

CMPT 983

Grounded Natural Language Understanding

January 28, 2021

Attention

Today

- Attention
 - Review of attention mechanism
 - Types of attentions

Attention

Need for "attention"

- Uses encoding of entire input when generating each output token
- Maybe would be useful to **focus** on a part of the input as the output tokens are generated

Translation

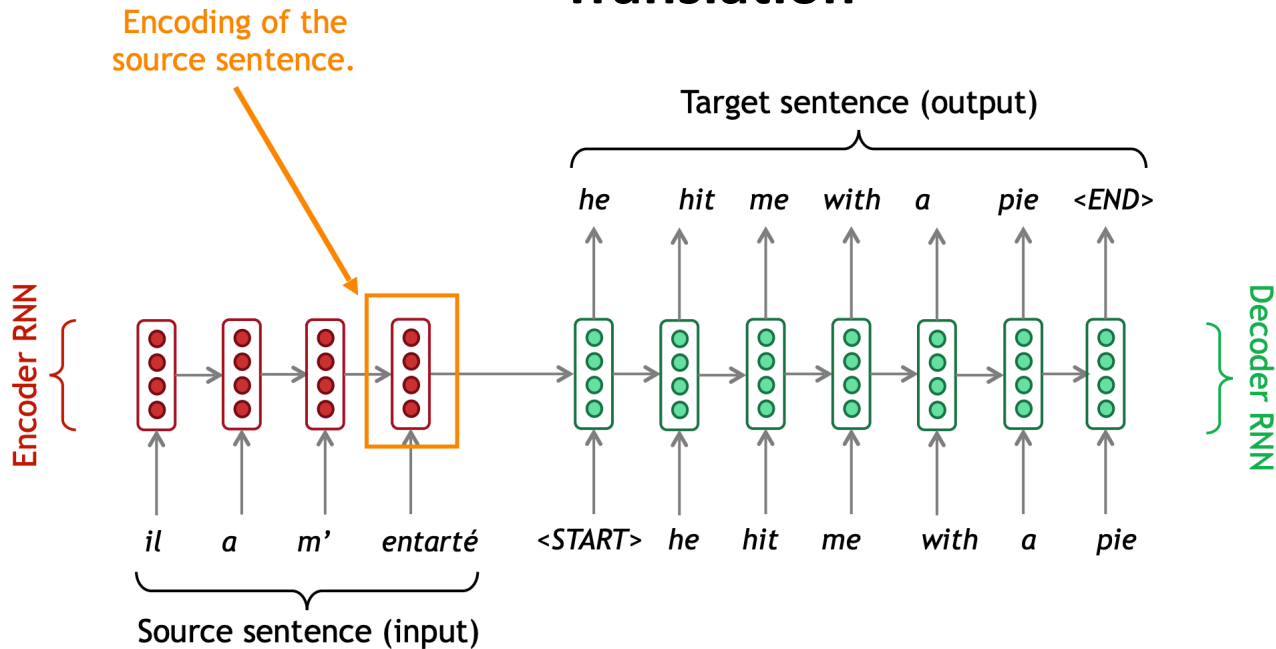


Image credit: Abigail See

Captioning

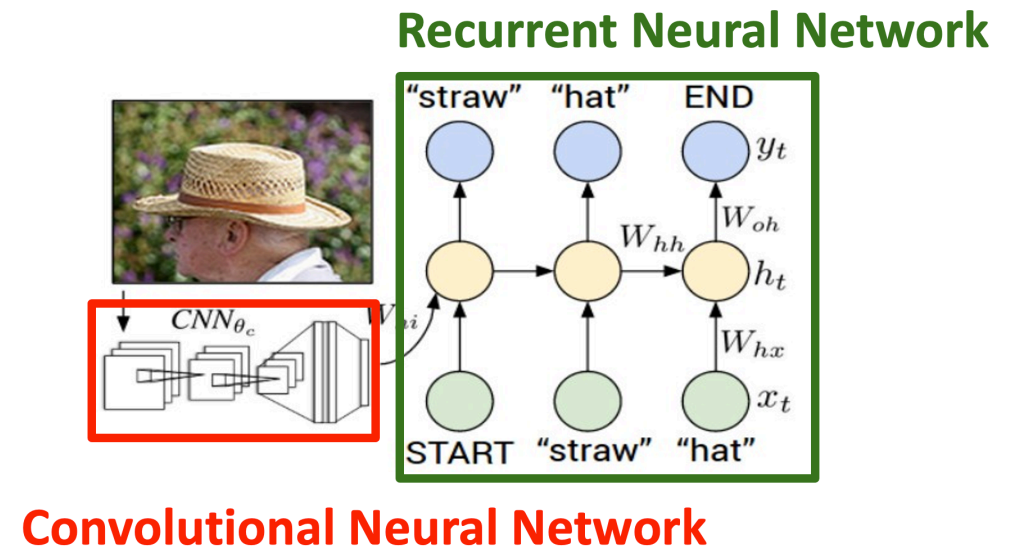
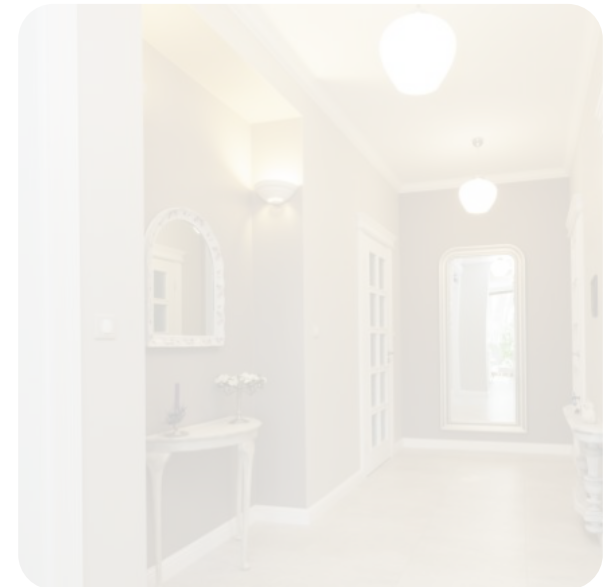


Image credit: Andrej Karpathy

Attention for VLN

Not every part of an input is important given the task context

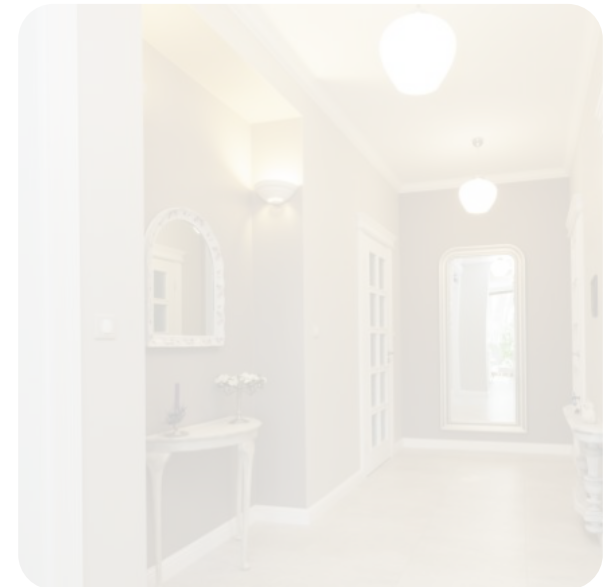
Exit the bathroom. Turn left and exit the room using the door on the left. Wait there.



Attention for VLN

Not every part of an input is important given the task context

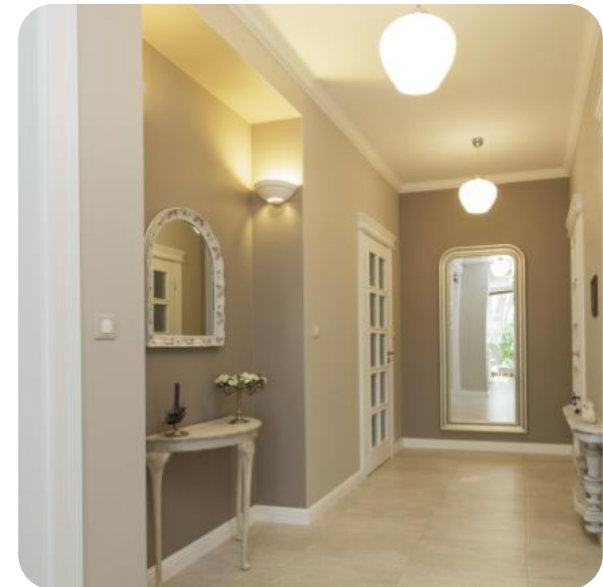
Exit the bathroom. Turn left and exit the room using the door on the left. Wait there.



Attention for VLN

Not every part of an input is important given the task context

Exit the bathroom. Turn left and exit the room using the door on the left. **Wait there.**



Attention for VLN

Not every part of an input is important given the task context

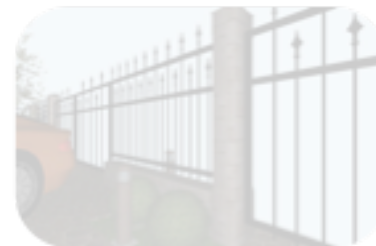
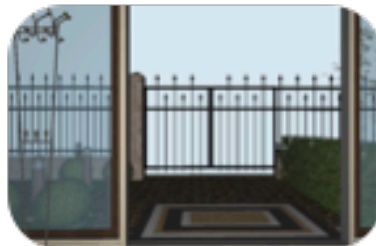
Take a right when you see the mirrored wardrobe.



Attention for VQA

Not every part of an input is important given the task context

What shape is the doormat?



Attention

- The concept of 'attention' has seen widespread use...
 - In many language and / or vision tasks, attention works extremely well!
 - Attention improves interpretability of neural networks

Q: What room are they in?



