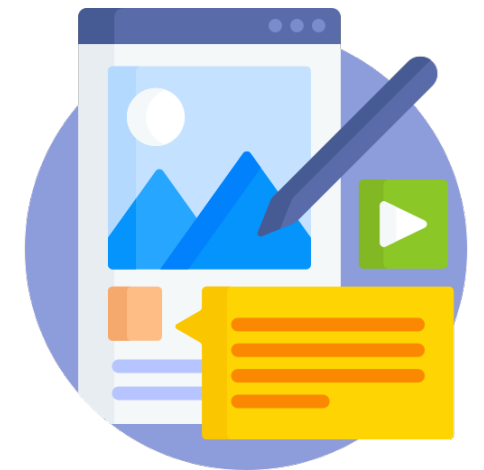


Image Captioning

Models Introduction

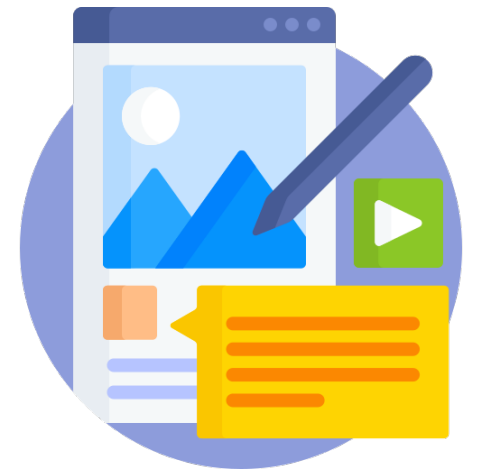
Date: Oct 27, 2021

Presenter: Everley Tseng



Outline

- Show Attend and Tell
 - Introduction to attention
 - NLP backgrounds
 - Model structure
 - Soft attention & hard attention
 - Loss function
 - Experimental results
- What's new in Image Captioning?



Why Attention?

Please write the caption for this image



Image source: <https://www.sandiegouniontribune.com/business/real-estate/story/2020-08-21/will-san-diego-stay-at-home-workers-leave-if-given-the-opportunity>

Why Attention?

- Important Components
 - Objects
 - Environments
- Relationships
 - Spatial
 - Interactive
- Details
 - Adjectives
 - Adverbs



Why Attention?

- Important Components
 - Objects
 - Environments
- Relationships
 - Spatial
 - Interactive
- Details
 - Adjectives
 - Adverbs



Why Attention?

- Important Components
 - Objects **Dog Dog**
 - Environments
- Relationships
 - Spatial
 - Interactive
- Details
 - Adjectives
 - Adverbs



Why Attention?

- Important Components

- Objects **Dog Dog**
- Environments

- Relationships

- Spatial
- Interactive

- Details

- Adjectives
- Adverbs

Brown
Colorful Collar
Happy



Why Attention?

- Important Components

- Objects **Dog Dog**
- Environments

- Relationships **Next to...**

- Spatial
- Interactive

Looking at...
Walking toward...

- Details

- Adjectives
- Adverbs

Brown
Colorful Collar
Happy



Why Attention?

- Important Components

- Objects **Dog Dog**
- Environments

- Relationships **Next to...
Looking at...
Walking toward...**

- Spatial
- Interactive

- Details

- Adjectives **Brown**
- Adverbs **Colorful Collar**
- Adjectives **Happy**



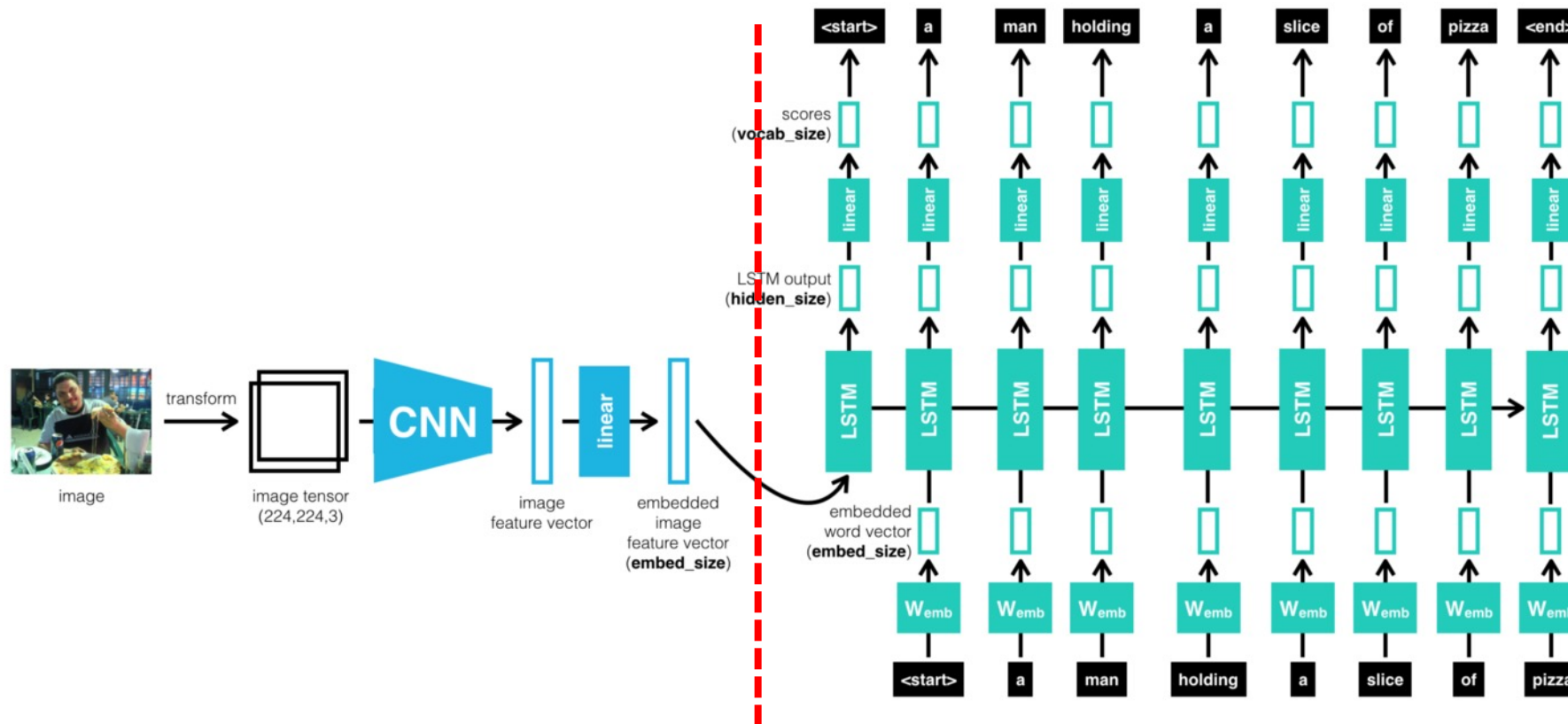
As we write down words, our attention moves across the image.

Show and Tell

O. Vinyals, A. Toshev, S. Bengio and D. Erhan, "Show and tell: A neural image caption generator," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015

Encoder

Decoder

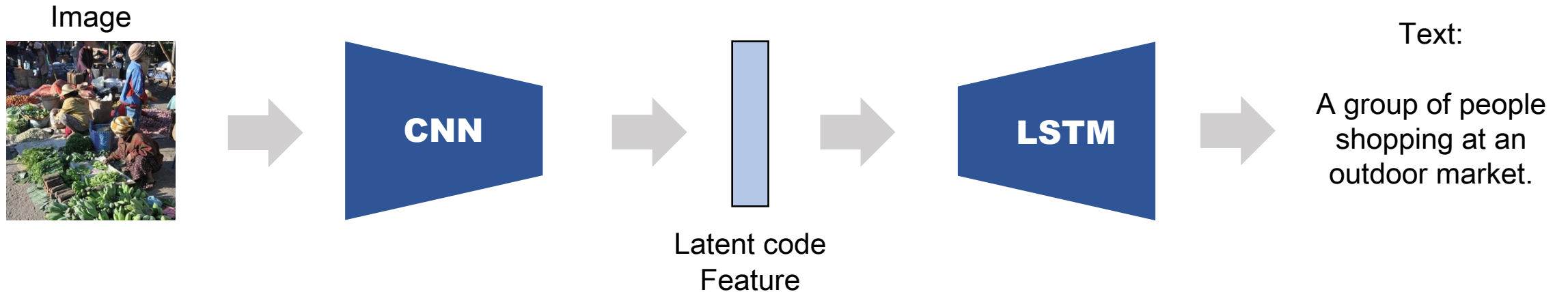


Show and Tell

O. Vinyals, A. Toshev, S. Bengio and D. Erhan, "Show and tell: A neural image caption generator," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015

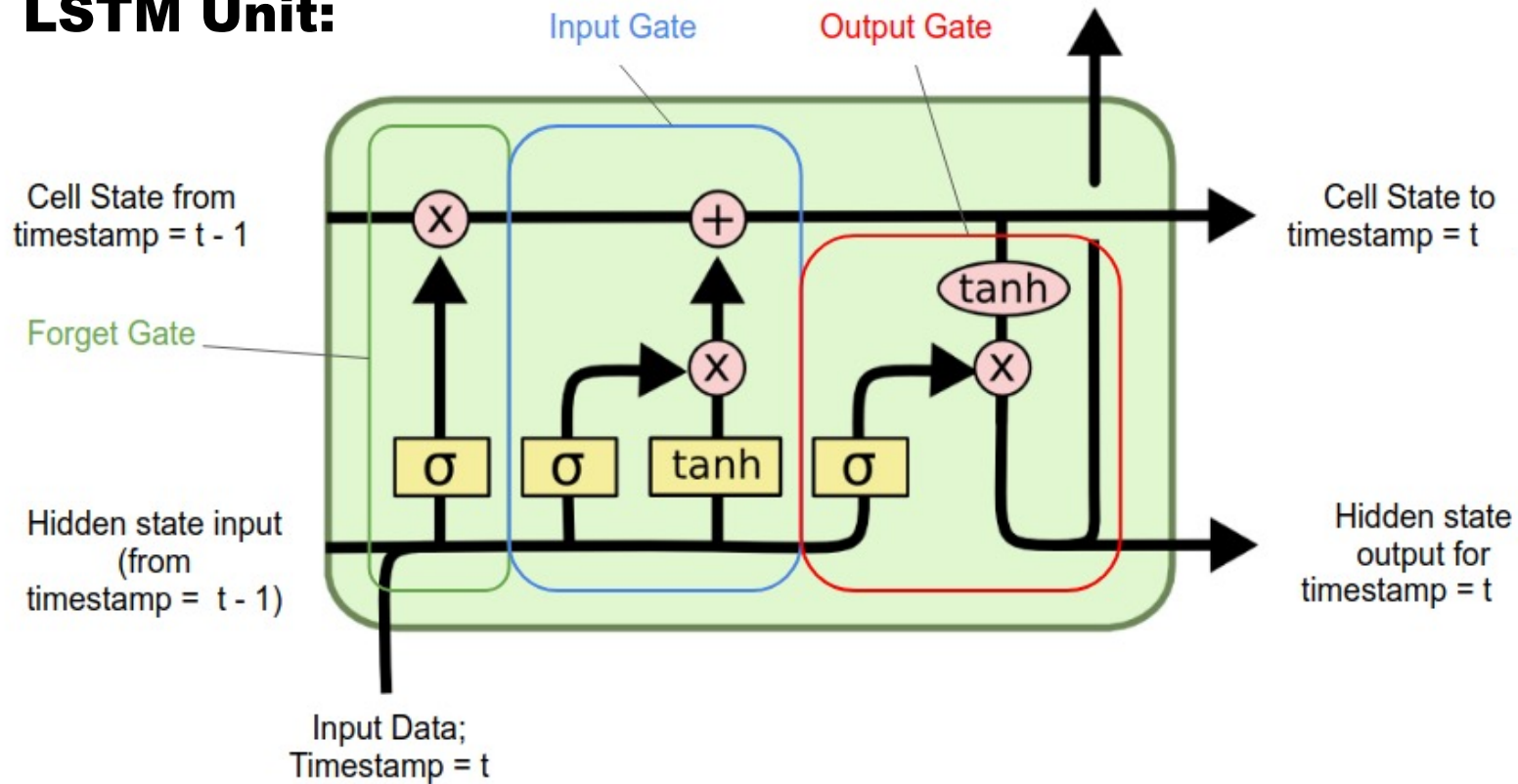
Encoder

Decoder



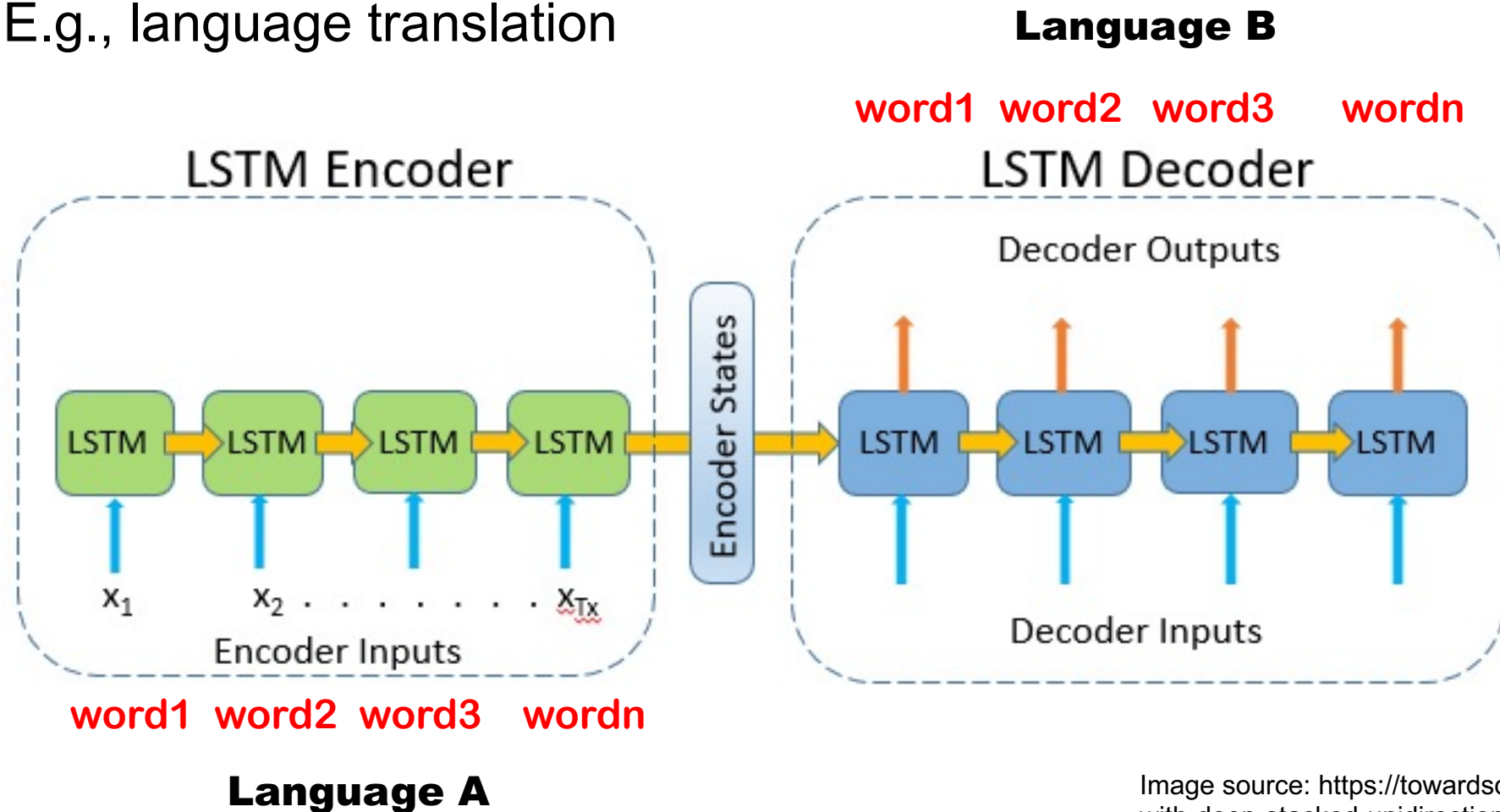
Recall – LSTM

LSTM Unit:



LSTM – Encoder-Decoder in NLP

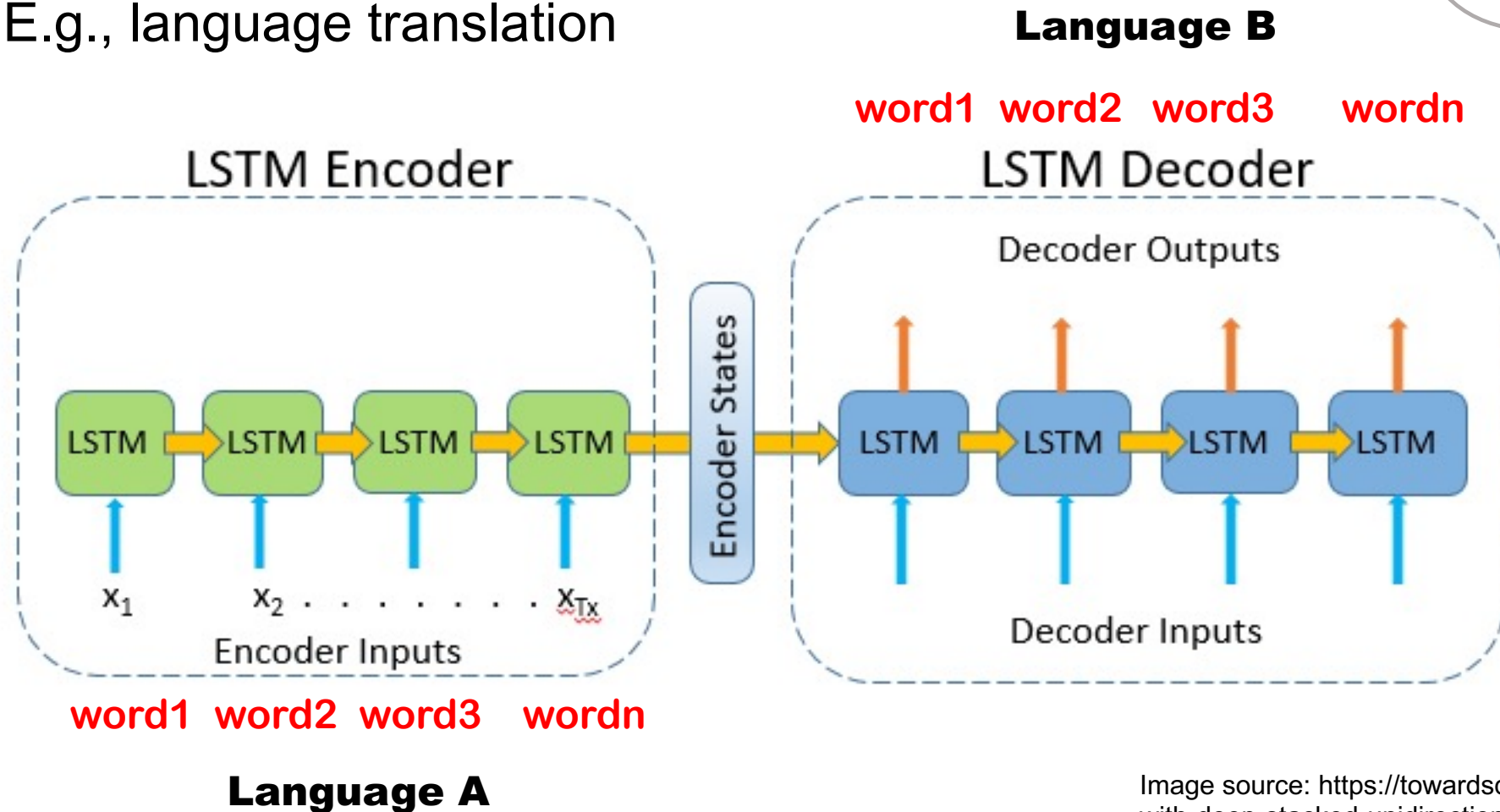
E.g., language translation



LSTM – Encoder-Decoder in NLP

How do we actually input/output words?

