

Course material at:

<https://github.com/lyeskhalil/mlbootcamp2022/>

# Advanced Topics

Friday

FASE ML Bootcamp

Based on material from various sources:

Data augmentation, Hyperparameter tuning, CNN understanding, Attention, GANs:

<http://cs231n.stanford.edu/schedule.html>

Uncertainty in Deep Learning: <https://www.gatsby.ucl.ac.uk/~balaji/balaji-uncertainty-talk-cifar-dlrl.pdf>

Graph Neural Networks: <http://tkipf.github.io/misc/SlidesCambridge.pdf>

Reinforcement Learning: [https://www.davidsilver.uk/wp-content/uploads/2020/03/deep\\_rl\\_tutorial\\_small\\_compressed.pdf](https://www.davidsilver.uk/wp-content/uploads/2020/03/deep_rl_tutorial_small_compressed.pdf)

2022 Free PDF book on all of ML by Kevin Murphy (Google): <https://probml.github.io/pml-book/book1.html>

Where we are now...

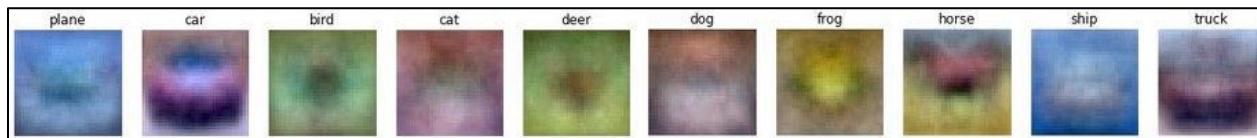
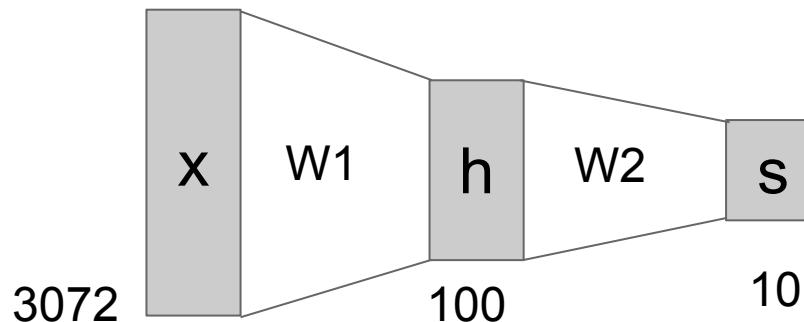
# Neural Networks

Linear score function:

$$f = Wx$$

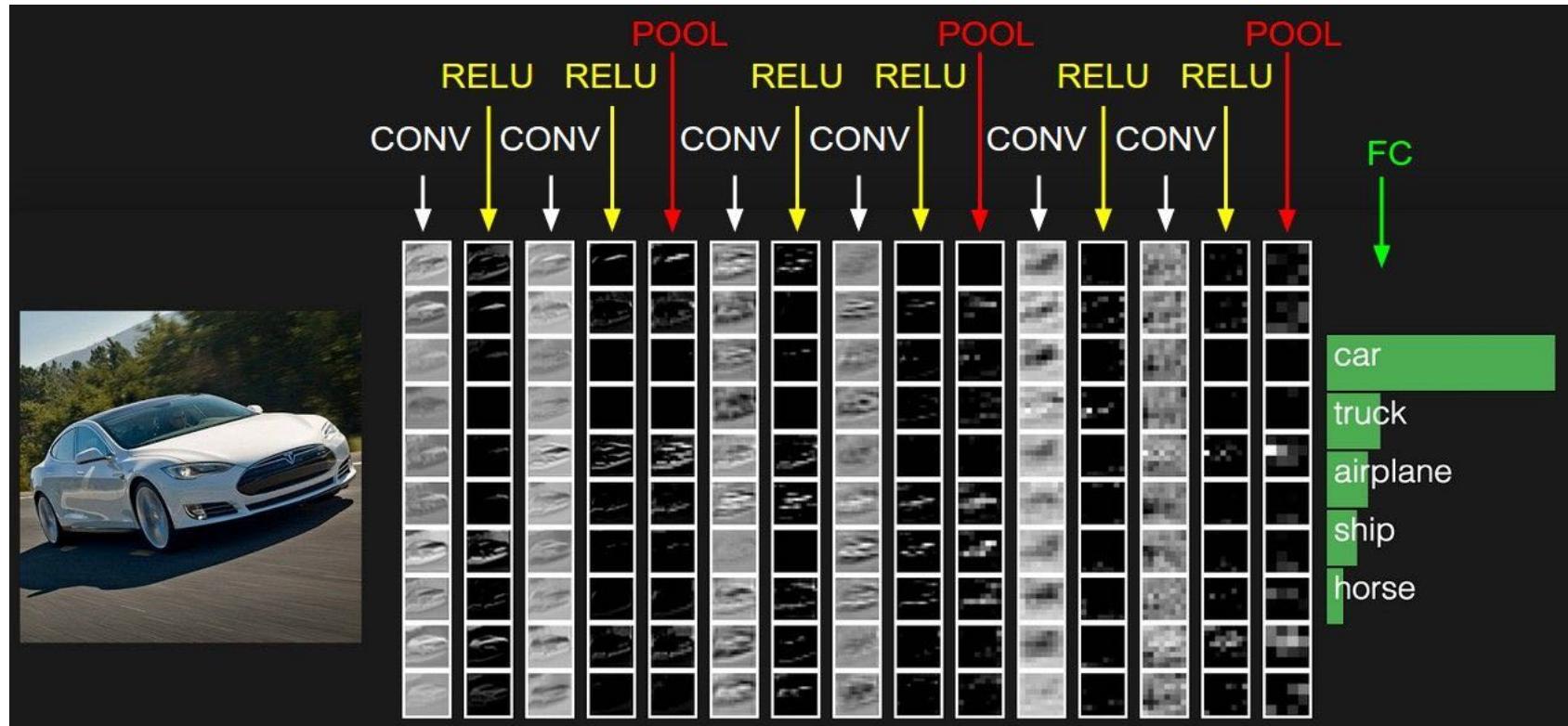
2-layer Neural Network

$$f = W_2 \max(0, W_1 x)$$



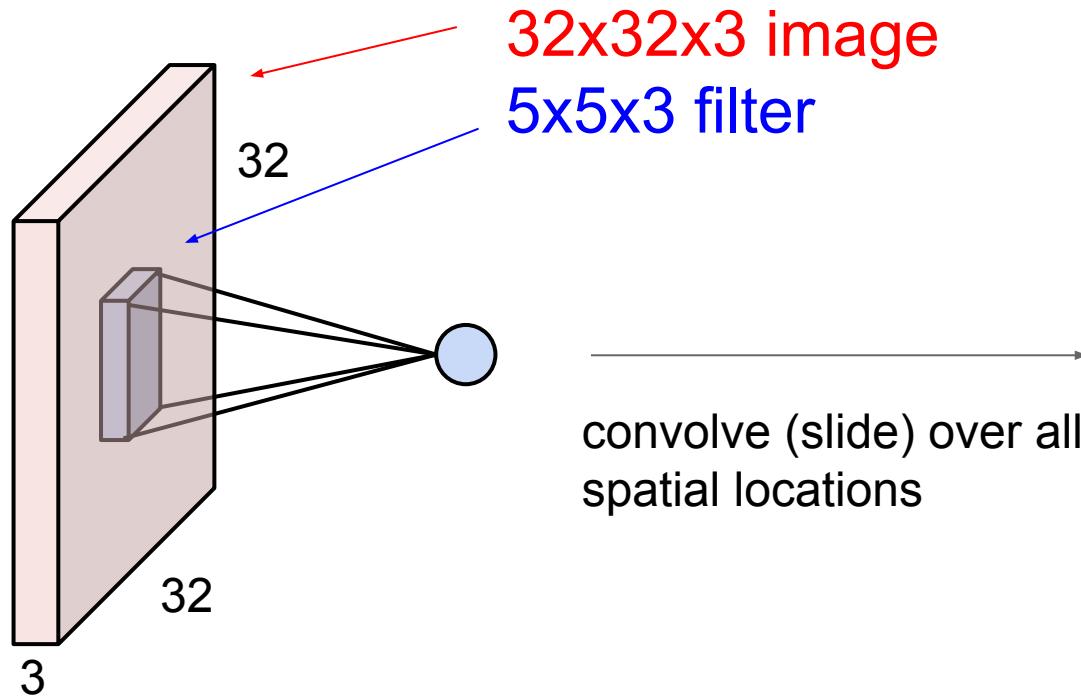
Where we are now...

# Convolutional Neural Networks

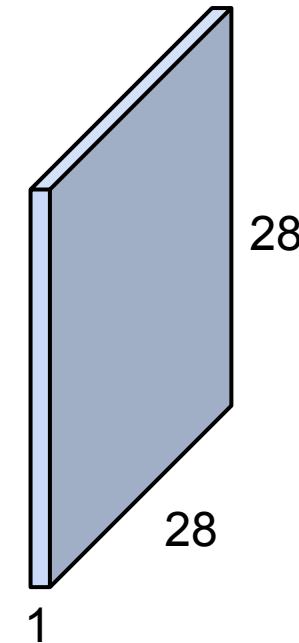


Where we are now...

## Convolutional Layer

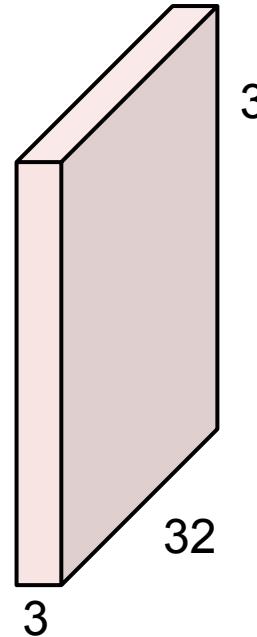


activation map



Where we are now...

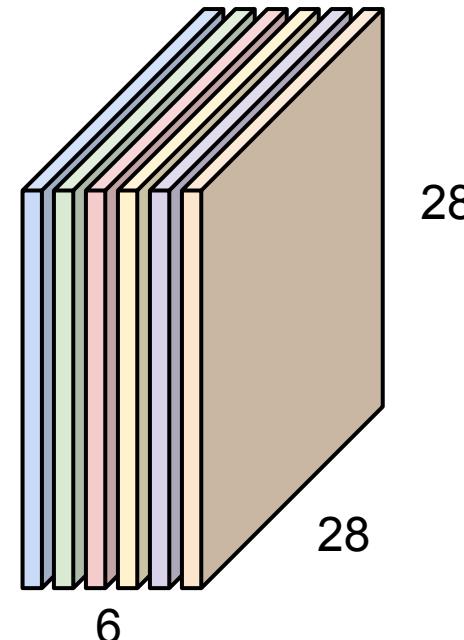
## Convolutional Layer



Convolution Layer

For example, if we had 6  $5 \times 5$  filters, we'll get 6 separate activation maps:

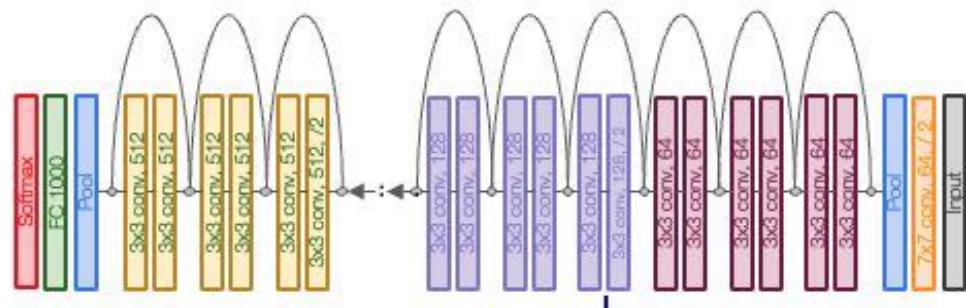
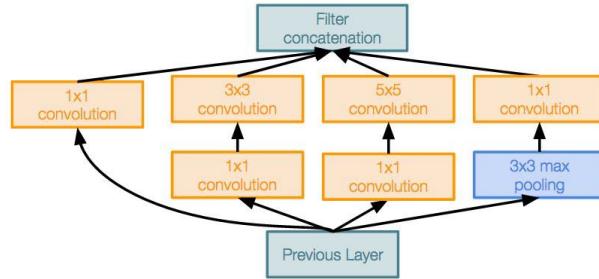
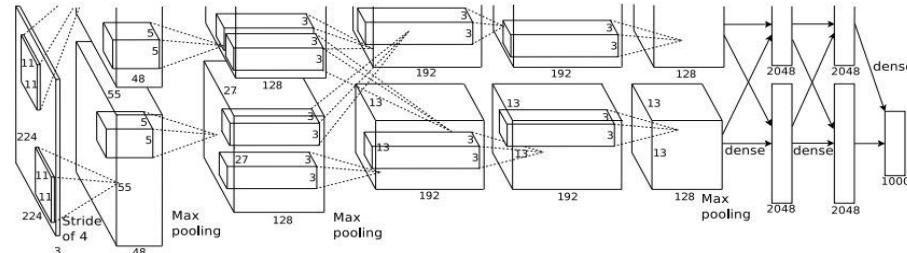
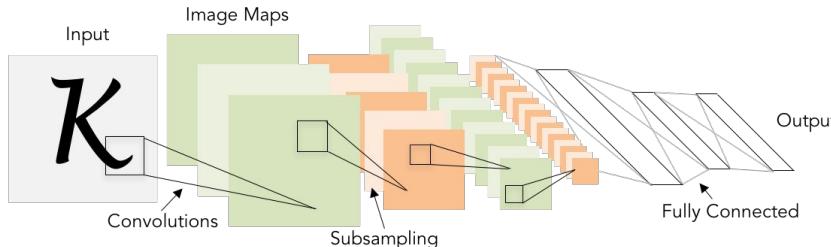
activation maps



We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

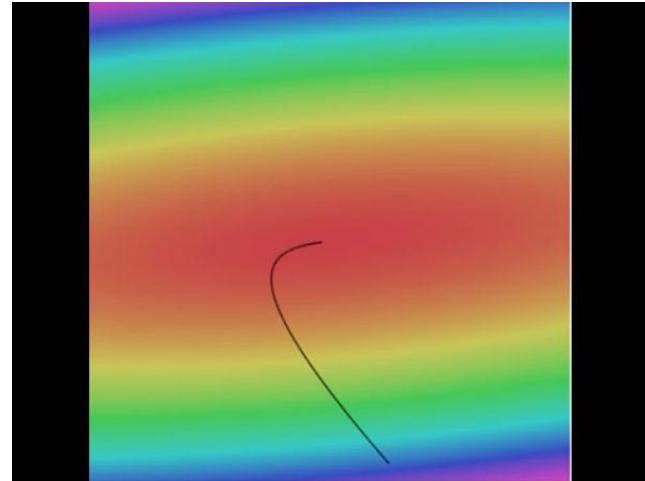
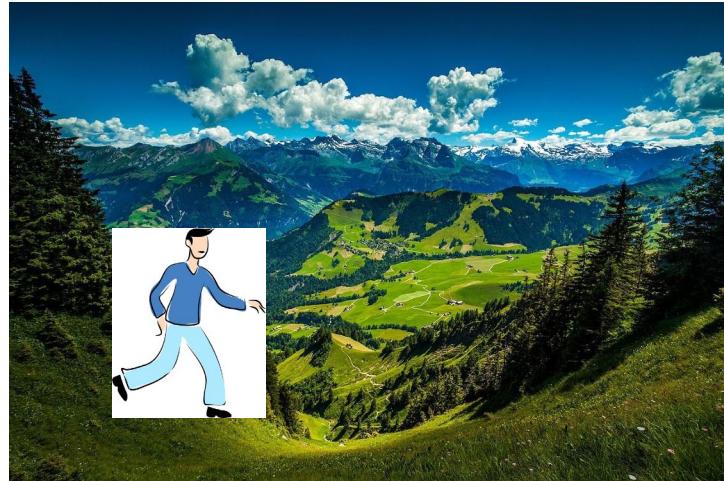
# Where we are now...

## CNN Architectures



Where we are now...

# Learning network parameters through optimization



```
# Vanilla Gradient Descent  
  
while True:  
    weights_grad = evaluate_gradient(loss_fun, data, weights)  
    weights += - step_size * weights_grad # perform parameter update
```

[Landscape image](#) is CC0 1.0 public domain  
[Walking man image](#) is CC0 1.0 public domain

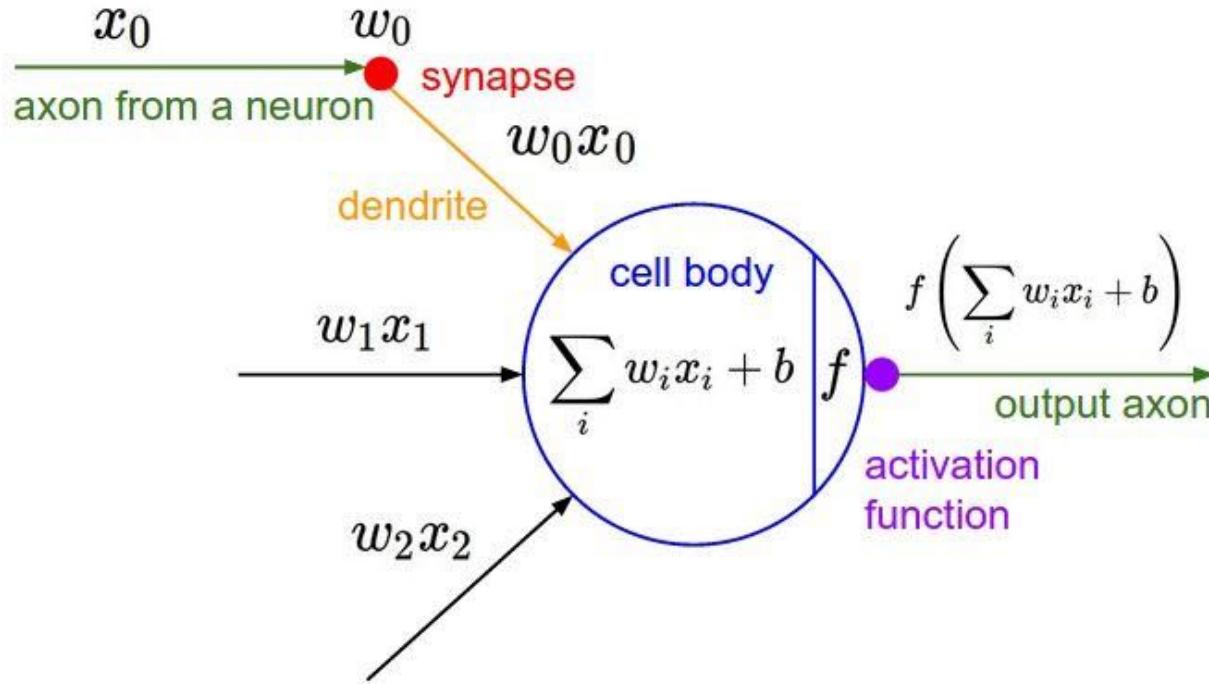
Where we are now...

# Mini-batch SGD

Loop:

1. **Sample** a batch of data
2. **Forward** prop it through the graph  
(network), get loss
3. **Backprop** to calculate the gradients
4. **Update** the parameters using the gradient

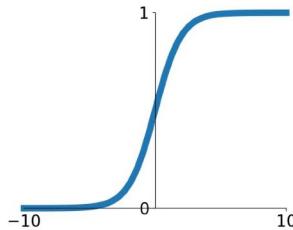
# Activation Functions



# Activation Functions

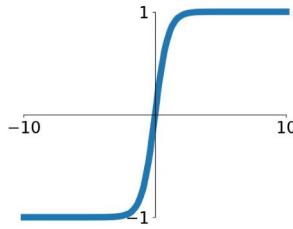
## Sigmoid

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



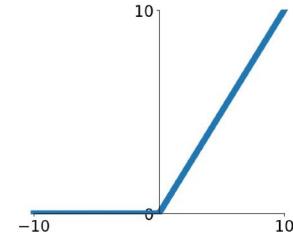
## tanh

$$\tanh(x)$$



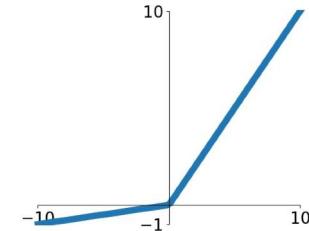
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

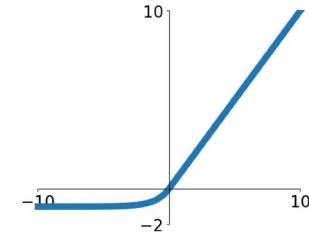


## Maxout

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

## ELU

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



# TLDR: In practice:

- Use ReLU. Be careful with your learning rates
- Try out Leaky ReLU / Maxout / ELU / SELU
  - To squeeze out some marginal gains
- Don't use sigmoid or tanh

# Data Preprocessing

# TLDR: In practice for Images: center only

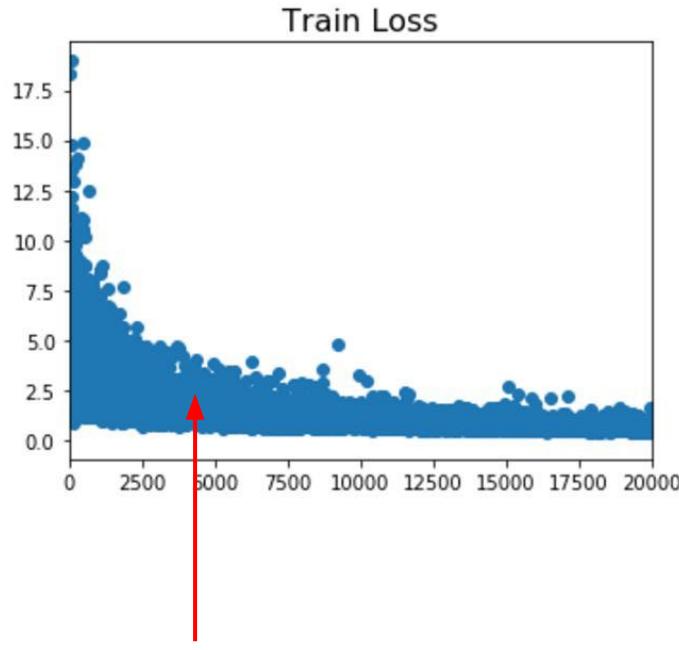
e.g. consider CIFAR-10 example with [32,32,3] images

- Subtract the mean image (e.g. AlexNet)  
(mean image = [32,32,3] array)
- Subtract per-channel mean (e.g. VGGNet)  
(mean along each channel = 3 numbers)
- Subtract per-channel mean and  
Divide by per-channel std (e.g. ResNet)  
(mean along each channel = 3 numbers)

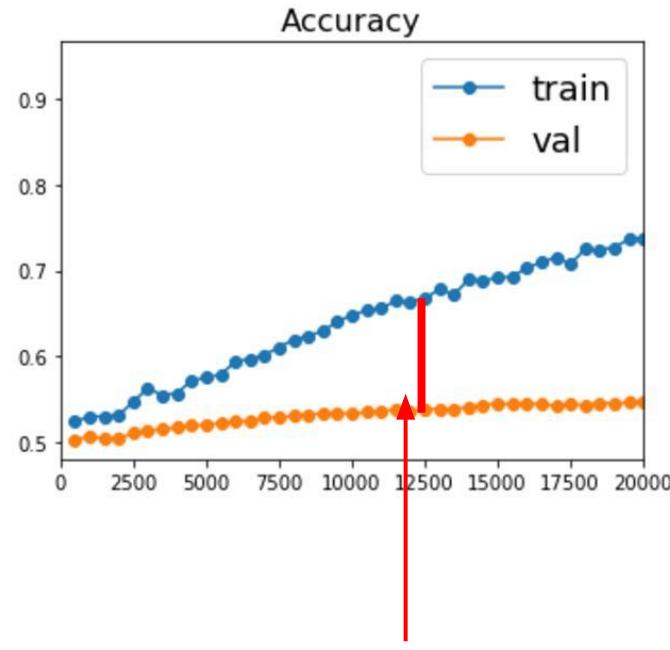
Not common  
to do PCA or  
whitening

# Training vs. Testing Error

# Beyond Training Error

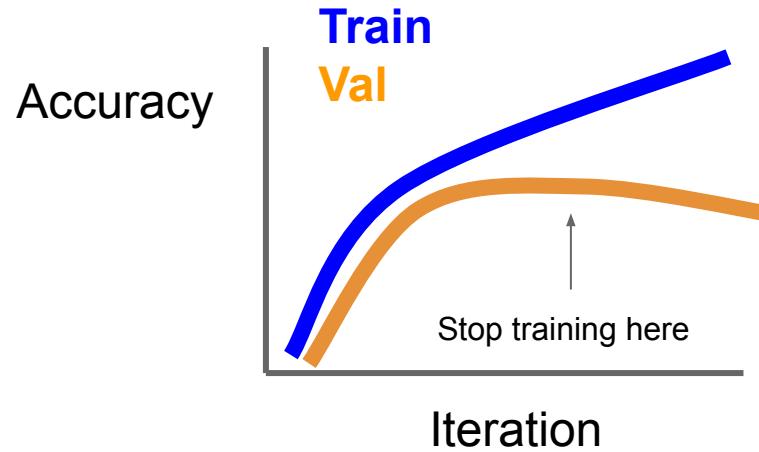
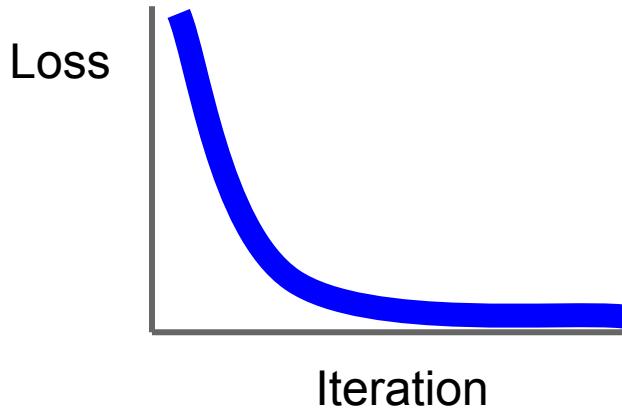


Better optimization algorithms help reduce training loss



But we really care about error on new data - how to reduce the gap?

