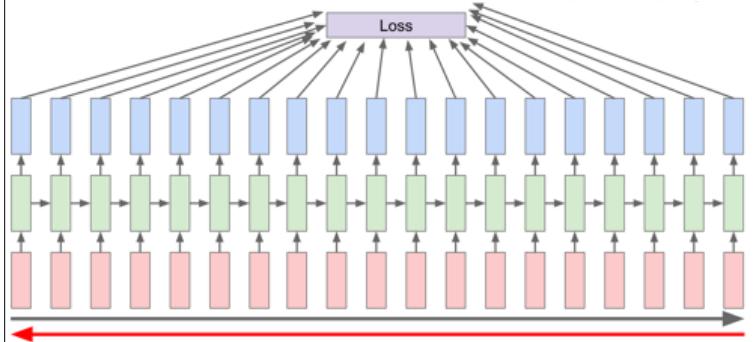


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Backpropagation for RNN

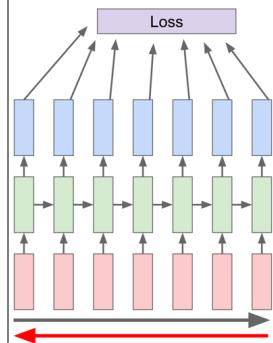


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Truncated Backpropagation



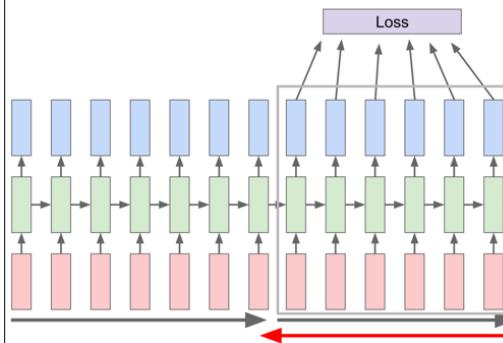
Run forward and backward through chunks of the sequence instead of whole sequence

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Truncated Backpropagation



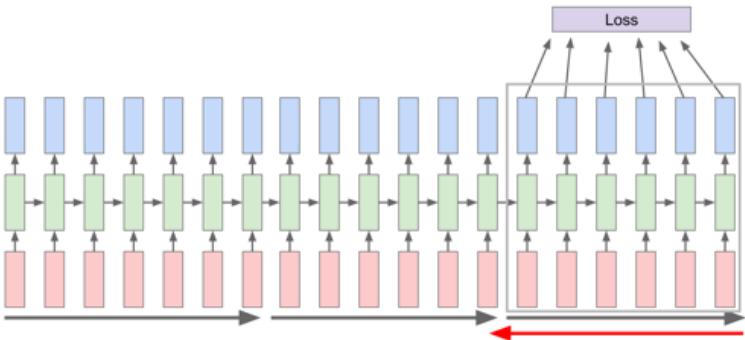
Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Truncated Backpropagation



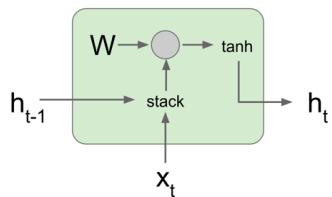
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

RNNs aren't as stable as CNNs

Vanilla RNN Gradient Flow



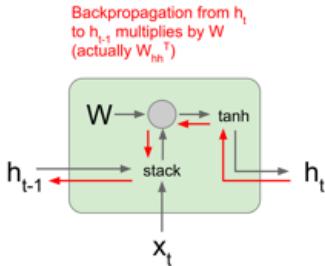
$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Vanilla RNN Gradient Flow



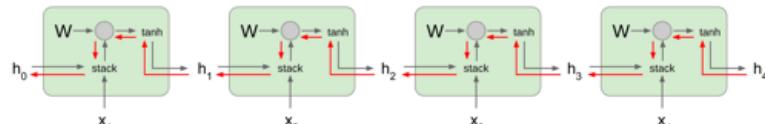
$$\begin{aligned} h_t &= \tanh(W_{hh}h_{t-1} + W_{xh}x_t) \\ &= \tanh \left(\begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \\ &= \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right) \end{aligned}$$

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Vanilla RNN Gradient Flow



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

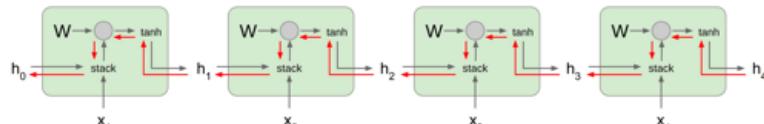
Largest singular value < 1 :
Vanishing gradients

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Vanilla RNN Gradient Flow



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

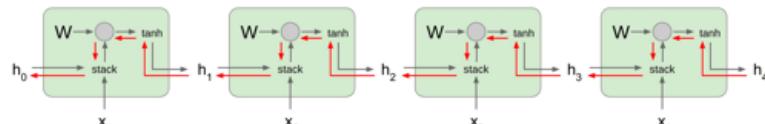
Gradient clipping: scale gradient if its norm is too big

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Vanilla RNN Gradient Flow



Computing gradient of h_0 involves many factors of W (and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

Change RNN architecture

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Long Short Term Memory (LSTM)

- Learn **long-term** dependencies
 - Remember information for long periods of time is practically their default behavior

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997
Colah's blog: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Long Short Term Memory (LSTM)

Vanilla RNN

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

LSTM

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

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

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

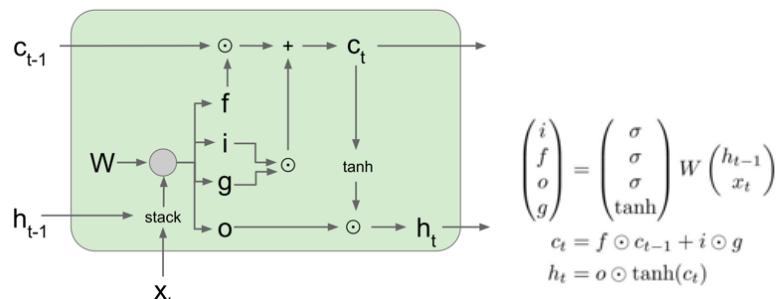
Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Long Short Term Memory (LSTM)



Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

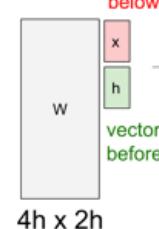
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Long Short Term Memory (LSTM)

vector from below (x)