E.g., language translation



NLP – Word Representation

- One-hot encoding
- Word embedding

Rome = [1, 0, 0, 0, 0, 0, ..., 0]

Paris = [0, 1, 0, 0, 0, 0, ..., 0]

Italy = [0, 0, 1, 0, 0, 0, ..., 0]

France = [0, 0, 0, 1, 0, 0, ..., 0]

NLP – Word Representation

- One-hot encoding
- Word embedding

My Corpus:

1. I live in Rome
2. I live in Paris
3. I live in Italy
4. I live in France

Rome = [1, 0, 0, 0, 0, 0, ..., 0]

Paris = [0, 1, 0, 0, 0, 0, ..., 0]

Italy = [0, 0, 1, 0, 0, 0, ..., 0]

France = [0, 0, 0, 1, 0, 0, ..., 0]

What is the length of each word vector?

NLP – Word Representation

- One-hot encoding
- Word embedding

My Corpus:

1. I live in Rome
2. I live in Paris
3. I live in Italy
4. I live in France

Rome = [1, 0, 0, 0, 0, 0, ..., 0]

Paris = [0, 1, 0, 0, 0, 0, ..., 0]

Italy = [0, 0, 1, 0, 0, 0, ..., 0]

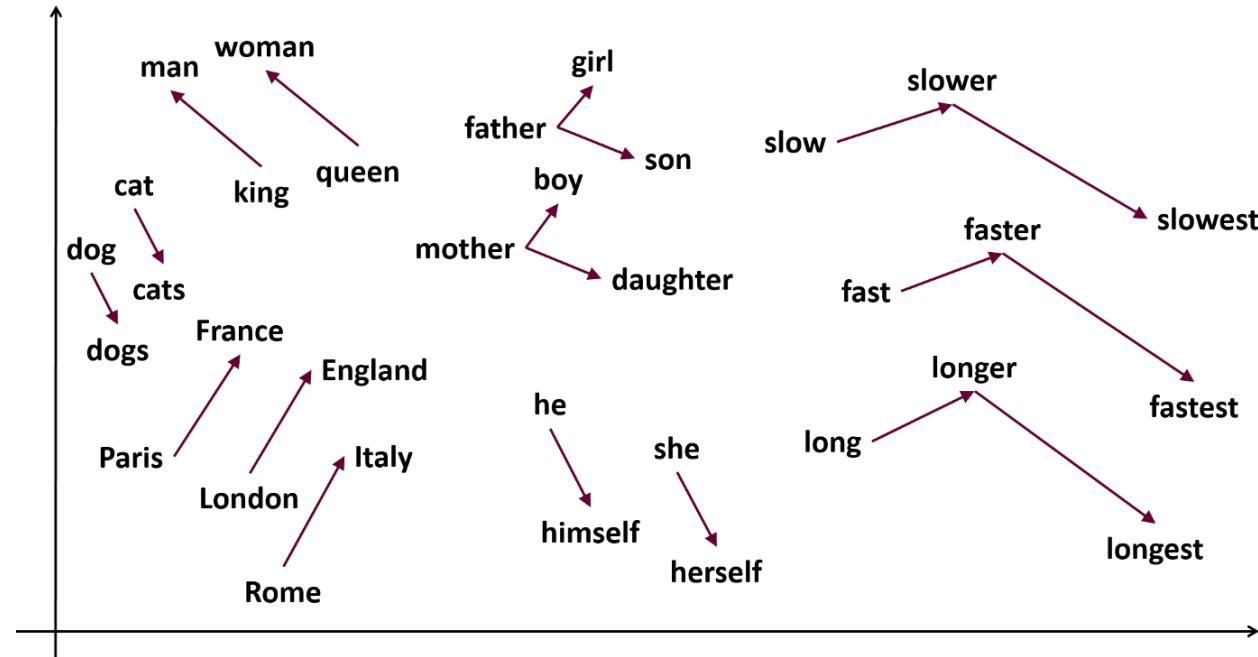
France = [0, 0, 0, 1, 0, 0, ..., 0]

What is the length of each word vector? **Length = 7**

Dictionary: ['I', 'live', 'in', 'Rome', 'Paris', 'Italy', 'France']

NLP – Word Representation

- One-hot encoding
- Word embedding
 - word2vec
 - BERT
 - GLoVe



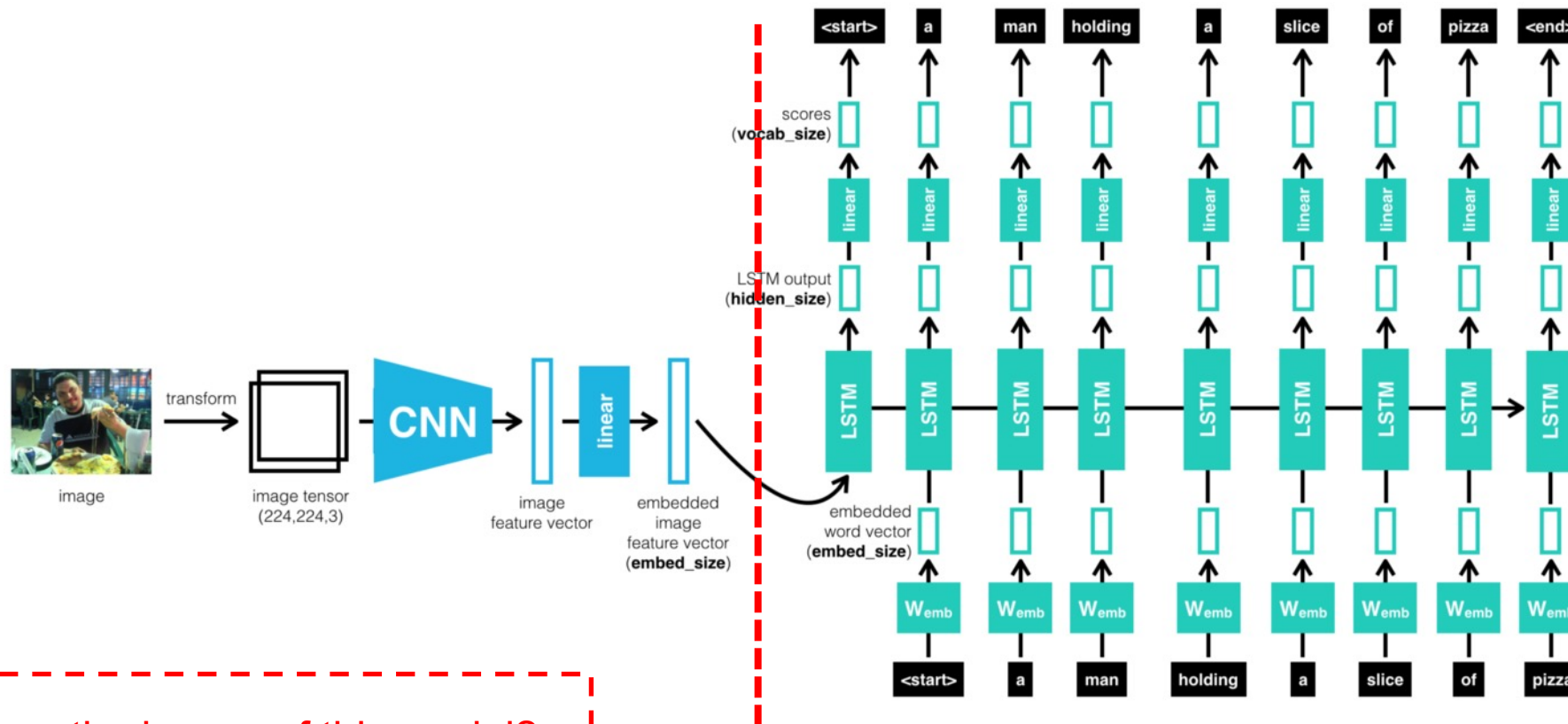
The semantic meaning of the words are embedded in the latent space.

Show and Tell

O. Vinyals, A. Toshev, S. Bengio and D. Erhan, "Show and tell: A neural image caption generator," *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015

Encoder

Decoder



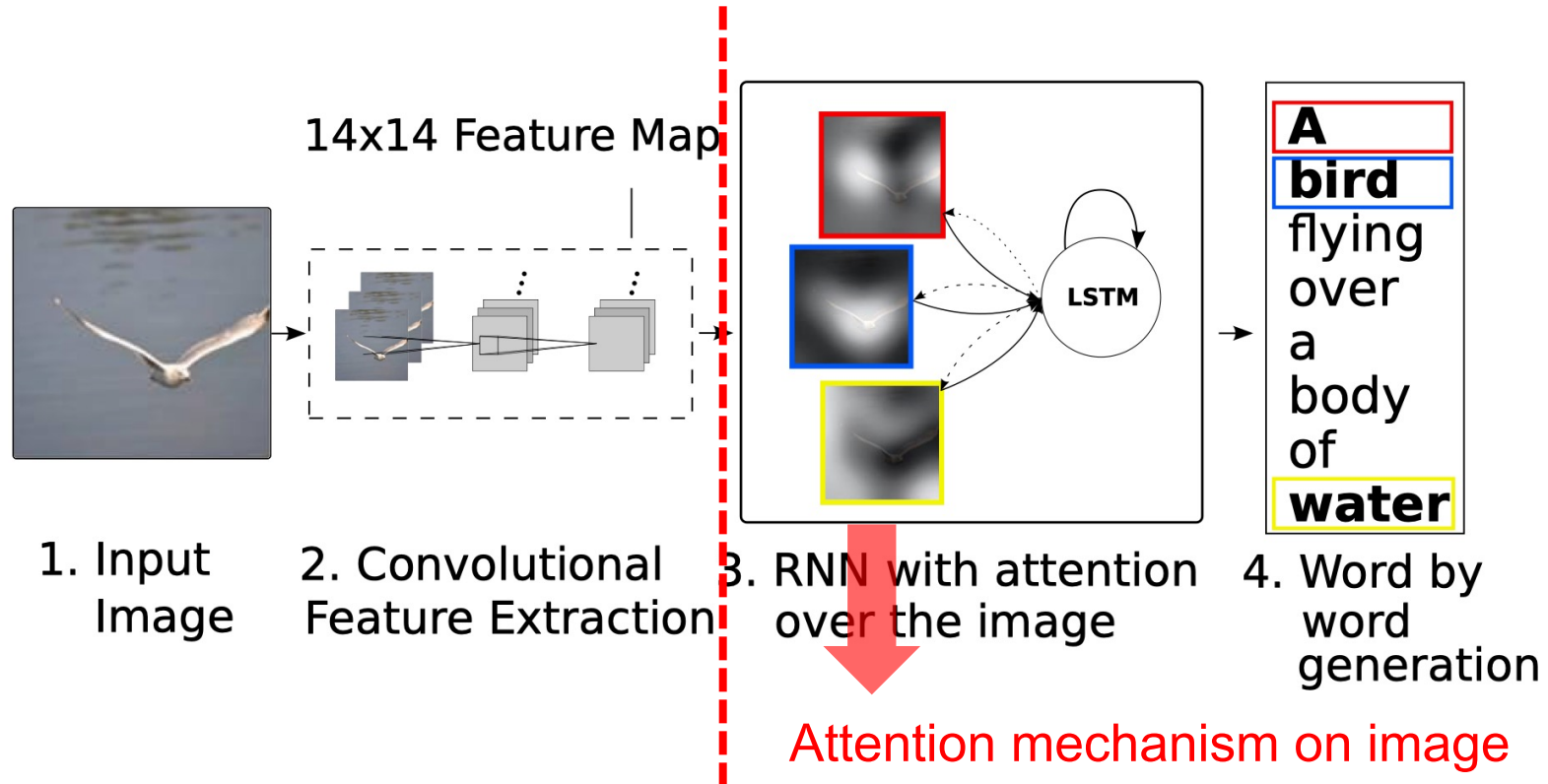
What are the issues of this model?

Image source:
<https://medium.com/swlh/image-caption-generation-with-visual-attention-c782dfc0634b>

Show, **Attend** and Tell

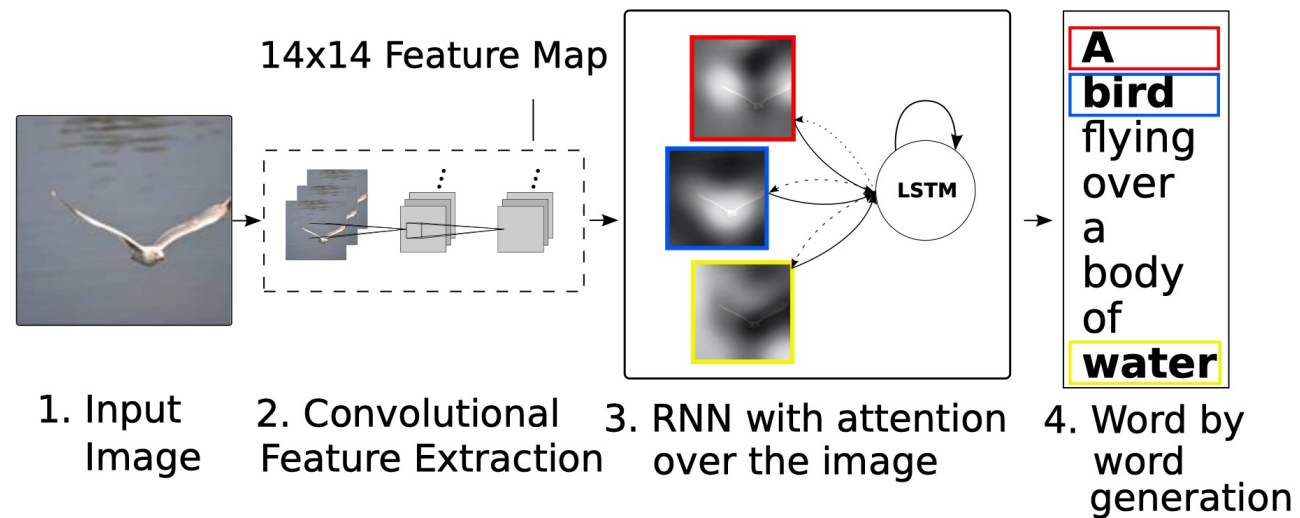
Encoder

Decoder

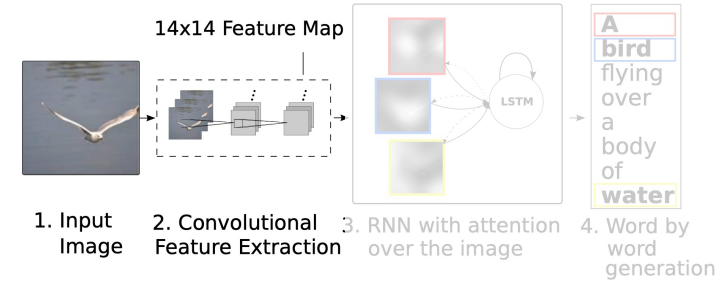


Show Attend and Tell

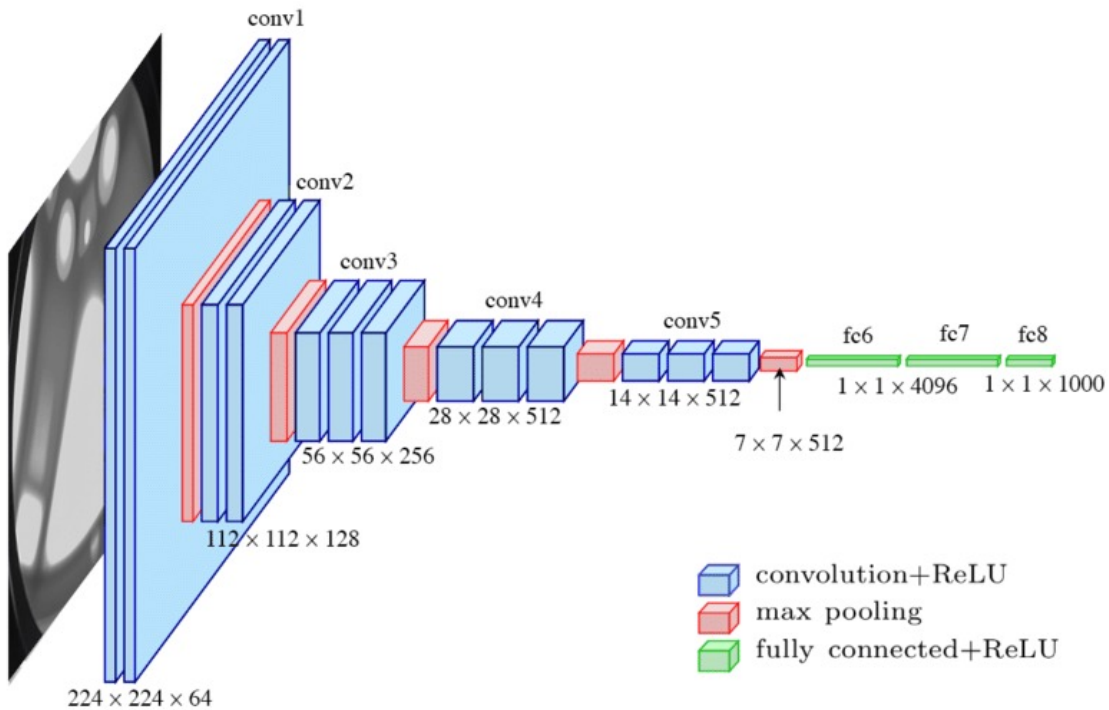
Model Structure



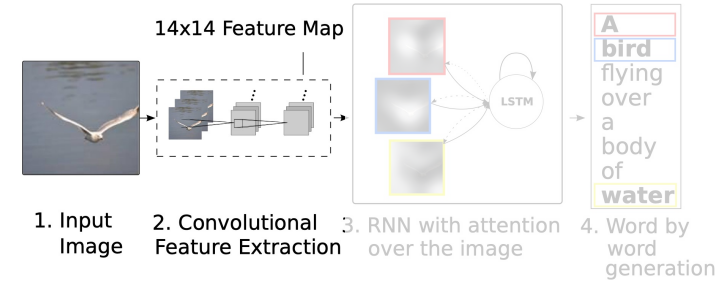
CNN Encoder



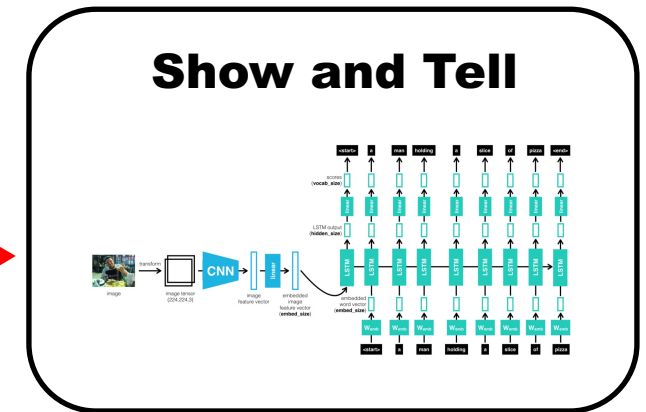
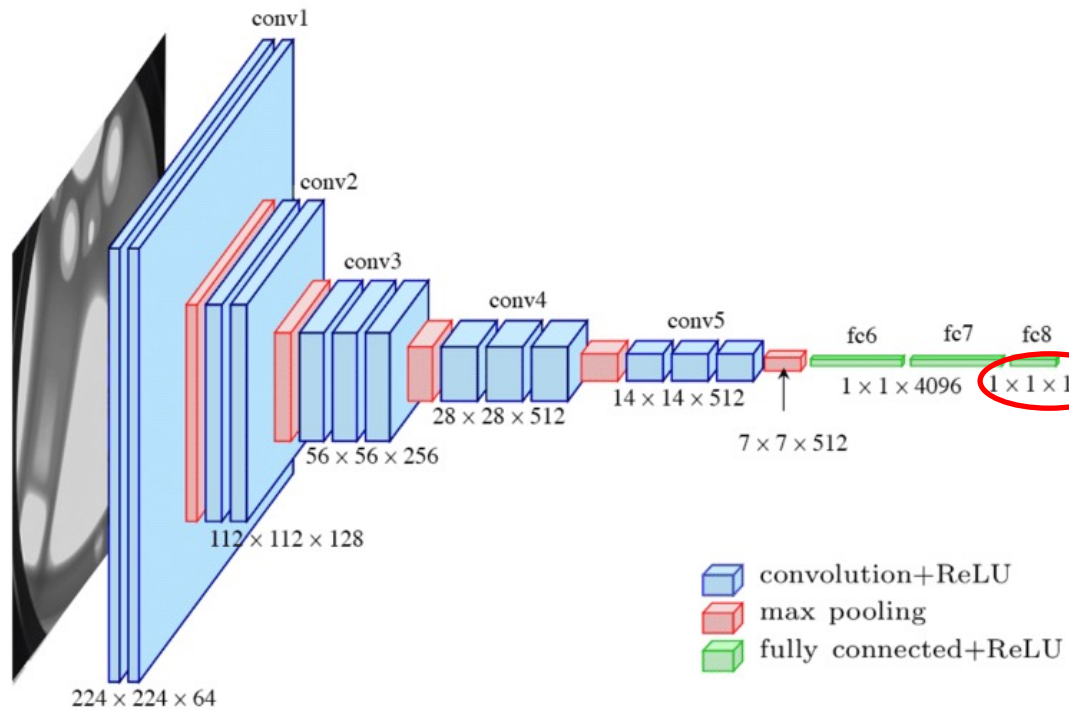
VGG-16



