

Attention Models

Jaegul Choo (주재걸)
Korea University

<https://sites.google.com/site/jaegulchoo/>

Slide made by Cheonbok Park

Contents

1. Recurrent Neural Network

1-1. RNN의 여러가지 형태

1-2. Character-level LM

1-3. Vanishing Gradient Problem

2. LSTM과 GRU

2-1. LSTM이란 무엇인가

2-2. GRU란 무엇인가

Image Captioning

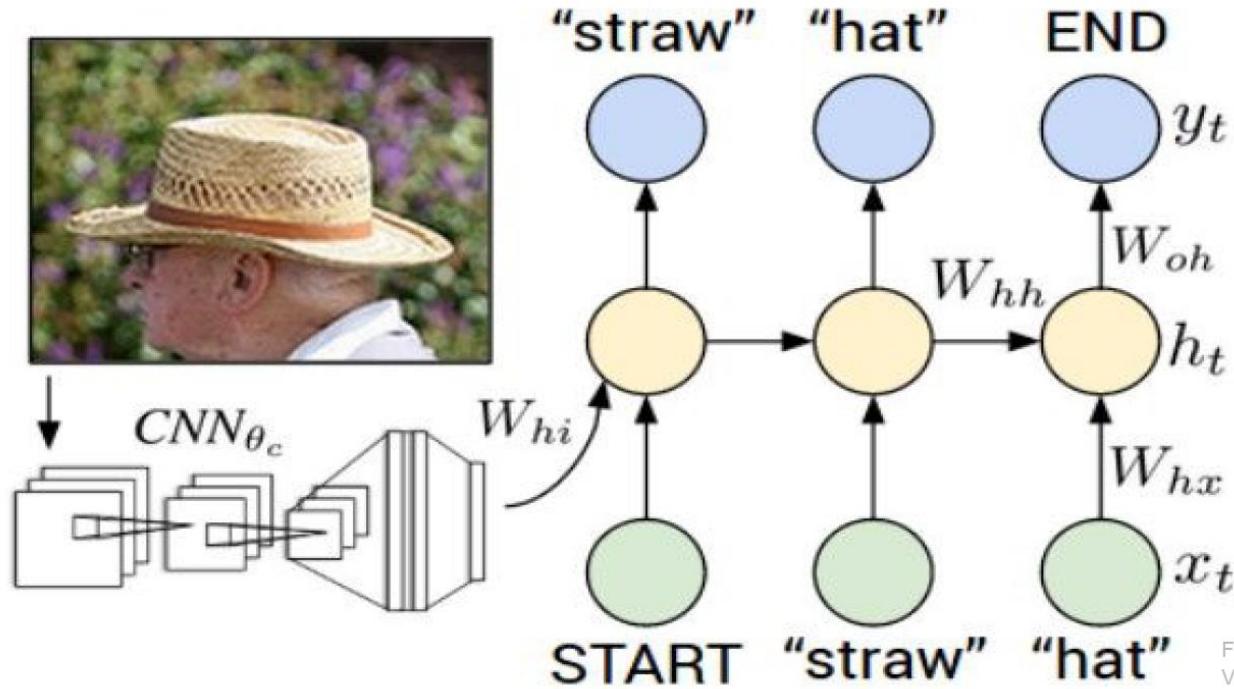


Figure from Karpathy et al, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015; figure copyright IEEE, 2015.
Reproduced for educational purposes.

Explain Images with Multimodal Recurrent Neural Networks, Mao et al.

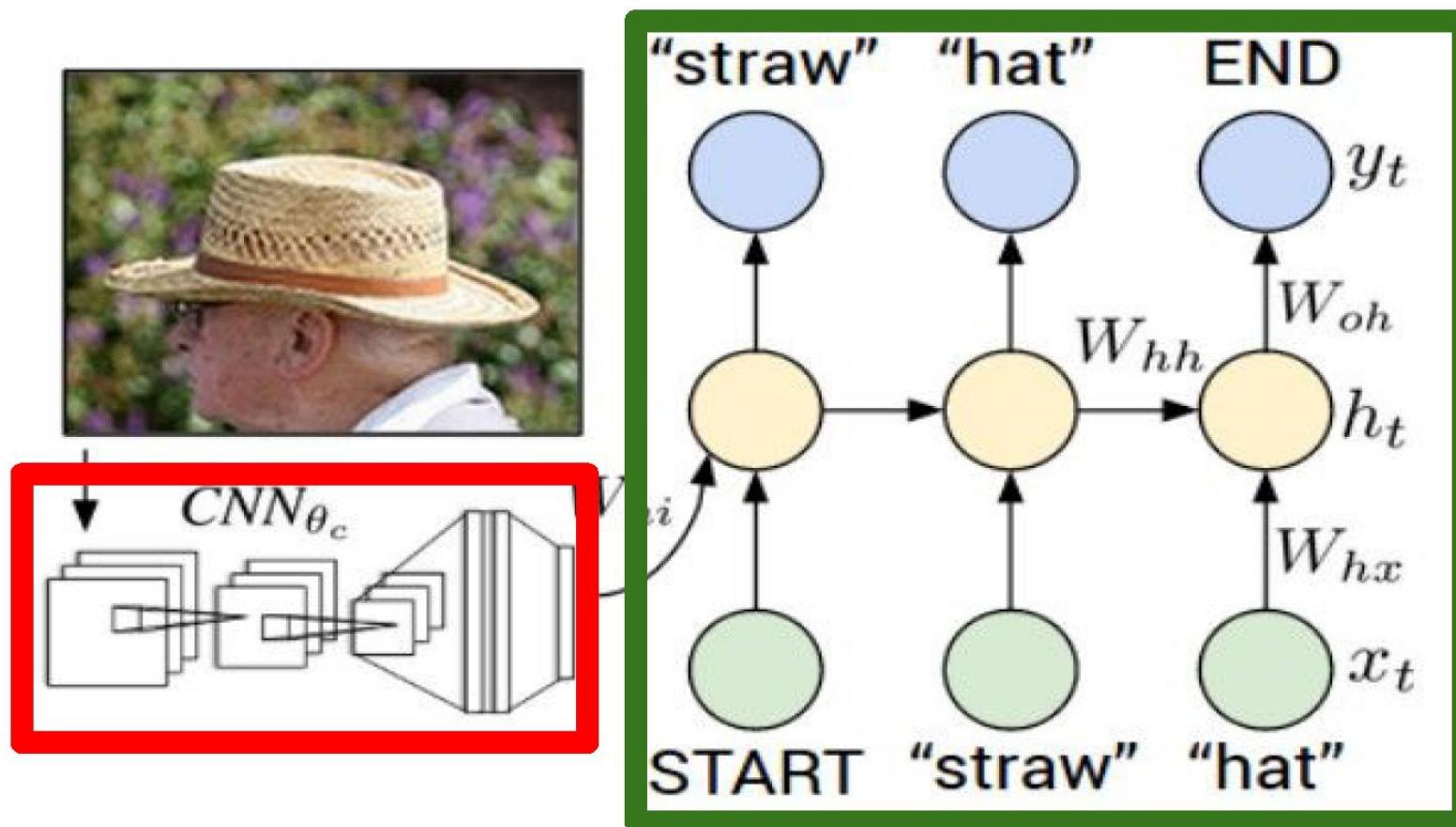
Deep Visual-Semantic Alignments for Generating Image Descriptions, Karpathy and Fei-Fei

Show and Tell: A Neural Image Caption Generator, Vinyals et al.

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

Learning a Recurrent Visual Representation for Image Caption Generation, Chen and Zitnick

Recurrent Neural Network



Convolutional Neural Network

image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

FC-4096

FC-4096

FC-1000

softmax

image



test image



conv-64

conv-64

maxpool

conv-128

conv-128

maxpool

conv-256

conv-256

maxpool

conv-512

conv-512

maxpool

conv-512

conv-512

maxpool

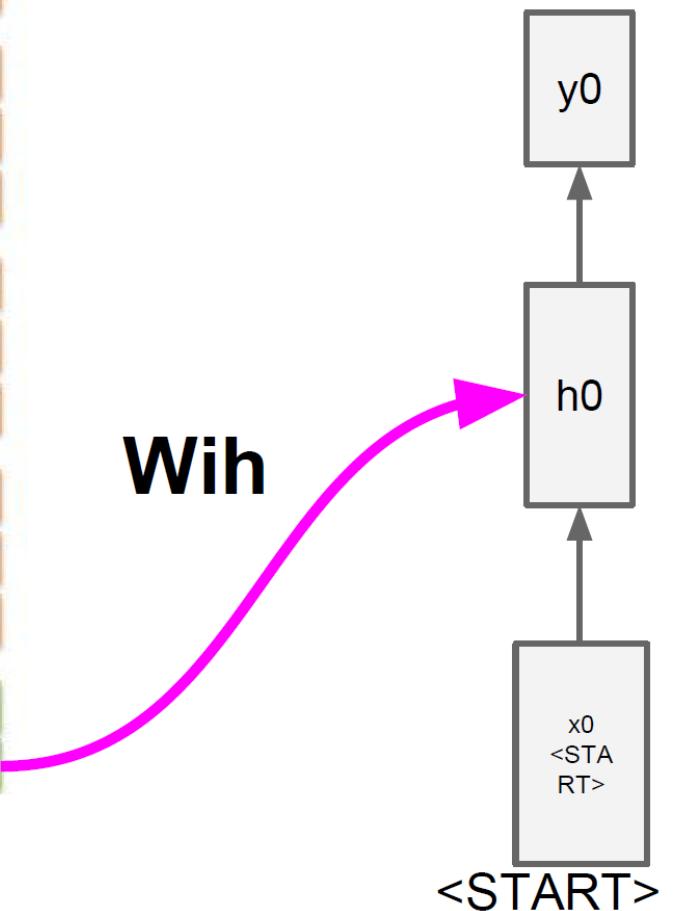
FC-4096

FC-4096

FC-1000

softmax

X



before:

$$h = \tanh(W_{xh} * x + W_{hh} * h)$$

now:

$$h = \tanh(W_{xh} * x + W_{hh} * h + W_{ih} * v)$$

Image Captioning: Example Results

Captions generated using [neuraltalk2](#)
All images are CC0 Public domain:
[cat suitcase](#), [cat tree](#), [dog](#), [bear](#),
[surfers](#), [tennis](#), [giraffe](#), [motorcycle](#)



A cat sitting on a suitcase on the floor



A cat is sitting on a tree branch



A dog is running in the grass with a frisbee



A white teddy bear sitting in the grass



Two people walking on the beach with surfboards



A tennis player in action on the court



Two giraffes standing in a grassy field



A man riding a dirt bike on a dirt track

Image Captioning: Failure Cases

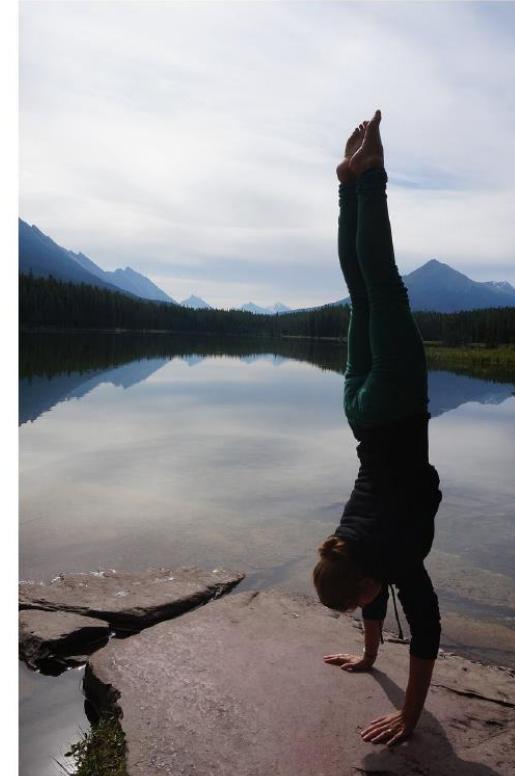
Captions generated using [neuraltalk2](#)
All images are CC0 Public domain: [fur coat](#), [handstand](#), [spider web](#), [baseball](#)



A woman is holding a cat in her hand



A person holding a computer mouse on a desk



A woman standing on a beach holding a surfboard



A bird is perched on a tree branch



A man in a baseball uniform throwing a ball

Image Captioning with Attention

RNN focuses its attention at a different spatial location when generating each word