- f: Forget gate, Whether to erase cell
- i: Input gate, whether to write to cell
- g: Gate gate (?), How much to write to cell
- o: Output gate, How much to reveal cell

Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

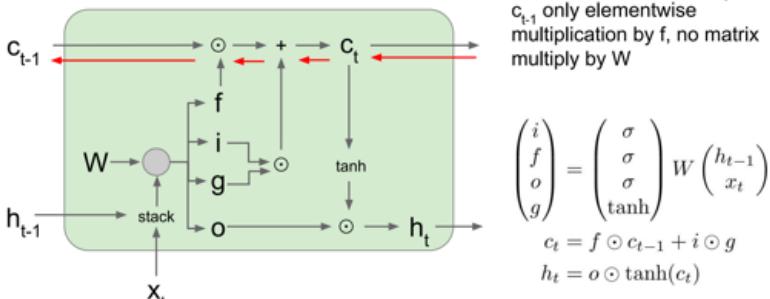
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Long Short Term Memory (LSTM)

Backpropagation from c_t to c_{t-1} only elementwise multiplication by f , no matrix multiply by W



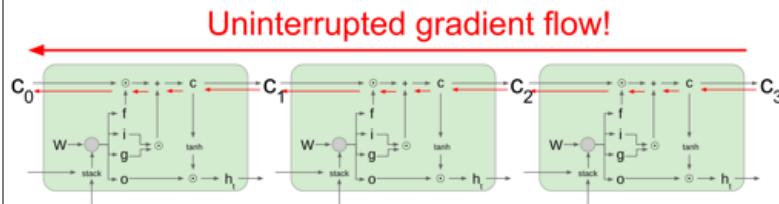
Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Long Short Term Memory (LSTM)



Hochreiter and Schmidhuber, "Long Short Term Memory", Neural Computation 1997

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Other variants

[An Empirical Exploration of Recurrent Network Architectures, Jozefowicz et al., 2015]

GRU [Learning phrase representations using rnn encoder-decoder for statistical machine translation, Cho et al. 2014]

$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$

$$\tilde{h}_t = \tanh(W_{xh}x_t + W_{hh}(r_t \odot h_{t-1}) + b_h)$$

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

MUT1:

$$\begin{aligned} z &= \sigma(W_{xz}x_t + W_{zh}h_t + b_z) \\ r &= \sigma(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{zh}(r \odot h_t) + W_{sh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

MUT2:

$$\begin{aligned} z &= \sigma(W_{xz}x_t + W_{zh}h_t + b_z) \\ r &= \sigma(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{zh}(r \odot h_t) + W_{sh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

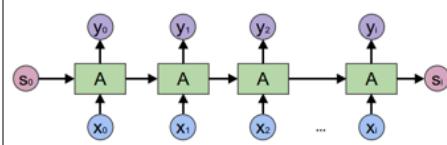
MUT3:

$$\begin{aligned} z &= \sigma(W_{xz}x_t + W_{zh}\tanh(h_t) + b_z) \\ r &= \sigma(W_{xr}x_t + W_{hr}h_t + b_r) \\ h_{t+1} &= \tanh(W_{zh}(r \odot h_t) + W_{sh}x_t + b_h) \odot z \\ &+ h_t \odot (1 - z) \end{aligned}$$

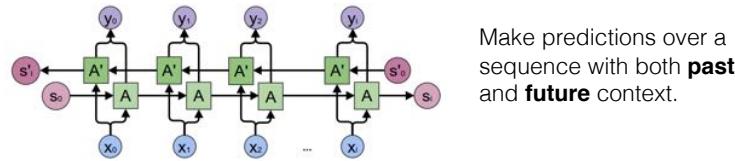
[LSTM: A Search Space Odyssey, Greff et al., 2015]

Other variants

RNN



Bidirectional RNN



Make predictions over a sequence with both **past** and **future** context.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

RNN applications

- Image Captioning
- Image Captioning with Attention
- Question Answering
- Visual Question Answering
- Speech Recognition
- Action Recognition in Videos
- Text Parsing
- Machine Translation

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning



I look better than Tom Cruise.

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei
Show and Tell: A Neural Image Caption Generator, Vinyals et al.
Long-term Recurrent Convolutional Networks for Visual Recognition and Description, Donahue et al.
Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

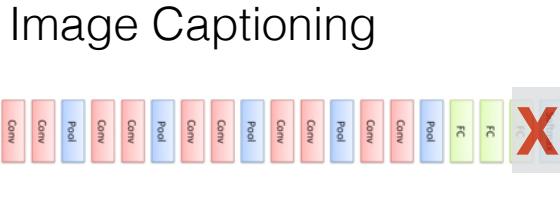
Image Captioning



I look better than Tom Cruise.



Image Captioning



Joseph J. Lim

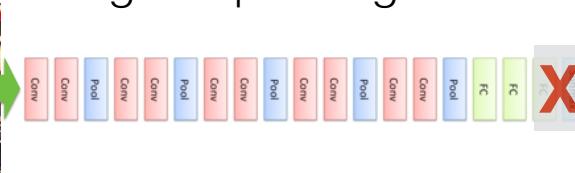
CSCI 599 @ USC

Lecture 8

Image Captioning



I look better than Tom Cruise.