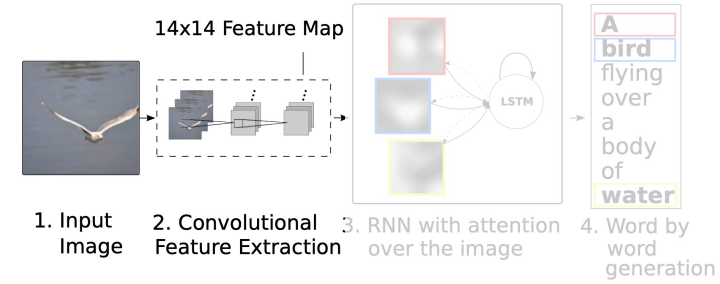
CNN Encoder



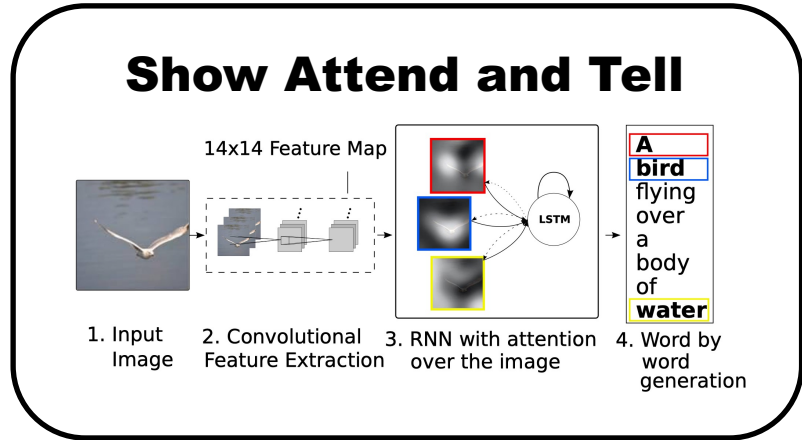
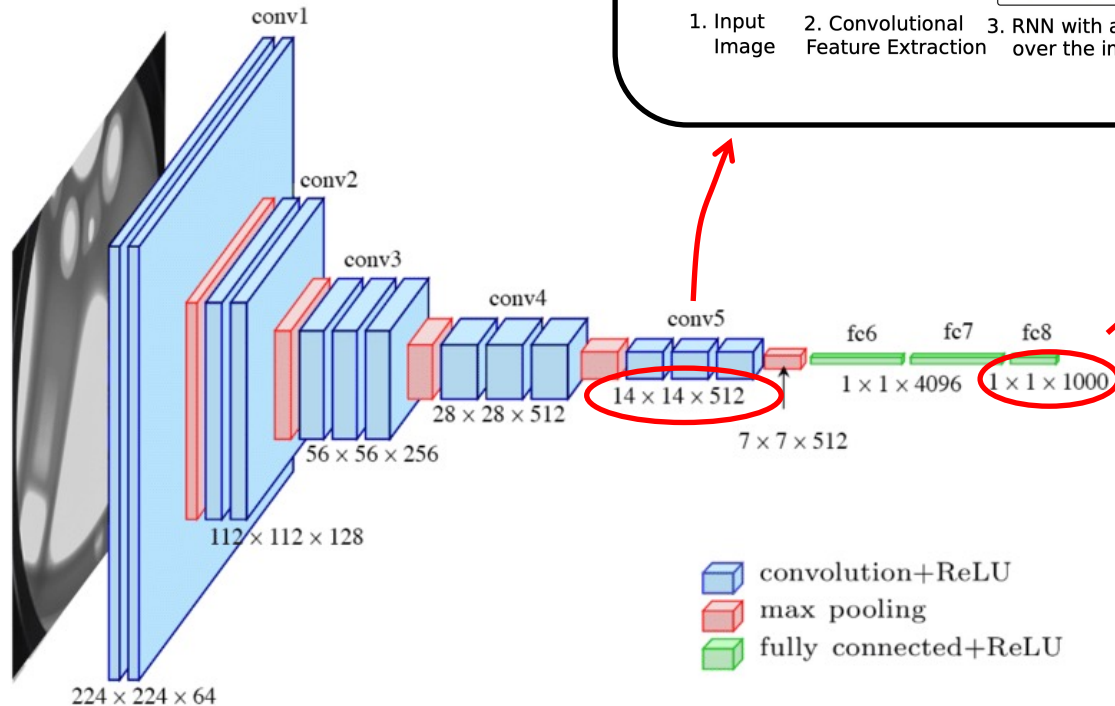
VGG-16



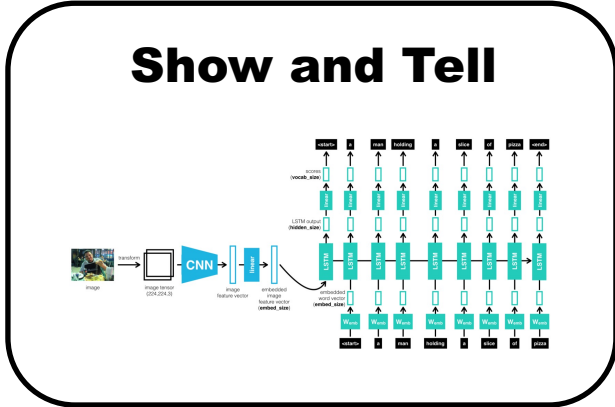
CNN Encoder



VGG-16

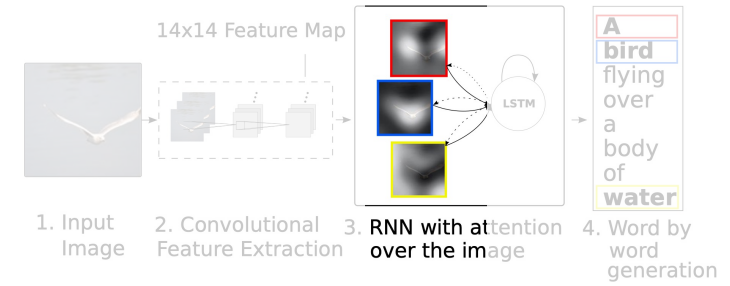


Spatial information is preserved for spatial attention



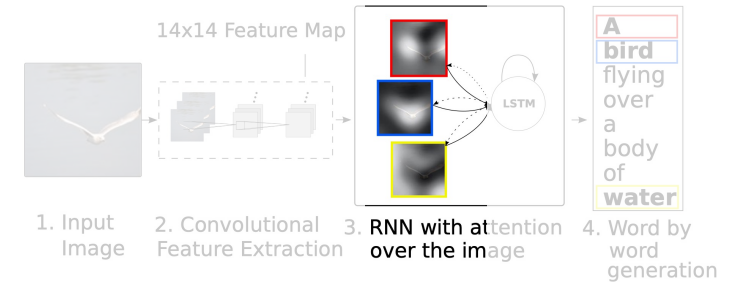
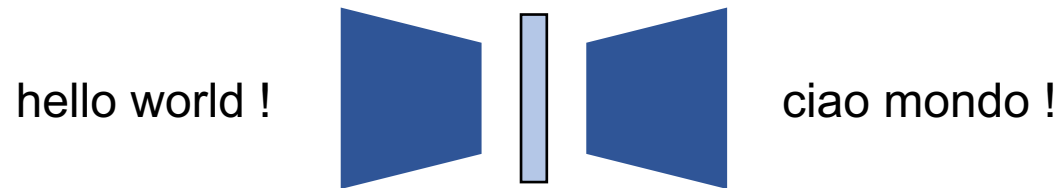
Attention

- Soft attention
- Hard attention



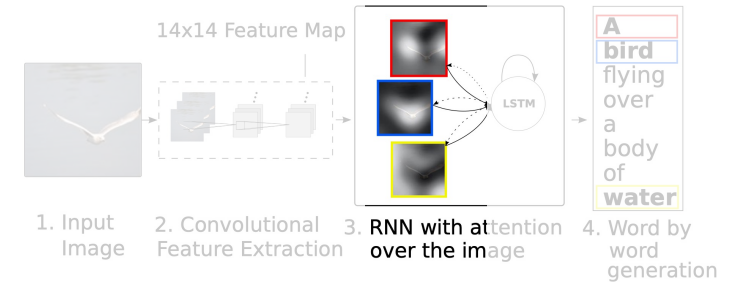
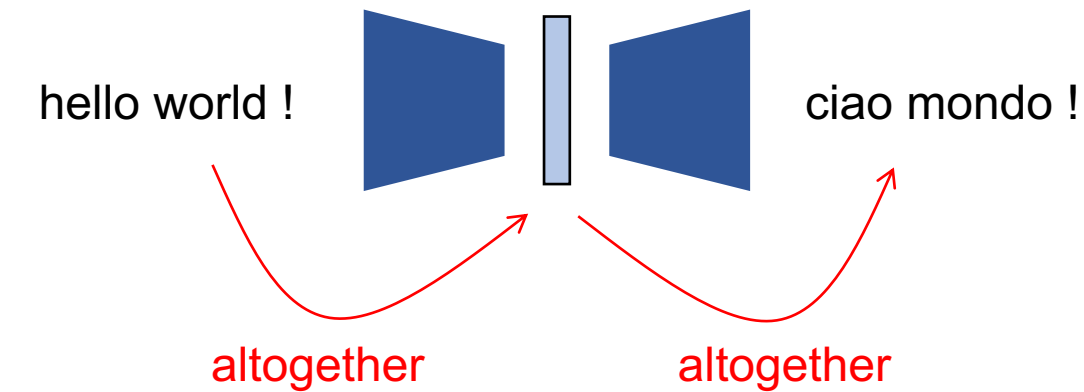
Soft Attention

- Attention Is All You Need, 2017
 - Natural Language Processing (NLP)
 - Language translation



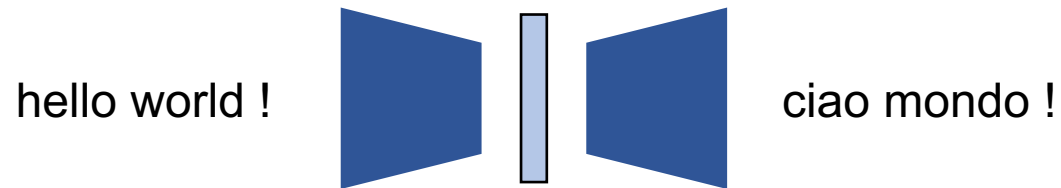
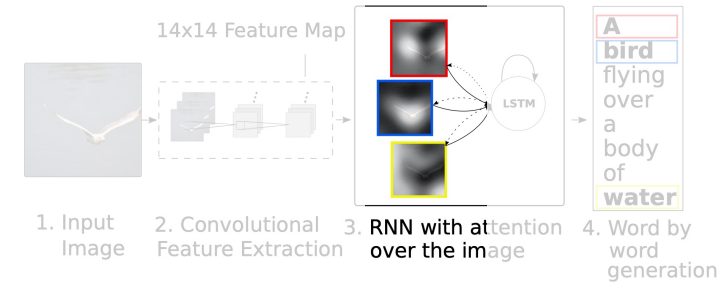
Soft Attention

- Attention Is All You Need, 2017
 - Natural Language Processing (NLP)
 - Language translation



Soft Attention

- Attention Is All You Need, 2017
 - Natural Language Processing (NLP)
 - Language translation

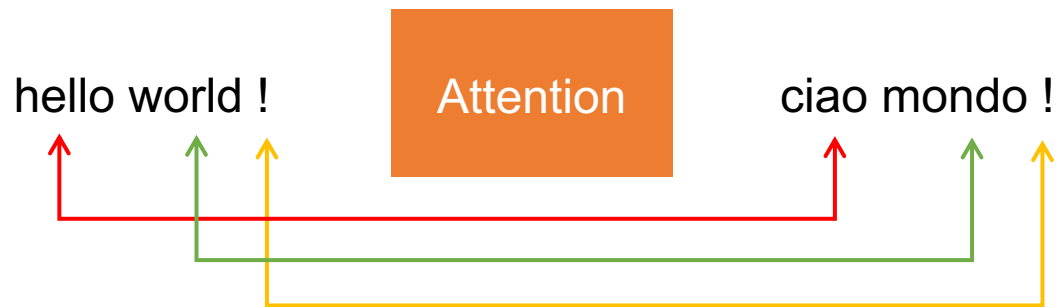
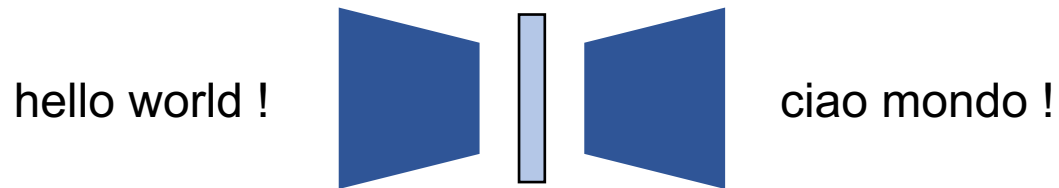
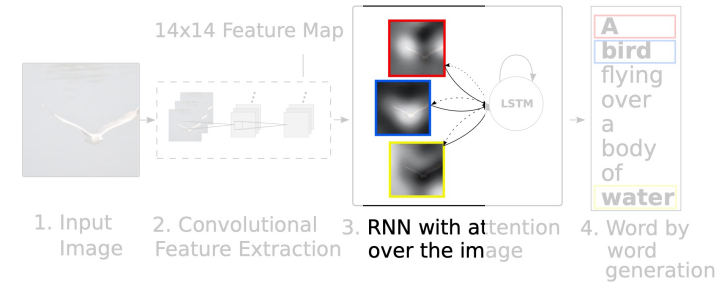


Should the contribution be equal?

Which word do you think should contribute to "ciao" more?

Soft Attention

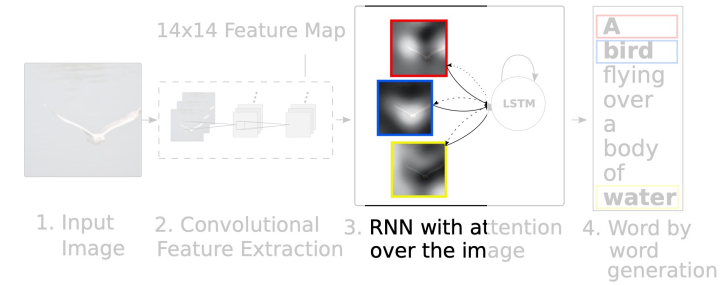
- Attention Is All You Need, 2017
 - Natural Language Processing (NLP)
 - Language translation



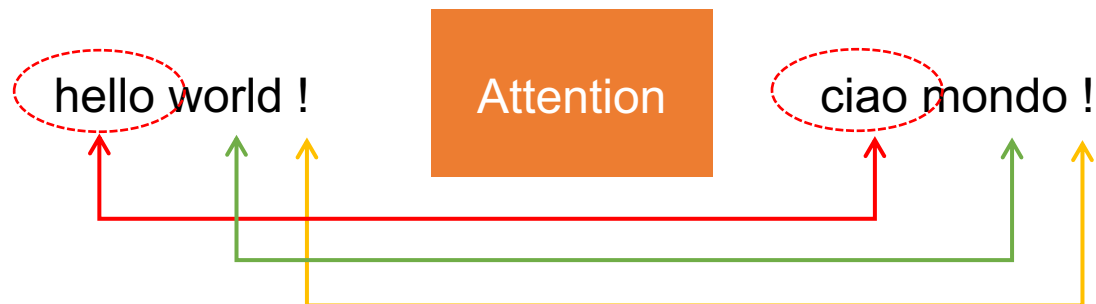
Attention helps the model to understand these relationships

Soft Attention

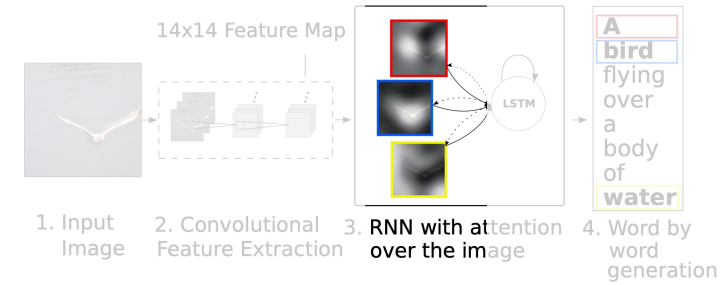
- Attention Is All You Need, 2017
 - Natural Language Processing (NLP)
 - Language translation



Attention mechanism:
How to calculate the influence on *ciao* from *hello*?

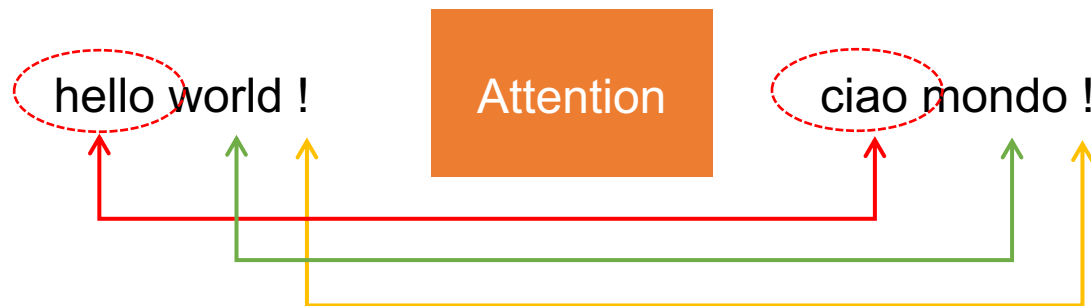


Soft Attention

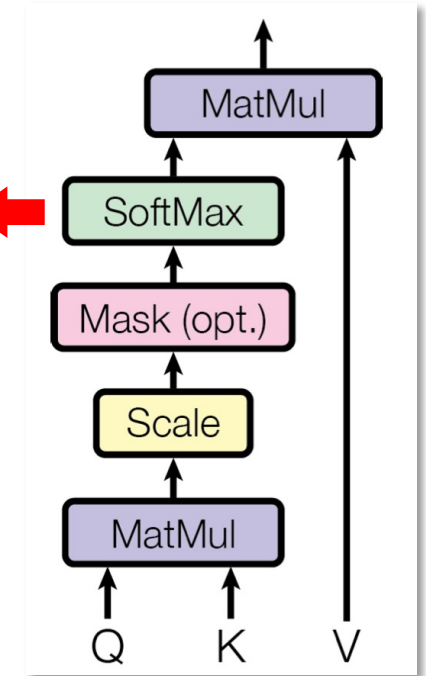


- Transformer: Attention Is All You Need, 2017
 - Natural Language Processing (NLP)
 - Language translation

Attention mechanism:
How to calculate the influence on *ciao* from *hello*?



