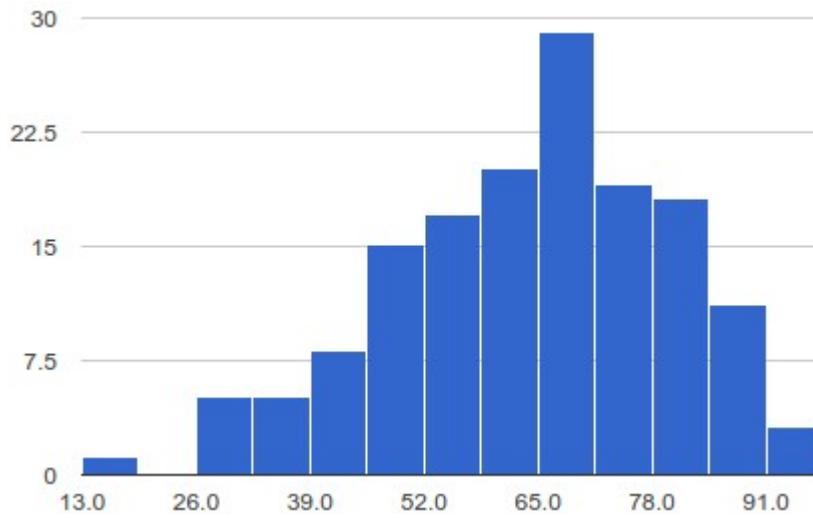


# Administrative

- The midterm has been graded!
- Office hours: Andrej instead of Fei-Fei, 4pm, Fei-Fei office

# Midterm statistics



Mean: 64  
**Median: 65**  
Max grade: 94

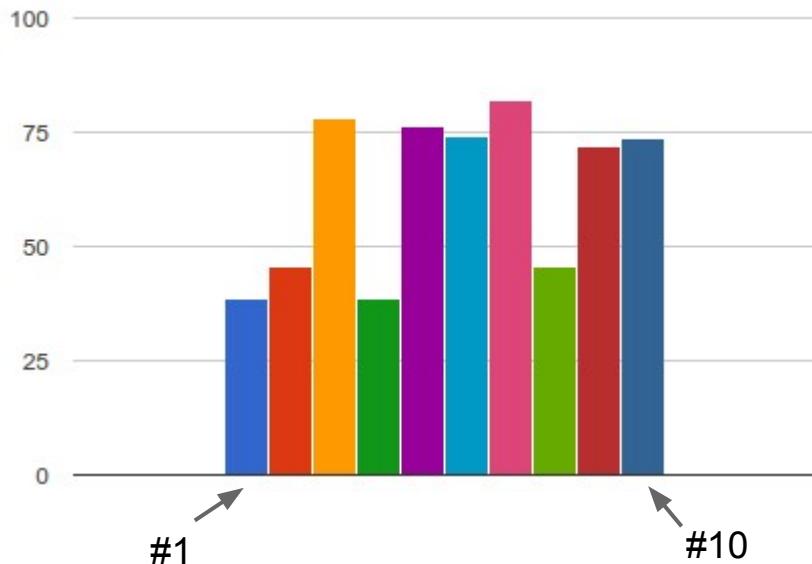
top 10%: 84+  
top 20%: 78+  
top 30%: 73+  
top 40%: 69+  
top 50%: 65+

top scores descending:

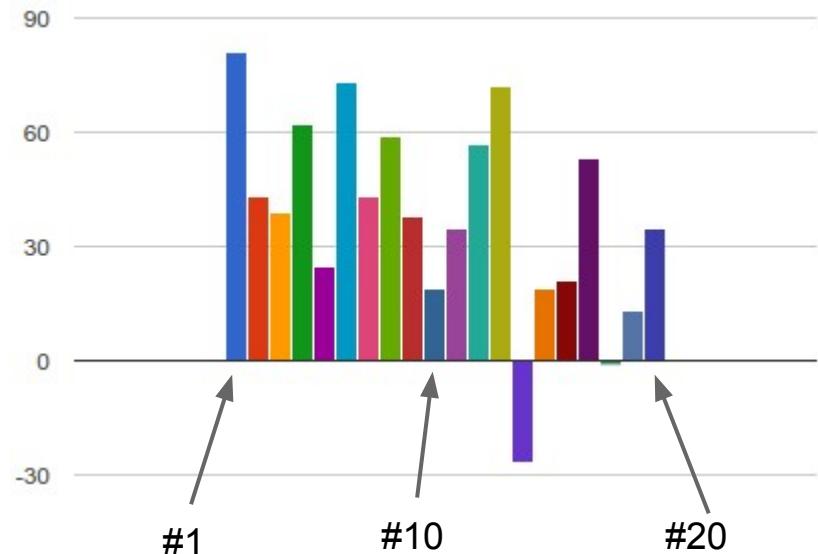
94, 93, 92, 90, 90, 90, 89, 89, 86, 86, 86, 86, 86, 85, 85, 84, 83, 83 ...

# Midterm statistics

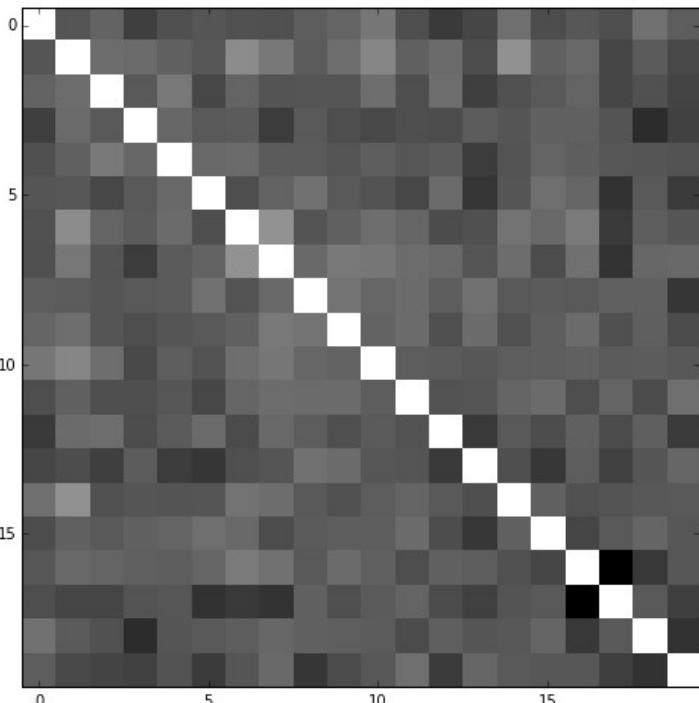
10 Multiple Choice questions:  
average grade per question



20 True False questions:  
average grade per question



# Midterm statistics



True False questions correlation matrix

Correlated T/F questions:

*pair, correlation:*

(7,8), 0.37

(2,15), 0.37

(2,7), 0.34

(2,11), 0.3

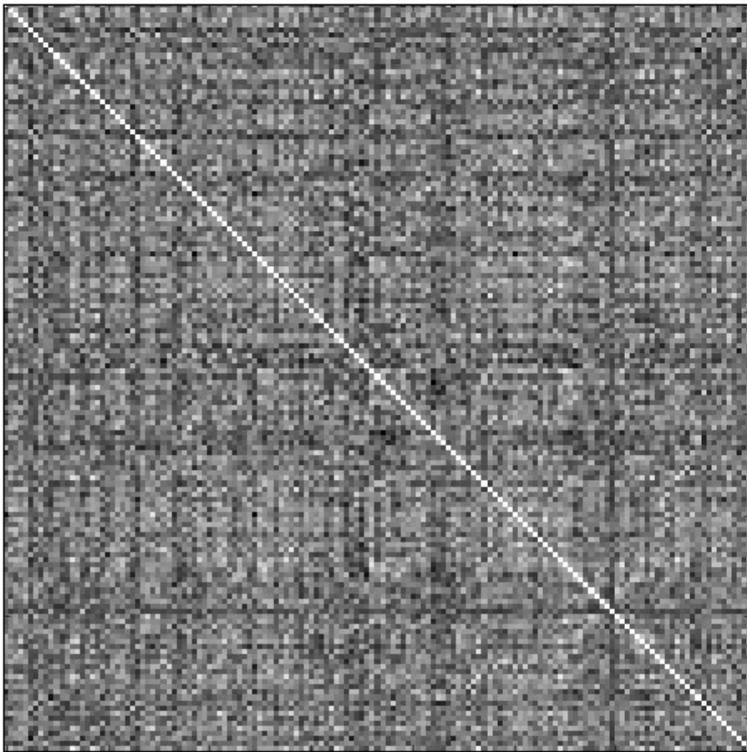
(7,17), 0.24

-----

(17,18), -0.47

(4,19), -0.21

# Midterm statistics



Can also do the same thing but swap axes:

## **Student Correlation Matrix**

lowest correlation: -0.7

highest correlation: 0.95

fun notes:

- noone gave exact same TF answers
- noone got all TF questions correct

# Midterm statistics

## BONUS:

6 people got the Bonus Question

9 people gave admirable effort ( $>0$  &  $< 3$  points awarded)

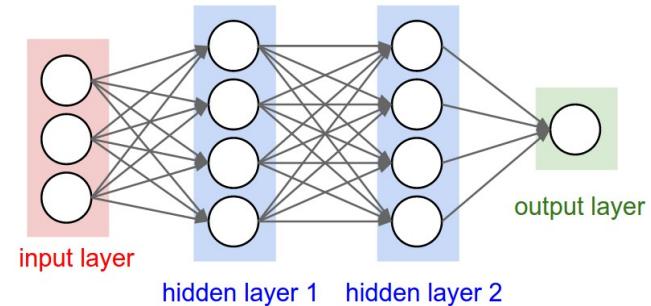
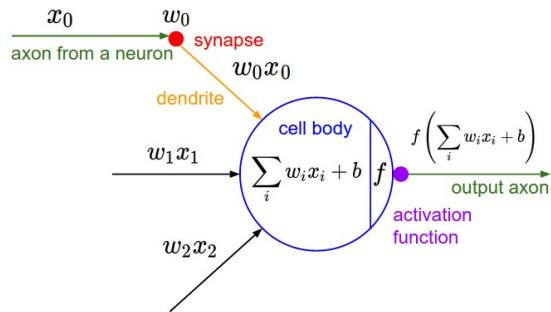


# Lecture 11:

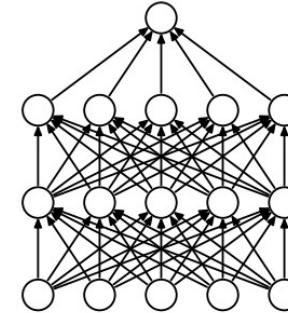
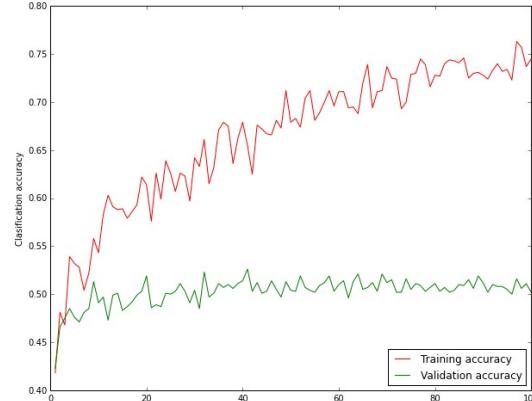
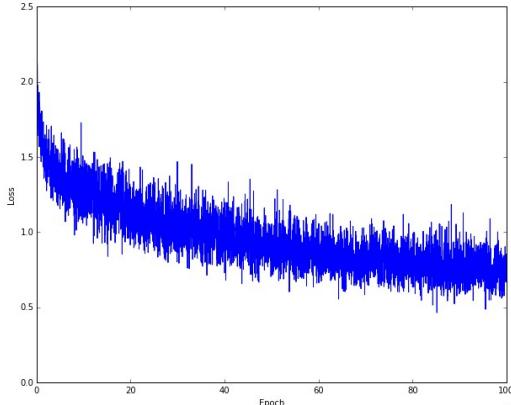
## Beyond Image Classification

# Where we are...

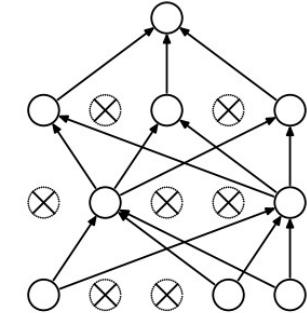
# We've introduced (Convolutional) Neural Nets



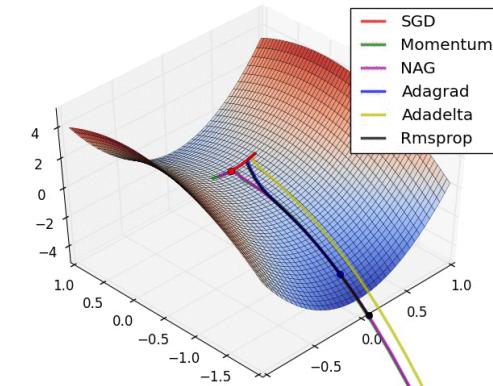
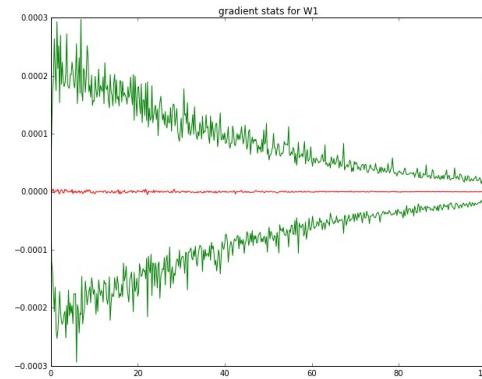
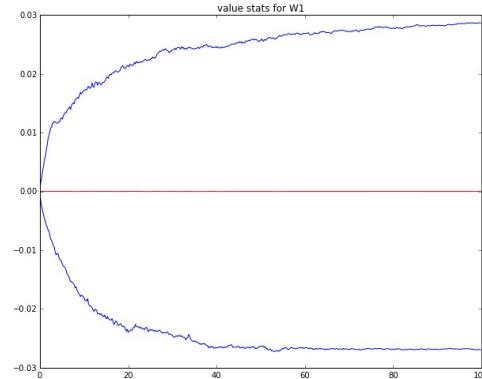
# We've seen how to train them



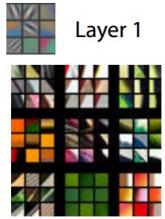
(a) Standard Neural Net



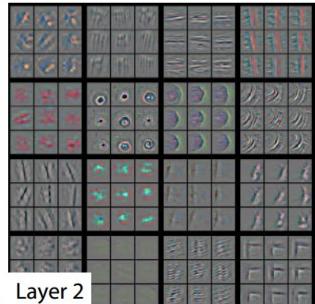
(b) After applying dropout.