# Early Stopping: Always do this



Stop training the model when accuracy on the validation set decreases  
Or train for a long time, but always keep track of the model snapshot  
that worked best on val

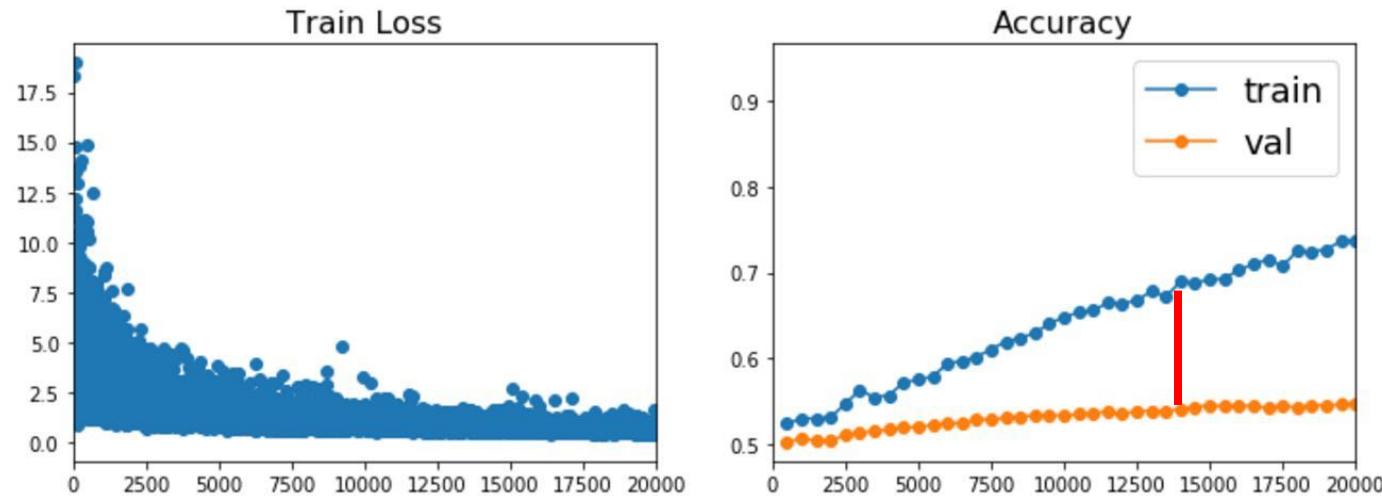
# Model Ensembles

1. Train multiple independent models
2. At test time average their results

(Take average of predicted probability distributions, then choose argmax)

Enjoy 2% extra performance

# How to improve single-model performance?



## Regularization

# Regularization: Add term to loss

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

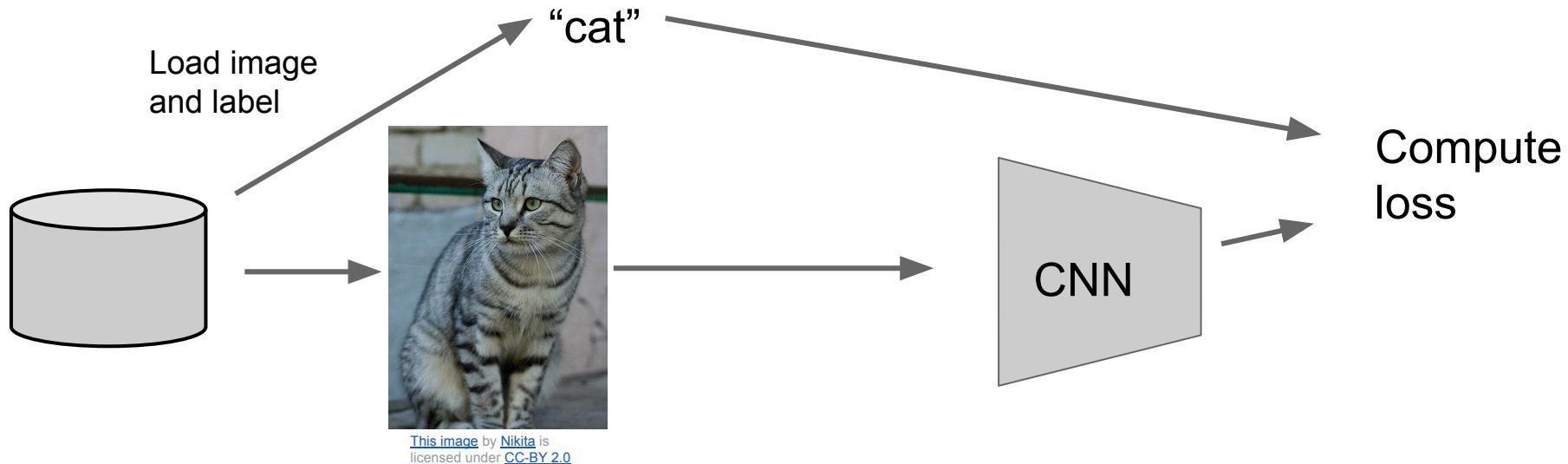
L1 regularization

$$R(W) = \sum_k \sum_l |W_{k,l}|$$

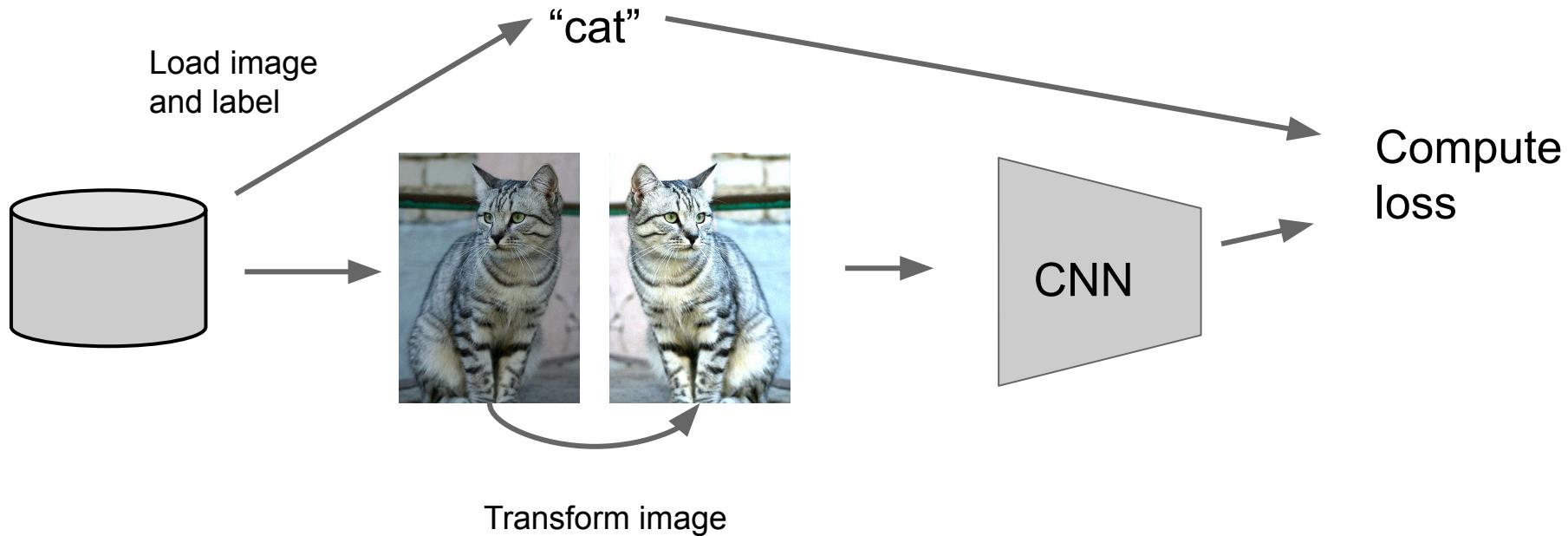
Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

# Regularization: Data Augmentation



# Regularization: Data Augmentation



# Data Augmentation

## Horizontal Flips



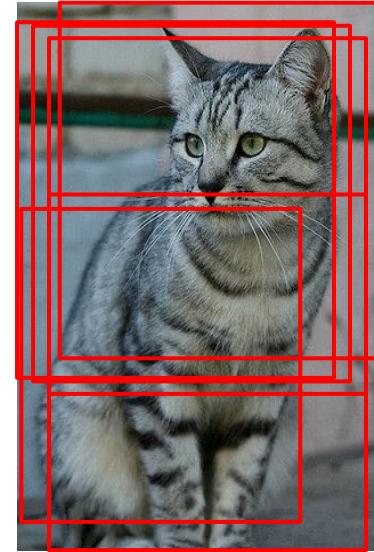
# Data Augmentation

## Random crops and scales

**Training:** sample random crops / scales

ResNet:

1. Pick random  $L$  in range [256, 480]
2. Resize training image, short side =  $L$
3. Sample random  $224 \times 224$  patch



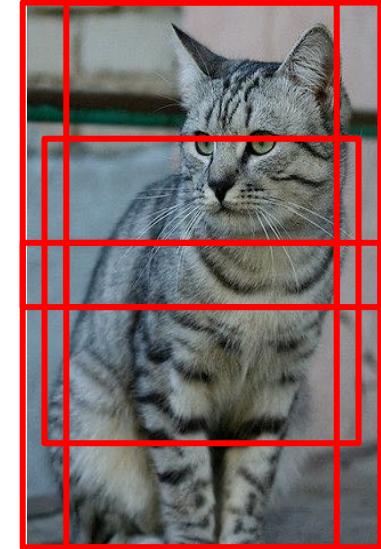
# Data Augmentation

## Random crops and scales

**Training:** sample random crops / scales

ResNet:

1. Pick random  $L$  in range [256, 480]
2. Resize training image, short side =  $L$
3. Sample random  $224 \times 224$  patch



**Testing:** average a fixed set of crops

ResNet:

1. Resize image at 5 scales: {224, 256, 384, 480, 640}
2. For each size, use 10  $224 \times 224$  crops: 4 corners + center, + flips

# Data Augmentation

## Color Jitter

Simple: Randomize  
contrast and brightness



# Data Augmentation

Get creative for your problem!

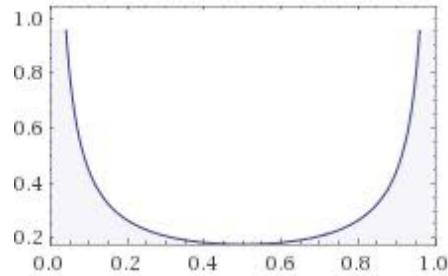
Examples of data augmentations:

- translation
- rotation
- stretching
- shearing,
- lens distortions, ... (go crazy)

# Regularization: Mixup

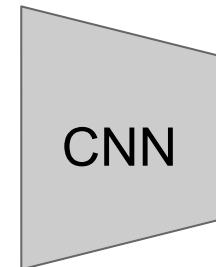
**Training:** Train on random blends of images

**Testing:** Use original images



## Examples:

- Dropout
- Batch Normalization
- Data Augmentation
- DropConnect
- Fractional Max Pooling
- Stochastic Depth
- Cutout / Random Crop
- Mixup



Target label:  
cat: 0.4  
dog: 0.6

Randomly blend the pixels  
of pairs of training images,  
e.g. 40% cat, 60% dog

Zhang et al, "mixup: Beyond Empirical Risk Minimization", ICLR 2018

# Choosing Hyperparameters

(without tons of GPUs)

# Choosing Hyperparameters

## Step 1: Check initial loss

Turn off weight decay, sanity check loss at initialization  
e.g.  $\log(C)$  for softmax with  $C$  classes

# Choosing Hyperparameters

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

Try to train to 100% training accuracy on a small sample of training data (~5-10 minibatches); fiddle with architecture, learning rate, weight initialization

Loss not going down? LR too low, bad initialization

Loss explodes to Inf or NaN? LR too high, bad initialization

# Choosing Hyperparameters

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

Use the architecture from the previous step, use all training data, turn on small weight decay, find a learning rate that makes the loss drop significantly within ~100 iterations

Good learning rates to try: 1e-1, 1e-2, 1e-3, 1e-4

# Choosing Hyperparameters

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

Choose a few values of learning rate and weight decay around what worked from Step 3, train a few models for ~1-5 epochs.

Good weight decay to try: 1e-4, 1e-5, 0

# Choosing Hyperparameters

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

**Step 5:** Refine grid, train longer

Pick best models from Step 4, train them for longer (~10-20 epochs) without learning rate decay

# Choosing Hyperparameters

**Step 1:** Check initial loss

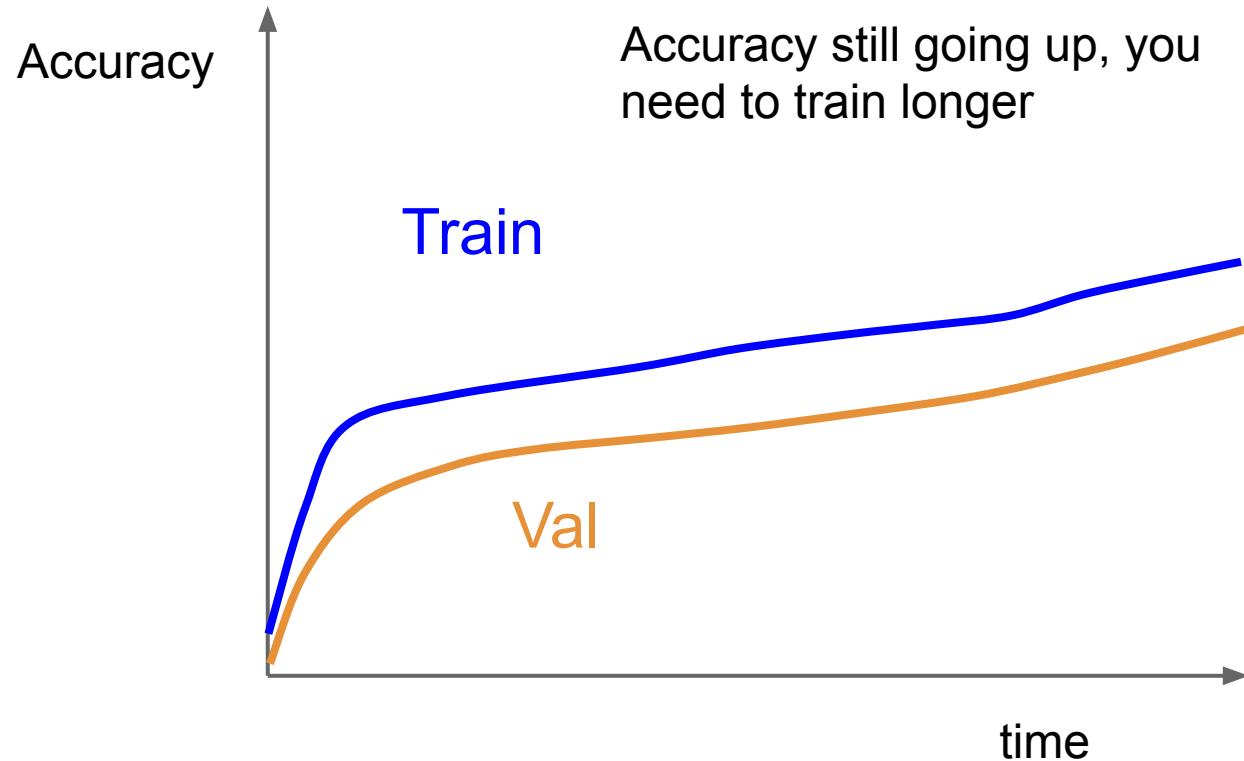
**Step 2:** Overfit a small sample

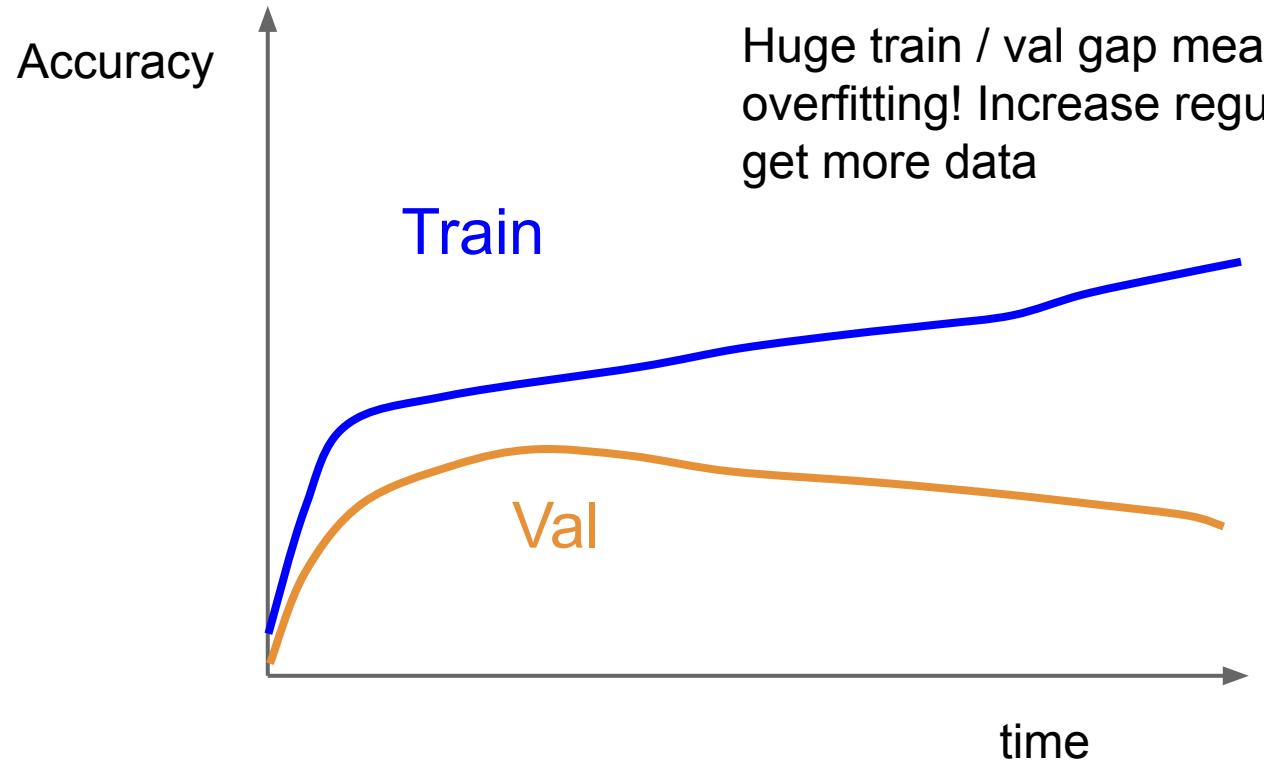
**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