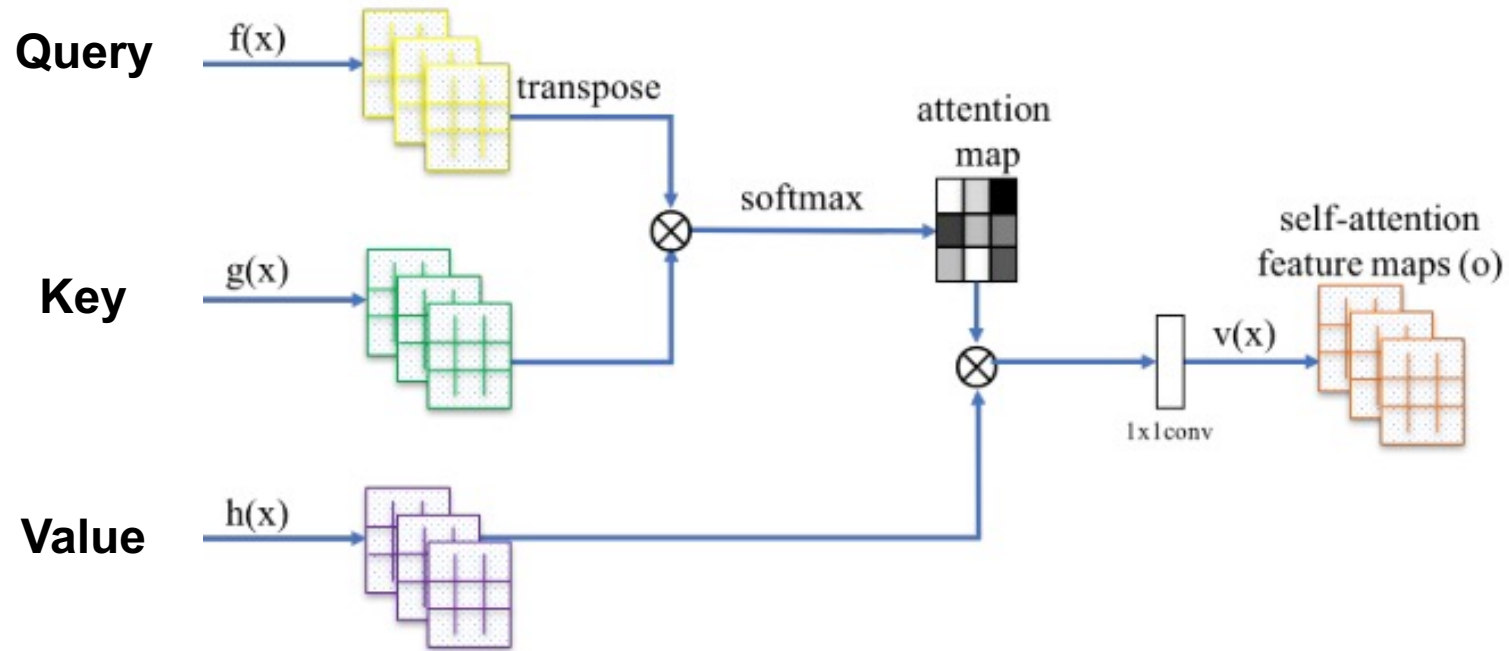
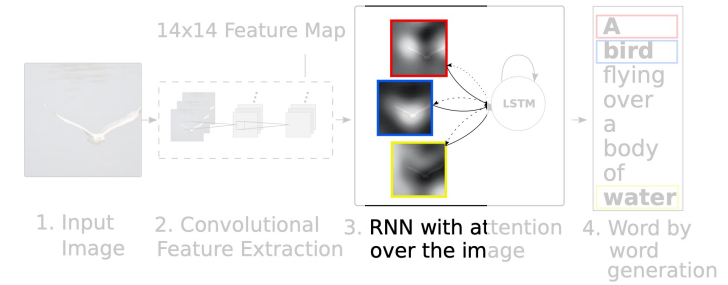
Attention values ←



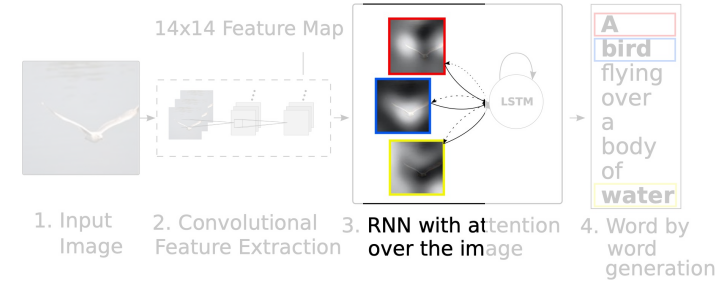
Query: information about *hello* as a feature candidate
Key: information about *ciao* as the output
Value: original feature vector of *hello*

Soft Attention

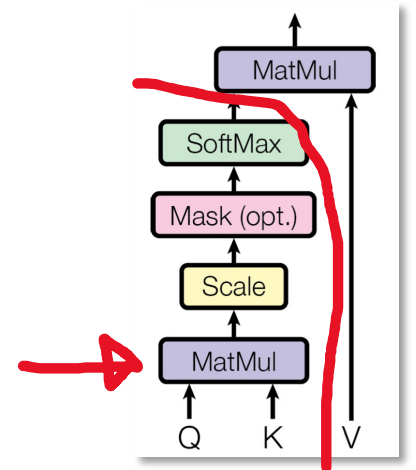
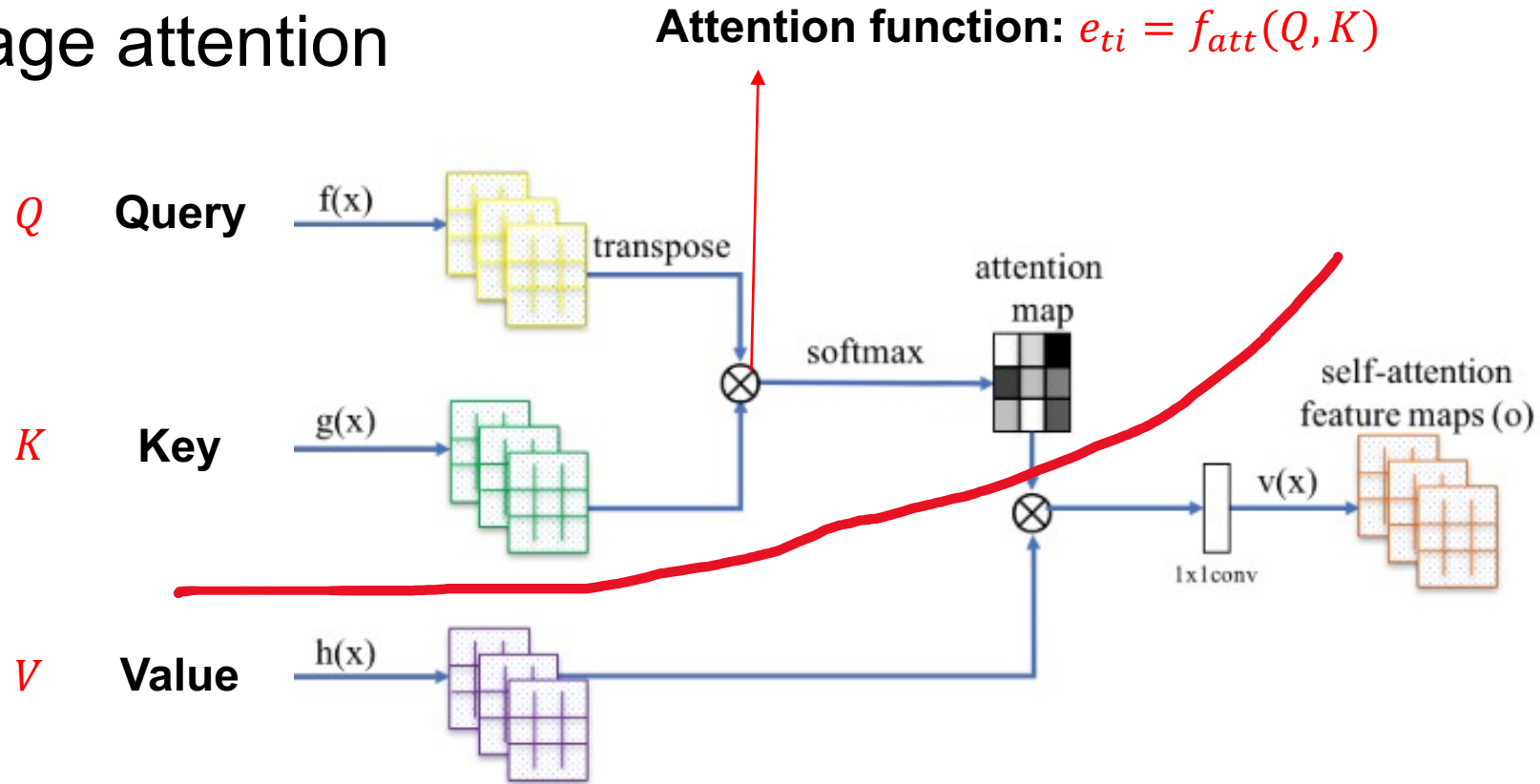
- Image attention



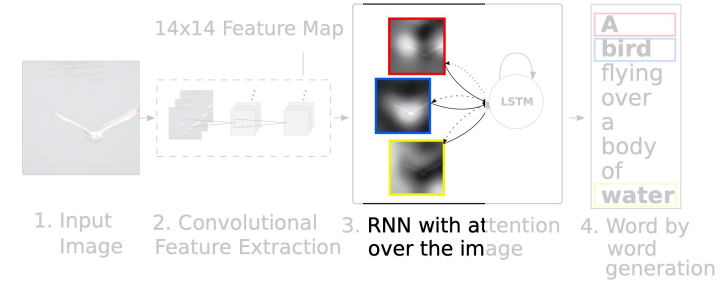
Soft Attention



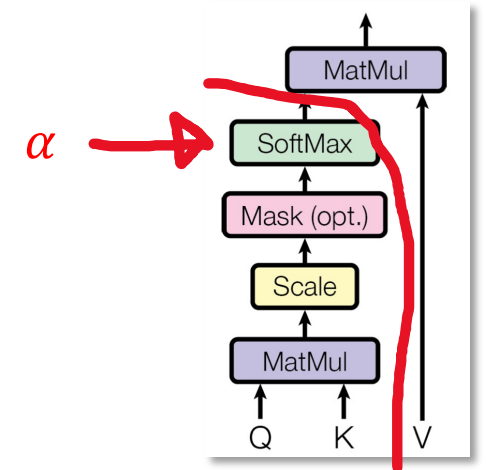
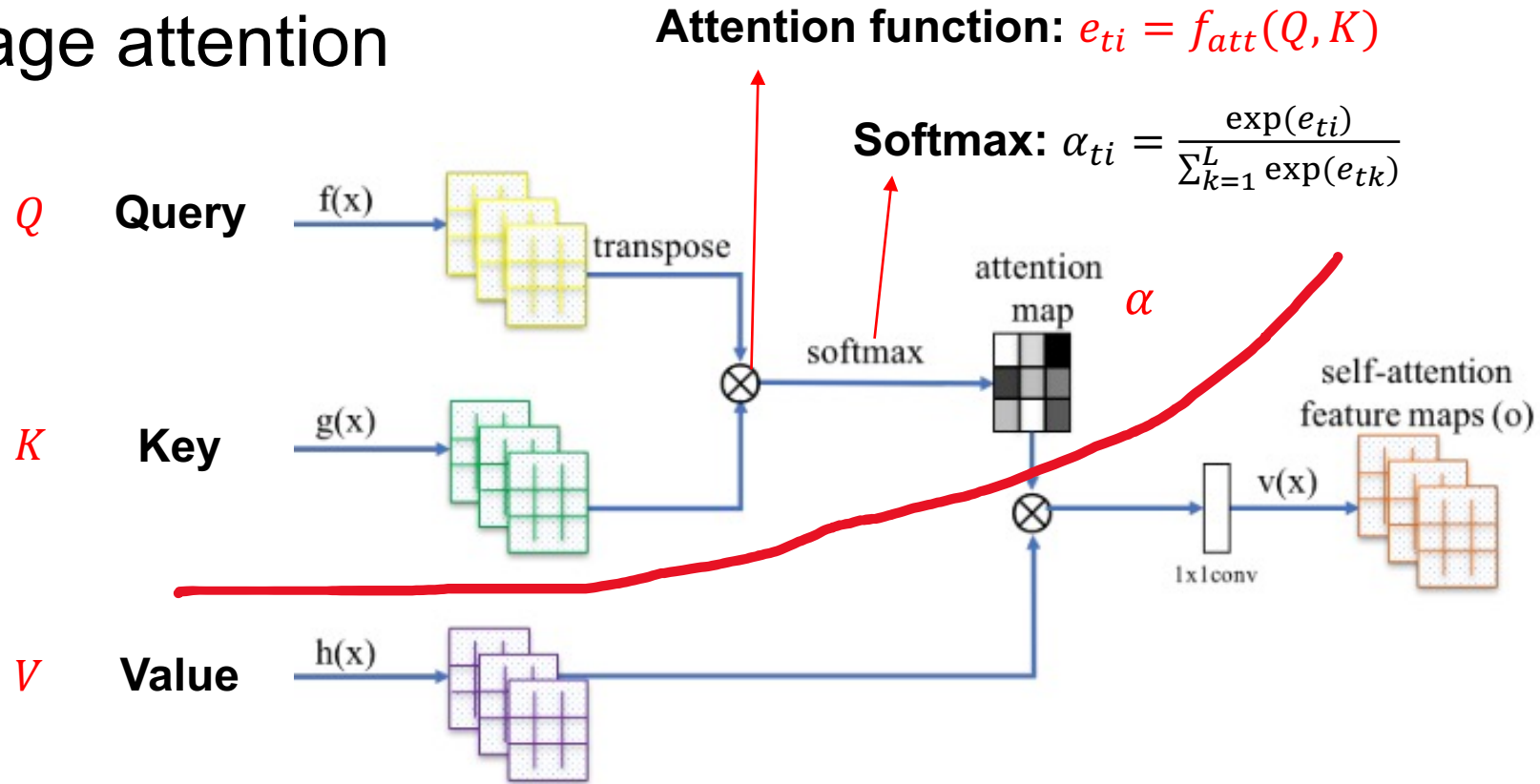
- Image attention



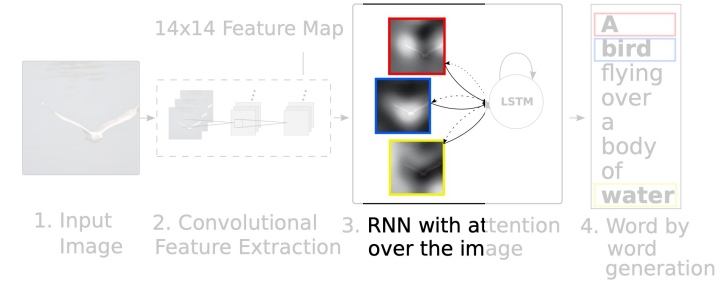
Soft Attention



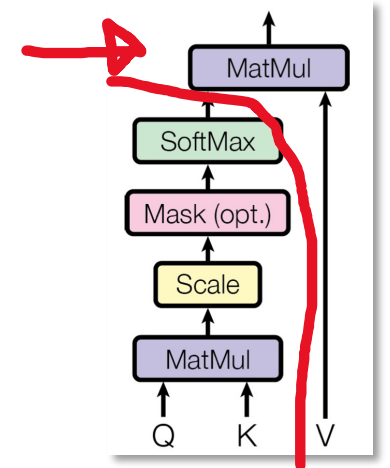
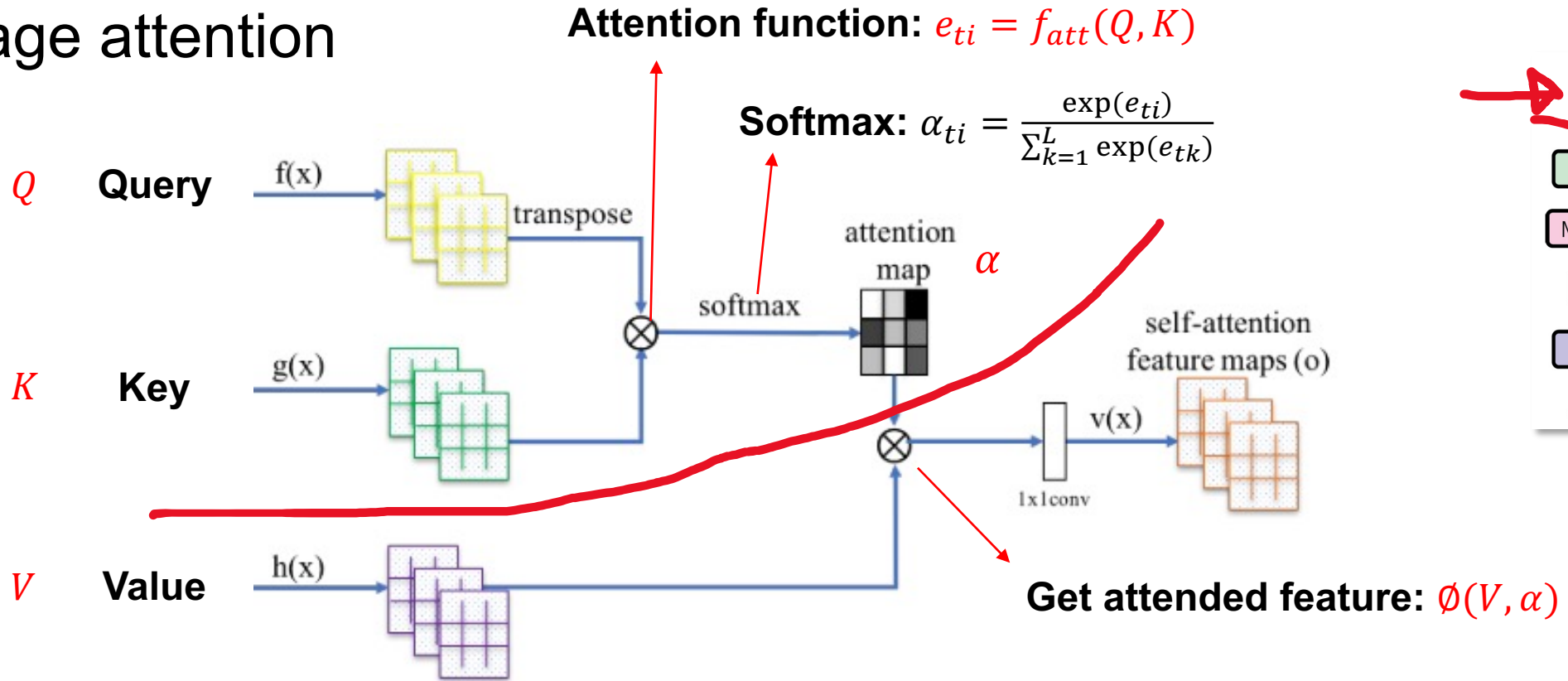
- Image attention