A: kitchen



Focus on part of input

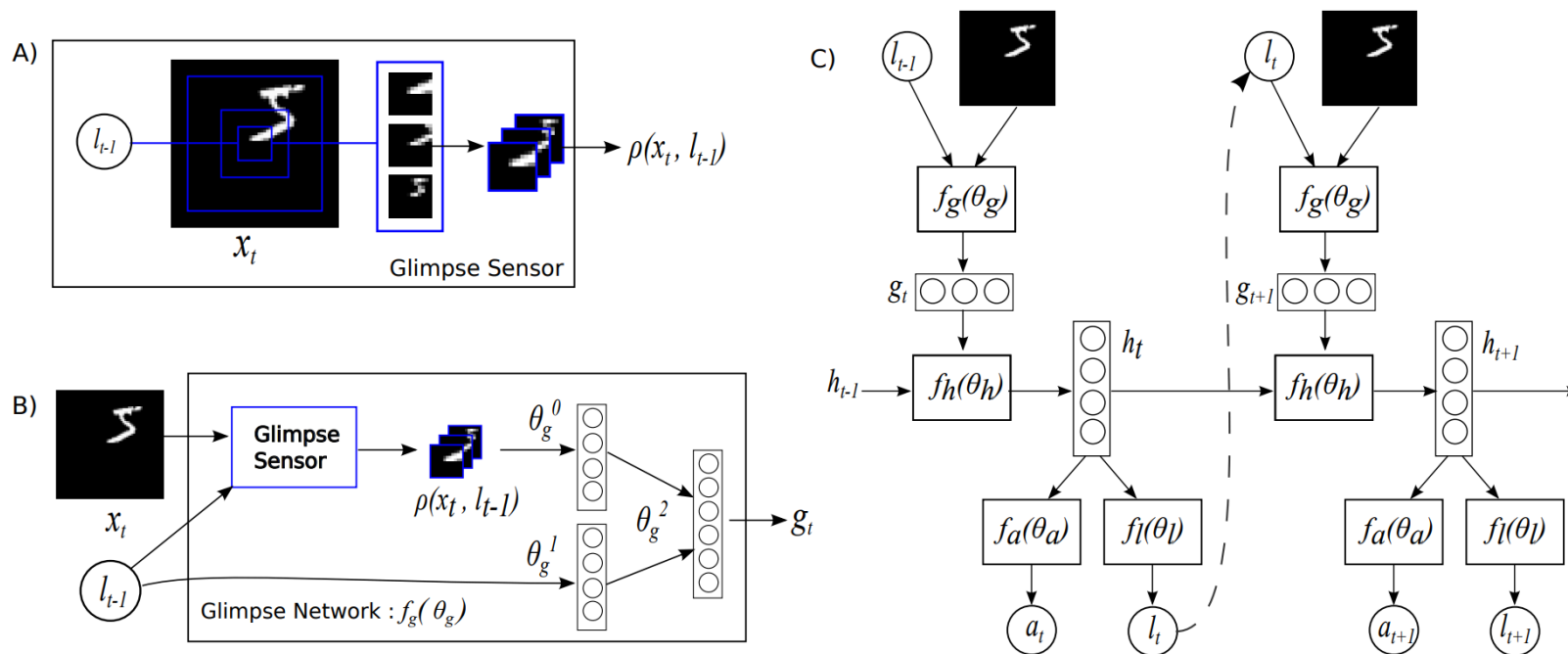
Attention mechanisms - summary

- Attention is probably one of the simplest and most effective ideas in deep learning – proven across many different domains
- Practically: focus on part of input by taking a weighted sum of different input parts
- With sufficient data, attention mechanisms can learn to ground language in visual content from 'distant' supervision
- Given the complexity of biological attention systems, assume there is still much to explore... particularly temporal aspects in context of LV&A

Attention

- Major impact on computer vision and NLP from 2014/5

Recurrent Models of Visual Attention

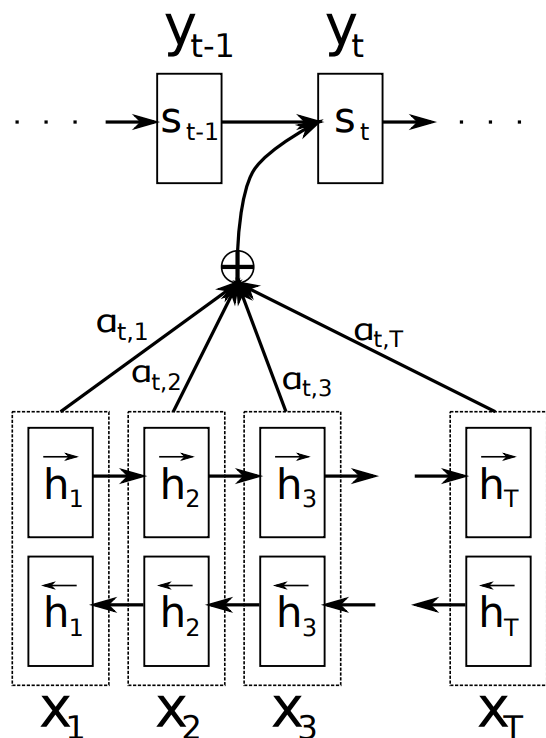


Mnih et al. NIPS 2014

Attention

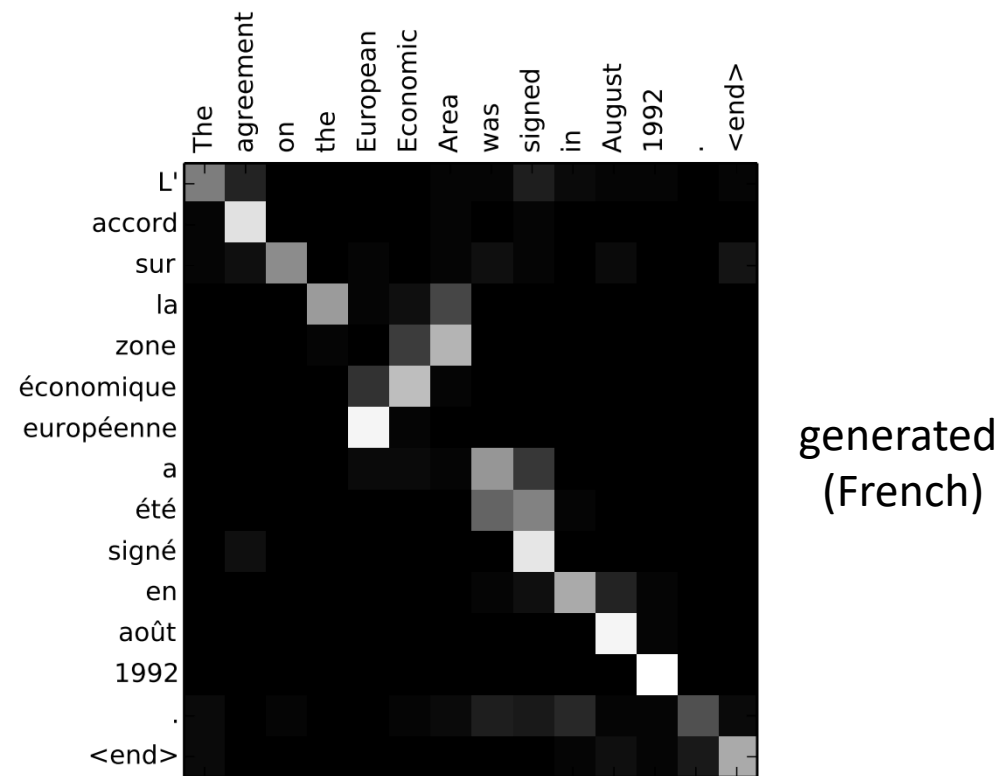
- Major impact on computer vision and NLP from 2014/5

Neural Machine Translation by Jointly Learning to Align and Translate



Bahdanau et al. ICLR 2015

Slides Credit: Peter Anderson



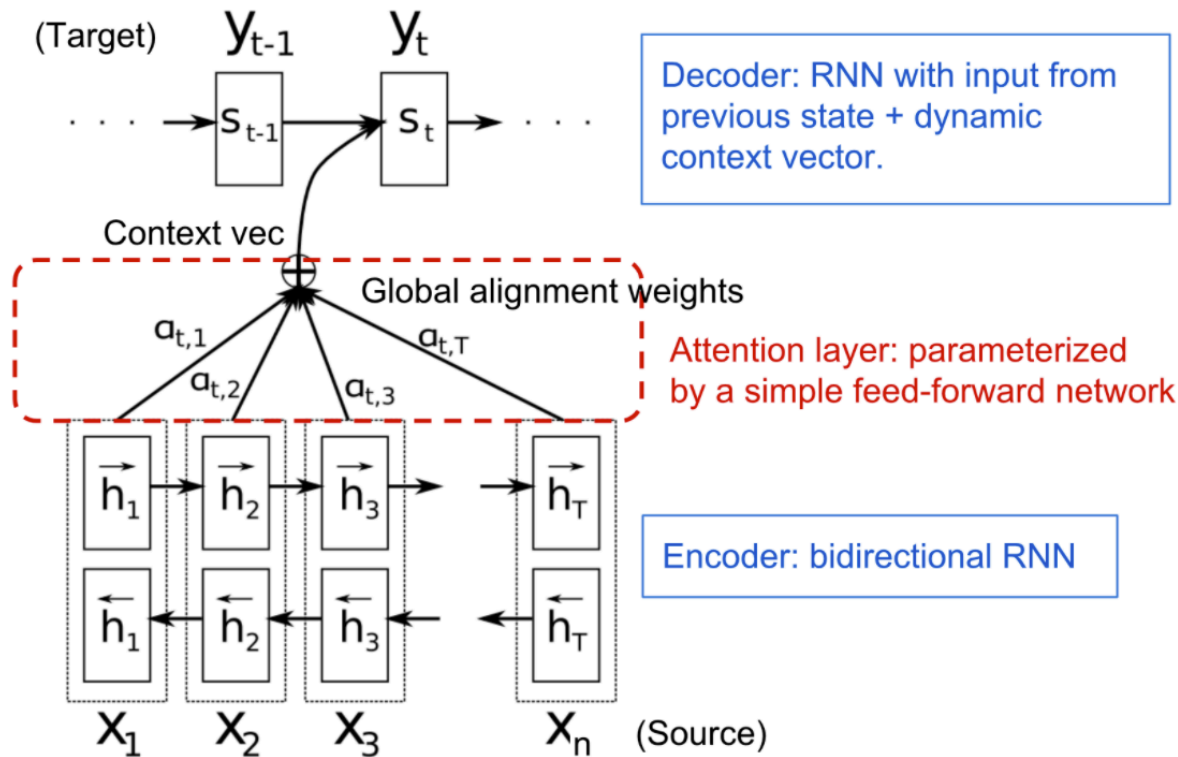
source (English)