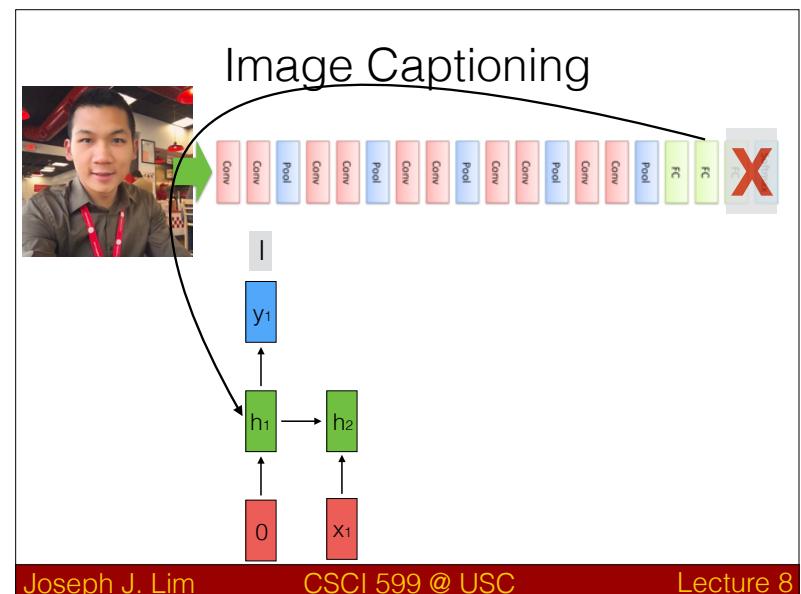
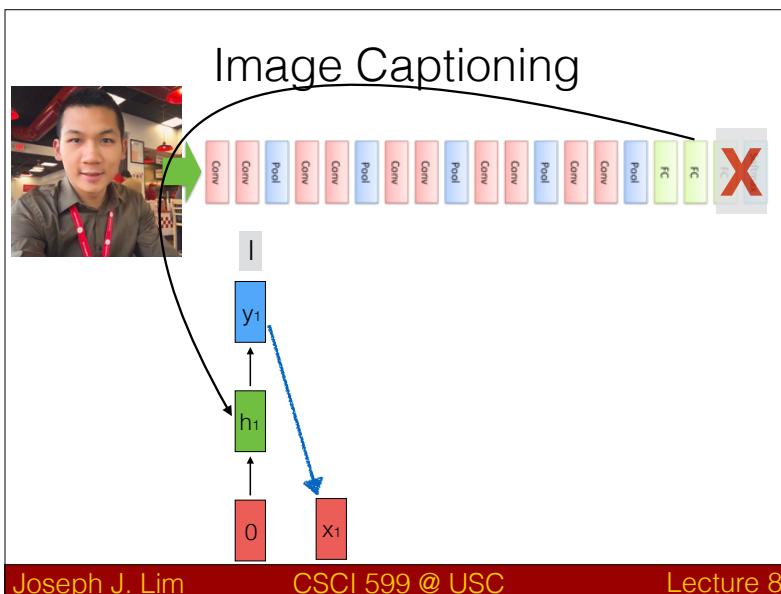
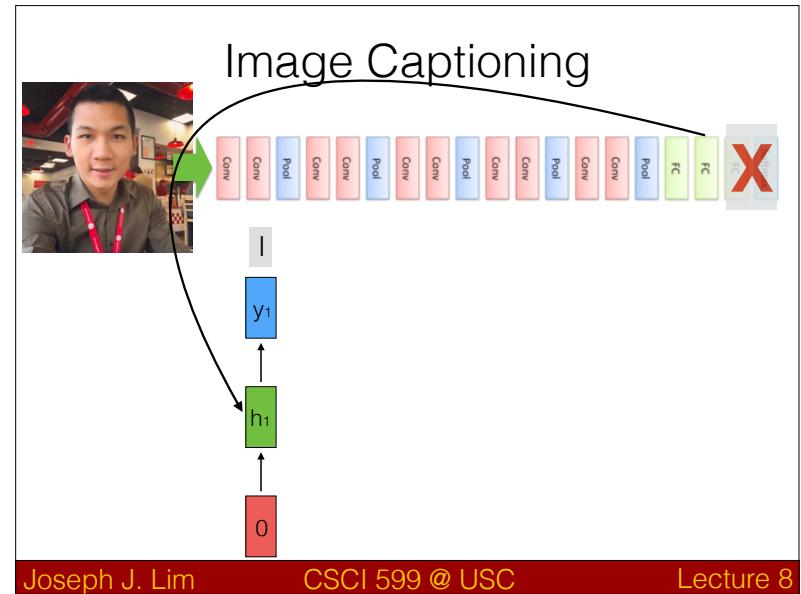
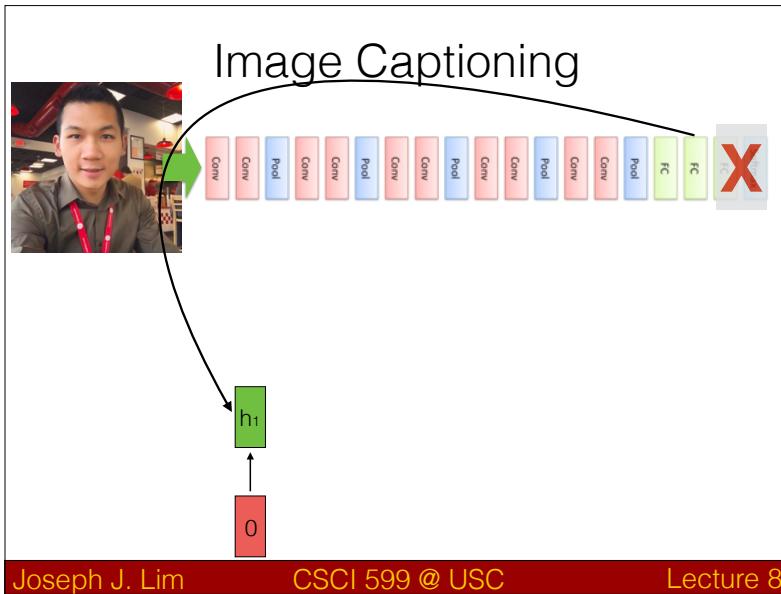