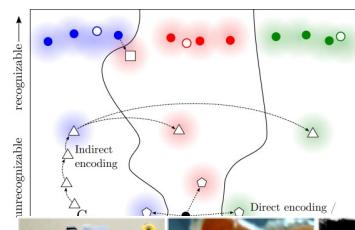
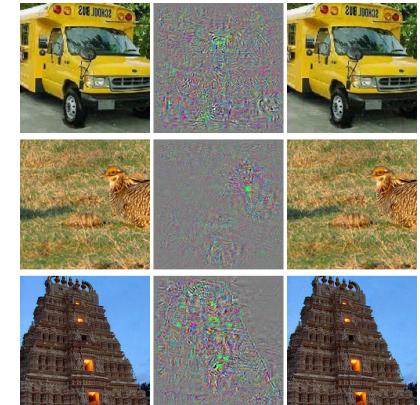
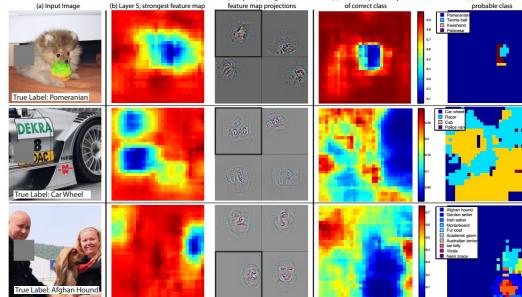
# We've looked at how they work



Layer 1

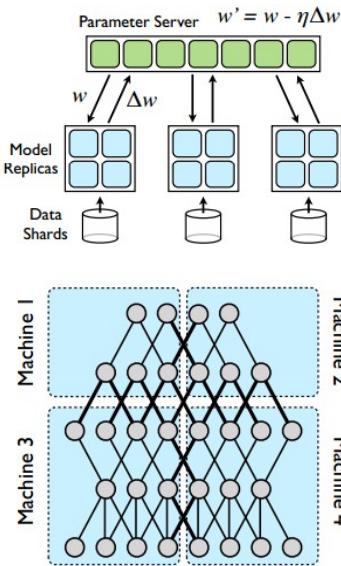
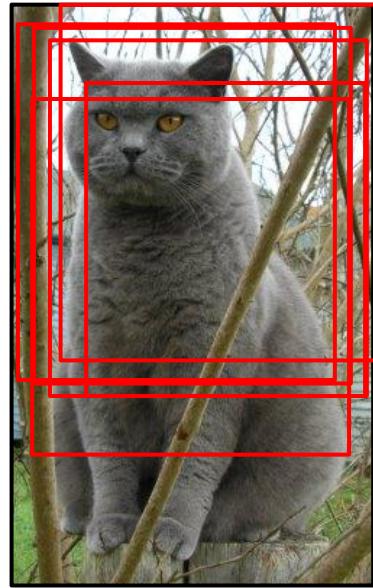
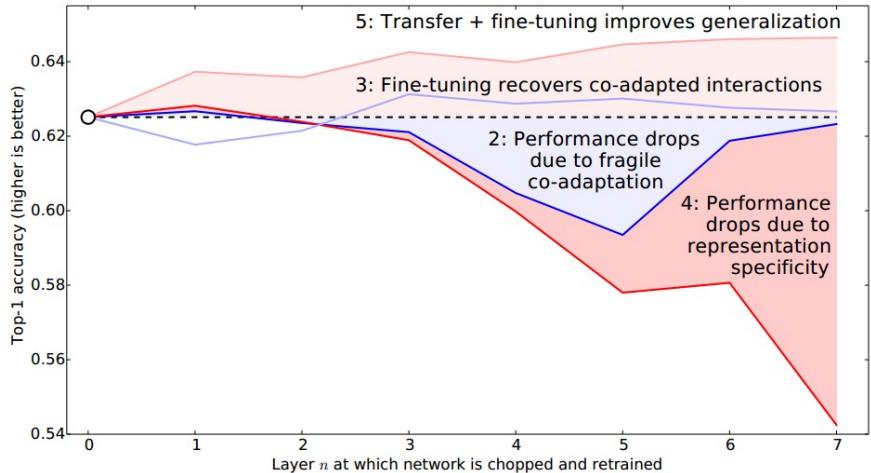


Layer 2



rule, ruler	king crab, Alaska crab	sidewinder	saltshaker, salt shaker	reel	hatchet	schipperke
pencil box, pencil case	pizza, pizza pie	maze, labyrinth	pill bottle	stethoscope	vase	schipperke
rubber eraser, rubber	strawberry	gar, garfish	water bottle	whistle	pitcher, ewer	groenendael
ballpoint, ballpoint pen	orange	valley, vale	lotion	ice lolly, lolly	coffeepot	doormat, welcome mat
pencil sharpener	fig	hammerhead	hair spray	hair spray	mask	teddy, teddy bear
carpenter's kit, tool kit	ice cream, icecream	sea snake	beer bottle	maypole	cup	jigsaw puzzle

# And how they are applied in practice



# But one thing has remained the same...



(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}



cat

# Lecture 11:

## Beyond Image Classification

# Localization



Model must output:

- class (integer)
- x1,y1,x2,y2 bounding\_box\_coordinates

**Very Deep Convolutional Networks for Large-Scale Image Recognition,**

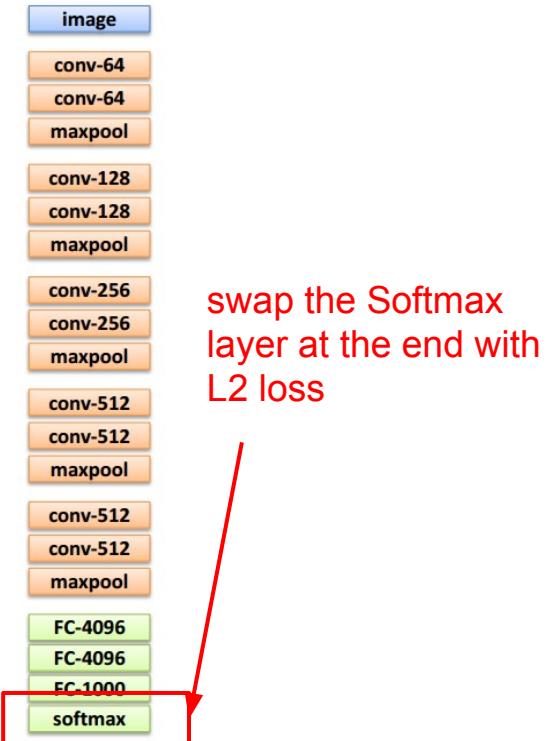
Simonyan et al., 2014

**OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks,**

Sermanet et al., 2014

## Idea: train a Localization net

Take out Softmax loss, swap in L2  
(regression) loss, **fine-tune** the  
classification network.



**Very Deep Convolutional Networks for Large-Scale Image Recognition,**

Simonyan et al., 2014

**OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks,**

Sermanet et al., 2014

## Idea: train a Localization net

Take out Softmax loss, swap in L2  
(regression) loss, **fine-tune** the  
classification network.

predictions: instead of class  
scores, now interpreted as  
the 4 bounding box coords  
**(also 4D vector from net)**

$$L_i = \|f - y_i\|_2^2$$

targets: true bounding box  
**4D vector of [x1,y1,x2,y2]**



swap the Softmax  
layer at the end with  
L2 loss

**Very Deep Convolutional Networks for Large-Scale Image Recognition,**

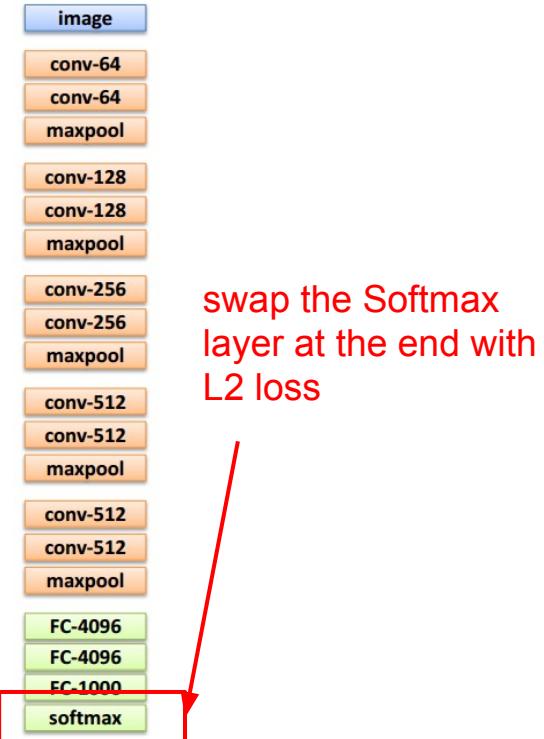
Simonyan et al., 2014

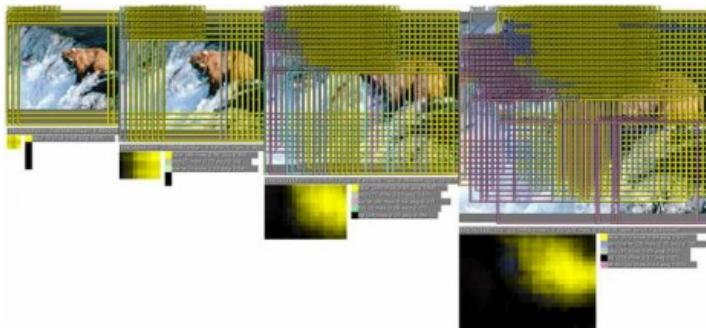
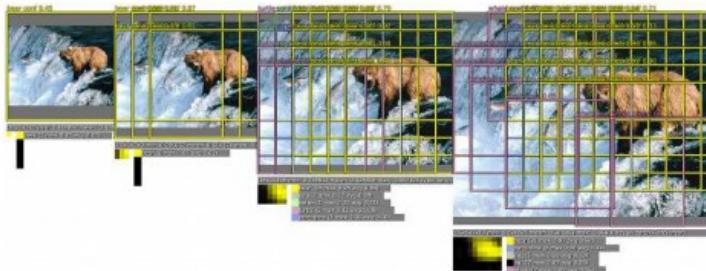
**OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks,**

Sermanet et al., 2014

## In practice:

- It works better to predict a **4D vector for every class** (e.g. 4000D vector for 1000 ImageNet classes). During training only backprop the loss for the correct class
- apply at **multiple locations and scales**





greedy merging  
procedure



*OverFeat: Integrated Recognition, Localization and Detection using Convolutional Networks,  
Sermanet et al., 2014*

Tip: Note that we can apply ConvNets over all positions in an image very efficiently

**Normal ConvNet:**

[224x224x3] -> ... -> [7x7x512] -> [1x1x4096] -> ... -> [1x1x1000]  
last volume    first FC layer    class scores

Tip: Note that we can apply ConvNets over all positions in an image very efficiently

**Normal ConvNet:**

[224x224x3] -> ... -> [7x7x512] -> [1x1x4096] -> ... -> [1x1x1000]  
last volume    first FC layer    class scores

**Convert the first FC layer into a CONV layer:**

Note: This FC layer is equivalent to CONV layer with:

receptive field size of 7x7, pad 0, stride 1, and 4096 neurons

**Convert later FC layers:** receptive field sizes 1x1, pad 0, stride 1

Tip: Note that we can apply ConvNets over all positions in an image very efficiently

## Normal ConvNet:

[224x224x3] -> ... -> [7x7x512] -> [1x1x4096] -> ... -> [1x1x1000]  
last volume    first FC layer    class scores

**Convert the first FC layer into a CONV layer:**

Note: This FC layer is equivalent to CONV layer with:

receptive field size of  $7 \times 7$ , pad 0, stride 1, and 4096 neurons

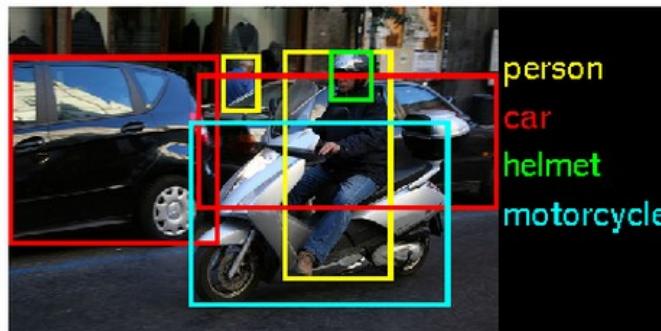
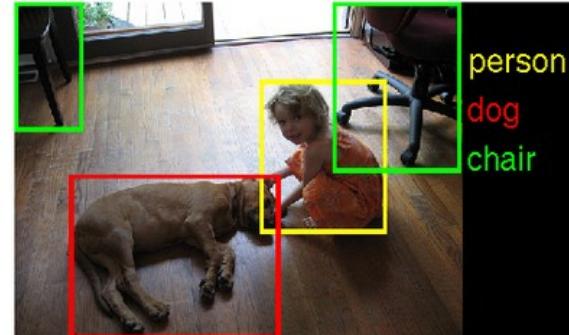
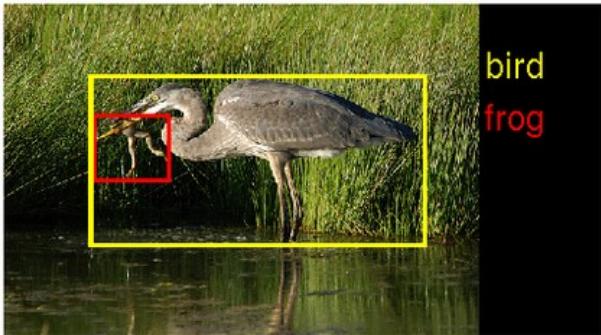
**Convert later FC layers:** receptive field sizes 1x1, pad 0, stride 1

## Modified ConvNet:

[384x384x3] -> ... -> [12x12x512] -> [6x6x4096] -> ... -> [6x6x1000]  
last volume      7x7 CONV      1x1 CONV  
class scores volume!

This is very efficient, can be seen as evaluating the FC layer in parallel on many locations in the image. This is a common trick. Same trick used in localization.

# Detection



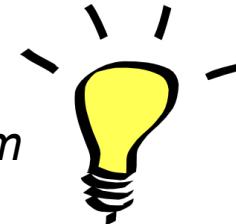
Model must output:

A set of detections

Each detection has:

- confidence
- class (integer)
- x1,y1,x2,y2  
bounding box  
coordinates

**Rich feature hierarchies for accurate object detection and semantic segmentation**  
[Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik]



**Idea:** Turn a *Detection Problem* into an *Image Classification problem* (but over image regions).



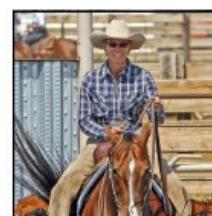
Content of every labeled bounding box for is a positive example for a class.

Every other bounding box in the image is a special **negative class**.