Soft Attention

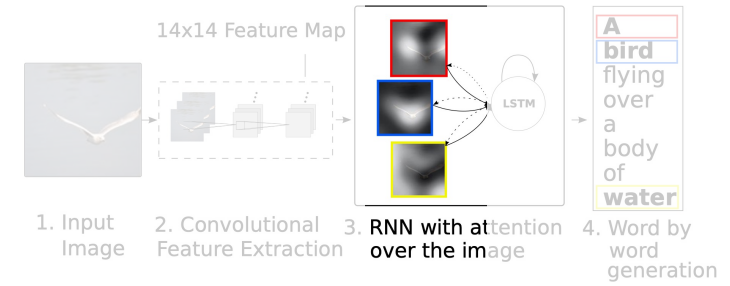


- Image attention



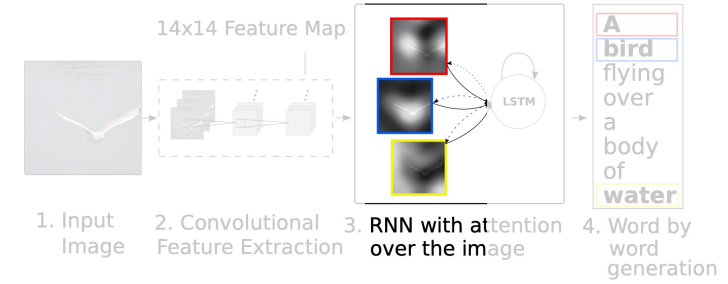
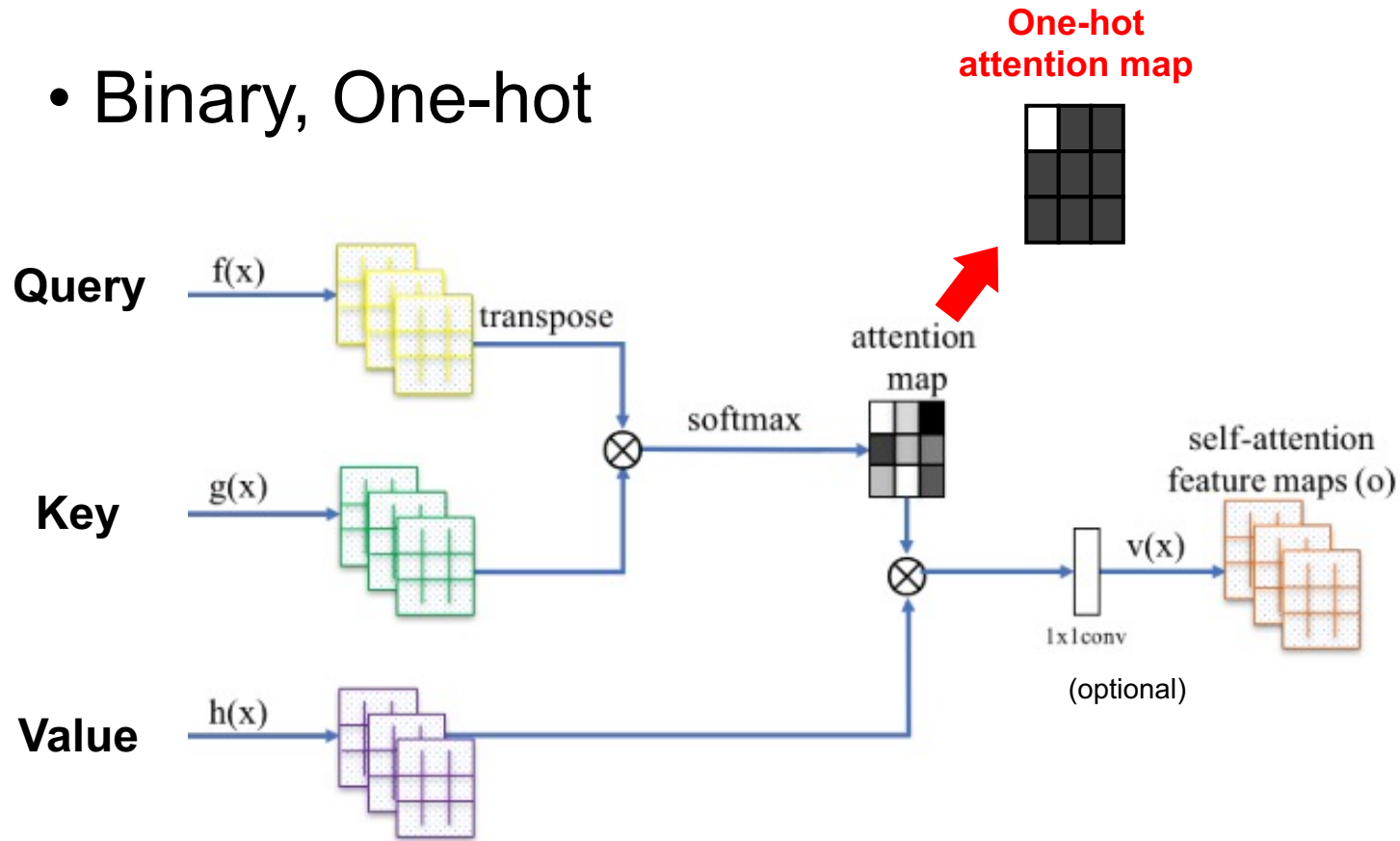
Hard Attention

- Binary, One-hot



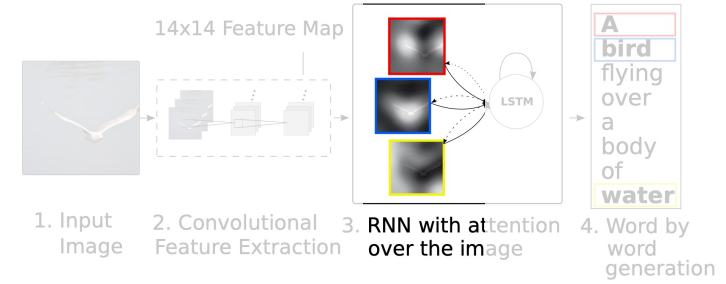
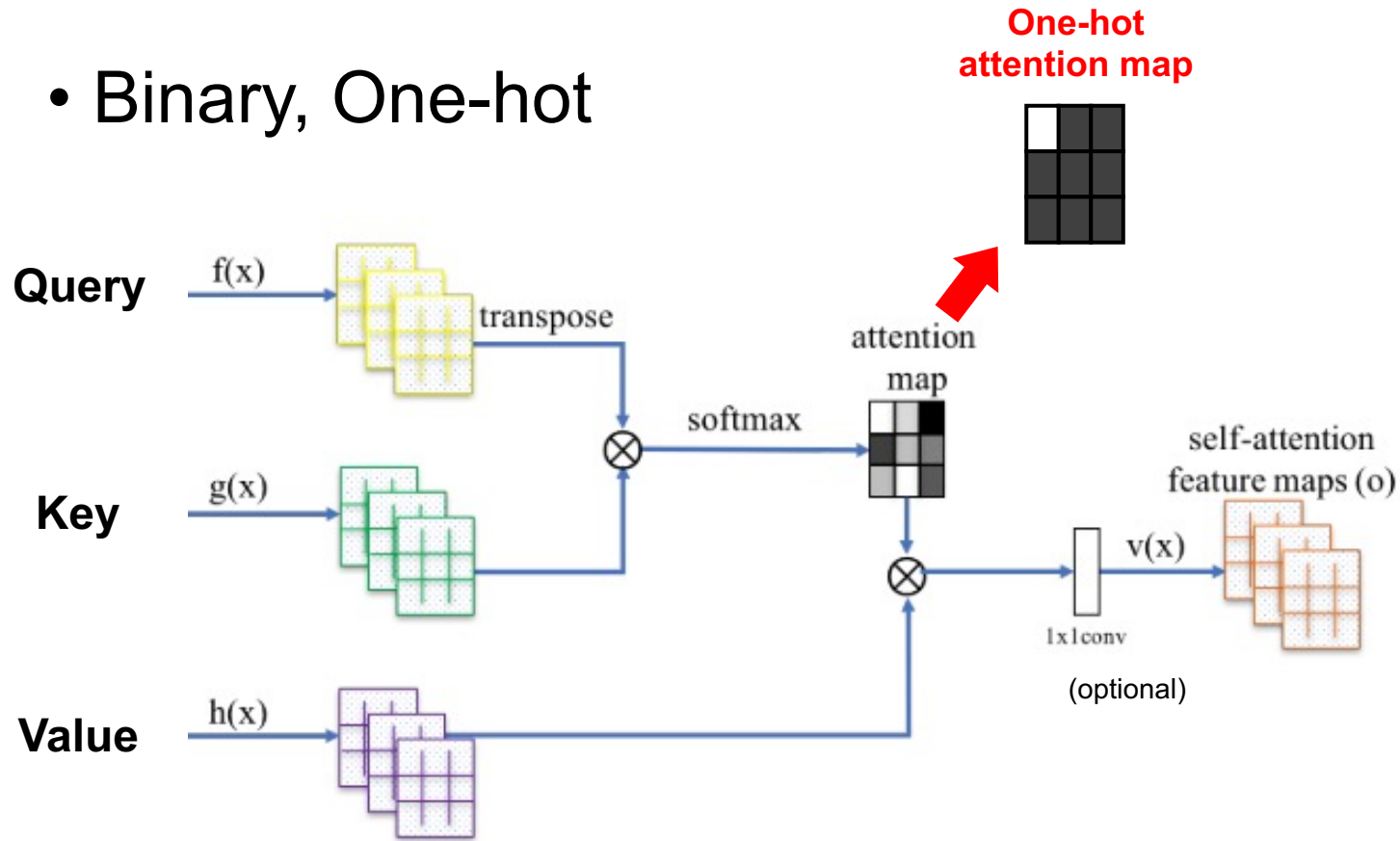
Hard Attention

- Binary, One-hot

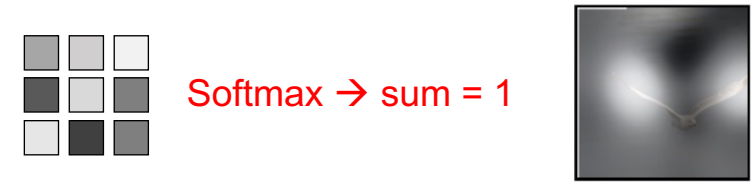


Hard Attention

- Binary, One-hot



Soft Attention:



Hard Attention:

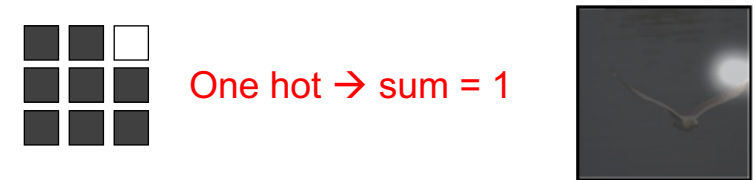
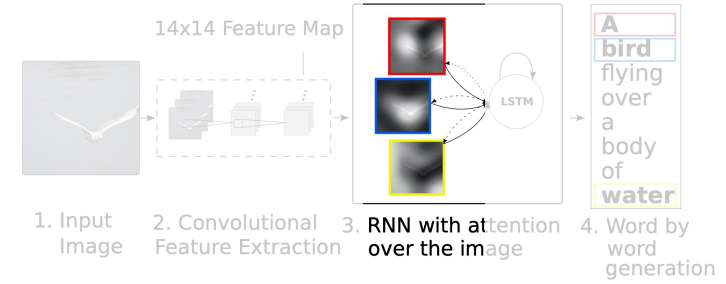
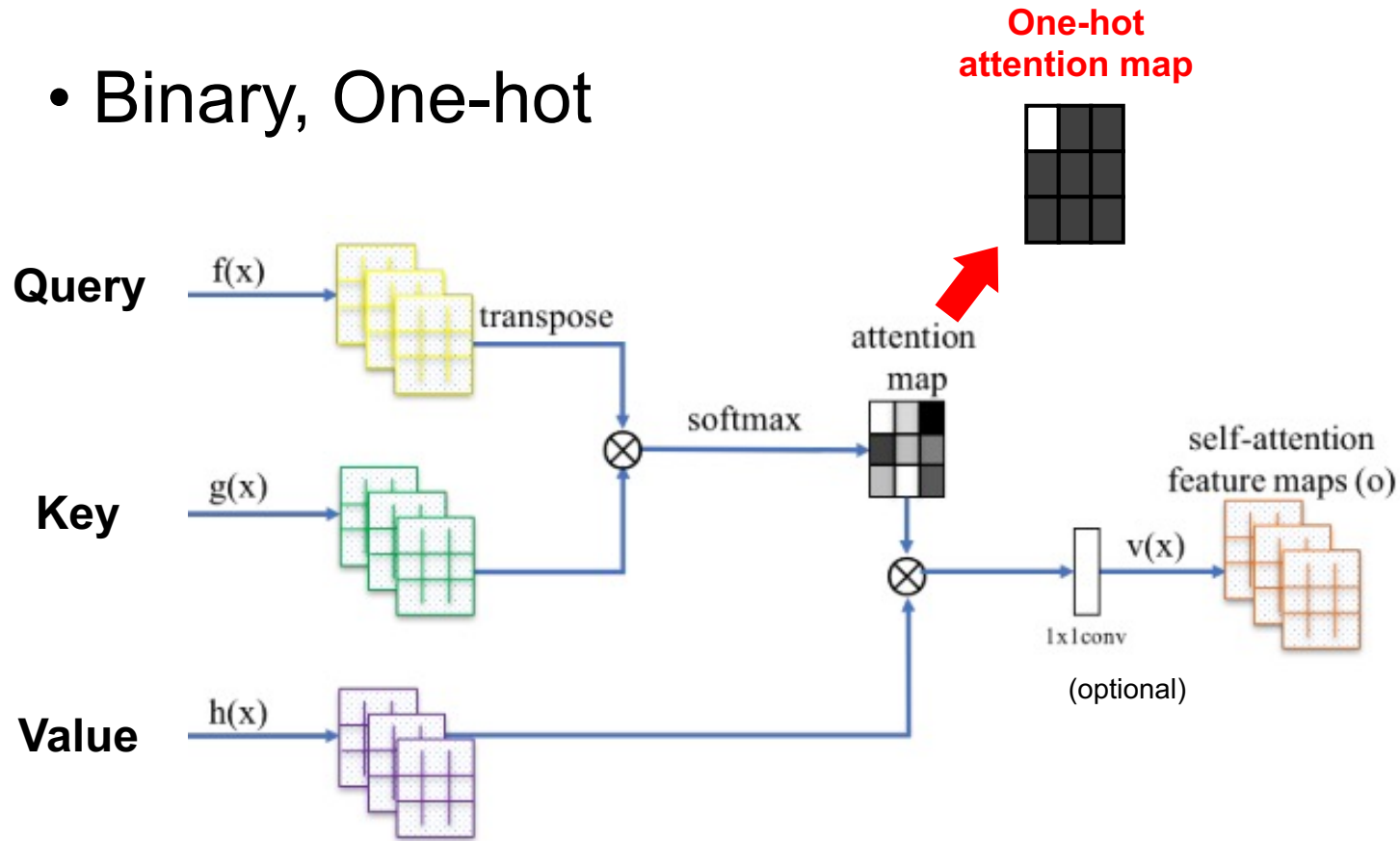


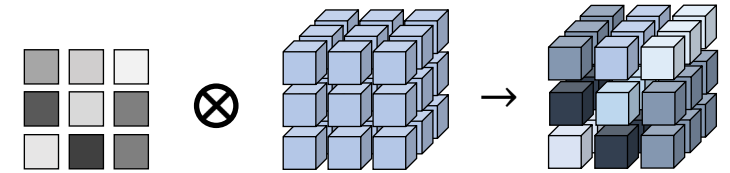
Image source:
<https://medium.com/mllearning-ai/self-attention-in-convolutional-neural-networks-172d947afc00>

Hard Attention

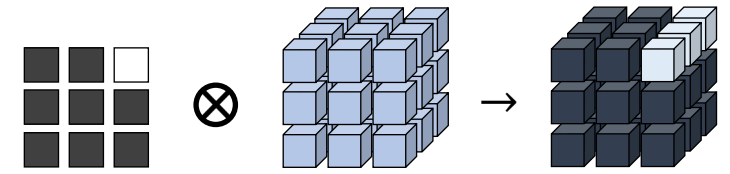
- Binary, One-hot



Soft Attention:

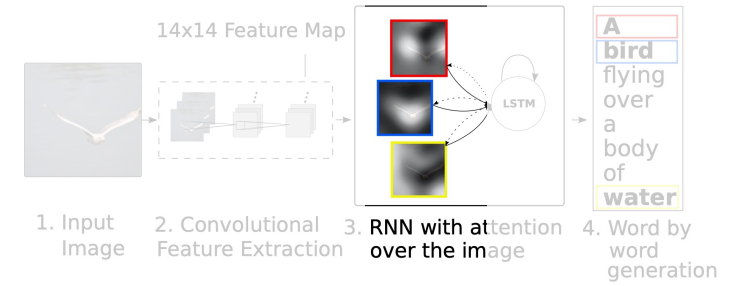


Hard Attention:

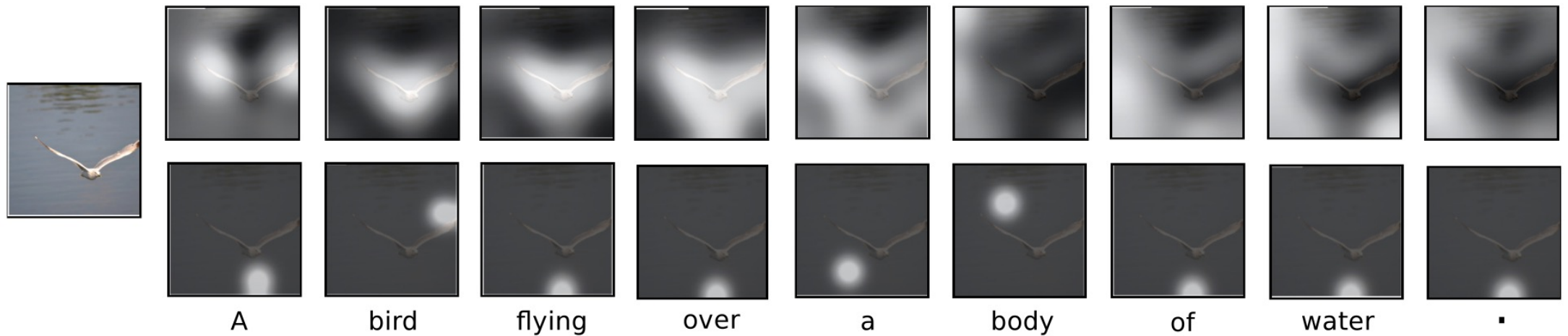


Attention

- Soft attention
- Hard attention

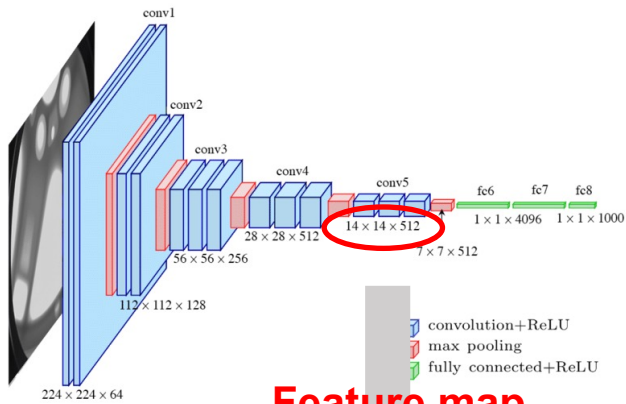


Soft attention: the attention map is a heatmap-like image



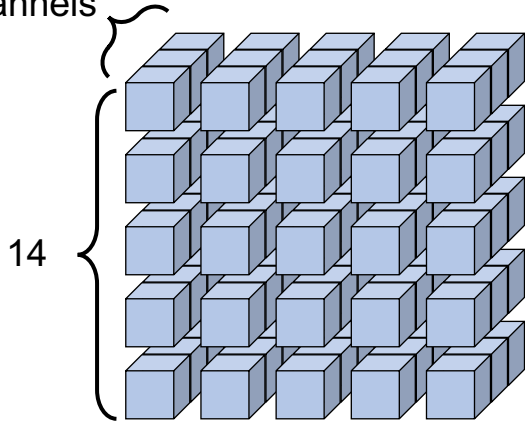
Hard attention: a single attention spot is selected on the map