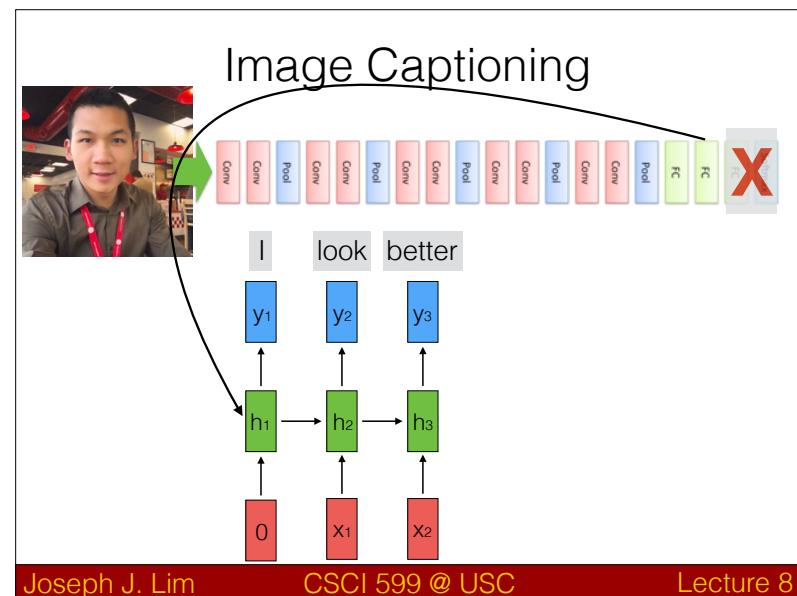
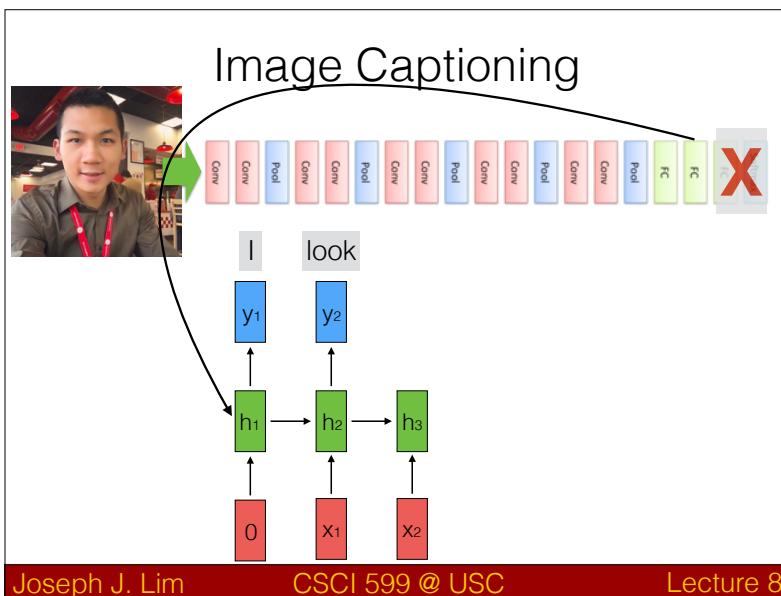
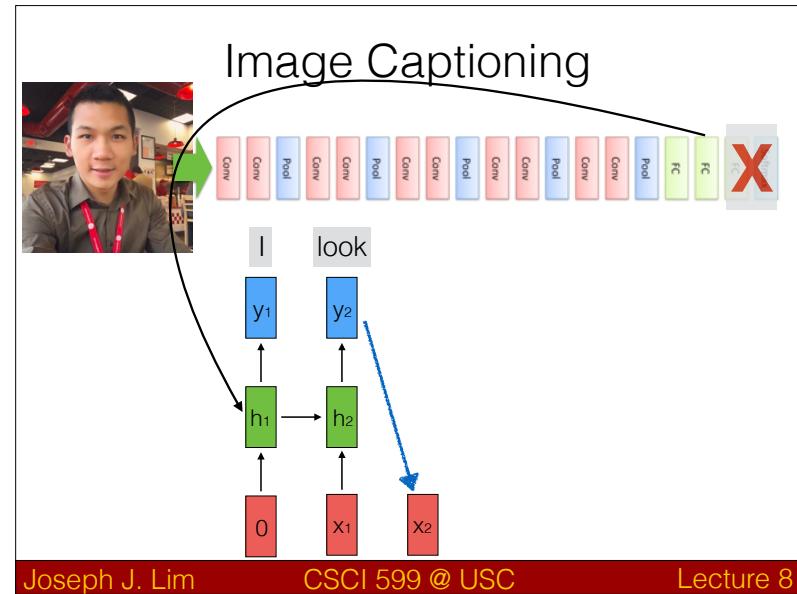
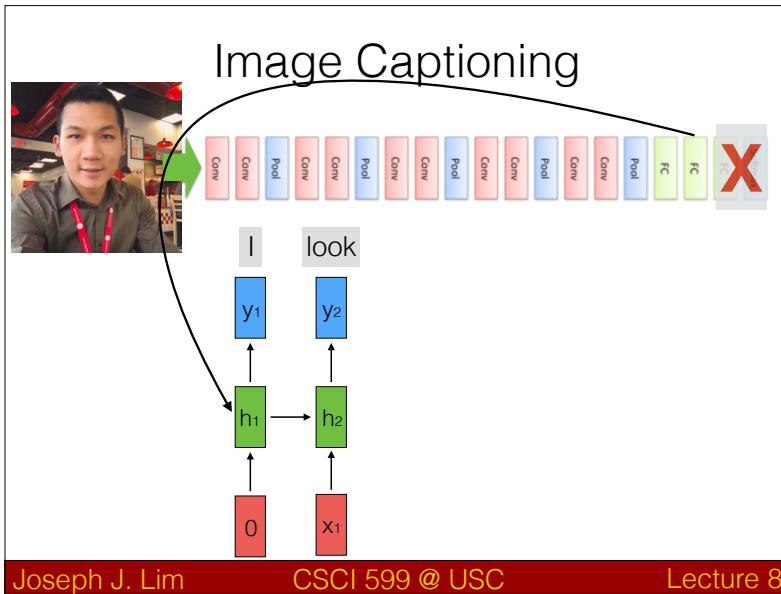
0