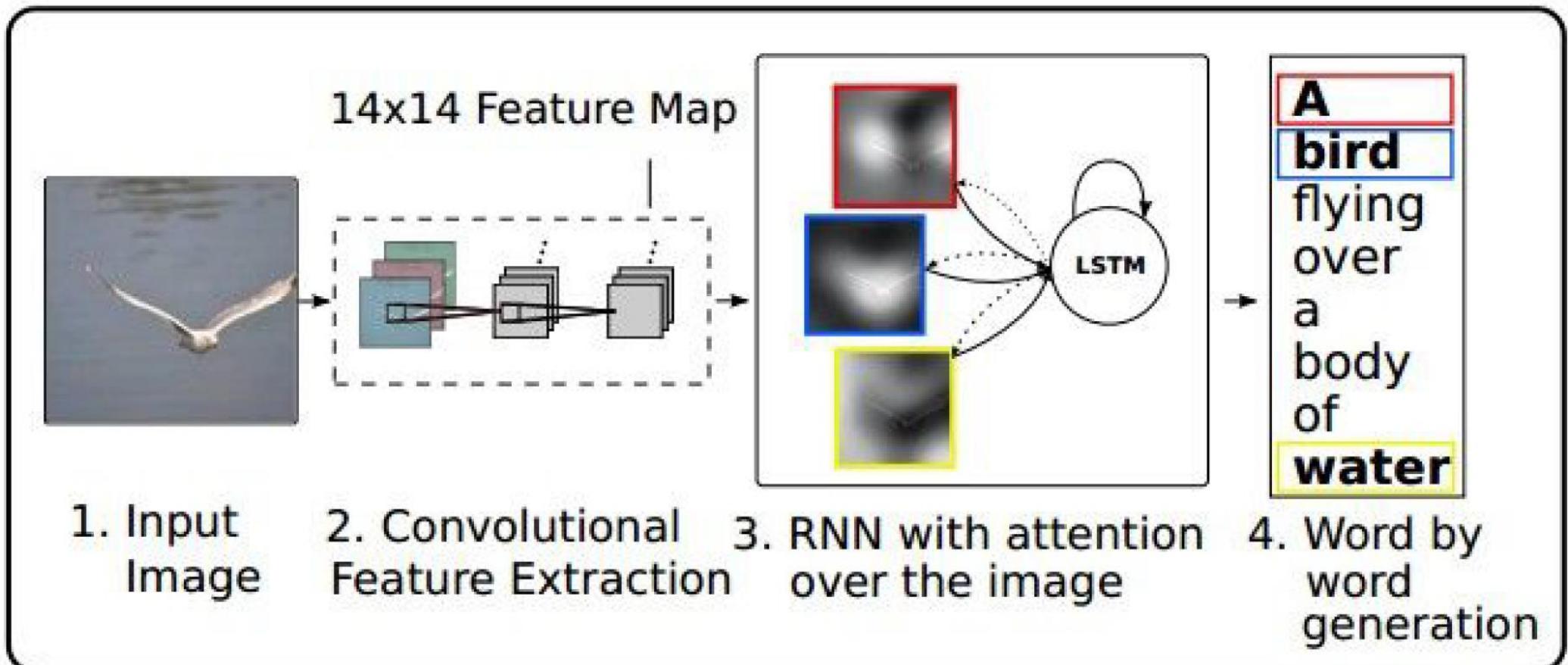


Image Captioning

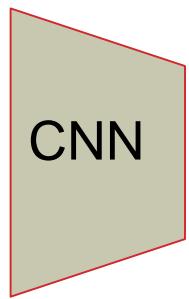


Image:
 $H \times W \times 3$

Image Captioning



Image:
 $H \times W \times 3$



Features:
 D



Image Captioning

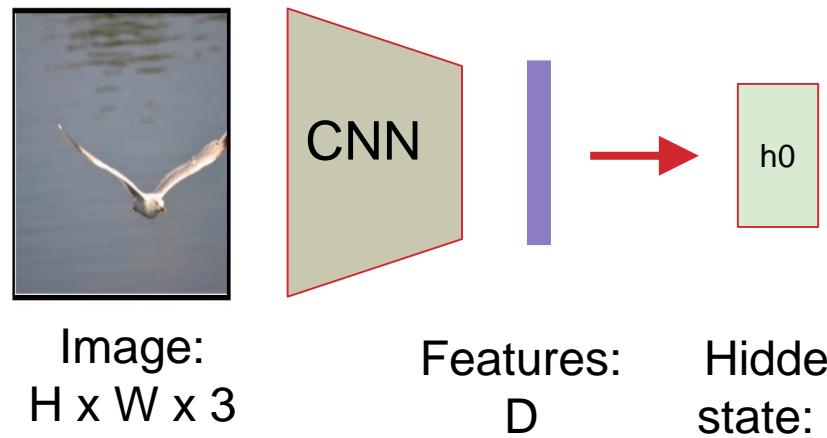


Image Captioning

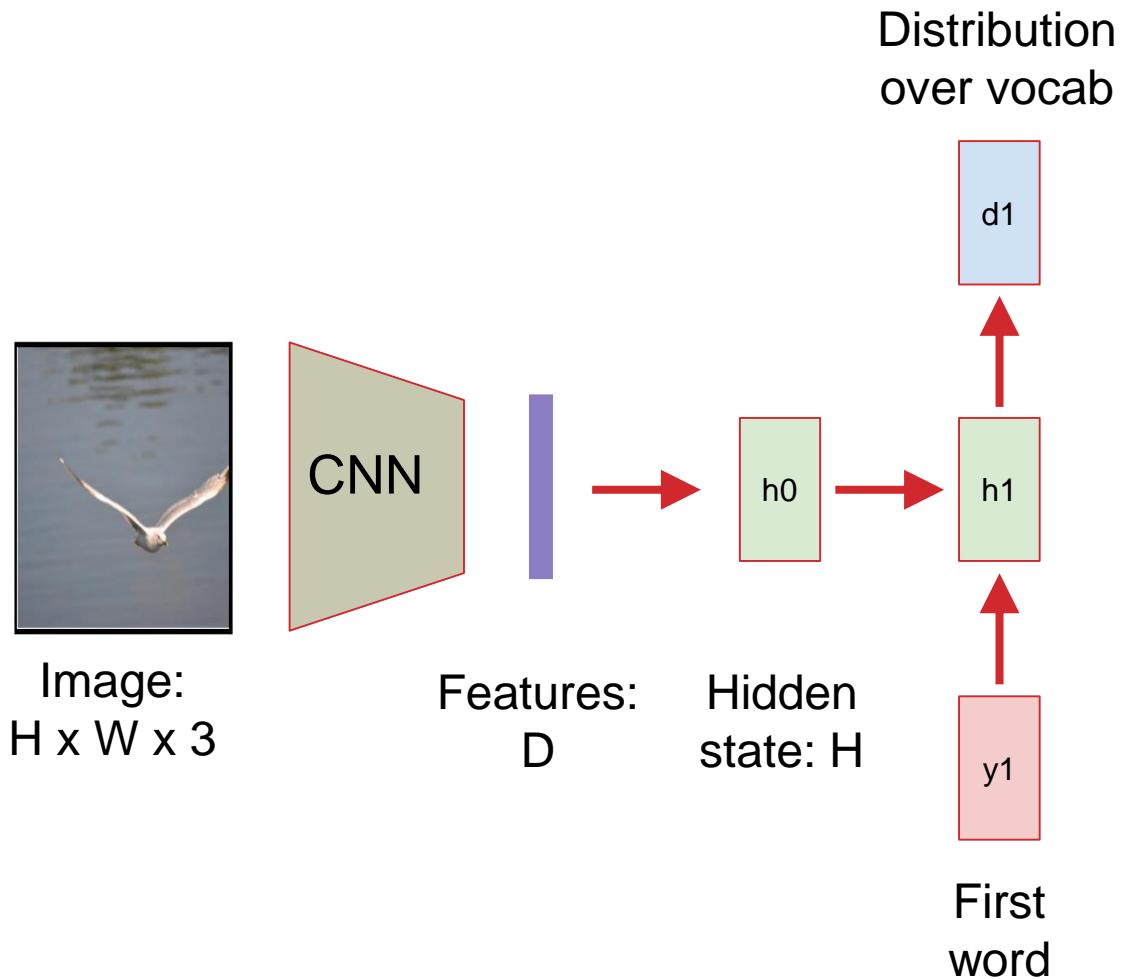


Image Captioning

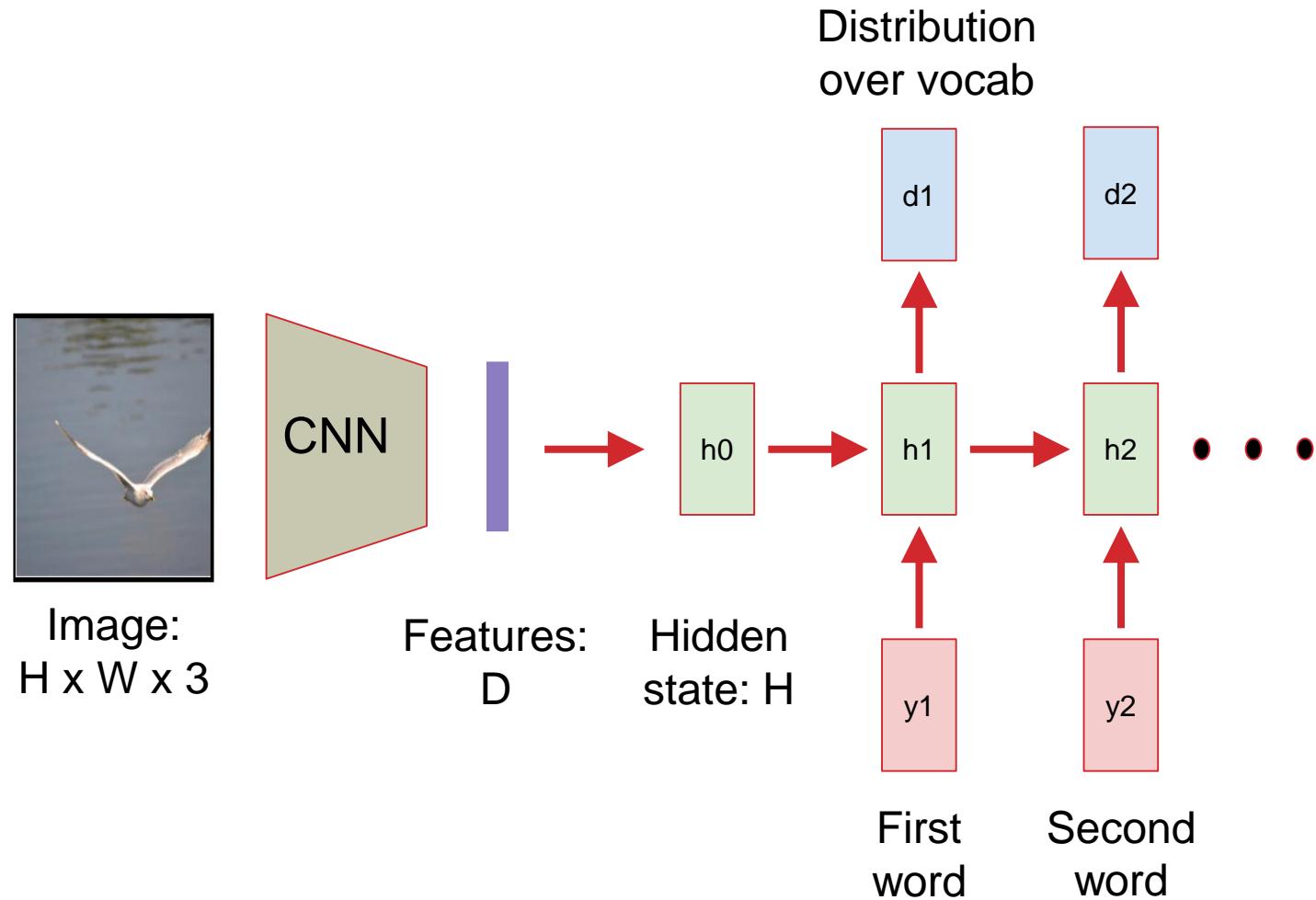
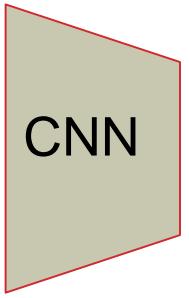


Image Captioning

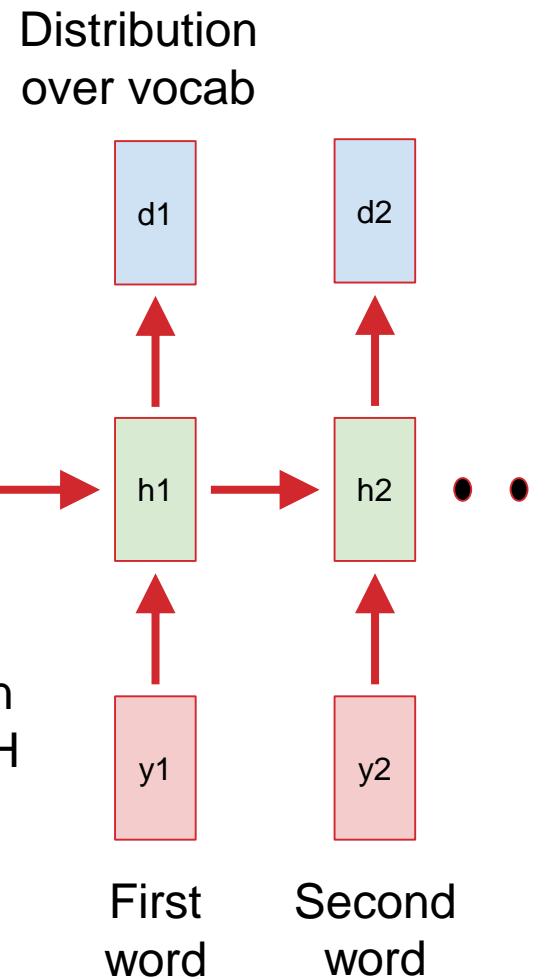


Image:
 $H \times W \times 3$



Features:
 D

Hidden
state: H

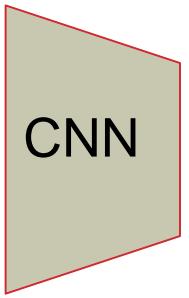


RNN only looks
at whole image,
once

Image Captioning

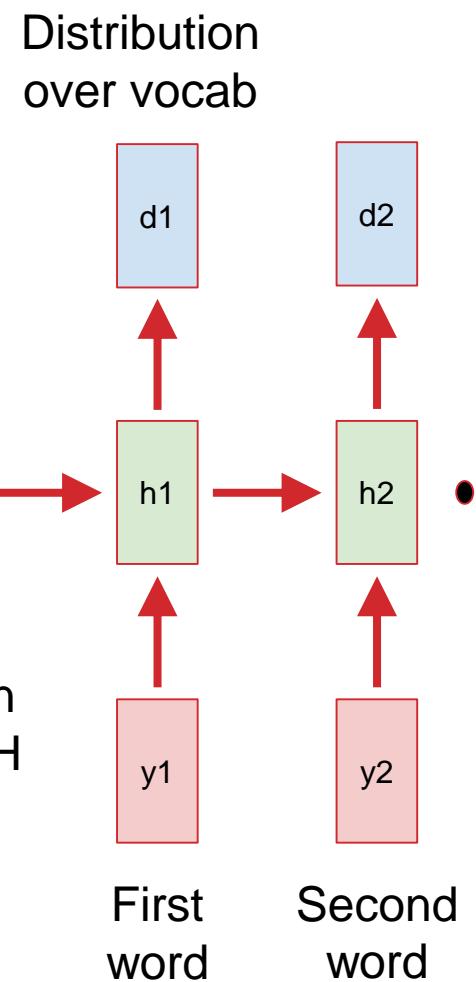


Image:
 $H \times W \times 3$



Features:
 D

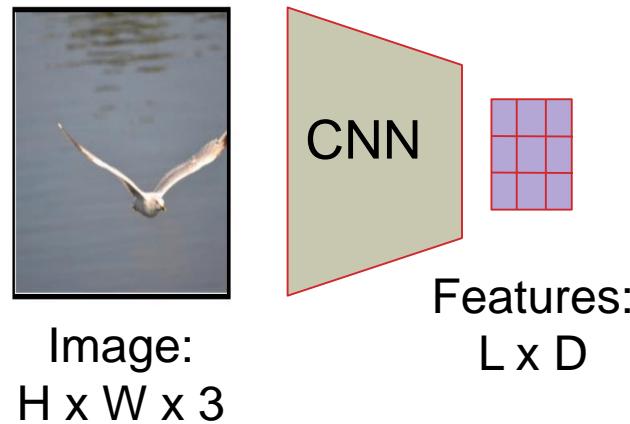
Hidden
state: H



RNN only looks
at whole image,
once

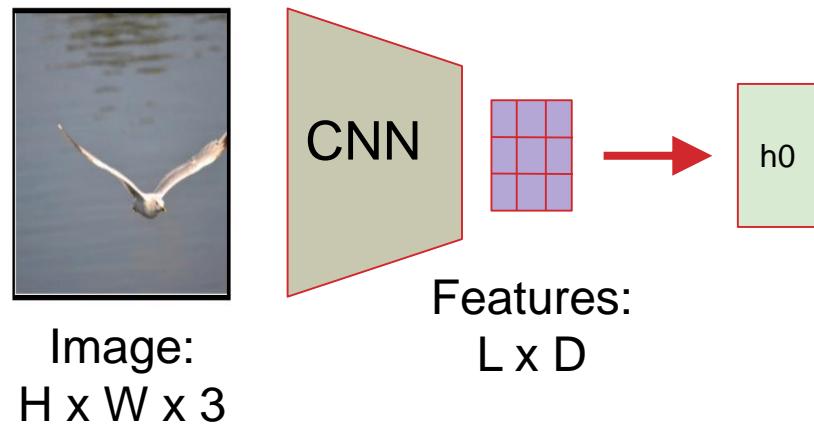
What if the
RNN looks at
different parts
of the image at
each timestep?

Image Captioning with Attention



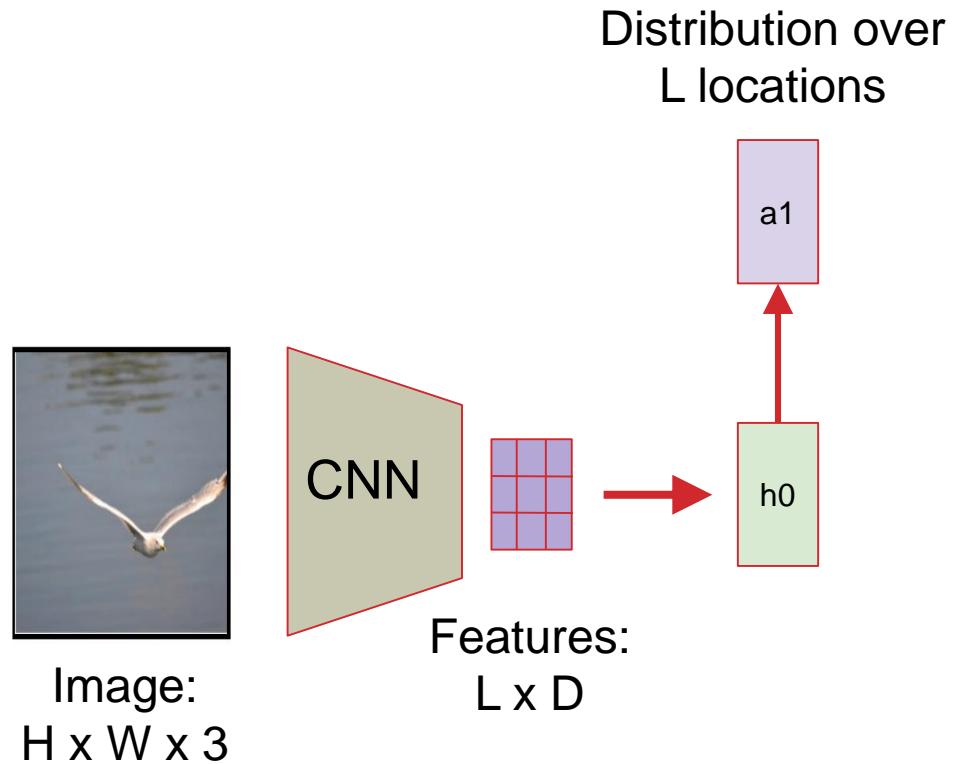
Xu et al, “Show, Attend and Tell:
Neural Image Caption Generation
with Visual Attention”, ICML 2015