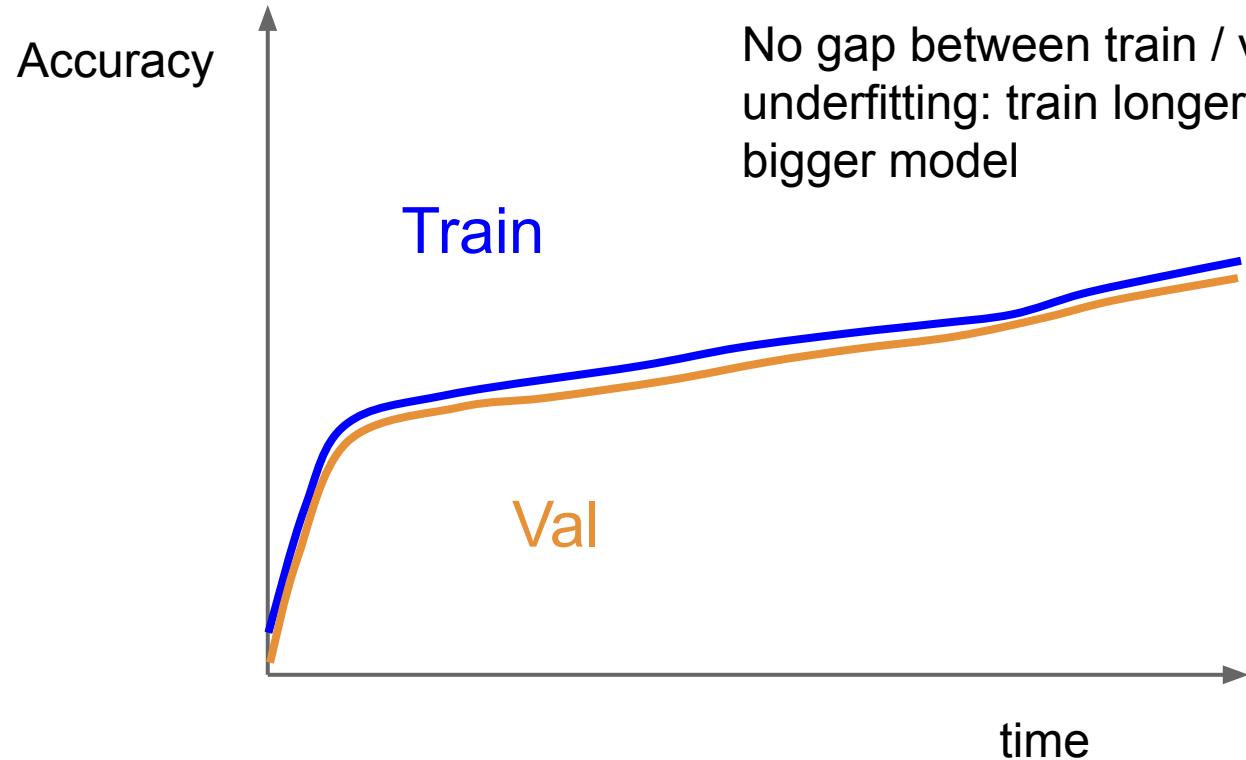
**Step 5:** Refine grid, train longer

**Step 6:** Look at loss and accuracy curves



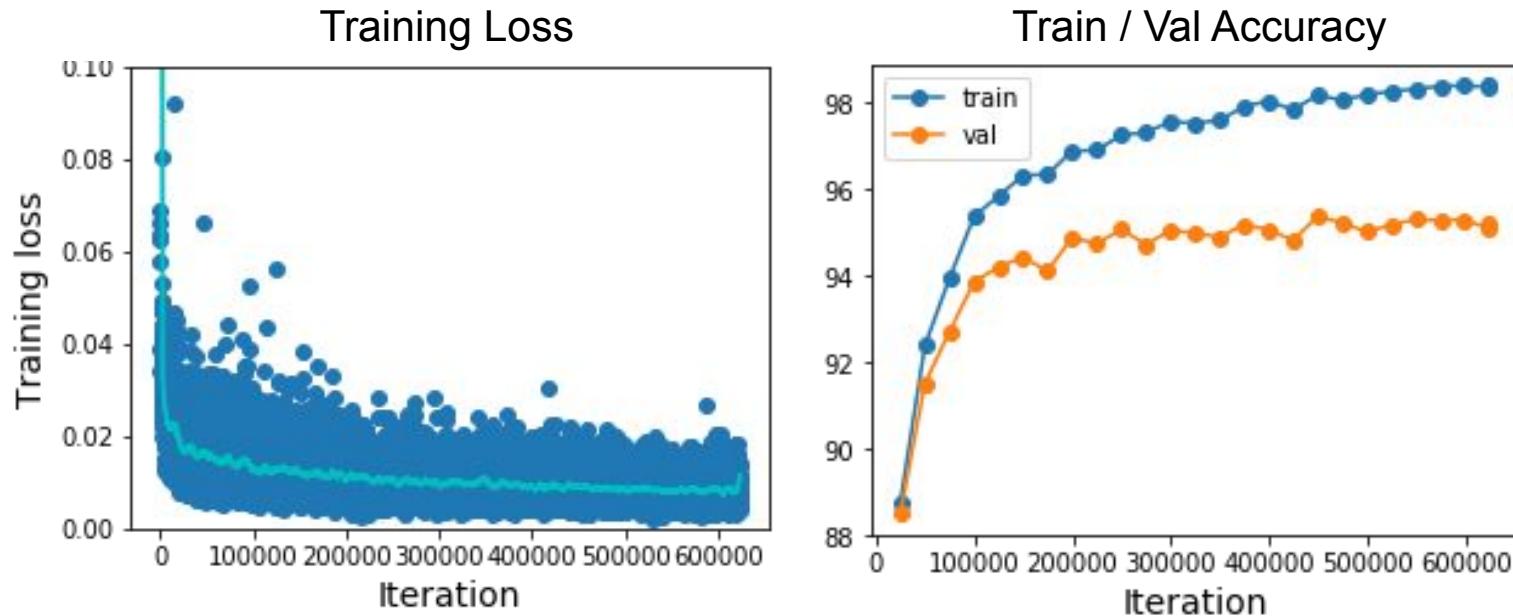


Huge train / val gap means overfitting! Increase regularization, get more data



No gap between train / val means underfitting: train longer, use a bigger model

# Look at learning curves!

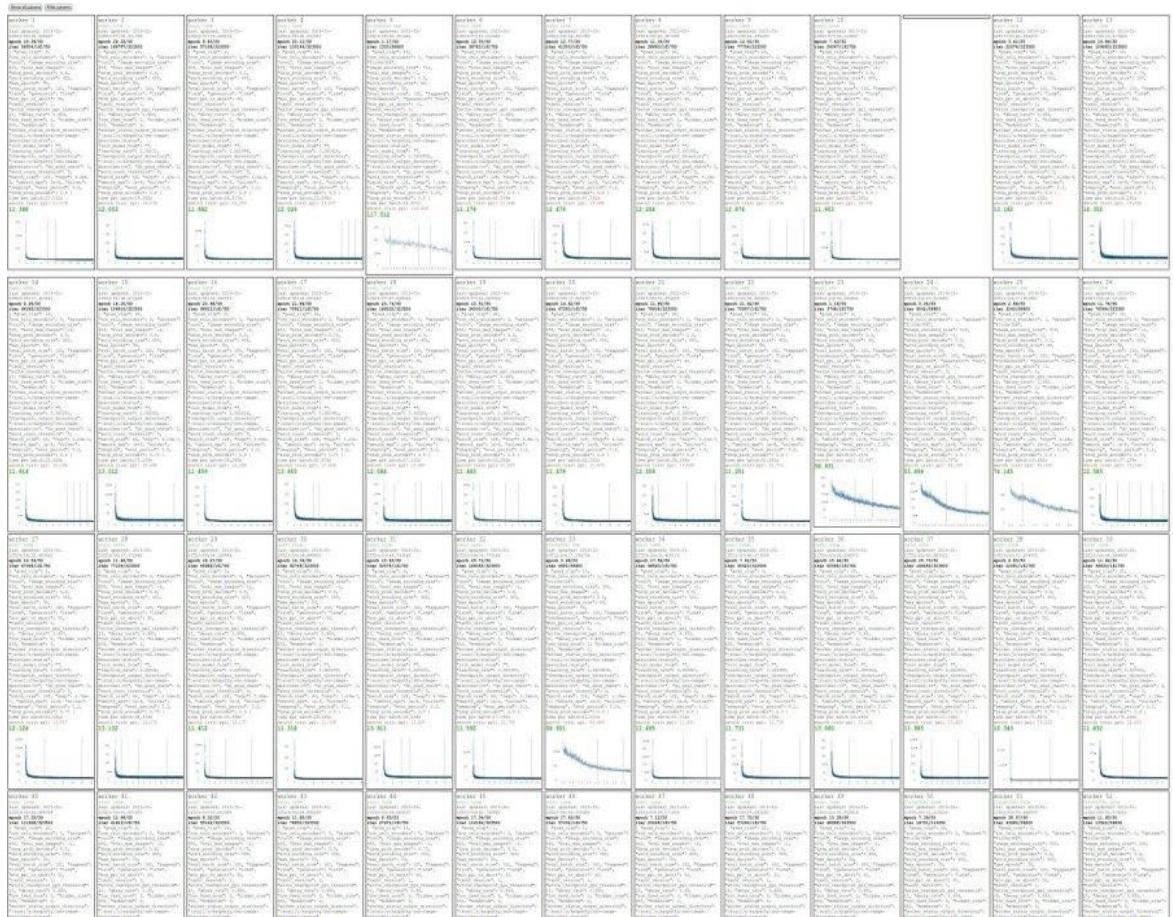


Losses may be noisy, use a scatter plot and also plot moving average to see trends better

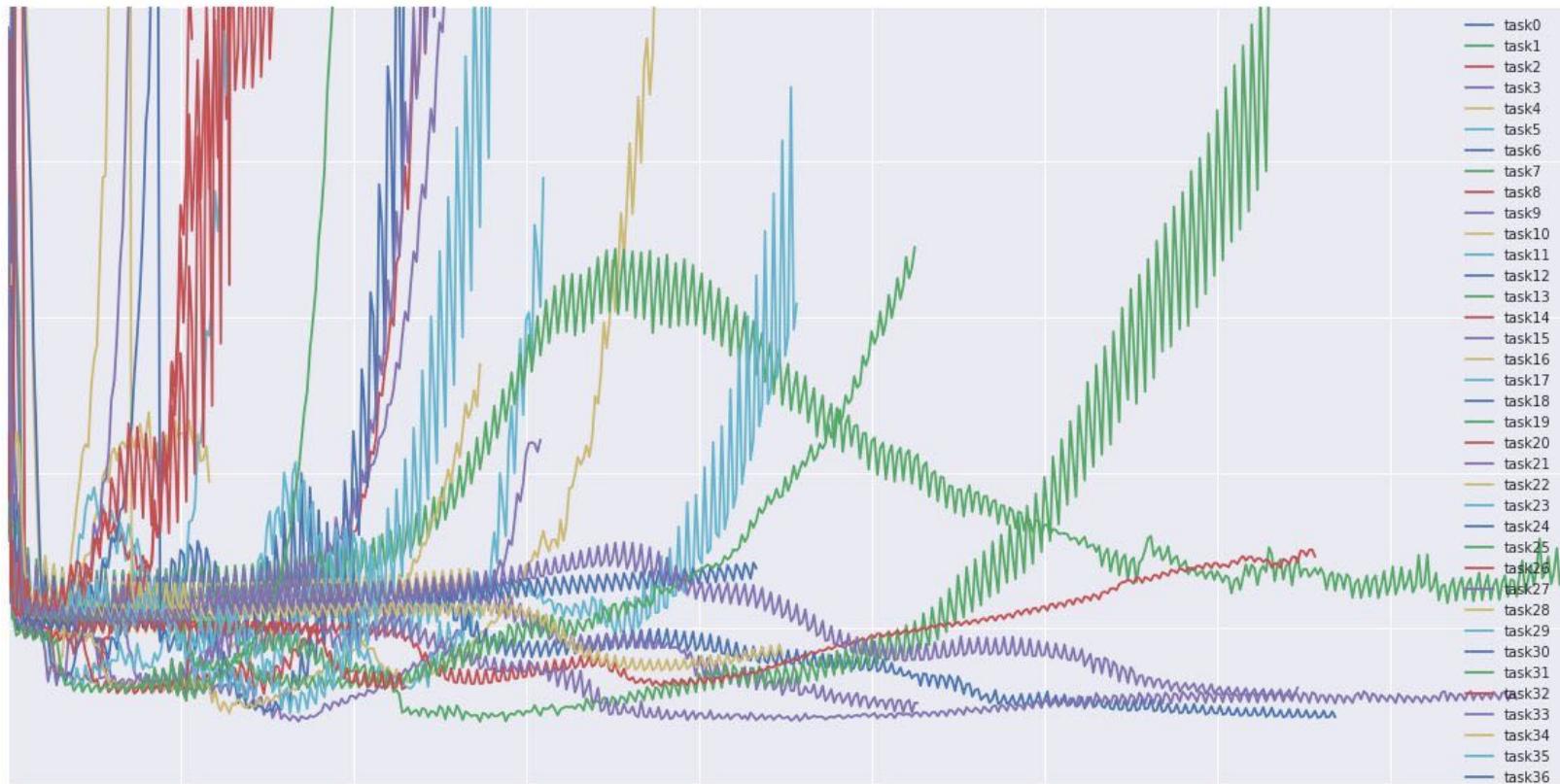
# Cross-validation

We develop "command centers" to visualize all our models training with different hyperparameters

check out [weights and biases](#)



You can plot all your loss curves for different hyperparameters on a single plot



# Choosing Hyperparameters

**Step 1:** Check initial loss

**Step 2:** Overfit a small sample

**Step 3:** Find LR that makes loss go down

**Step 4:** Coarse grid, train for ~1-5 epochs

**Step 5:** Refine grid, train longer

**Step 6:** Look at loss and accuracy curves

**Step 7:** GOTO step 5

# Random Search vs. Grid Search

*Random Search for  
Hyper-Parameter Optimization  
Bergstra and Bengio, 2012*

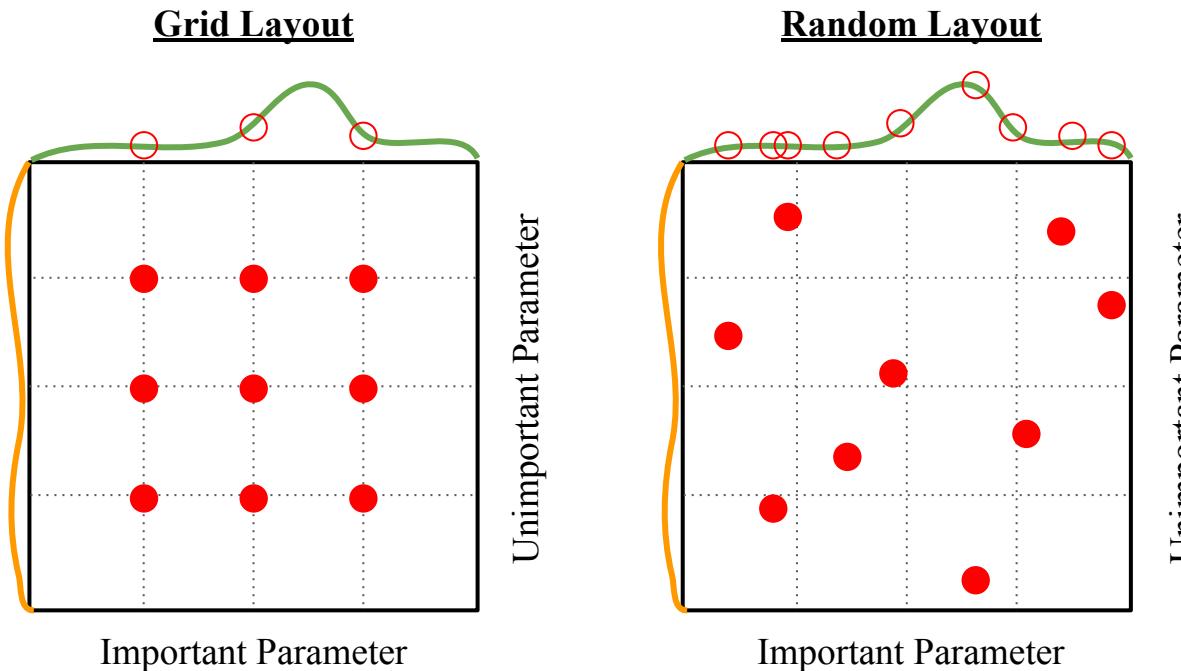
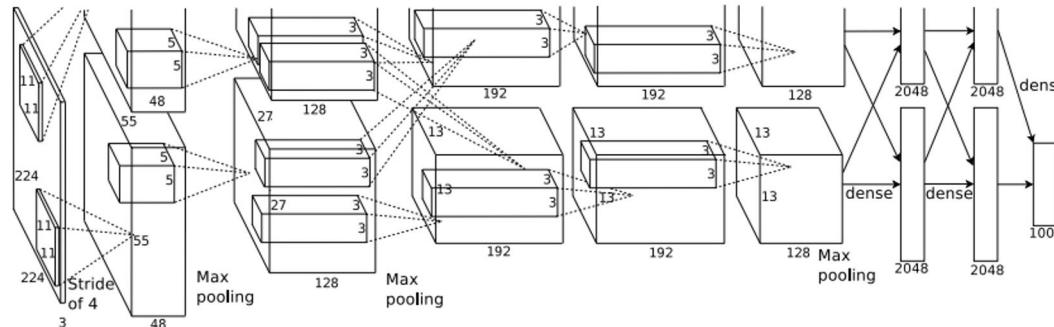


Illustration of Bergstra et al., 2012 by Shayne  
Longpre, copyright CS231n 2017

# Lecture 8: Visualizing and Understanding

# Today: What's going on inside ConvNets?

This image is CC0 public domain



Input Image:  
3 x 224 x 224

What are the intermediate features looking for?

Class Scores:  
1000 numbers

Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.  
Figure reproduced with permission.

## **Visualizing what models have learned:**

- Visualizing filters
- Visualizing final layer features
- Visualizing activations

## **Understanding input pixels**

- Identifying important pixels
- Saliency via backprop
- Guided backprop to generate images
- Gradient ascent to visualize features

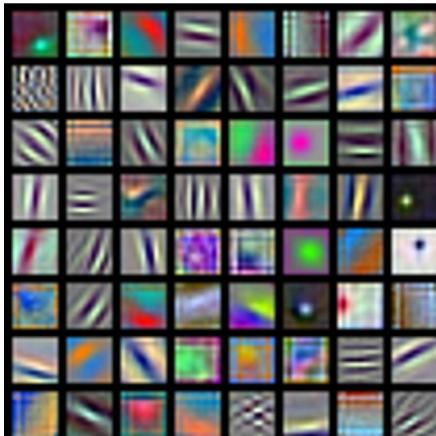
## **Adversarial perturbations**

## **Style transfer**

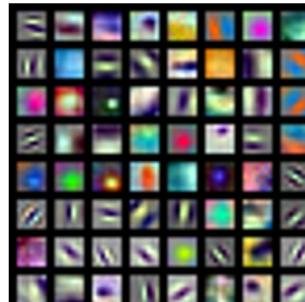
- Features inversion
- Deep dream
- Texture synthesis
- Neural style transfer

## **Today's agenda**

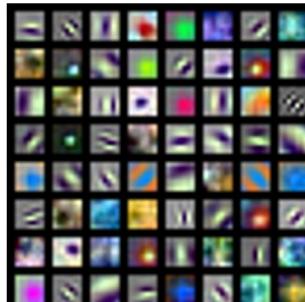
# First Layer: Visualize Filters



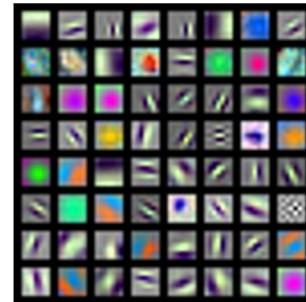
AlexNet:  
 $64 \times 3 \times 11 \times 11$



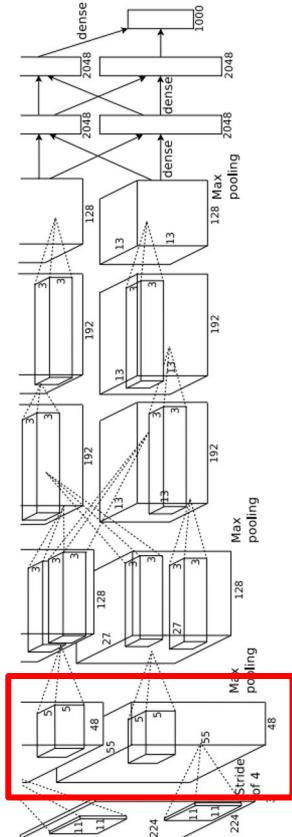
ResNet-18:  
 $64 \times 3 \times 7 \times 7$



ResNet-101:  
 $64 \times 3 \times 7 \times 7$



DenseNet-121:  
 $64 \times 3 \times 7 \times 7$



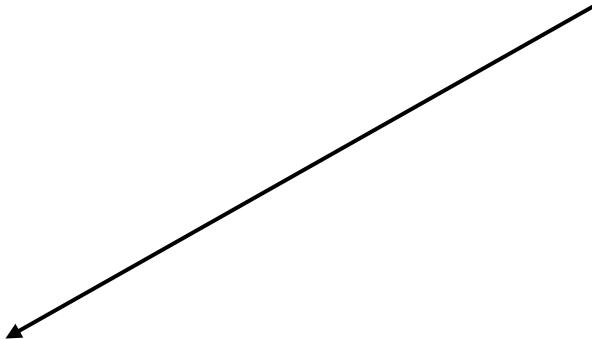
Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014

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

Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

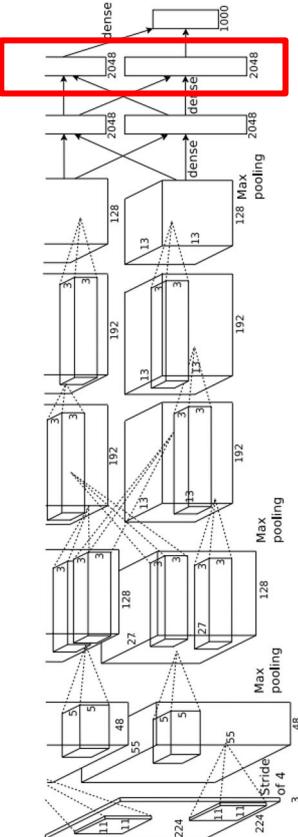
# Last Layer

FC7 layer



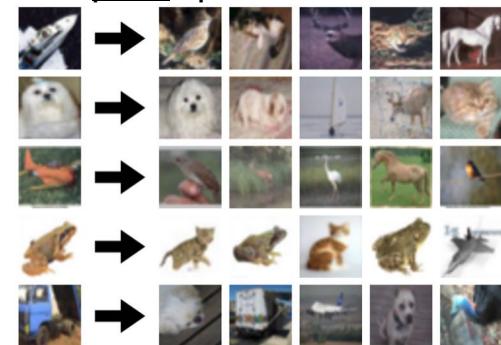
4096-dimensional feature vector for an image  
(layer immediately before the classifier)

Run the network on many images, collect the  
feature vectors



# Last Layer: Nearest Neighbors

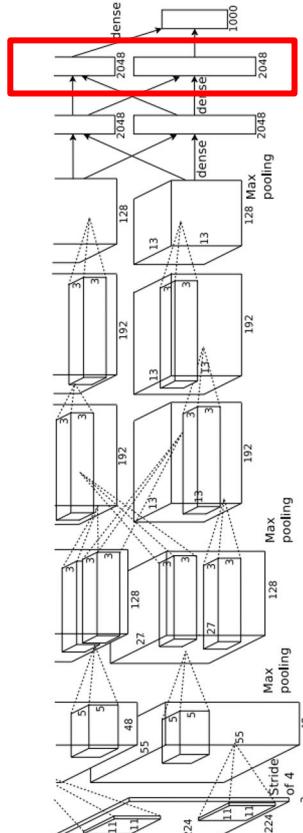
**Recall:** Nearest neighbors  
in pixel space



Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.  
Figures reproduced with permission.