Joseph J. Lim

CSCI 599 @ USC

Lecture 8





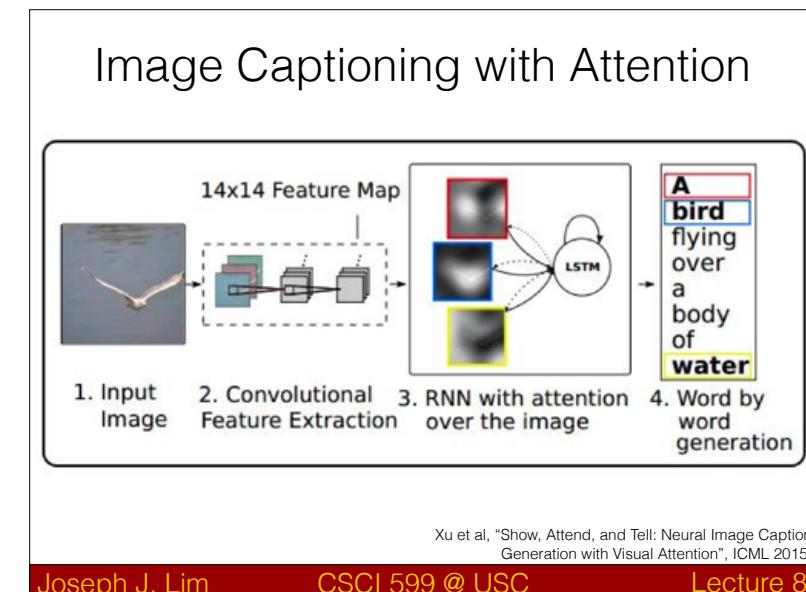
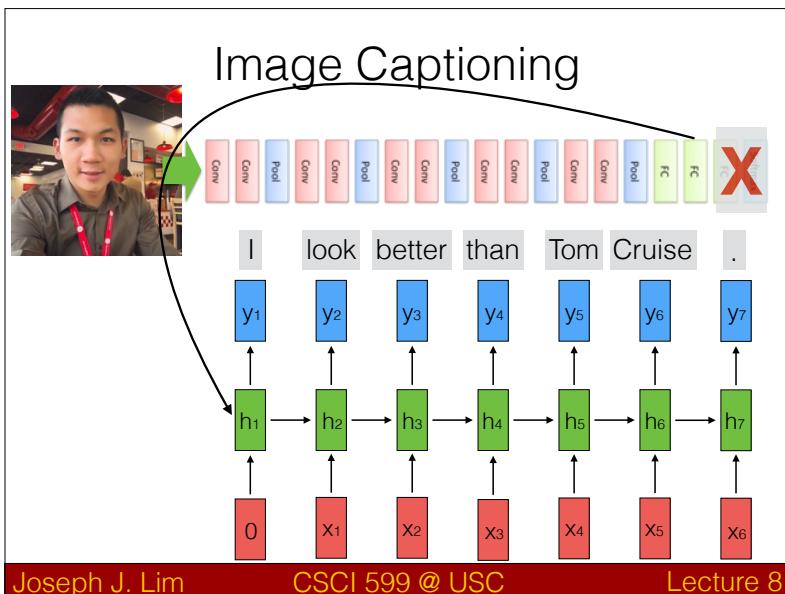
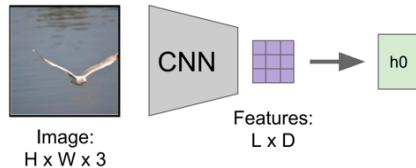


Image Captioning with Attention

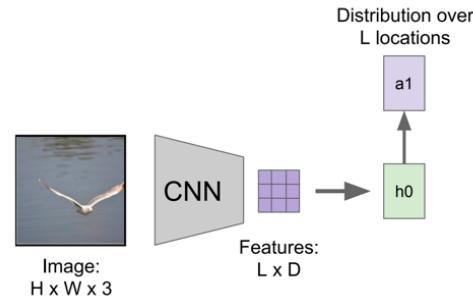


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning with Attention

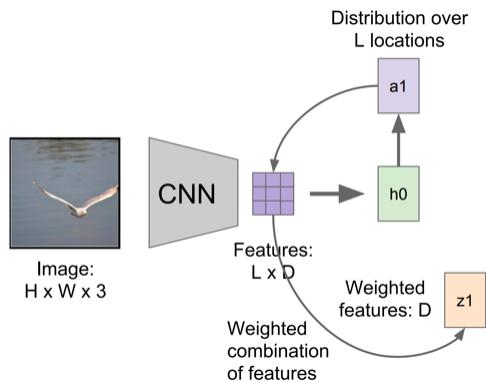


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning with Attention

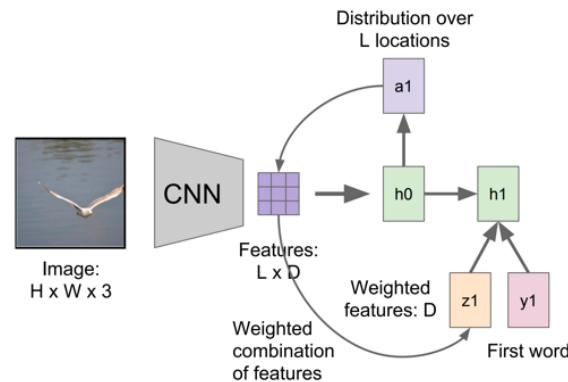


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning with Attention

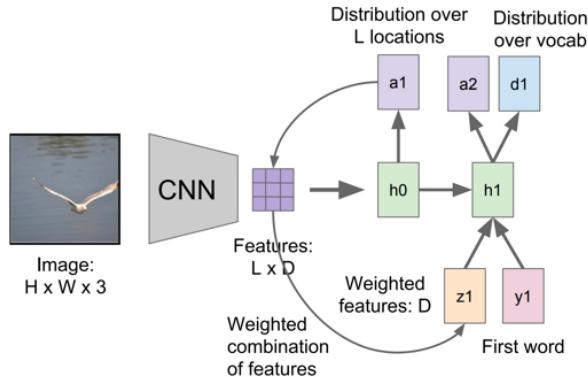


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning with Attention

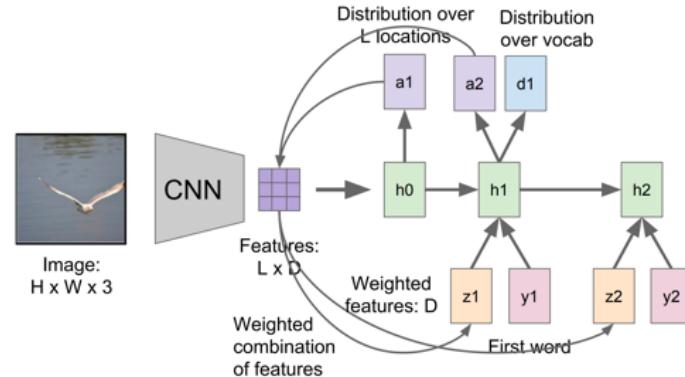


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning with Attention

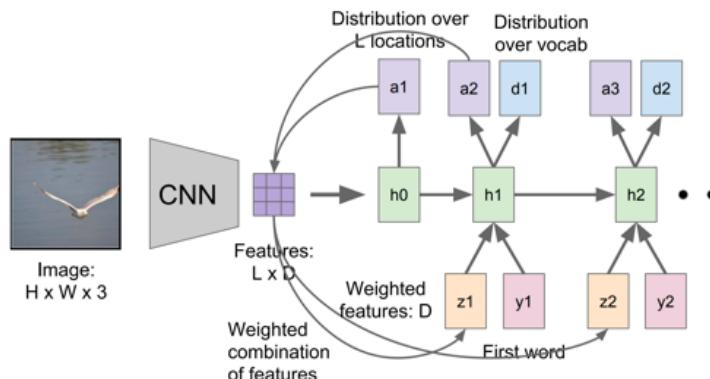


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning with Attention

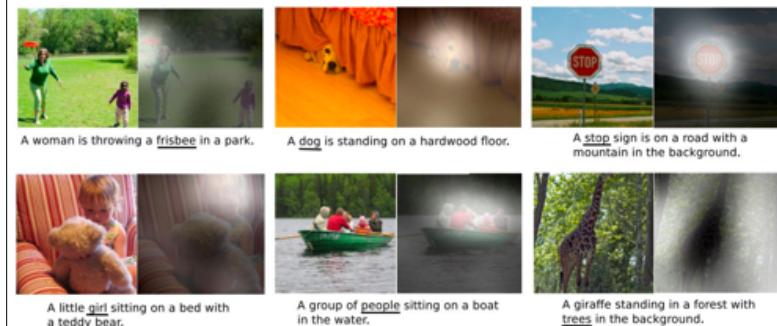


Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Image Captioning with Attention