CNN + Attention + LSTM



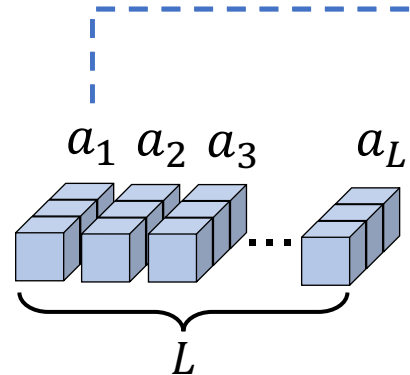
Feature map

512 channels



14

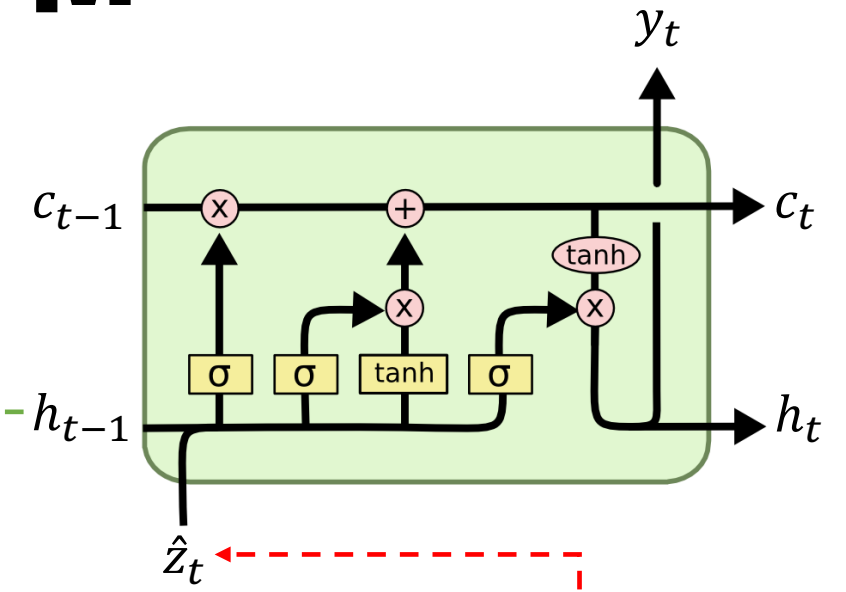
14



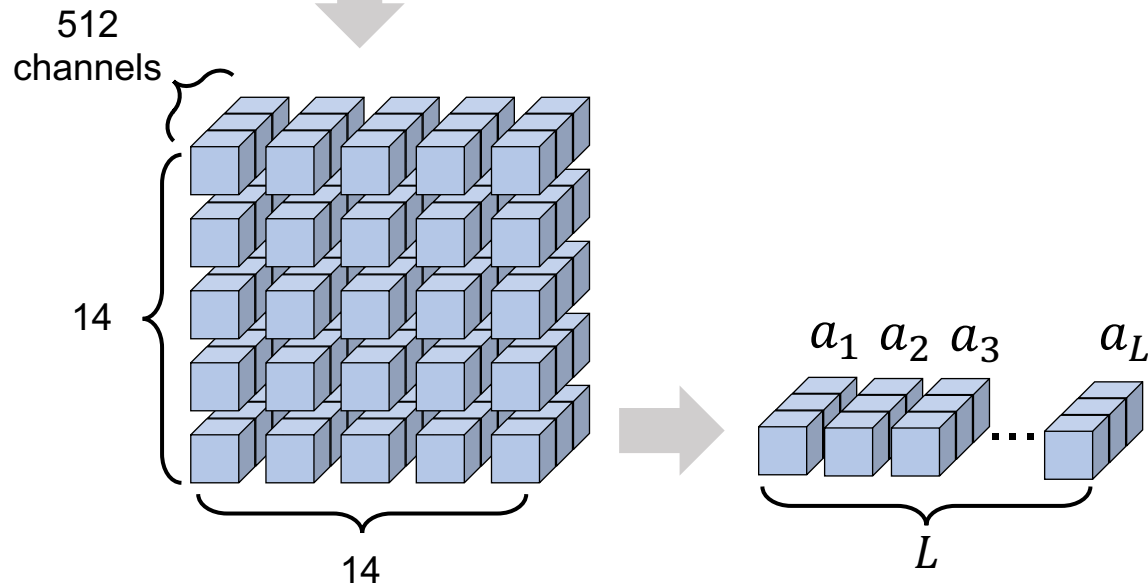
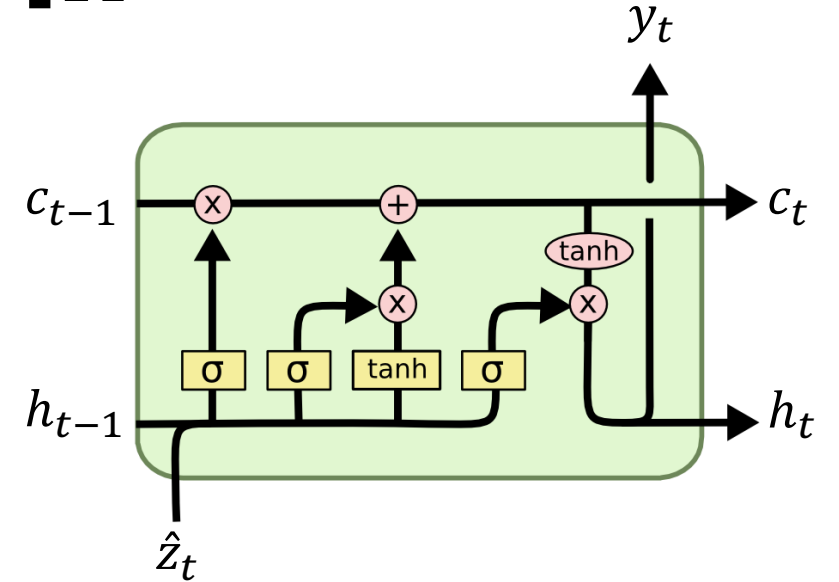
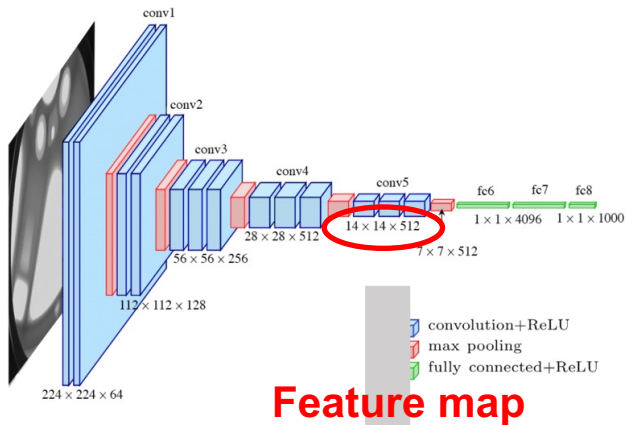
L

$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$

$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$



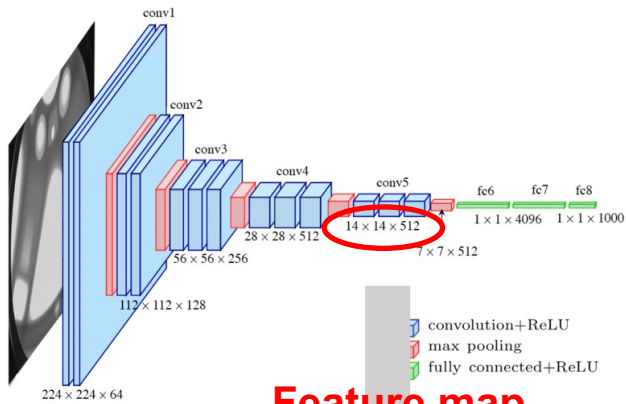
CNN + Attention + LSTM



$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$

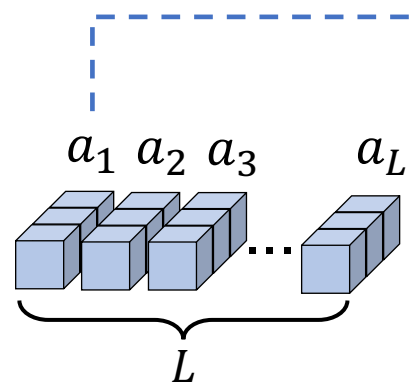
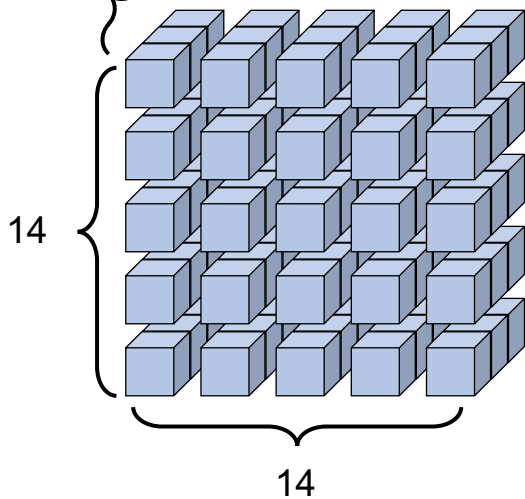
$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$

CNN + Attention + LSTM



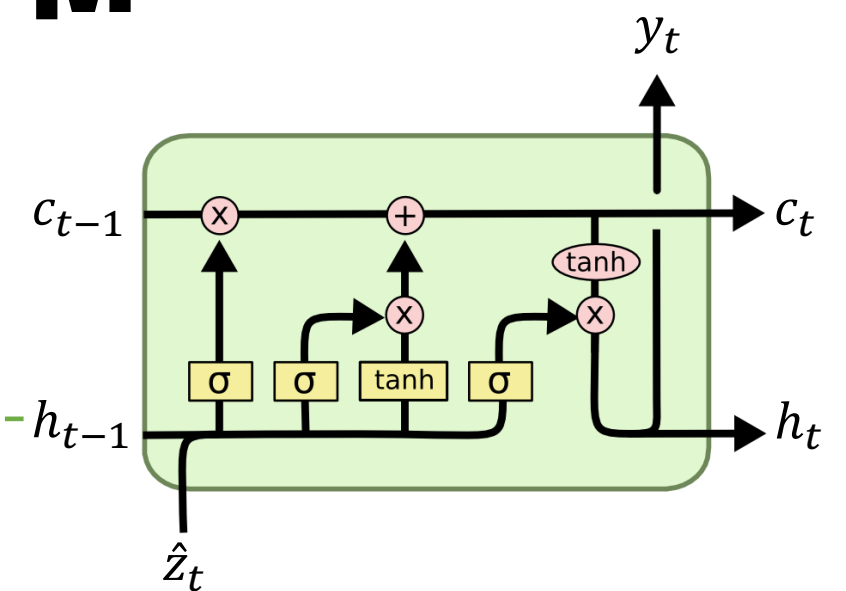
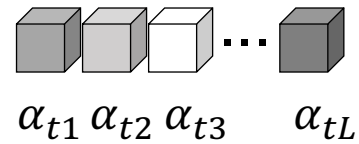
Feature map

512 channels

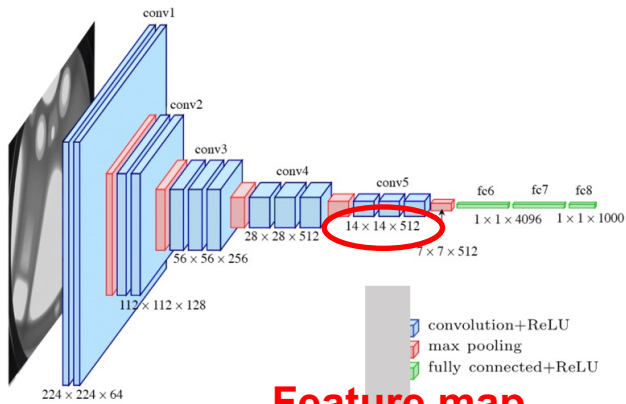


$f_{att}(a_i, h_{t-1}) \rightarrow softmax \rightarrow \alpha_{ti}$

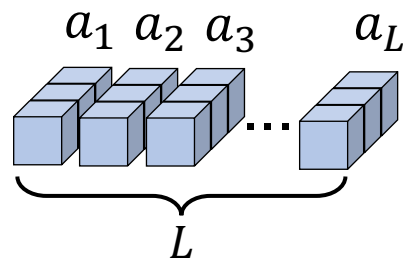
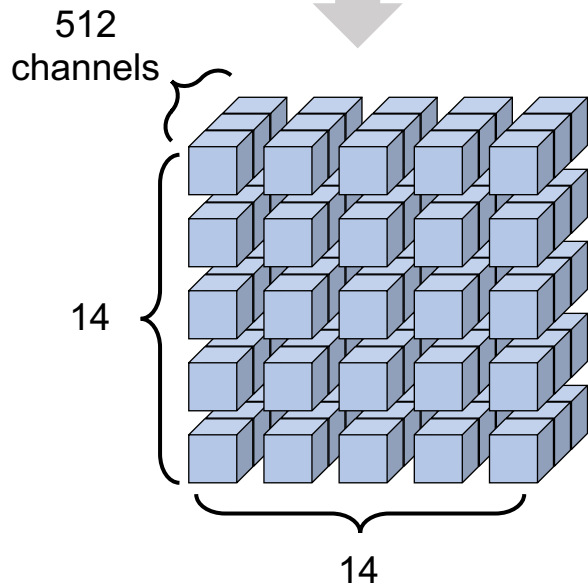
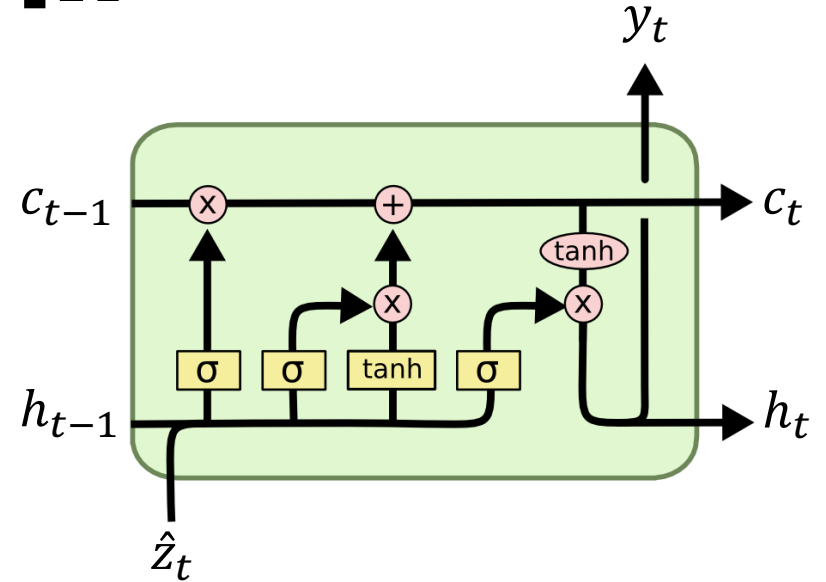
$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$



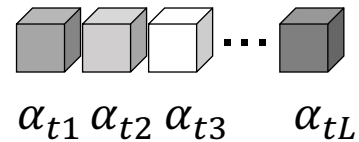
CNN + Attention + LSTM



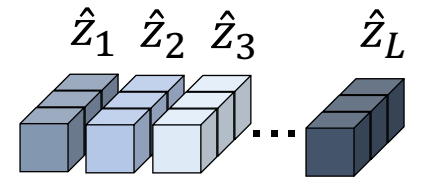
Feature map



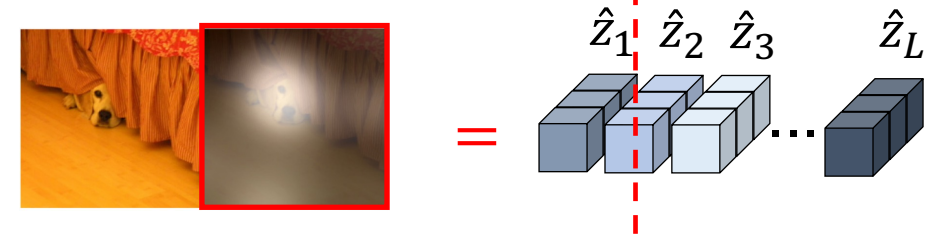
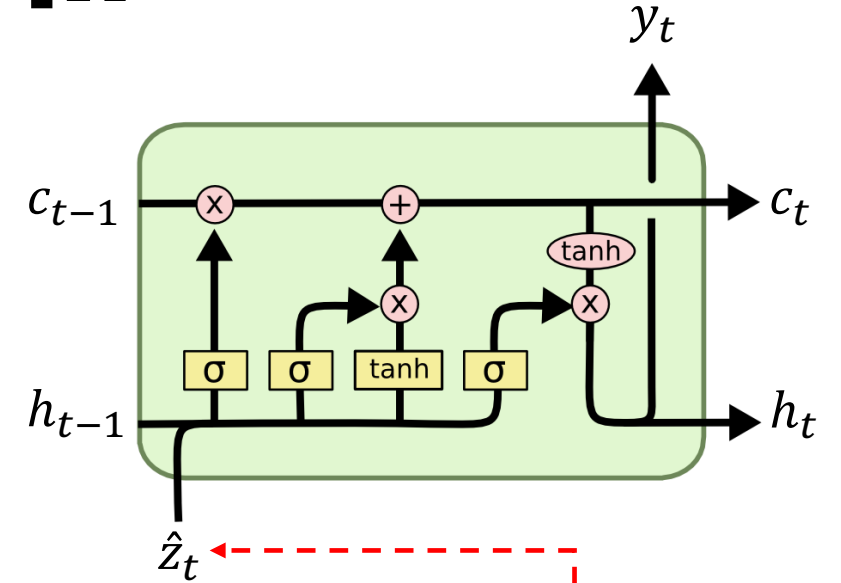
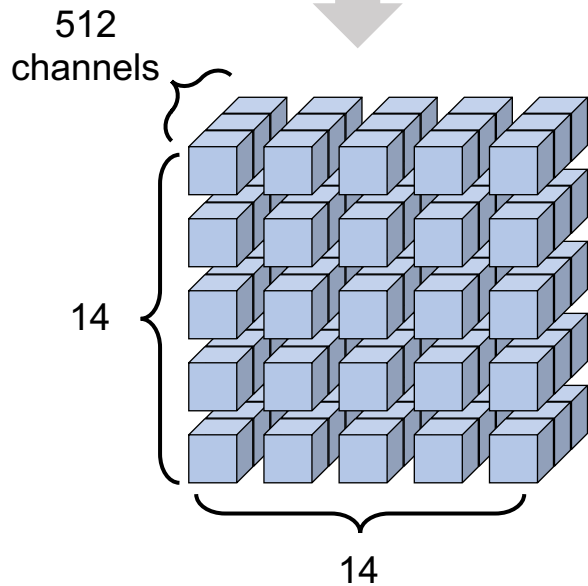
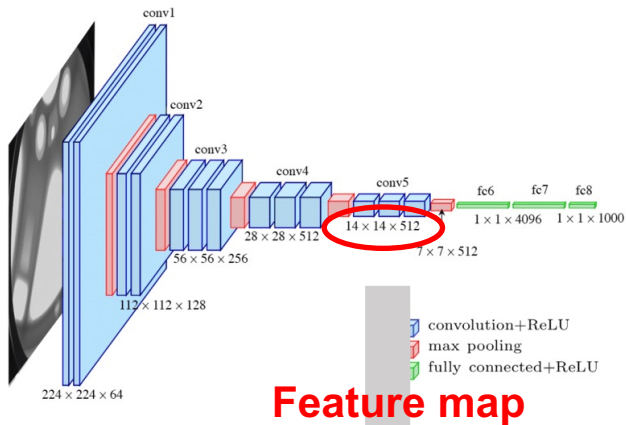
$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$



$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$



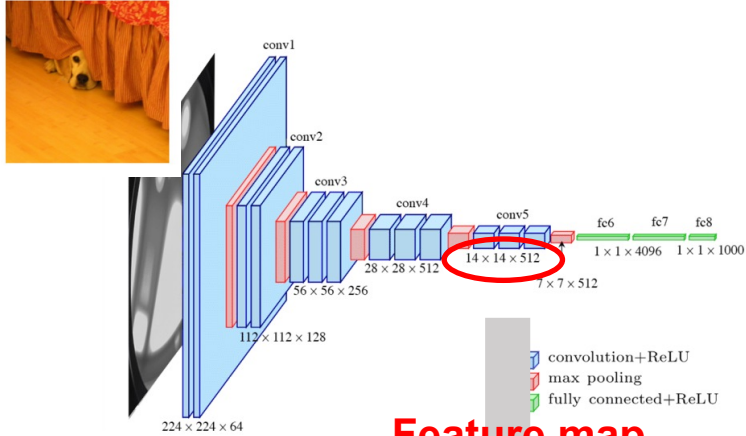
CNN + Attention + LSTM



$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$

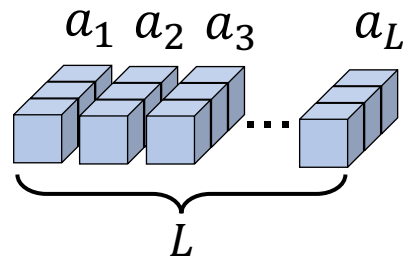
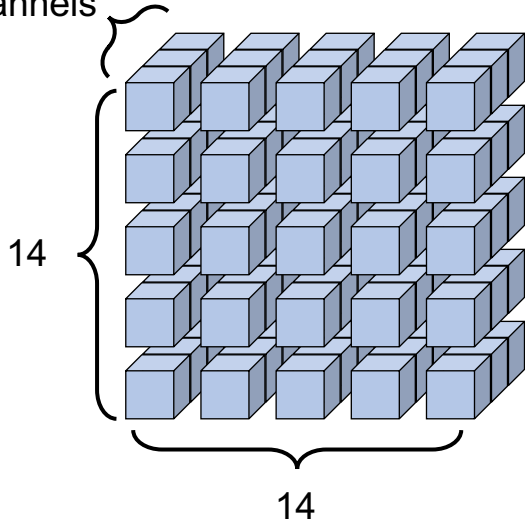
$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$

