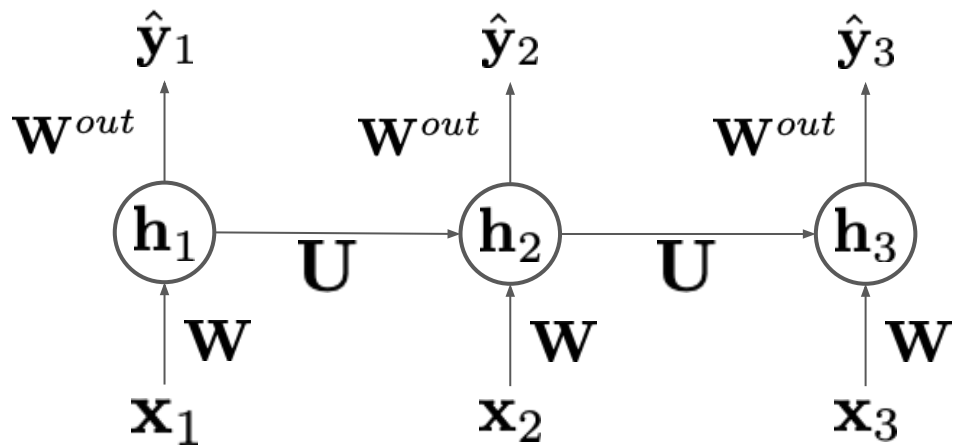# Language and Vision

Instructor: Seunghoon Hong

# Course logistics

- No lecture on the next week

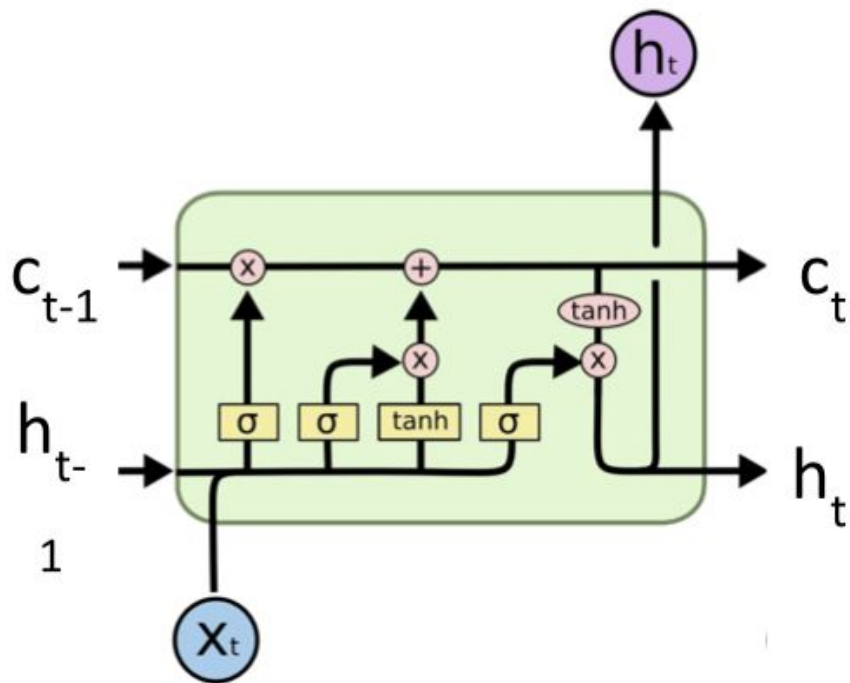# Recap: (Vanilla) Recurrent Neural Network



In general, for any $t \geq 1$,

$$\mathbf{h}_t = \sigma(\mathbf{U}\mathbf{h}_{t-1} + \mathbf{W}\mathbf{x}_t + \mathbf{b})$$

$$\hat{\mathbf{y}}_t = \mathbf{W}^{out}\mathbf{h}_t$$

$$\mathbf{h}_0 = \mathbf{0}$$

# Recap: Long-Short Term Memory



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

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

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

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

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

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

4

# Today's agenda

- Language modeling using RNNs
- Image captioning
    - Naive image captioning, image captioning with attention
- Visual question answering
    - Naive visual question answering, memory network

# Today's agenda

- **Language modeling using RNNs**
- Image captioning
    - Naive image captioning, image captioning with attention
- Visual question answering
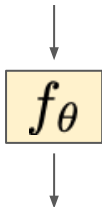    - Naive visual question answering, memory network

# Modeling language

**Sentence generation**

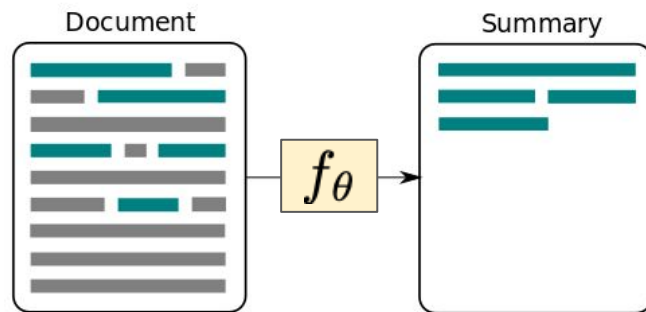$f_\theta$ → I am very hungry at the end of every class

**Machine translation**

The agreement on the European Economic
Area was signed in August 1992.

↓

$f_\theta$

↓

L'accord sur l'Espace économique
européen a été signé en août 1992.

**Text summarization**

Document → $f_\theta$ → Summary

**And so many more...**

# Modeling language

- Sentence = a sequence of discrete symbols

| A | Man | is | holding | a | bat |
|---|---|---|---|---|---|

$\mathbf{x}_t \in \mathbb{R}^n$

| [00001000] | [01000000] | [00000010] | [00000100] | [00001000] | [10000000] |
|---|---|---|---|---|---|
| $\mathbf{x}_0$ | $\mathbf{x}_1$ | $\mathbf{x}_2$ | $\mathbf{x}_3$ | $\mathbf{x}_4$ | $\mathbf{x}_5$ |

**One-hot encoding of discrete symbols (tokenization)**

# RNN as a language model

- Sentence generation = predicting a next token

A

$\hat{\mathbf{x}}_{t+1} \in \mathbb{R}^n$  [00001000]  Classification layer  $\hat{\mathbf{x}}_{t+1} = \mathbf{W}_d \mathbf{h}_t + \mathbf{b}_d$

$\mathbf{W}_d \in \mathbb{R}^{n \times d}$

$\mathbf{h}_t \in \mathbb{R}^d$  Continuous output embedding  $\mathbf{h}_t = \text{LSTM}(\mathbf{z}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}; \theta)$

LSTM  Recurrent neural network

$\mathbf{z}_t \in \mathbb{R}^d$  Continuous word embedding  $\mathbf{z}_t = \mathbf{W}_e \mathbf{x}_t + \mathbf{b}_e$

$\mathbf{W}_e \in \mathbb{R}^{d \times n}$

$\mathbf{x}_t \in \mathbb{R}^n$  [00100000]  One-hot encoding

<SOS>

9

# RNN as a language model

- Sentence generation = predicting a next token

A          Man

$\hat{\mathbf{x}}_{t+1} \in \mathbb{R}^n$  **[00001000]**  **[01000000]**

$\mathbf{W}_d \in \mathbb{R}^{n \times d}$

$\mathbf{h}_t \in \mathbb{R}^d$

| LSTM | → | LSTM |

$\mathbf{z}_t \in \mathbb{R}^d$

$\mathbf{W}_e \in \mathbb{R}^{d \times n}$

$\mathbf{x}_t \in \mathbb{R}^n$  **[00100000]**  **[00001000]**

<SOS>          A

# RNN as a language model

- Sentence generation = predicting a next token



A    Man    is    holding    a    bat    <EOS>

$\hat{\mathbf{x}}_{t+1} \in \mathbb{R}^n$ [00001000] [01000000] [00000010] [00000100] [00001000] [10000000] [00001000]

$\mathbf{W}_d \in \mathbb{R}^{n \times d}$

$\mathbf{h}_t \in \mathbb{R}^d$

LSTM    LSTM    LSTM    LSTM    LSTM    LSTM    LSTM

$\mathbf{z}_t \in \mathbb{R}^d$

$\mathbf{W}_e \in \mathbb{R}^{d \times n}$

$\mathbf{x}_t \in \mathbb{R}^n$ [00100000] [00001000] [01000000] [00000010] [00000100] [00001000] [10000000]