Image Captioning with Attention



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

Image Captioning with Attention



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

Image Captioning with Attention

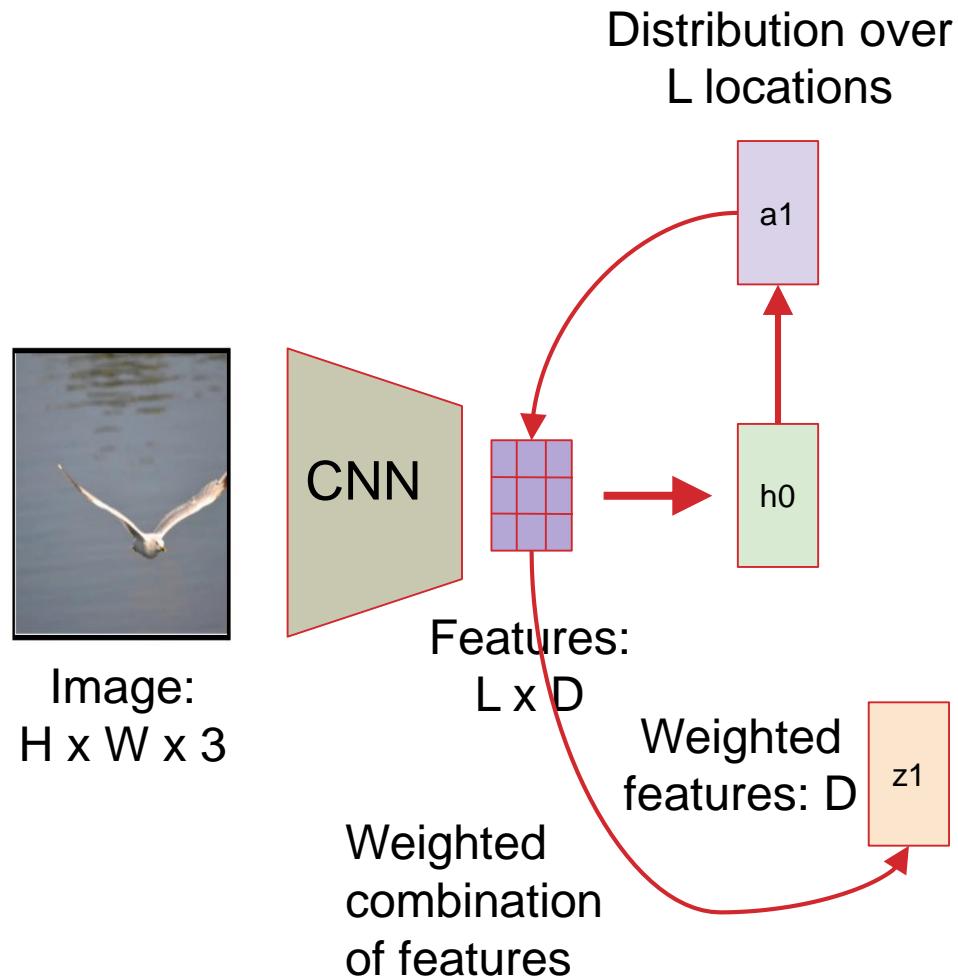


Image Captioning with Attention

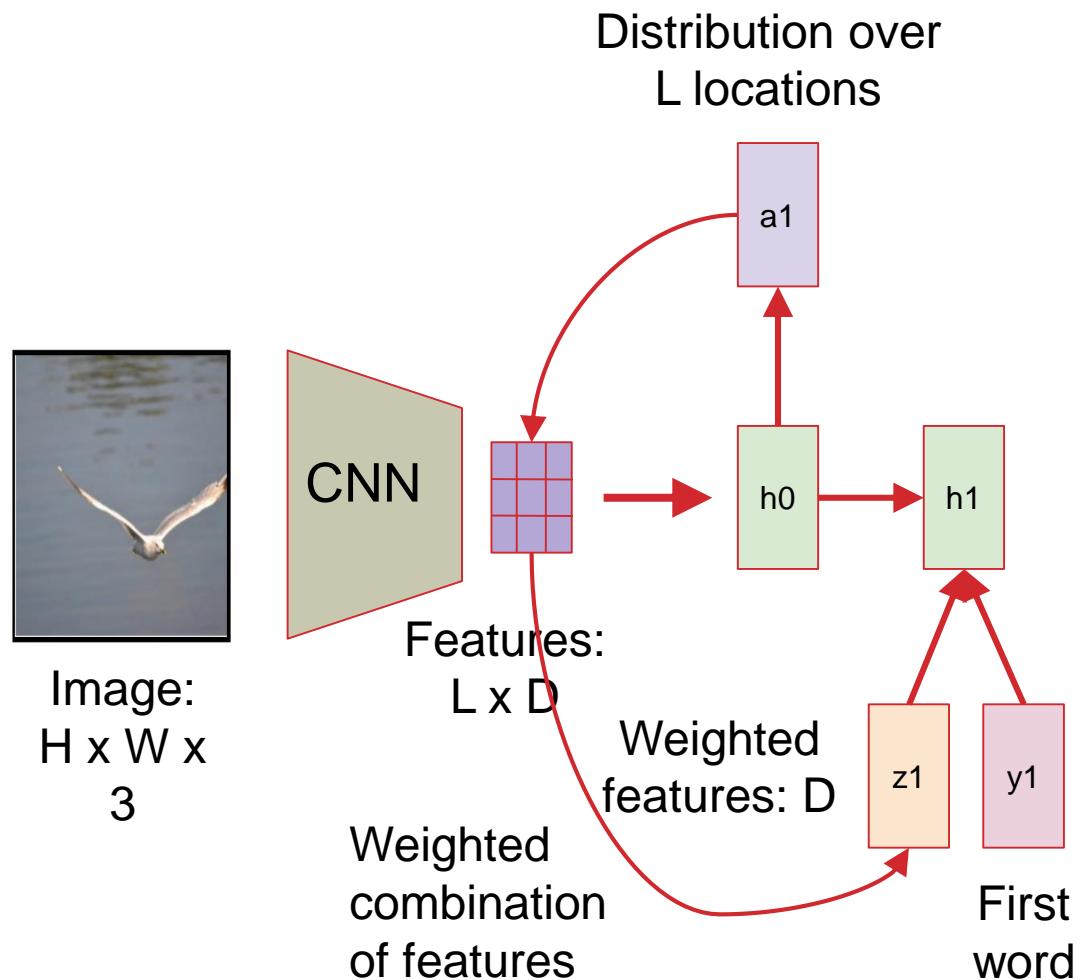


Image Captioning with Attention

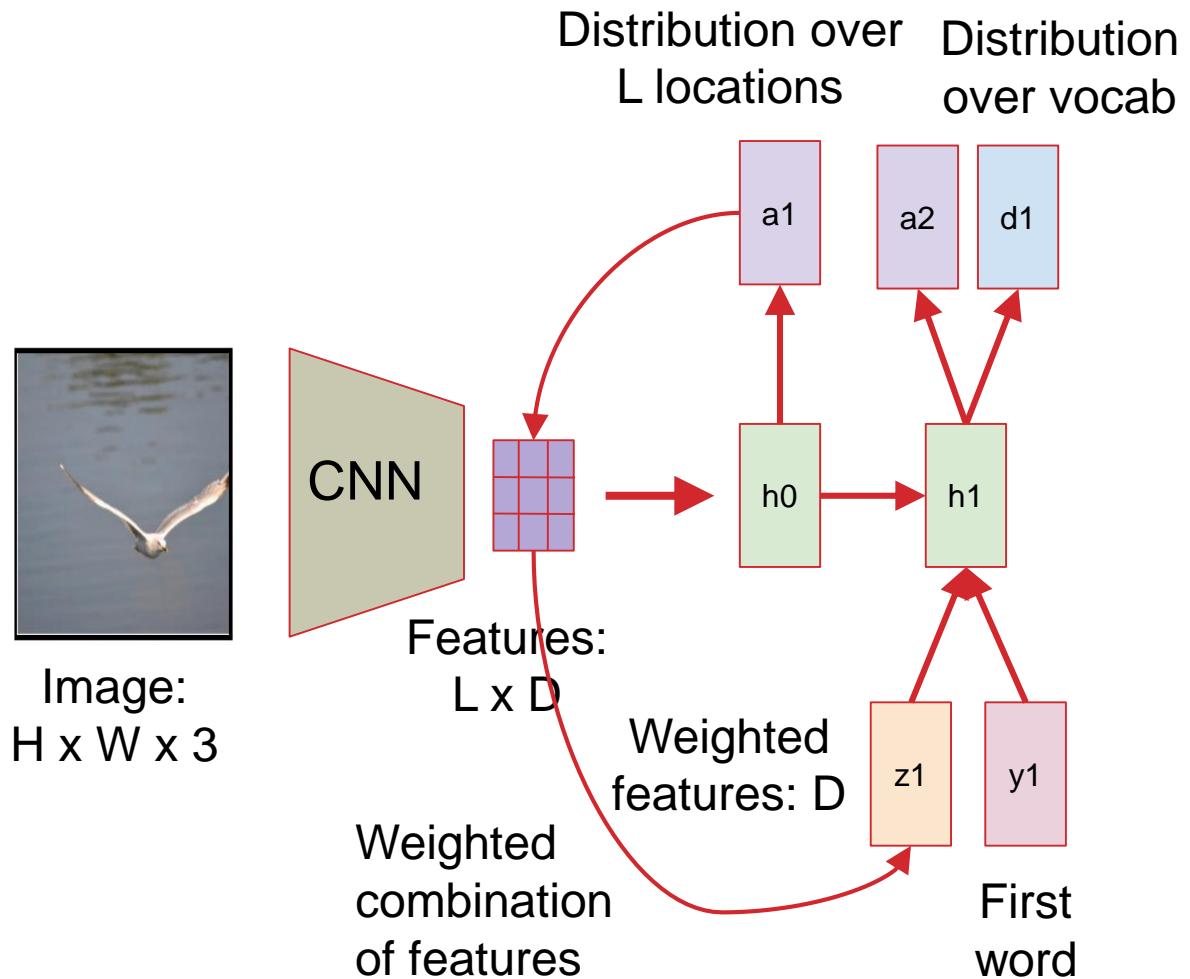


Image Captioning with Attention

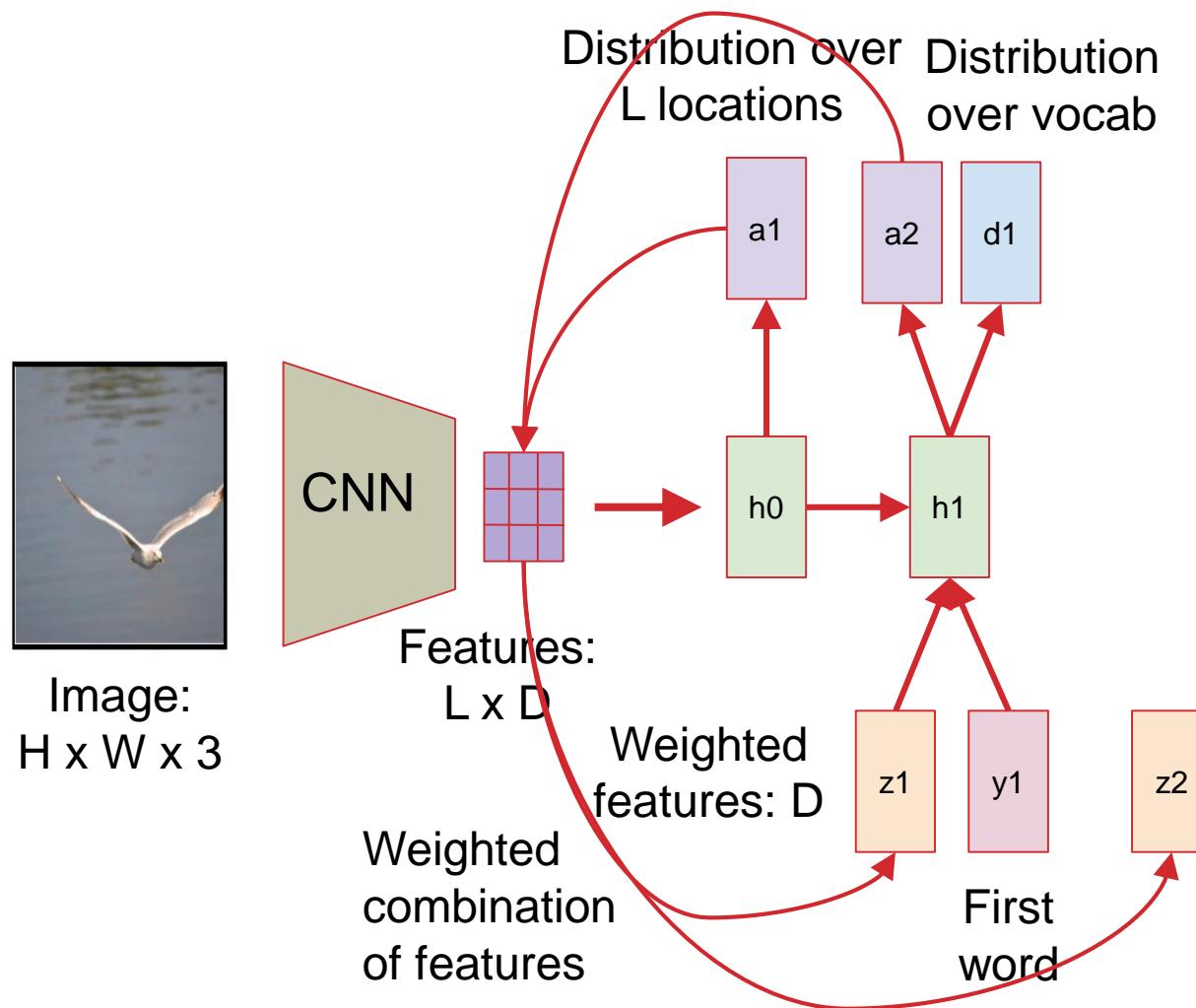


Image Captioning with Attention

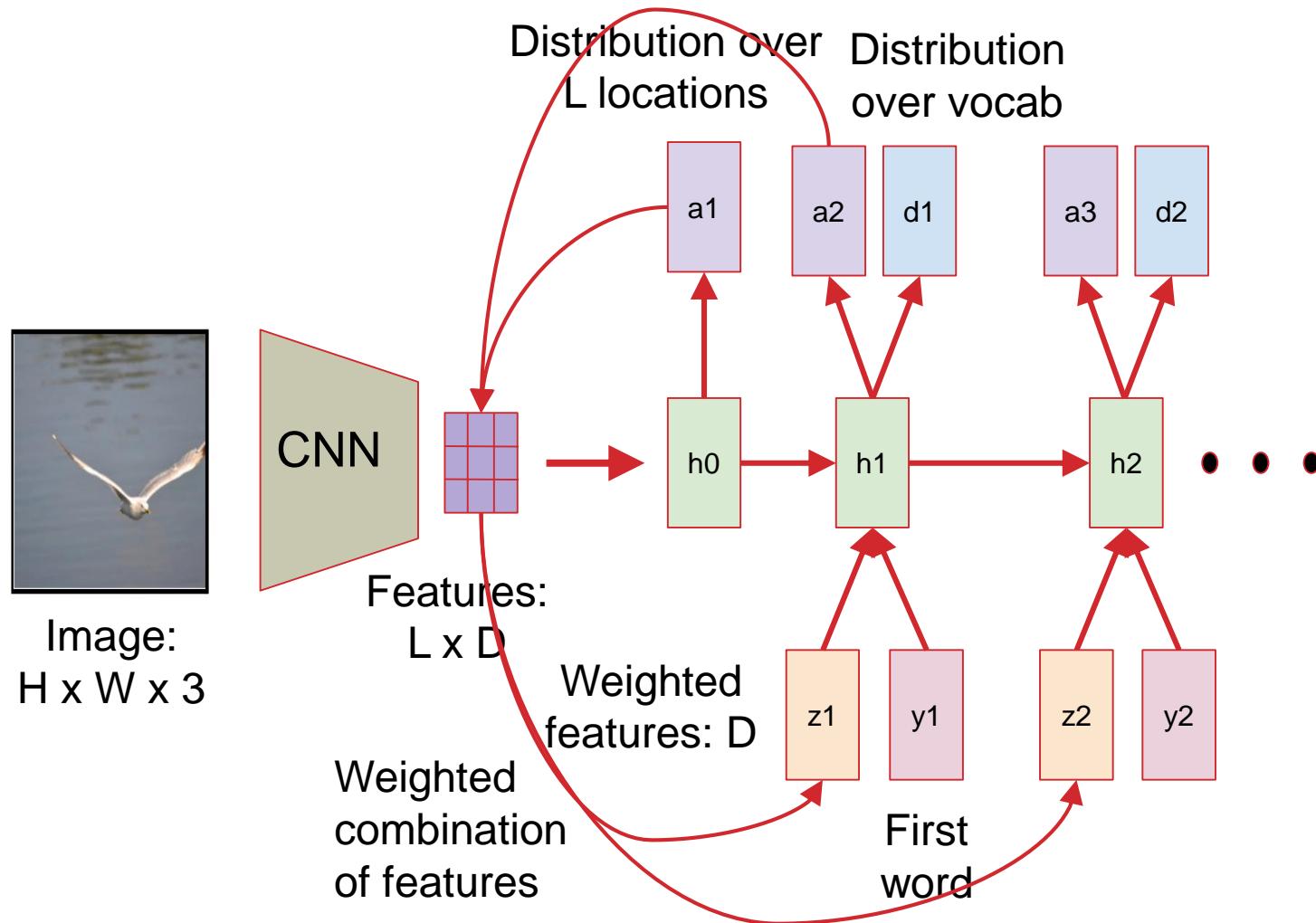


Image Captioning with Attention

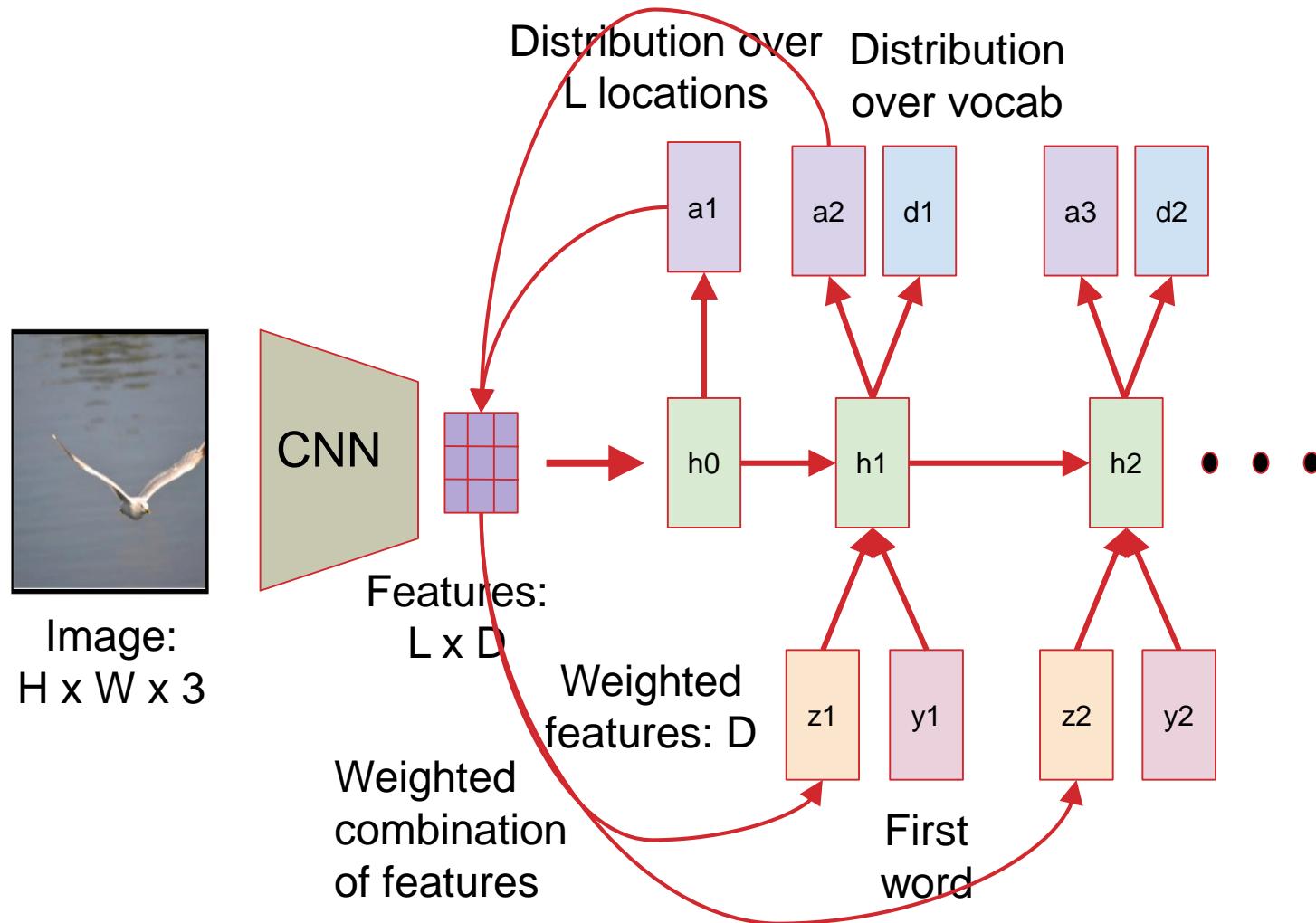


Image Captioning with Attention

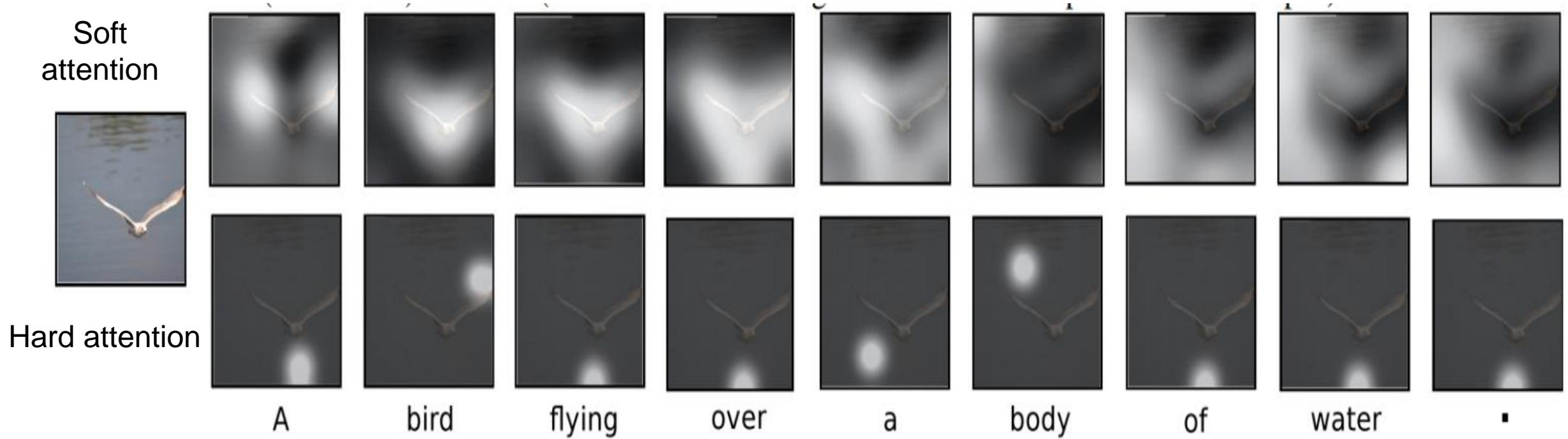


Image Captioning with Attention

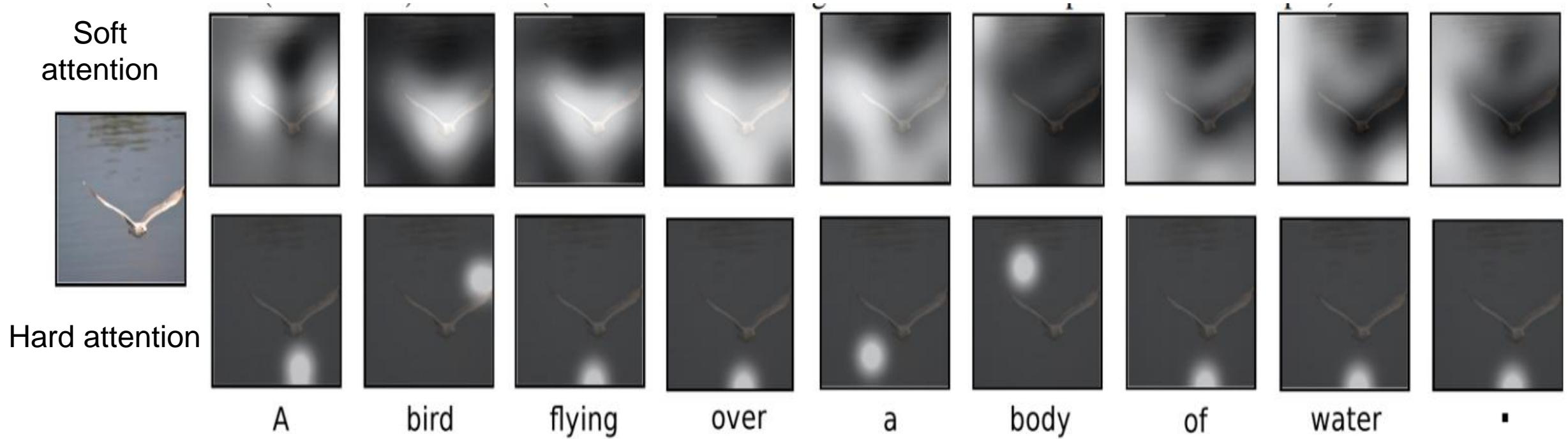


Image Captioning with Attention



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



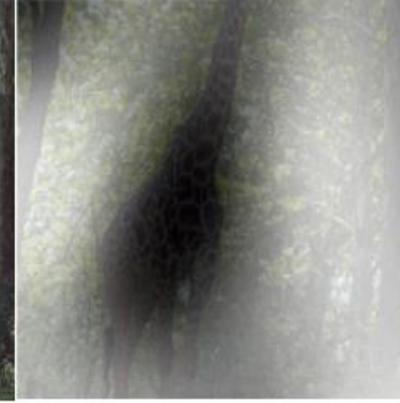
A stop sign is on a road with a mountain in the background.