Attention weights α_i (0 = black, 1 = white)

Attention

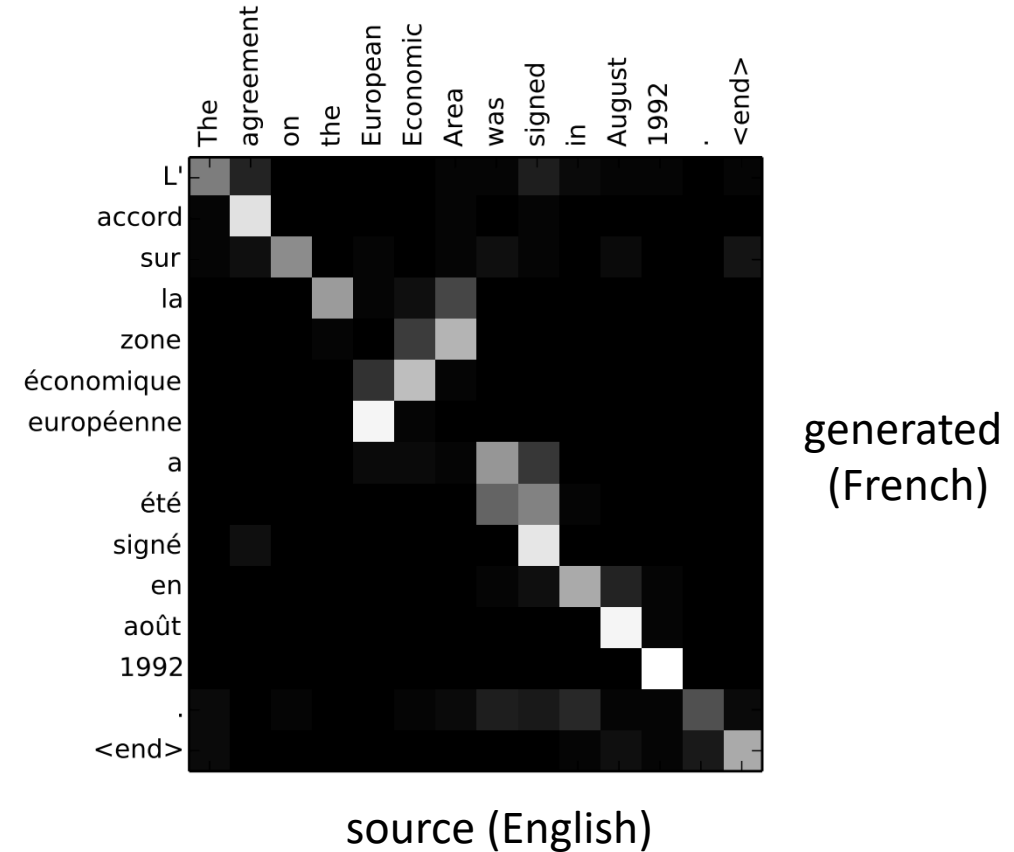
- Major impact on computer vision and NLP from 2014/5

Neural Machine Translation by Jointly Learning to Align and Translate



Bahdanau et al. ICLR 2015

Slides Credit: Peter Anderson

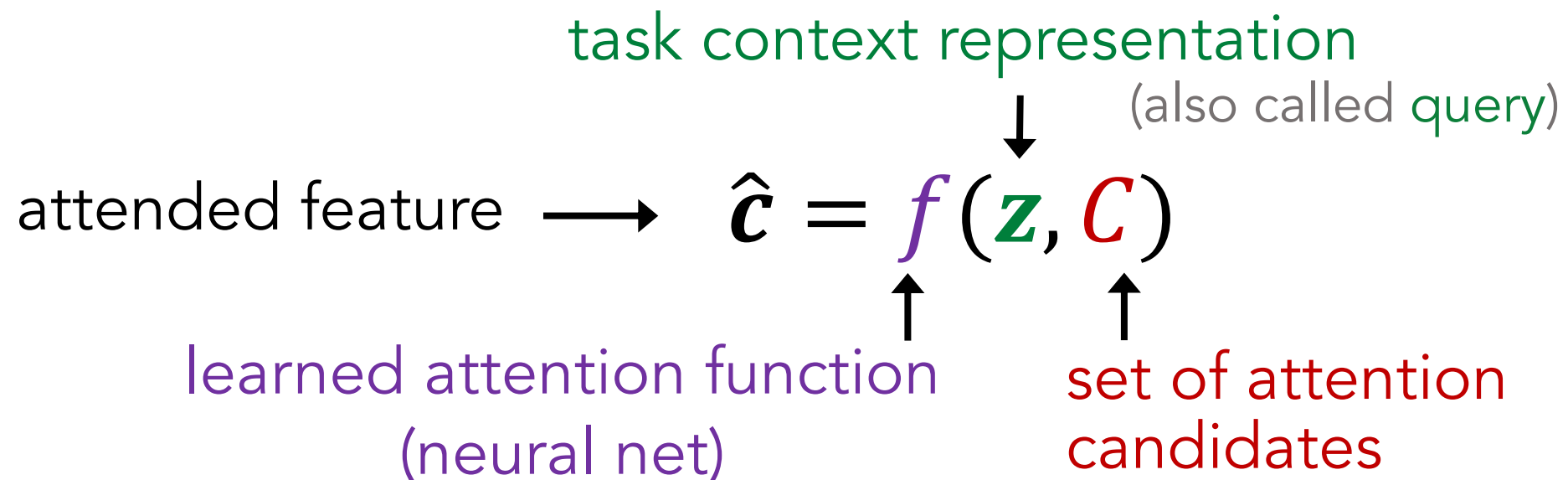


Attention weights α_i (0 = black, 1 = white)

Attention

Attention in Neural Networks:

A learned mechanism that **learns to focus** on a subset of the **input** that is most **relevant to the current task**.



Computing attention

Attention function, f

$$a_i = g(\mathbf{c}_i, \mathbf{z})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$

Attention scores: \mathbf{a} (unnormalized)

Attention weights: $\boldsymbol{\alpha}$ (normalized)

Final attention output

Weighted sum of context features

Attention score $a_i = g(\mathbf{c}_i, \mathbf{z})$

how well does the attention candidate \mathbf{c}_i match the query \mathbf{z}

- Dot-product attention:

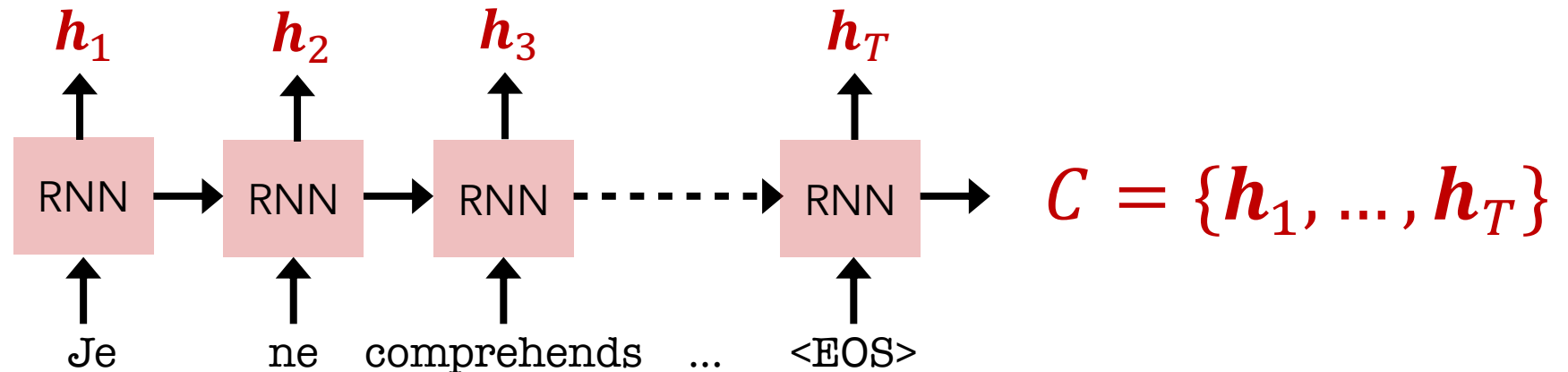
$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{c}_i$$

- Neural network

$$g(\mathbf{c}_i, \mathbf{z}) = v^\top \tanh(W_1 \mathbf{c}_i + W_2 \mathbf{z})$$

Attention over Text

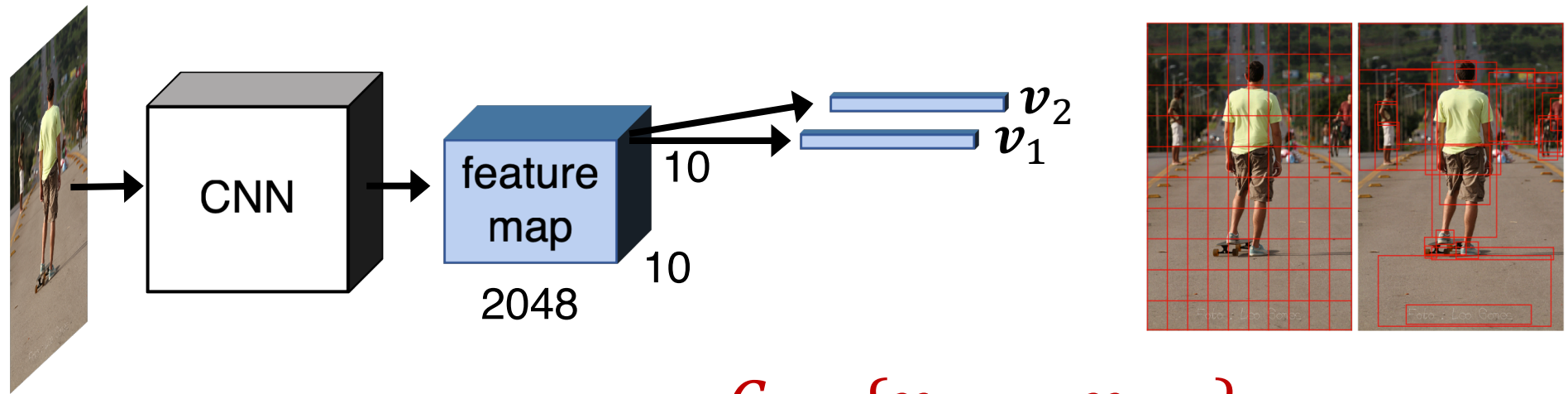
Attention candidates, C typically defined by the hidden states of an encoder (e.g. one feature vector for each word in the input text)