<SOS>    A    Man    is    holding    a    bat

# Training: RNN-based language model

$$\sum_{t=1}^{T} \mathcal{L}(\hat{\mathbf{x}}_t, \mathbf{x}_t)$$

Cross-entropy loss (word-level classification)



| $\hat{\mathbf{x}}_1$ | $\hat{\mathbf{x}}_2$ | $\hat{\mathbf{x}}_3$ | $\hat{\mathbf{x}}_4$ | $\hat{\mathbf{x}}_5$ | $\hat{\mathbf{x}}_6$ | $\hat{\mathbf{x}}_7$ |
|---|---|---|---|---|---|---|
| A | Man | is | holding | a | bat | <EOS> |
| LSTM | LSTM | LSTM | LSTM | LSTM | LSTM | LSTM |
| <SOS> | A | Man | is | holding | a | bat |
| $\mathbf{x}_0$ | $\mathbf{x}_1$ | $\mathbf{x}_2$ | $\mathbf{x}_3$ | $\mathbf{x}_4$ | $\mathbf{x}_5$ | $\mathbf{x}_6$ |

We feed ground-truth words as inputs (also known as **teacher forcing**)

12

# Inference: RNN-based language model



- **For each step, sample a word from the output score**
- **Sampling methods:**

  - Take the word with maximum score (greedy, deterministic)
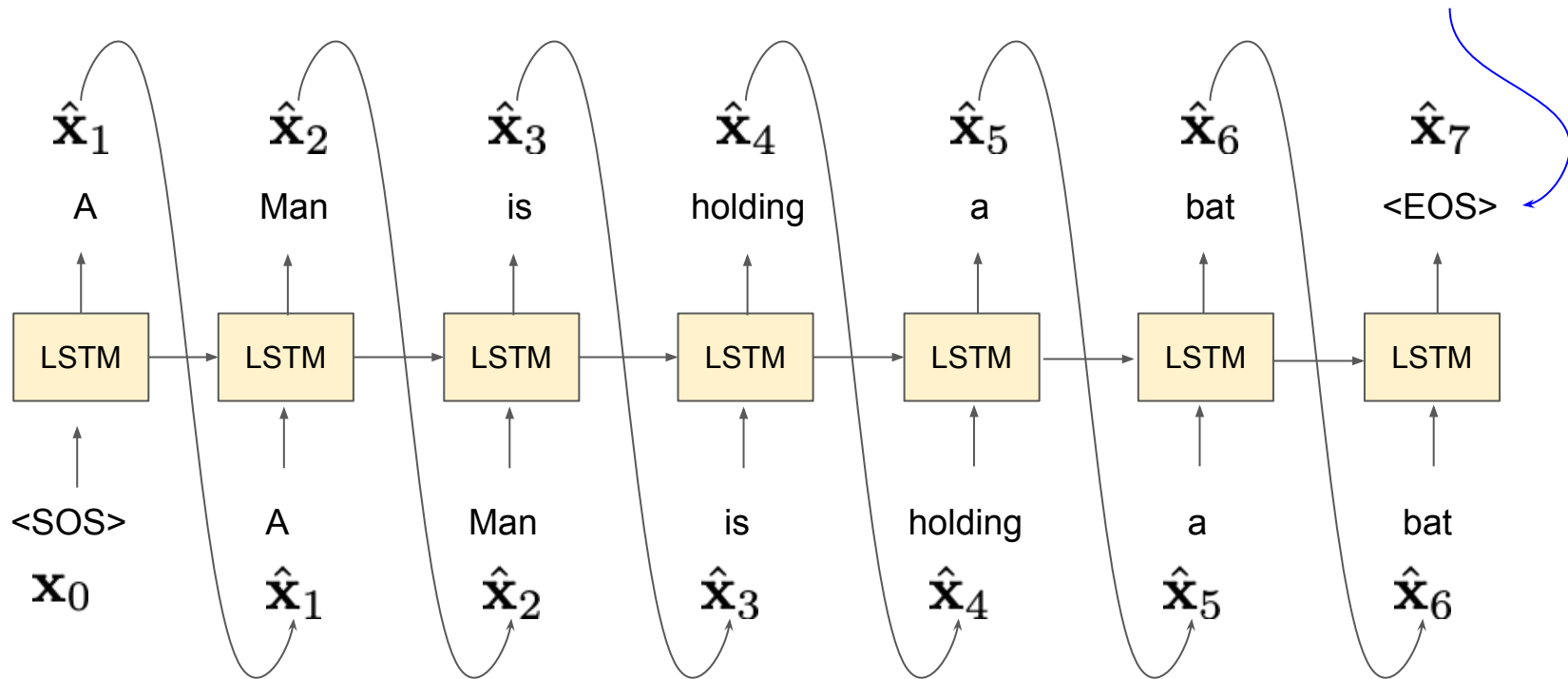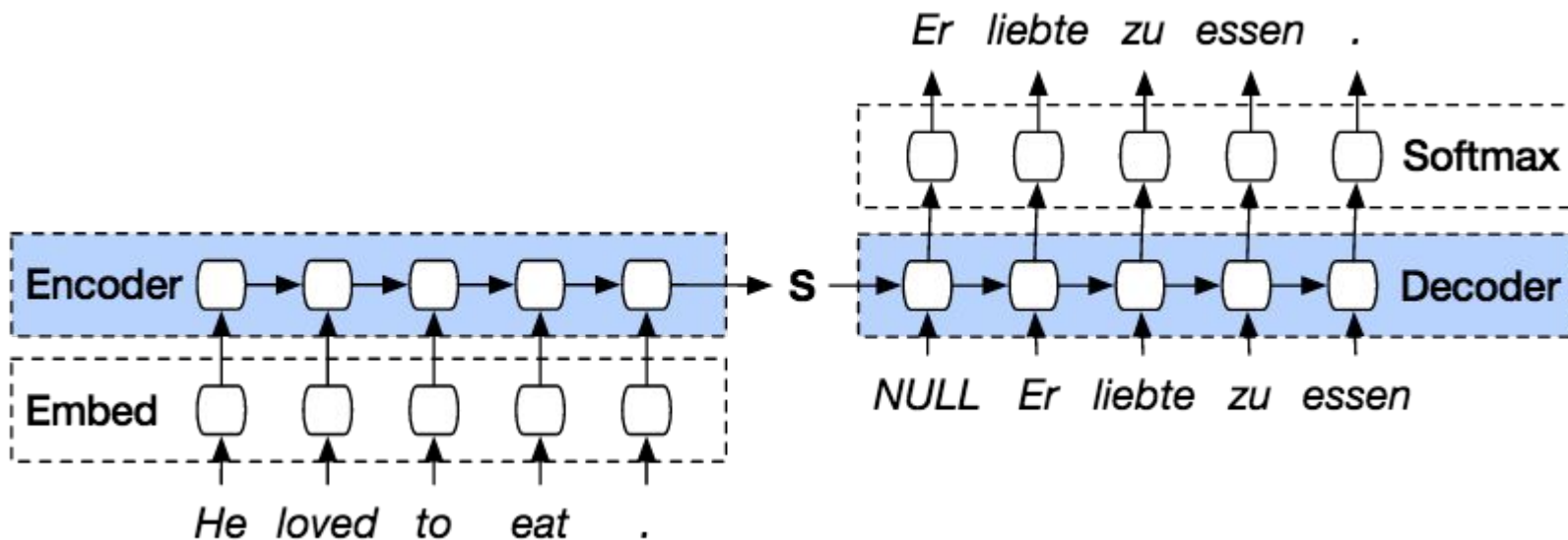  - Sample a word according to score probability (stochastic)

$\hat{\mathbf{x}}_1$

LSTM

$\mathbf{x}_0$

$\text{argmax } S(x_t)$

1 2 3 4 5 6 …    N

**Greedy method**

$x_t \sim S(x_t)$

1 2 3 4 5 6 …    N

**Probabilistic method**

# Inference: RNN-based language model

Stop sampling when it samples the end-of-sentence symbol

# Machine translation

- Translate a sentence in one language to another

# Summary: LSTM-based language model

- Sentence = a sequence of discrete symbols
- RNN (e.g. LSTM) for modeling a sequence of discrete symbols
  - Each word: an one-hot encoding
  - Sentence generation: prediction of the next word given the previous words
  - Training: sequential classification (classification of each word at a time)
  - Inference: sequentially predict a word and use it as an input to the next step

# Today's agenda

- Language modeling using RNNs
- **Image captioning**
  - Naive image captioning, image captioning with attention
- Visual question answering
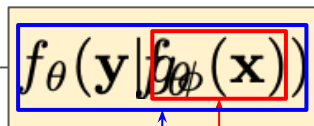  - Naive visual question answering, memory network

# Image captioning

- Task definition: describe an image using natural language (sentence)

**x**: image

**y**: sentence



$$f_\theta(\mathbf{y} \mid f_\phi(\mathbf{x}))$$

"man in black shirt is playing guitar."