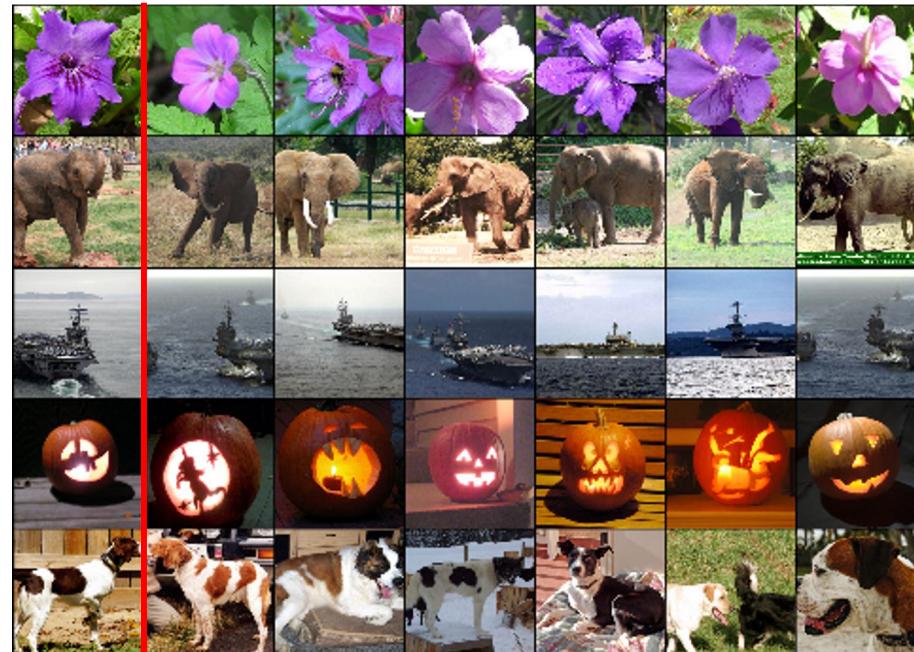
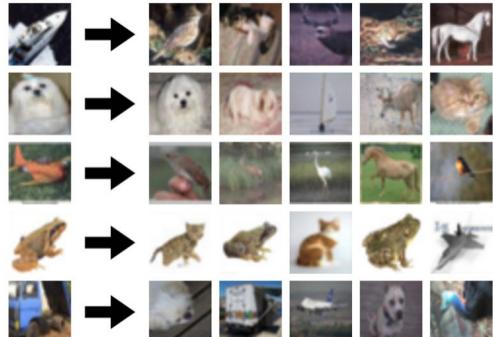
# Last Layer: Nearest Neighbors

4096-dim vector



Test image L2 Nearest neighbors in feature space

Recall: Nearest neighbors in pixel space



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

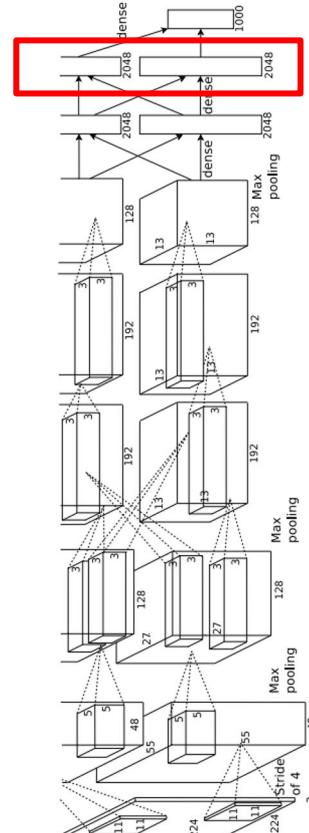
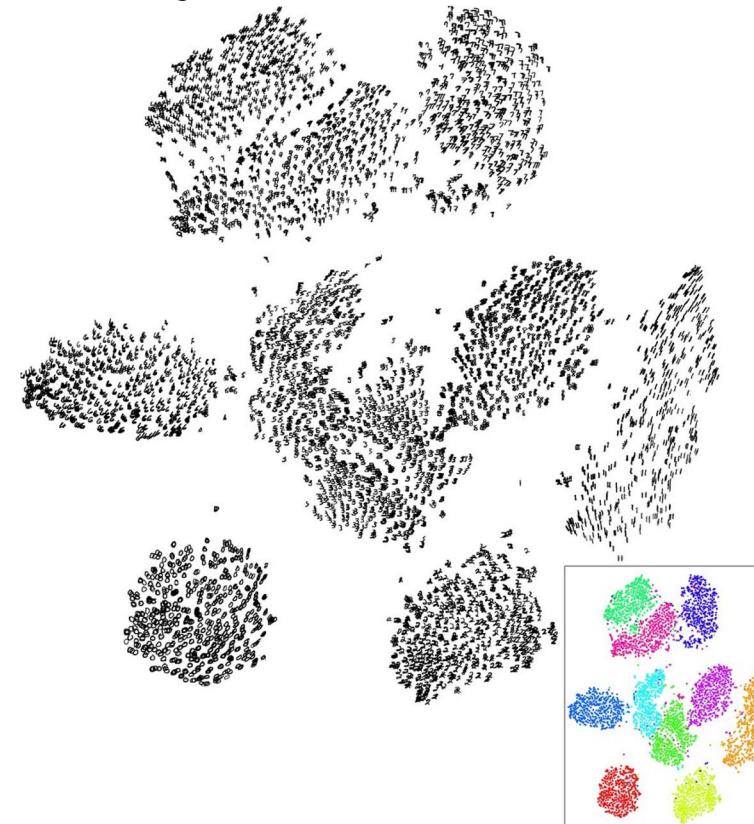
Figures reproduced with permission.

# Last Layer: Dimensionality Reduction

Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

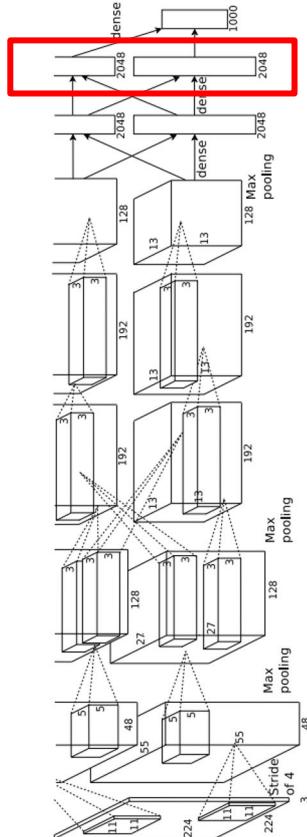
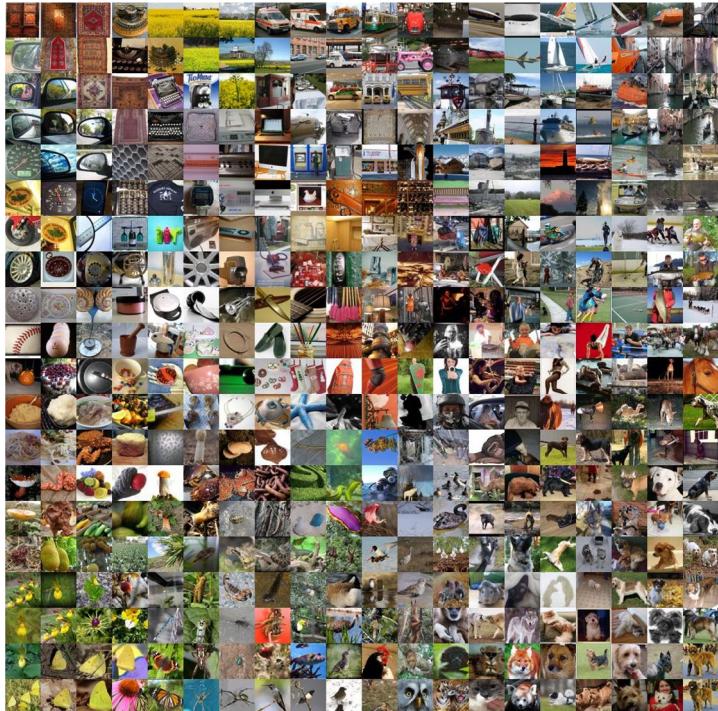
Simple algorithm: Principal Component Analysis (PCA)

More complex: t-SNE



Van der Maaten and Hinton, “Visualizing Data using t-SNE”, JMLR 2008  
Figure copyright Laurens van der Maaten and Geoff Hinton, 2008. Reproduced with permission.

# Last Layer: Dimensionality Reduction

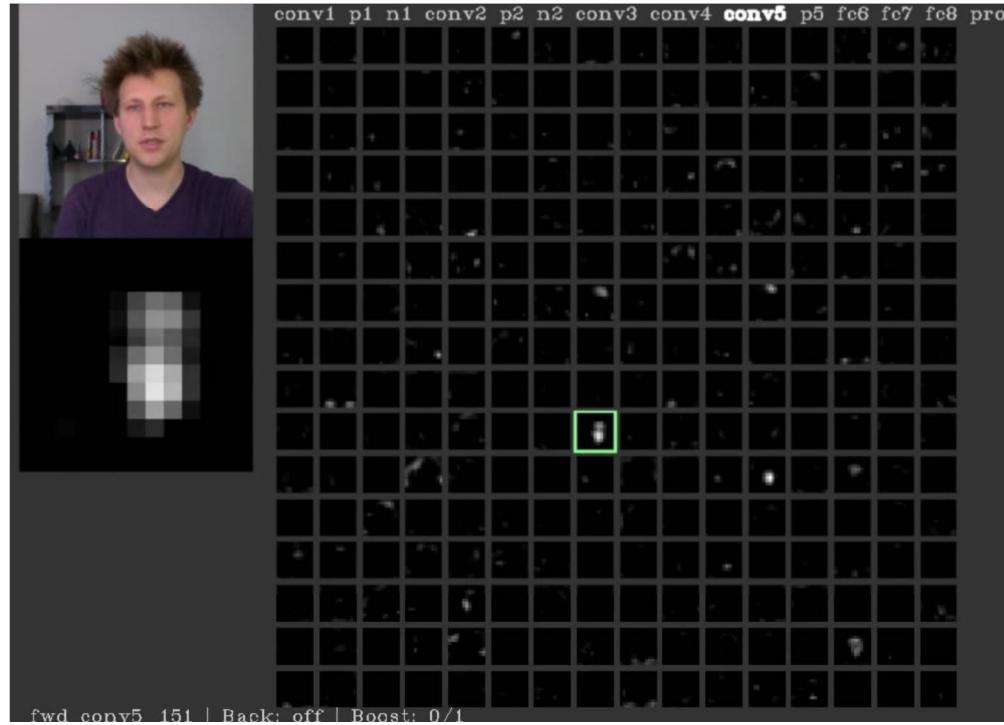


Van der Maaten and Hinton, "Visualizing Data using t-SNE", JMLR 2008  
Krizhevsky et al, "ImageNet Classification with Deep Convolutional Neural Networks", NIPS 2012.  
Figure reproduced with permission.

See high-resolution versions at  
<http://cs.stanford.edu/people/karpathy/cnnembed/>

# Visualizing Activations

conv5 feature map  
is  $128 \times 13 \times 13$ ;  
visualize as 128  
 $13 \times 13$  grayscale  
images



Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.  
Figure copyright Jason Yosinski, 2014. Reproduced with permission.

## Visualizing what models have learned:

- Visualizing filters
- Visualizing final layer features
- Visualizing activations

## Understanding input pixels

- Identifying important pixels
- Saliency via backprop
- Guided backprop to generate images
- Gradient ascent to visualize features

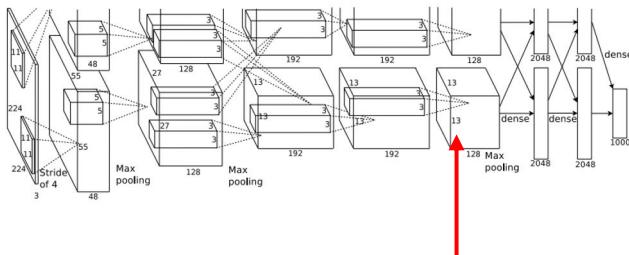
## Adversarial perturbations

## Style transfer

- Deep dream
- Features inversion
- Texture synthesis
- Neural style transfer

# Today's agenda

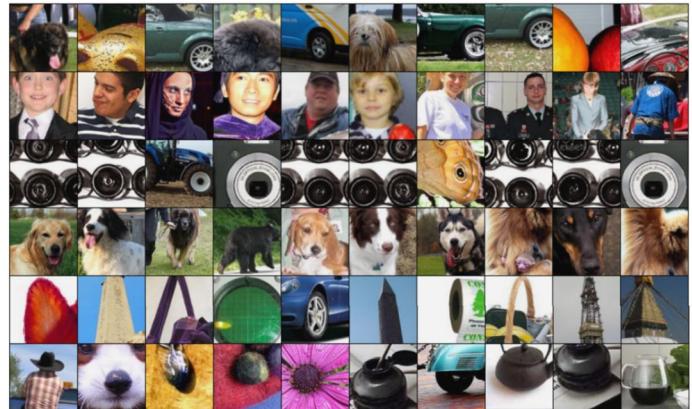
# Maximally Activating Patches



Pick a layer and a channel; e.g. conv5 is 128 x 13 x 13, pick channel 17/128

Run many images through the network, record values of chosen channel

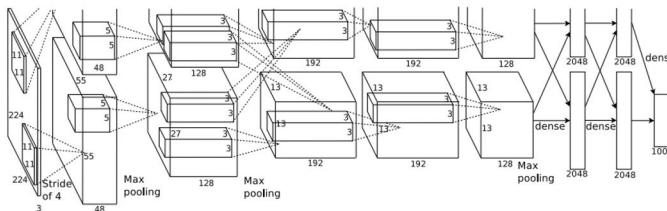
Visualize image patches that correspond to maximal activations



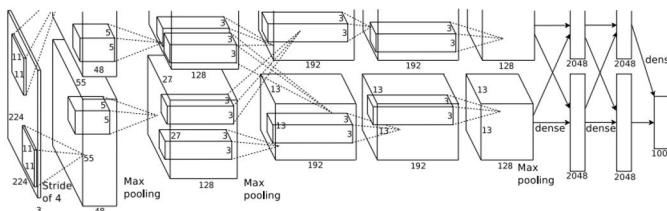
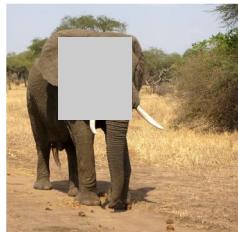
Springenberg et al., "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015  
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015;  
reproduced with permission.

# Which pixels matter: Saliency via Occlusion

Mask part of the image before feeding to CNN,  
check how much predicted probabilities change



$P(\text{elephant}) = 0.95$



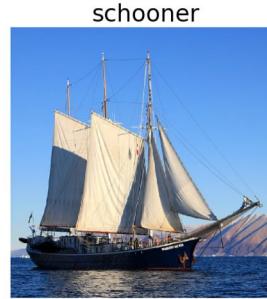
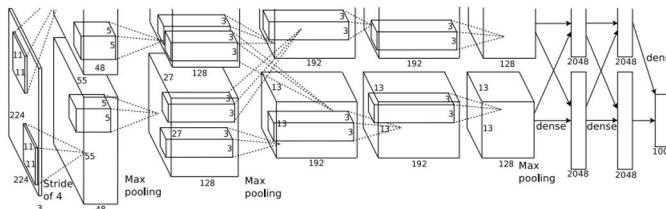
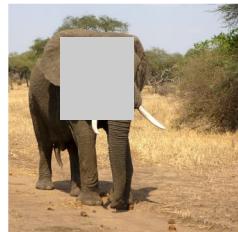
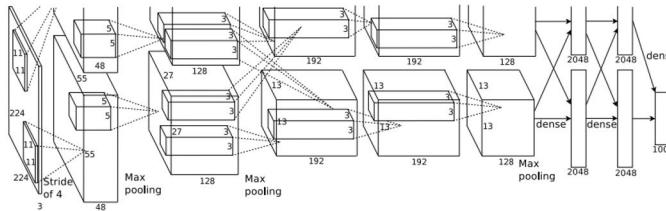
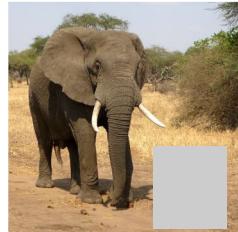
$P(\text{elephant}) = 0.75$

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

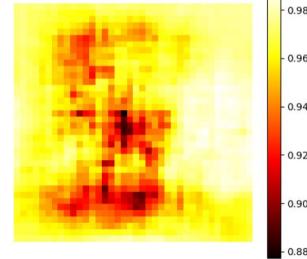
Boat image is CC0 public domain  
Elephant image is CC0 public domain  
Go-Karts image is CC0 public domain

# Which pixels matter: Saliency via Occlusion

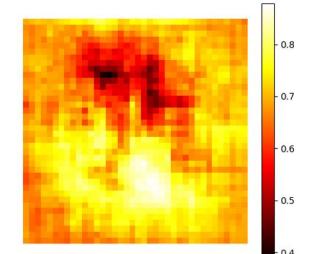
Mask part of the image before feeding to CNN,  
check how much predicted probabilities change



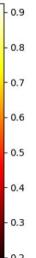
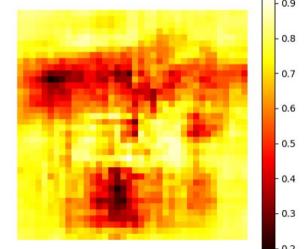
schooner



African elephant, *Loxodonta africana*



go-kart

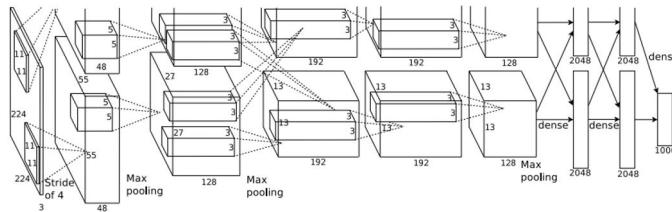


Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

Boat image is CC0 public domain  
Elephant image is CC0 public domain  
Go-Karts image is CC0 public domain

# Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities

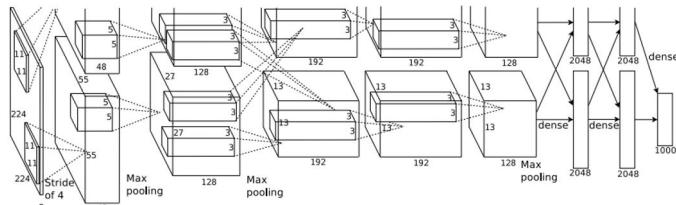


Dog

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.  
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

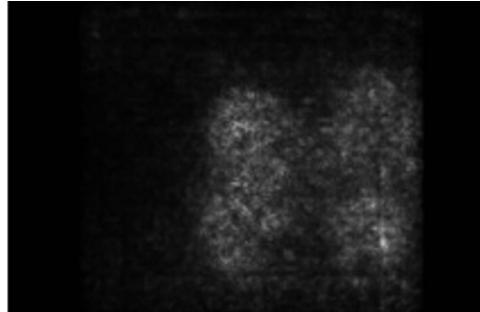
# Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



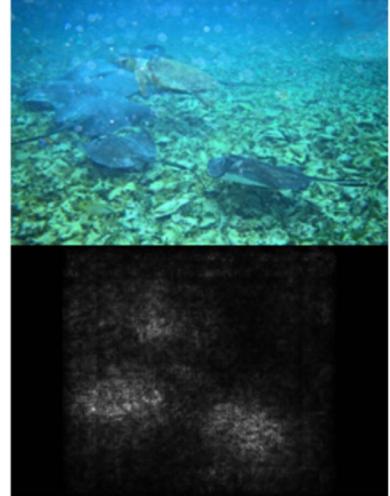
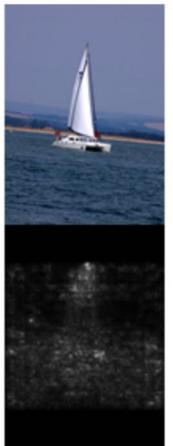
Dog

Compute gradient of **(unnormalized) class score** with respect to image pixels, take absolute value and max over RGB channels



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.  
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Saliency Maps

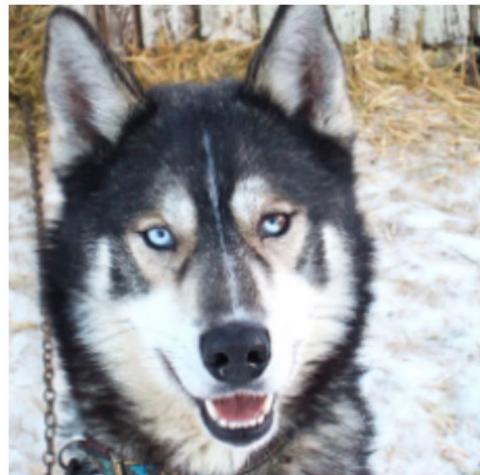


Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.  
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

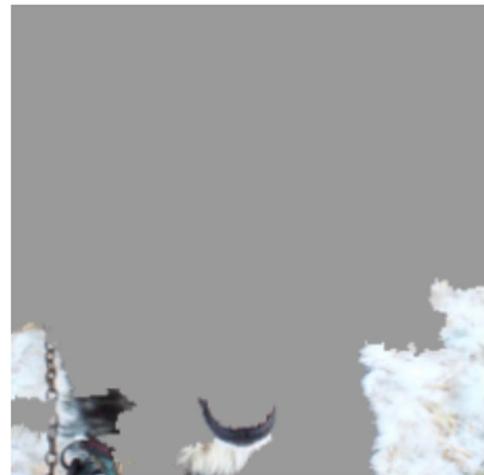
# Saliency maps: Uncovers biases

Such methods also find  
biases

wolf vs dog classifier looks  
is actually a snow vs no-  
snow classifier



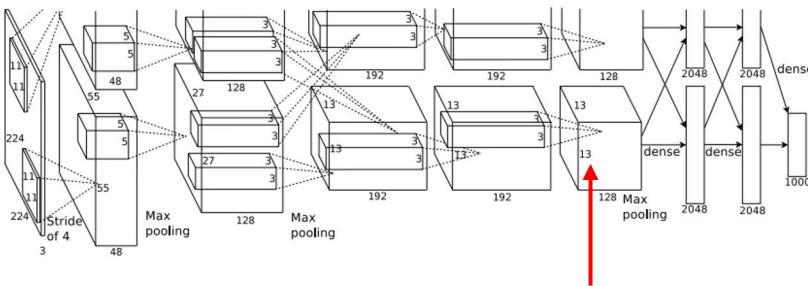
(a) Husky classified as wolf



(b) Explanation

Figures copyright Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin, 2016; reproduced with permission.  
Ribeiro et al, "Why Should I Trust You?" Explaining the Predictions of Any Classifier", ACM KDD 2016

# Intermediate Features via (guided) backprop

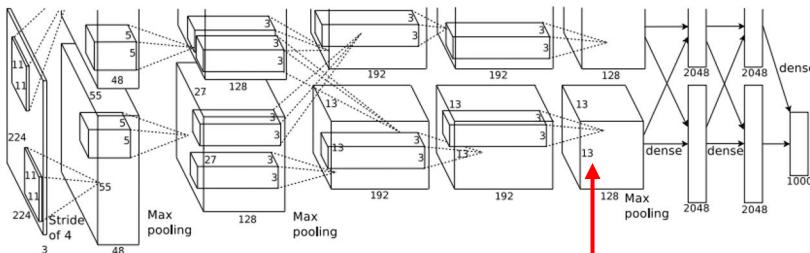


Pick a single intermediate channel, e.g. one value in  $128 \times 13 \times 13$  conv5 feature map

Compute gradient of activation value with respect to image pixels

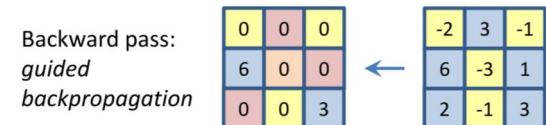
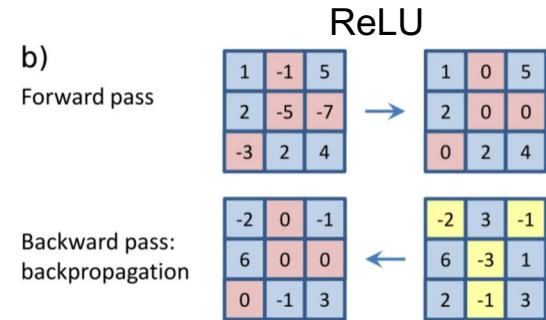
Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014  
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015

# Intermediate Features via (guided) backprop



Pick a single intermediate neuron, e.g. one value in  $128 \times 13 \times 13$  conv5 feature map

Compute gradient of neuron value with respect to image pixels



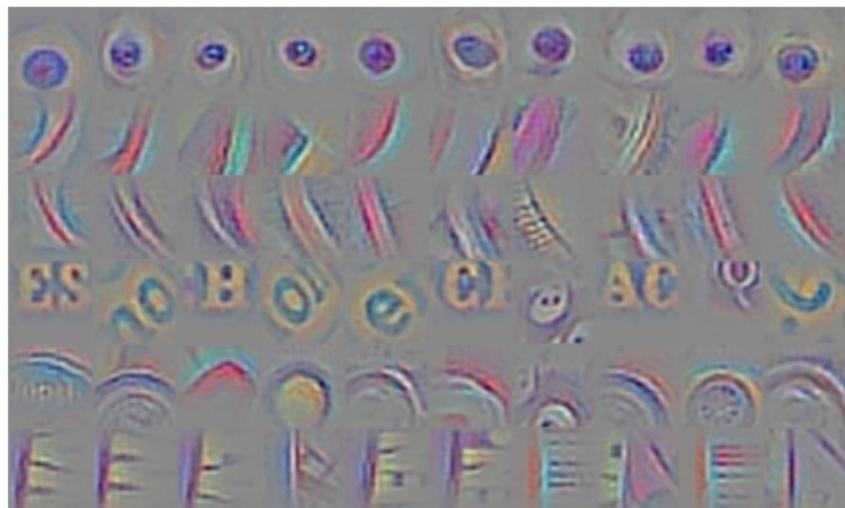
Images come out nicer if you only backprop positive gradients through each ReLU (guided backprop)

Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

# Intermediate features via (guided) backprop



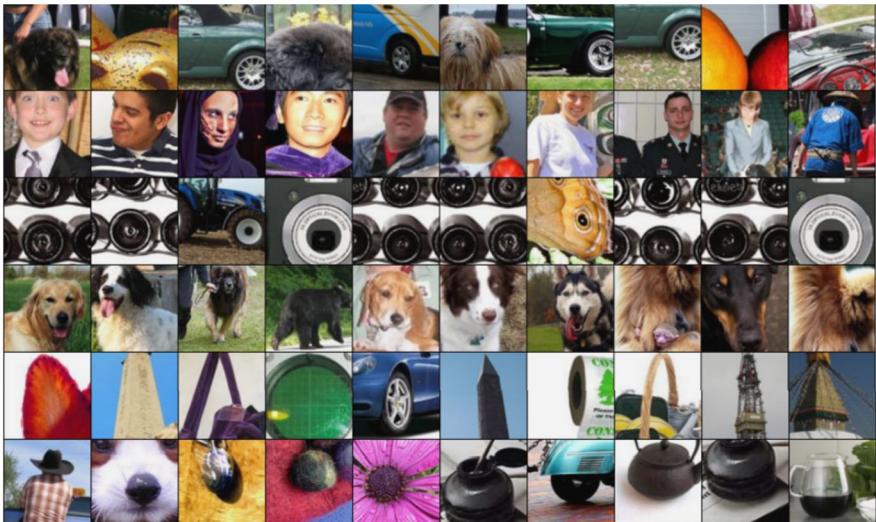
Maximally activating patches  
(Each row is a different neuron)



Guided Backprop

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014  
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015  
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

# Intermediate features via (guided) backprop



Maximally activating patches  
(Each row is a different neuron)



Guided Backprop

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014  
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015  
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

# Visualizing CNN features: Gradient Ascent

**(Guided) backprop:**

Find the part of an image that a neuron responds to

**Gradient ascent:**

Generate a synthetic image that maximally activates a neuron

$$I^* = \arg \max_I f(I) + R(I)$$

Neuron value

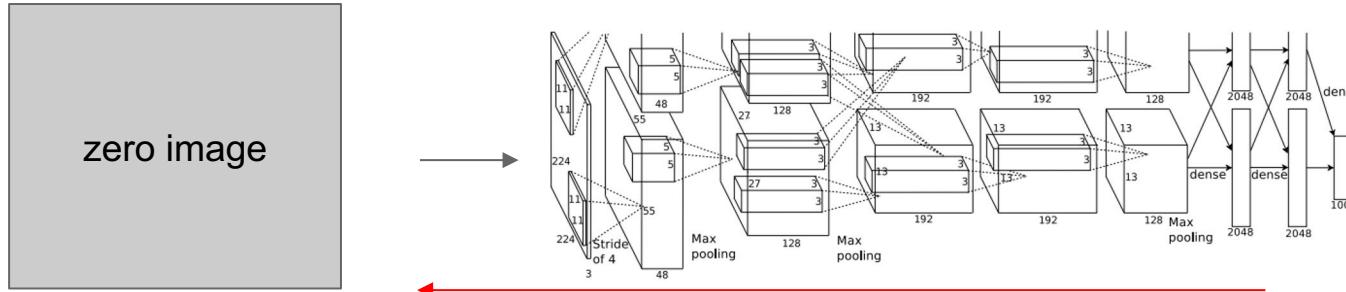
Natural image regularizer

# Visualizing CNN features: Gradient Ascent

1. Initialize image to zeros

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

score for class c (before Softmax)



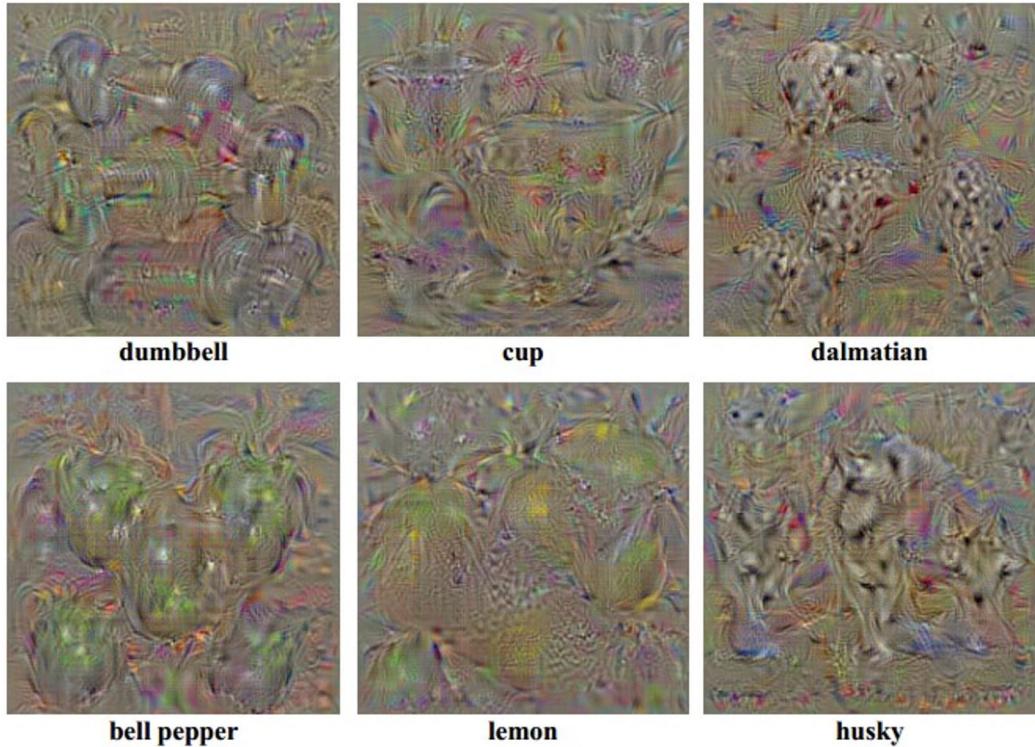
Repeat:

2. Forward image to compute current scores
3. Backprop to get gradient of neuron value with respect to image pixels
4. Make a small update to the image

# Visualizing CNN features: Gradient Ascent

$$\arg \max_I S_c(I) - \boxed{\lambda \|I\|_2^2}$$

Simple regularizer: Penalize L2  
norm of generated image

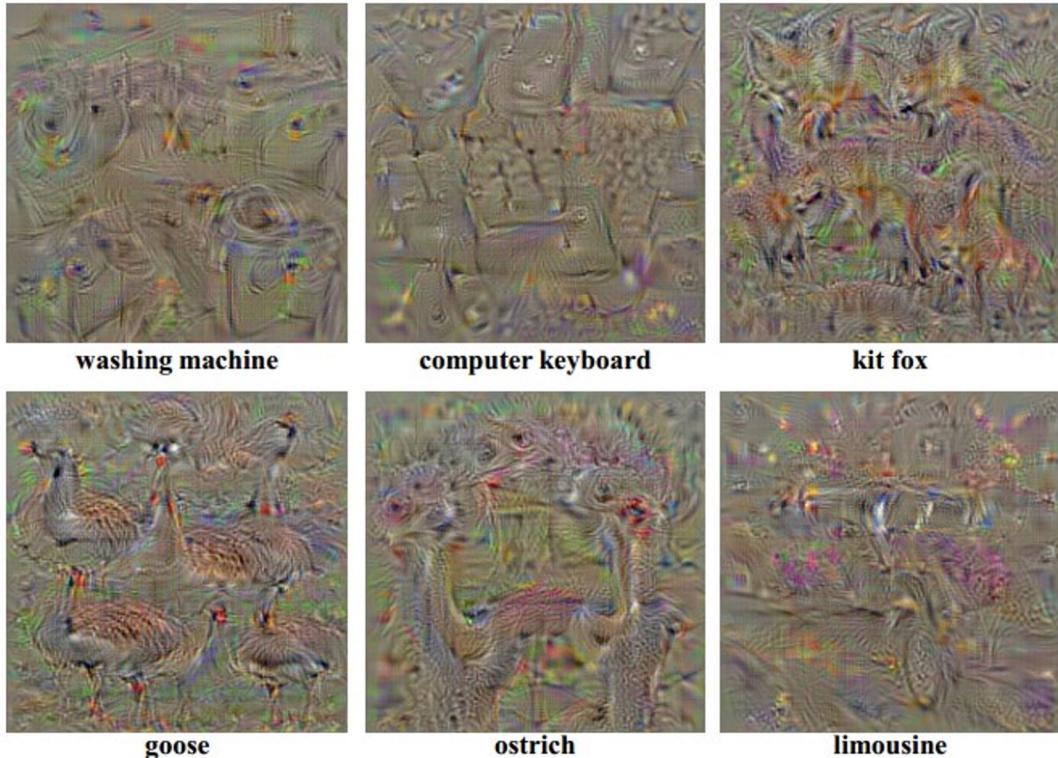


Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.  
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Visualizing CNN features: Gradient Ascent

$$\arg \max_I S_c(I) - \boxed{\lambda \|I\|_2^2}$$

Simple regularizer: Penalize L2  
norm of generated image



Yosinski et al., "Understanding Neural Networks Through Deep Visualization", ICML DL Workshop 2014.  
Figure copyright Jason Yosinski, Jeff Clune, Anh Nguyen, Thomas Fuchs, and Hod Lipson, 2014.  
Reproduced with permission.

# Visualizing CNN features: Gradient Ascent

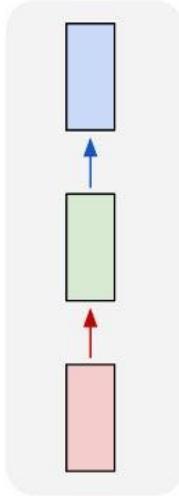


Nguyen et al., "Multifaceted Feature Visualization: Uncovering the Different Types of Features Learned By Each Neuron in Deep Neural Networks", ICML Visualization for Deep Learning Workshop 2016.  
Figures copyright Anh Nguyen, Jason Yosinski, and Jeff Clune, 2016; reproduced with permission.

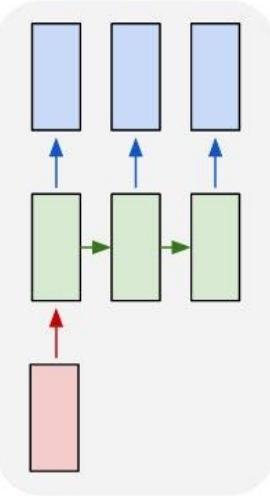
# Lecture 11: Attention and Transformers

# Last Time: Recurrent Neural Networks

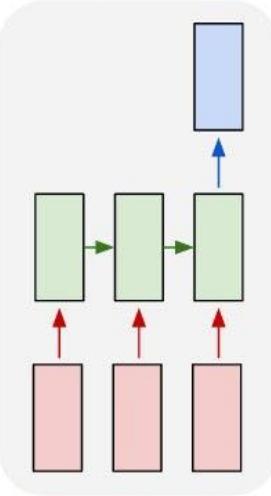
one to one



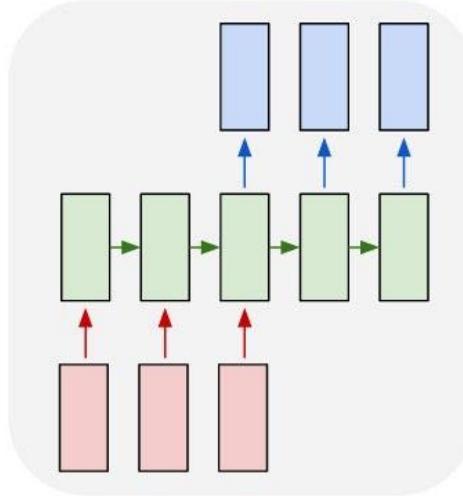
one to many



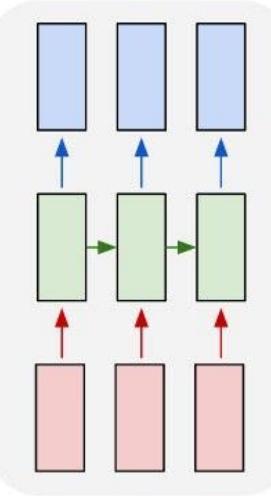
many to one



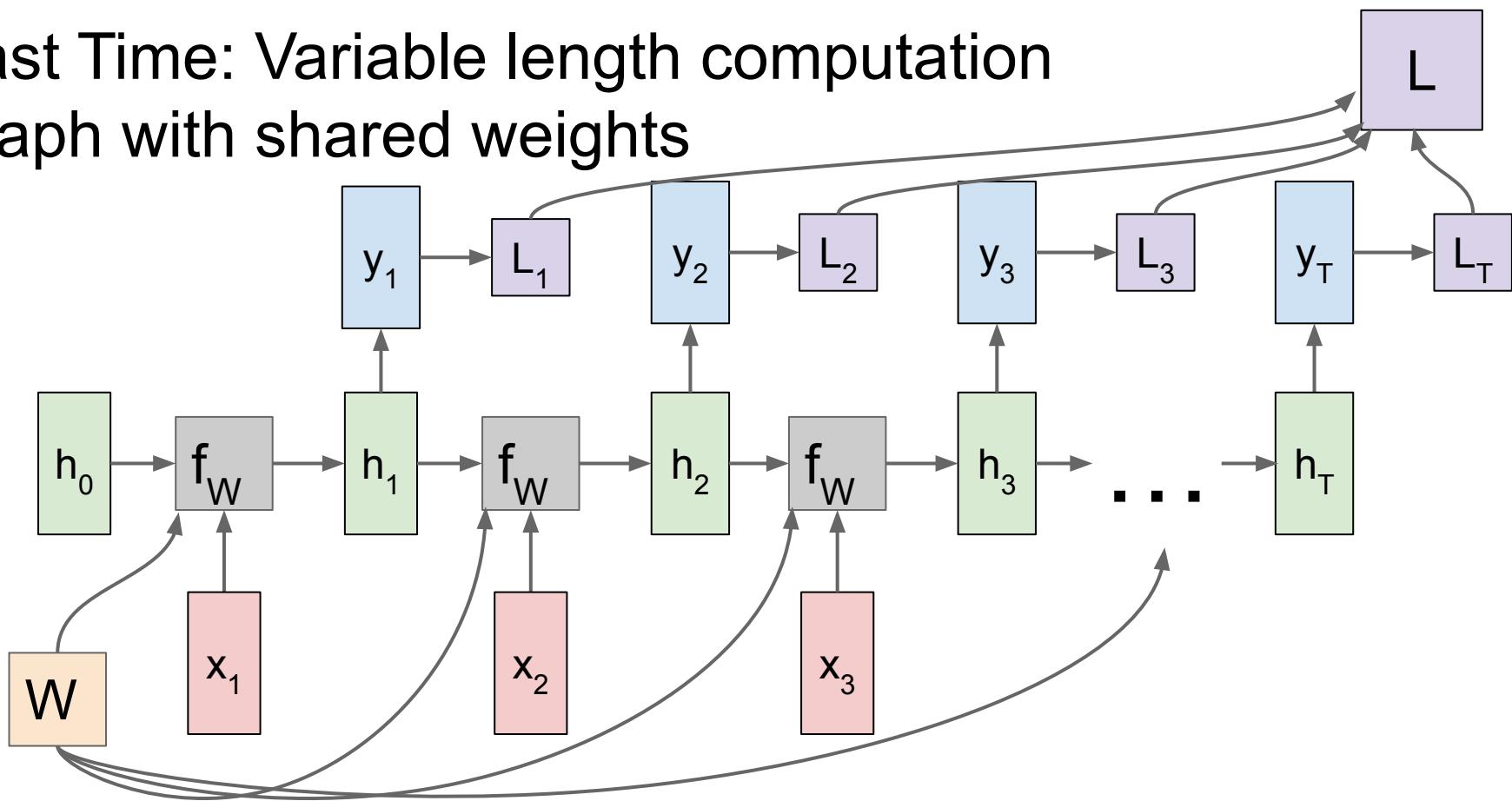
many to many



many to many



# Last Time: Variable length computation graph with shared weights

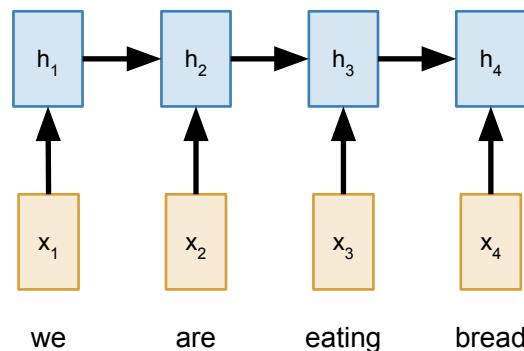


# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

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



Sutskever et al, "Sequence to sequence learning with neural networks", NeurIPS 2014

# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

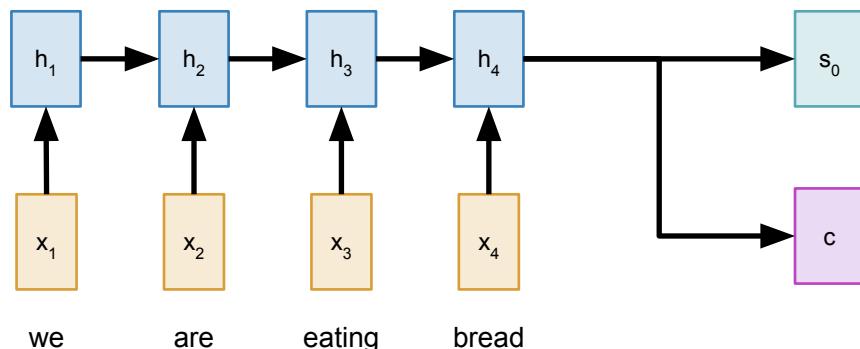
**Output:** Sequence  $y_1, \dots, y_T$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



Sutskever et al., "Sequence to sequence learning with neural networks", NeurIPS 2014

# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

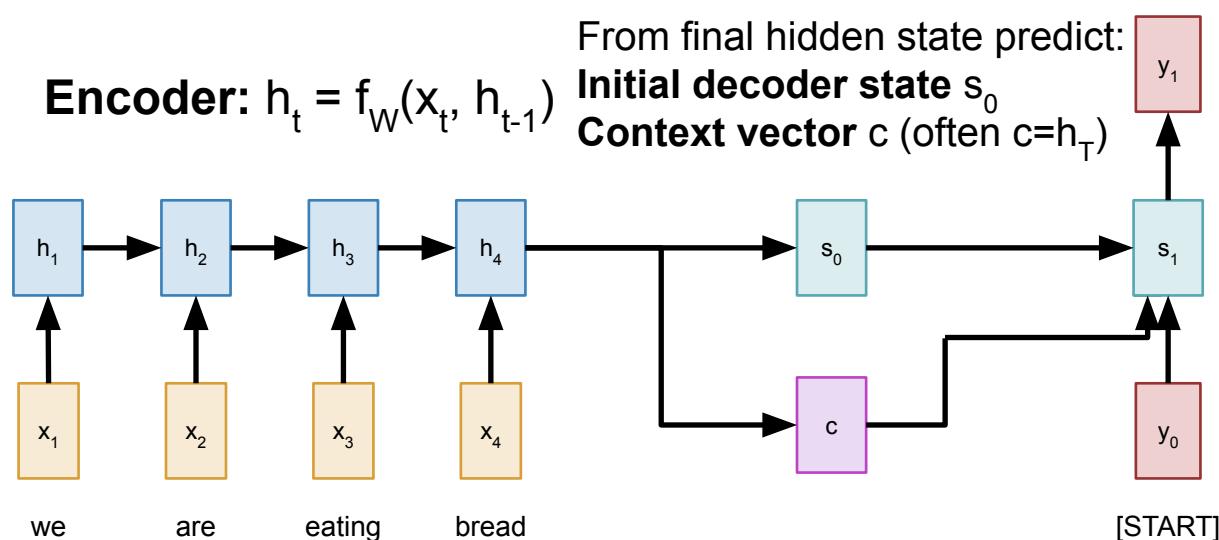
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

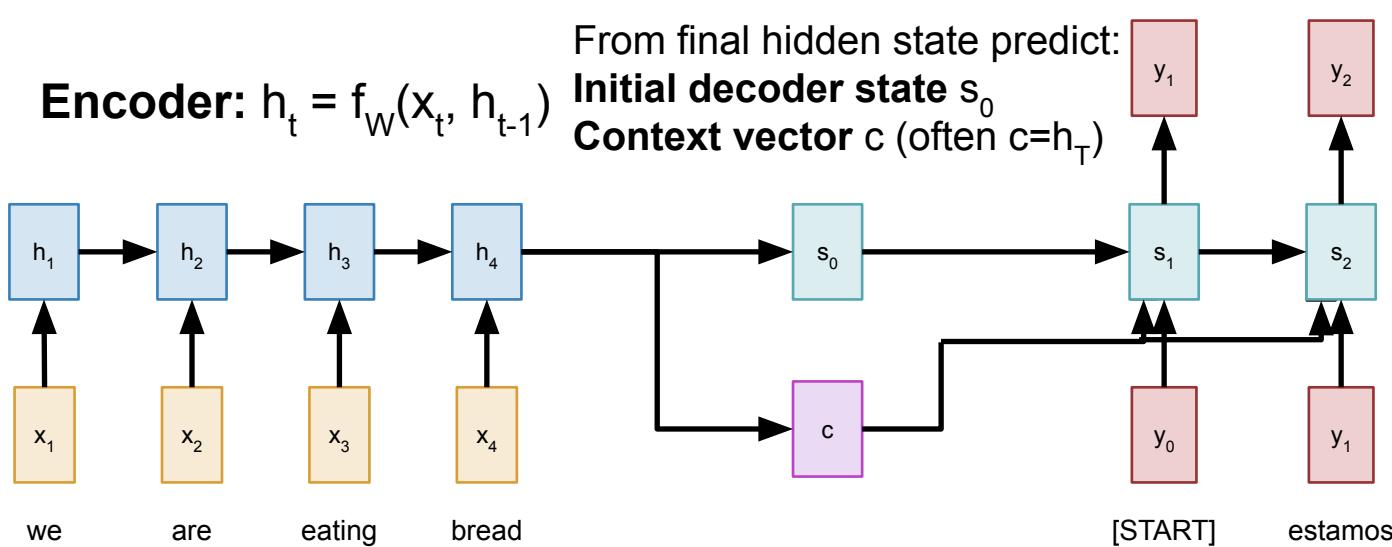
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

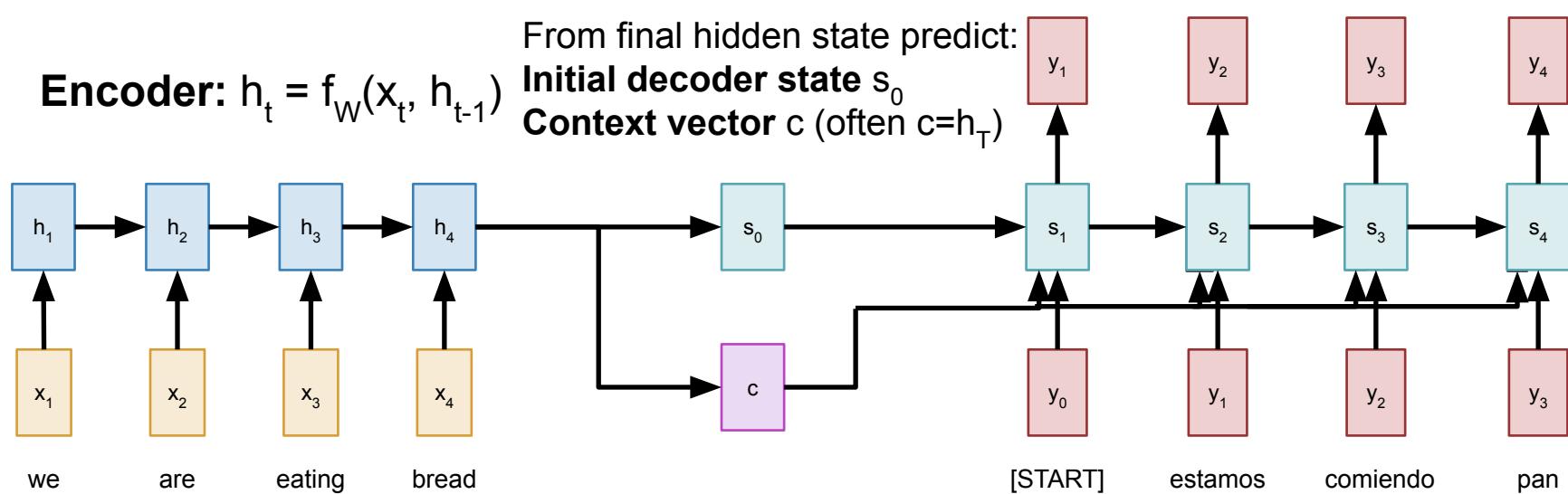
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