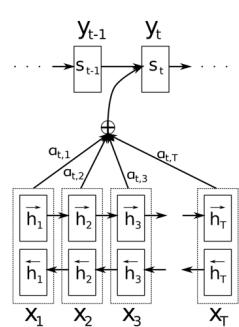


Joseph J. Lim

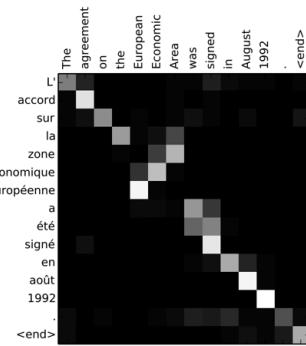
CSCI 599 @ USC

Lecture 8

Machine Translation with Attention



Self-Attention Mechanism
a: the soft attention coefficients



Visualizing a sample attention

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." arXiv preprint arXiv:1409.0473 (2014).

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Question Answering

Mary moved to the bathroom.

John went to the hallway.

Where is Mary? bathroom

<https://allenai.github.io/bi-att-flow/demo/>
Dialogue Learning With Human-in-the-loop, Jiwei et al.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Question Answering

Text

Mary moved to the bathroom.
John went to the hallway.

Question

Where is Mary?

Answer

bathroom

<https://allenai.github.io/bi-att-flow/demo/>
Dialogue Learning With Human-in-the-loop, Jiwei et al.

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Question Answering

Text

Mary moved to the bathroom.
John went to the hallway.

T Text embedding

Question

- Where is Mary? Q Question embedding

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer

X1
Mary

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer

h₁
X1
Mary

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer

h₁
X1
X2
Mary moved

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer

h₁ → h₂
X1 X2
Mary moved

Joseph J. Lim

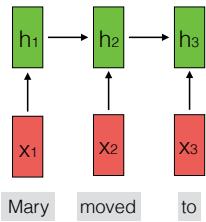
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Mary moved to

Joseph J. Lim

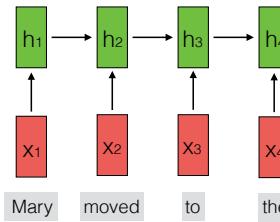
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Mary moved to the

Joseph J. Lim

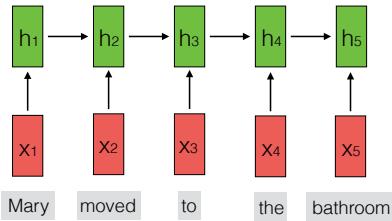
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Mary moved to the bathroom

Joseph J. Lim

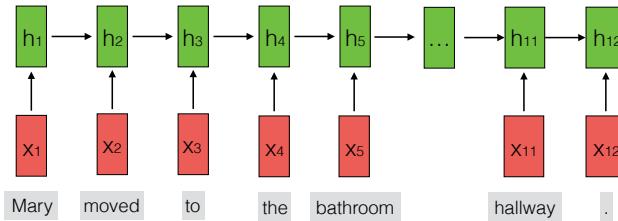
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Mary moved to the bathroom, John went to the hallway, .

Joseph J. Lim

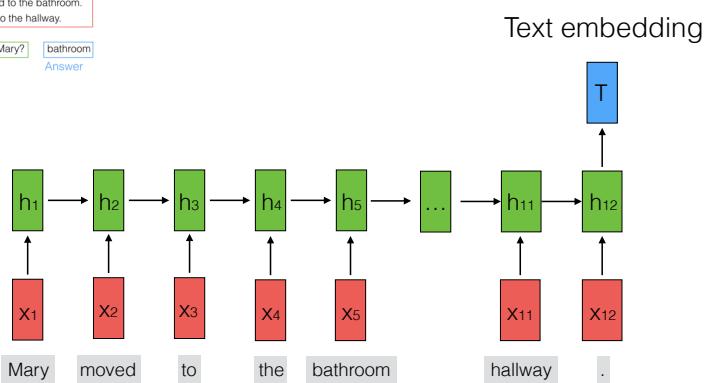
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Joseph J. Lim

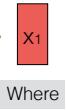
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Joseph J. Lim

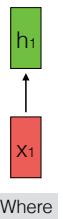
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Joseph J. Lim

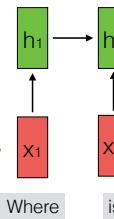
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? [bathroom]
Question Answer



Joseph J. Lim

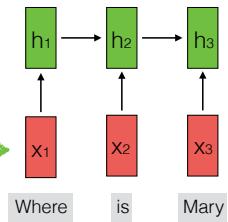
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? bathroom
Question Answer



Joseph J. Lim

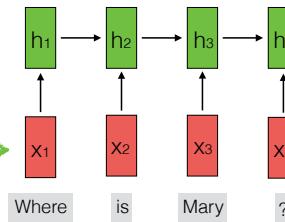
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? bathroom
Question Answer



Joseph J. Lim

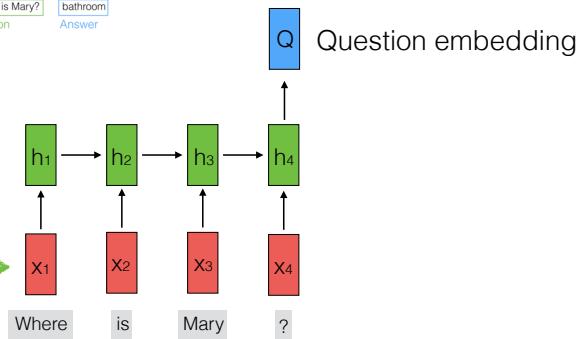
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? bathroom
Question Answer