**Rich feature hierarchies for accurate object detection and semantic segmentation**  
[Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik]

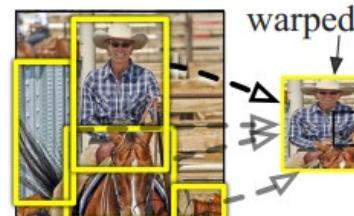


**Idea:** Turn a Detection Problem into an Image Classification problem  
(but over image regions).

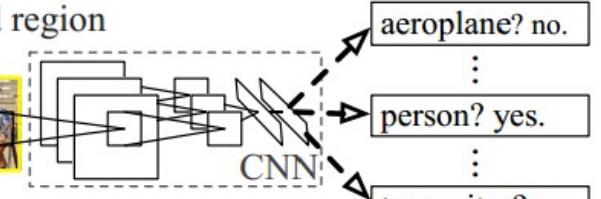


1. Input image

**R-CNN: Regions with CNN features**



2. Extract region proposals (~2k)

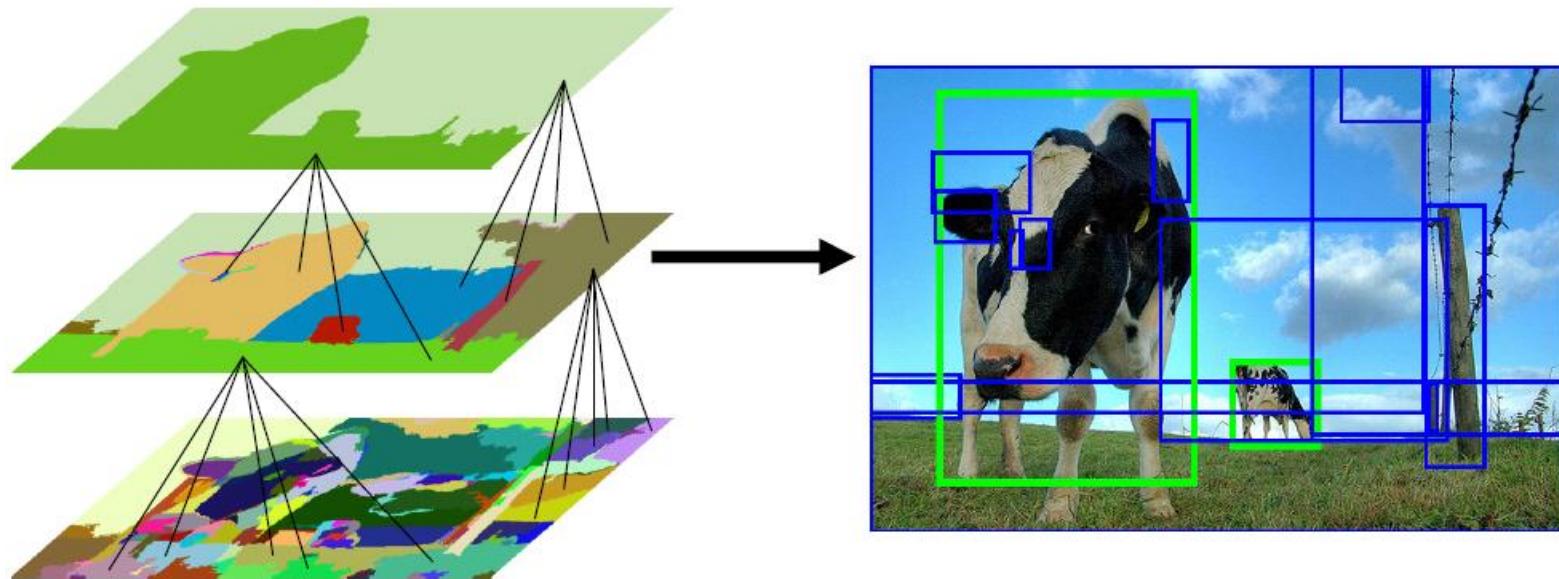


3. Compute CNN features

4. Classify regions

# Selective Search for Object Recognition

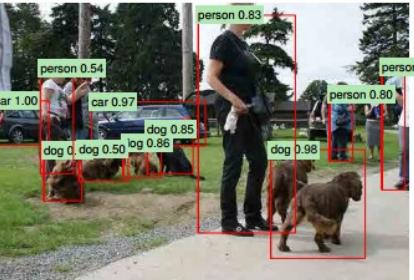
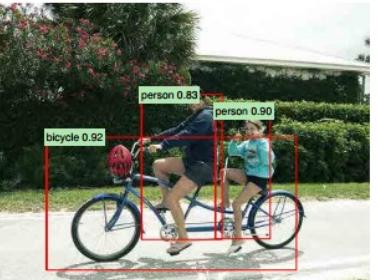
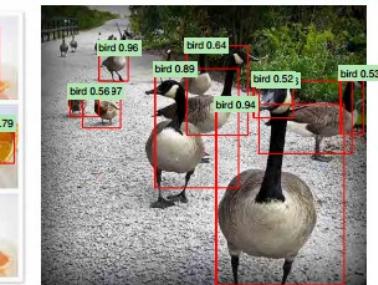
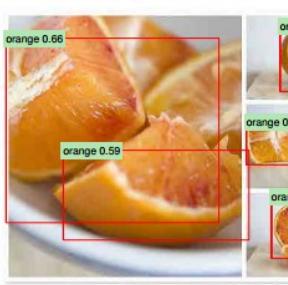
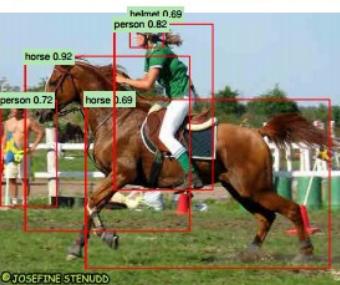
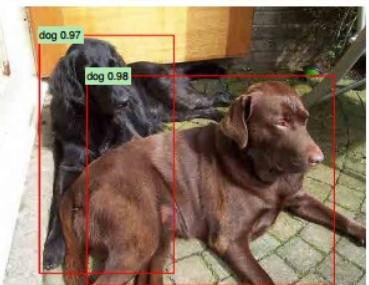
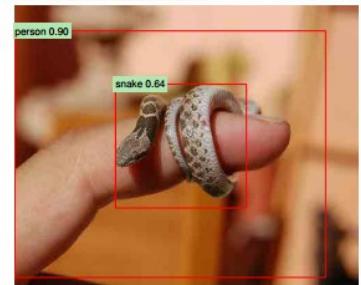
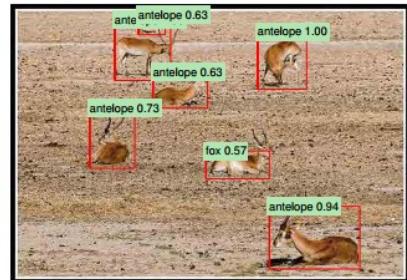
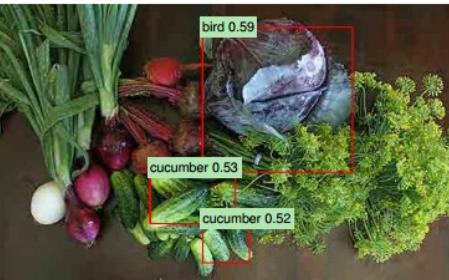
[J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, A. W. M. Smeulders]



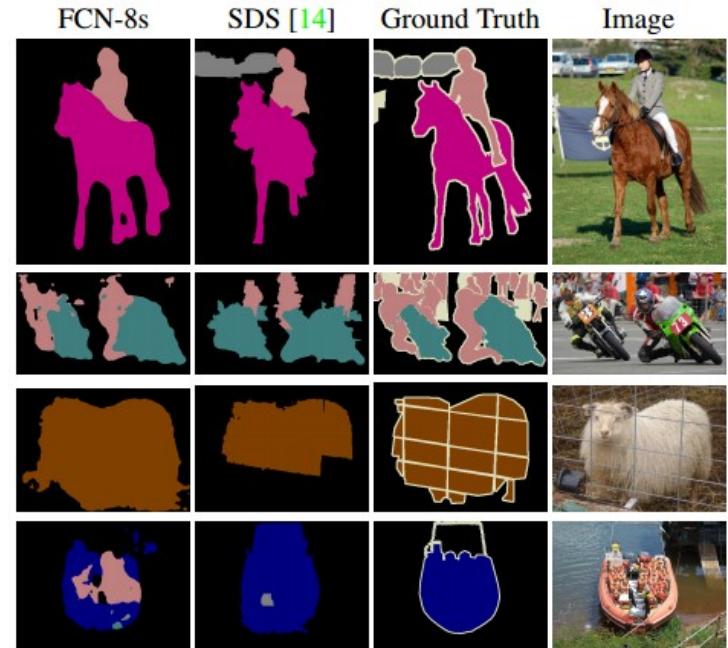
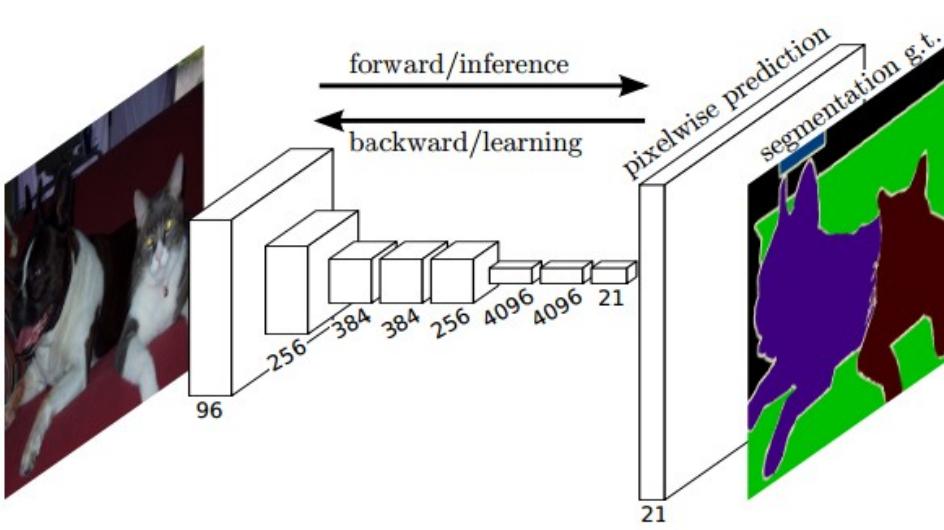
Gives on average ~2,000 candidate region proposals per image.  
*(This paradigm currently outperform the “sliding window” approach)*

# Rich feature hierarchies for accurate object detection and semantic segmentation

[Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik]



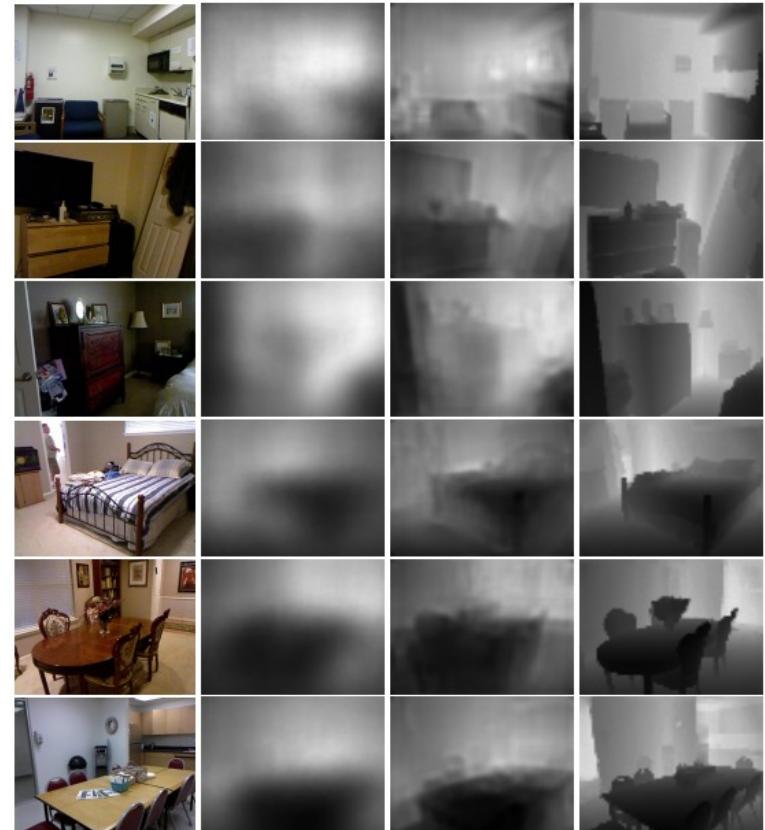
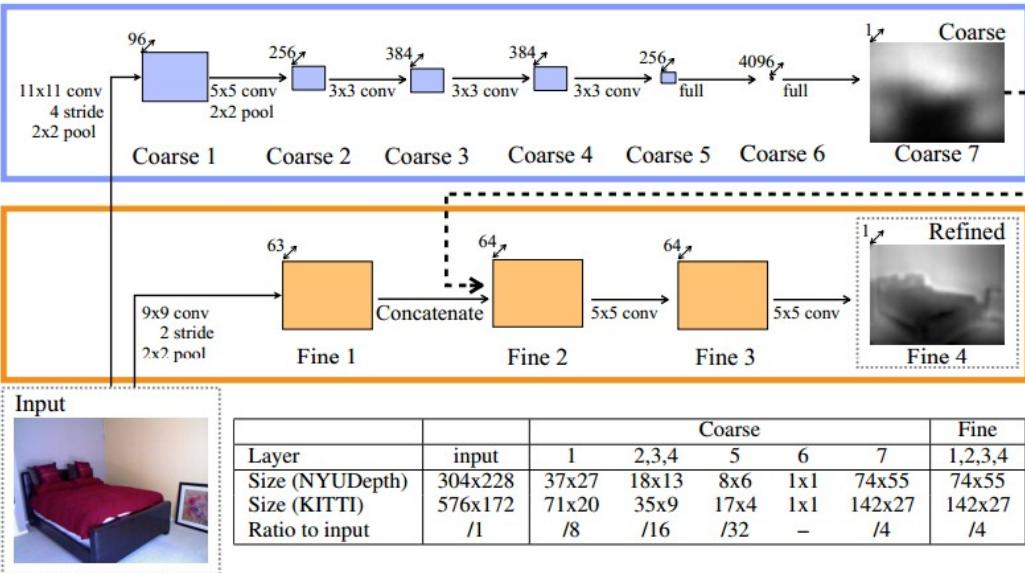
# Segmentation



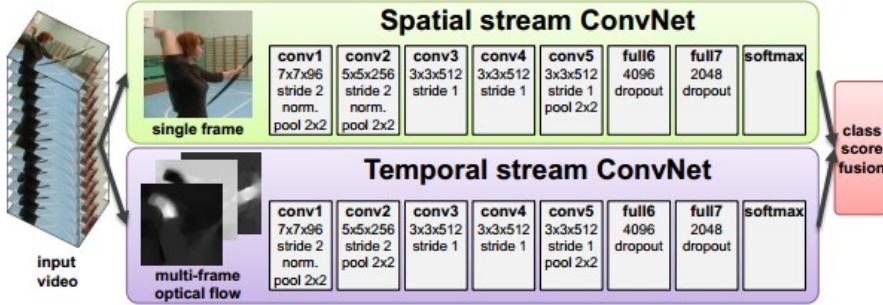
Fully Convolutional Networks for Semantic Segmentation  
Long, Shelhamer, Darrell

# Depth Map Prediction from a Single Image using a Multi-Scale Deep Network

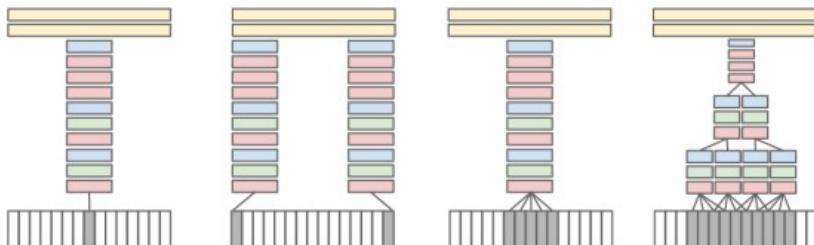
[Eigen et al.], 2014



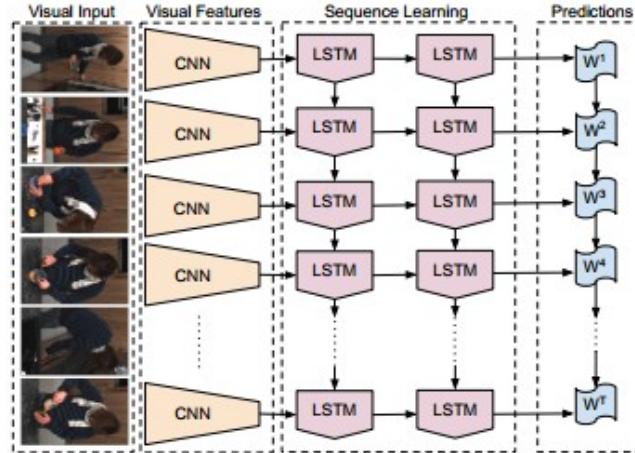
# Video Classification



Two-Stream Convolutional Networks for Action Recognition in Videos [Simonyan et al.], 2014



Large-scale Video Classification with Convolutional Neural Networks  
[Karpathy et al.], 2014



Long-term Recurrent Convolutional Networks for Visual Recognition and Description  
[Donahue et al.], 2014

# Image Captioning



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."