Attention over Visual Features

- Attention candidates, \mathcal{C} typically defined by the spatial output of a CNN (feature vectors for different parts of the image)

Grid based or over object proposals

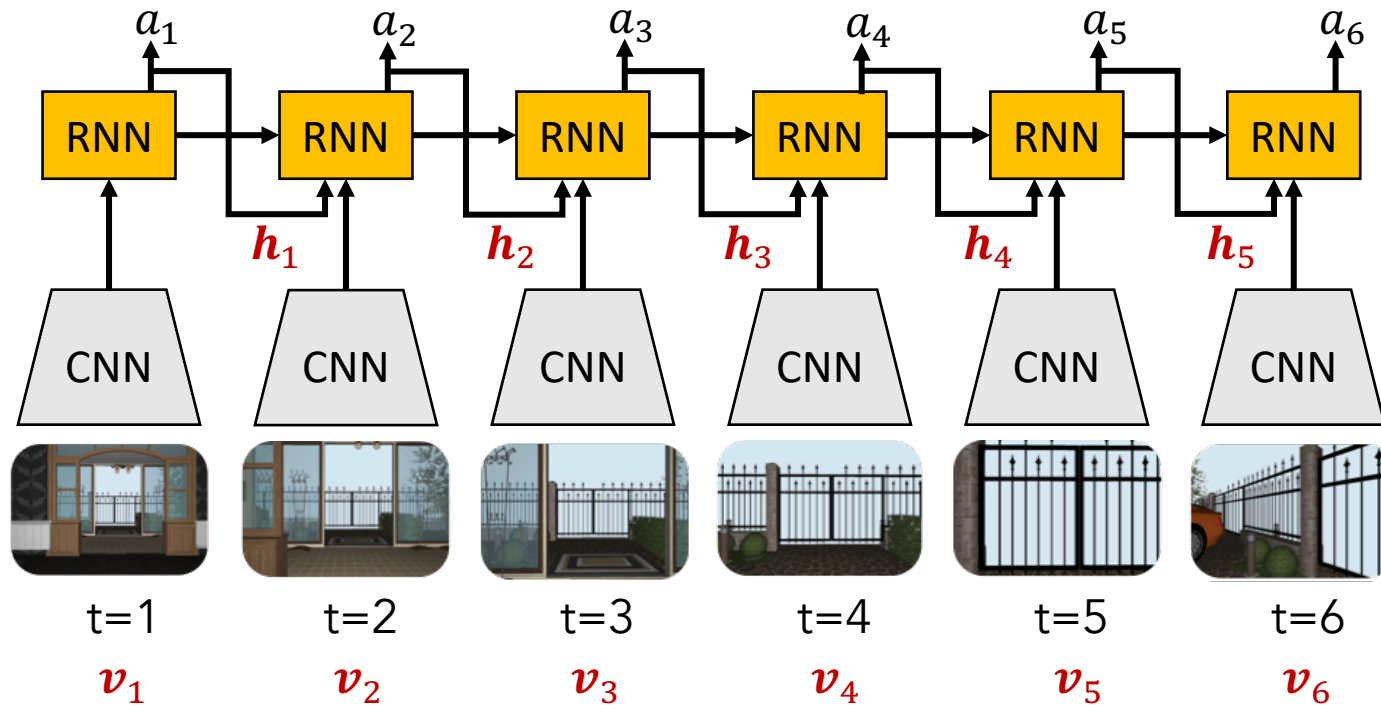


$$\mathcal{C} = \{v_1, \dots, v_{100}\}$$

Attention over Agent Experience

Embodied AI (visual language navigation)

- Attention candidates, \mathcal{C} as agent hidden state or visual vectors



$$\mathcal{C} = \{h_1, \dots, h_5\}$$

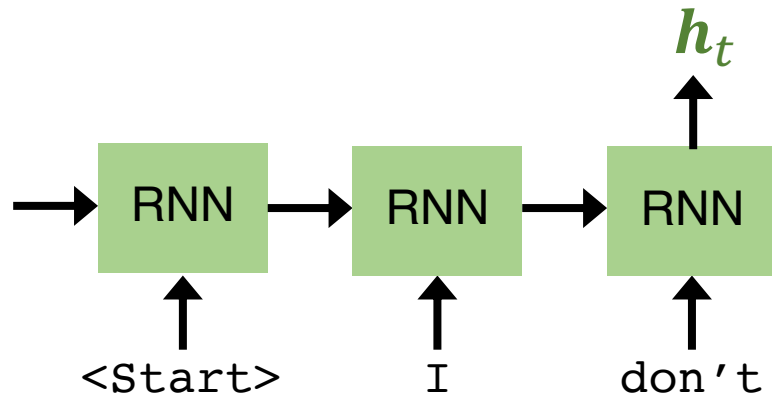
or

$$\mathcal{C} = \{v_1, \dots, v_6\}$$

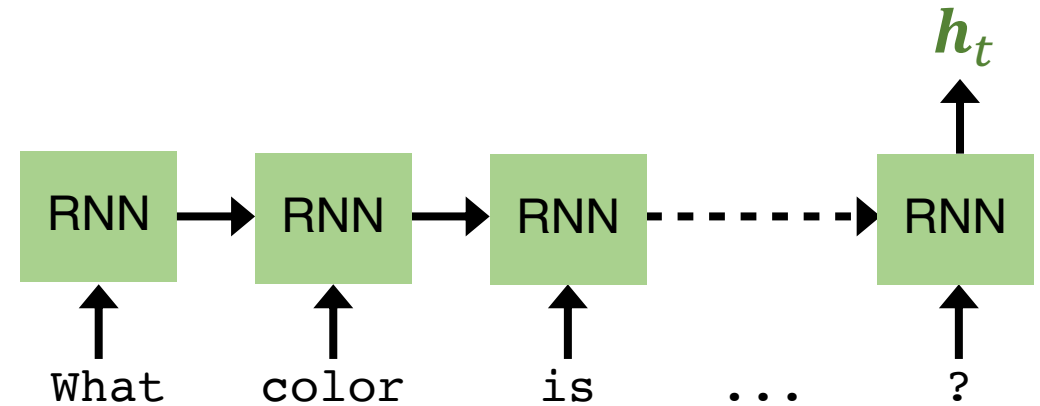
Task context for Attention

- Task context representation, z , is often an RNN encoding

Machine translation / image captioning:
Decoder hidden state

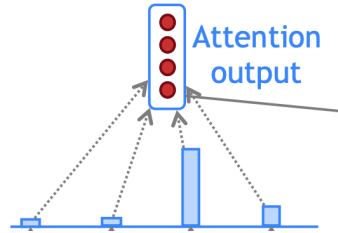


VQA: Question encoding
(final encoder hidden state)