<span style="color:red">1. Recognizing visual patterns in an image</span>

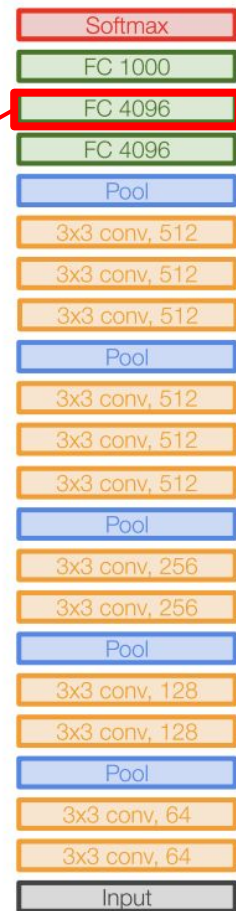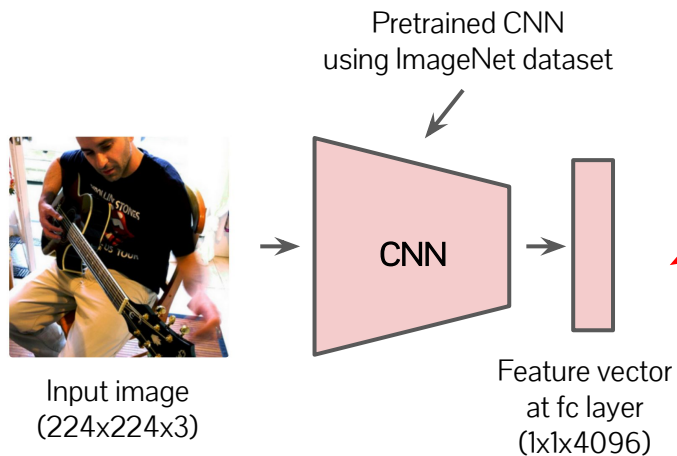<span style="color:blue">2. Generate a sentence conditioned on image</span>
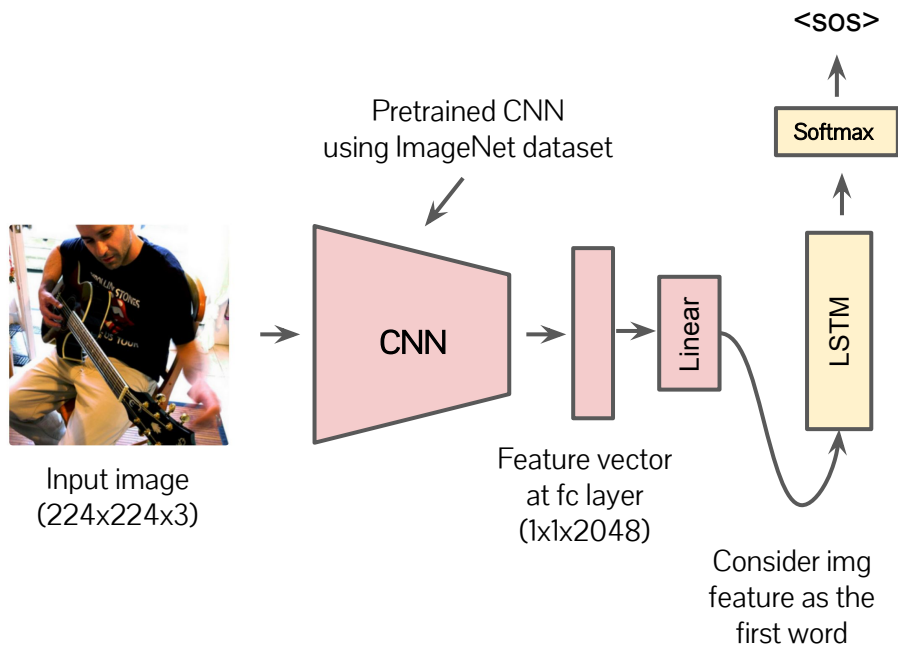
<span style="color:red">**CNN**</span>

<span style="color:blue">**RNN**</span>

18

# Naive image captioning

$$f_\theta(\mathbf{y}|\boxed{g_\phi(\mathbf{x})})$$



Pretrained CNN
using ImageNet dataset

CNN

Input image
(224x224x3)

Feature vector
at fc layer
(1x1x4096)

| Softmax |
| FC 1000 |
| FC 4096 |
| FC 4096 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| Pool |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| Pool |
| 3x3 conv, 128 |
| 3x3 conv, 128 |
| Pool |
| 3x3 conv, 64 |
| 3x3 conv, 64 |
| Input |

VGG16

# Naive image captioning

$$f_\theta(\mathbf{y}|g_\phi(\mathbf{x}))$$



<SOS>

Softmax

Pretrained CNN
using ImageNet dataset

CNN

Linear

LSTM

Input image
(224x224x3)

Feature vector
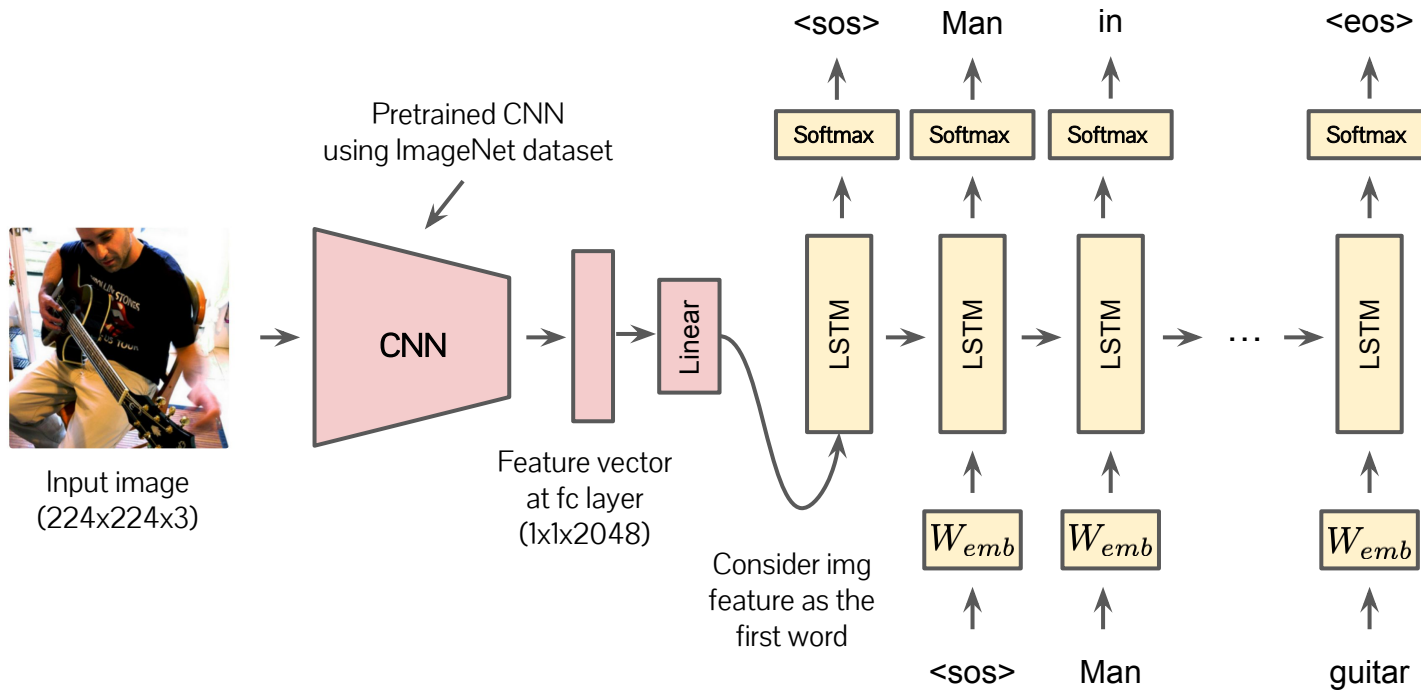at fc layer
(1x1x2048)