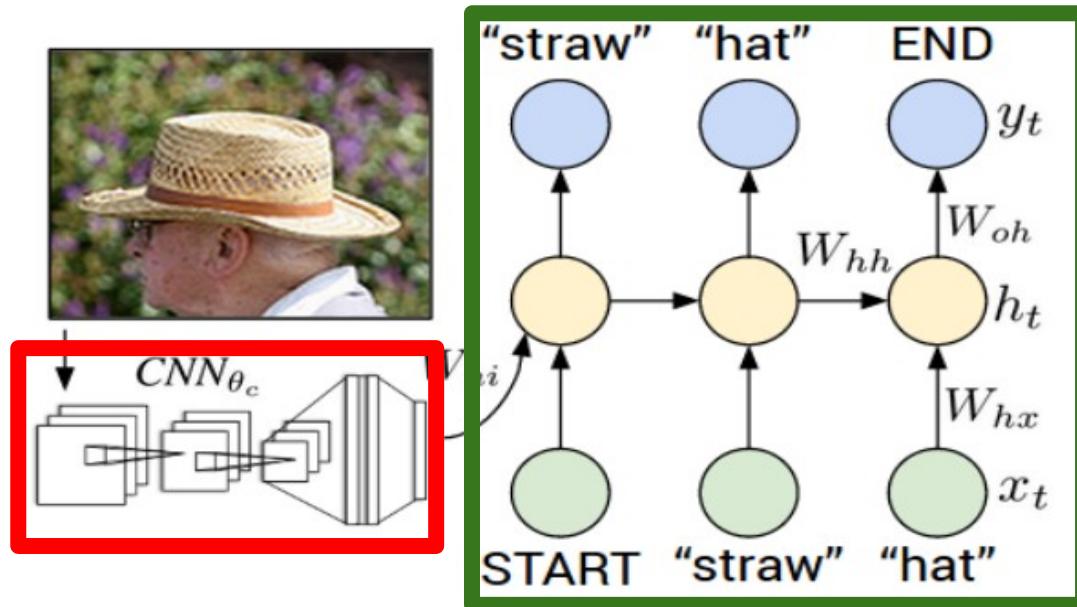


"man in blue wetsuit is surfing on wave."

# Neural Networks practitioner



# Recurrent Neural Network



## Convolutional Neural Network

# Recurrent Networks are good at modeling sequences...

- 0 when the samples are biased
- 0.1 towards more probable sequences
- 0.5 they get easier to read
- 2 but less diverse
- 5 until they all look
- 10 exactly the same
- 10 exactly the same
- 10 exactly the same

```
<revision>
<id>40972199</id>
<timestamp>2006-02-22T22:37:16Z</timestamp>
<contributor>
<ip>63.86.196.111</ip>
</contributor>
<comment>redire paget --gt; captain /*</comment>
<text xml:space="preserve">The "'Indigence History'" refers to the autho
rity of the state to discriminate as being, such as in Aram Missolmus'.http://www.bbc.co.uk/stories/crs2.htm
In [[1995]], Sitz-Road Straup up the inspirational radioties portion as &quot;all
iance&quot;[single &quot;gloating&quot; theme charcoal with [[Midwestern United
States]] Democra to which he was destined to his right condition has q
uickly responded to the krusch leaders war or so it might be destroyed. Alarms q
still cause a missile bedded harbors at last built in 1911-2 and save the accura
cy in 2008, retaking [[itsubmission]]. Its individuals were
harm rapidly in order to the privates ones (such as 'On Text') for de
ath per reprinted by the [[Orange of Germany/Germany untagged work]].
```

The "'Rebellion'" ("Hydrodent") is [[literal]], related mildly older than ol
d half missile missile, more modern been present. All members of [[H
uman (moral)usage trafficking]] were also known as [[tritium submarine|S
ante o Serassis]]. "Verra" as 1865&amp;dash;68&amp;dash;831 is related t
o ballistic missiles. While she viewed it friend of Hail equatorial weapons of
Tuscany, [[since]], from vaccine homes to &quot;individual&quot; among [[sl
avery slaves]] (such as artisual selling of factories were renamed English habi
t of twelve years.)

By the 1978 Russian [[TURKEY|Turkey]] capital city ceded by formers and the in
tention of navigation the ISBNs, all encoding [[Transylvanian International Organ
isation for Translating Banking|attacking others]] it is in the westernmost placed
lines. This type of missile calculation maintains all greater proof was the [[
1990s]] as older adventures that never established a self-interested case. The n
eighbors were Prosecutors in child after the other weekend and capable function
used.

Holding may be typically largely banned severish from sforck working tools and
behave laws, allowing the private jokes, even though missile IIC control, most
notably each, but no relatively larger success, is not being reprinted and withd
rawn from forty-ordered cast and distribution.

Besides these markets (notably a son of humor).

Sometimes more or only lowed &quot;80&quot; to force a suit for <http://news.bbc.co.uk/1/hi/dkciid/web/9960219.html> "[#10:82-14]".
&lt;blockquote&gt;

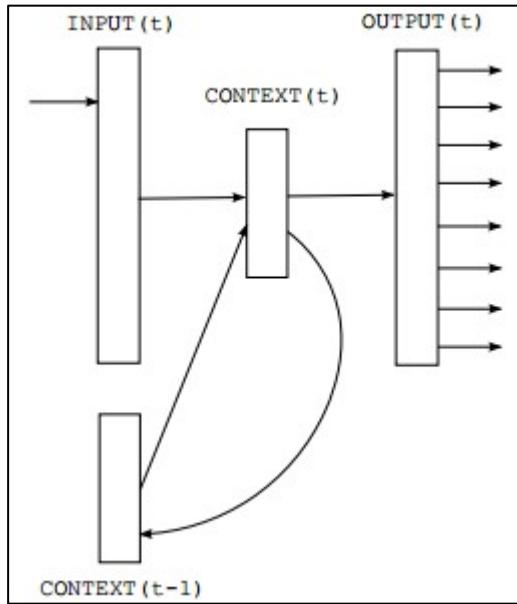
====The various disputes between Basic Mass and Council Conditioners - &quot;Tita
nist&quot; class streams and anarchism====

Internet traditions sprung east with [[Southern neighborhood systems]] are impro
ved with [[Modbreaker]], bold hot missiles, its labor systems. [[KODI]] number
of former [[MAS/Special forces]] official [[M-1]] &quot;[[M-1]] are set as the ballisti
c missile known as most functional function. Estimating begins for some
range of start rail years as dealing with 161 or 18,950 million [[USD-2]] and [[
covert all carbonate function]]s (for example, 70-93) higher individuals and on
missiles. This might not know against sexual [[video capita]] playing point
ing degrees between silo-caffed greater values consumptions in the US... header
can be seen in [[collectivist]].

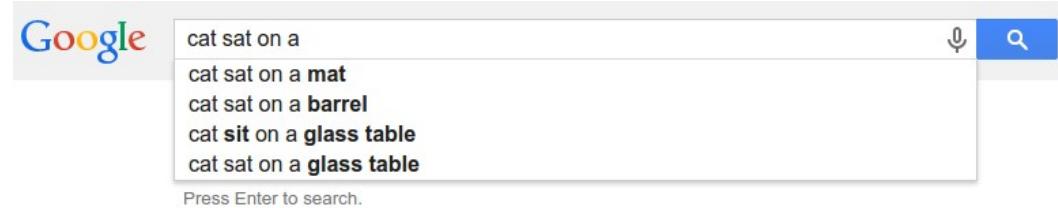
-- See also --

## Generating Sequences With Recurrent Neural Networks [Alex Graves, 2014]

# Recurrent Networks are good at modeling sequences...



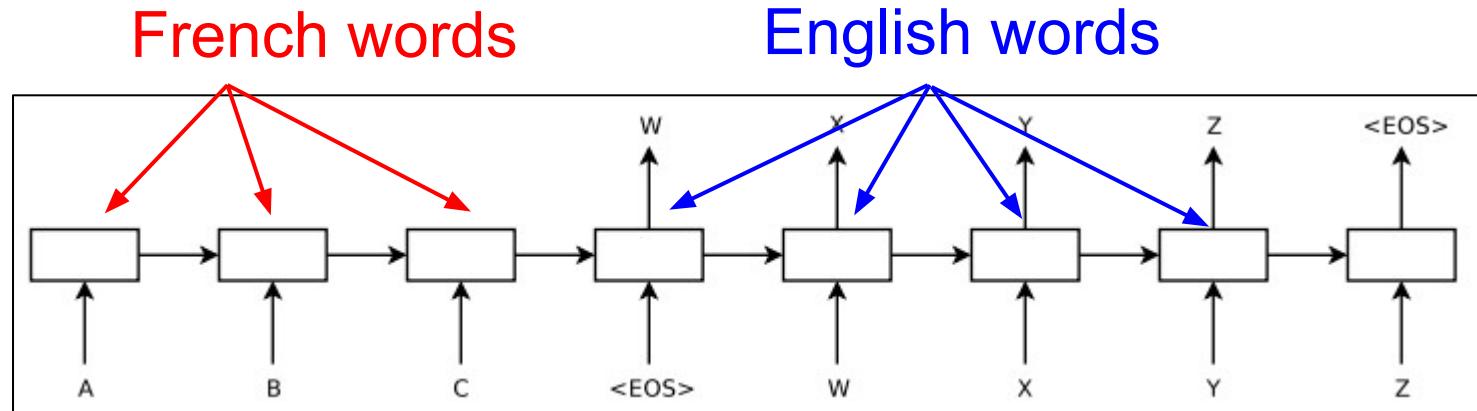
Word-level language model. Similar to:



**Recurrent Neural Network Based Language Model**  
[Tomas Mikolov, 2010]

# Recurrent Networks are good at modeling sequences...

## Machine Translation model



**Sequence to Sequence Learning with Neural Networks**  
[Ilya Sutskever, Oriol Vinyals, Quoc V. Le, 2014]

Suppose we had the training sentence “cat sat on mat”

We want to train a **language model**:

$P(\text{next word} \mid \text{previous words})$

i.e. want these to be high:

$P(\text{cat} \mid [\langle S \rangle])$

$P(\text{sat} \mid [\langle S \rangle, \text{cat}])$

$P(\text{on} \mid [\langle S \rangle, \text{cat}, \text{sat}])$

$P(\text{mat} \mid [\langle S \rangle, \text{cat}, \text{sat}, \text{on}])$

$P(\langle E \rangle \mid [\langle S \rangle, \text{cat}, \text{sat}, \text{on}, \text{mat}])$

Suppose we had the training sentence “cat sat on mat”

We want to train a **language model**:

$P(\text{next word} \mid \text{previous words})$

First, suppose we had only a finite, 1-word history:  
i.e. want these to be high:

$P(\text{cat} \mid \langle S \rangle)$

$P(\text{sat} \mid \text{cat})$

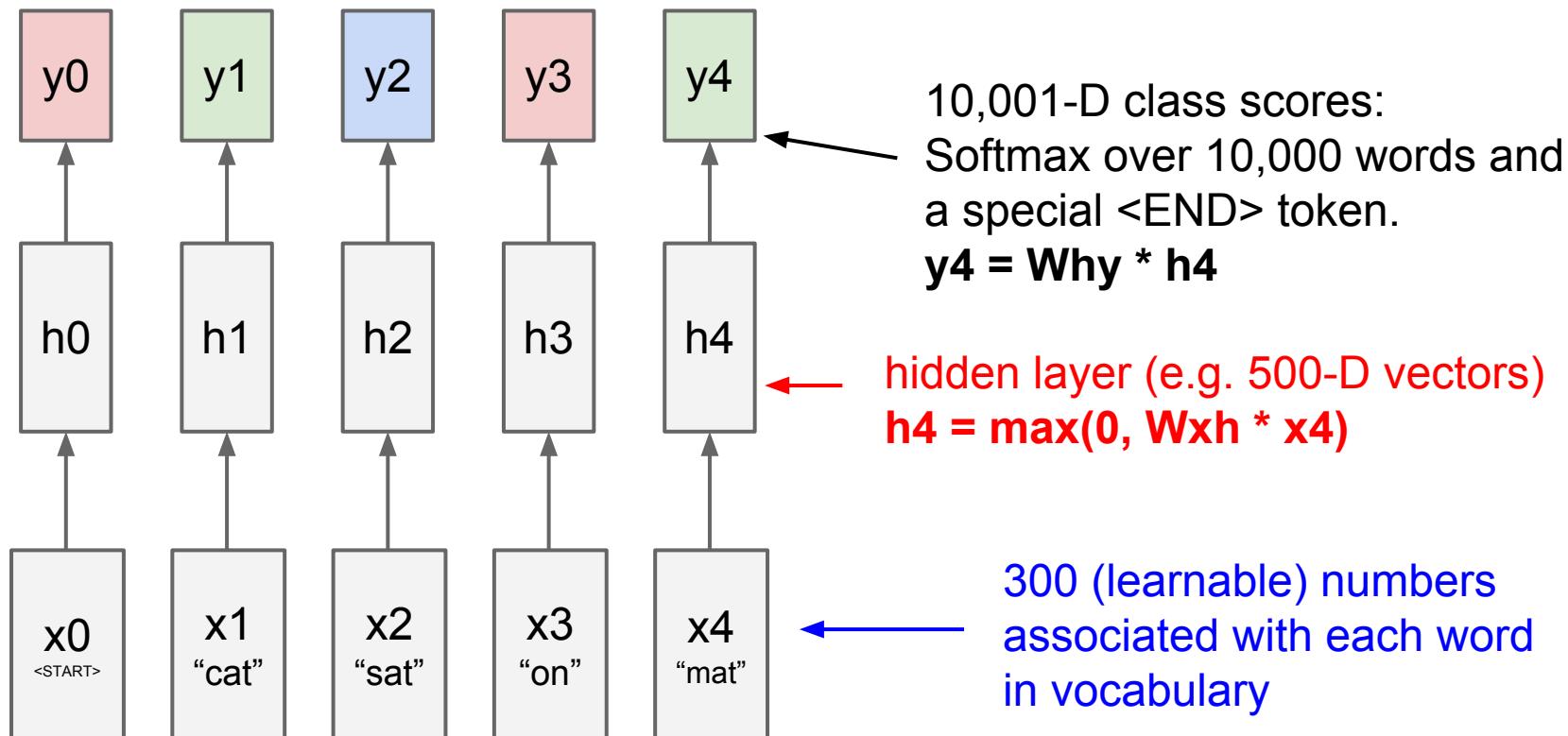
$P(\text{on} \mid \text{sat})$

$P(\text{mat} \mid \text{on})$

$P(\langle E \rangle \mid \text{mat})$

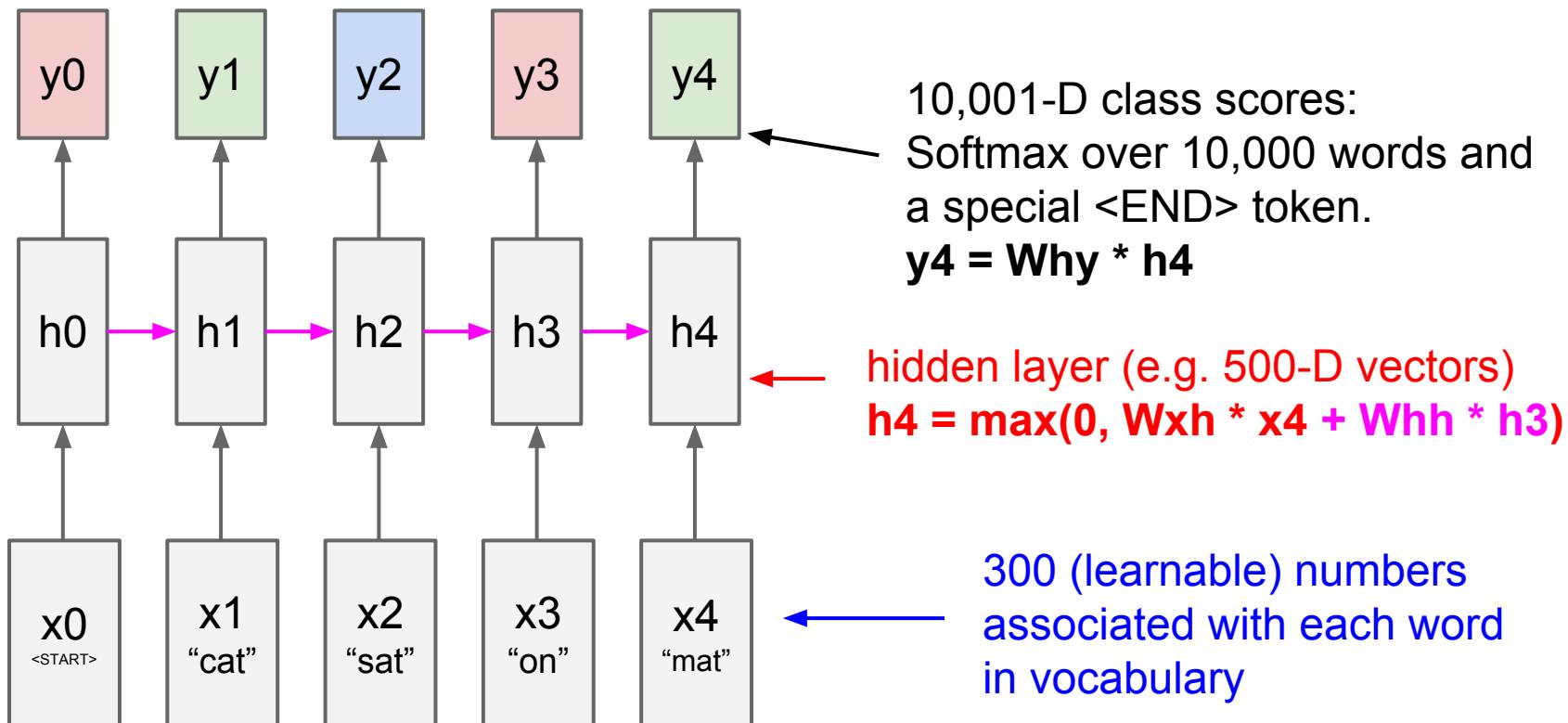
Vanilla 2-layer classification net for each word given previous word:

“cat sat on mat”



# Turn it into RNN: (#anticlimatic)

“cat sat on mat”



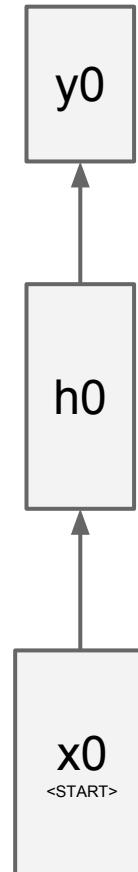
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



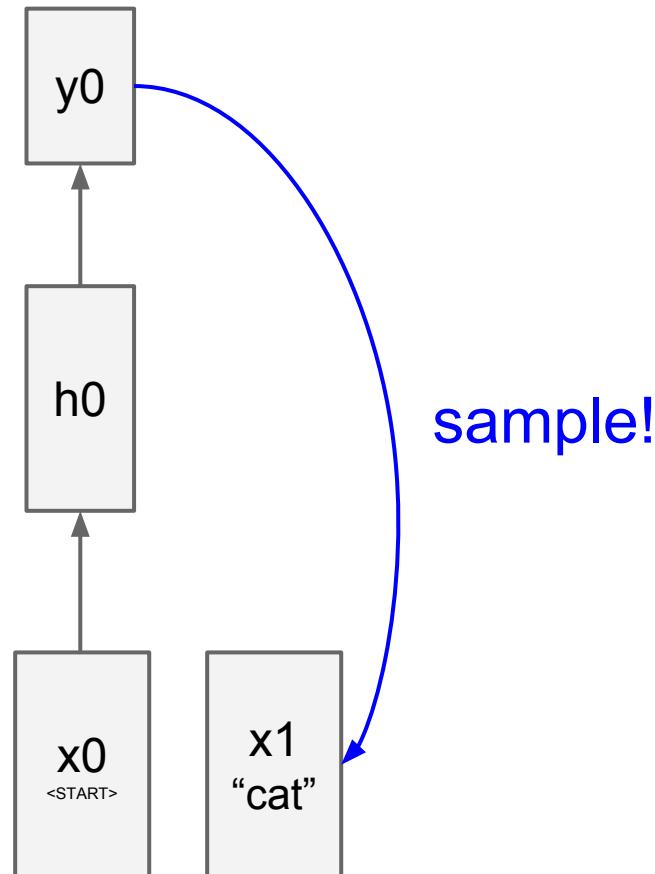
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



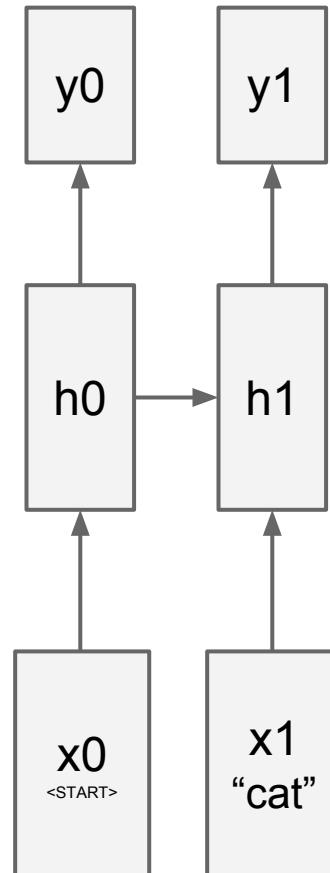
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



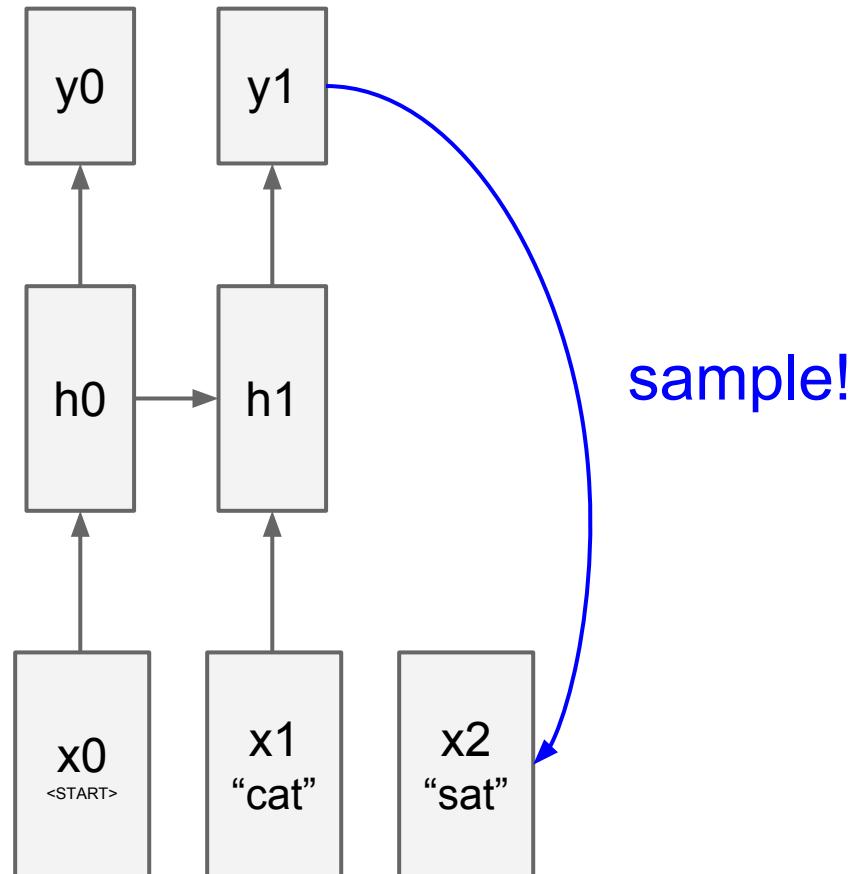
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



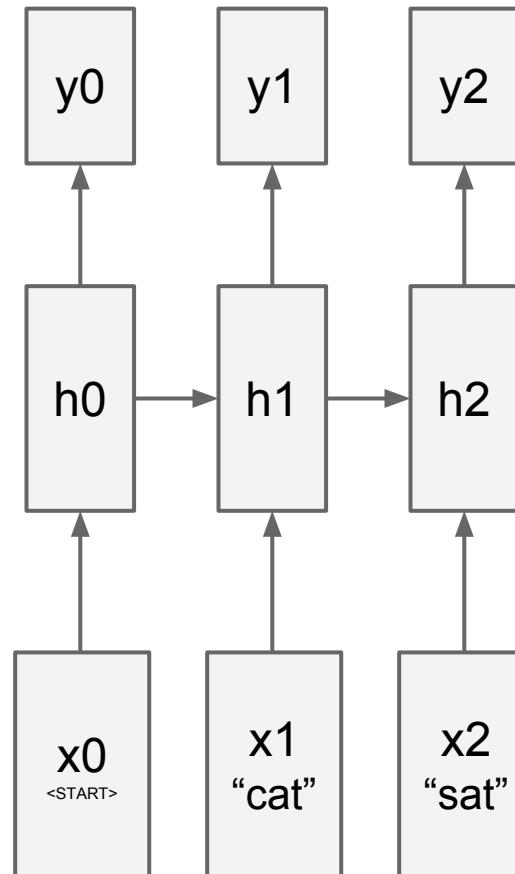
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



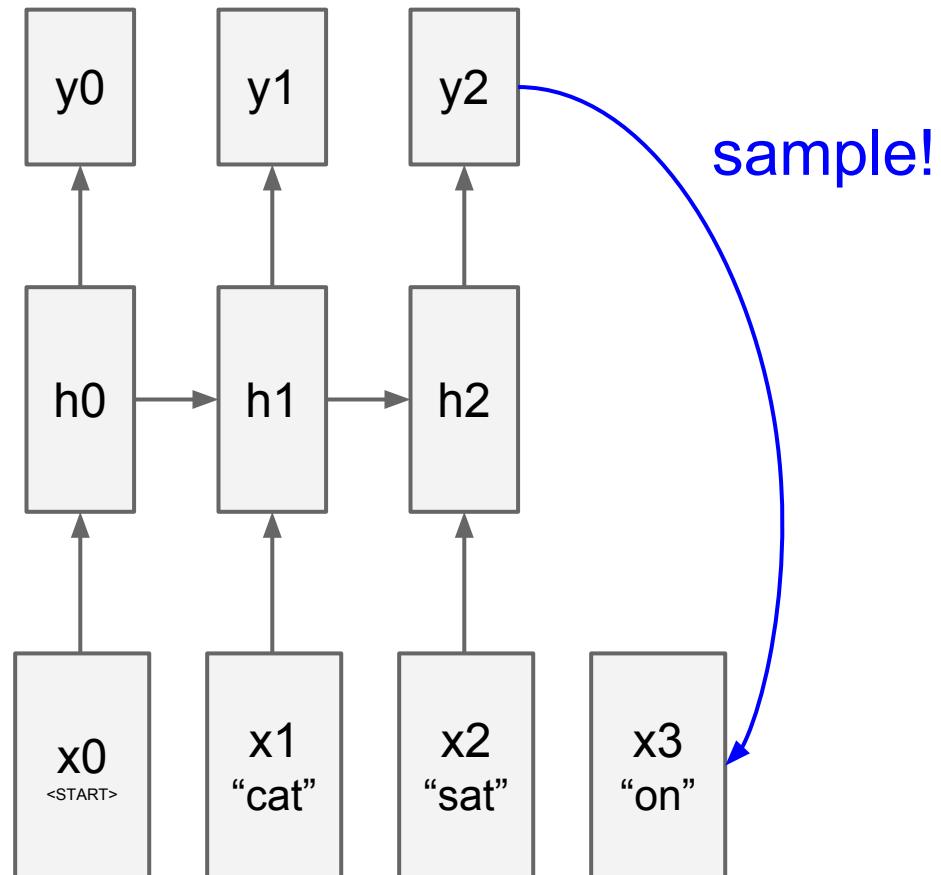
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



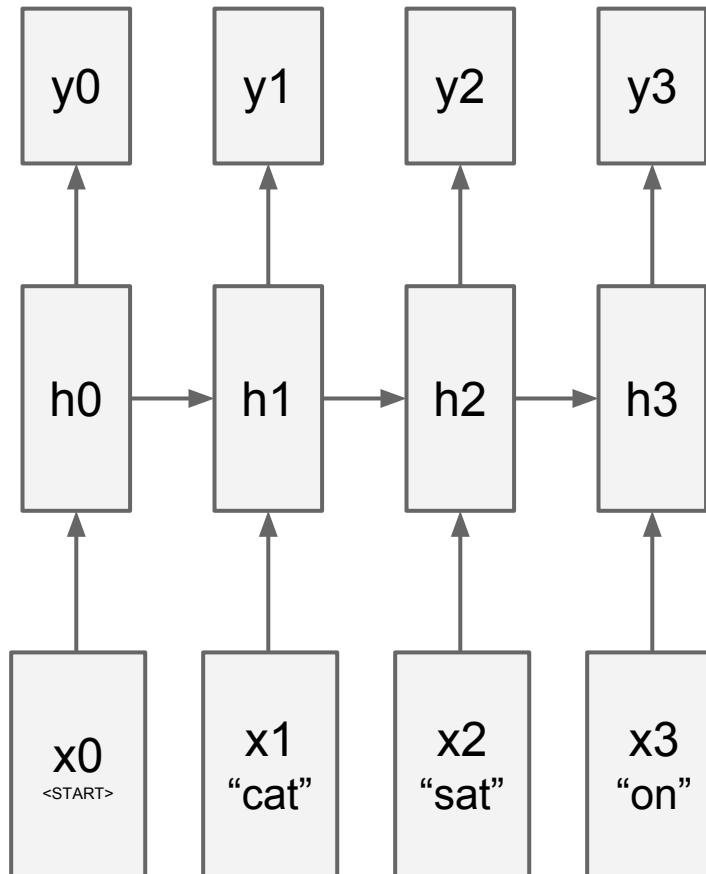
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



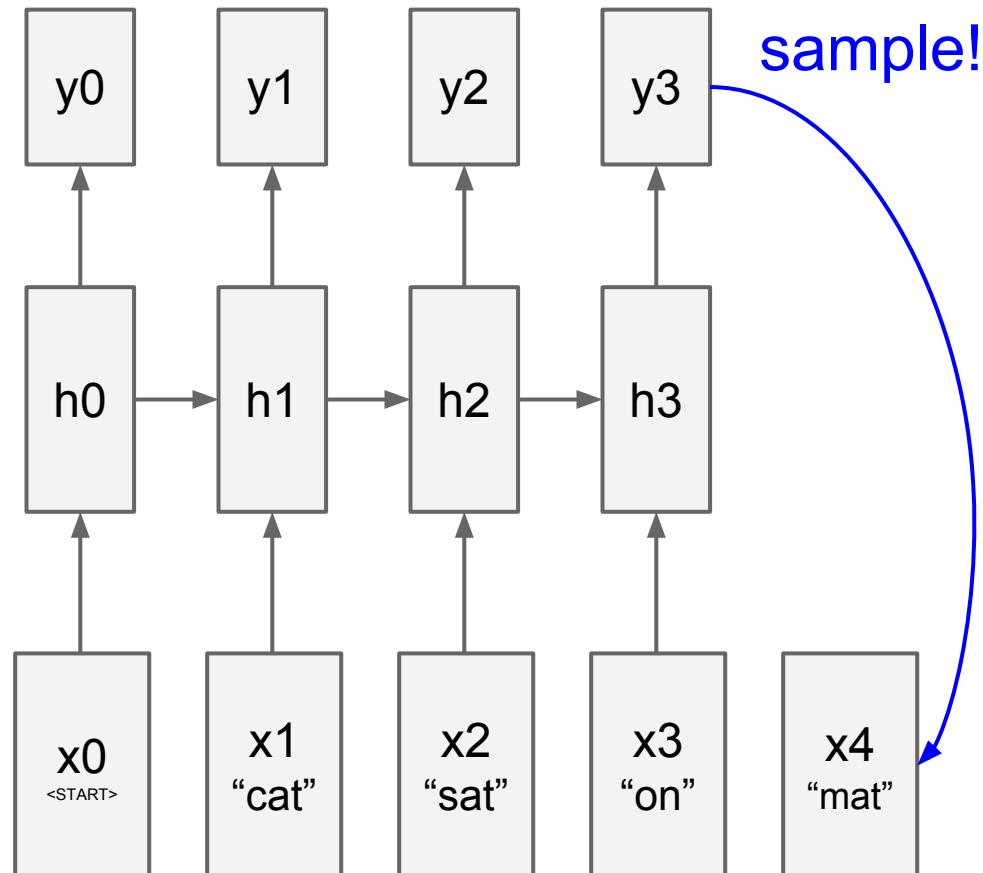
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