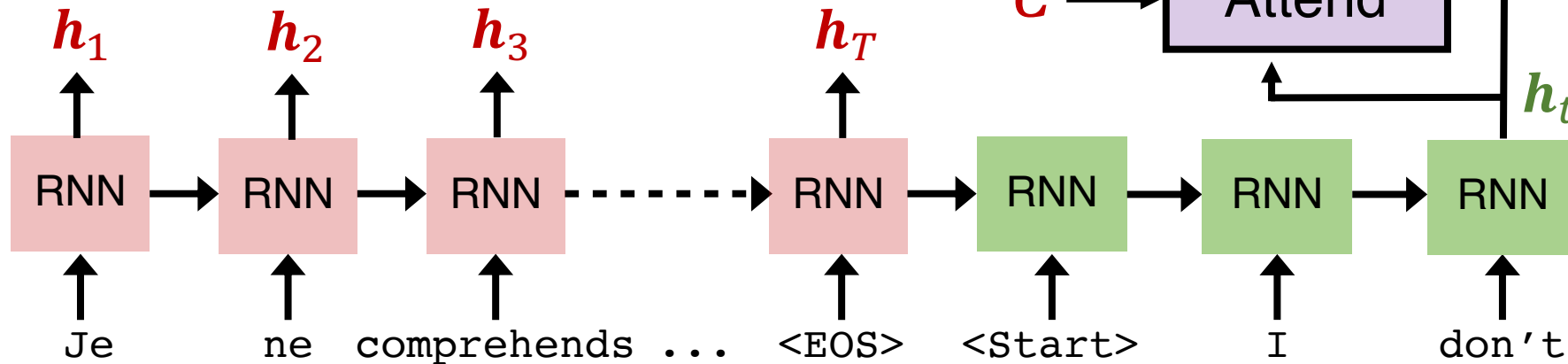
Putting it all together

- Attentive machine translation model



\hat{c}_t is the weighted sum of hidden states

$$C = \{h_1, \dots, h_T\}$$



understand

Predict next word to generate

Softmax

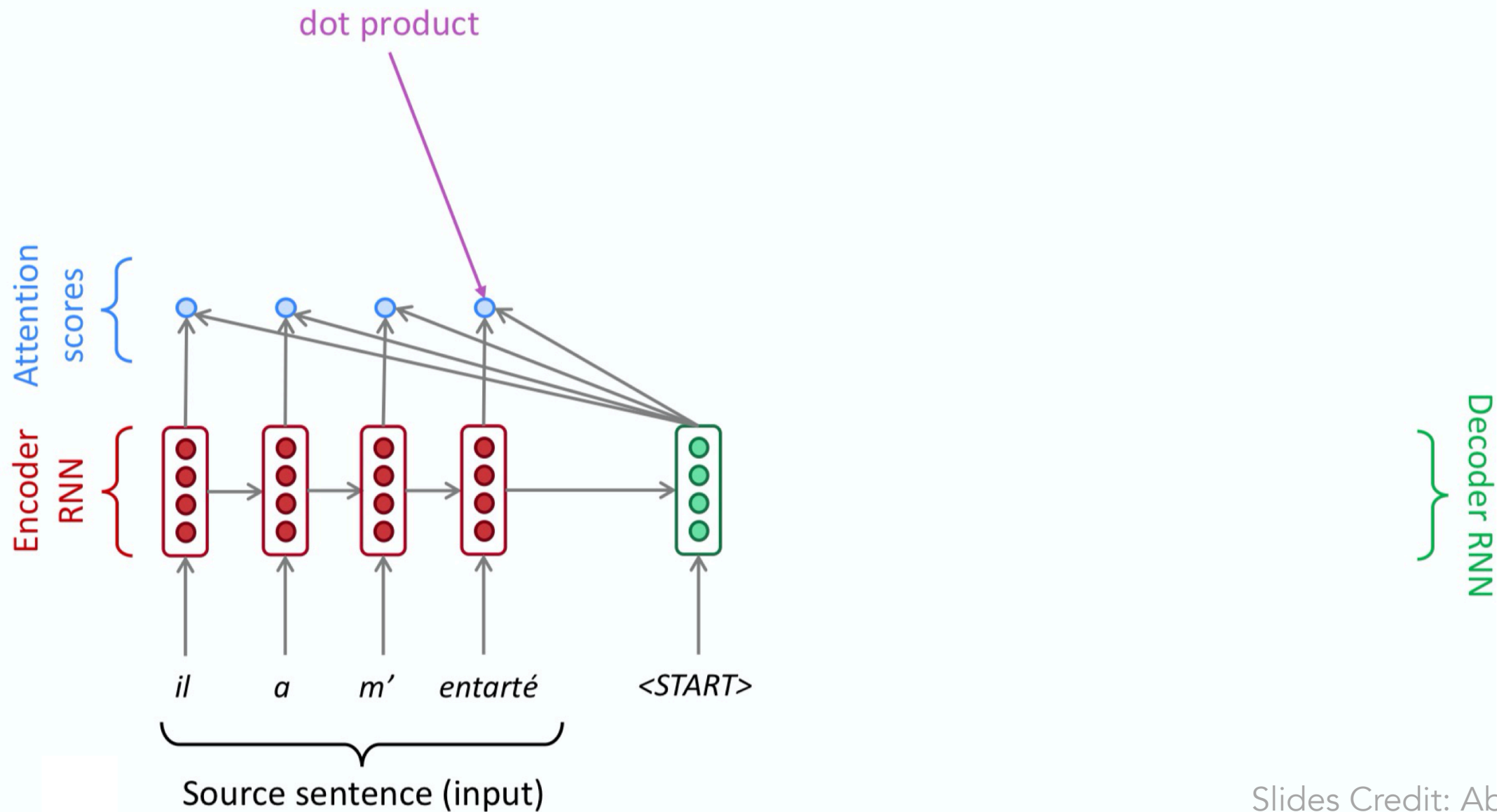
Feedforward Net

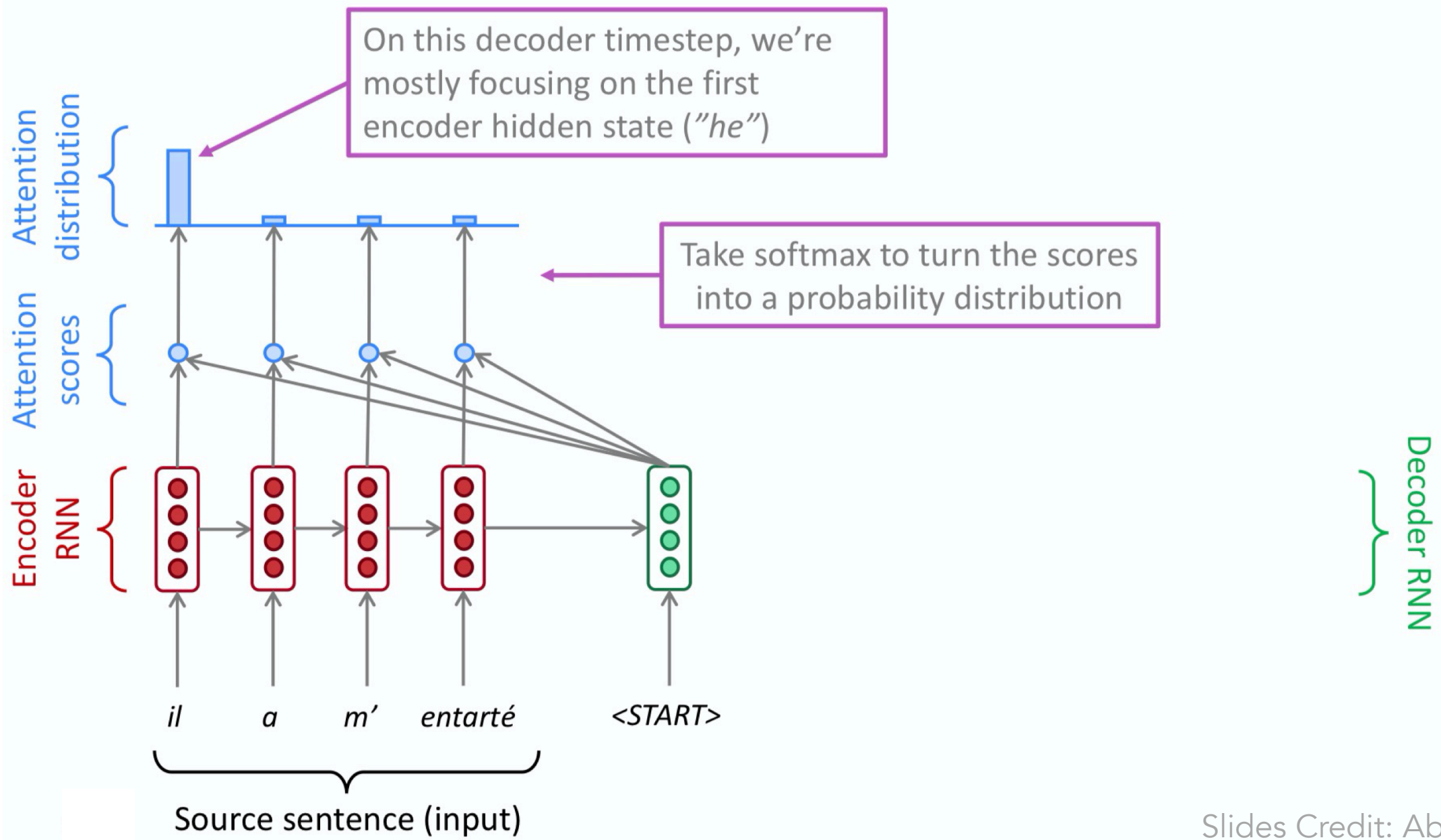
Concatenation

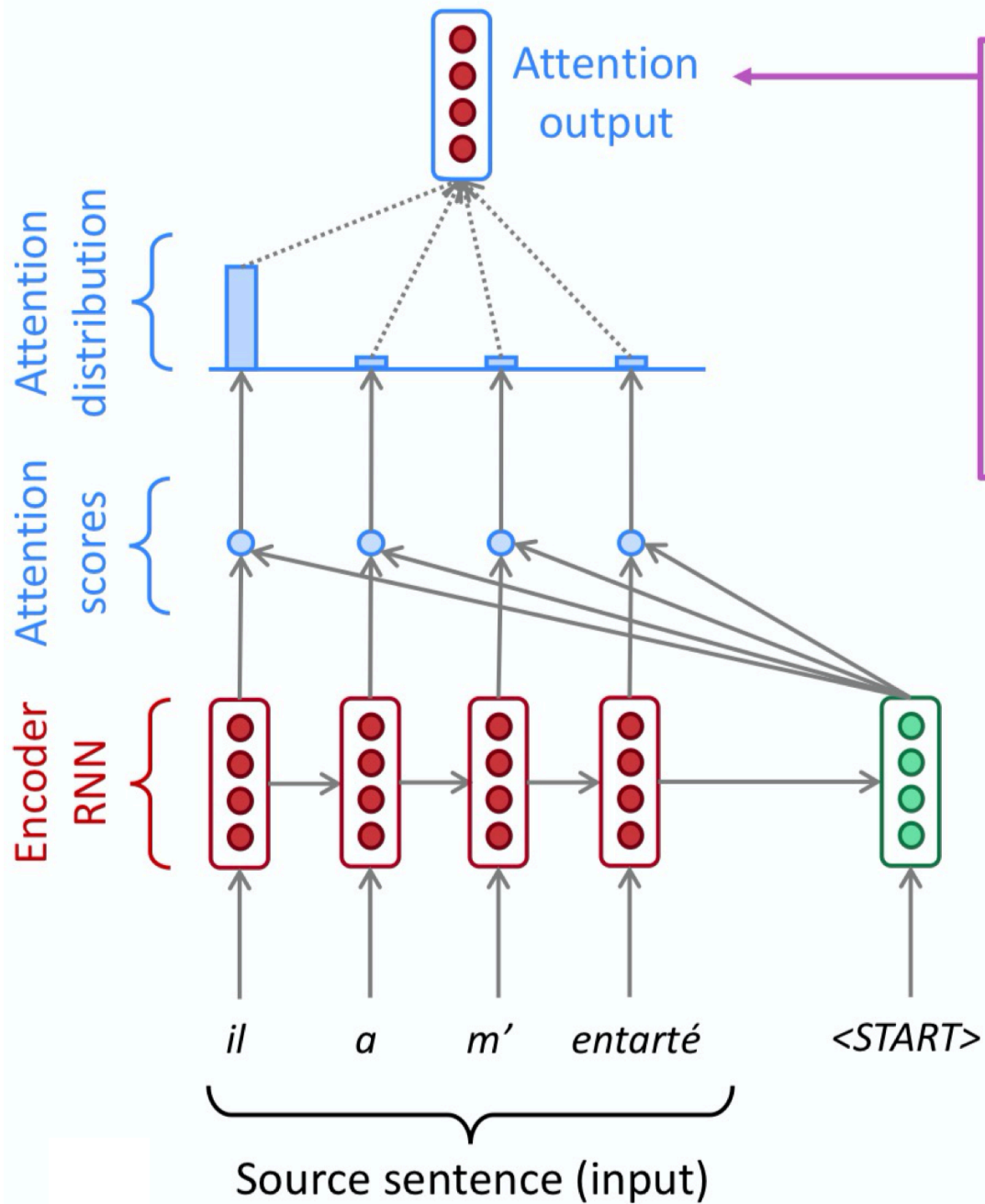
Attend