Consider img
feature as the
first word

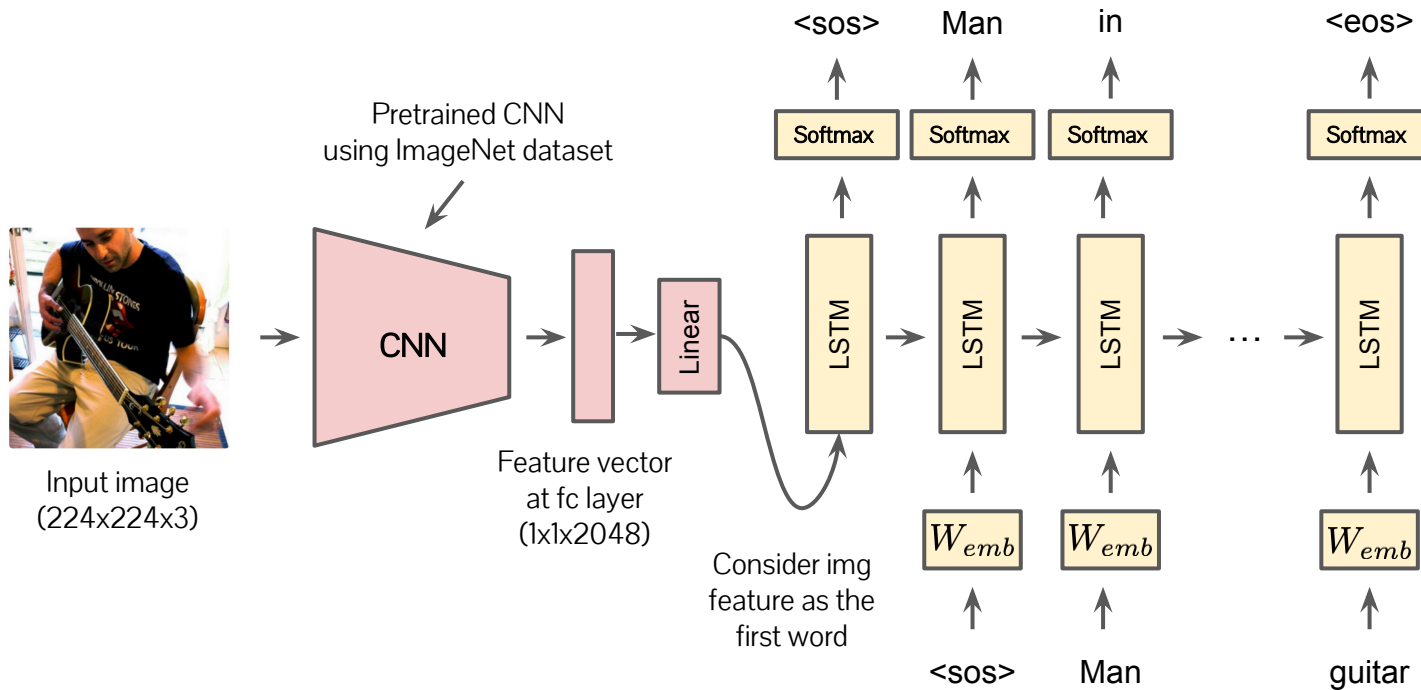# Naive image captioning

$$f_\theta(\mathbf{y}|g_\phi(\mathbf{x}))$$

# Naive image captioning

$$f_\theta(\mathbf{y}|g_\phi(\mathbf{x}))$$



Input image
(224x224x3)

Pretrained CNN
using ImageNet dataset

CNN

Feature vector
at fc layer
(1x1x2048)

Linear

Consider img
feature as the
first word

LSTM

Softmax

<sos>

$W_{emb}$

<sos>

LSTM

Softmax

Man

$W_{emb}$

Man

LSTM

Softmax

in

…

LSTM

Softmax

<eos>

$W_{emb}$

guitar

# Naive image captioning: Training

**Training data**
(image, sentence) pairs

$\mathbf{x}$      $\mathbf{y}$

A person on a beach
flying a kite.

A black and white photo of
a train on a train track.

A person skiing down a
snow covered slope.

Cross-entropy loss
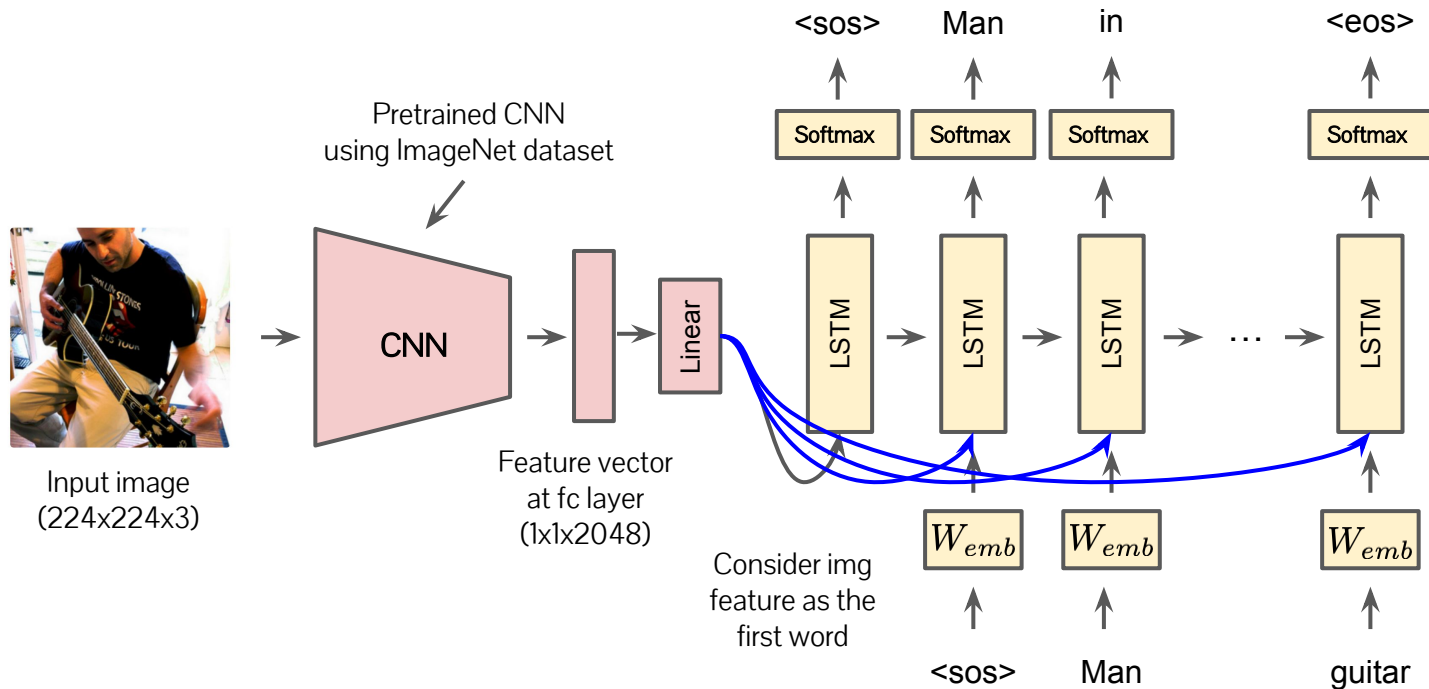(same as sentence generation)

$$\sum_{t=1}^{T} \mathcal{L}(\hat{\mathbf{y}}_t, \mathbf{y}_t)$$

<sos>    Man    in    <eos>

Softmax    Softmax    Softmax    Softmax

$\mathbf{x}$

Pretrained CNN
using ImageNet dataset

CNN

Linear

LSTM   LSTM   LSTM   ...   LSTM

Input image
(224x224x3)

Feature vector
at fc layer
(1x1x2048)

Consider img
feature as the
first word

$W_{emb}$    $W_{emb}$    $W_{emb}$

<sos>    Man    guitar

# Practical issue



Pretrained CNN using ImageNet dataset

Input image (224x224x3)

Feature vector at fc layer (1x1x2048)

Consider img feature as the first word

<sos>   A   Man   **?**

$W_{emb}$   $W_{emb}$   $W_{emb}$

<sos>   A   Man

Can you guess what would be the following word?

- If RNN is strong enough, it can generate reasonable sentences **conditioned only on previous words** (not an image)

- This is especially prominent since the previous words has more direct impact on prediction of next words (shorter dependency length)

- In order to make captioning conditioned on image content, we have to **strengthen** the conditioning to image