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



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



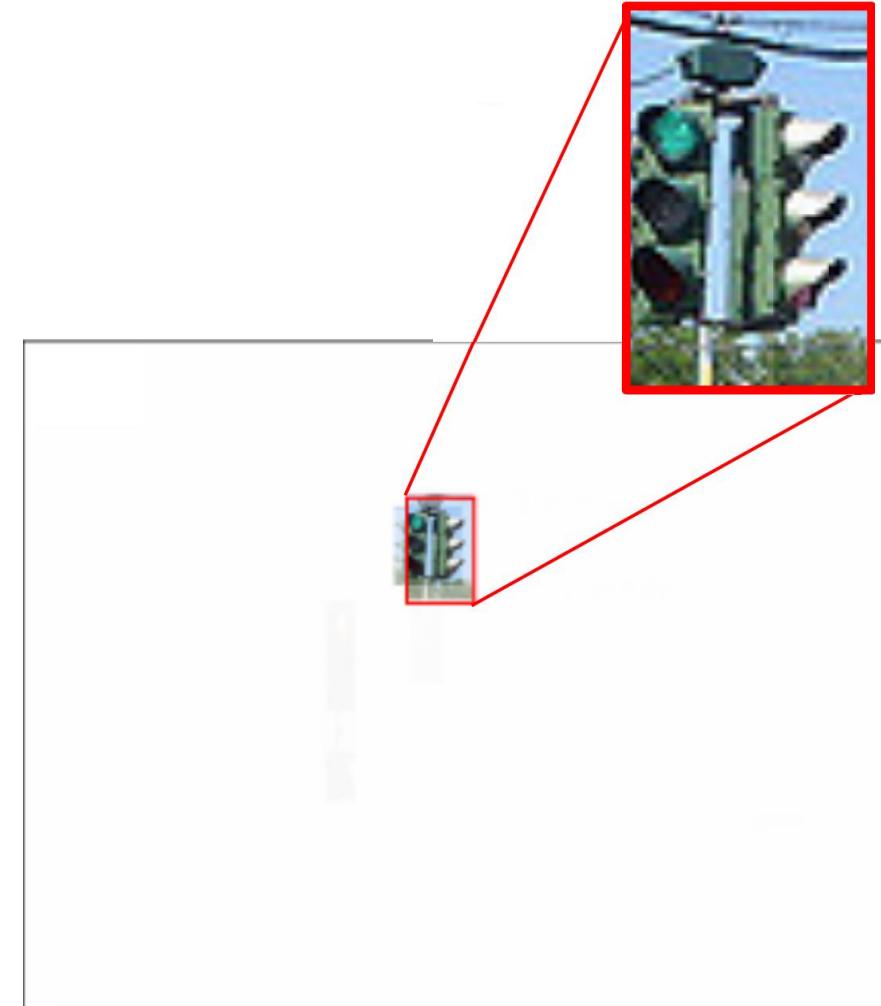
A giraffe standing in a forest with trees in the background.

RECAP, Visual Attention

- Vision and language tasks often require fine grained visual processing, e.g. VQA:

Q: What color is illuminated on the traffic light ?

A: green

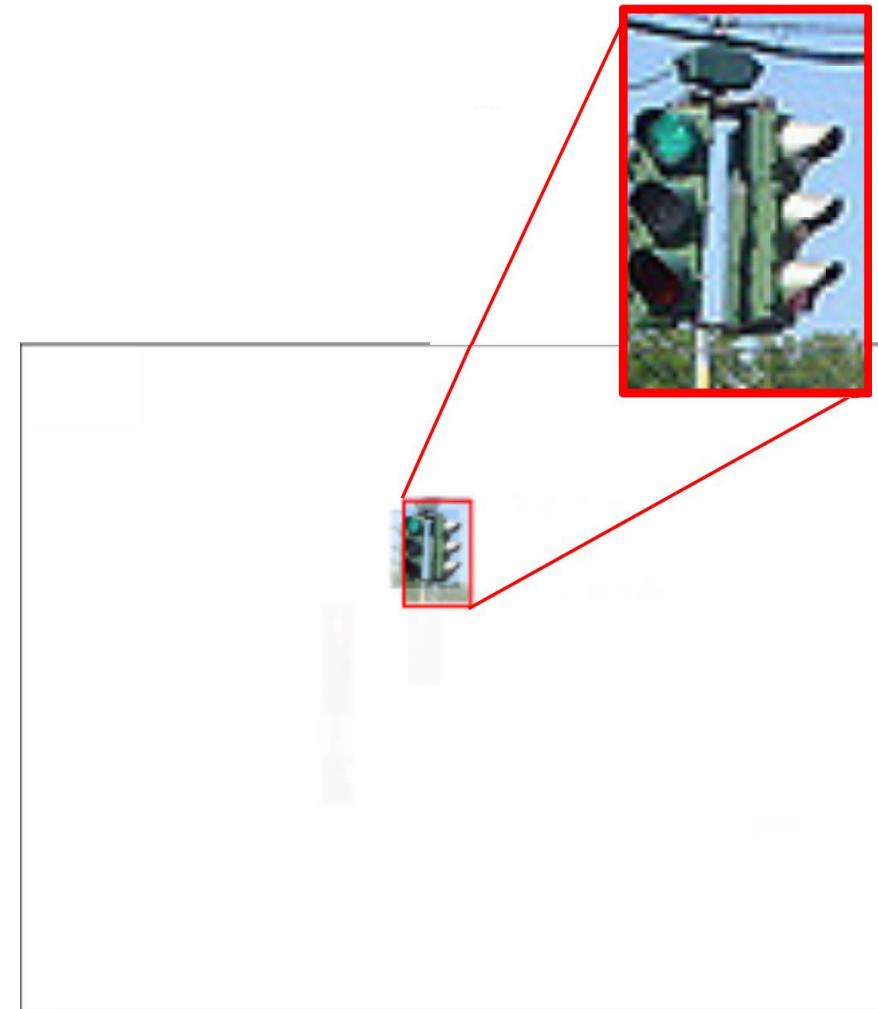


RECAP, Visual Attention

- Visual attention mechanisms learn to focus on image regions that are relevant to the task

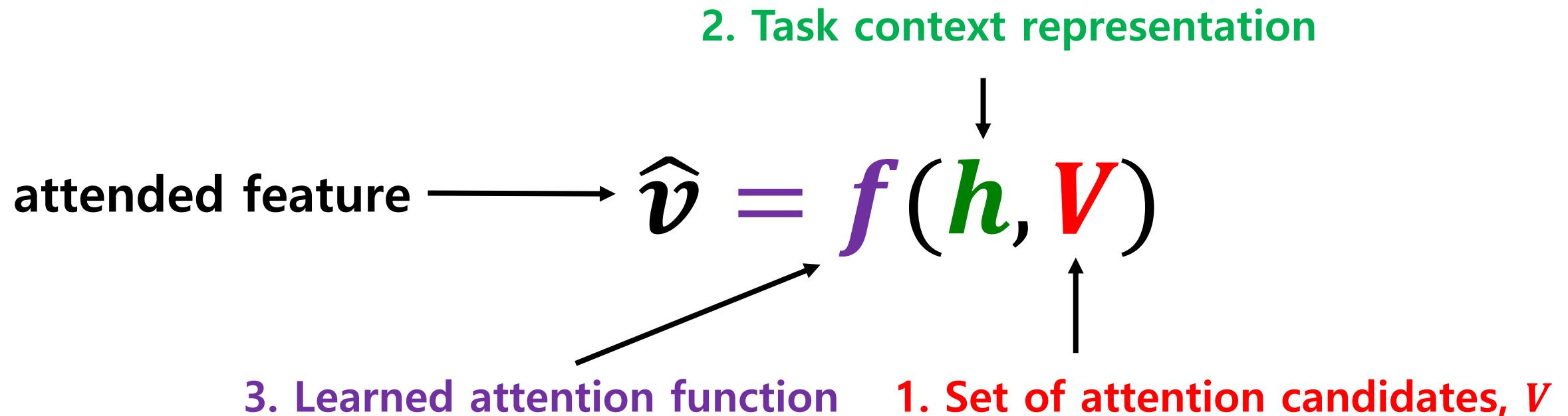
Q: What color is illuminated on the traffic light ?

A: green



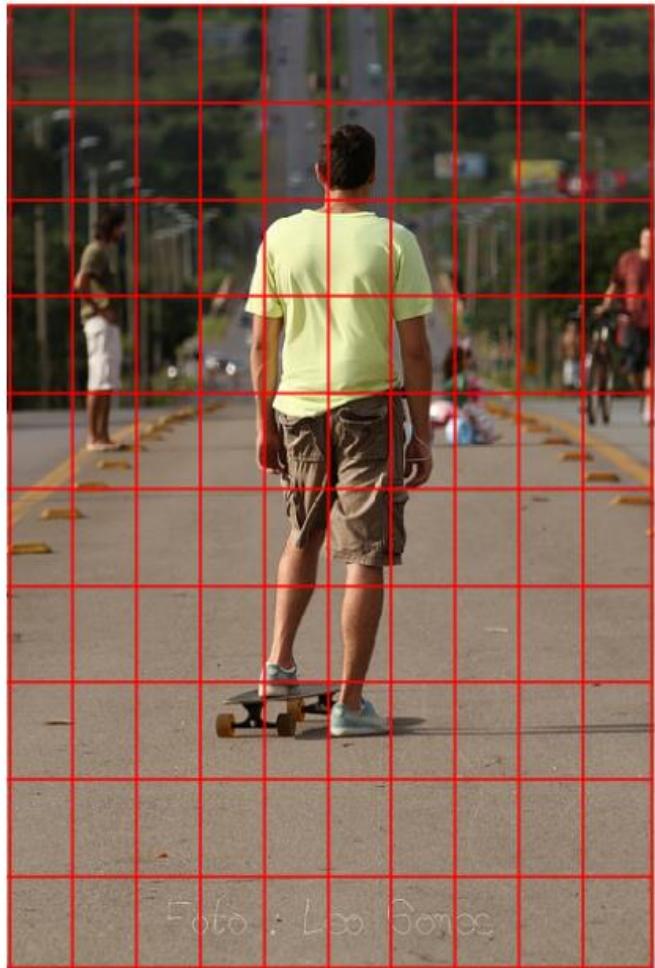
Components of visual attention

- Visual attention mechanisms **learn to focus** on **image regions** that are **relevant to the task**



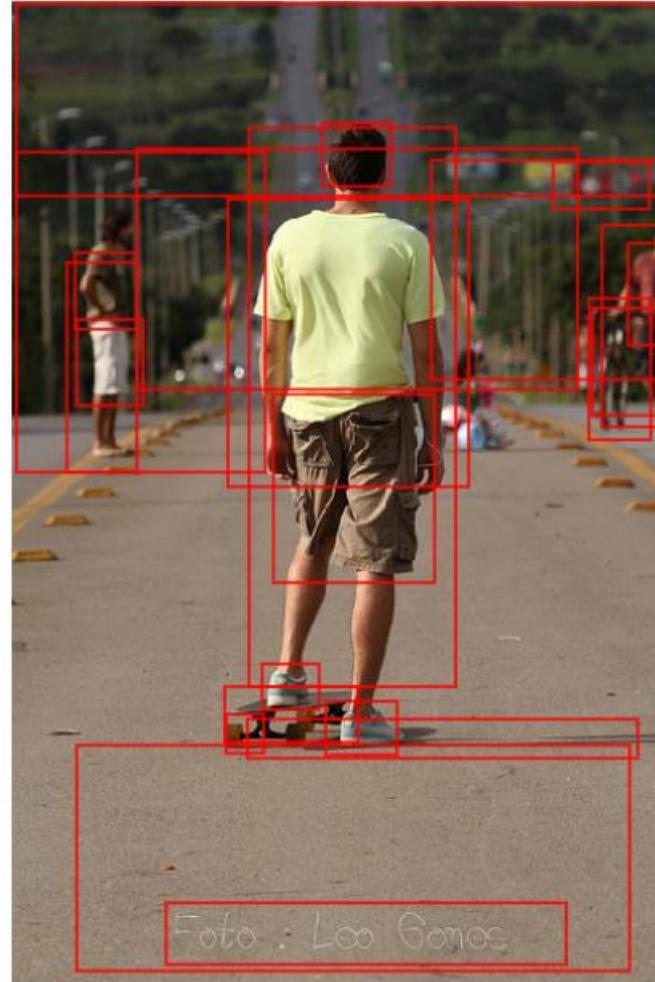
Attention Candidates , V

10



10

Standard approach:
Spatial output of a CNN

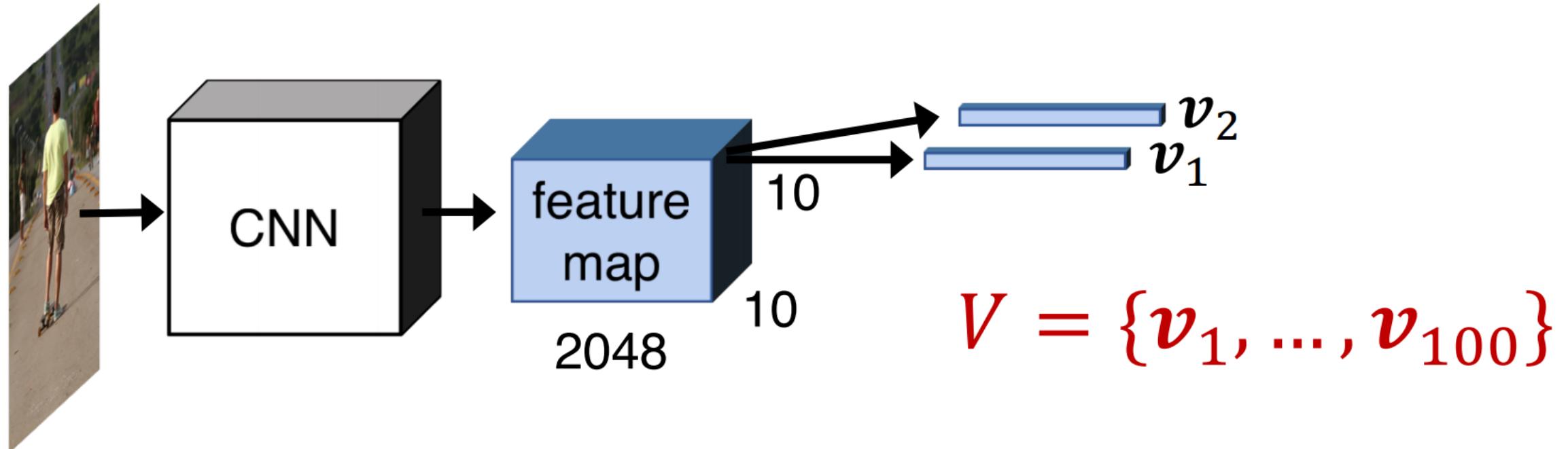


k regions

Bottom-up attention approach:
Object-based attention

Attention Candidates , V

Standard approach: use the spatial output of a CNN to extract vectors for each position in a grid



Objects are a natural basis for attention

- Human visual attention can select discrete objects, not just spatial regions
- Image captioning and VQA are concerned with objects

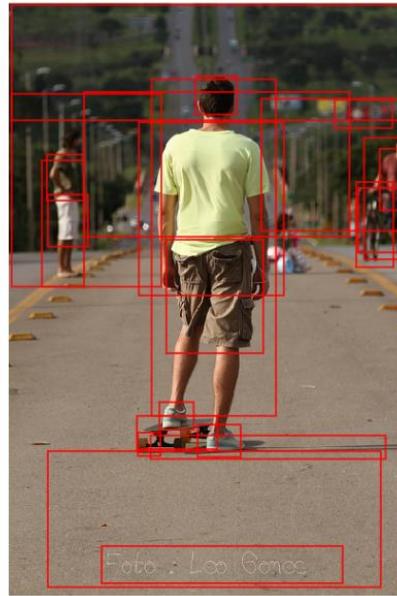


A **young man** on a **skateboard** looking down **street** with **people** watching.

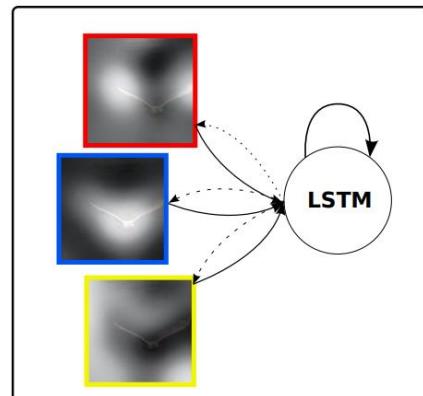
Q: Is the **boy** in the **yellow shirt** wearing **head protective gear**? **A:** No

Bottom-up and top-down attention

V



Bottom-up process: Extract all objects and other salient regions from the image (independent of the question / partially-completed caption)

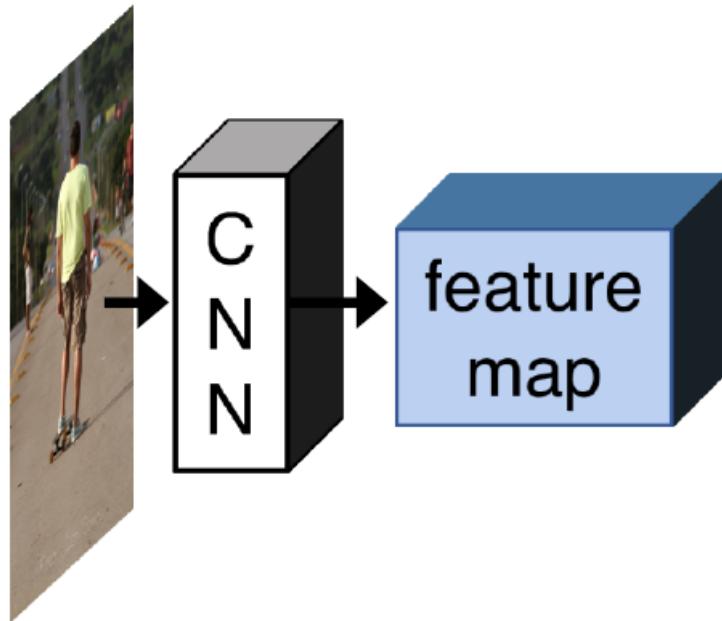


i. RNN with attention over the image

Top-down process: Given task context, weight the attention candidates (i.e., use existing VQA / captioning models)