Interlude: Attention for machine translation

Attention for machine translation

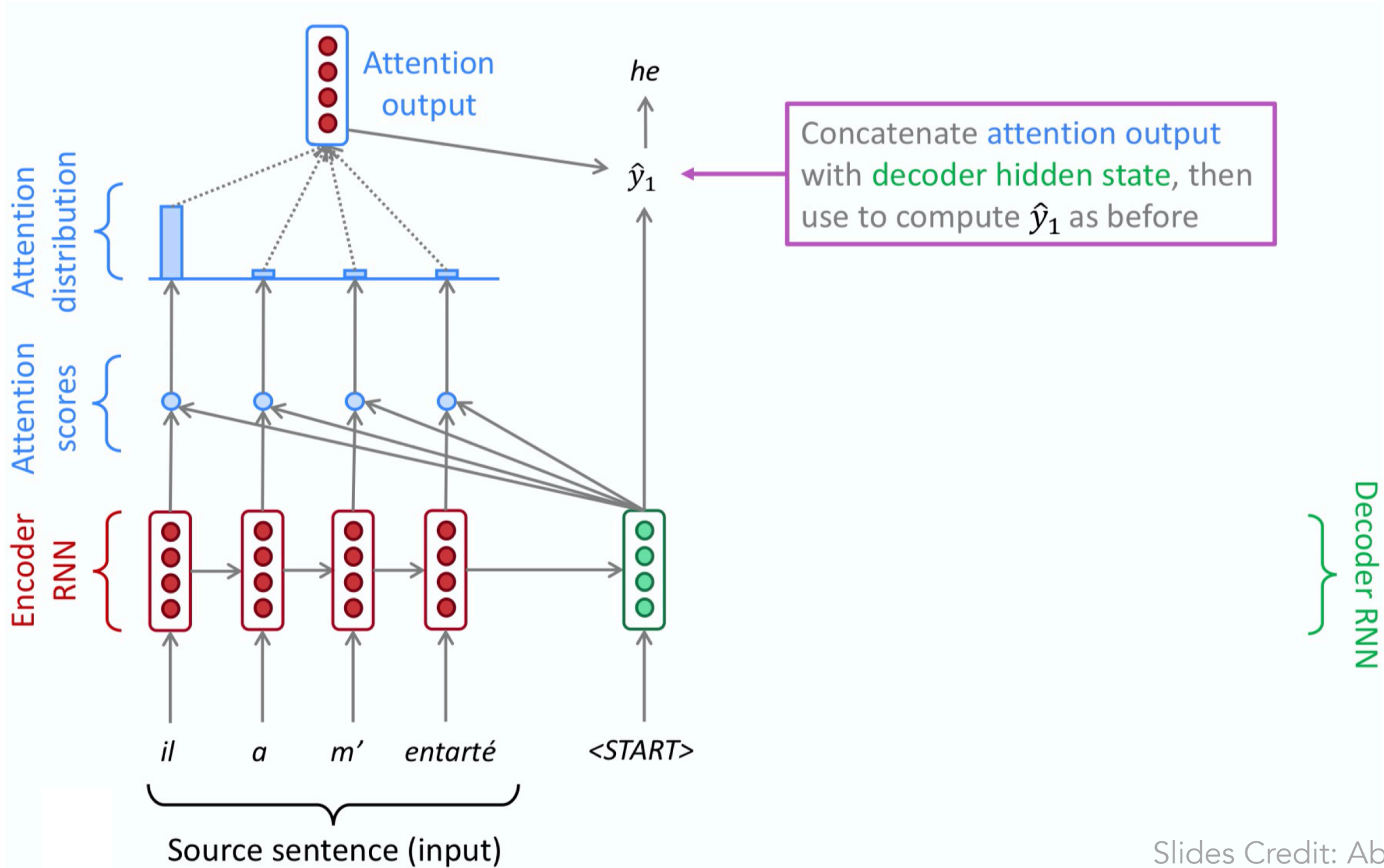


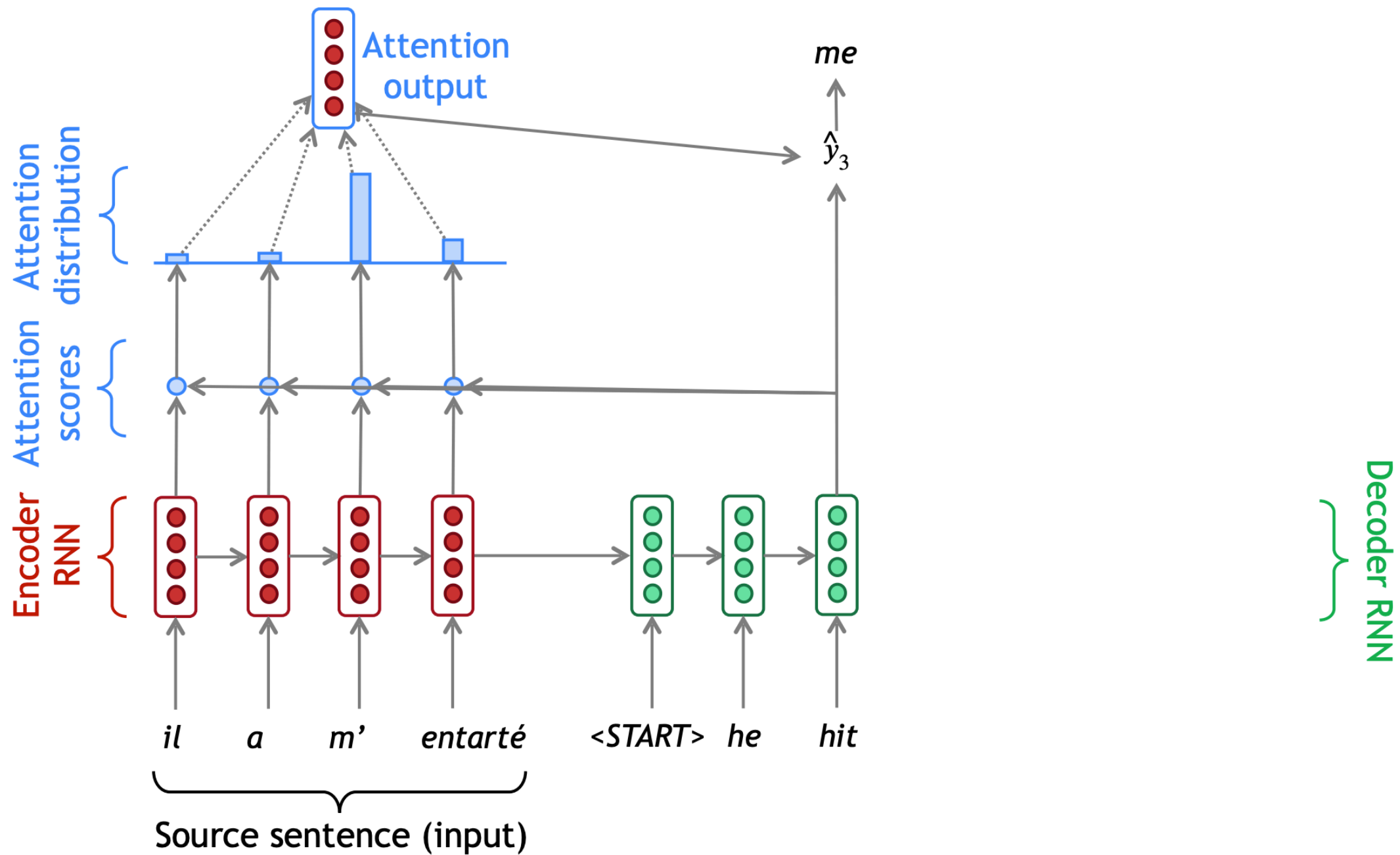


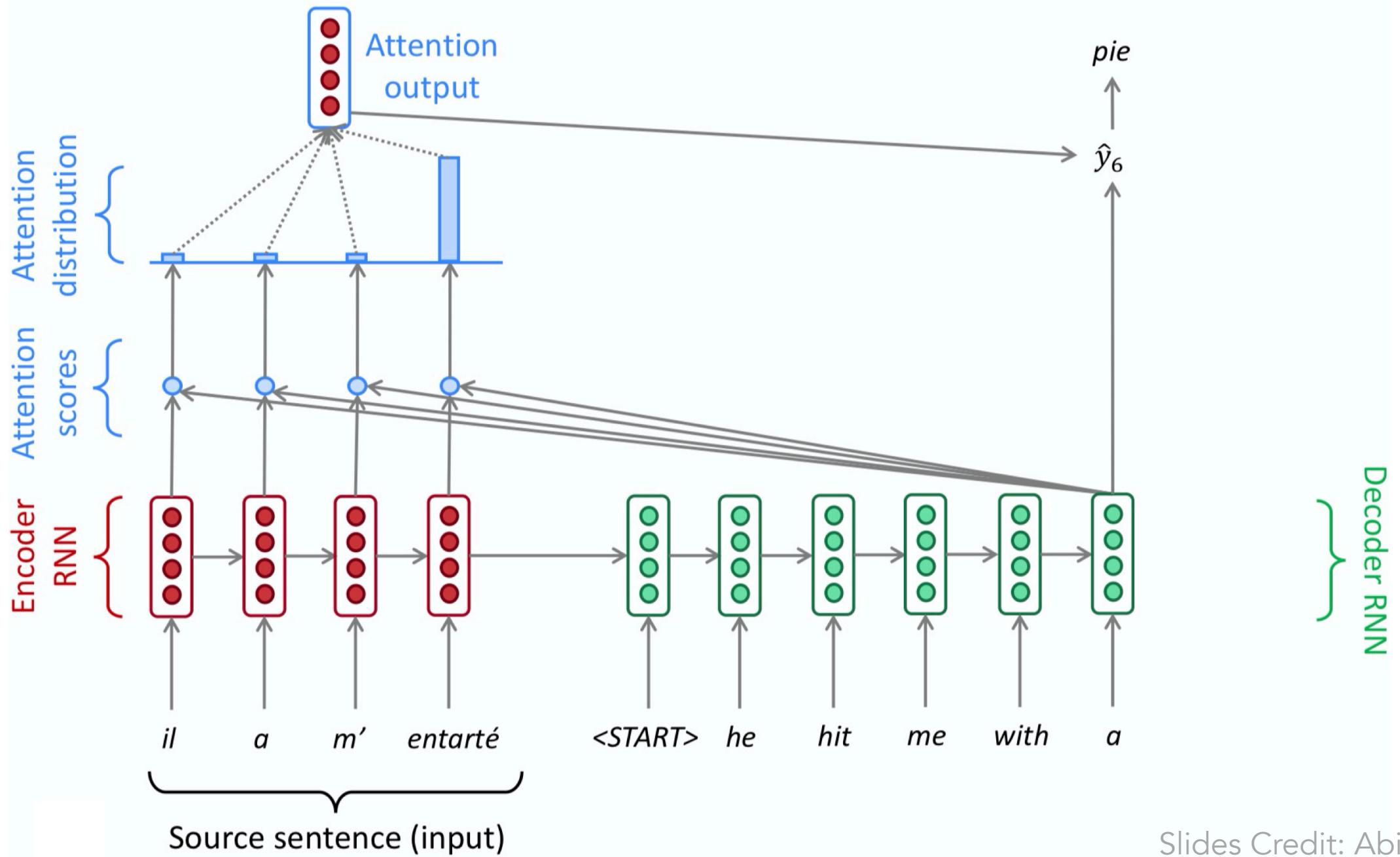


Use the attention distribution to take a **weighted sum** of the **encoder hidden states**.

The **attention output** mostly contains information from the **hidden states** that received **high attention**.