Joseph J. Lim

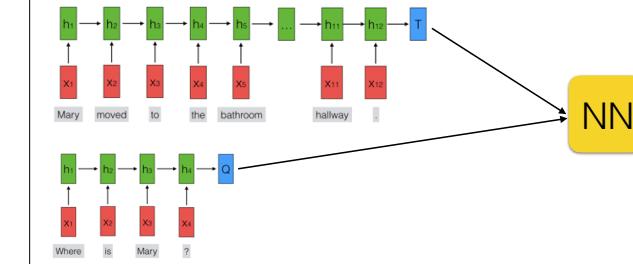
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

- Where is Mary? bathroom
Question Answer



Joseph J. Lim

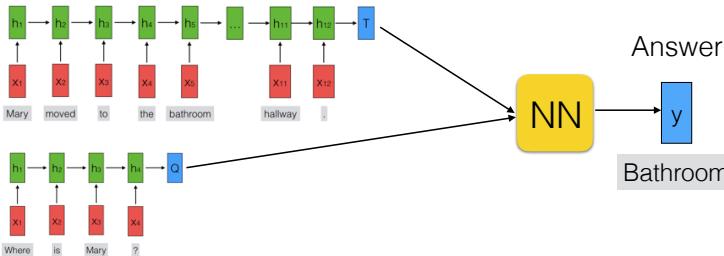
CSCI 599 @ USC

Lecture 8

Question Answering

Text
Mary moved to the bathroom.
John went to the hallway.

Where is Mary? bathroom
Question Answer



Joseph J. Lim

CSCI 599 @ USC

Lecture 8

DEMO

Paragraph

Southern California, often abbreviated SoCal, is a geographic and cultural region that generally comprises California's southernmost 10 counties. The region is traditionally described as "eight counties", based on demographics and economic ties: Imperial, Los Angeles, Orange, Riverside, San Bernardino, San Diego, Santa Barbara, and Ventura. The more extensive 10-county definition, including Kern and San Luis Obispo counties, is also used based on historical political divisions. Southern California is a major economic center for the state of California and the United States.

Question

What is Southern California often abbreviated as?
new question

Answer

SoCal

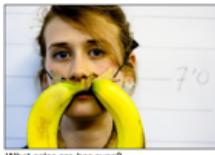
- <https://allenai.github.io/bi-att-flow/demo/>

Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Visual Question Answering



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

VQA: Visual Question Answering, Agrawal et al
A simple neural network module for relational reasoning, Santoro et al

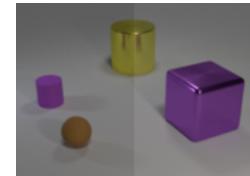
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Visual Question Answering

Image



Question

What is the size of
the brown sphere?

Answer

Small

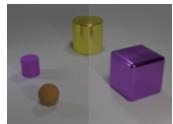
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Visual Question Answering

Image



Question

What is the size of
the brown sphere?

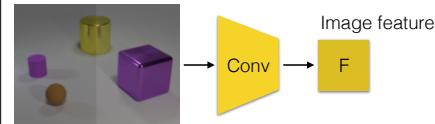
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Visual Question Answering

Image



Question

What is the size of
the brown sphere?

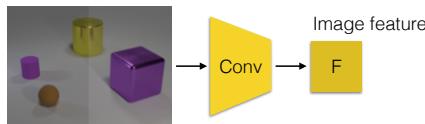
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Visual Question Answering

Image



Question

What is the size of
the brown sphere?

h_1

What

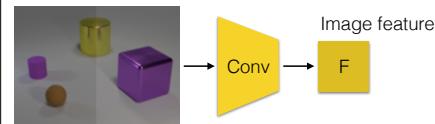
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Visual Question Answering

Image



Question

What is the size of
the brown sphere?