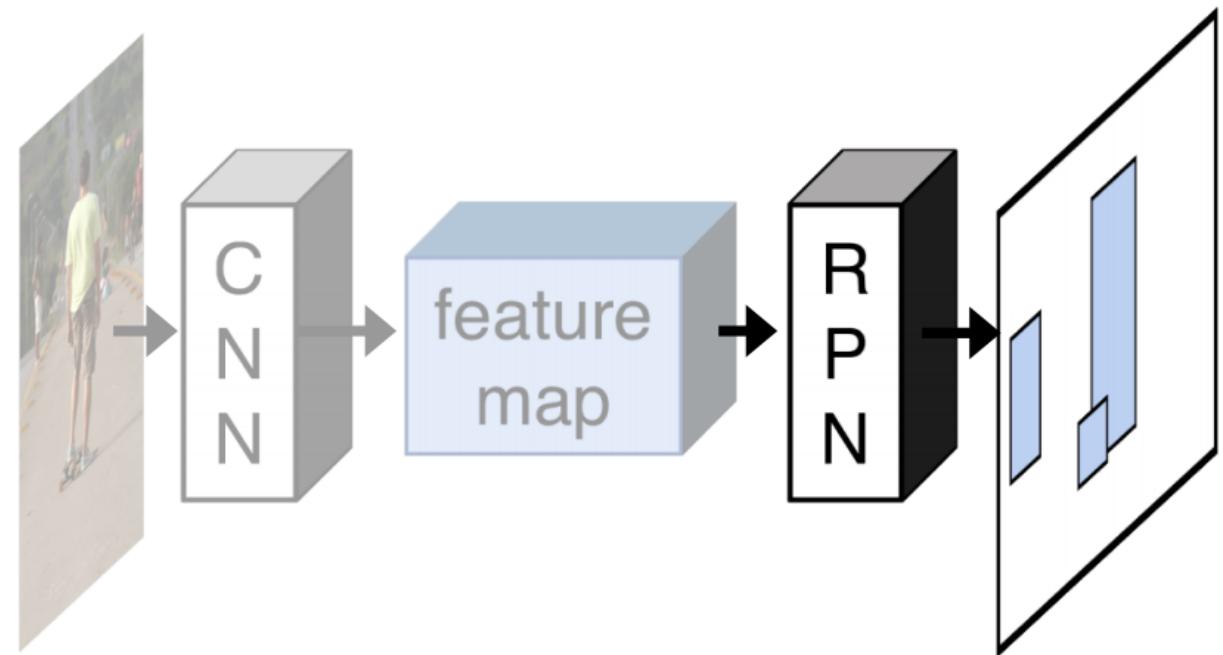
Attention Candidates , V

Bottom-up attention (using Faster R-CNN₂)



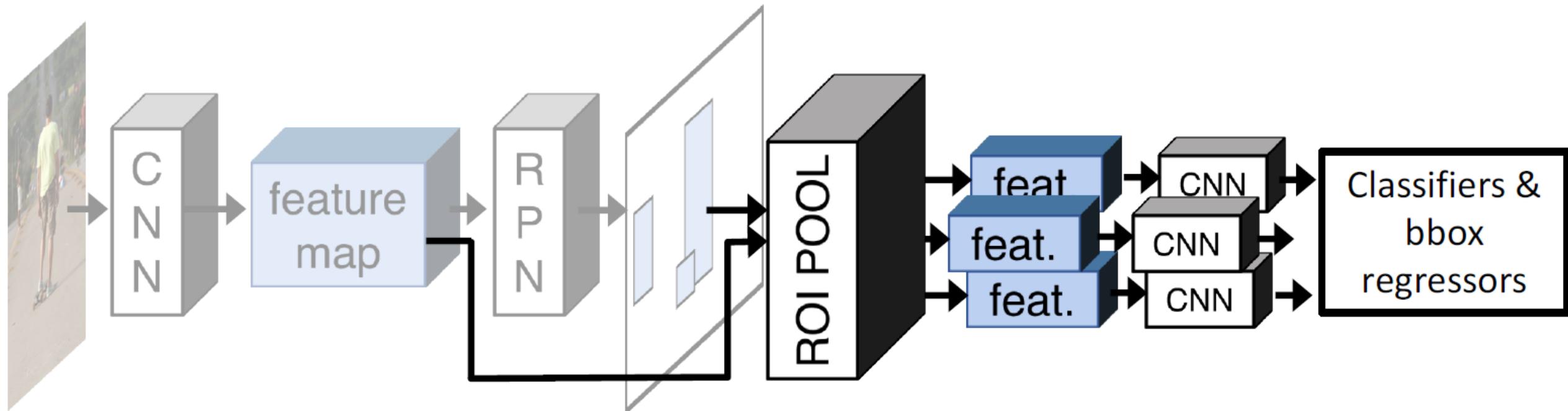
Attention Candidates , V

Bottom-up attention (using Faster R-CNN₂)



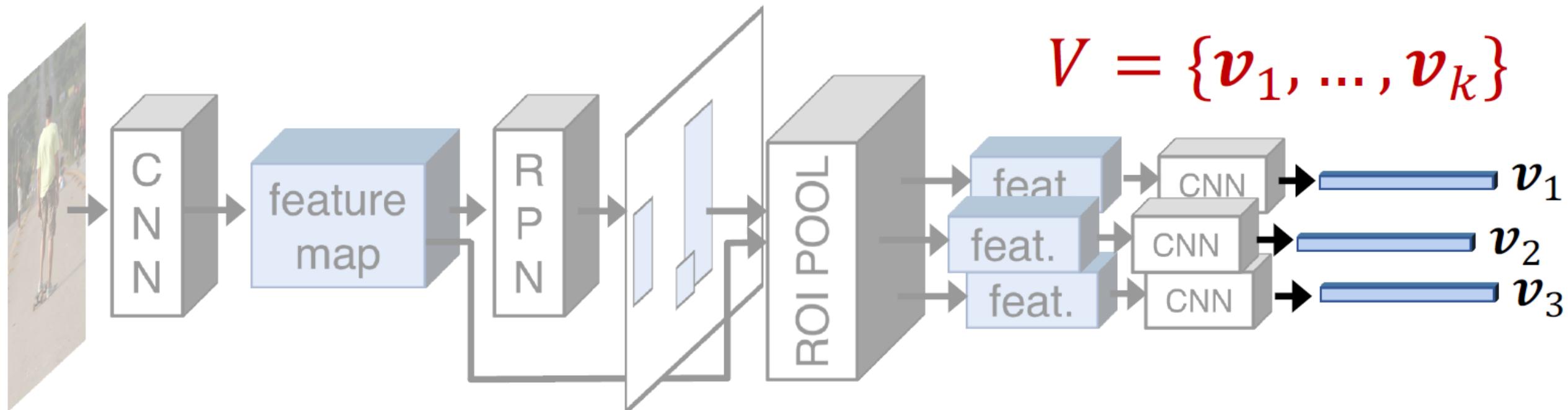
Attention Candidates , V

Bottom-up attention (using Faster R-CNN₂)



Attention Candidates , V

Bottom-up attention (using Faster R-CNN₂)

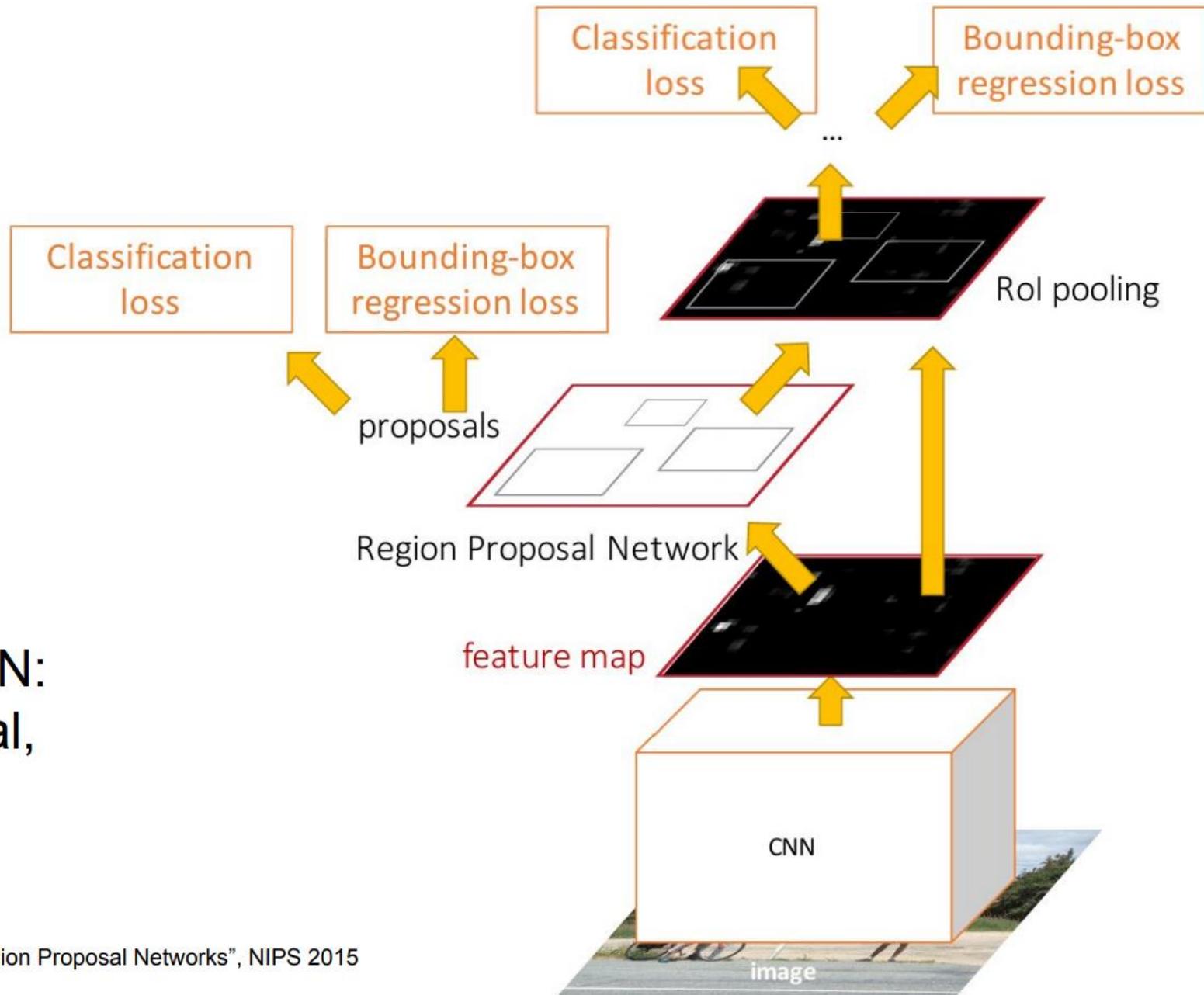


Faster R-CNN:

Make CNN do proposals!

Insert **Region Proposal Network (RPN)** to predict proposals from features

Otherwise same as Fast R-CNN:
Crop features for each proposal,
classify each one

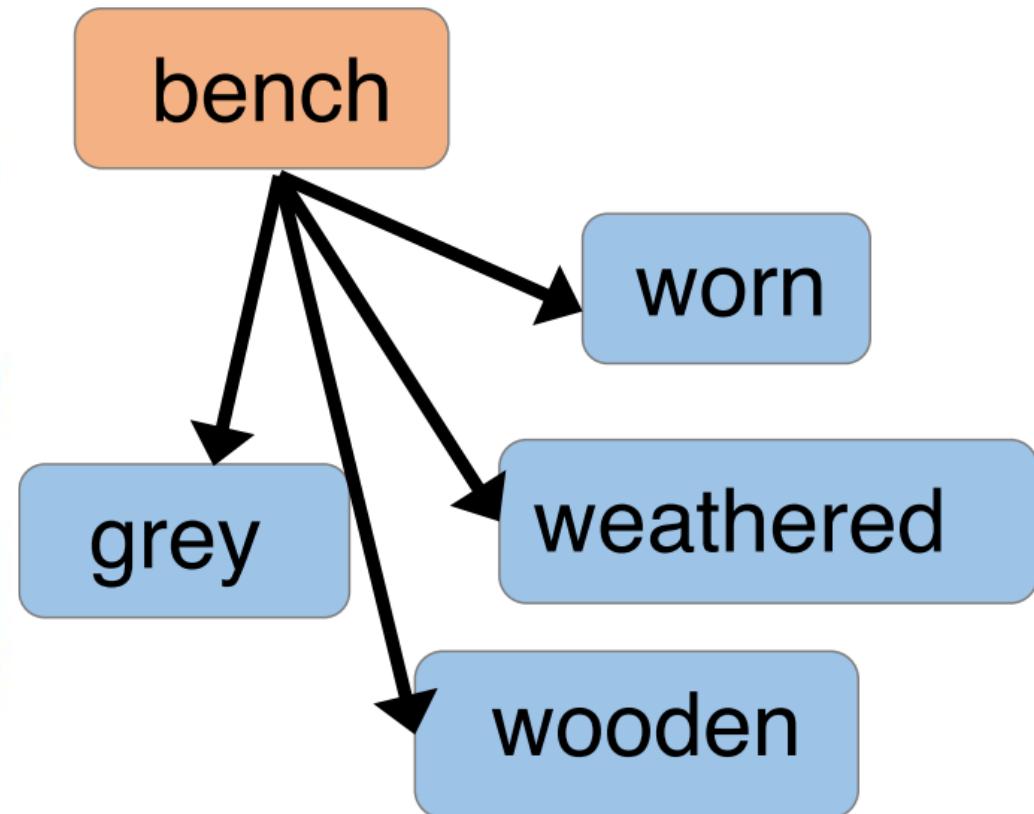


Ren et al, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks", NIPS 2015
Figure copyright 2015, Ross Girshick; reproduced with permission

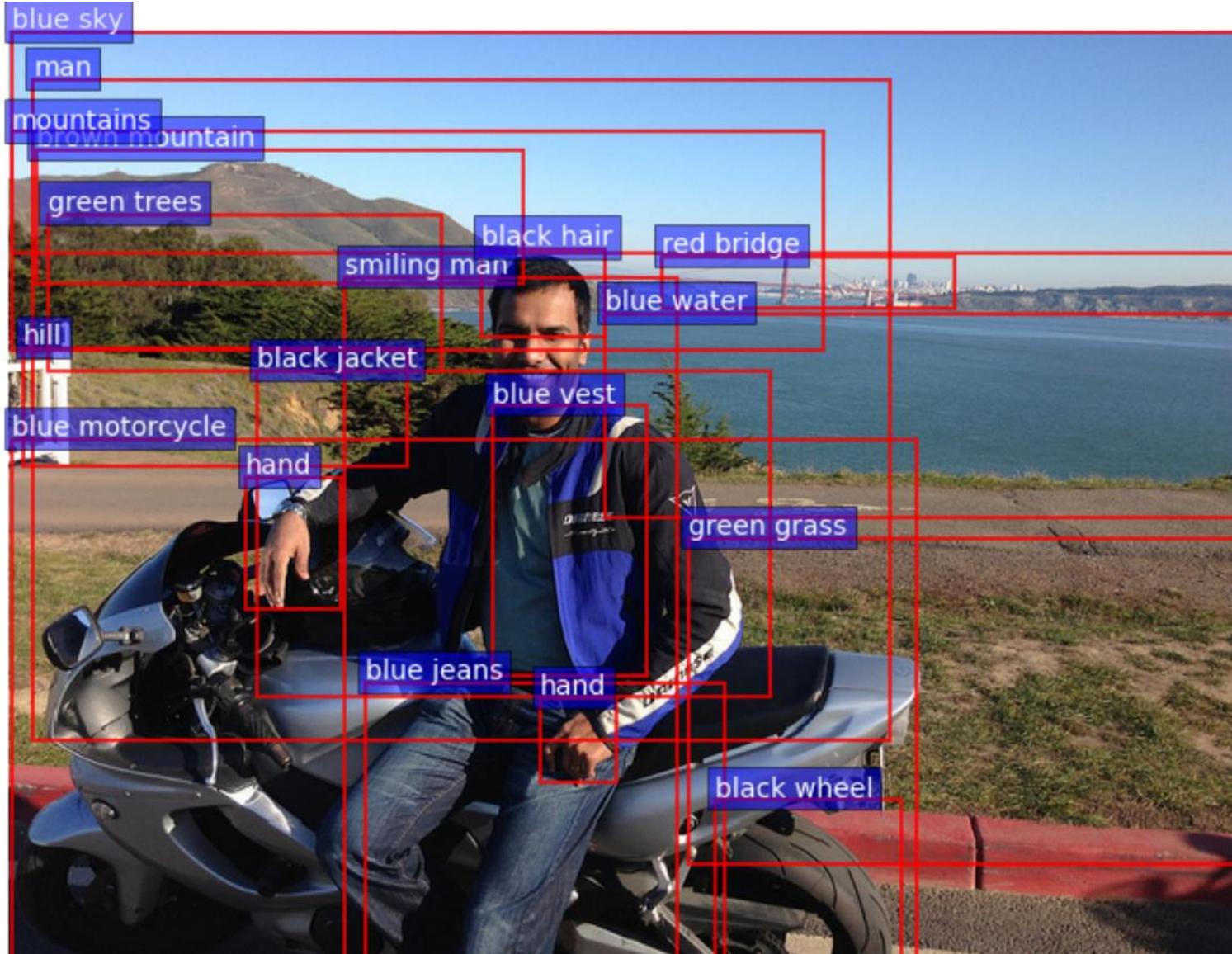
Faster RCNN pre-training

Using Visual Genome with:

- 1600 filtered object classes
- 400 filtered attribute classes

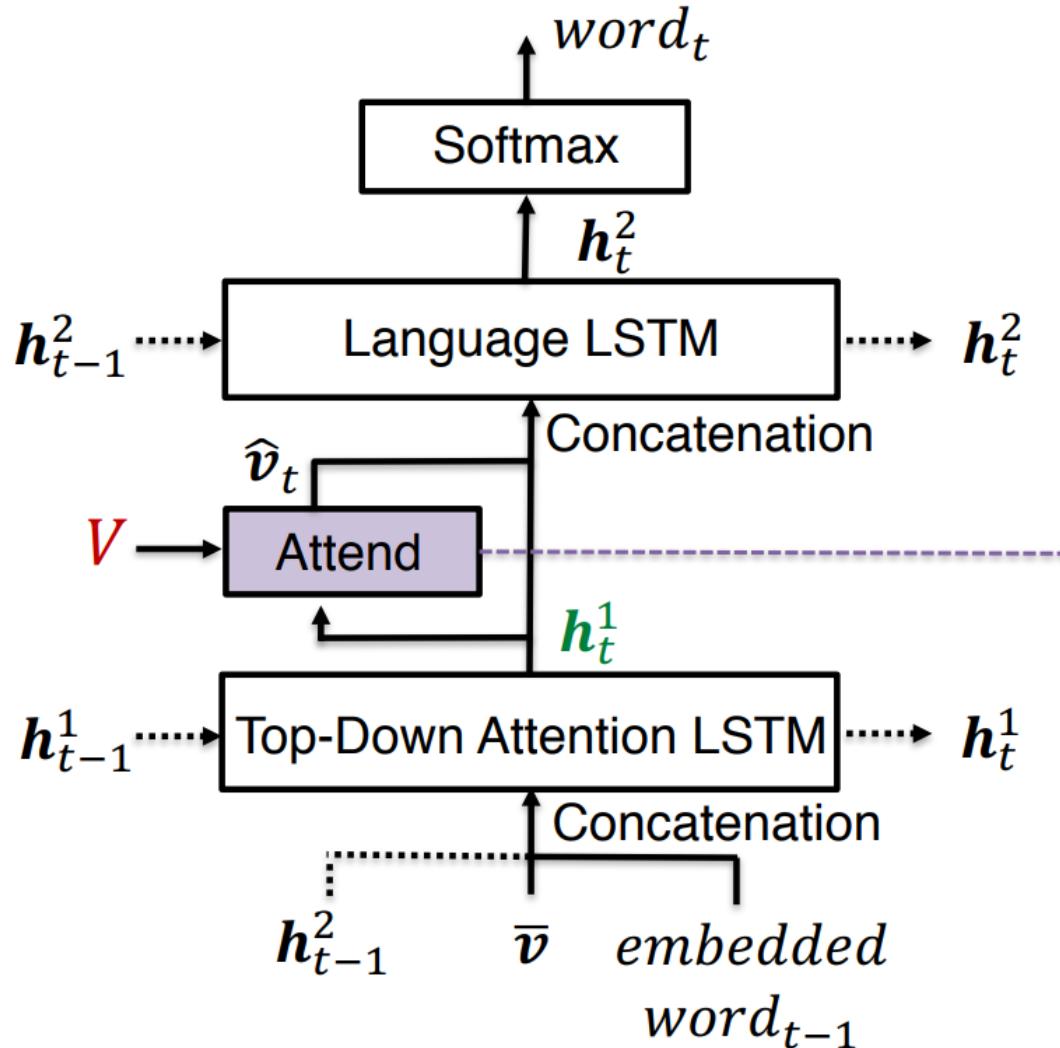


Result of Faster RCNN



Result of Faster RCNN

Image captioning model:



Attend block

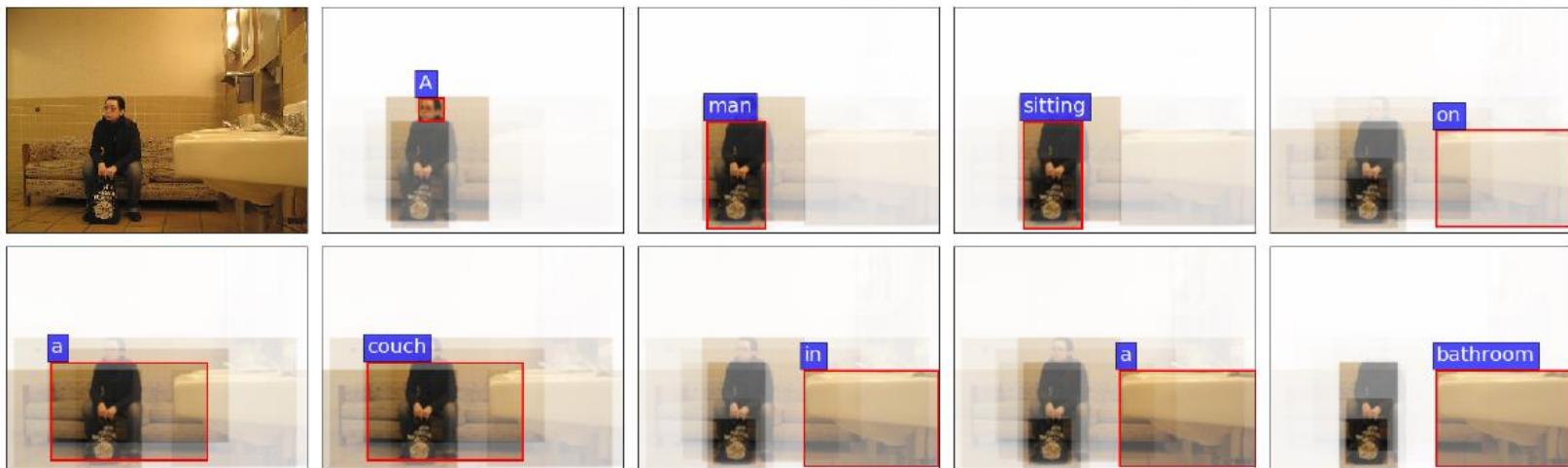
$$a_i = \mathbf{w}^T \tanh(W_v \mathbf{v}_i + W_h \mathbf{h})$$
$$\alpha = \text{softmax}(\mathbf{a})$$
$$\hat{\mathbf{v}} = \sum_{i=1}^k \alpha_i \mathbf{v}_i$$

Results of Up-Down attention model

ResNet (10x10): A man sitting on a **toilet** in a bathroom.

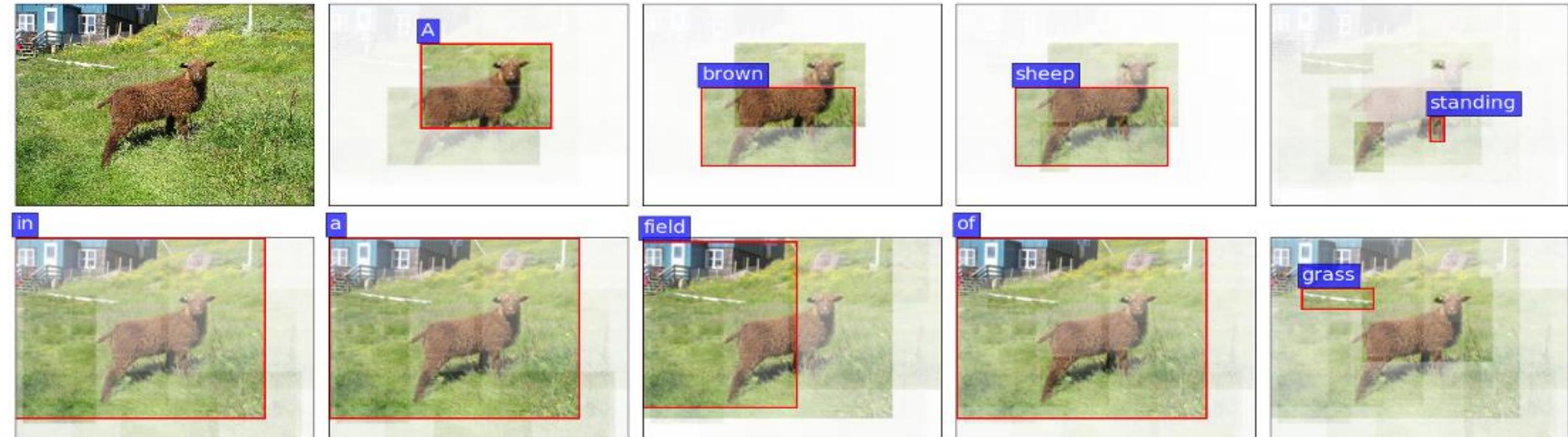


Up-Down : A man sitting on a **couch** in a bathroom.



Results of Up-Down attention model

Up-Down : A brown sheep standing in a field of grass.



Results of Up-Down attention model

1st COCO Captions leaderboard (July 2017)

COCO Captions “Karpathy” test set (single-model):

	BLEU-4	METEOR	CIDEr	SPICE
ResNet (10×10)	34.0	26.5	111.1	20.2
Up-Down (Ours)	36.3	27.7	120.1	21.4

+8% **+6%**

Neural Machine Translation

2014

(dramatic reenactment)

Neural Machine Translation

2014



Neural Machine Translation

- Neural Machine Translation (NMT) is a way to do Machine Translation with a *single neural network*
- The neural network architecture is called sequence-to-sequence (aka seq2seq) and it involves *two RNNs*.

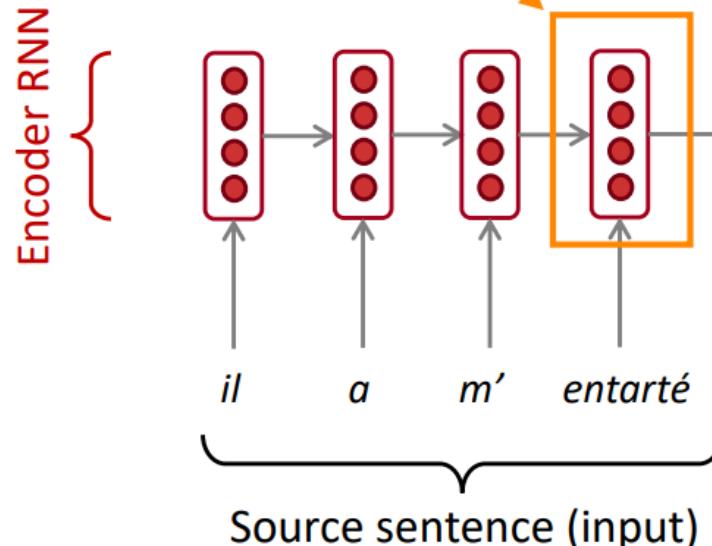
```
# i go school      가      가  
#                 가  
#  
#  
#      :  
#      가  
#  
#  
#  
#  
#
```

Neural Machine Translation[#]

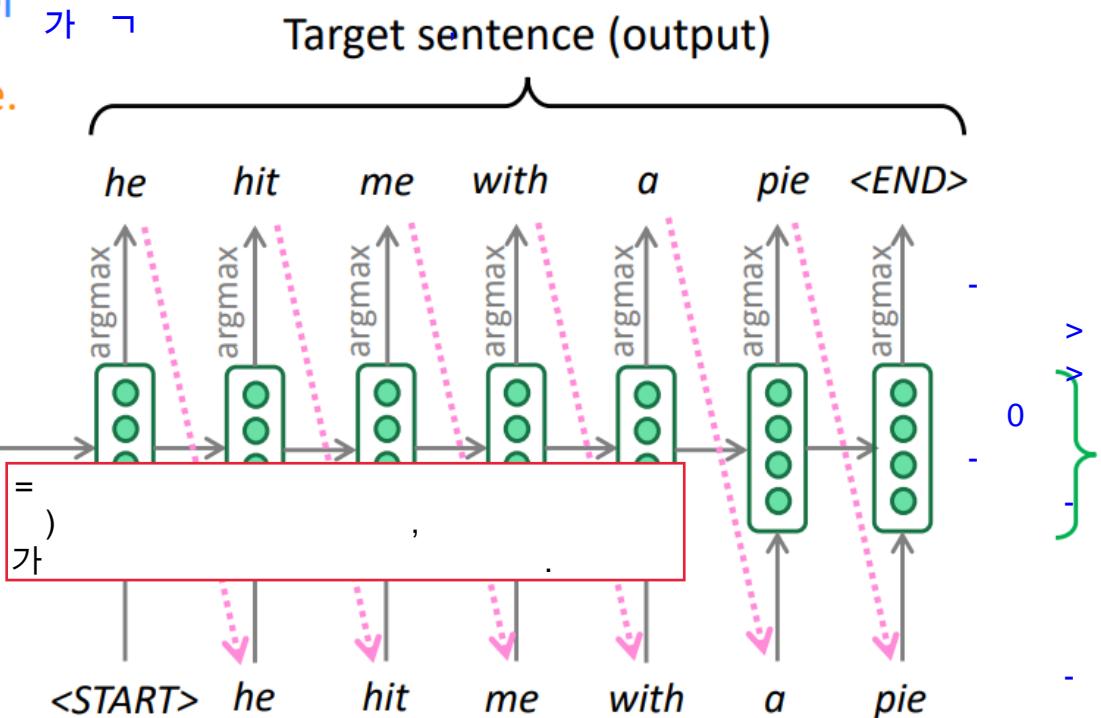
The sequence-to-sequence model

Encoding of the source sentence.

Provides initial hidden state
for Decoder RNN.



Encoder RNN produces
an **encoding** of the
source sentence.



Decoder RNN is a Language Model that generates
target sentence, *conditioned on encoding*.[#]

- sos, eos, pad

Note: This diagram shows **test time behavior**:
decoder output is fed in → as next step's input

Neural Machine Translation

- Sequence-to-sequence is useful for *more than just MT*
- Many NLP tasks can be phrased as sequence-to-sequence:
 - **Summarization** (long text → short text)
 - **Dialogue** (previous utterances → next utterance)
 - **Parsing** (input text → output parse as sequence)
 - **Code generation** (natural language → Python code)

- 가
- .
- :
- .
- .
- .
- .
#

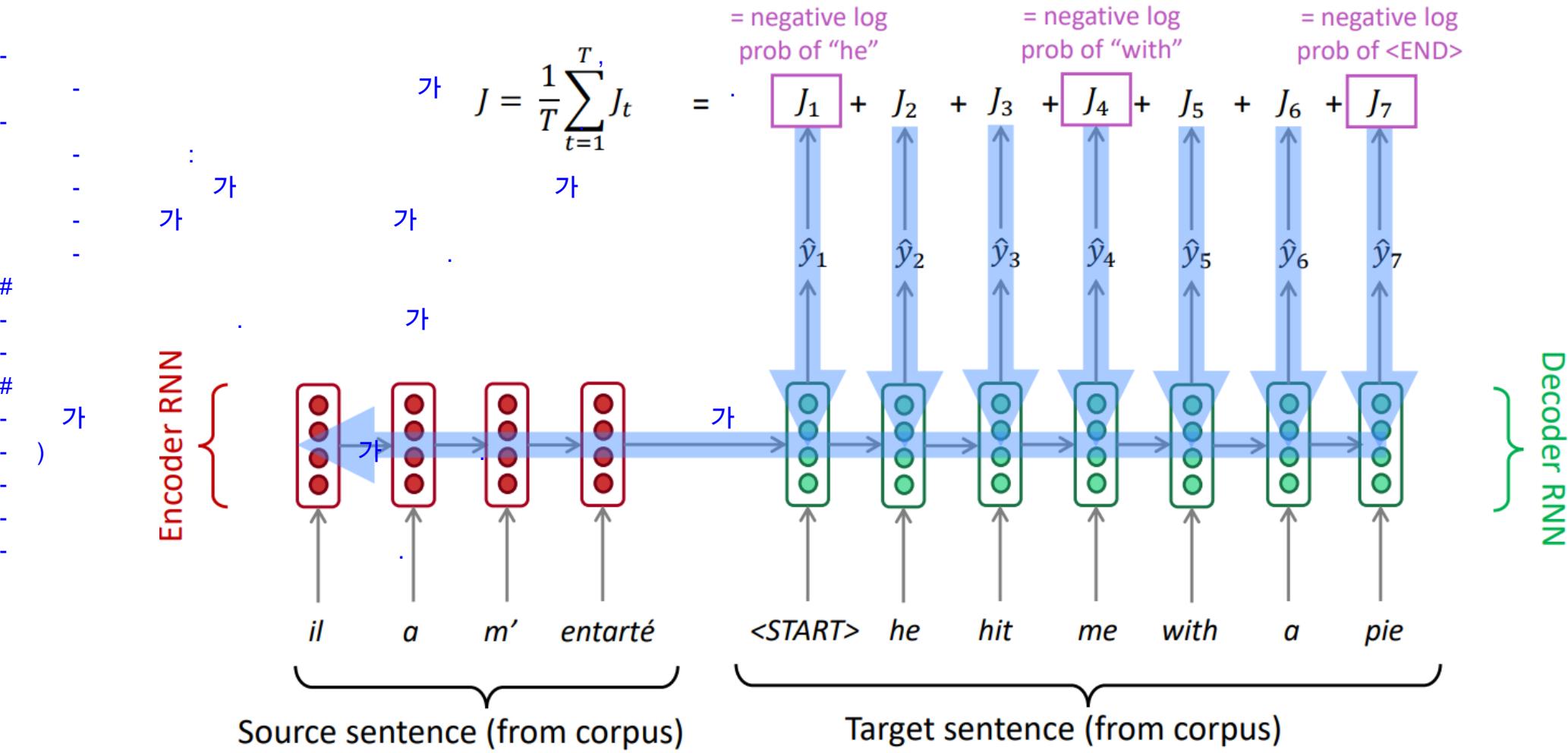
Neural Machine Translation

- The sequence-to-sequence model is an example of a **Conditional Language Model**.
 - **Language Model** because the decoder is predicting the next word of the target sentence y
 - **Conditional** because its predictions are *also* conditioned on the source sentence x
- NMT directly calculates $P(y|x)$:
$$P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots P(y_T|y_1, \dots, y_{T-1}, x)$$

NMT #MT

Probability of next target word, given target words so far and source sentence x
- **Question:** How to train a NMT system?
- **Answer:** Get a big parallel corpus...

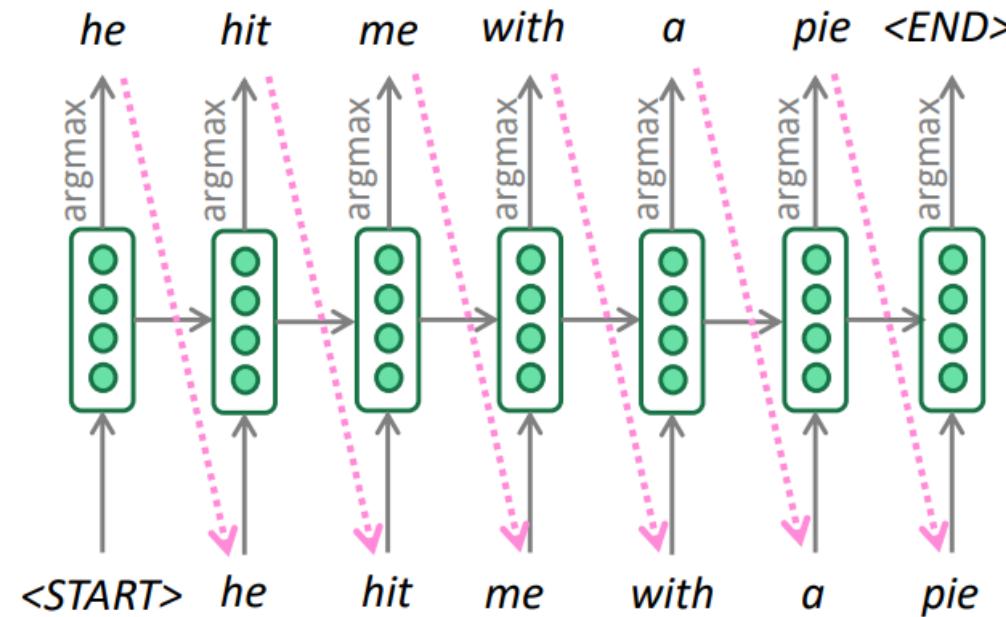
Neural Machine Translation



Seq2seq is optimized as a single system.
Backpropagation operates “end-to-end”.