24

# Improving image conditioning: shortcuts

# Improving image conditioning: attention

- Make the model "gaze" on salient objects for generating corresponding words

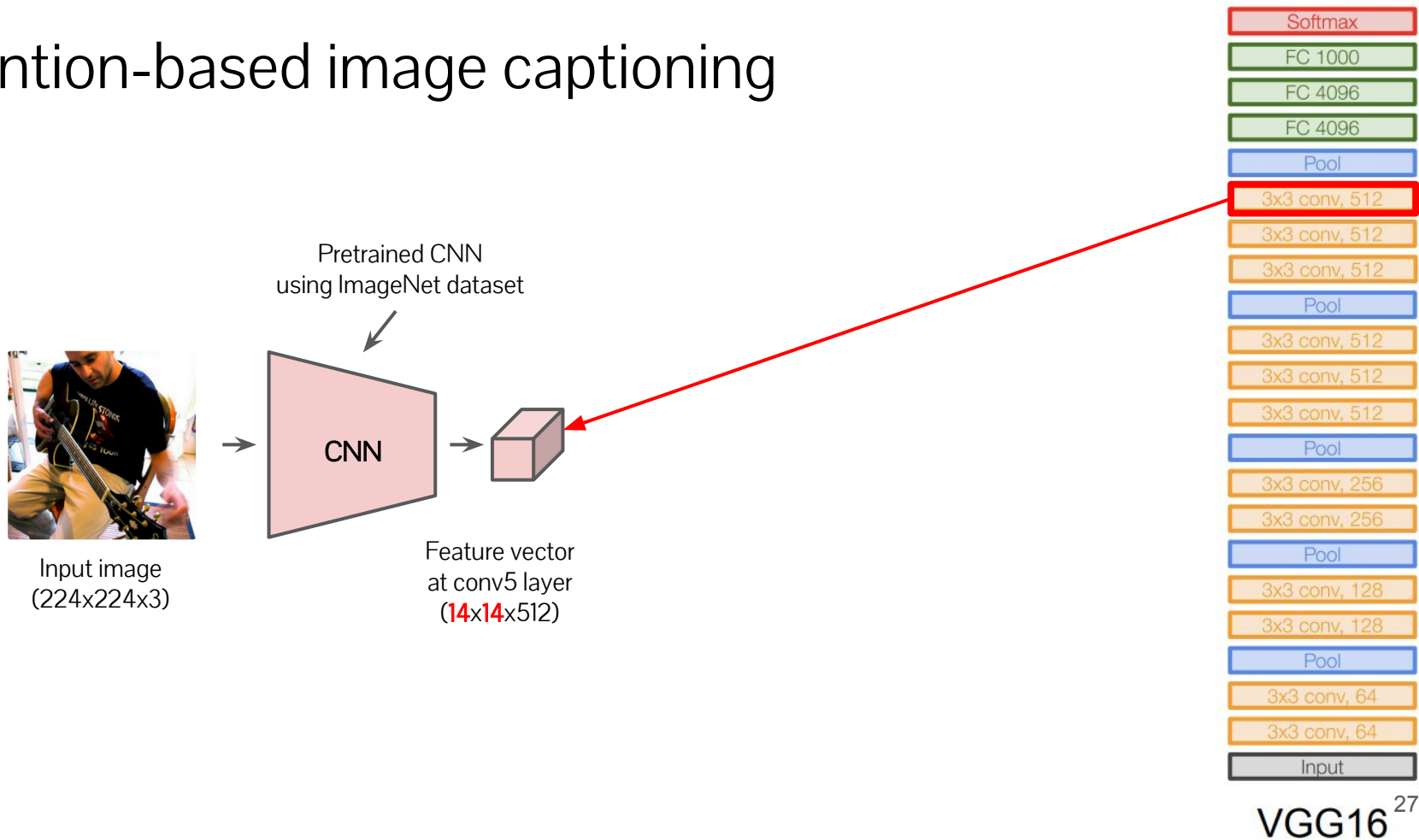**No attention**   (= uniform attention)
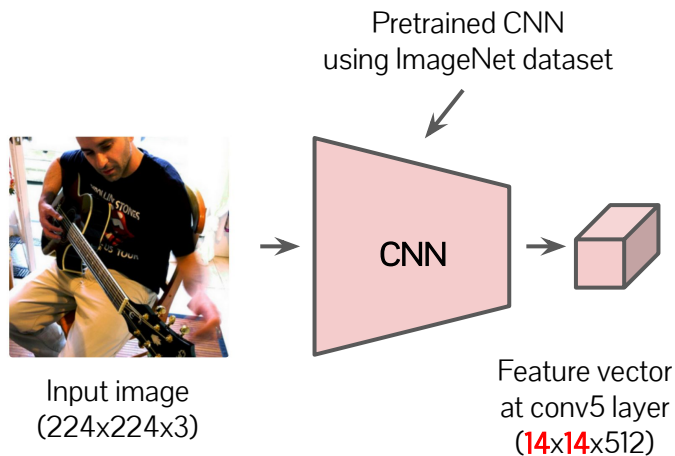
A bird flying over a body of water.

Word prediction is conditioned on parts of an image

**Soft attention**



A    bird    flying    over    a    body    of    water    .

# Attention-based image captioning



Pretrained CNN
using ImageNet dataset

CNN

Input image
(224x224x3)

Feature vector
at conv5 layer
(14x14x512)

Softmax
FC 1000
FC 4096
FC 4096
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 512
3x3 conv, 512
3x3 conv, 512
Pool
3x3 conv, 256
3x3 conv, 256
Pool
3x3 conv, 128
3x3 conv, 128
Pool
3x3 conv, 64
3x3 conv, 64
Input

VGG16 <sup>27</sup>

# Attention-based image captioning

Pretrained CNN
using ImageNet dataset

CNN

Input image
(224x224x3)

Feature vector
at conv5 layer
(14x14x512)

**Attention**:
- A positive matrix that has same spatial dimension as feature map

$$\alpha^{(t)} \in \mathbb{R}^{W \times H} \qquad \sum_{i,j} \alpha_{i,j}^{(t)} = 1$$

- We want to compute attention for each word

a   bird   flying     .

14   $\alpha^{(1)}$   $\alpha^{(2)}$   $\alpha^{(3)}$     $\alpha^{(N)}$

14

- Attention is used to abstract image feature

$$\mathbf{z}^{(t)} = \sum_{i,j} \alpha_{i,j}^{(t)} \mathbf{x}_{i,j} \in \mathbb{R}^C$$
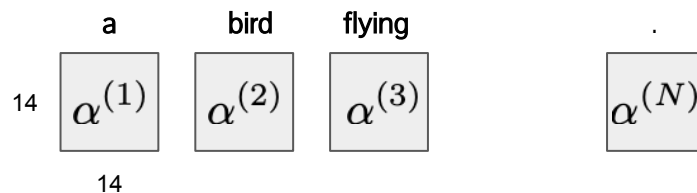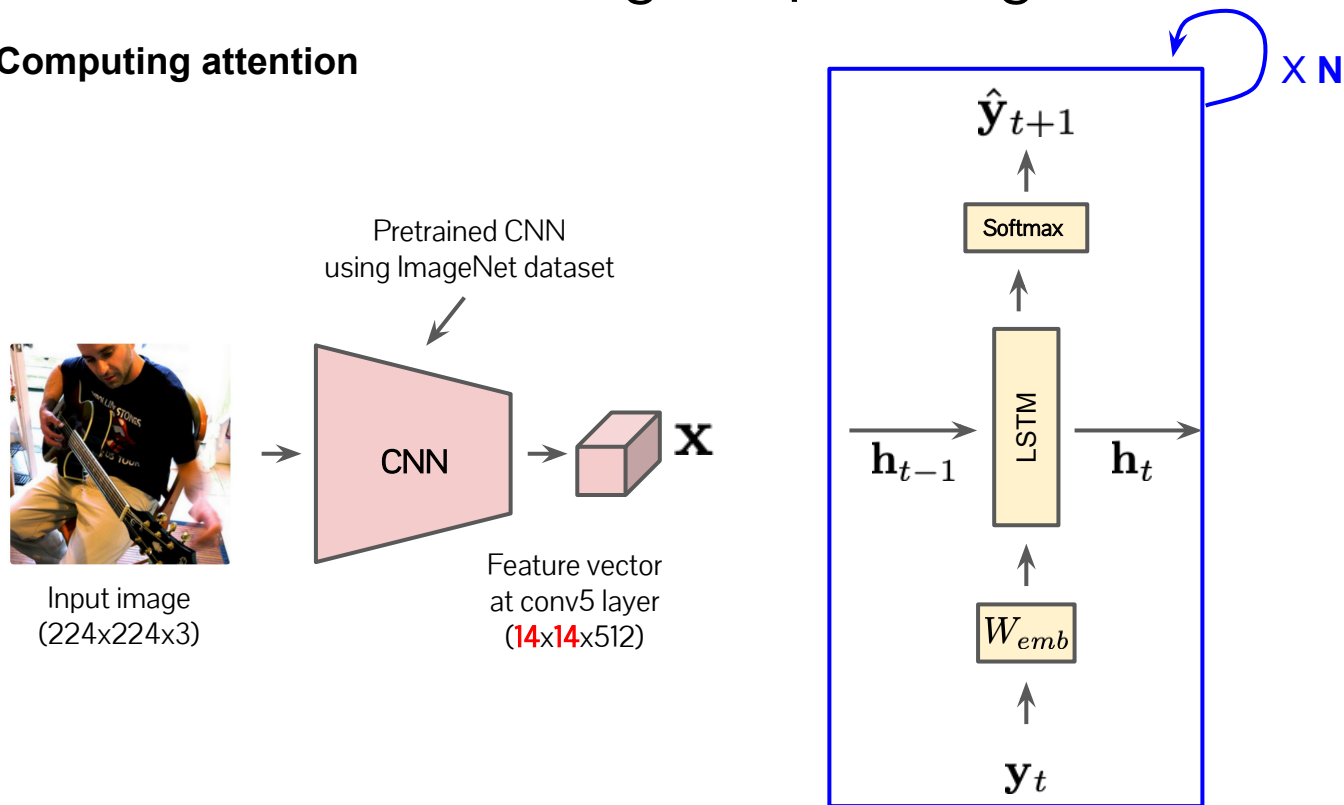
# Attention-based image captioning

**Challenges:**

**Attention:**

- A positive matrix that has same spatial dimension as feature map

$$\alpha^{(t)} \in \mathbb{R}^{W \times H} \qquad \sum_{i,j} \alpha_{i,j}^{(t)} = 1$$

- How do we compute the attention? $\longrightarrow$
- We want to compute attention for each word

a      bird     flying         .