$h_1 \rightarrow h_2$

What is

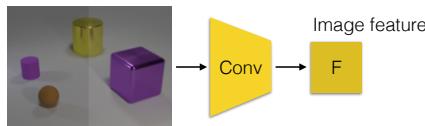
Joseph J. Lim

CSCI 599 @ USC

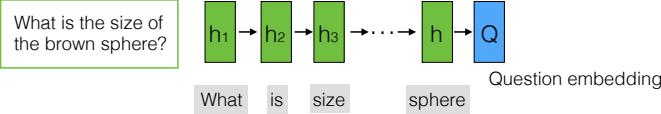
Lecture 8

Visual Question Answering

Image



Question



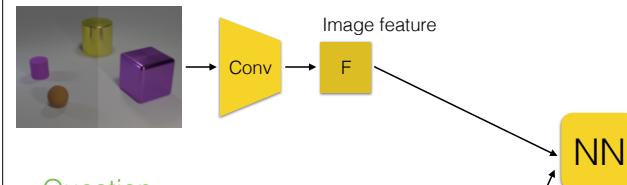
Joseph J. Lim

CSCI 599 @ USC

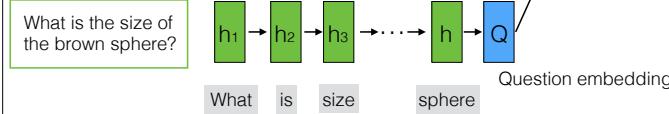
Lecture 8

Visual Question Answering

Image



Question



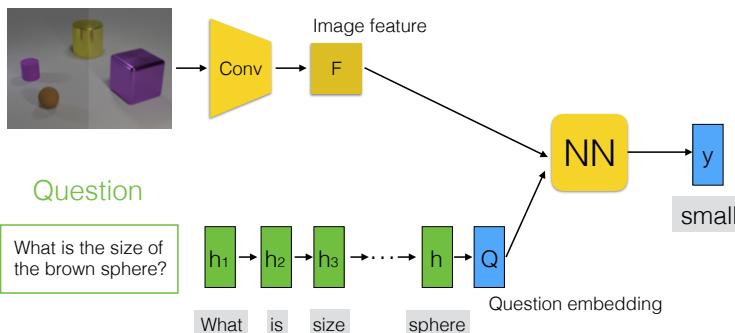
Joseph J. Lim

CSCI 599 @ USC

Lecture 8

Visual Question Answering

Image



Joseph J. Lim

CSCI 599 @ USC

Lecture 8

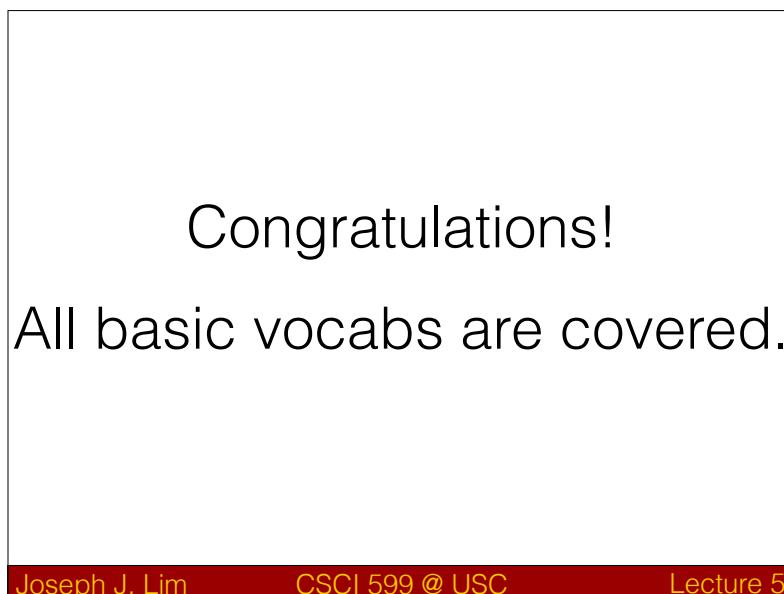
Summary

- RNNs deal with various sequential data or dependencies.
- LSTM or GRU are go-to models (rather than vanilla RNNs)
- RNNs are still under development (relative to CNNs)

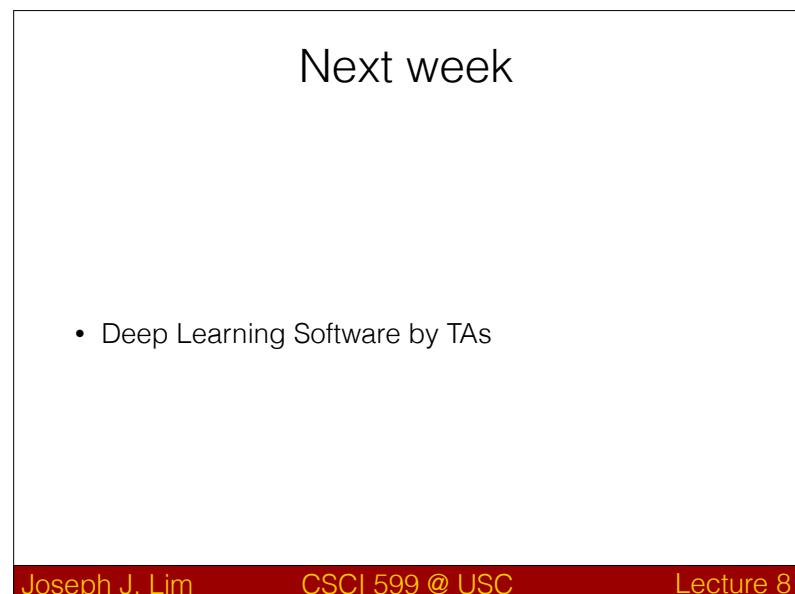
Joseph J. Lim

CSCI 599 @ USC

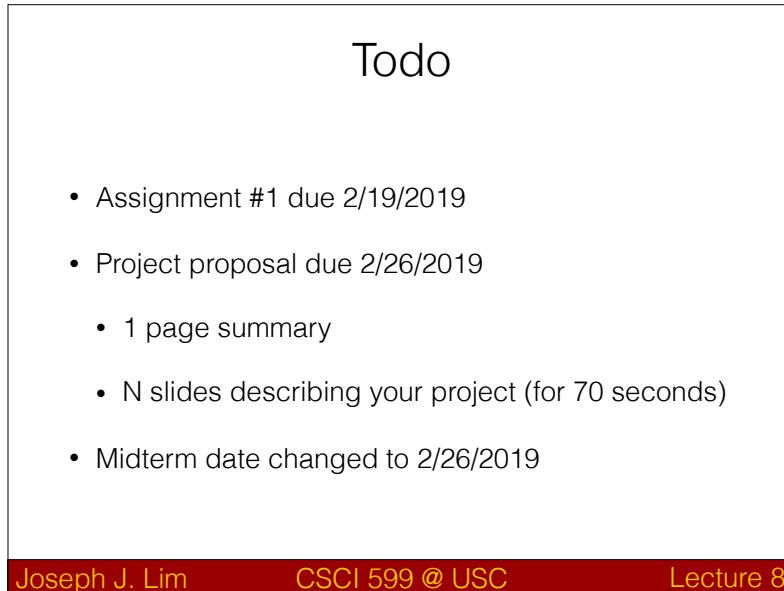
Lecture 8



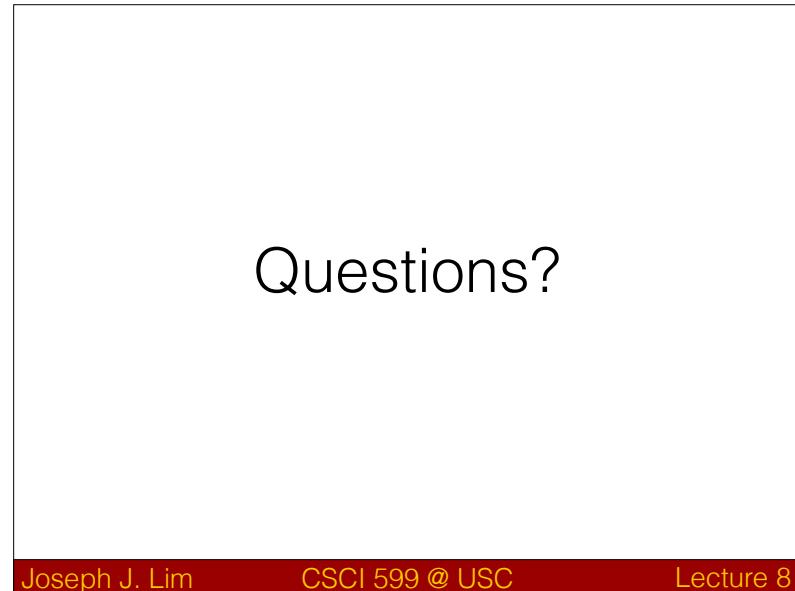
Joseph J. Lim CSCI 599 @ USC Lecture 5



Joseph J. Lim CSCI 599 @ USC Lecture 8



Joseph J. Lim CSCI 599 @ USC Lecture 8



Joseph J. Lim CSCI 599 @ USC Lecture 8