Sutskever et al., "Sequence to sequence learning with neural networks", NeurIPS 2014

# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

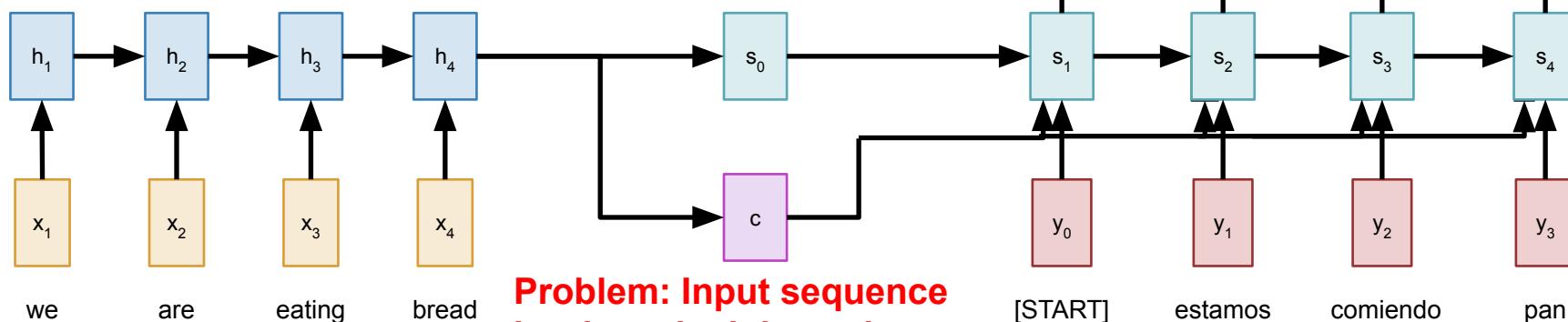
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



**Problem: Input sequence bottlenecked through fixed-sized vector. What if  $T=1000$ ?**

Sutskever et al, "Sequence to sequence learning with neural networks", NIPS 2014

# Sequence to Sequence with RNNs

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_T$

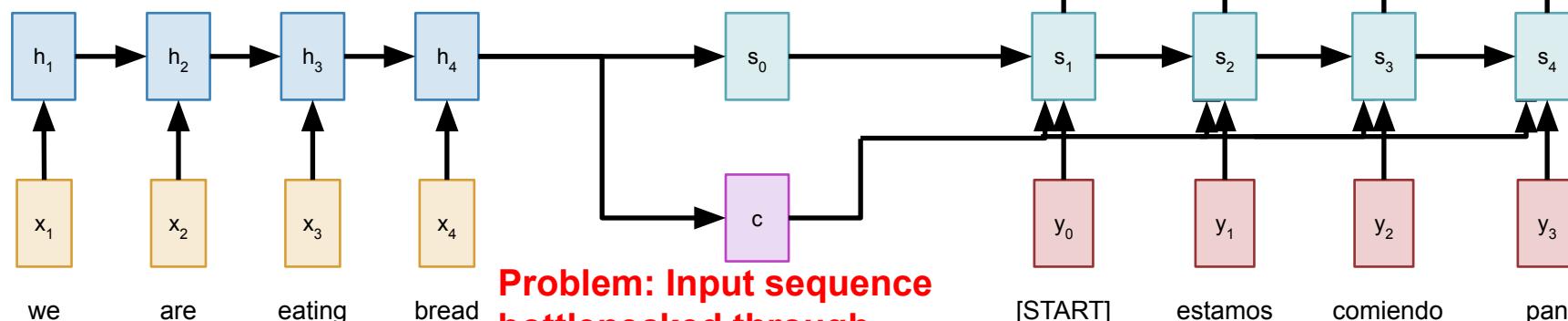
**Decoder:**  $s_t = g_U(y_{t-1}, s_{t-1}, c)$

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

From final hidden state predict:

**Initial decoder state**  $s_0$

**Context vector**  $c$  (often  $c=h_T$ )



**Problem:** Input sequence  
bottlenecked through  
fixed-sized vector. What if  
 $T=1000$ ?

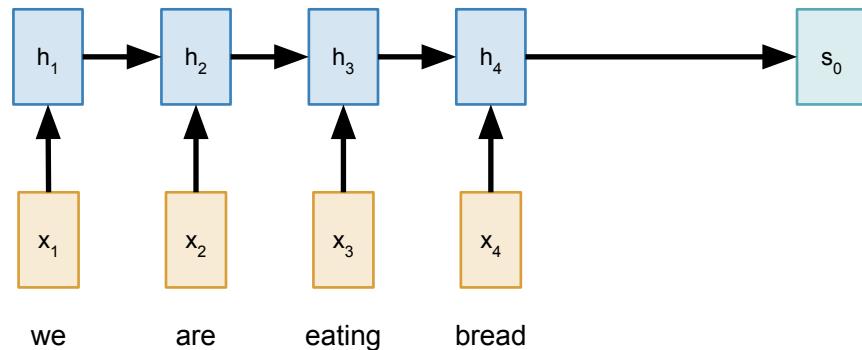
**Idea:** use new context vector  
at each step of decoder!

# Sequence to Sequence with RNNs and Attention

**Input:** Sequence  $x_1, \dots, x_T$

**Output:** Sequence  $y_1, \dots, y_{T'}$

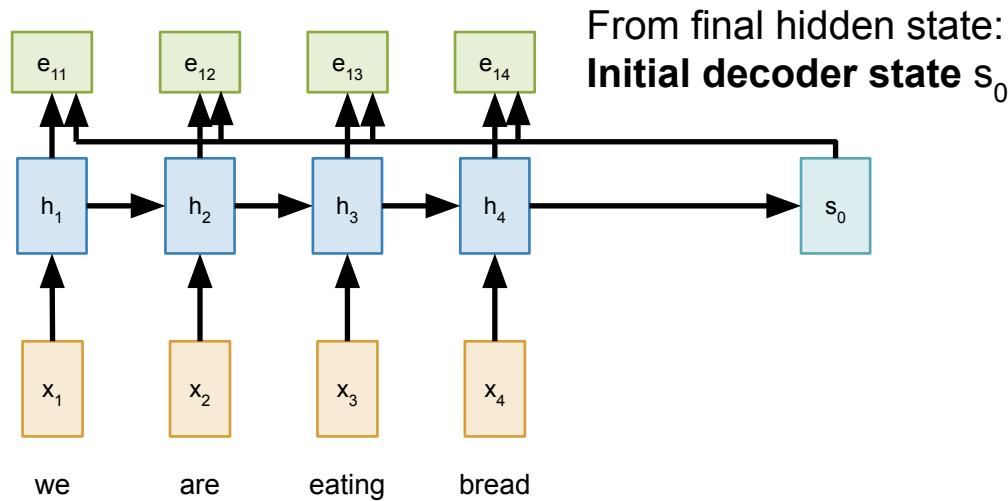
**Encoder:**  $h_t = f_W(x_t, h_{t-1})$  From final hidden state:  
**Initial decoder state**  $s_0$



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

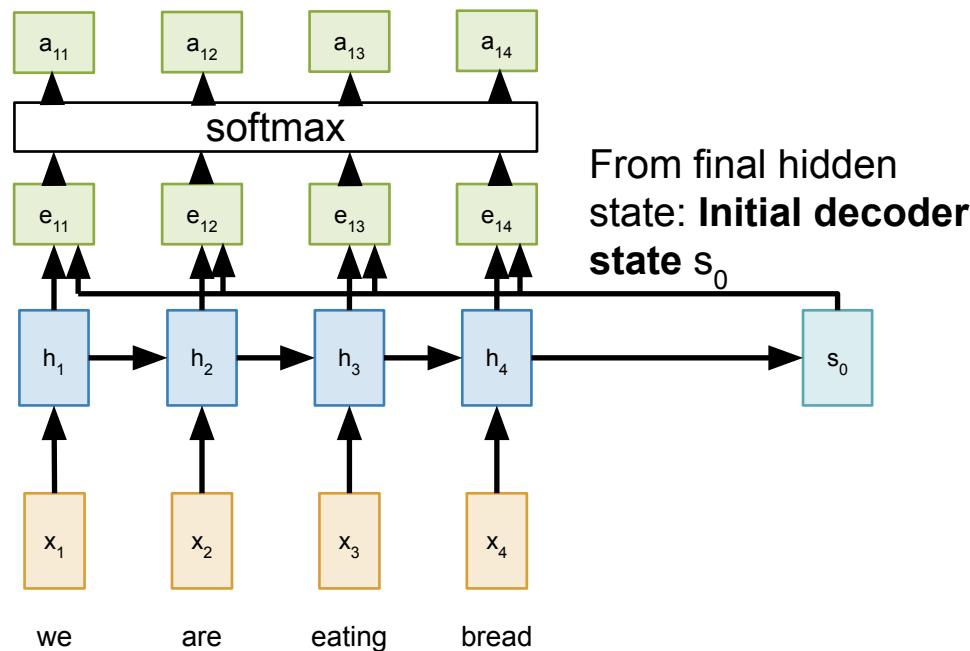
# Sequence to Sequence with RNNs and Attention

Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$       ( $f_{\text{att}}$  is an MLP)



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence to Sequence with RNNs and Attention



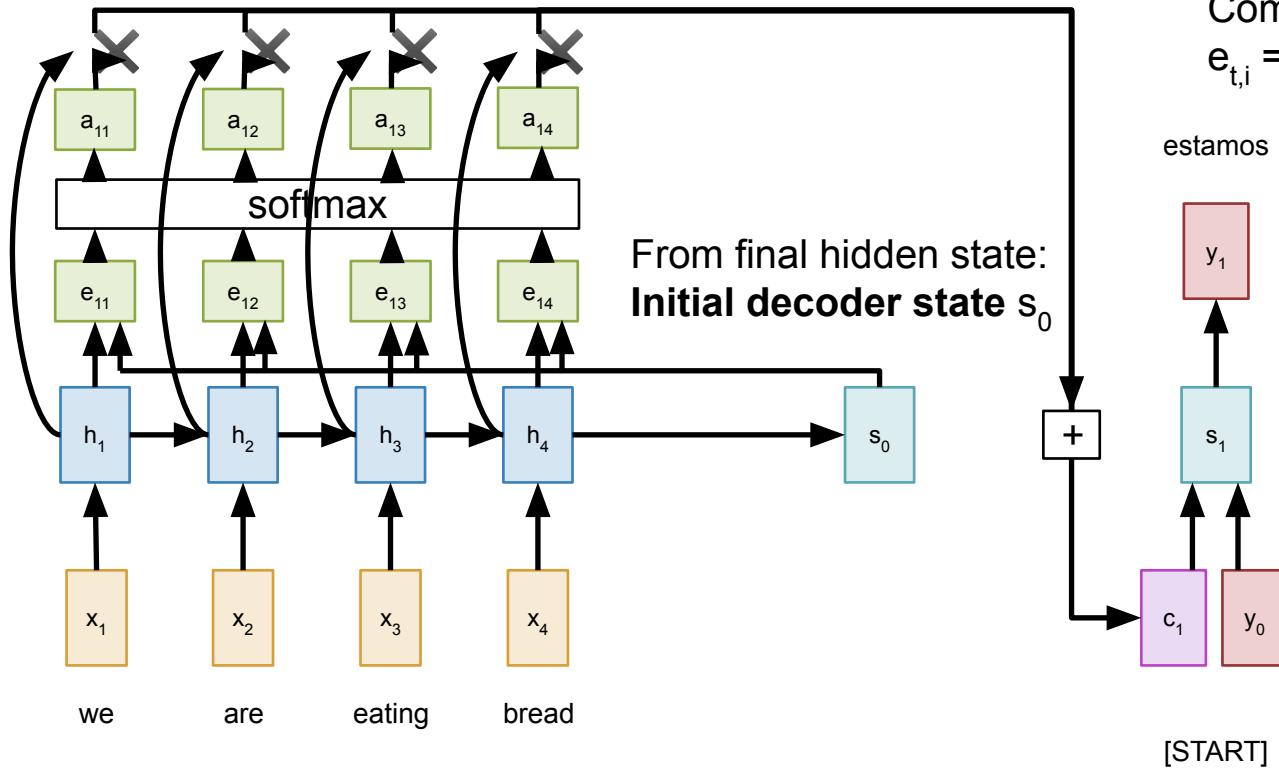
Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$       ( $f_{att}$  is an MLP)

Normalize alignment scores  
to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

From final hidden state: **Initial decoder state  $s_0$**

Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence to Sequence with RNNs and Attention



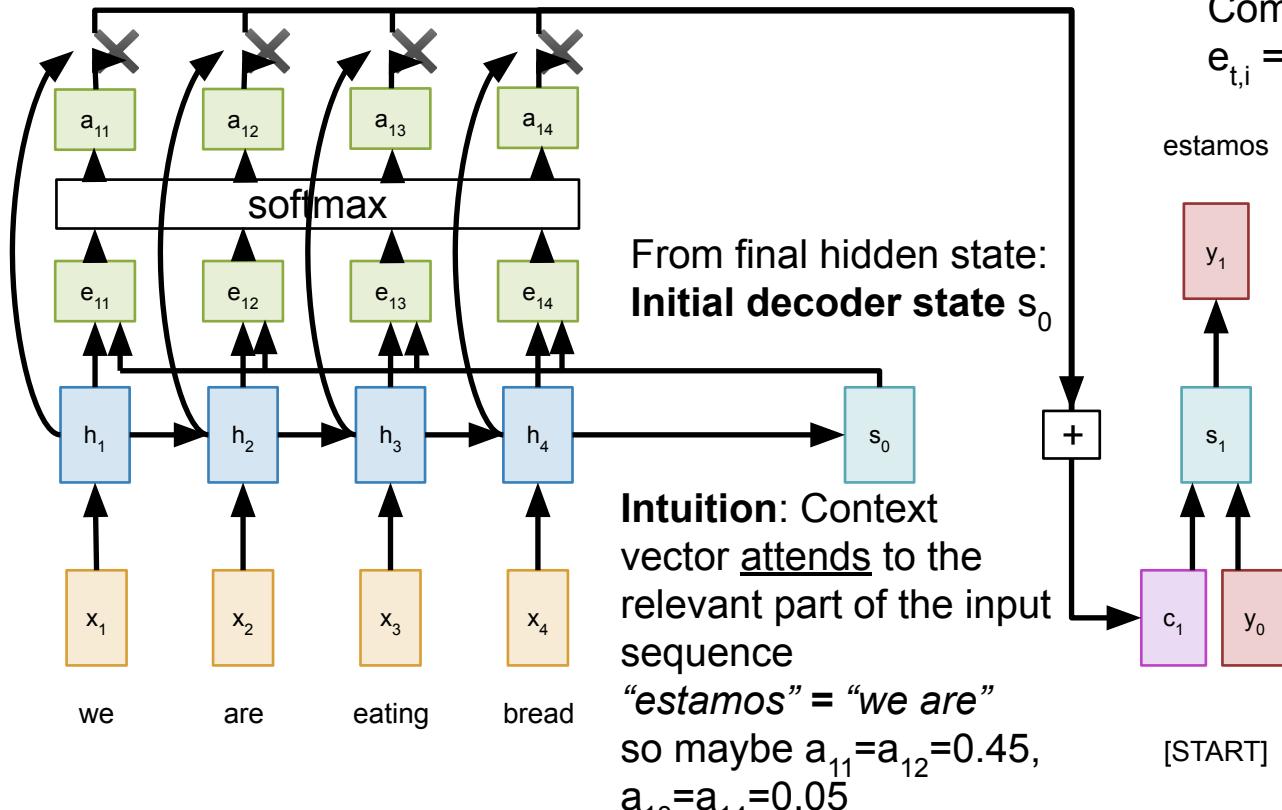
Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{att}(s_{t-1}, h_i)$  ( $f_{att}$  is an MLP)

estamos

Normalize alignment scores  
to get **attention weights**  
 $0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$

Compute context vector as  
linear combination of hidden  
states  
 $c_t = \sum_i a_{t,i} h_i$

# Sequence to Sequence with RNNs and Attention



Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$  ( $f_{\text{att}}$  is an MLP)

estamos

Normalize alignment scores  
to get **attention weights**

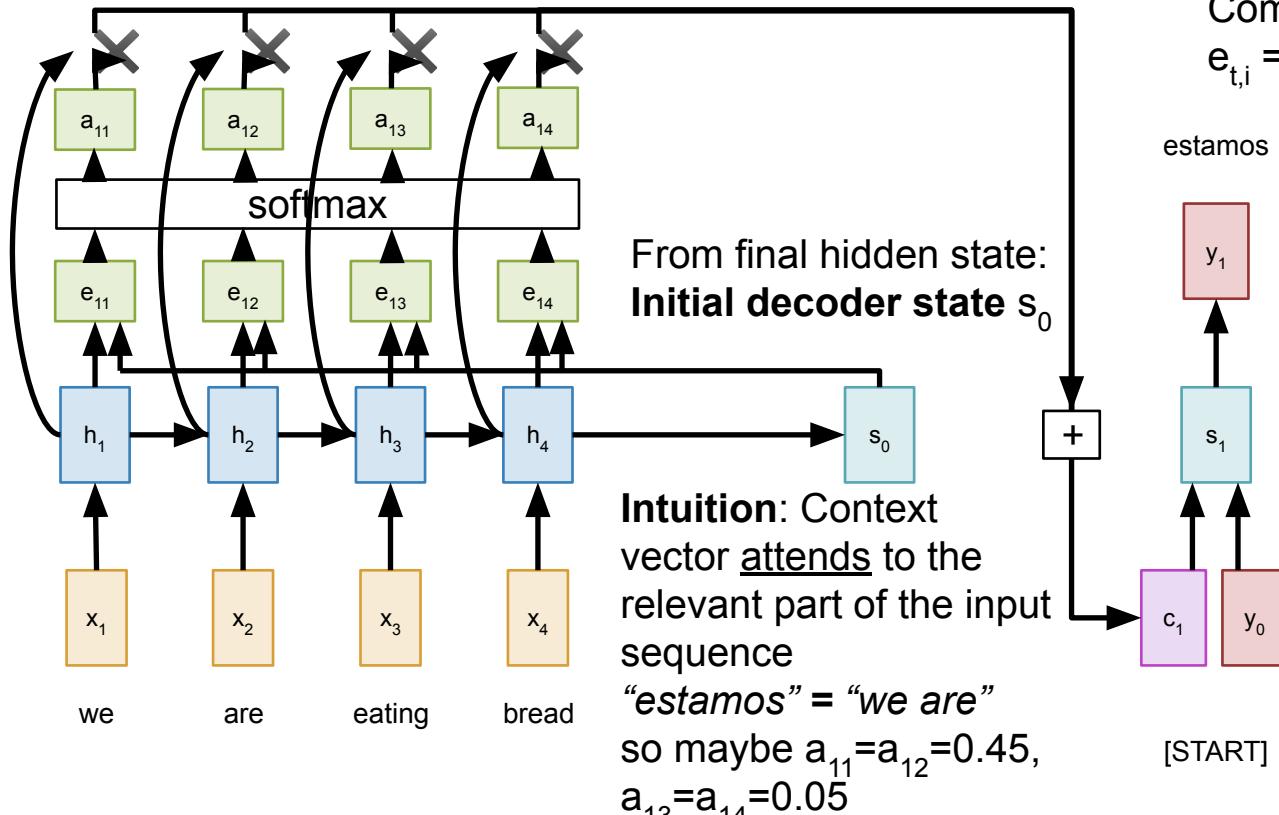
$$0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$$

Compute context vector as  
linear combination of hidden  
states

$$c_t = \sum_i a_{t,i} h_i$$

Use context vector in  
decoder:  $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

# Sequence to Sequence with RNNs and Attention



Compute (scalar) **alignment scores**  
 $e_{t,i} = f_{\text{att}}(s_{t-1}, h_i)$  ( $f_{\text{att}}$  is an MLP)

estamos

Normalize alignment scores  
to get **attention weights**

$$0 < a_{t,i} < 1 \quad \sum_i a_{t,i} = 1$$

Compute context vector as  
linear combination of hidden  
states

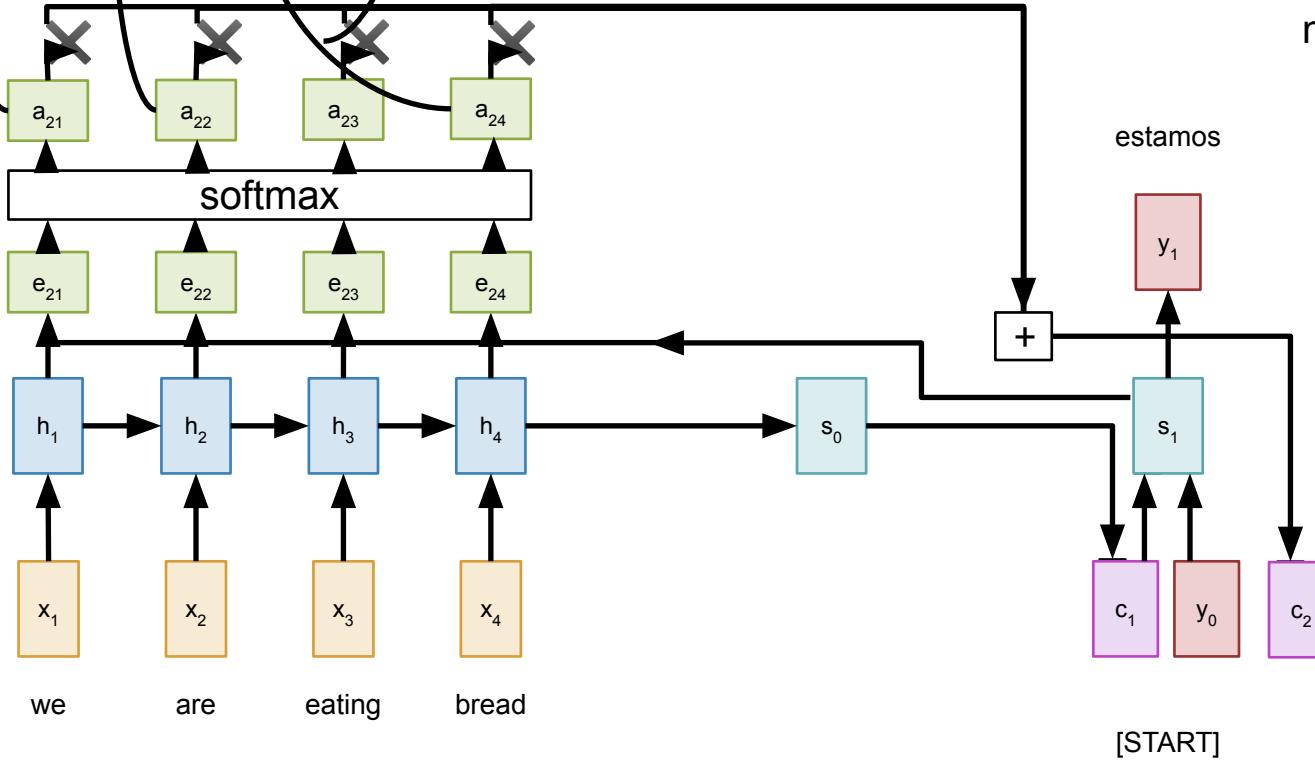
$$c_t = \sum_i a_{t,i} h_i$$

Use context vector in  
decoder:  $s_t = g_U(y_{t-1}, s_{t-1}, c_t)$

This is all differentiable! No  
supervision on attention  
weights – backprop through  
everything