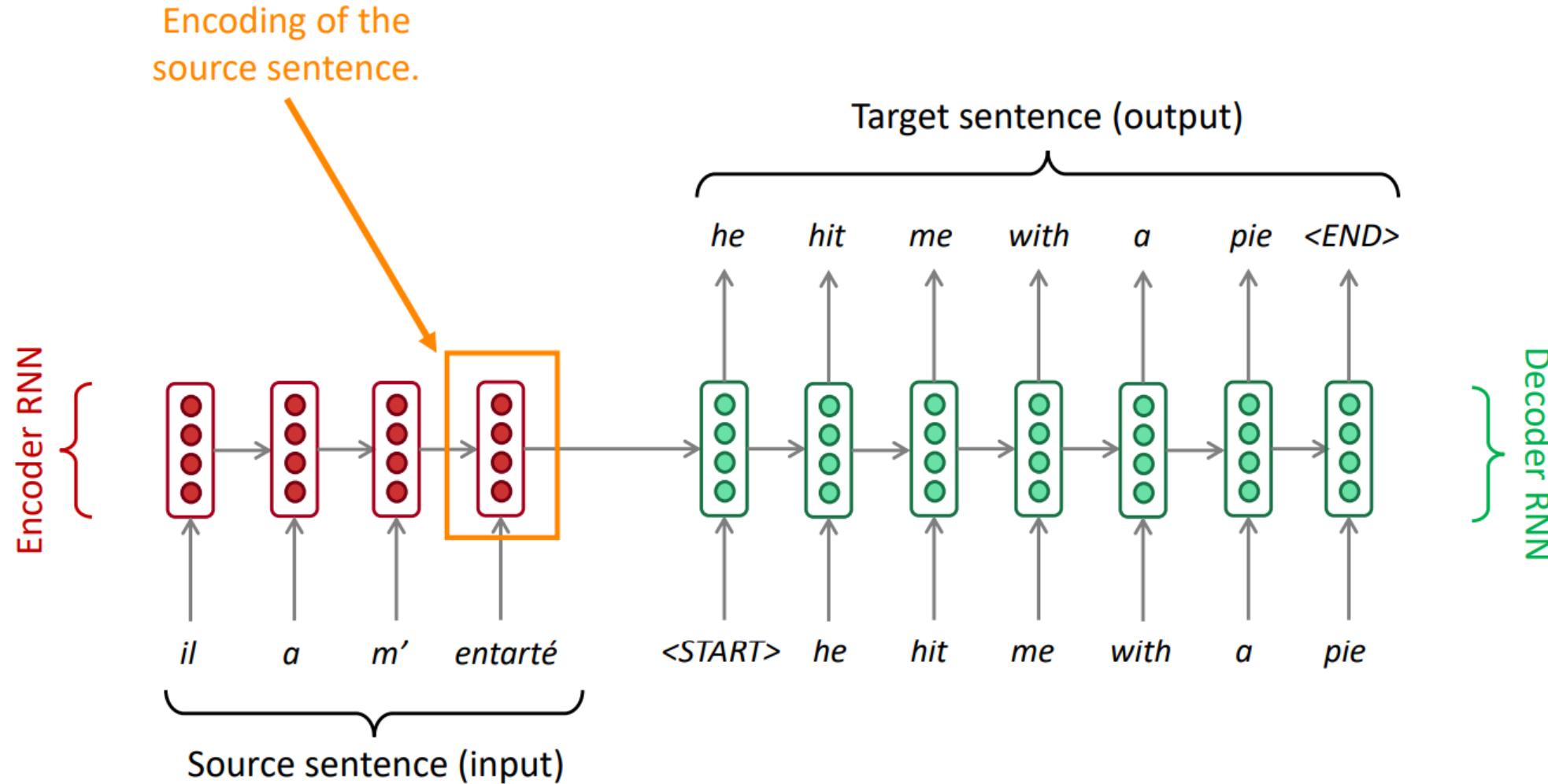
Greedy Decoding

- We saw how to generate (or “decode”) the target sentence by taking argmax on each step of the decoder



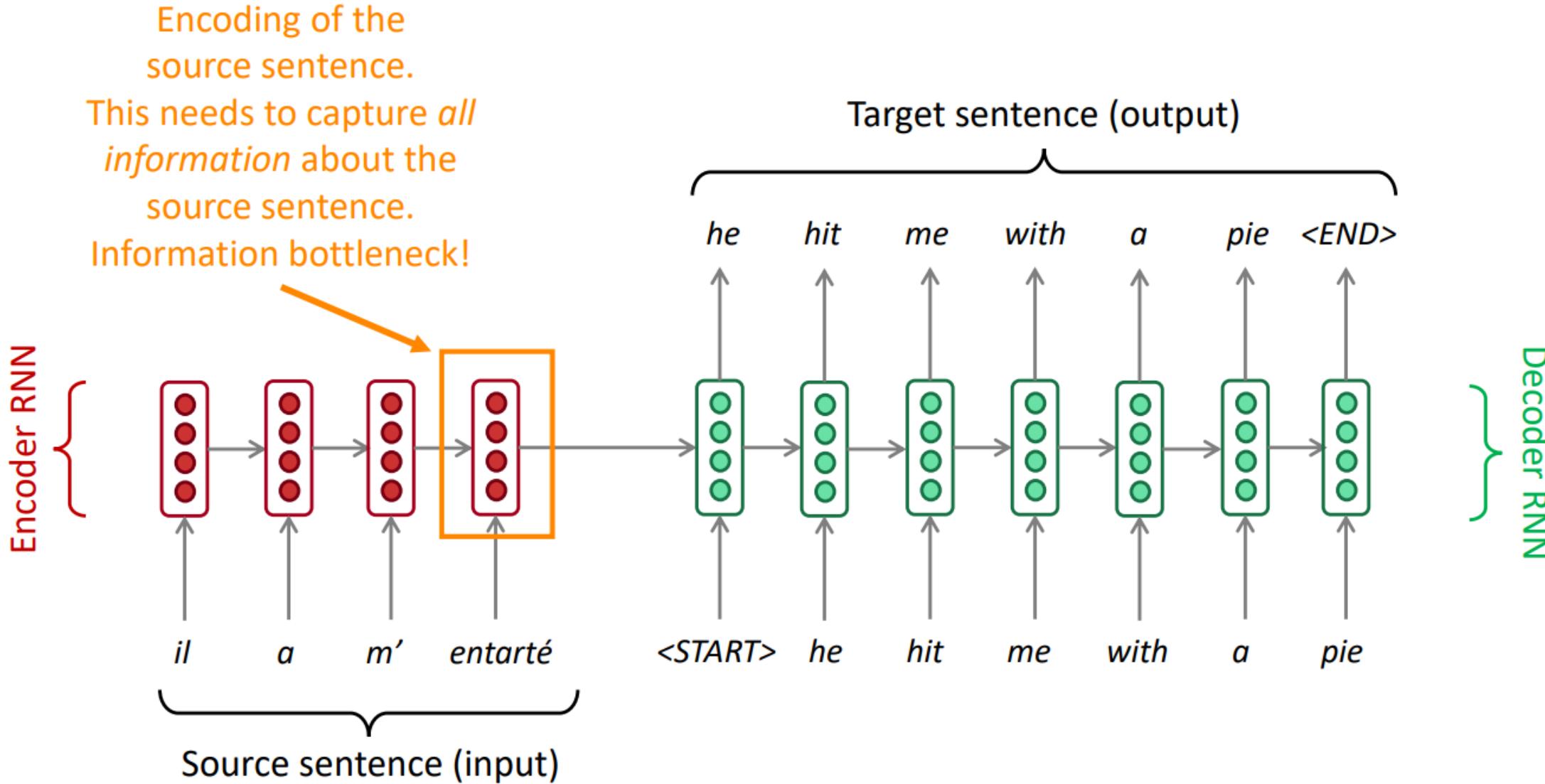
- This is **greedy decoding** (take most probable word on each step)
- **Problems with this method?**

Neural Machine Translation



Problems with this architecture?

Neural Machine Translation



Attention

- **Attention** provides a solution to the bottleneck problem.
- Core idea: on each step of the decoder, use *direct connection to the encoder* to *focus on a particular part* of the source sequence

-
-
-



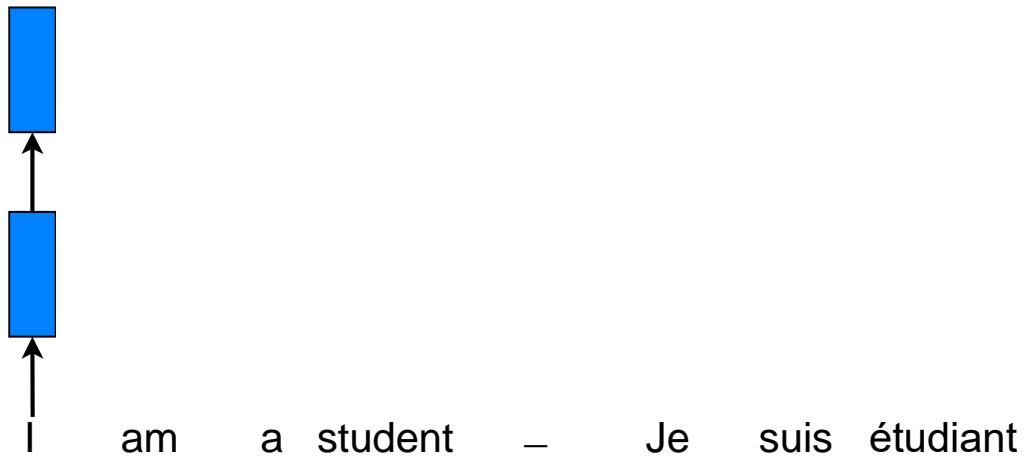
- First we will show via diagram (no equations), then we will show with equations

Neural Machine Translation (NMT)

I am a student – Je suis étudiant

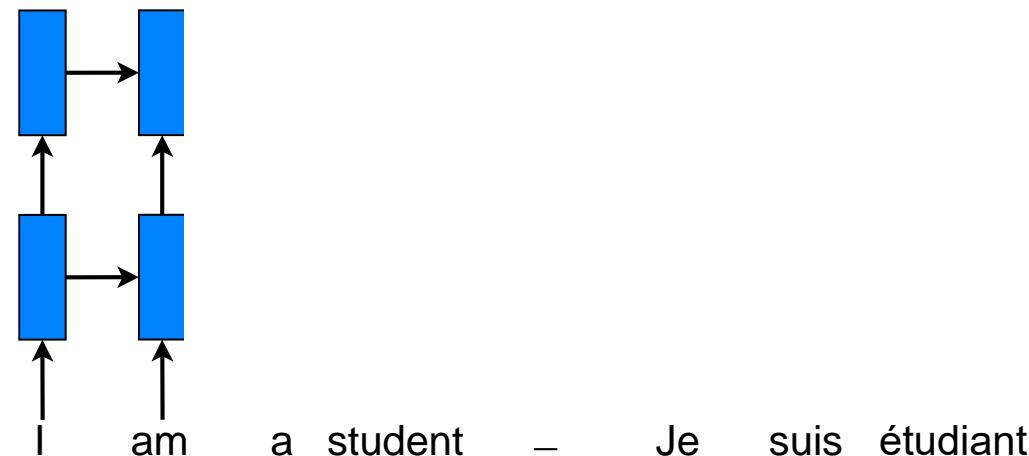
- Big RNNs trained **end-to-end**.

Neural Machine Translation (NMT)



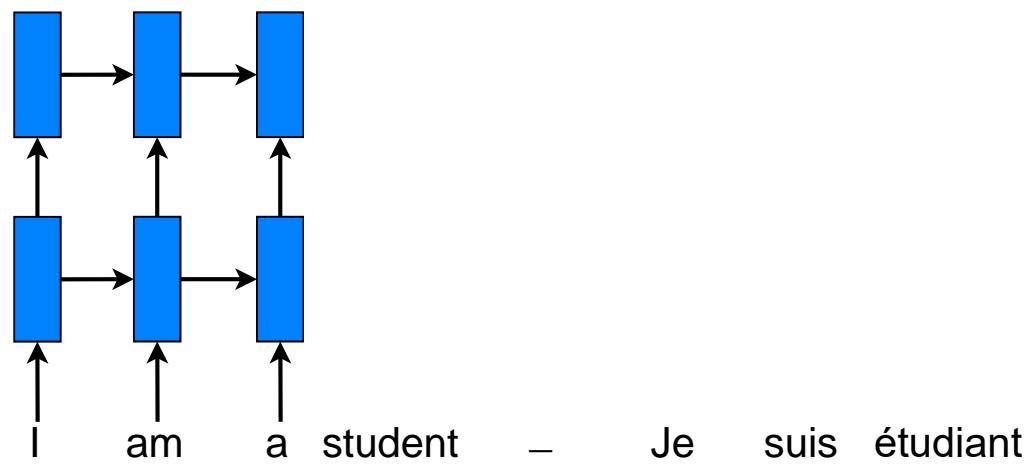
- Big RNNs trained **end-to-end**.

Neural Machine Translation (NMT)



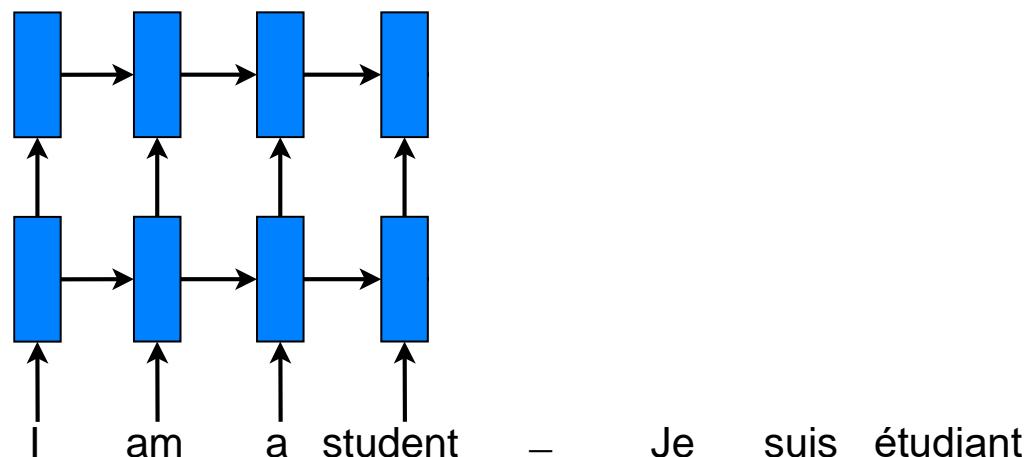
- Big RNNs trained **end-to-end**.

Neural Machine Translation (NMT)



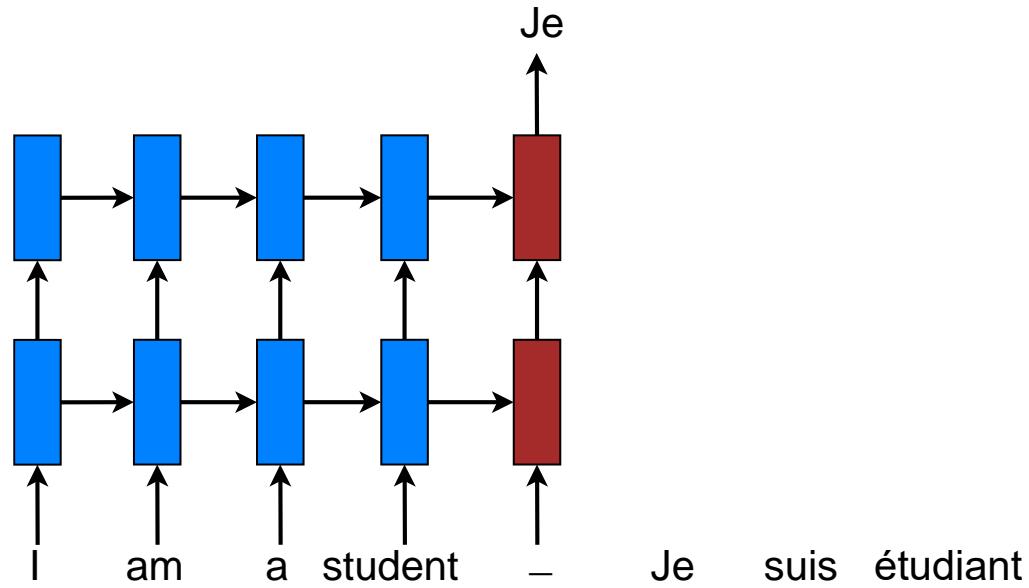
- Big RNNs trained end-to-end.

Neural Machine Translation (NMT)



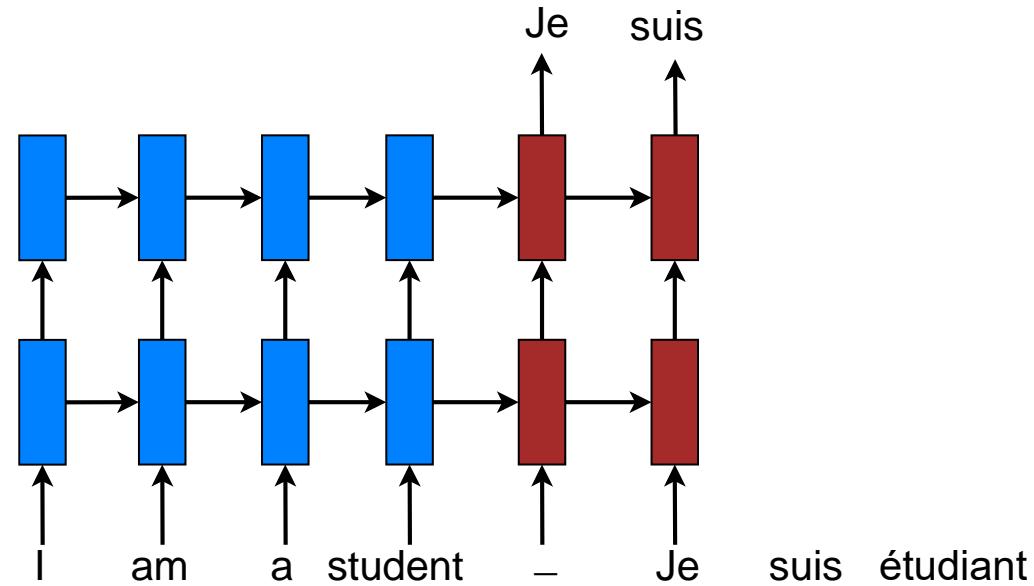
- Big RNNs trained end-to-end.

Neural Machine Translation (NMT)



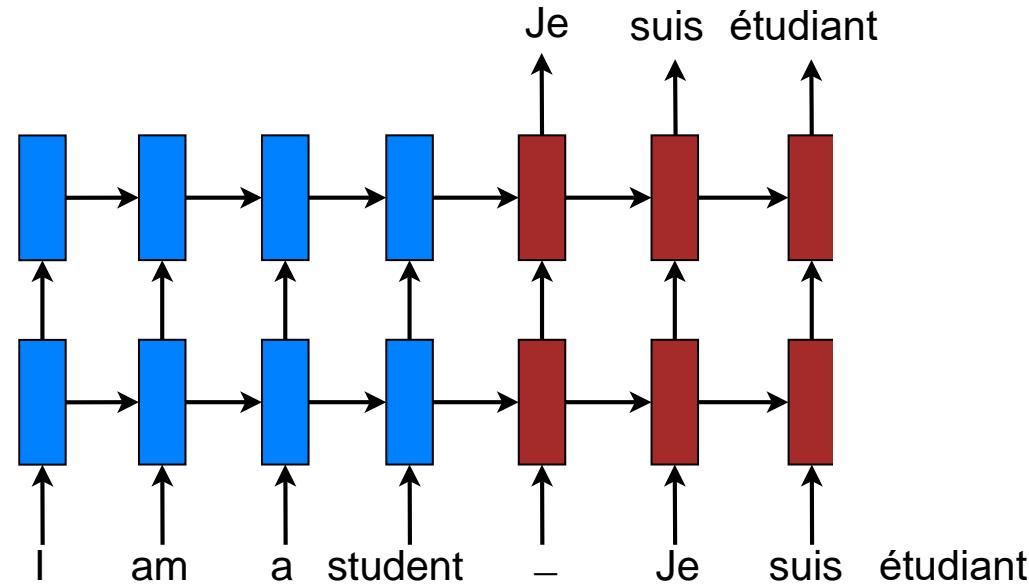
- Big RNNs trained end-to-end: **encoder-decoder**.

Neural Machine Translation (NMT)



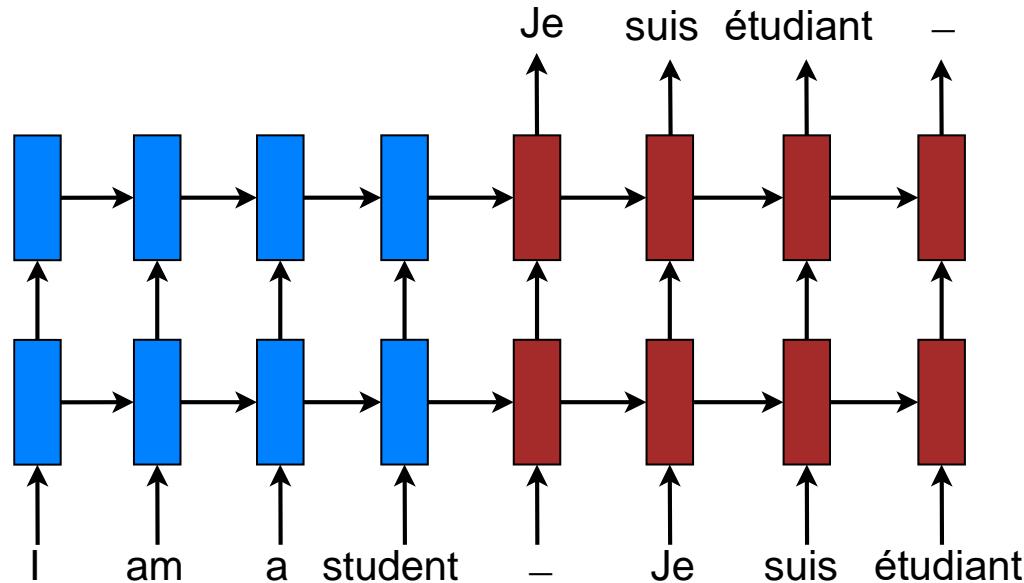
- Big RNNs trained end-to-end: **encoder-decoder**.