14   $\alpha^{(1)}$   $\alpha^{(2)}$   $\alpha^{(3)}$       $\alpha^{(N)}$

14

- How do we use it to predict the word? $\longrightarrow$
- Attention is used to abstract image feature

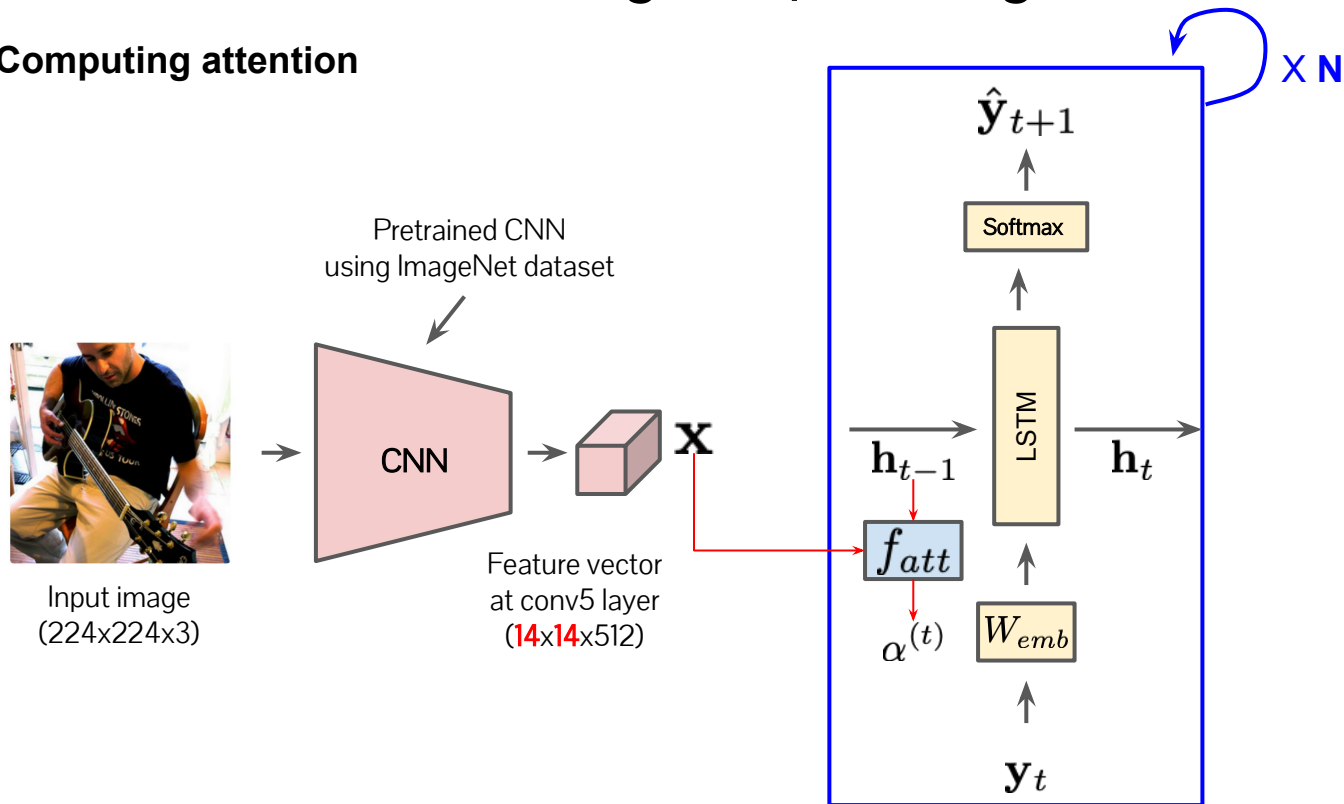$$\mathbf{z}^{(t)} = \sum_{i,j} \alpha_{i,j}^{(t)} \mathbf{x}_{i,j} \in \mathbb{R}^{C}$$

29

# Attention-based image captioning

**Challenges:**

**Attention:**

- A positive matrix that has same spatial dimension as feature map

$$\alpha^{(t)} \in \mathbb{R}^{W \times H} \qquad \sum_{i,j} \alpha^{(t)}_{i,j} = 1$$

- **How do we compute the attention?** ⟶

- We want to compute attention for each word

| a | bird | flying | . |

14  $\alpha^{(1)}$  $\alpha^{(2)}$  $\alpha^{(3)}$  $\alpha^{(N)}$

14

- How do we use it to predict the word? ⟶

- Attention is used to abstract image feature

$$\mathbf{z}^{(t)} = \sum_{i,j} \alpha^{(t)}_{i,j} \mathbf{x}_{i,j} \in \mathbb{R}^C$$

30

# Attention-based image captioning

**Computing attention**

Pretrained CNN
using ImageNet dataset

CNN

$\mathbf{X}$

Input image
(224x224x3)

Feature vector
at conv5 layer
(14x14x512)

$\hat{\mathbf{y}}_{t+1}$

Softmax

LSTM

$\mathbf{h}_{t-1}$ → LSTM → $\mathbf{h}_t$

$W_{emb}$

$\mathbf{y}_t$

X **N**

# Attention-based image captioning

**Computing attention**



Pretrained CNN
using ImageNet dataset

CNN

$\mathbf{X}$

Input image
(224x224x3)

Feature vector
at conv5 layer
(14x14x512)