CNN + Attention + LSTM

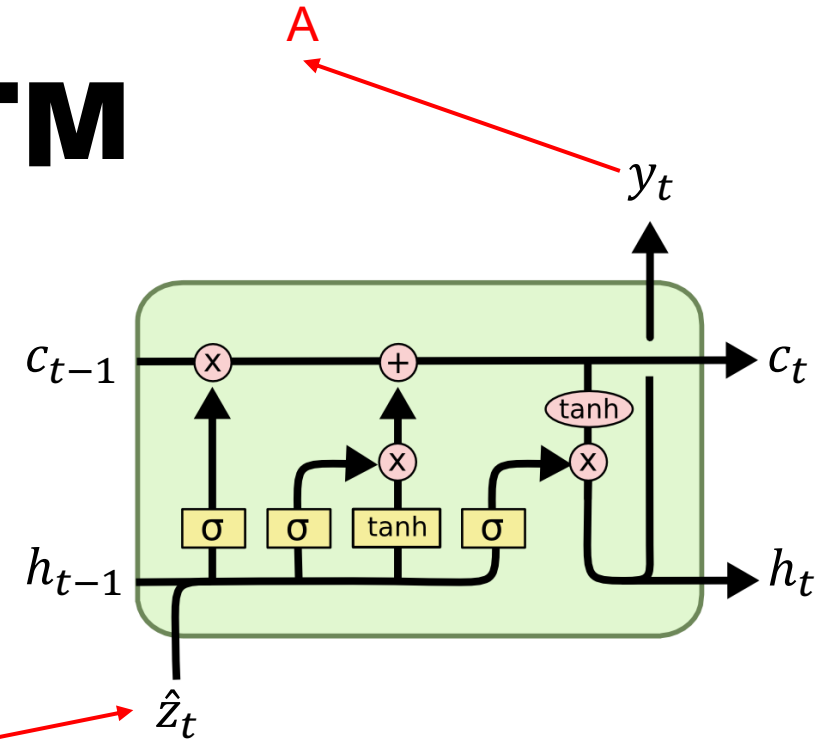


Feature map

512 channels

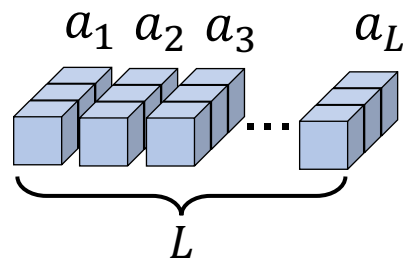
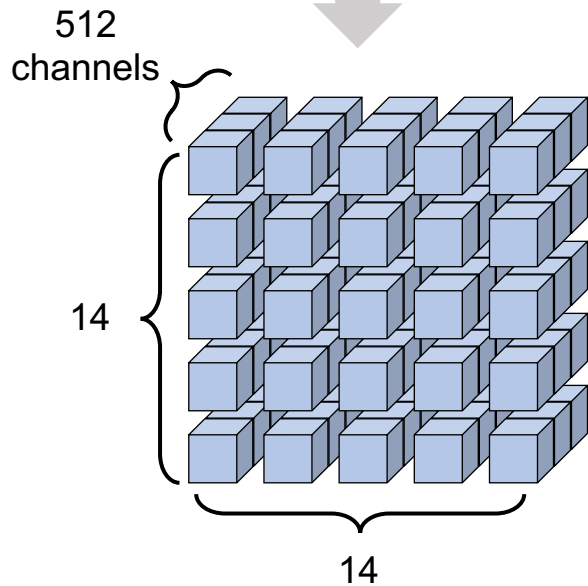
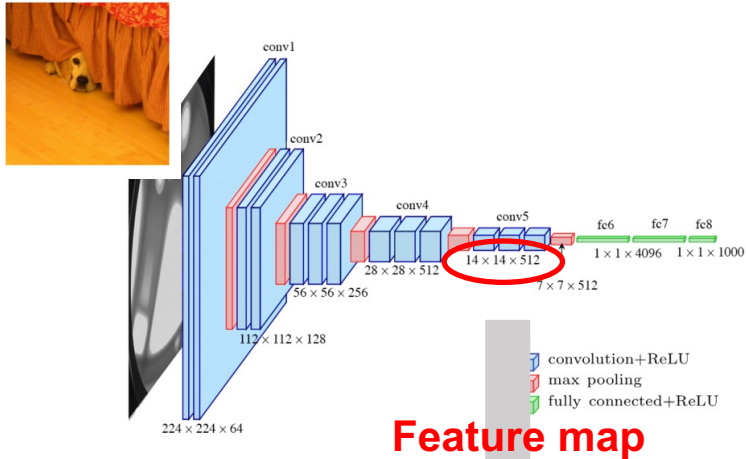


$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$



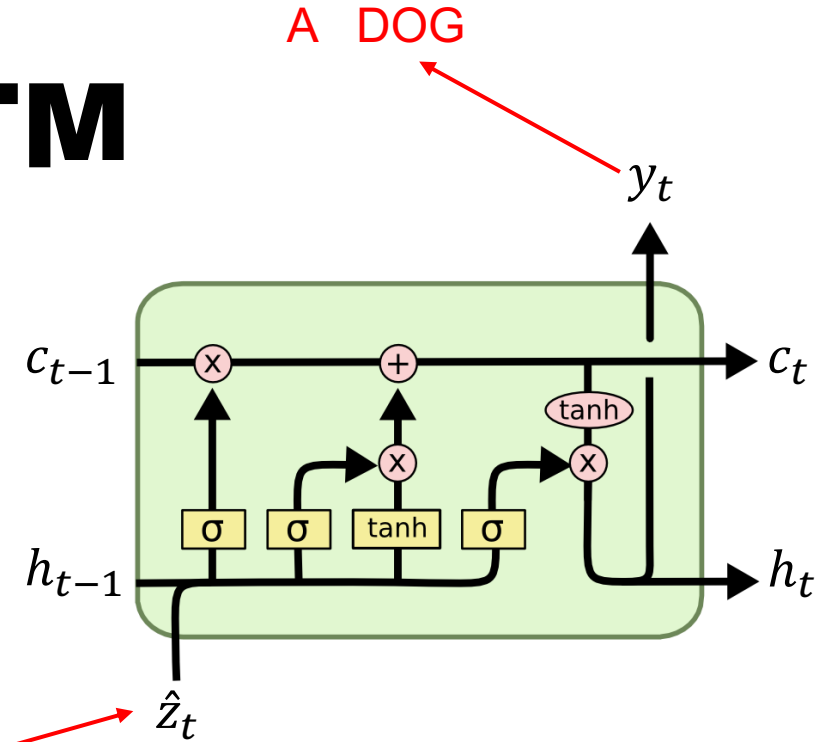
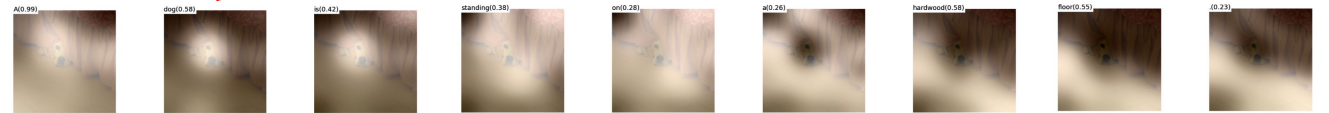
$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$

CNN + Attention + LSTM

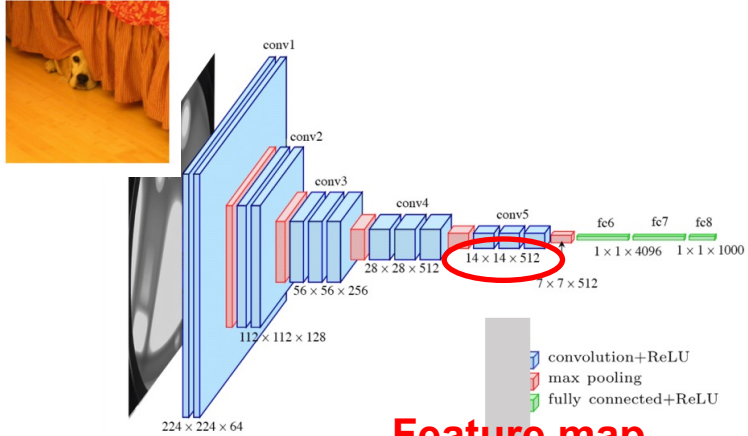


$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$

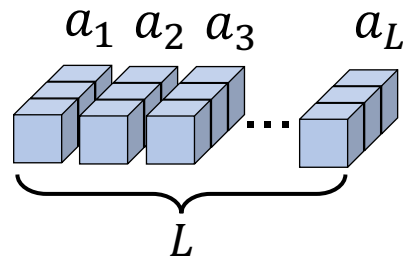
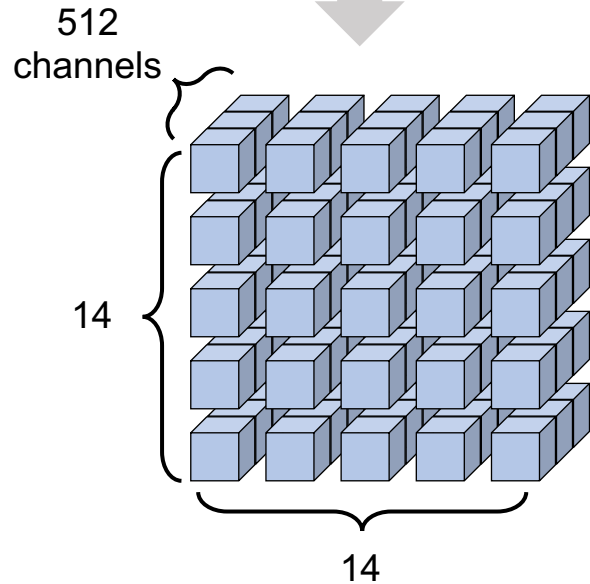
$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$



CNN + Attention + LSTM

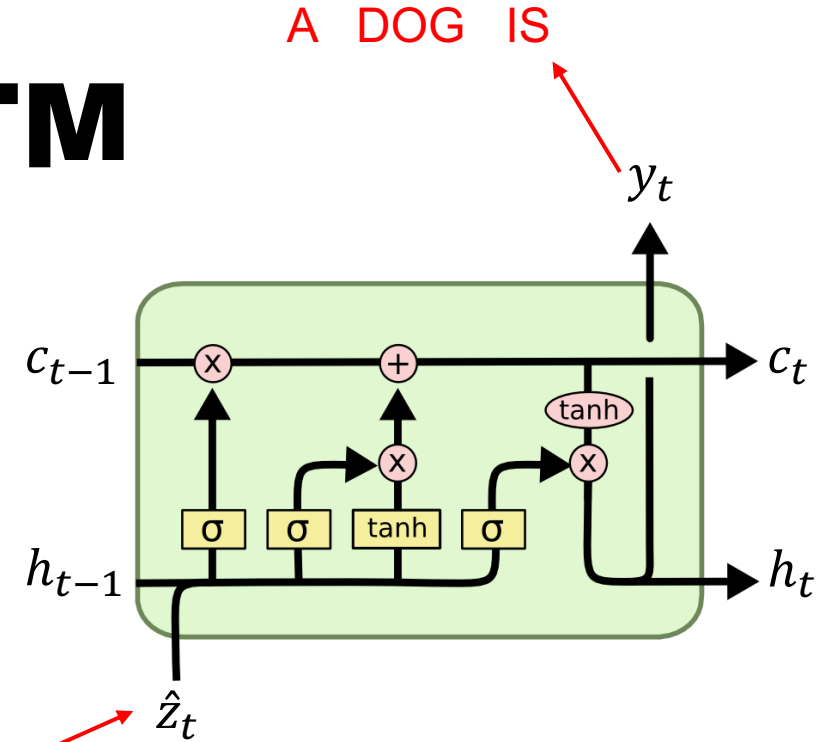


Feature map



$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$

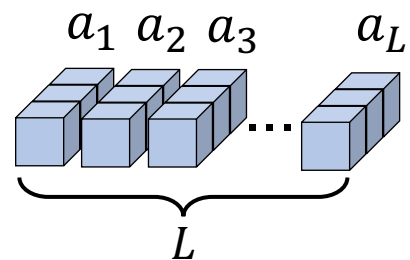
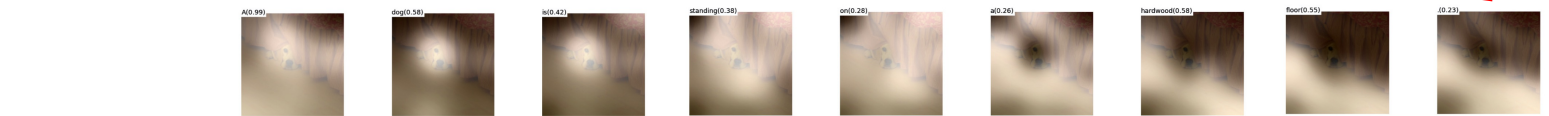
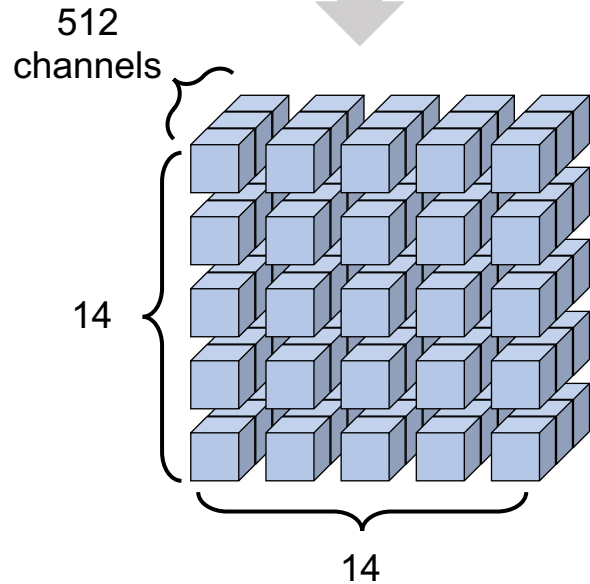
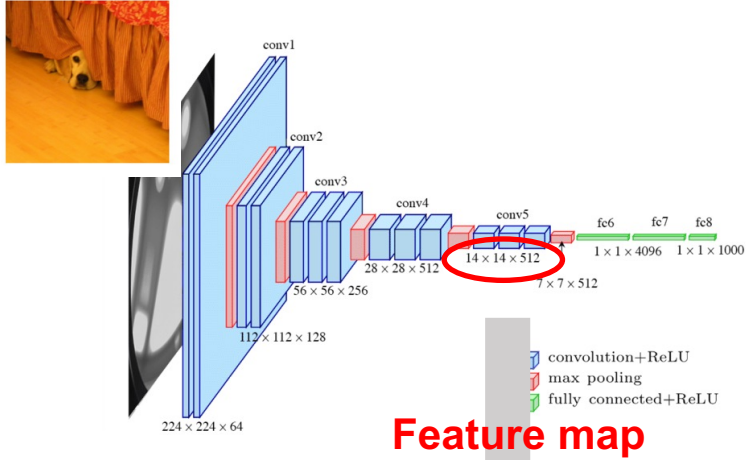
$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$



A DOG IS

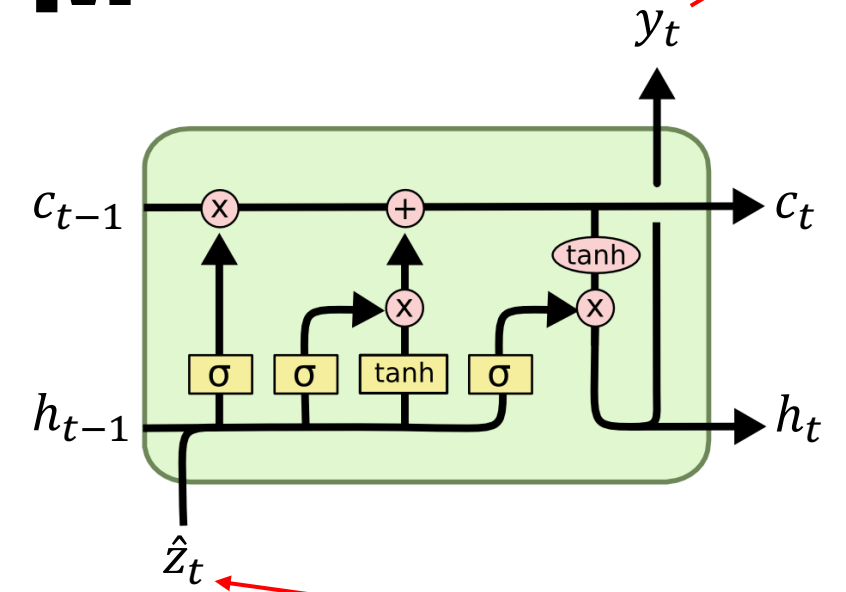
CNN + Attention + LSTM

A DOG IS FLOOR .



$$f_{att}(a_i, h_{t-1}) \rightarrow \text{softmax} \rightarrow \alpha_{ti}$$

$$\Phi(\{a_i\}, \{\alpha_{ti}\}) = \hat{z}_t$$



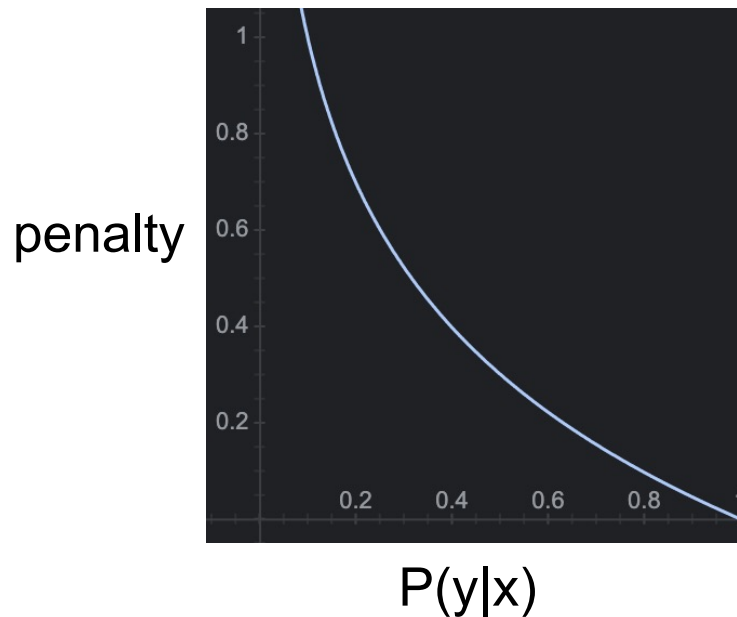
↑

↑

Loss Function

$$L_d = \underline{-\log(P(\mathbf{y}|\mathbf{x}))} + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$

negative log-likelihood



x:



y: The dog is laying under a bed

Loss Function

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \underbrace{\sum_i^L (1 - \sum_t^C \alpha_{ti})^2}$$

$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Softmax, so $\sum_i \alpha_{ti} = 1$

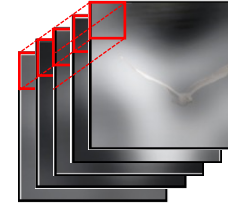


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Softmax, so $\sum_i \alpha_{ti} = 1$

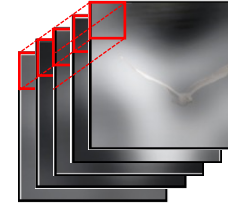


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i (1 - \sum_t \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

α_t : α_1 α_2 α_3

0.1	0.4	0.1	0.1	0.0	0.4
0.2	0.3	0.1	0.7	0.3	0.3

Softmax, so $\sum_i \alpha_{ti} = 1$

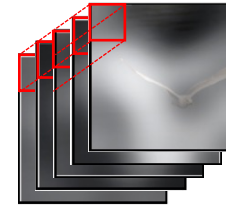


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i (1 - \sum_t \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

α_t : α_1 α_2 α_3

0.1	0.4
0.2	0.3

0.1	0.1
0.1	0.7

0.0	0.4
0.3	0.3

$1 - \sum_t \alpha_{ti}$

?	

Softmax, so $\sum_i \alpha_{ti} = 1$

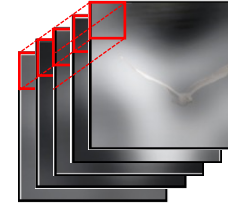


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i (1 - \sum_t \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

$\alpha_t:$	α_1	α_2	α_3	$1 - \sum_t \alpha_{ti}$
	0.1	0.1	0.0	0.8
	0.4	0.1	0.4	?
	0.2	0.7	0.3	
	0.3	0.3	0.3	

Softmax, so $\sum_i \alpha_{ti} = 1$

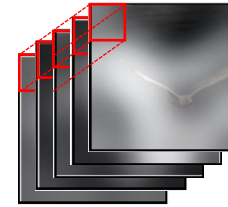


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i (1 - \sum_t \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

α_t : α_1

0.1	0.4
0.2	0.3

α_2

0.1	0.1
0.1	0.7

α_3

0.0	0.4
0.3	0.3

$1 - \sum_t \alpha_{ti}$

0.8	0.1
?	