Neural Machine Translation (NMT)



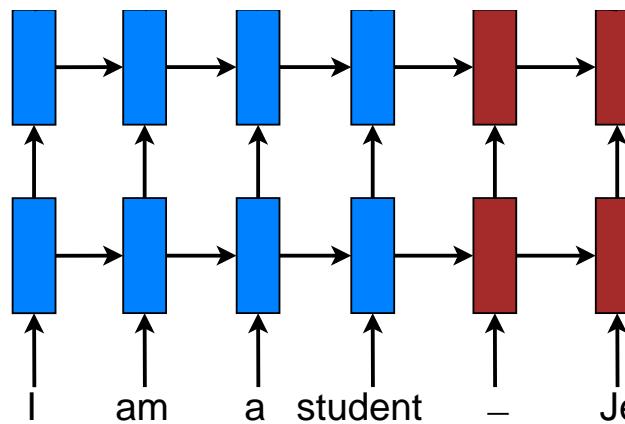
- Big RNNs trained end-to-end: **encoder-decoder**.

Neural Machine Translation (NMT)



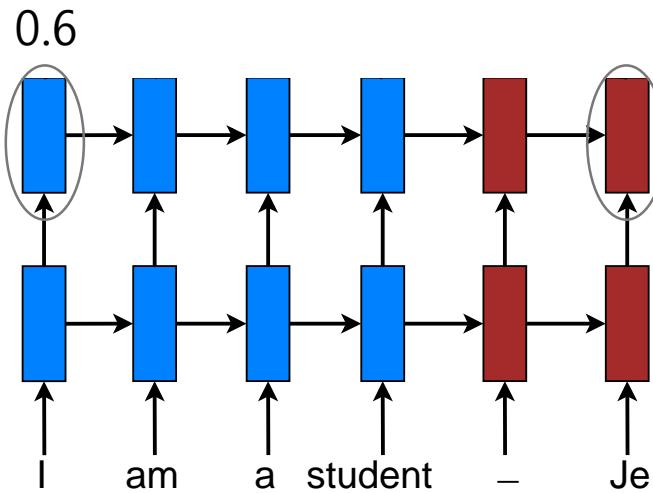
- Big RNNs trained end-to-end: **encoder-decoder**.
 - Generalize well to long sequences.
 - Small memory footprint.
 - Simple decoder.

Attention Mechanism



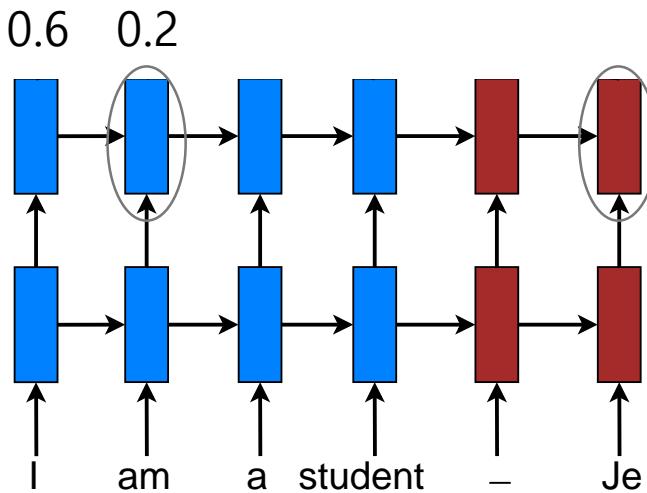
- Maintain a **memory** of source hidden states

Attention Mechanism



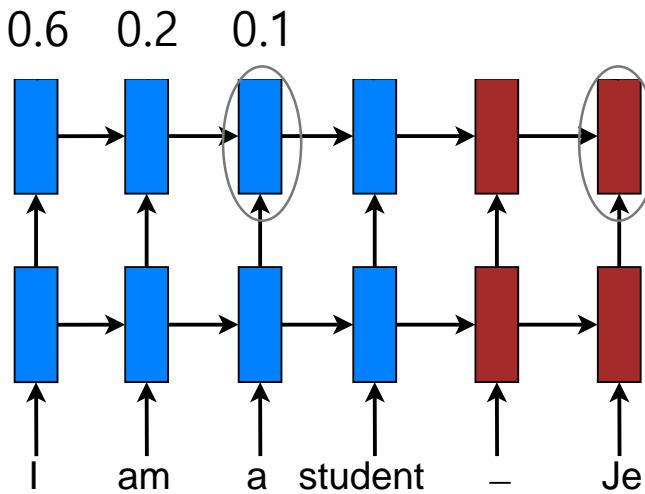
- Maintain a **memory** of source hidden states
 - Compare target and source hidden states

Attention Mechanism



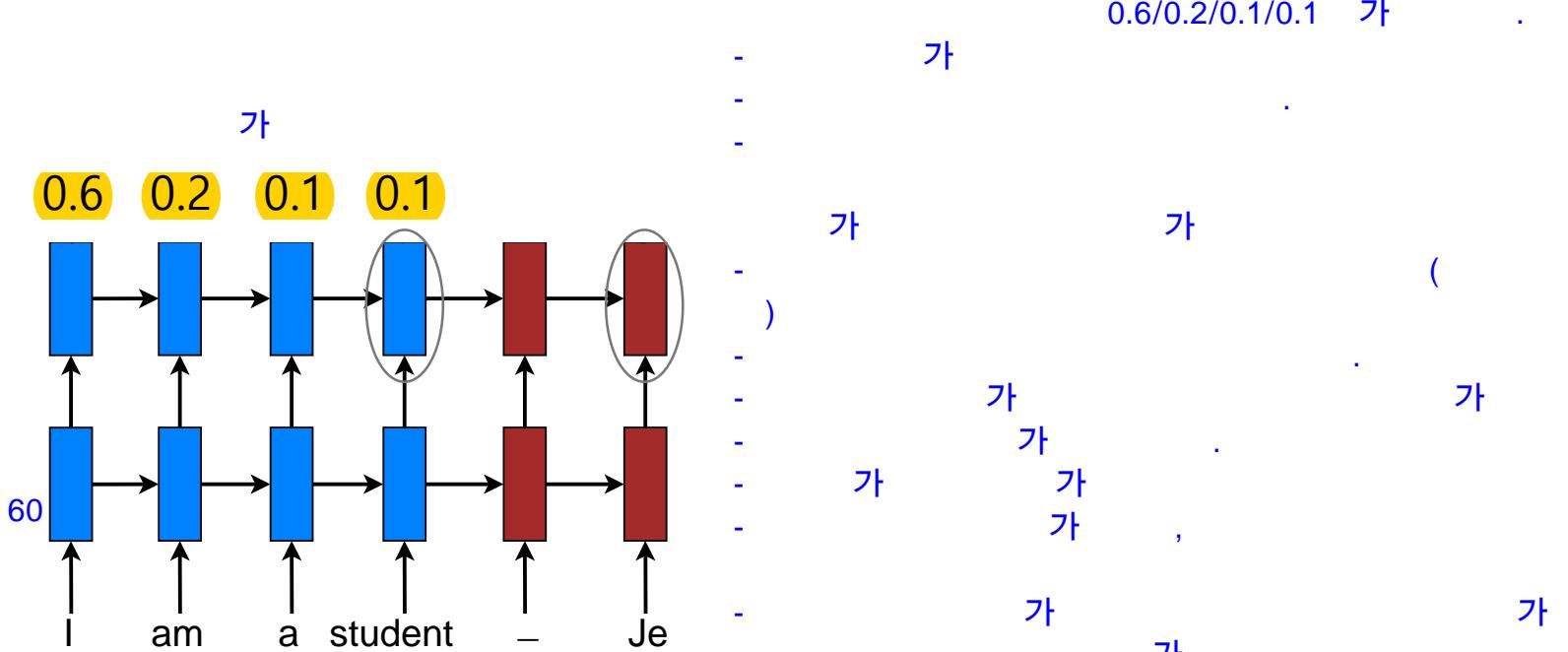
- Maintain a **memory** of source hidden states
 - Compare target and source hidden states

Attention Mechanism



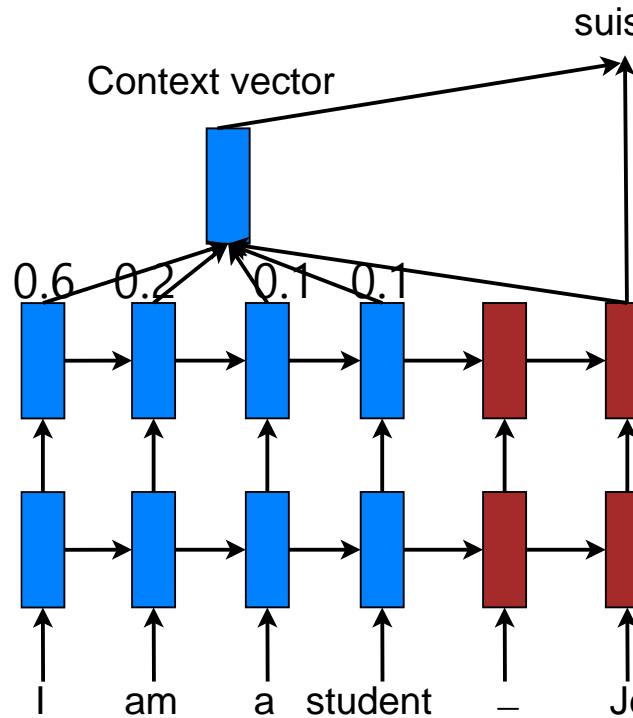
- Maintain a **memory** of source hidden states
 - Compare target and source hidden states

Attention Mechanism



- Maintain a **memory** of source hidden states
 - Compare target and source hidden states

Attention Mechanism



- Maintain a **memory** of source hidden states
 - Able to translate long sentences.

No other attention architectures b
eside (Bahdanau et al., 2015)

Attention

- We have encoder hidden states $h_1, \dots, h_N \in \mathbb{R}^h$
- On timestep t , we have decoder hidden state $s_t \in \mathbb{R}^h$
- We get the attention scores e^t for this step:

$$e^t = [s_t^T h_1, \dots, s_t^T h_N] \in \mathbb{R}^N$$

- We take softmax to get the attention distribution α^t for this step (this is a probability distribution and sums to 1)

$$\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$$

- We use α^t to take a weighted sum of the encoder hidden states to get the attention output a_t

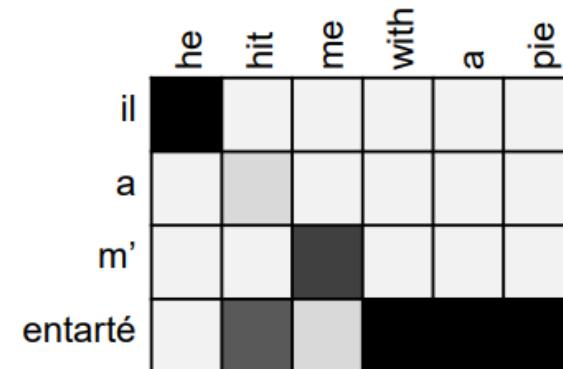
$$a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$$

- Finally we concatenate the attention output a_t with the decoder hidden state s_t and proceed as in the non-attention seq2seq model

$$[a_t; s_t] \in \mathbb{R}^{2h}$$

Attention is great

- Attention significantly improves NMT performance
 - It's very useful to allow decoder to focus on certain parts of the source
- Attention solves the bottleneck problem
 - Attention allows decoder to look directly at source; bypass bottleneck
- Attention helps with vanishing gradient problem
 - Provides shortcut to faraway states
- Attention provides some interpretability
 - By inspecting attention distribution, we can see what the decoder was focusing on
 - We get (soft) alignment for free!
 - This is cool because we never explicitly trained an alignment system
 - The network just learned alignment by itself



Attention Models

- We've seen that attention is a great way to improve the sequence-to-sequence model for Machine Translation.
 - However: You can use attention in **many architectures** (not just seq2seq) and **many tasks** (not just MT)
-
- More general definition of attention:
 - Given a set of vector *values*, and a vector *query*, attention is a technique to compute a weighted sum of the values, dependent on the query.
 - We sometimes say that the *query attends to the values*.
 - For example, in the seq2seq + attention model, each decoder hidden state (*query*) *attends to* all the encoder hidden states (*values*).

Attention is a general Deep Learning technique

More general definition of attention:

Given a set of vector *values*, and a vector *query*, attention is a technique to compute a weighted sum of the values, dependent on the query.

Intuition:

- The weighted sum is a *selective summary* of the information contained in the values, where the query determines which values to focus on.
- Attention is a way to obtain a *fixed-size representation of an arbitrary set of representations* (the values), dependent on some other representation (the query).

There are several attention variants

- We have some *values* $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$ and a *query* $\mathbf{s} \in \mathbb{R}^{d_2}$
- Attention always involves:
 1. Computing the *attention scores* $\mathbf{e} \in \mathbb{R}^N$
 2. Taking softmax to get *attention distribution* α :

There are
multiple ways
to do this

$$\alpha = \text{softmax}(\mathbf{e}) \in \mathbb{R}^N$$

- 3. Using attention distribution to take weighted sum of values:

$$\mathbf{a} = \sum_{i=1}^N \alpha_i \mathbf{h}_i \in \mathbb{R}^{d_1}$$

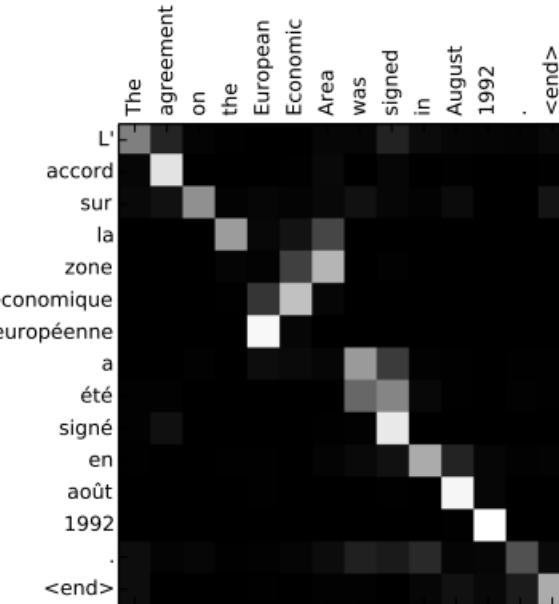
thus obtaining the *attention output* \mathbf{a} (sometimes called the *context vector*)

Attention variants

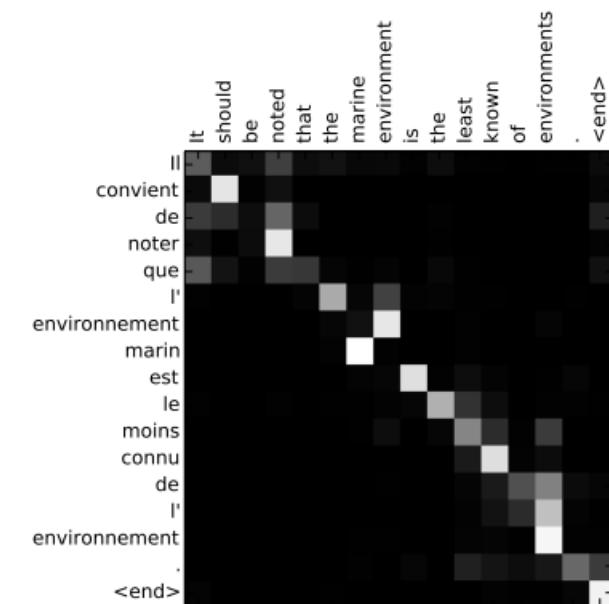
There are **several ways** you can compute $e \in \mathbb{R}^N$ from $\mathbf{h}_1, \dots, \mathbf{h}_N \in \mathbb{R}^{d_1}$ and $\mathbf{s} \in \mathbb{R}^{d_2}$:

- Basic dot-product attention: $e_i = \mathbf{s}^T \mathbf{h}_i \in \mathbb{R}$
 - Note: this assumes $d_1 = d_2$
 - This is the version we saw earlier
- Multiplicative attention: $e_i = \mathbf{s}^T \mathbf{W} \mathbf{h}_i \in \mathbb{R}$
 - Where $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$ is a weight matrix
- Additive attention: $e_i = \mathbf{v}^T \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s}) \in \mathbb{R}$
 - Where $\mathbf{W}_1 \in \mathbb{R}^{d_3 \times d_1}$, $\mathbf{W}_2 \in \mathbb{R}^{d_3 \times d_2}$ are weight matrices and $\mathbf{v} \in \mathbb{R}^{d_3}$ is a weight vector.
 - d_3 (the attention dimensionality) is a hyperparameter

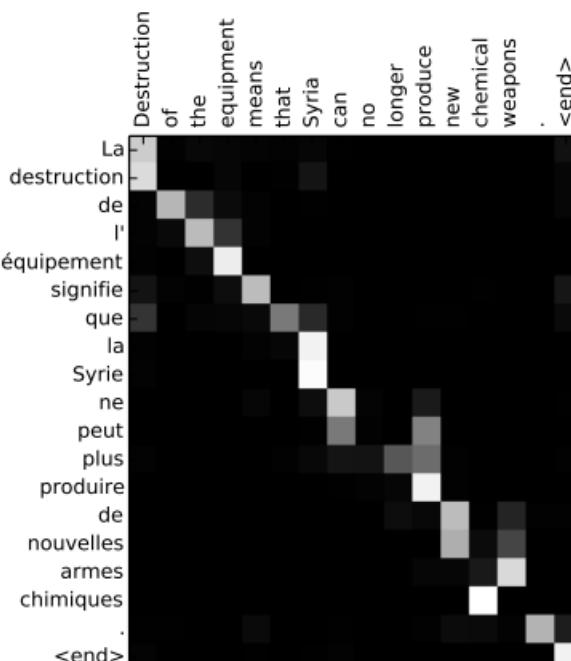
Attention Example in Machine Translation



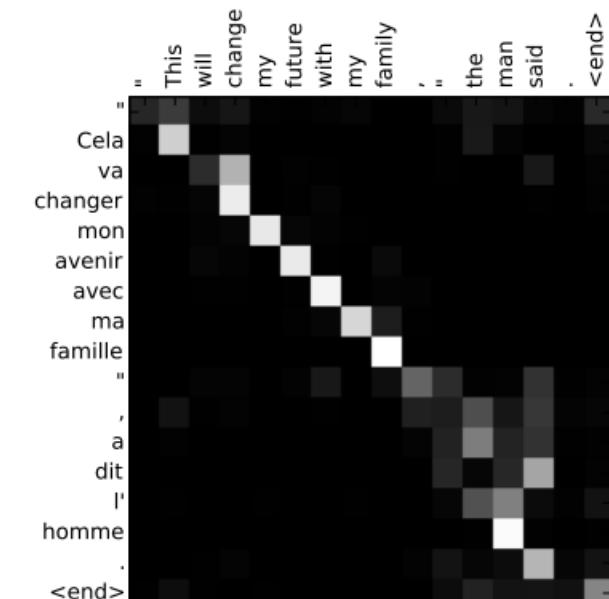
(a)



(b)



(c)



(d)

References

[Stanford University CS231n: Convolutional Neural Networks for Visual Recognition](#)

[Deep Learning Summer School, Montreal 2016 - VideoLectures.NET](#)

[Understanding LSTM Networks -- colah's blog](#)

<https://panderson.me/images/CVPR-Up-Down-talk.pdf>