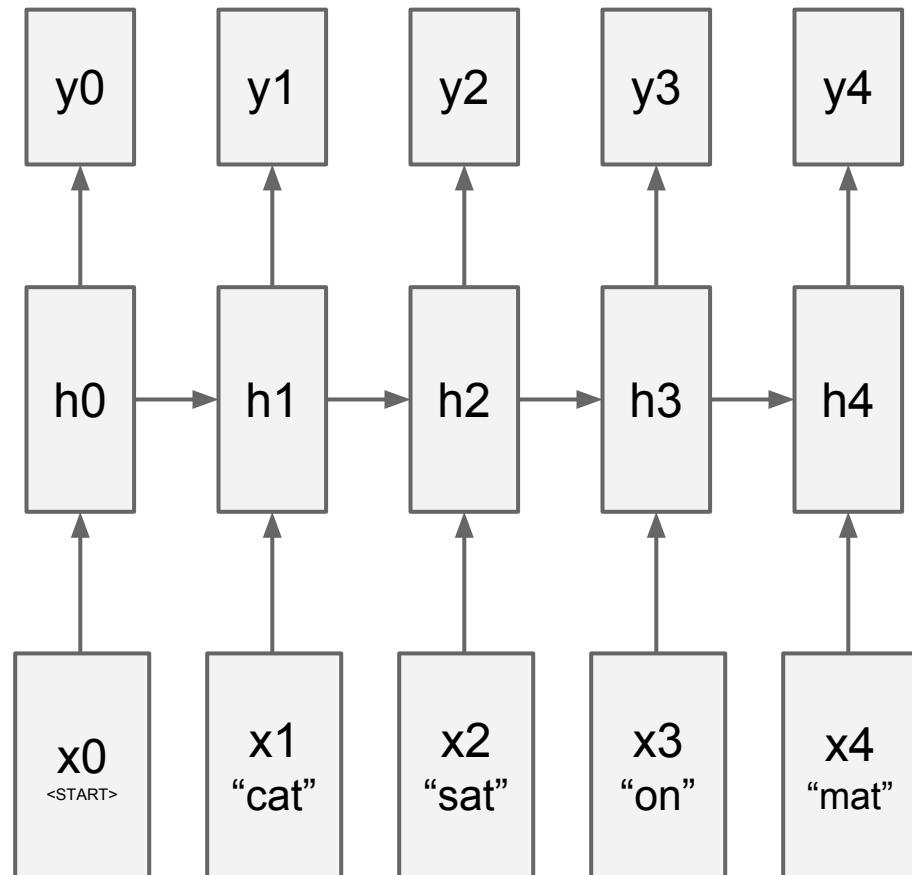
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$



Training this on a lot of sentences would give us a language model. A way to predict

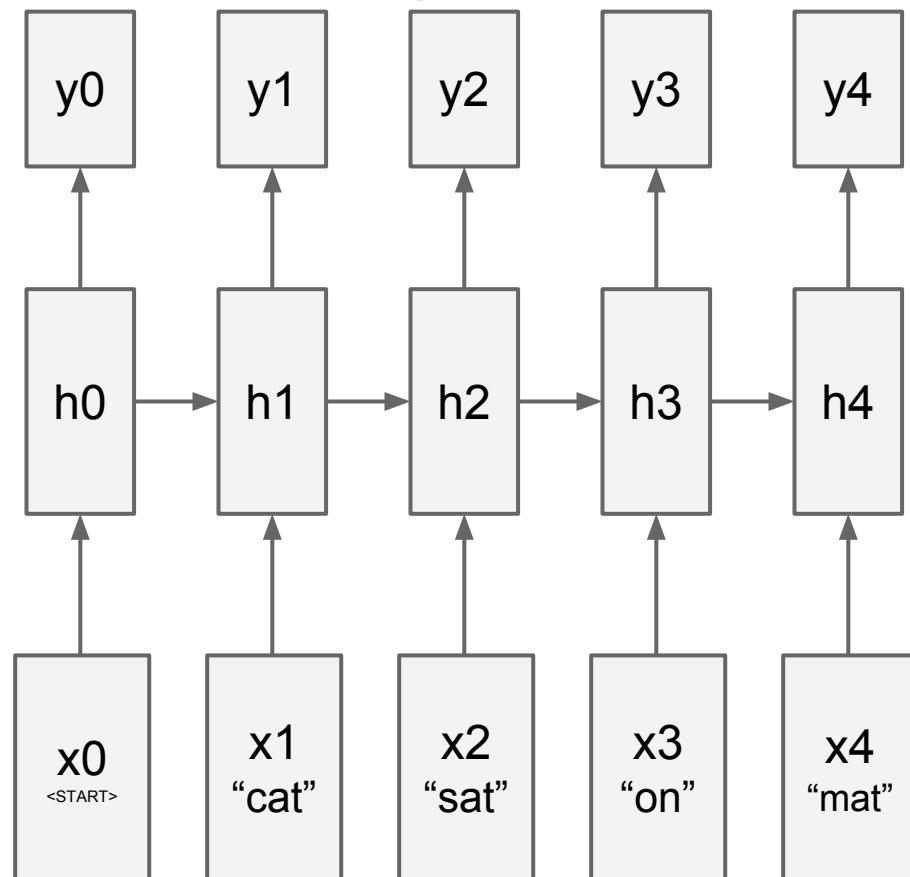
$P(\text{next word} \mid \text{previous words})$



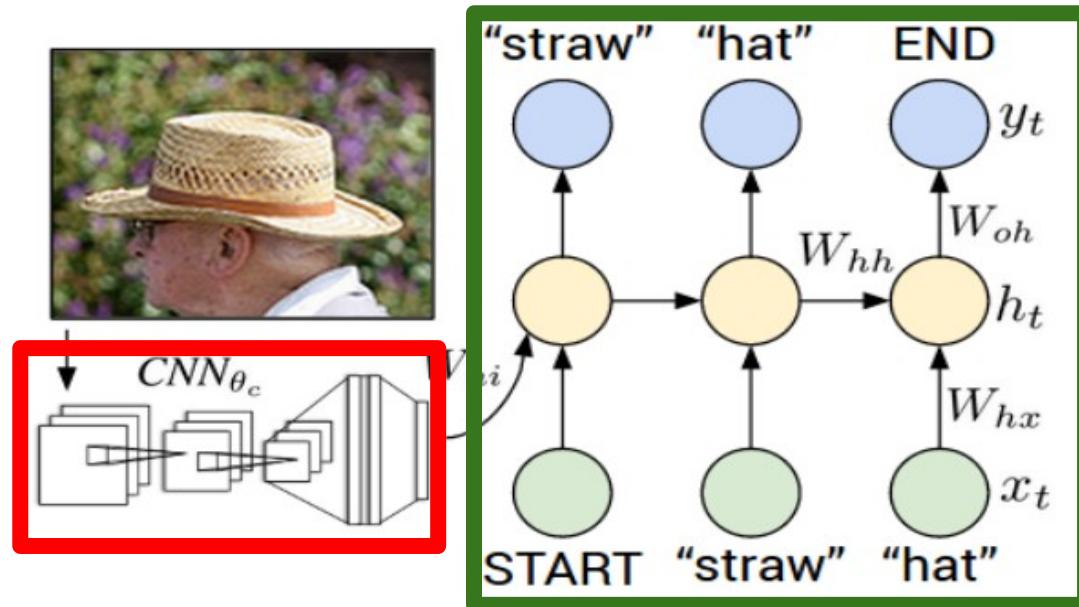
Training this on a lot of sentences would give us a language model. A way to predict

$P(\text{next word} \mid \text{previous words})$

samples <END>? done.



# Recurrent Neural Network



## Convolutional Neural Network

image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

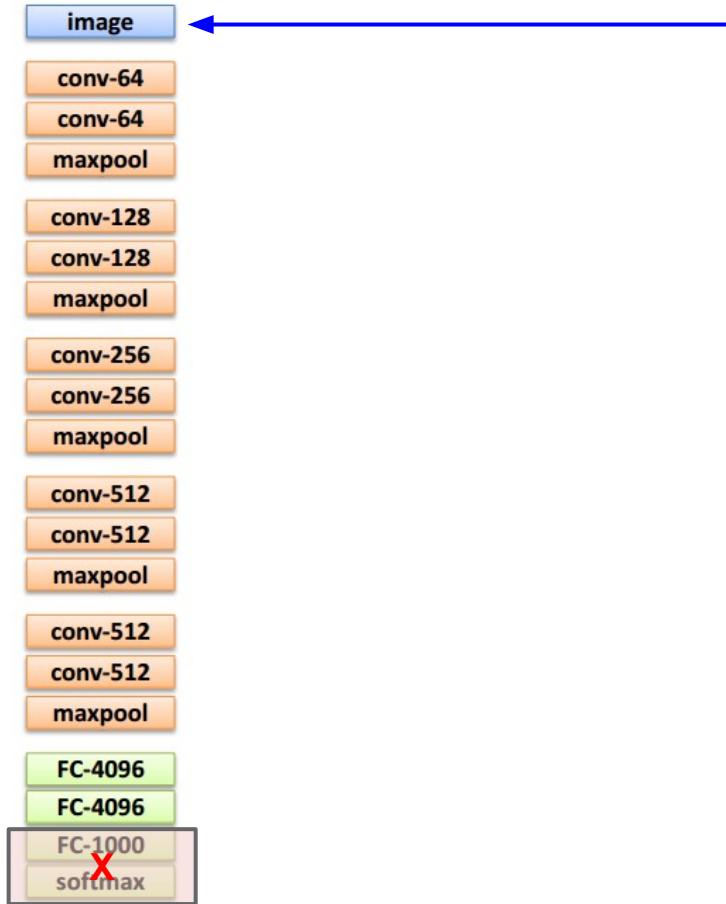
FC-1000

softmax



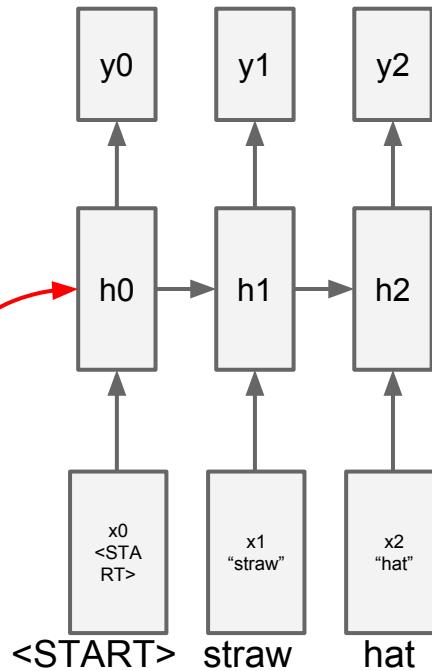
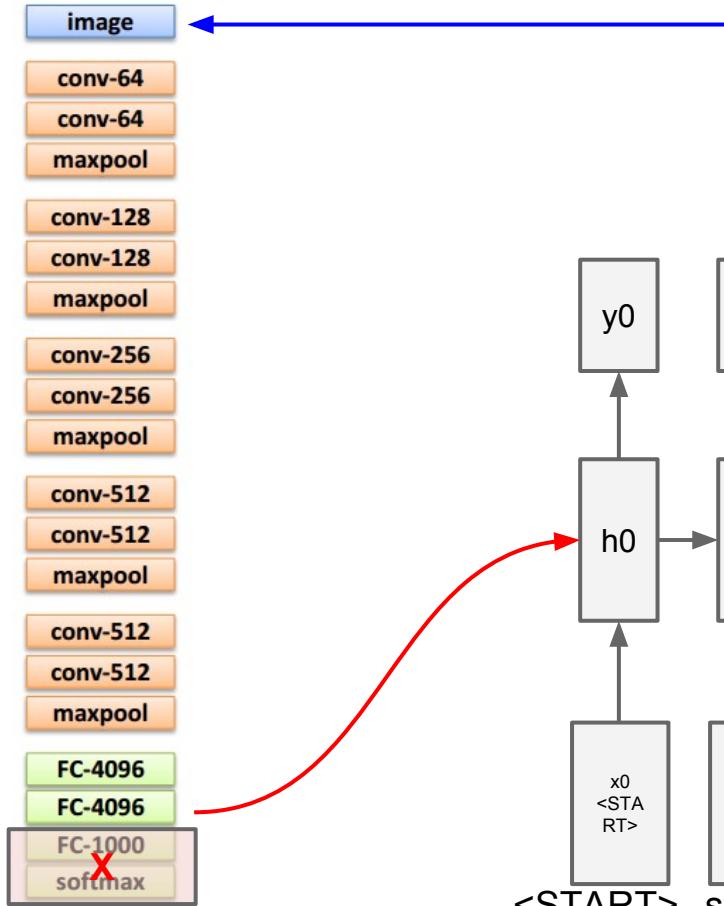
“straw hat”

training example

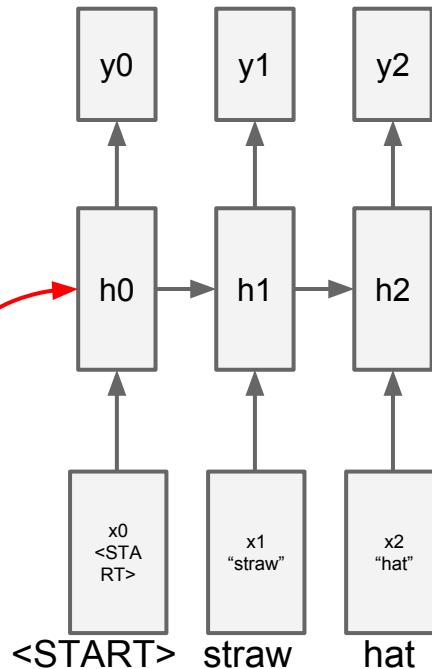
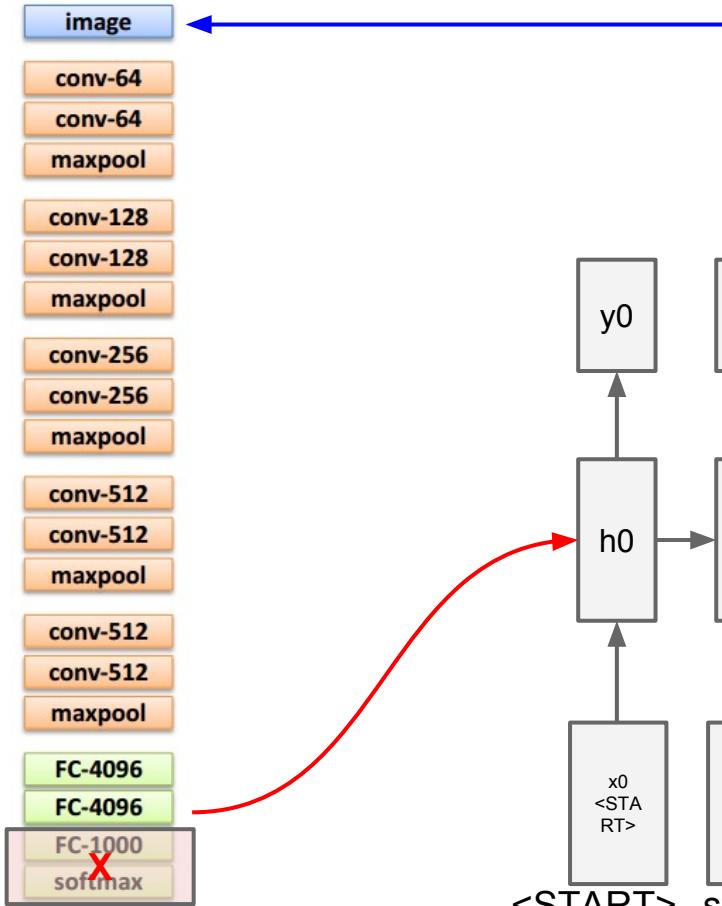