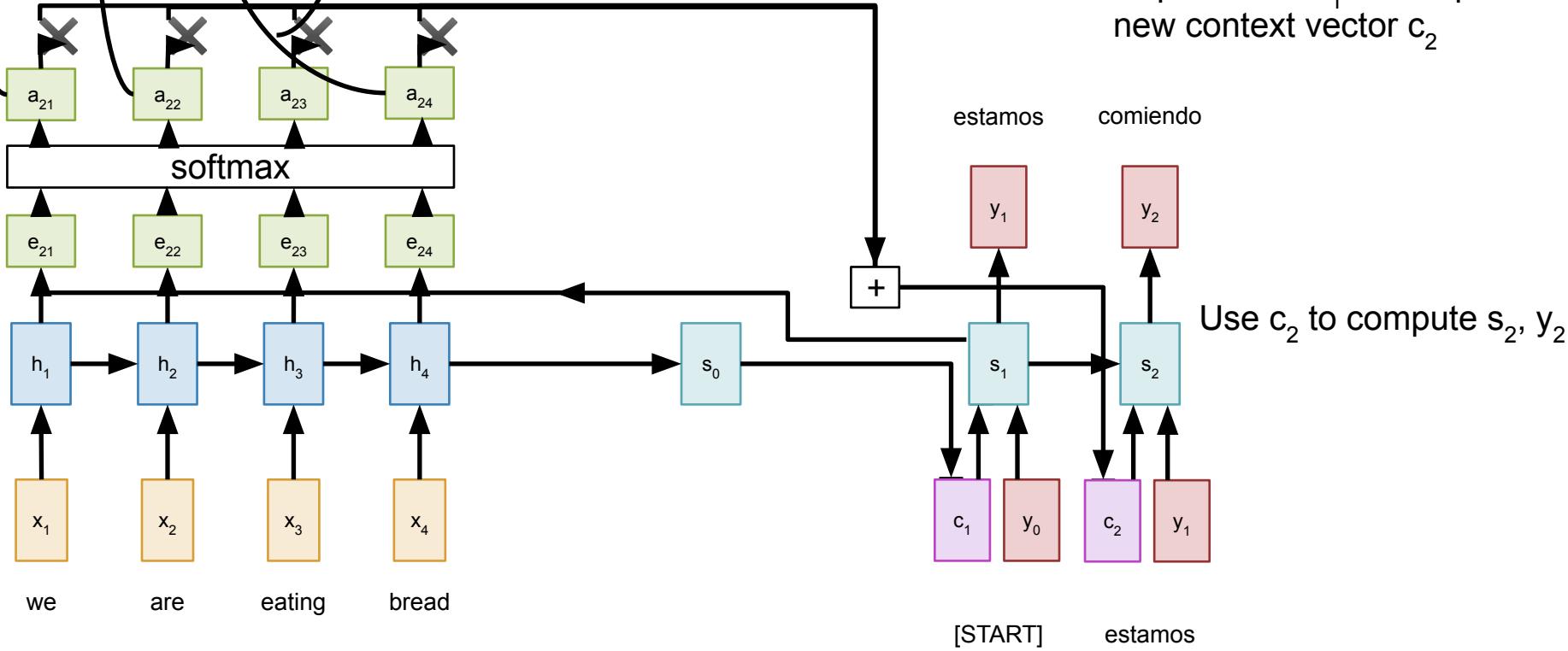
# Sequence to Sequence with RNNs and Attention

Repeat: Use  $s_1$  to compute new context vector  $c_2$



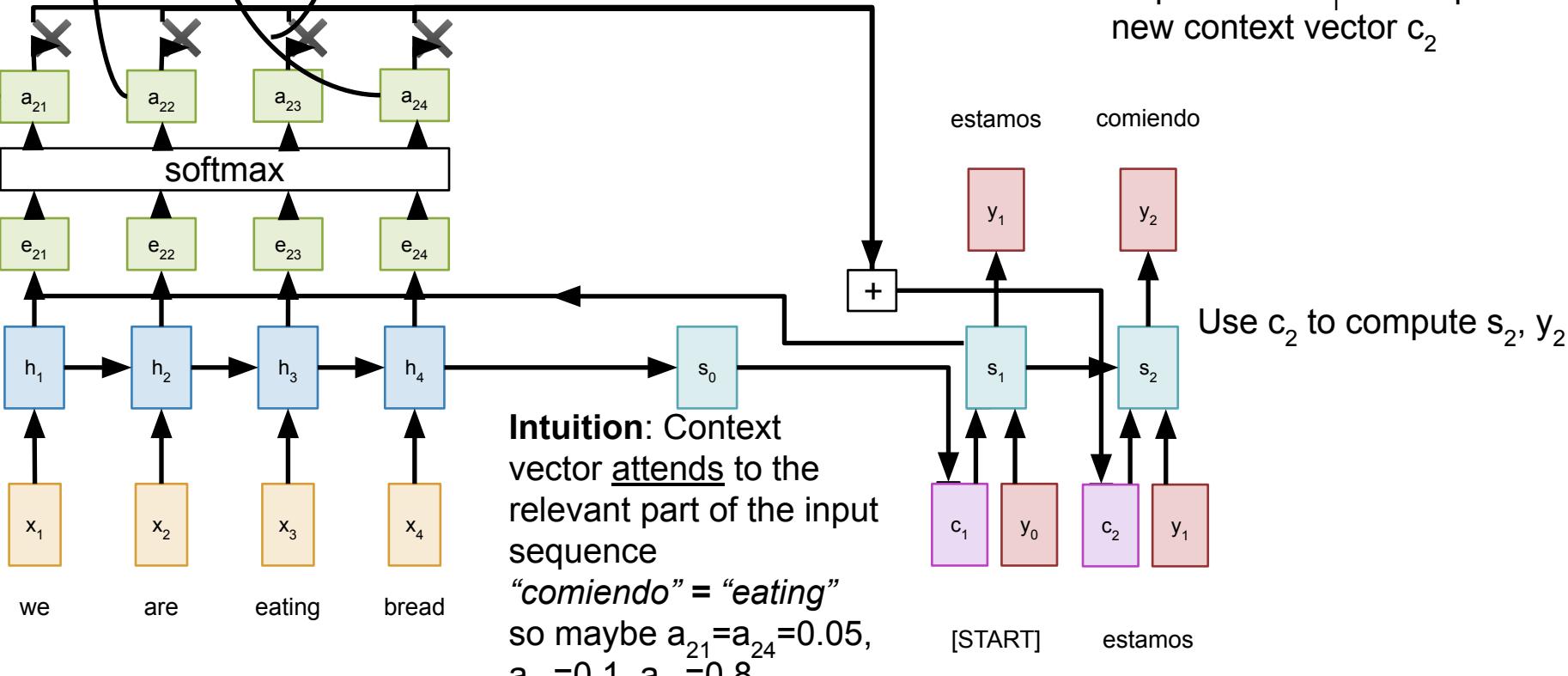
Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence to Sequence with RNNs and Attention



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence to Sequence with RNNs and Attention

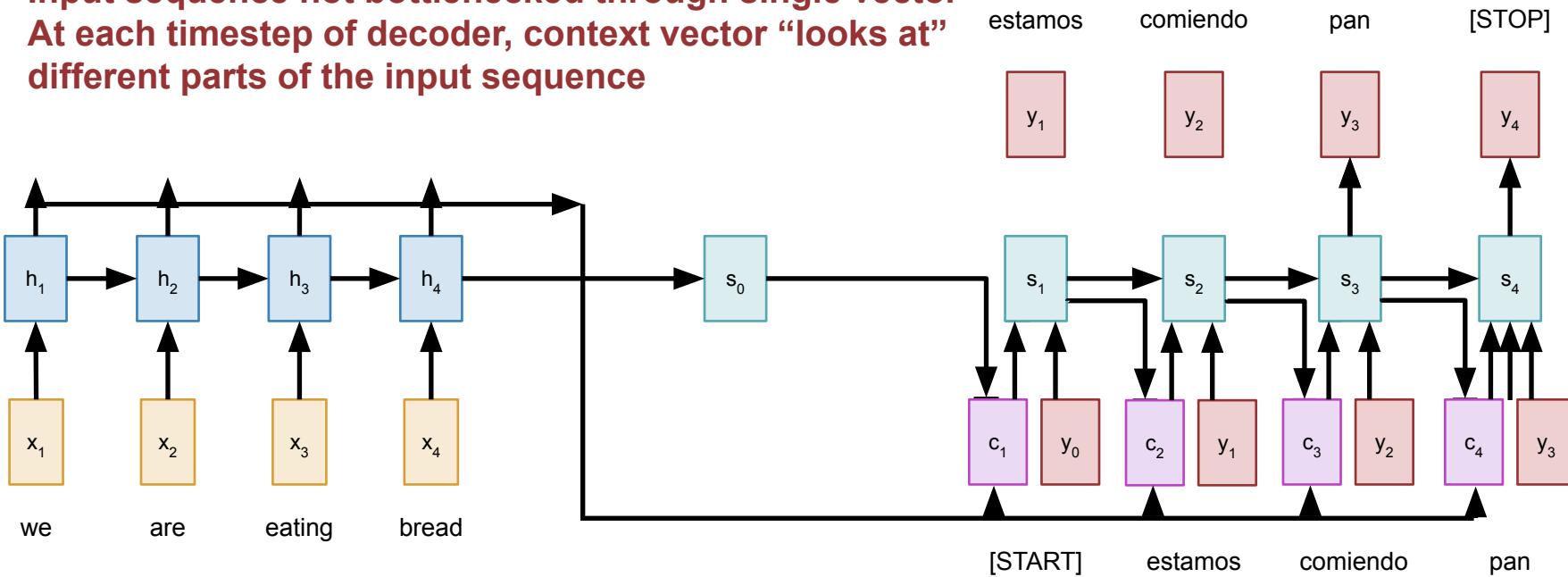


Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Sequence to Sequence with RNNs and Attention

Use a different context vector in each timestep of decoder

- Input sequence not bottlenecked through single vector
- At each timestep of decoder, context vector “looks at” different parts of the input sequence



Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

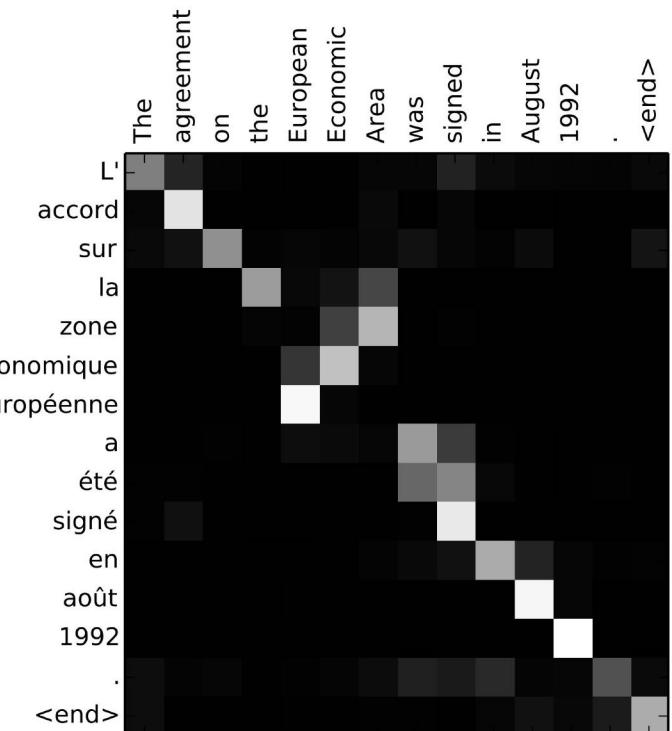
# Sequence to Sequence with RNNs and Attention

**Example:** English to French translation

**Input:** “The agreement on the European Economic Area was signed in August 1992.”

**Output:** “L'accord sur la zone économique européenne a été signé en août 1992.”

Visualize attention weights  $a_{t,i}$



Bahdanau et al, “Neural machine translation by jointly learning to align and translate”, ICLR 2015

# Sequence to Sequence with RNNs and Attention

**Example:** English to French translation

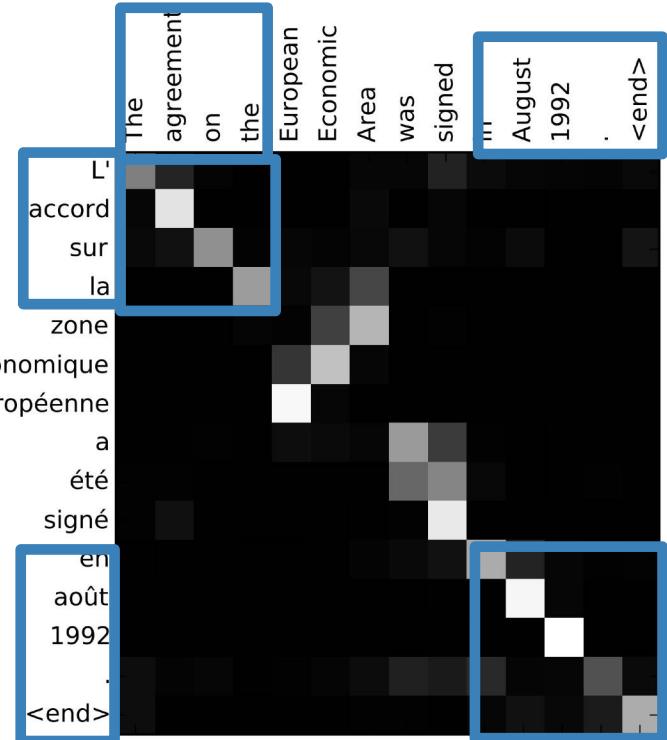
**Input:** “**The agreement on the European Economic Area was signed in August 1992.**”

**Output:** “**L'accord sur la zone économique européenne a été signé en août 1992.**”

Diagonal attention means words correspond in order

Diagonal attention means words correspond in order

Visualize attention weights  $a_{t,i}$



# Sequence to Sequence with RNNs and Attention

**Example:** English to French translation

**Input:** “**The agreement on the European Economic Area** was signed in **August 1992**.”

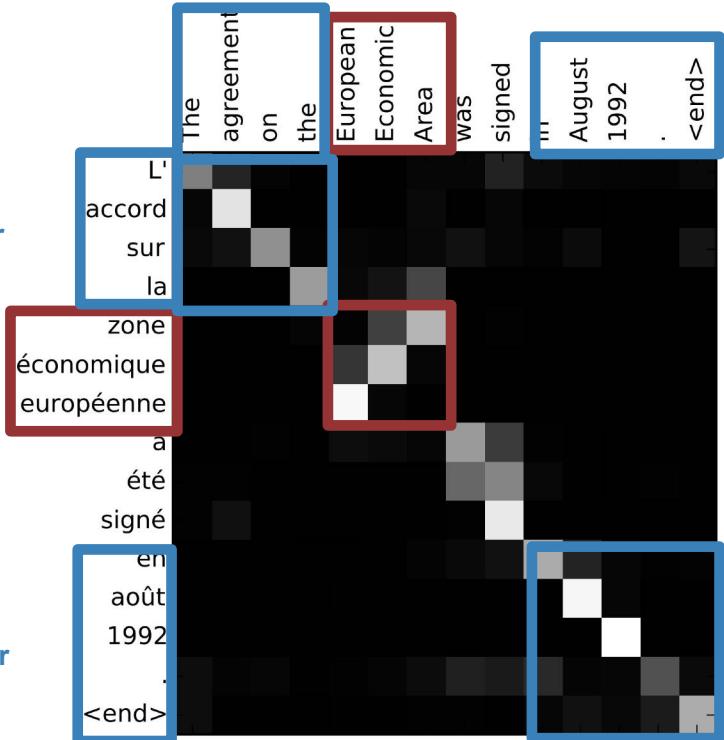
**Output:** “**L'accord sur la zone économique européenne** a été signé en **août 1992**.”

Diagonal attention means words correspond in order

Attention figures out different word orders

Diagonal attention means words correspond in order

Visualize attention weights  $a_{t,i}$

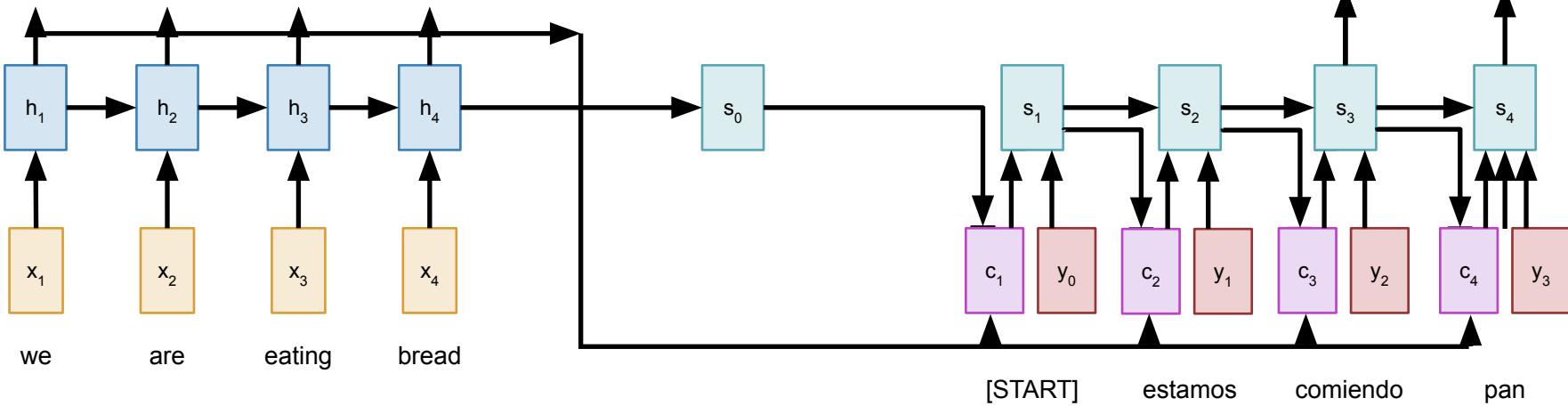


Bahdanau et al, “Neural machine translation by jointly learning to align and translate”, ICLR 2015

# Sequence to Sequence with RNNs and Attention

The decoder doesn't use the fact that  $h_i$  form an ordered sequence – it just treats them as an unordered set  $\{h_i\}$

Can use similar architecture given any set of input hidden vectors  $\{h_i\}$ !

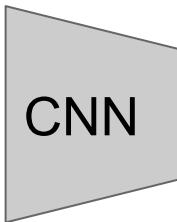
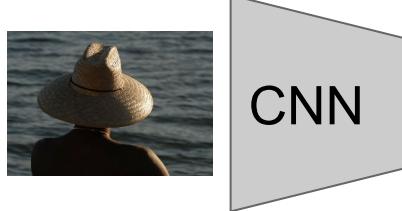


Bahdanau et al, "Neural machine translation by jointly learning to align and translate", ICLR 2015

# Image Captioning using spatial features

**Input:** Image  $I$

**Output:** Sequence  $\mathbf{y} = y_1, y_2, \dots, y_T$



$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Extract spatial  
features from a  
pretrained CNN

Features:  
 $H \times W \times D$

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning using spatial features

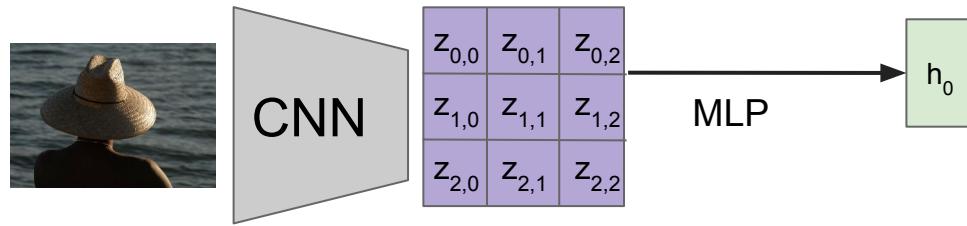
**Input:** Image  $I$

**Output:** Sequence  $\mathbf{y} = y_1, y_2, \dots, y_T$

**Encoder:**  $h_0 = f_w(\mathbf{z})$

where  $\mathbf{z}$  is spatial CNN features

$f_w(\cdot)$  is an MLP



Extract spatial  
features from a  
pretrained CNN

Features:  
 $H \times W \times D$

Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning using spatial features

**Input:** Image  $I$

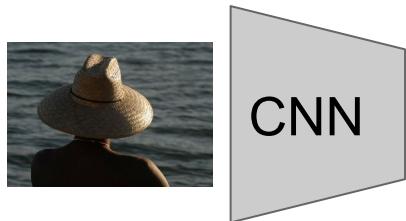
**Output:** Sequence  $\mathbf{y} = y_1, y_2, \dots, y_T$

**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c)$   
where context vector  $c$  is often  $c = h_0$

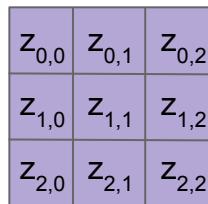
**Encoder:**  $h_0 = f_w(\mathbf{z})$

where  $\mathbf{z}$  is spatial CNN features

$f_w(\cdot)$  is an MLP

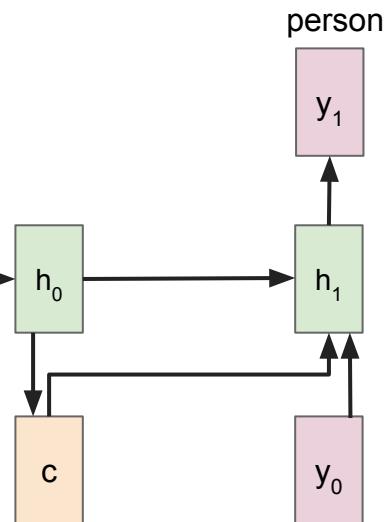


Extract spatial  
features from a  
pretrained CNN



Features:  
 $H \times W \times D$

MLP



[START]

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning using spatial features

**Input:** Image  $I$

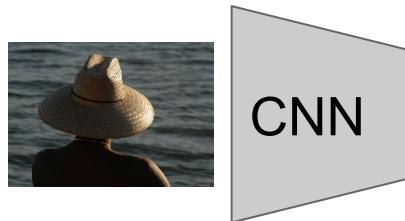
**Output:** Sequence  $\mathbf{y} = y_1, y_2, \dots, y_T$

**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c)$   
where context vector  $c$  is often  $c = h_0$

**Encoder:**  $h_0 = f_w(\mathbf{z})$

where  $\mathbf{z}$  is spatial CNN features

$f_w(\cdot)$  is an MLP

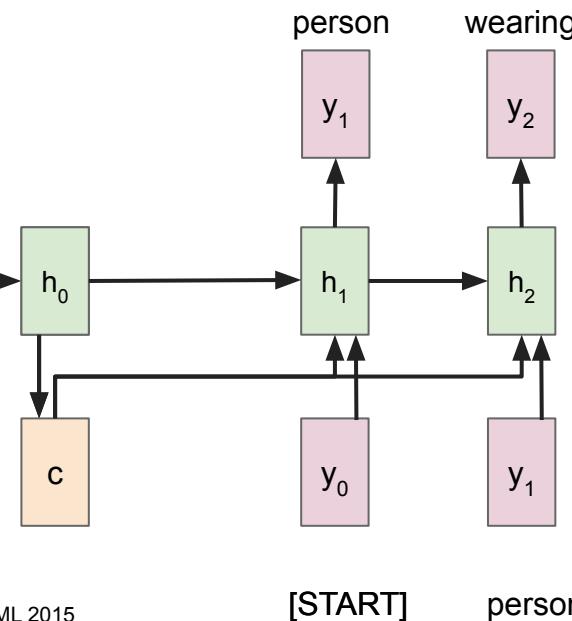


Extract spatial  
features from a  
pretrained CNN

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

MLP



Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

[START] person

# Image Captioning using spatial features

**Input:** Image  $I$

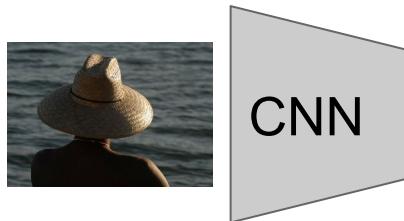
**Output:** Sequence  $\mathbf{y} = y_1, y_2, \dots, y_T$