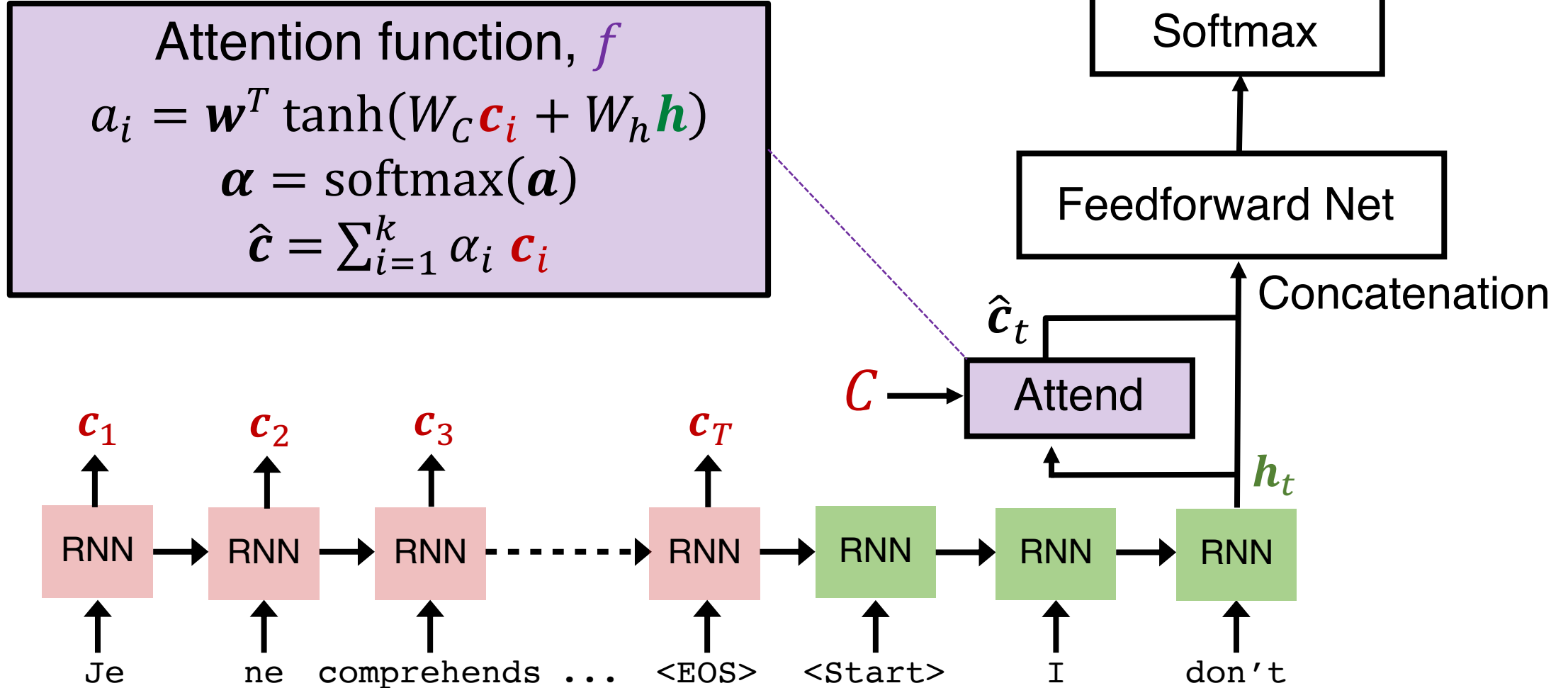






Putting it all together

- Attentive machine translation model



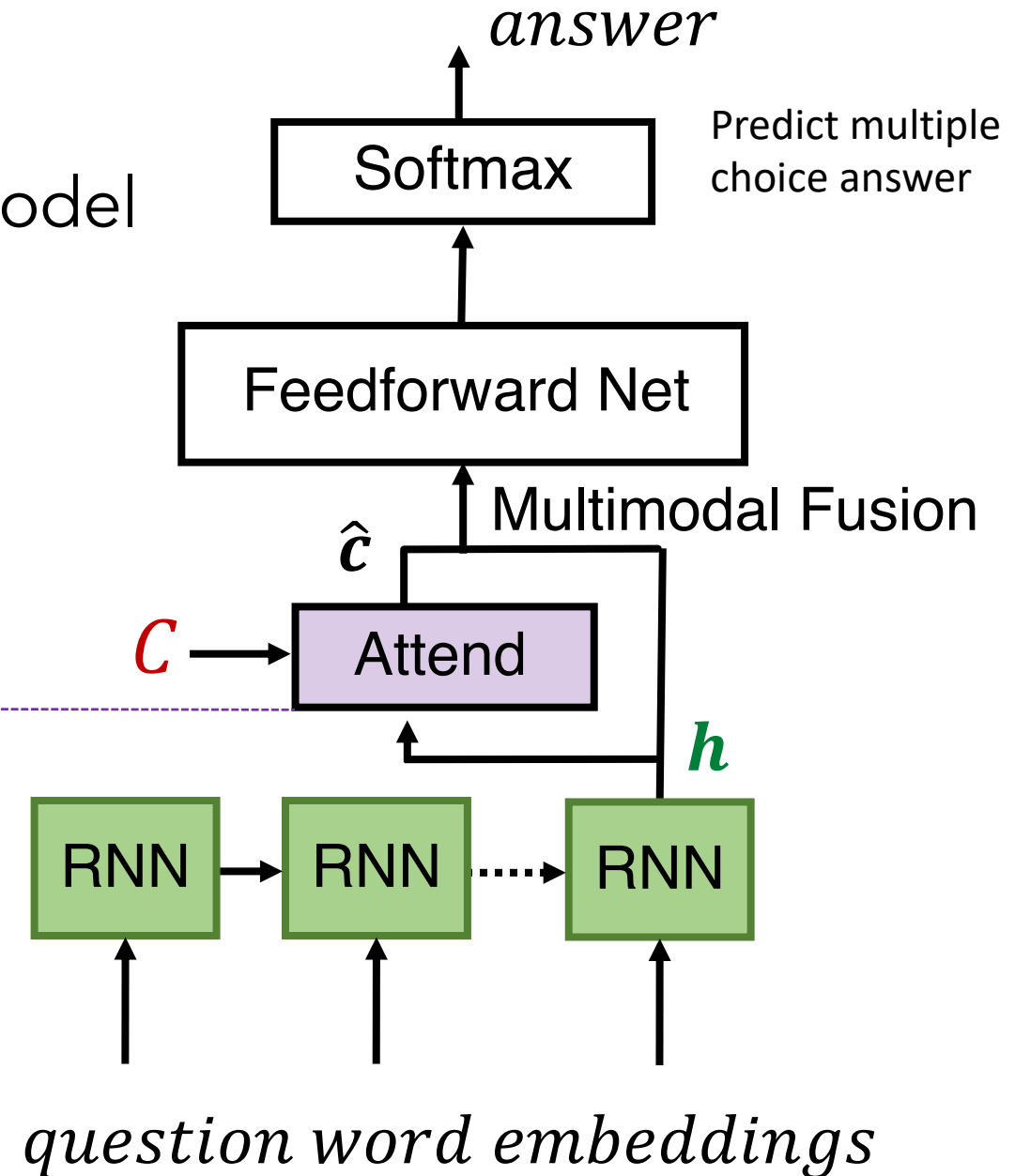
Putting it all together

- Attentive Visual Question Answering model

Attention over visual regions!
 $C = \{v_1, \dots, v_{100}\}$

Attention function, f

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



Putting it all together

- Attentive Image Captioning model

Attention over visual regions!
 $C = \{v_1, \dots, v_{100}\}$

Attention function, f

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

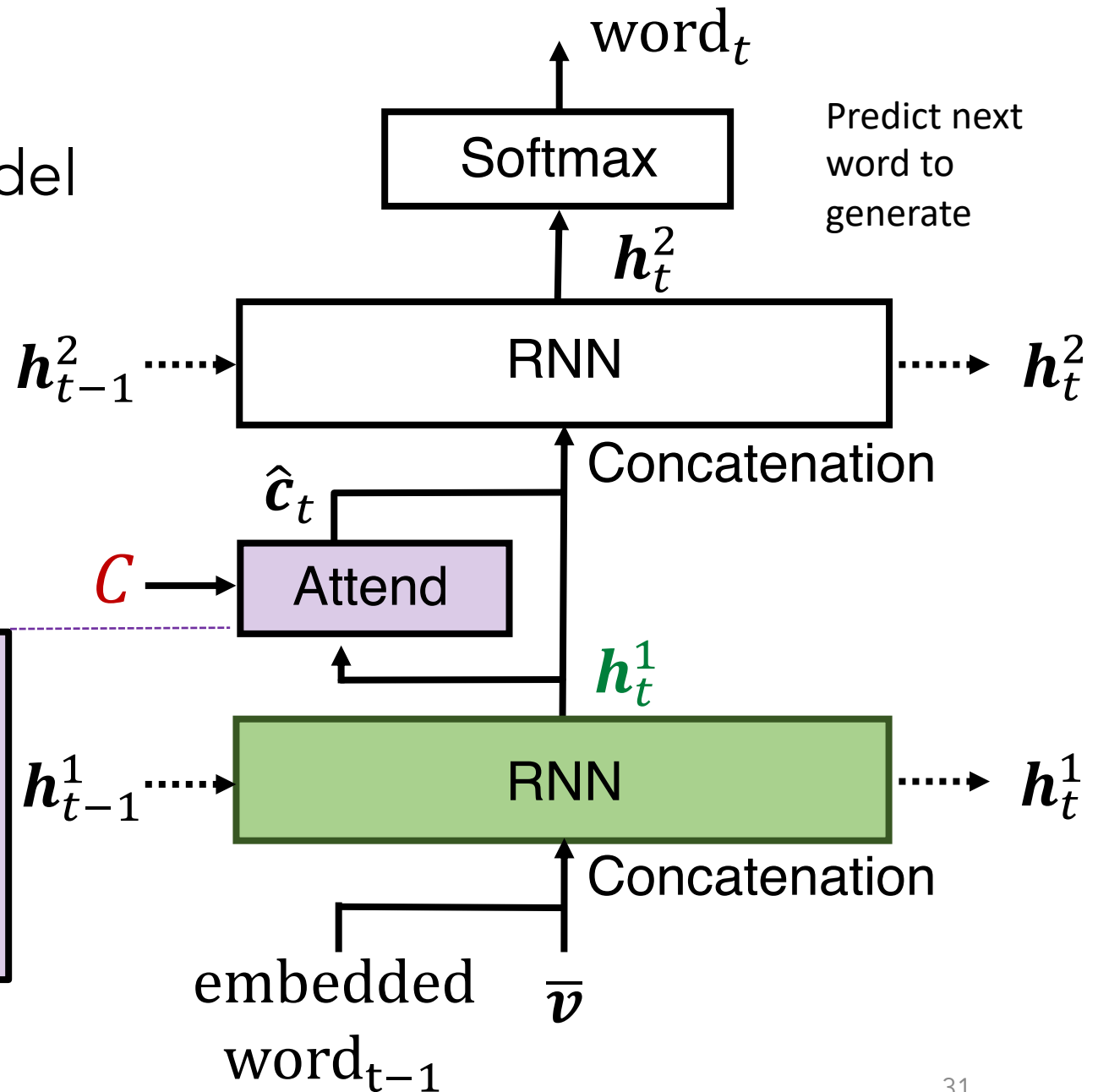