“straw hat”

training example



training example



training example

before:

$$h_0 = \max(0, W_{xh} * x_0)$$

now:

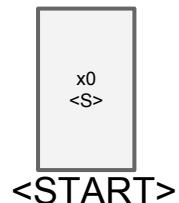
$$h_0 = \max(0, W_{xh} * x_0 + W_{ih} * v)$$

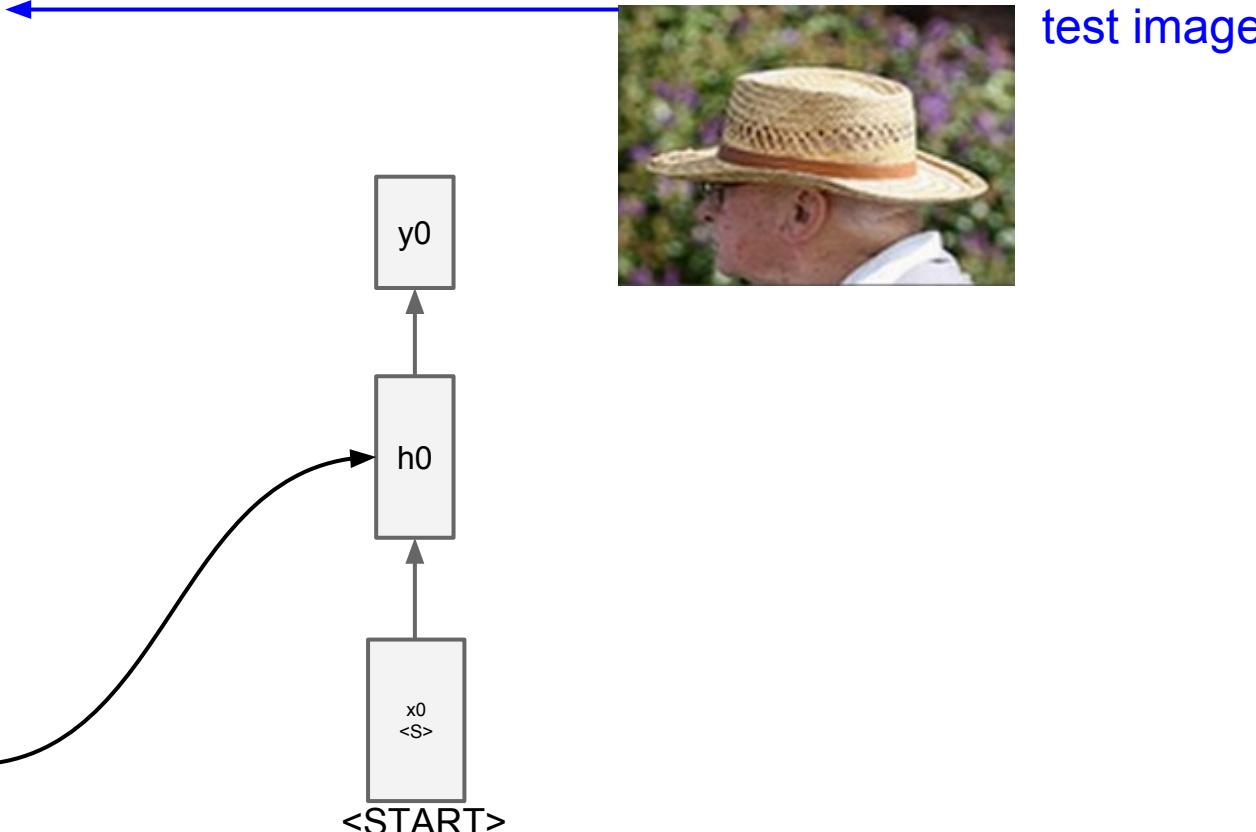
test image

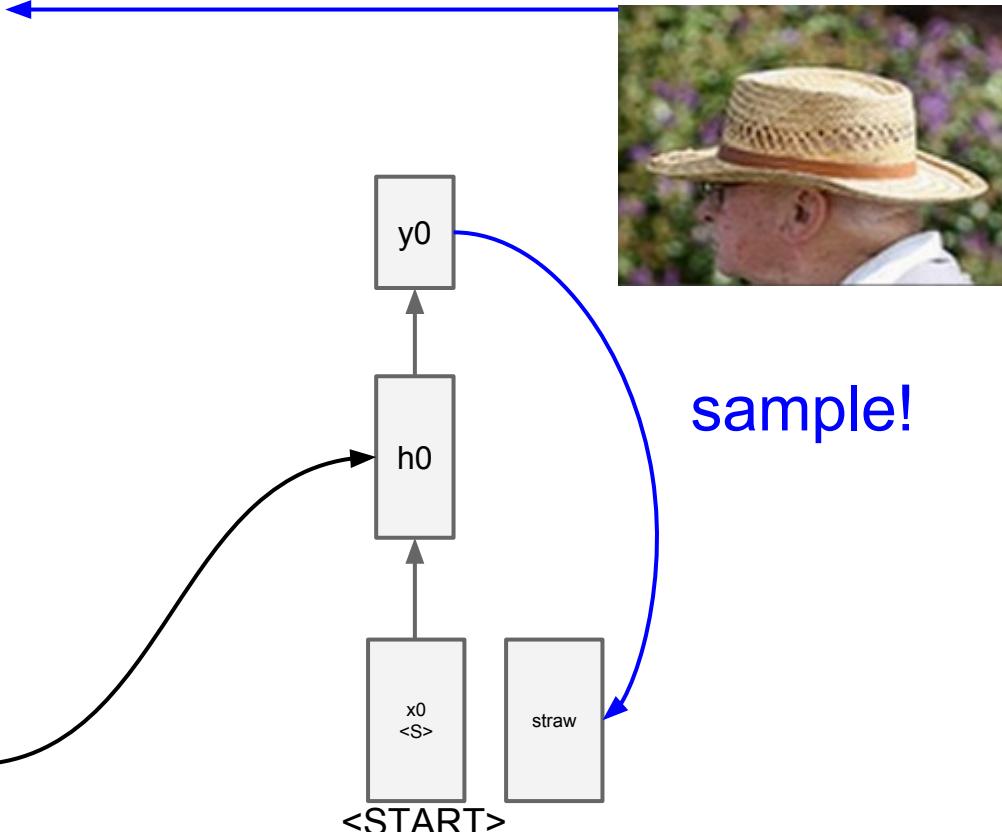


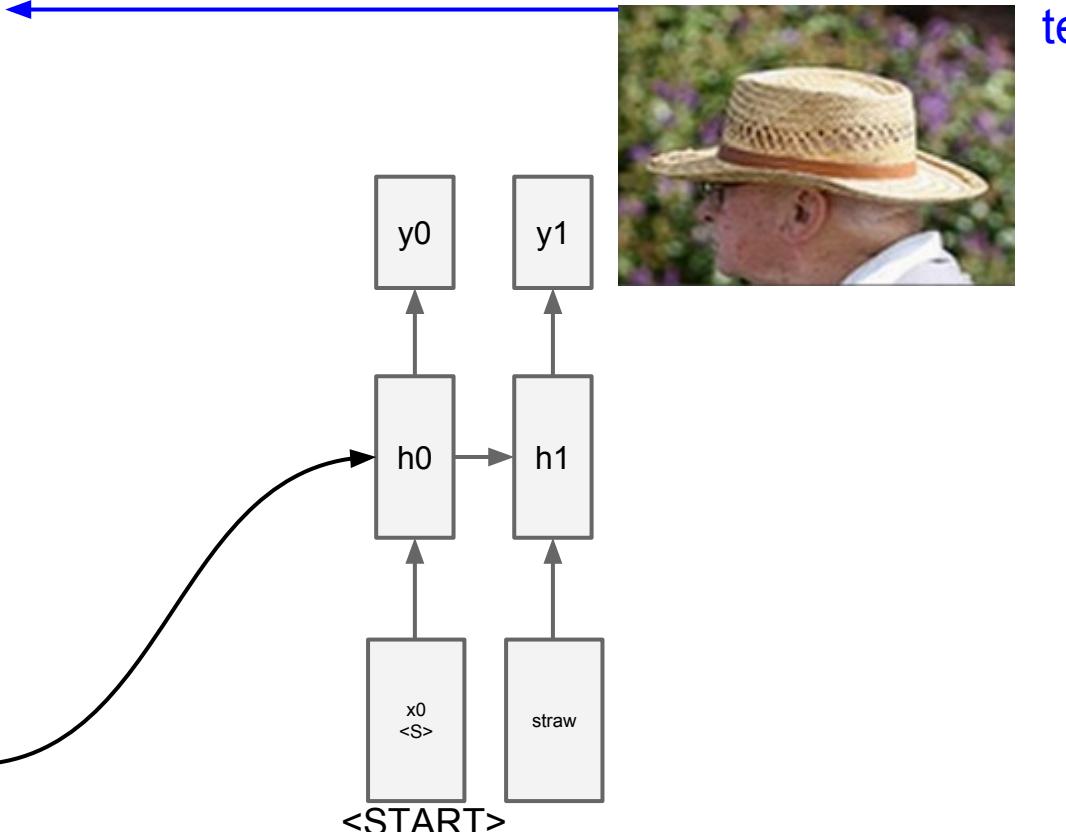


test image

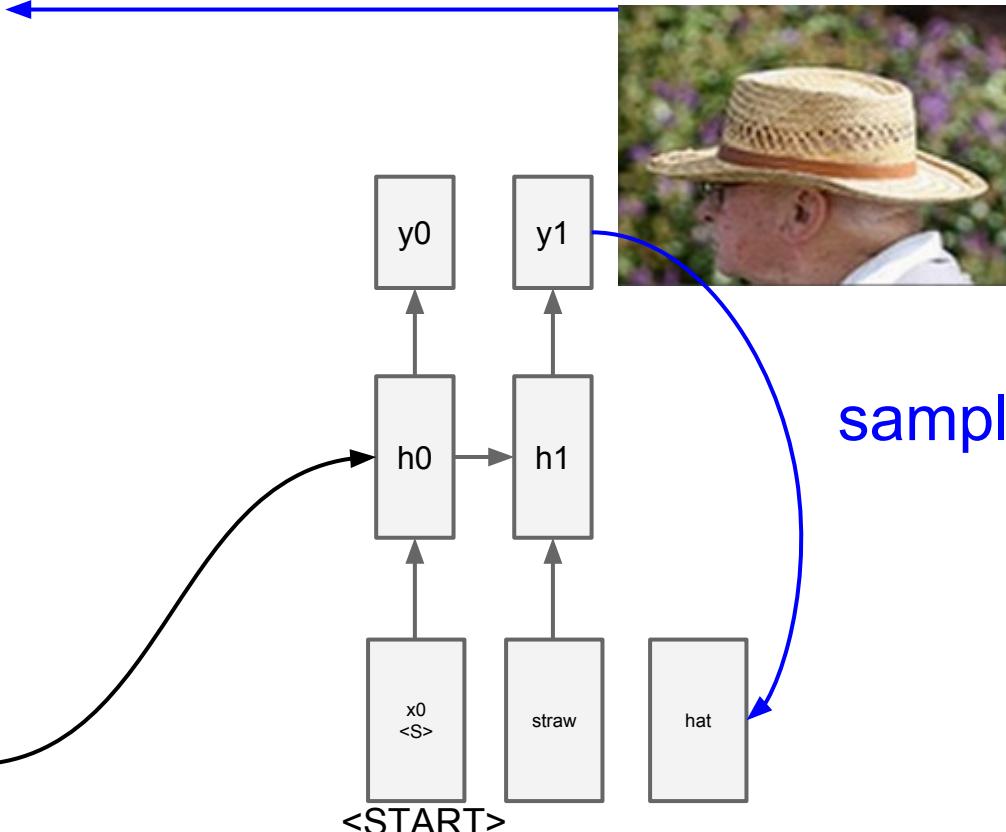






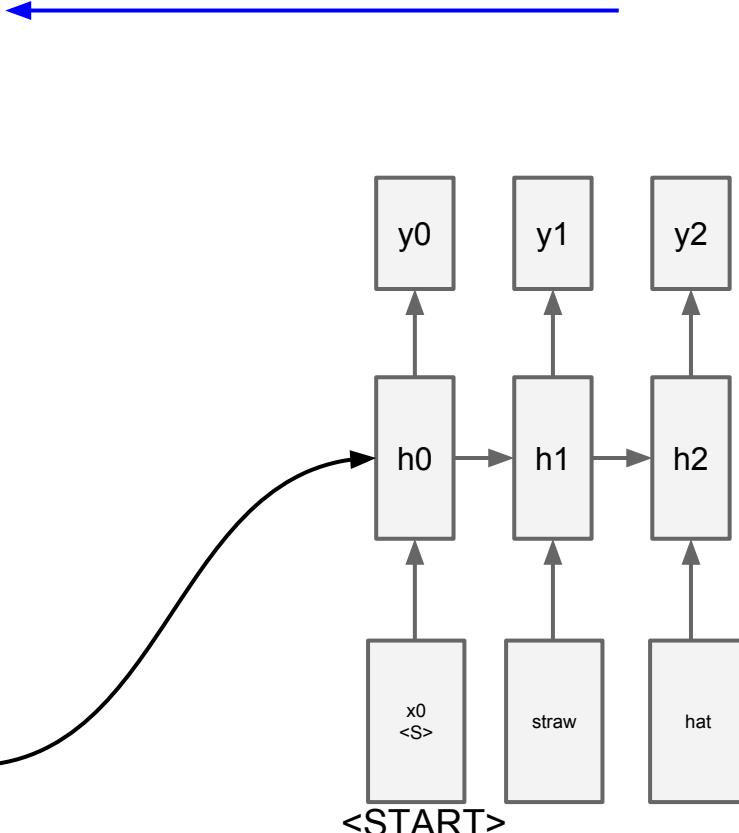


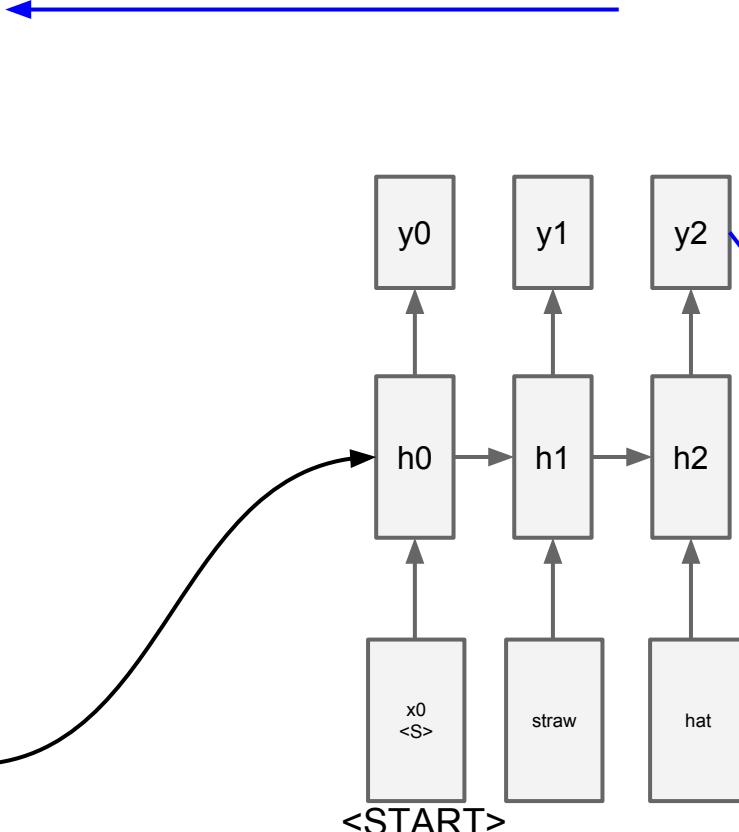
test image



test image

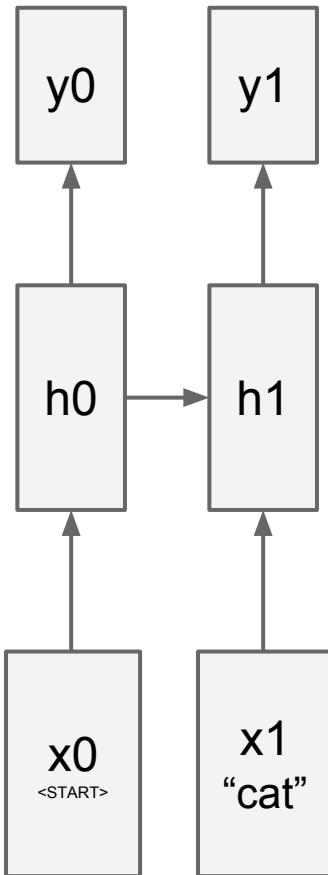
sample!





test image

sample!  
 <END> token  
 => finish.