Softmax, so $\sum_i \alpha_{ti} = 1$

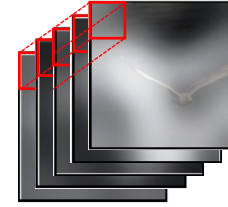


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i (1 - \sum_t \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

$\alpha_t:$	α_1	α_2	α_3	$1 - \sum_t \alpha_{ti}$
	0.1 0.4	0.1 0.1	0.0 0.4	0.8 0.1
	0.2 0.3	0.1 0.7	0.3 0.3	0.4 ?

Softmax, so $\sum_i \alpha_{ti} = 1$

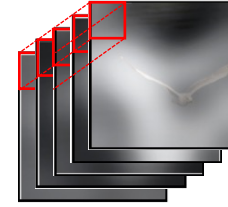


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i (1 - \sum_t \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

α_t : α_1

0.1	0.4
0.2	0.3

α_2

0.1	0.1
0.1	0.7

α_3

0.0	0.4
0.3	0.3

$1 - \sum_t \alpha_{ti}$

0.8	0.1
0.4	-0.3

Softmax, so $\sum_i \alpha_{ti} = 1$

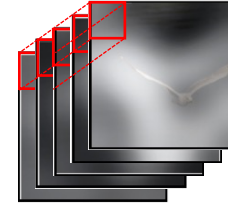


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

$\alpha_t:$	α_1	α_2	α_3	$1 - \sum_t^C \alpha_{ti}$																
	<table border="1"><tr><td>0.1</td><td>0.4</td></tr><tr><td>0.2</td><td>0.3</td></tr></table>	0.1	0.4	0.2	0.3	<table border="1"><tr><td>0.1</td><td>0.1</td></tr><tr><td>0.1</td><td>0.7</td></tr></table>	0.1	0.1	0.1	0.7	<table border="1"><tr><td>0.0</td><td>0.4</td></tr><tr><td>0.3</td><td>0.3</td></tr></table>	0.0	0.4	0.3	0.3	<table border="1"><tr><td>0.8</td><td>0.1</td></tr><tr><td>0.4</td><td>-0.3</td></tr></table>	0.8	0.1	0.4	-0.3
0.1	0.4																			
0.2	0.3																			
0.1	0.1																			
0.1	0.7																			
0.0	0.4																			
0.3	0.3																			
0.8	0.1																			
0.4	-0.3																			

Softmax, so $\sum_i \alpha_{ti} = 1$



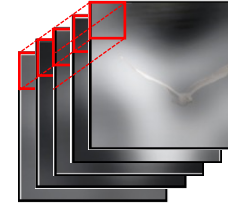
→ Sum over **positions** = 1

$$\sum_i^L (1 - \sum_t^C \alpha_{ti})^2 = ?$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: Calculate the loss for these attention maps.

$\alpha_t:$	α_1	α_2	α_3	$1 - \sum_t^C \alpha_{ti}$																
	<table border="1"><tr><td>0.1</td><td>0.4</td></tr><tr><td>0.2</td><td>0.3</td></tr></table>	0.1	0.4	0.2	0.3	<table border="1"><tr><td>0.1</td><td>0.1</td></tr><tr><td>0.1</td><td>0.7</td></tr></table>	0.1	0.1	0.1	0.7	<table border="1"><tr><td>0.0</td><td>0.4</td></tr><tr><td>0.3</td><td>0.3</td></tr></table>	0.0	0.4	0.3	0.3	<table border="1"><tr><td>0.8</td><td>0.1</td></tr><tr><td>0.4</td><td>-0.3</td></tr></table>	0.8	0.1	0.4	-0.3
0.1	0.4																			
0.2	0.3																			
0.1	0.1																			
0.1	0.7																			
0.0	0.4																			
0.3	0.3																			
0.8	0.1																			
0.4	-0.3																			

Softmax, so $\sum_i \alpha_{ti} = 1$



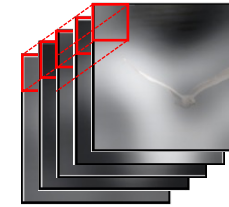
→ Sum over **positions** = 1

$$\begin{aligned} \sum_i^L (1 - \sum_t^C \alpha_{ti})^2 &= 0.8^2 + 0.1^2 + 0.4^2 + 0.3^2 \\ &= 0.9 \end{aligned}$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i (1 - \sum_t \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i (1 - \sum_t \alpha_{ti})^2$?

$\alpha_t:$	α_1	α_2	α_3
?	?	?	?
?	?	?	?

Softmax, so $\sum_i \alpha_{ti} = 1$

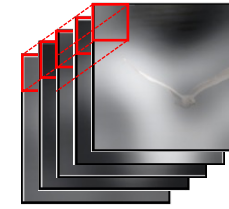


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

$\alpha_t:$	α_1	α_2	α_3												
	<table border="1" style="border-collapse: collapse; width: 40px; height: 40px; text-align: center;"> <tr><td>?</td><td>?</td></tr> <tr><td>?</td><td>?</td></tr> </table>	?	?	?	?	<table border="1" style="border-collapse: collapse; width: 40px; height: 40px; text-align: center;"> <tr><td>?</td><td>?</td></tr> <tr><td>?</td><td>?</td></tr> </table>	?	?	?	?	<table border="1" style="border-collapse: collapse; width: 40px; height: 40px; text-align: center;"> <tr><td>?</td><td>?</td></tr> <tr><td>?</td><td>?</td></tr> </table>	?	?	?	?
?	?														
?	?														
?	?														
?	?														
?	?														
?	?														

$$\sum_t^C \alpha_{ti}$$

Softmax, so $\sum_i \alpha_{ti} = 1$

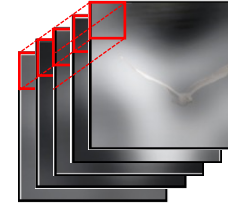


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

$\alpha_t:$

α_1
? ?
? ?

α_2
? ?
? ?

α_3
? ?
? ?

$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

Softmax, so $\sum_i \alpha_{ti} = 1$

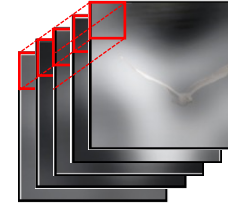


→ Sum over **positions** = 1

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

α_t : α_1 α_2 α_3

?	?	?	?	?	?
?	?	?	?	?	?

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

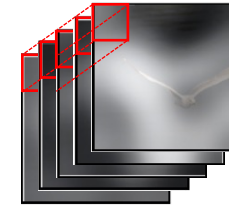
$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

$$|1 - \sum_t^C \alpha_{ti}|$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

α_t : α_1 α_2 α_3

?	?	?	?	?	?
?	?	?	?	?	?

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

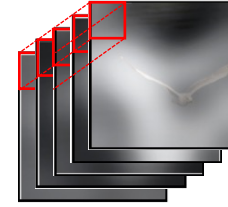
$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

$$0 \leq |1 - \sum_t^C \alpha_{ti}| \leq 2$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

α_t : α_1 α_2 α_3

?	?	?	?	?	?
?	?	?	?	?	?

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

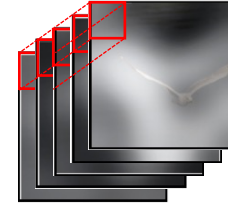
$$0 \leq |1 - \sum_t^C \alpha_{ti}| \leq 2$$

$$(1 - \sum_t^C \alpha_{ti})^2$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

α_t : α_1 α_2 α_3

?	?	?	?	?	?
?	?	?	?	?	?

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

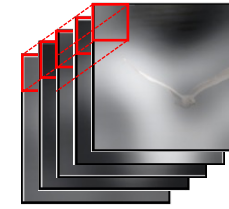
$$0 \leq |1 - \sum_t^C \alpha_{ti}| \leq 2$$

$$0 \leq (1 - \sum_t^C \alpha_{ti})^2 \leq 4$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

α_t :

α_1	α_2	α_3
0 1	0 1	0 1
0 0	0 0	0 0

$$1 - \sum_t^C \alpha_{ti}$$

1	-2
1	1

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

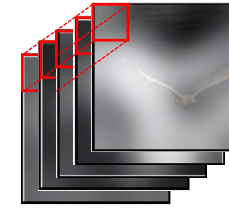
$$0 \leq |1 - \sum_t^C \alpha_{ti}| \leq 2$$

$$0 \leq (1 - \sum_t^C \alpha_{ti})^2 \leq 4$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

α_t : α_1 α_2 α_3

0	1	0	1	0	1
0	0	0	0	0	0

$$1 - \sum_t^C \alpha_{ti}$$

$$(1 - \sum_t^C \alpha_{ti})^2$$

1	-2
1	1

1	4
1	1

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

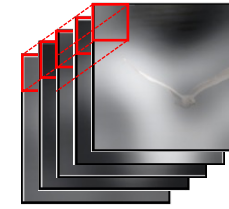
$$0 \leq |1 - \sum_t^C \alpha_{ti}| \leq 2$$

$$0 \leq (1 - \sum_t^C \alpha_{ti})^2 \leq 4$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{att}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Practice: What is the max value of $\sum_i^L (1 - \sum_t^C \alpha_{ti})^2$?

α_t :

α_1	α_2	α_3
0 1	0 1	0 1
0 0	0 0	0 0

$1 - \sum_t^C \alpha_{ti}$	$(1 - \sum_t^C \alpha_{ti})^2$
1 -2	1 4
1 1	1 1

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

$$0 \leq \sum_t^C \alpha_{ti} \leq 3$$

$$0 \leq |1 - \sum_t^C \alpha_{ti}| \leq 2$$

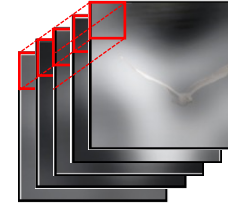
$$0 \leq (1 - \sum_t^C \alpha_{ti})^2 \leq 4$$

$$\sum_i^L (1 - \sum_t^C \alpha_{ti})^2 = 7$$

Loss Function

encourage $\sum_t \alpha_{ti} \approx 1$

$$L_d = -\log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_i^L (1 - \sum_t^C \alpha_{ti})^2$$



→ Sum over **time** = 1

$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

Softmax, so $\sum_i \alpha_{ti} = 1$



→ Sum over **positions** = 1

Why do we need the second term?

- Encourage the model to **pay equal attention** to **every part** of the image over the course of generation
- This penalty was important quantitatively to improving overall **BLEU score** and that qualitatively this leads to more **rich** and **descriptive** captions

Experimental Results

Dataset	Model	BLEU				METEOR
		BLEU-1	BLEU-2	BLEU-3	BLEU-4	
Flickr8k	Google NIC(Vinyals et al., 2014) ^{†Σ}	63	41	27	—	—
	Log Bilinear (Kiros et al., 2014a) [◦]	65.6	42.4	27.7	17.7	17.31
	Soft-Attention	67	44.8	29.9	19.5	18.93
	Hard-Attention	67	45.7	31.4	21.3	20.30
Flickr30k	Google NIC ^{†$\circ\Sigma$}	66.3	42.3	27.7	18.3	—
	Log Bilinear	60.0	38	25.4	17.1	16.88
	Soft-Attention	66.7	43.4	28.8	19.1	18.49
	Hard-Attention	66.9	43.9	29.6	19.9	18.46
COCO	CMU/MS Research (Chen & Zitnick, 2014) ^a	—	—	—	—	20.41
	MS Research (Fang et al., 2014) ^{†a}	—	—	—	—	20.71
	BRNN (Karpathy & Li, 2014) [◦]	64.2	45.1	30.4	20.3	—
	Google NIC ^{†$\circ\Sigma$}	66.6	46.1	32.9	24.6	—
	Log Bilinear [◦]	70.8	48.9	34.4	24.3	20.03
	Soft-Attention	70.7	49.2	34.4	24.3	23.90
	Hard-Attention	71.8	50.4	35.7	25.0	23.04

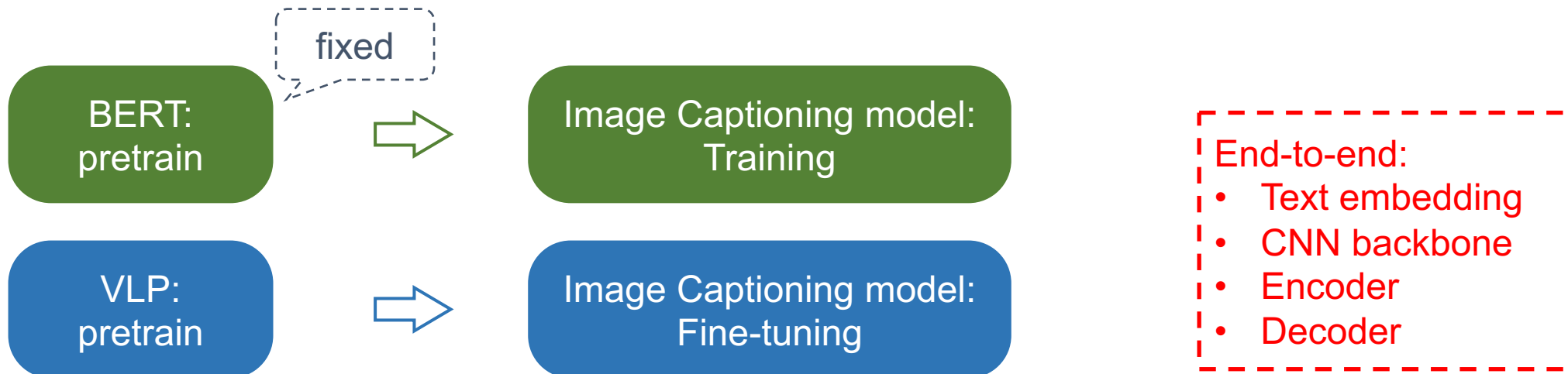
What's new in Image Captioning?



- Vision-Language Pre-training (VLP)
- Object Anchors
- Visual VOcabulary pre-training (VIVO)
- Generative Adversarial Networks (GAN)
- Deep Reinforcement Learning (RL) + Meta Learning

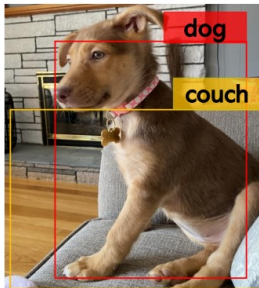
What's new in Image Captioning?

- Vision-Language Pre-training (VLP)
- Object Anchors
- Visual VOcabulary pre-training (VIVO)
- Generative Adversarial Networks (GAN)
- Deep Reinforcement Learning (RL) + Meta Learning



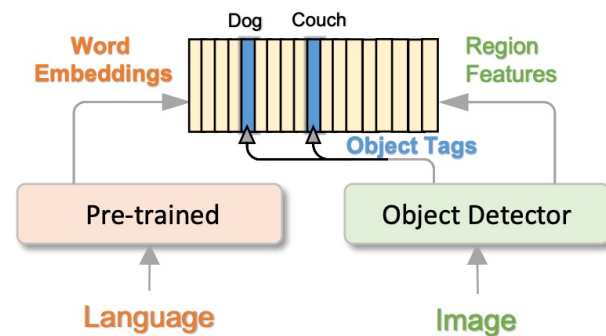
What's new in Image Captioning?

- Vision-Language Pre-training (VLP)
- Object Anchors
- Visual Vocabulary pre-training (VIVO)
- Generative Adversarial Networks (GAN)
- Deep Reinforcement Learning (RL) + Meta Learning

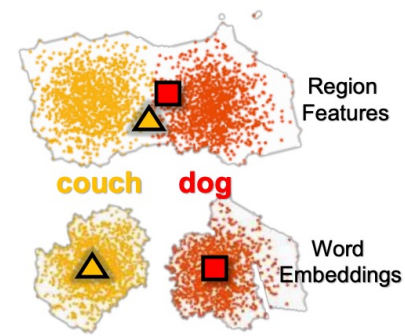


A **dog** is sitting on a **couch**

(a) Image-text pair



(b) Objects as anchor points

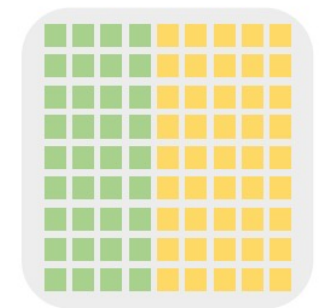
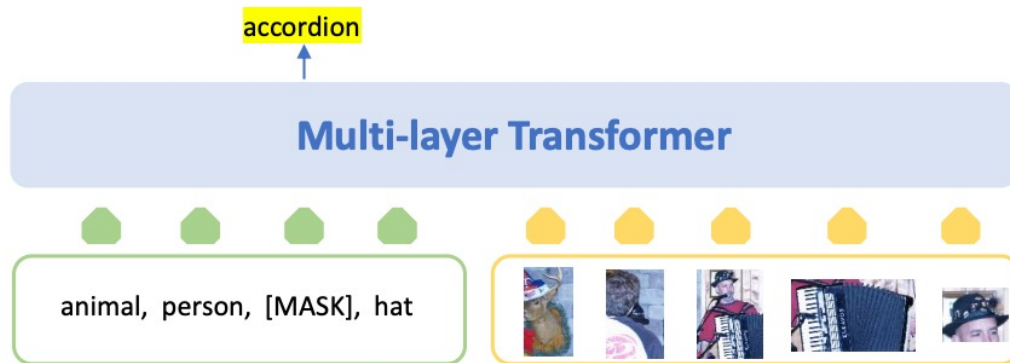
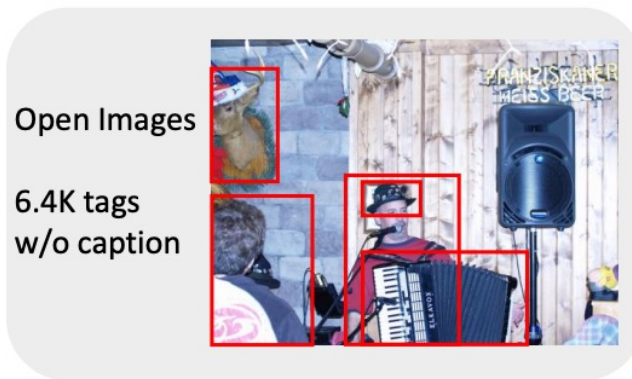


(c) Semantics spaces

What's new in Image Captioning?

- Vision-Language Pre-training (VLP)
- Object Anchors
- Visual VOcabulary pre-training (VIVO)
- Generative Adversarial Networks (GAN)
- Deep Reinforcement Learning (RL) + Meta Learning

Some challenges require models trained without other image captioning dataset



attention mask

What's new in Image Captioning?

- Vision-Language Pre-training (VLP)
- Object Anchors
- Visual VOcabulary pre-training (VIVO)
- **Generative Adversarial Networks (GAN)**
- Deep Reinforcement Learning (RL) + Meta Learning

What's new in Image Captioning?

- Vision-Language Pre-training (VLP)
- Object Anchors
- Visual VOcabulary pre-training (VIVO)
- Generative Adversarial Networks (GAN)
- Deep Reinforcement Learning (RL) + Meta Learning

Thank you!



References

- Papers

- <https://arxiv.org/abs/1411.4555>
- <https://arxiv.org/abs/1502.03044>

- Others

- <https://becominghuman.ai/only-numpy-deriving-partial-forward-feed-lstm-on-show-attend-and-tell-neural-image-caption-4e44aa2b966d>
- <https://www.youtube.com/watch?v=y1S3Ri7myMg>
- <https://cv-tricks.com/artificial-intelligence/show-attend-tell-image-captioning-explained/>

Image Sources

- <https://www.sandiegouniontribune.com/business/real-estate/story/2020-08-21/will-san-diego-stay-at-home-workers-leave-if-given-the-opportunity>
- <https://medium.com/swlh/image-caption-generation-with-visual-attention-c782dfc0634b>
- <https://towardsdatascience.com/time-series-forecasting-with-deep-stacked-unidirectional-and-bidirectional-lstms-de7c099bd918>
- <https://www.usna.edu/Users/cs/nchamber/courses/nlp/f20/labs/lab5/index.html>
- <https://medium.com/intelligentmachines/word-embedding-and-one-hot-encoding-ad17b4bbe111>
- <https://medium.com/mlearning-ai/self-attention-in-convolutional-neural-networks-172d947afc00>
- https://www.ecva.net/papers/eccv_2020/papers_ECCV/papers/123750120.pdf
- <https://arxiv.org/pdf/2009.13682.pdf>
- <https://www.animatorisland.com/story-101-why-do-we-tell-stories/>