**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c)$   
where context vector  $c$  is often  $c = h_0$

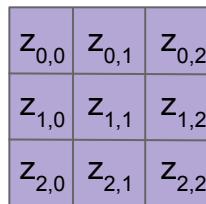
**Encoder:**  $h_0 = f_w(\mathbf{z})$

where  $\mathbf{z}$  is spatial CNN features

$f_w(\cdot)$  is an MLP

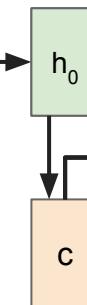


Extract spatial  
features from a  
pretrained CNN

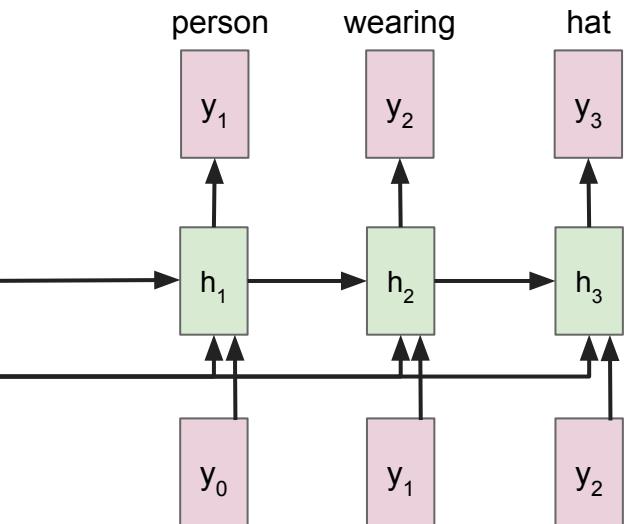


Features:  
 $H \times W \times D$

MLP



$c$



Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

[START]

person

wearing

# Image Captioning using spatial features

**Input:** Image  $I$

**Output:** Sequence  $\mathbf{y} = y_1, y_2, \dots, y_T$

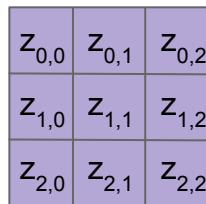
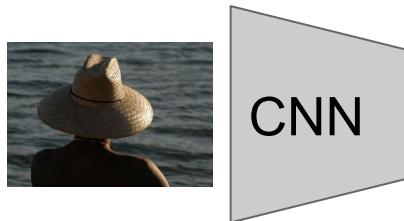
**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c)$

where context vector  $c$  is often  $c = h_0$

**Encoder:**  $h_0 = f_w(\mathbf{z})$

where  $\mathbf{z}$  is spatial CNN features

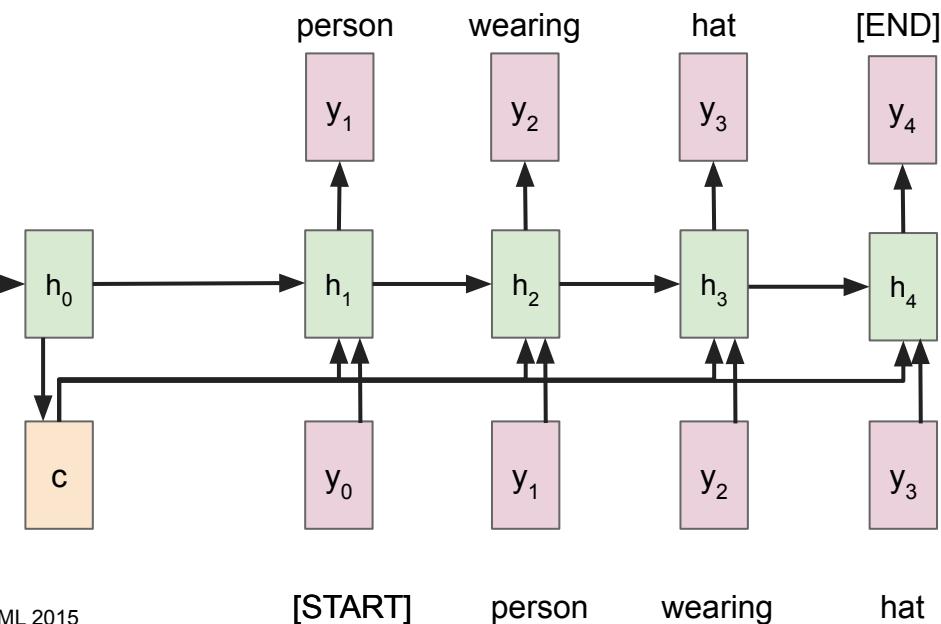
$f_w(\cdot)$  is an MLP



MLP

Extract spatial  
features from a  
pretrained CNN

Features:  
 $H \times W \times D$



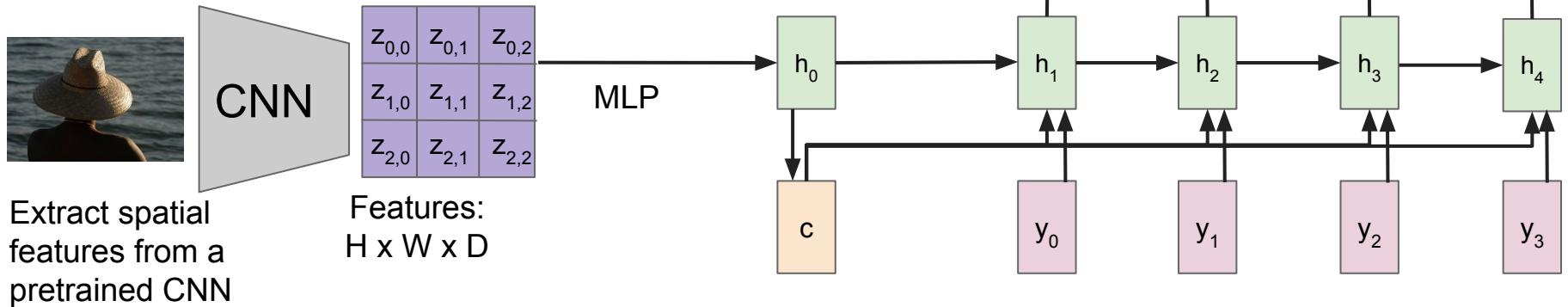
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning using spatial features

**Problem: Input is "bottlenecked" through c**

- Model needs to encode everything it wants to say within c

This is a problem if we want to generate really long descriptions? 100s of words long



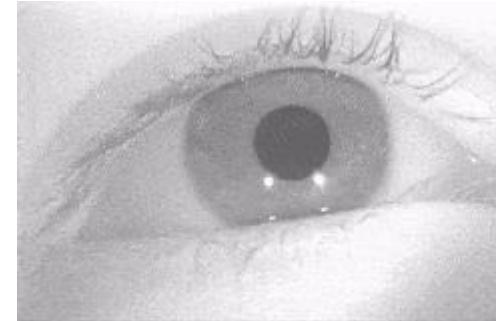
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with RNNs and **Attention**

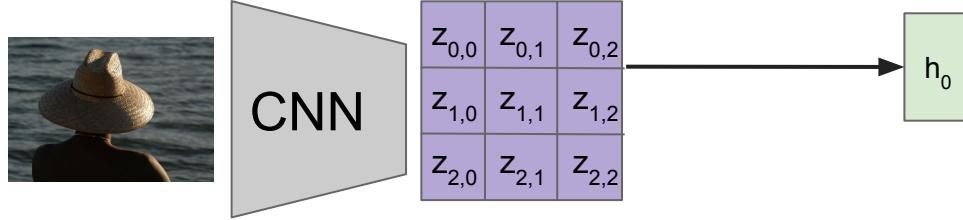
[gif source](#)

Attention idea: New context vector at every time step.

Each context vector will attend to different image regions



Attention Saccades in humans



Extract spatial  
features from a  
pretrained CNN

Features:  
 $H \times W \times D$

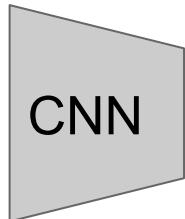
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with RNNs and Attention

Compute alignments scores (scalars):

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$f_{att}(\cdot)$  is an MLP



Extract spatial features from a pretrained CNN

Features:  
 $H \times W \times D$

Alignment scores:

$H \times W$

$e_{1,0,0}$	$e_{1,0,1}$	$e_{1,0,2}$
$e_{1,1,0}$	$e_{1,1,1}$	$e_{1,1,2}$
$e_{1,2,0}$	$e_{1,2,1}$	$e_{1,2,2}$

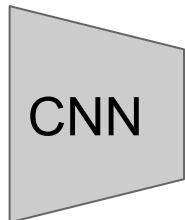


# Image Captioning with RNNs and Attention

Compute alignments scores (scalars):

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$f_{att}(\cdot)$  is an MLP



Extract spatial features from a pretrained CNN

Alignment scores:

$H \times W$

$e_{1,0,0}$	$e_{1,0,1}$	$e_{1,0,2}$
$e_{1,1,0}$	$e_{1,1,1}$	$e_{1,1,2}$
$e_{1,2,0}$	$e_{1,2,1}$	$e_{1,2,2}$

Attention:

$H \times W$

$a_{1,0,0}$	$a_{1,0,1}$	$a_{1,0,2}$
$a_{1,1,0}$	$a_{1,1,1}$	$a_{1,1,2}$
$a_{1,2,0}$	$a_{1,2,1}$	$a_{1,2,2}$

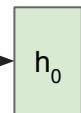
Normalize to get attention weights:

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$0 < a_{t,i,j} < 1$ ,  
attention values sum to 1

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$



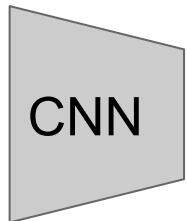
Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with RNNs and Attention

Compute alignments scores (scalars):

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$f_{att}(\cdot)$  is an MLP



Extract spatial features from a pretrained CNN

Alignment scores:

$$H \times W$$

$e_{1,0,0}$	$e_{1,0,1}$	$e_{1,0,2}$
$e_{1,1,0}$	$e_{1,1,1}$	$e_{1,1,2}$
$e_{1,2,0}$	$e_{1,2,1}$	$e_{1,2,2}$

Attention:

$$H \times W$$

$a_{1,0,0}$	$a_{1,0,1}$	$a_{1,0,2}$
$a_{1,1,0}$	$a_{1,1,1}$	$a_{1,1,2}$
$a_{1,2,0}$	$a_{1,2,1}$	$a_{1,2,2}$

Normalize to get attention weights:

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

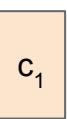
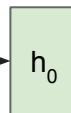
$0 < a_{t,i,j} < 1$ ,  
attention values sum to 1

Compute context vector:

$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

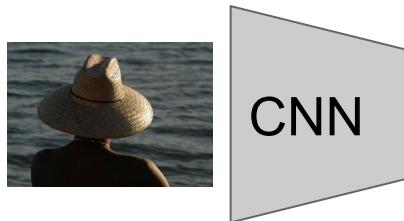
# Image Captioning with RNNs and Attention

Each timestep of decoder uses a different context vector that looks at different parts of the input image

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$

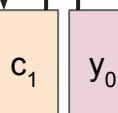
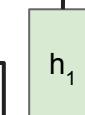
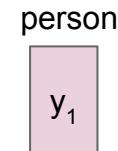


Extract spatial features from a pretrained CNN

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

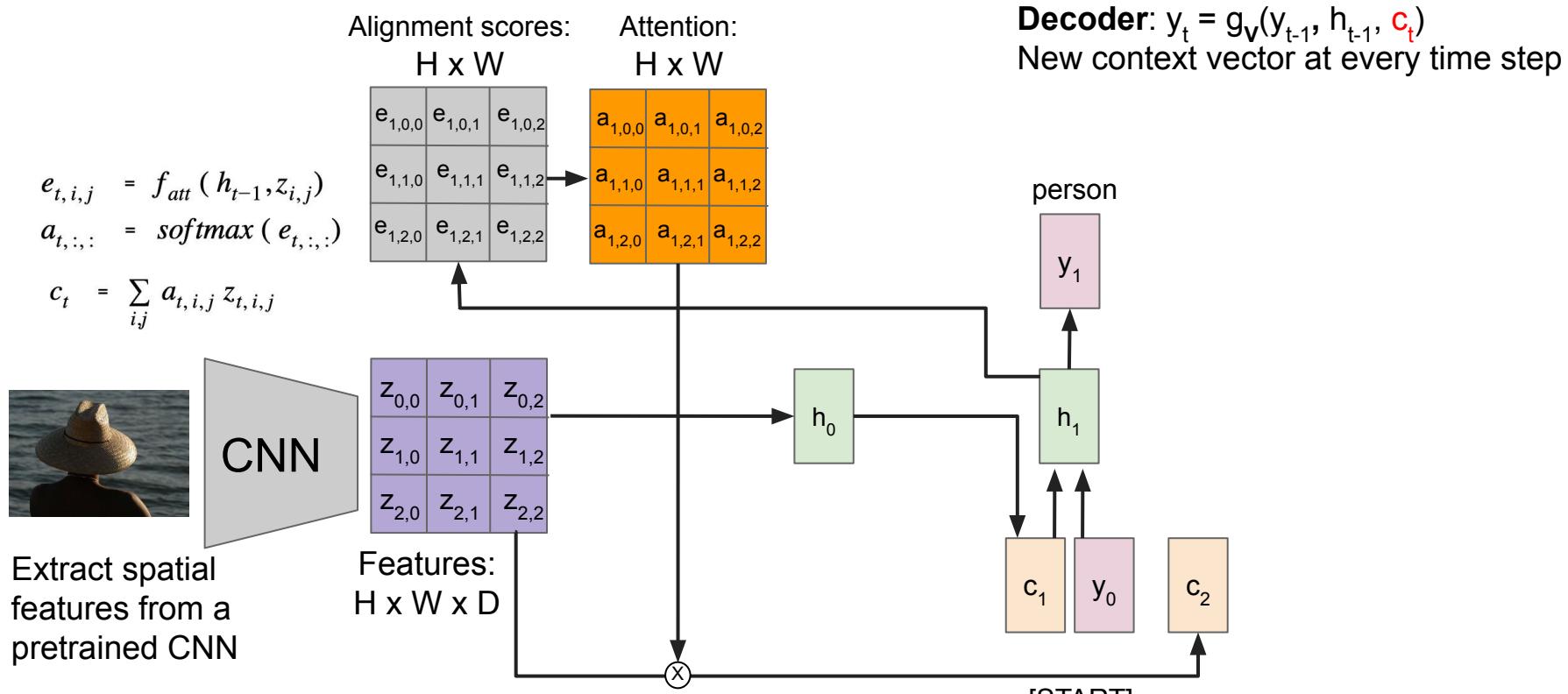
**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c_t)$   
New context vector at every time step



[START]

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with RNNs and Attention



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

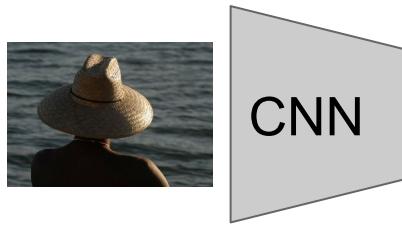
# Image Captioning with RNNs and Attention

Each timestep of decoder uses a different context vector that looks at different parts of the input image

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$

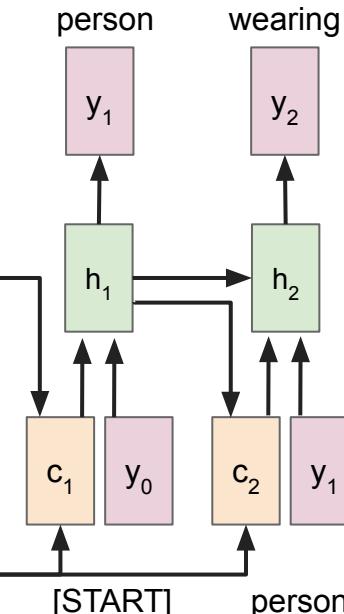


Extract spatial features from a pretrained CNN

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

**Decoder:**  $y_t = g_v(y_{t-1}, h_{t-1}, c_t)$   
New context vector at every time step



Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

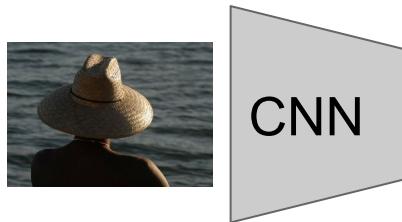
# Image Captioning with RNNs and Attention

Each timestep of decoder uses a different context vector that looks at different parts of the input image

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$



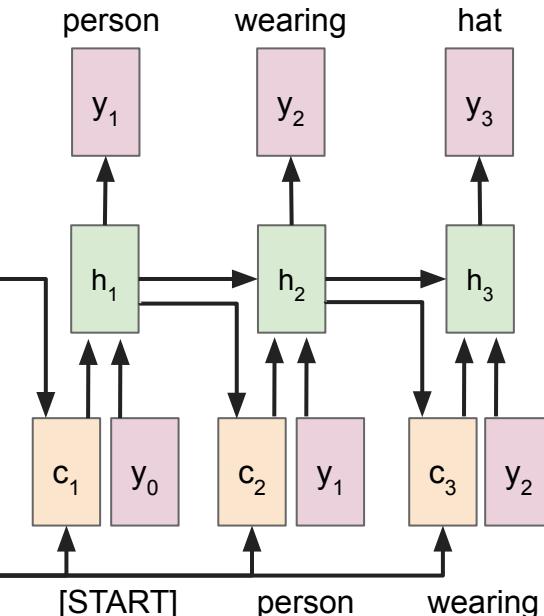
Extract spatial features from a pretrained CNN

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

Decoder:  $y_t = g_v(y_{t-1}, h_{t-1}, c_t)$

New context vector at every time step



Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

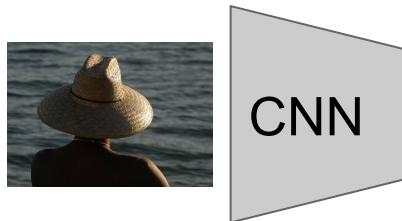
# Image Captioning with RNNs and Attention

Each timestep of decoder uses a different context vector that looks at different parts of the input image

$$e_{t,i,j} = f_{att}(h_{t-1}, z_{i,j})$$

$$a_{t,:,:} = \text{softmax}(e_{t,:,:})$$

$$c_t = \sum_{i,j} a_{t,i,j} z_{t,i,j}$$



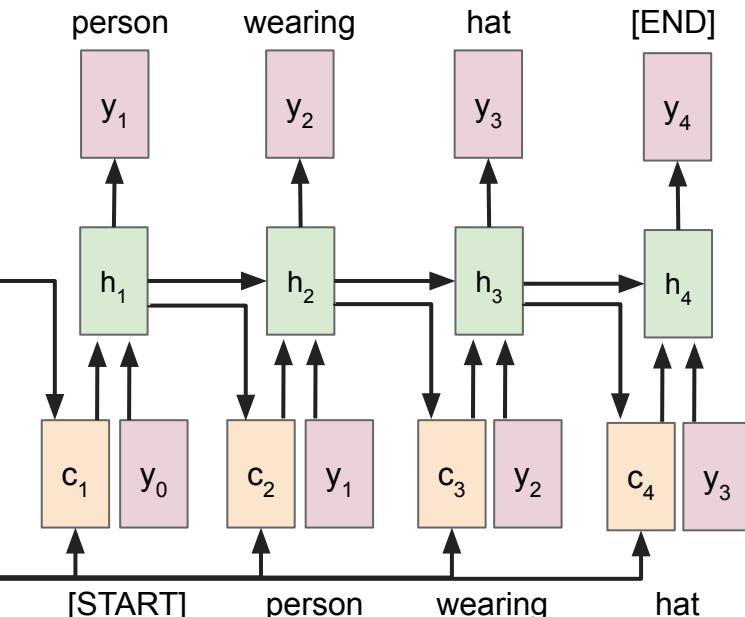
Extract spatial  
features from a  
pretrained CNN

$z_{0,0}$	$z_{0,1}$	$z_{0,2}$
$z_{1,0}$	$z_{1,1}$	$z_{1,2}$
$z_{2,0}$	$z_{2,1}$	$z_{2,2}$

Features:  
 $H \times W \times D$

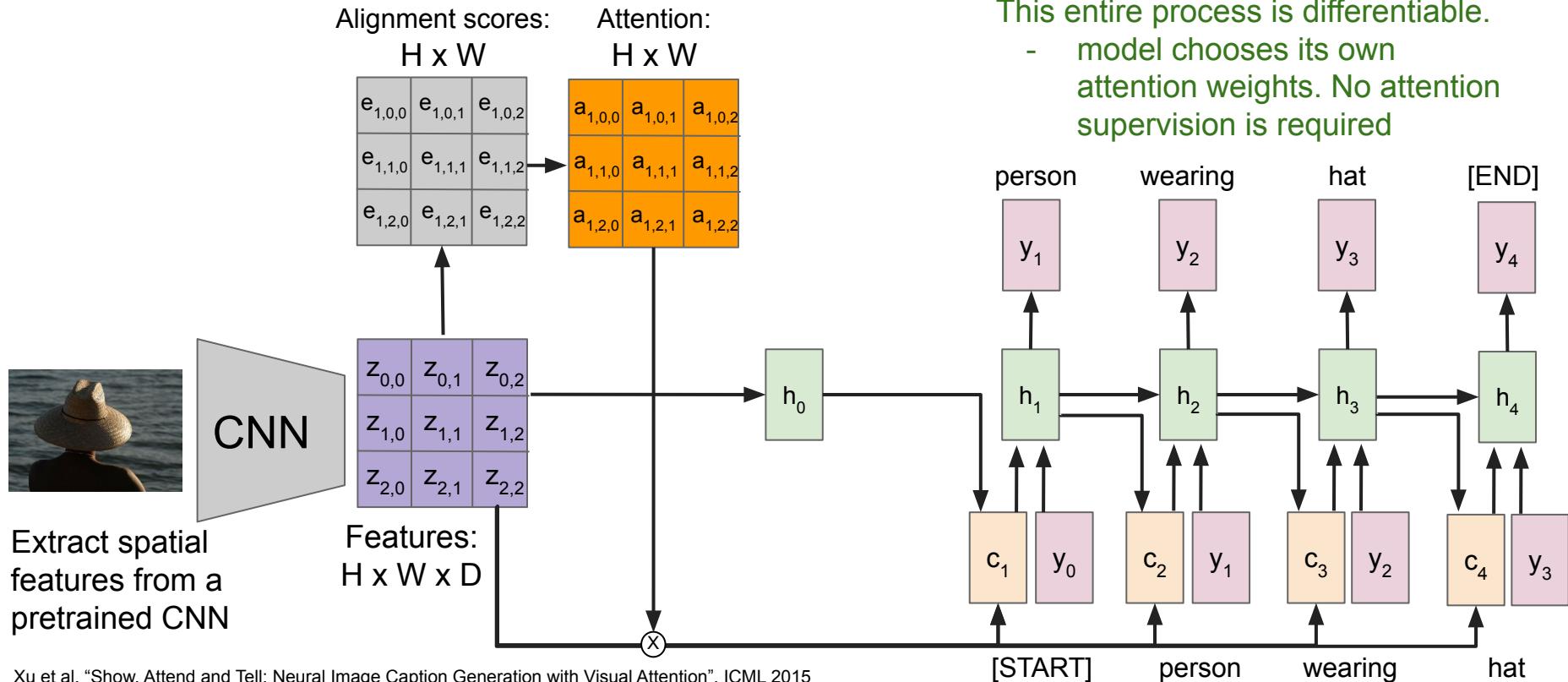
$$\text{Decoder: } y_t = g_v(y_{t-1}, h_{t-1}, c_t)$$

New context vector at every time step



Xu et al., "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with RNNs and Attention



Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

# Image Captioning with Attention



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



A giraffe standing in a forest with trees in the background.

Xu et al., "Show, Attend, and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015

Figure copyright Kelvin Xu, Jimmy Lei Ba, Jamie Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio, 2015. Reproduced with permission.