$\hat{\mathbf{y}}_{t+1}$

Softmax

$\mathbf{h}_{t-1}$

LSTM

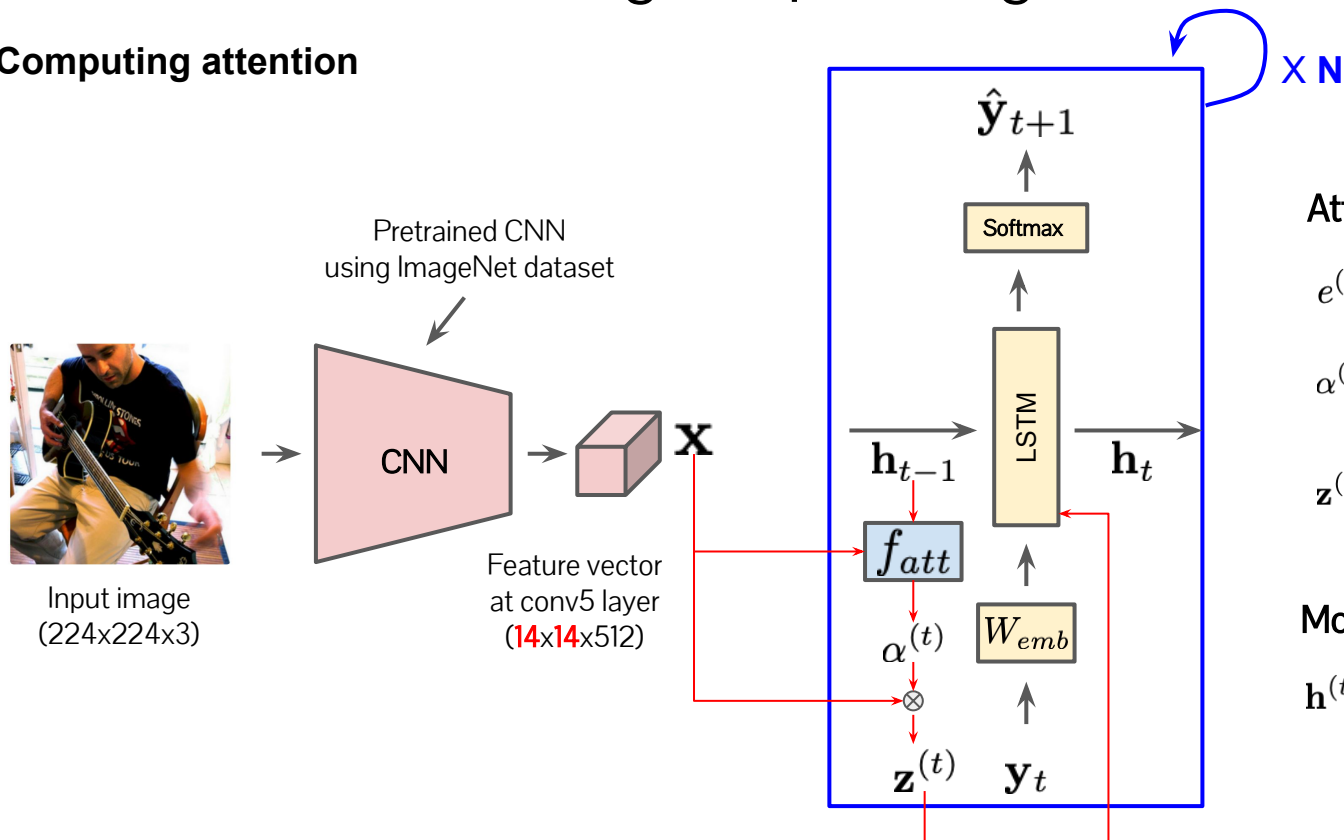$\mathbf{h}_t$

$f_{att}$

$\alpha^{(t)}$

$W_{emb}$

$\mathbf{y}_t$

X **N**

Attention module

$$e^{(t)} = f_{att}(\mathbf{x}, \mathbf{h}_{t-1})$$

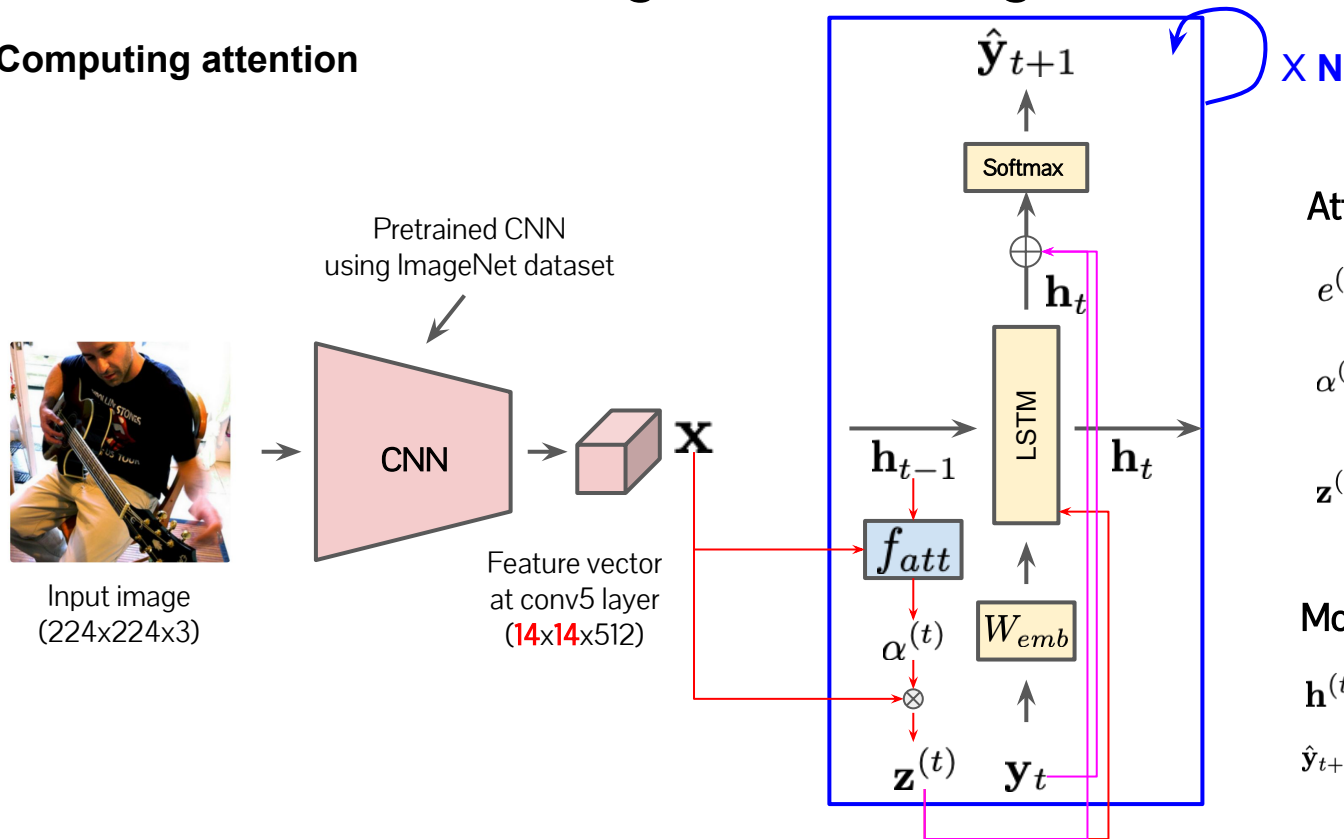$$\alpha^{(t)} = \frac{\exp(e^{(t)})}{\sum_{i,j} \exp(e_{i,j}^{(t)})}$$

32

# Attention-based image captioning

**Computing attention**



Pretrained CNN
using ImageNet dataset

Input image
(224x224x3)

CNN

Feature vector
at conv5 layer
(14x14x512)

$\mathbf{X}$

$\hat{\mathbf{y}}_{t+1}$

Softmax

LSTM

$\mathbf{h}_{t-1}$    $\mathbf{h}_t$

$f_{att}$

$\alpha^{(t)}$

$\otimes$

$W_{emb}$

$\mathbf{z}^{(t)}$    $\mathbf{y}_t$

X **N**

Attention module

$$e^{(t)} = f_{att}(\mathbf{x}, \mathbf{h}_{t-1})$$

$$\alpha^{(t)} = \frac{\exp(e^{(t)})}{\sum_{i,j} \exp(e_{i,j}^{(t)})}$$

$$\mathbf{z}^{(t)} = \sum_{i,j} \alpha_{i,j}^{(t)} \mathbf{x}_{i,j}$$

# Attention-based image captioning

**Computing attention**



Pretrained CNN
using ImageNet dataset

CNN

Input image
(224x224x3)

Feature vector
at conv5 layer
(14x14x512)

$\mathbf{X}$

$\hat{\mathbf{y}}_{t+1}$

Softmax

LSTM

$\mathbf{h}_{t-1}$

$\mathbf{h}_t$

$f_{att}$

$\alpha^{(t)}$

$W_{emb}$

$\otimes$

$\mathbf{z}^{(t)}$

$\mathbf{y}_t$

X **N**

## Attention module

$$e^{(t)} = f_{att}(\mathbf{x}, \mathbf{h}_{t-1})$$

$$\alpha^{(t)} = \frac{\exp(e^{(t)})}{\sum_{i,j} \exp(e_{i,j}^{(t)})}$$

$$\mathbf{z}^{(t)} = \sum_{i,j} \alpha_{i,j}^{(t)} \mathbf{x}_{i,j}$$

## Modified LSTM

$$\mathbf{h}^{(t)} = LSTM(\mathbf{h}_{t-1}, W_{emb}\mathbf{y}_t, \mathbf{z}^{(t)})$$

34

# Attention-based image captioning

**Computing attention**



Pretrained CNN
using ImageNet dataset

Input image
(224x224x3)

CNN

Feature vector
at conv5 layer
(14x14x512)

$\mathbf{X}$

Softmax

$\hat{\mathbf{y}}_{t+1}$

X **N**

$\mathbf{h}_t$

LSTM

$\mathbf{h}_{t-1}$

$\mathbf{h}_t$

$f_{att}$

$\alpha^{(t)}$

$W_{emb}$

$\mathbf{z}^{(t)}$

$\mathbf{y}_t$

**Attention module**

$$e^{(t)} = f_{att}(\mathbf{x}, \mathbf{h}_{t-1})$$

$$\alpha^{(t)} = \frac{\exp(e^{(t)})}{\sum_{i,j} \exp(e^{(t)}_{i,j})}$$

$$\mathbf{z}^{(t)} = \sum_{i,j} \alpha^{(t)}_{i,j} \mathbf{x}_{i,j}$$

**Modified LSTM**

$$\mathbf{h}^{(t)} = LSTM(\mathbf{h}_{t-1}, W_{emb}\mathbf{y}_t, \mathbf{z}^{(t)})$$

$$\hat{\mathbf{y}}_{t+1} = \exp(W^o(W_{emb}\mathbf{y}_t + W^h\mathbf{h}_t + W^z\mathbf{z}_t))$$

35

A(0.99)  large(0.49)  white(0.40)

bird(0.35)  standing(0.29)  in(0.27)  a(0.35)

forest(0.54)  .(0.46)