# RNN vs. LSTM

“hidden” representation  
(e.g. 200 numbers)

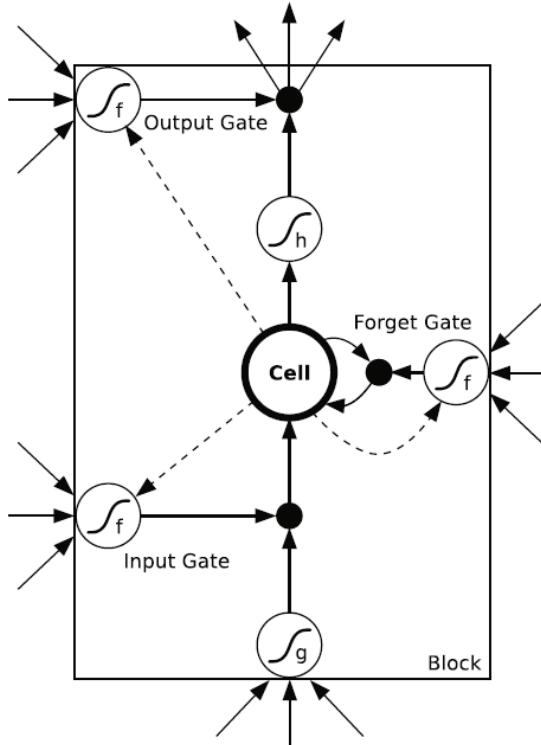
$$\mathbf{h}_1 = \max(0, \mathbf{W}_{xh} * \mathbf{x}_1 + \mathbf{Whh} * \mathbf{h}_0)$$

- LSTM just changes the form of the equation for  $\mathbf{h}$  such that:
1. more expressive multiplicative interactions
  2. gradients flow nicer
  3. network can explicitly decide to reset the hidden state

# RNN vs. LSTM

- Use LSTM over RNN
- Do not be intimidated by pictures that try to draw an LSTM, e.g.:

*(it's just a particular funny form of forward function, backpropagation as usual)*



$$\text{LSTM} : h_t^{l-1}, h_{t-1}^l, c_{t-1}^l \rightarrow h_t^l, c_t^l$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} T_{2n,4n} \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

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

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

Recurrent Neural Network Regularization  
[Zaremba, Sutskever, Vinyals]

# Image Sentence Datasets

a man riding a bike on a dirt path through a forest.

bicyclist raises his fist as he rides on desert dirt trail.

this dirt bike rider is smiling and raising his fist in triumph.

a man riding a bicycle while pumping his fist in the air.

a mountain biker pumps his fist in celebration.



Microsoft COCO

*[Tsung-Yi Lin et al. 2014]*

[mscoco.org](http://mscoco.org)

currently:

~120K images

~5 sentences each

# Wow I can't believe that worked



a group of people standing around a room with remotes  
logprob: -9.17



a young boy is holding a baseball bat  
logprob: -7.61



a cow is standing in the middle of a street  
logprob: -8.84

# Well, I can kind of see it



a baby laying on a bed with a stuffed bear  
logprob: -8.66

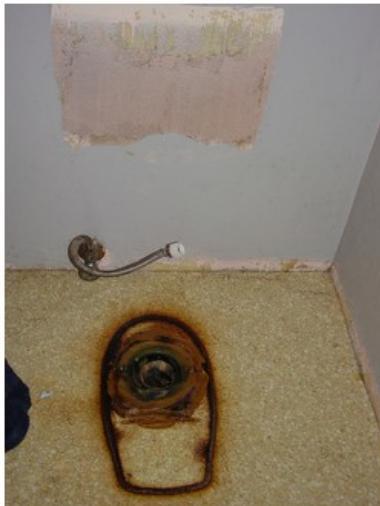


a young boy is holding a  
baseball bat  
logprob: -7.65



a cat is sitting on a couch with a remote control  
logprob: -12.45

# Not sure what happened there...



a toilet with a seat up in a bathroom  
logprob: -13.44



a woman holding a teddy bear in front of a mirror  
logprob: -9.65

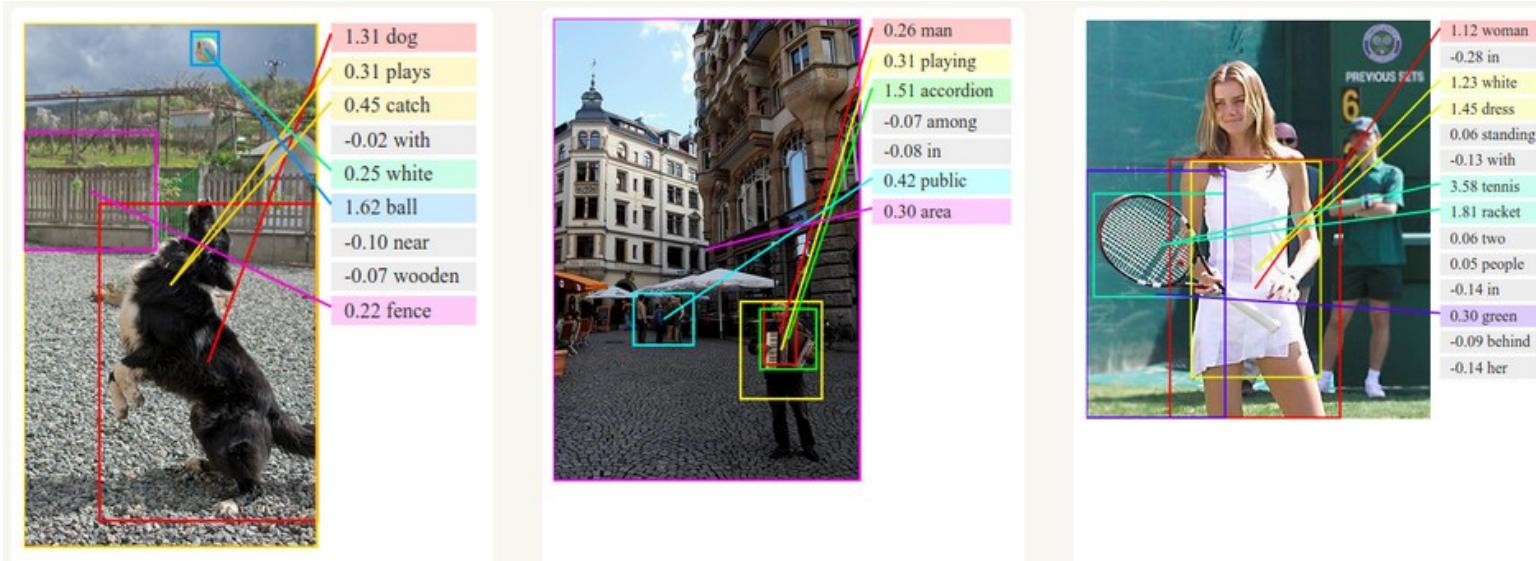


a horse is standing in the middle of a road  
logprob: -10.34

More examples in Web demo: <http://bit.ly/neuraltalkdemo>

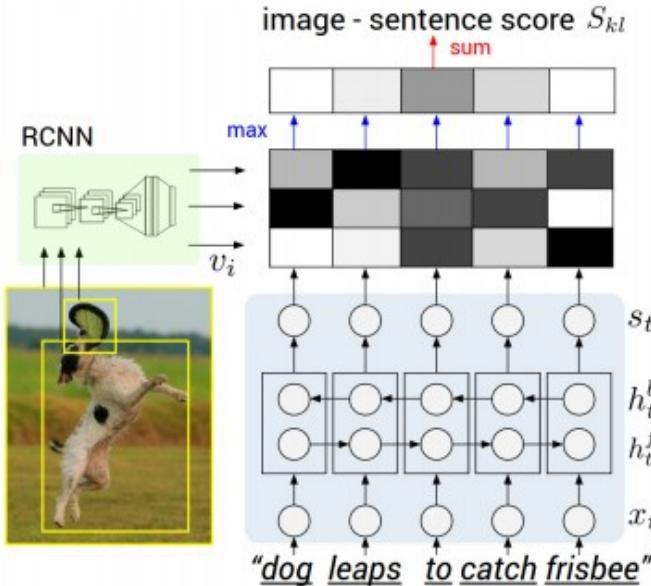
# Ranking and Retrieval

Each example is a query test sentence, the most likely retrieved image & the grounding:

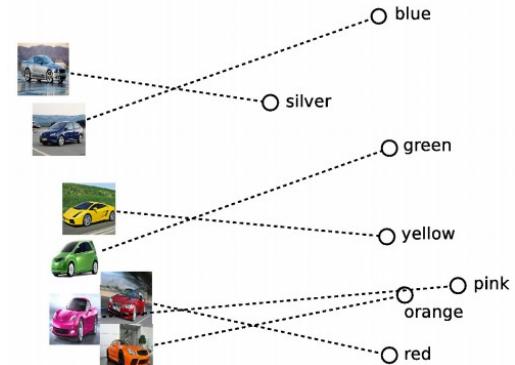
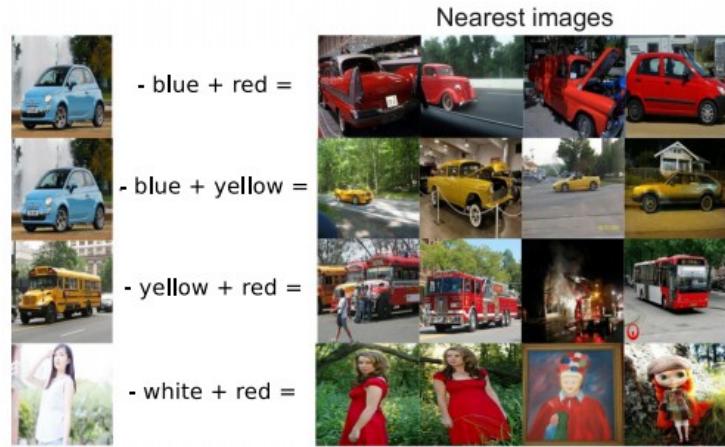
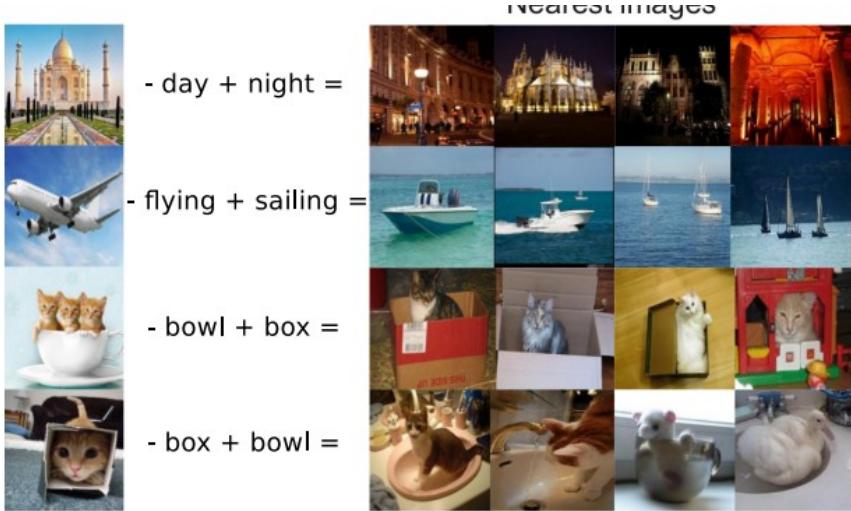
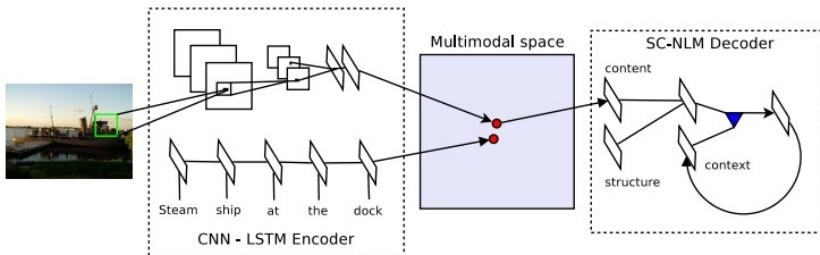


More examples in Web demo: <http://bit.ly/rankingdemo>

# Ranking and Retrieval

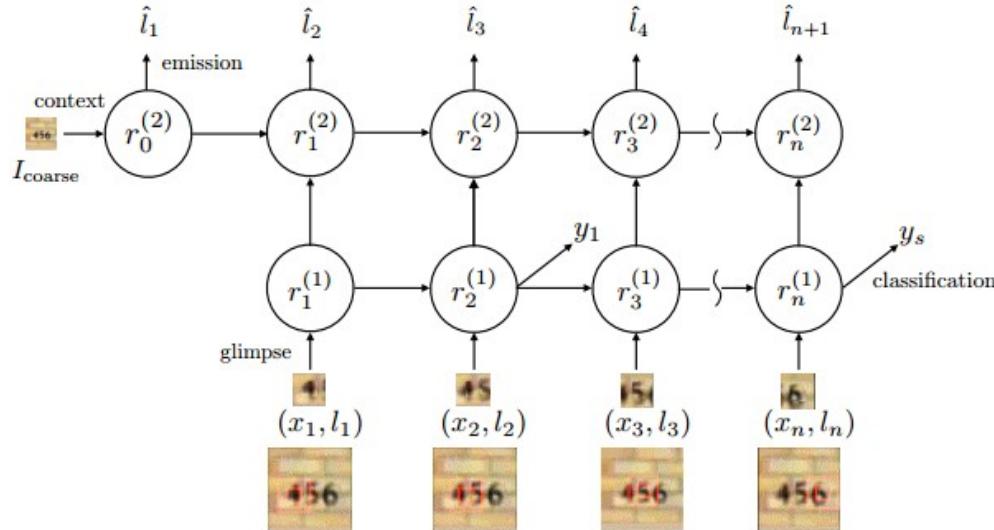


Deep Visual-Semantic Alignments for Generating Image Descriptions  
[Karpathy and Fei-Fei, 2015]



**Unifying Visual-Semantic Embeddings with Multimodal Neural Language Models**  
*[Kiros, Salakhutdinov, Zemel, 2014]*

# Recurrent Attention Models



web demo

<http://www.psi.toronto.edu/~jimmy/dram/>

also DRAW: <https://www.youtube.com/watch?v=Zt7MI9eKEo>

**Multiple Object Recognition with Visual Attention**  
[Jimmy Lei Ba, Volodymyr Mnih, Koray Kavukcuoglu], 2014

# Summary:

- We looked at many Computer Vision tasks beyond Image Classification and how they are addressed with Convolutional Neural Networks

Next: Guest Lecture:  
**Evan Shelhamer**  
*Caffe*