A(0.96)  group(0.27)  of(0.27)
people(0.21)  sitting(0.28)  on(0.22)  a(0.21)
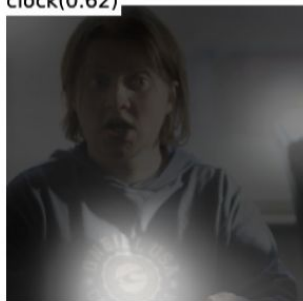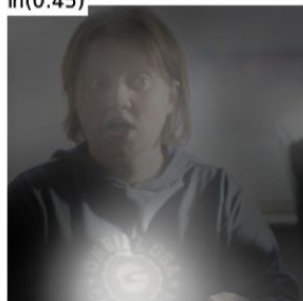boat(0.19)  in(0.13)  the(0.10)  water(0.30)
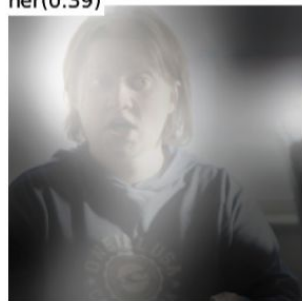
37

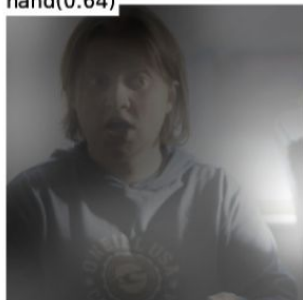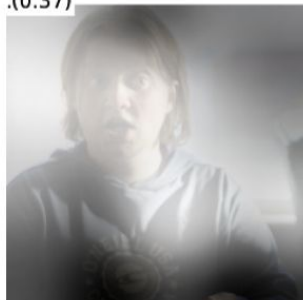A(1.00)  woman(0.80)  holding(0.68)  a(0.58)  clock(0.62)  in(0.45)  her(0.39)  hand(0.64)  .(0.37)
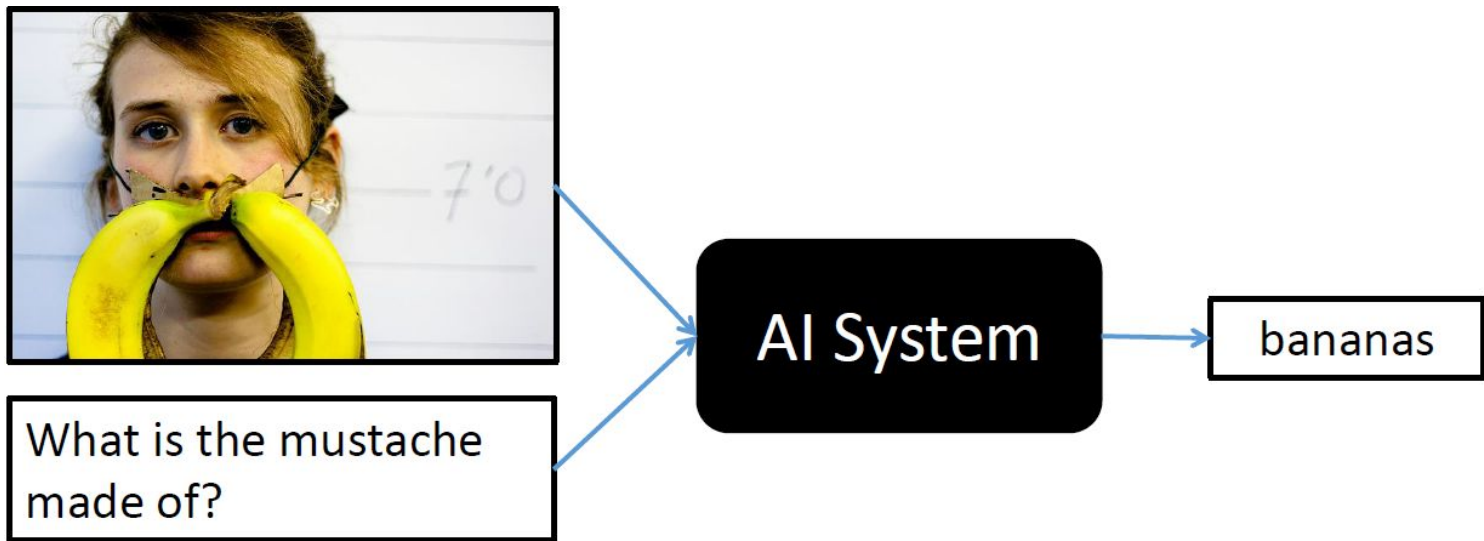
# Today's agenda

- Language modeling using RNNs
- Image captioning
    - Naive image captioning, image captioning with attention
- **Visual question answering**
    - Naive visual question answering, memory network

# Visual question answering

- Objective: given an image and a question about an image, predict an answer.

# Visual question answering

- Objective: given an image and a question about an image, predict an answer.



How do we design this system?
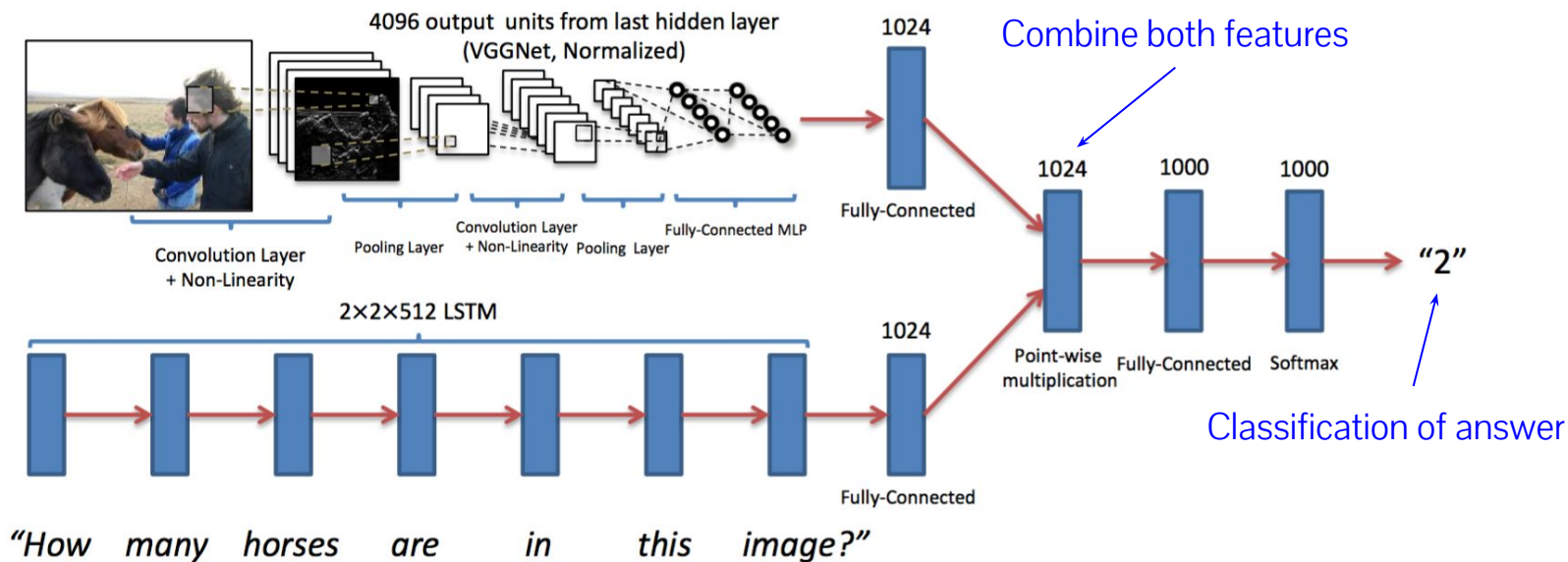
# Visual question answering



Encoding image using CNN

4096 output units from last hidden layer (VGGNet, Normalized)

1024
Fully-Connected

Convolution Layer + Non-Linearity
Pooling Layer
Convolution Layer + Non-Linearity
Pooling Layer
Fully-Connected MLP

Combine both features

1024    1000    1000

"2"

Point-wise multiplication    Fully-Connected    Softmax

Classification of answer

2×2×512 LSTM

1024
Fully-Connected

"How    many    horses    are    in    this    image?"

Encoding question using LSTM

42

# VQA with attention



feature vectors of different parts of image

**Question:**
What are sitting in the basket on a bicycle?

**Answer:**
dogs

Query

Attention layer 1

Attention layer 2

Yang et al., Stacked Attention Networks for Image Question Answering, In CVPR, 2016

# VQA with attention

Question: What are sitting in the basket on a bicycle?



**Original Image**  **First Attention Layer**  **Second Attention Layer**