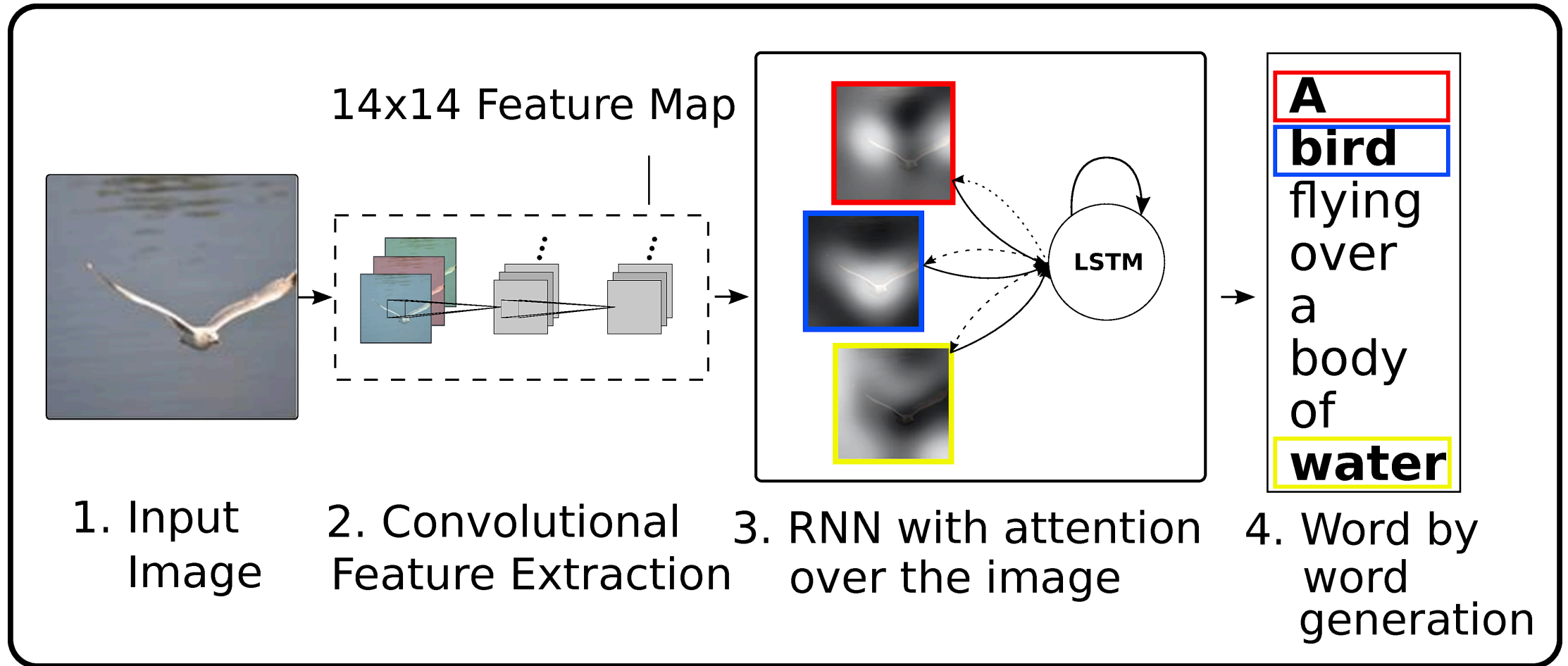


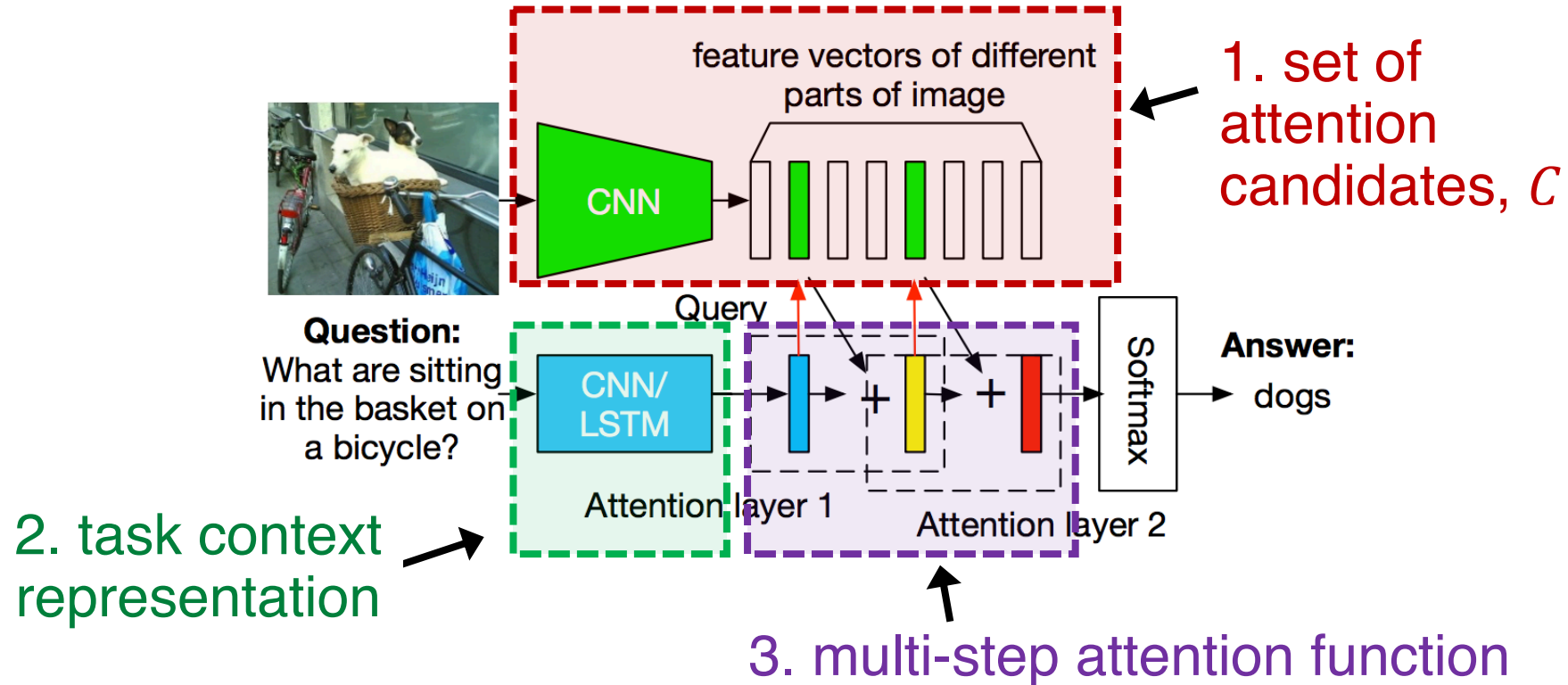
Image captioning example



Xu et al. ICML 2015

Stacked attention networks

- Multiple attentions steps for VQA



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

Stacked attention networks

(a) What are pulling a man on a wagon down on dirt road?
Answer: horses Prediction: horses



(b) What is the color of the box ?
Answer: red Prediction: red



(c) What next to the large umbrella attached to a table?
Answer: trees Prediction: tree



(d) How many people are going up the mountain with walking sticks?
Answer: four Prediction: four

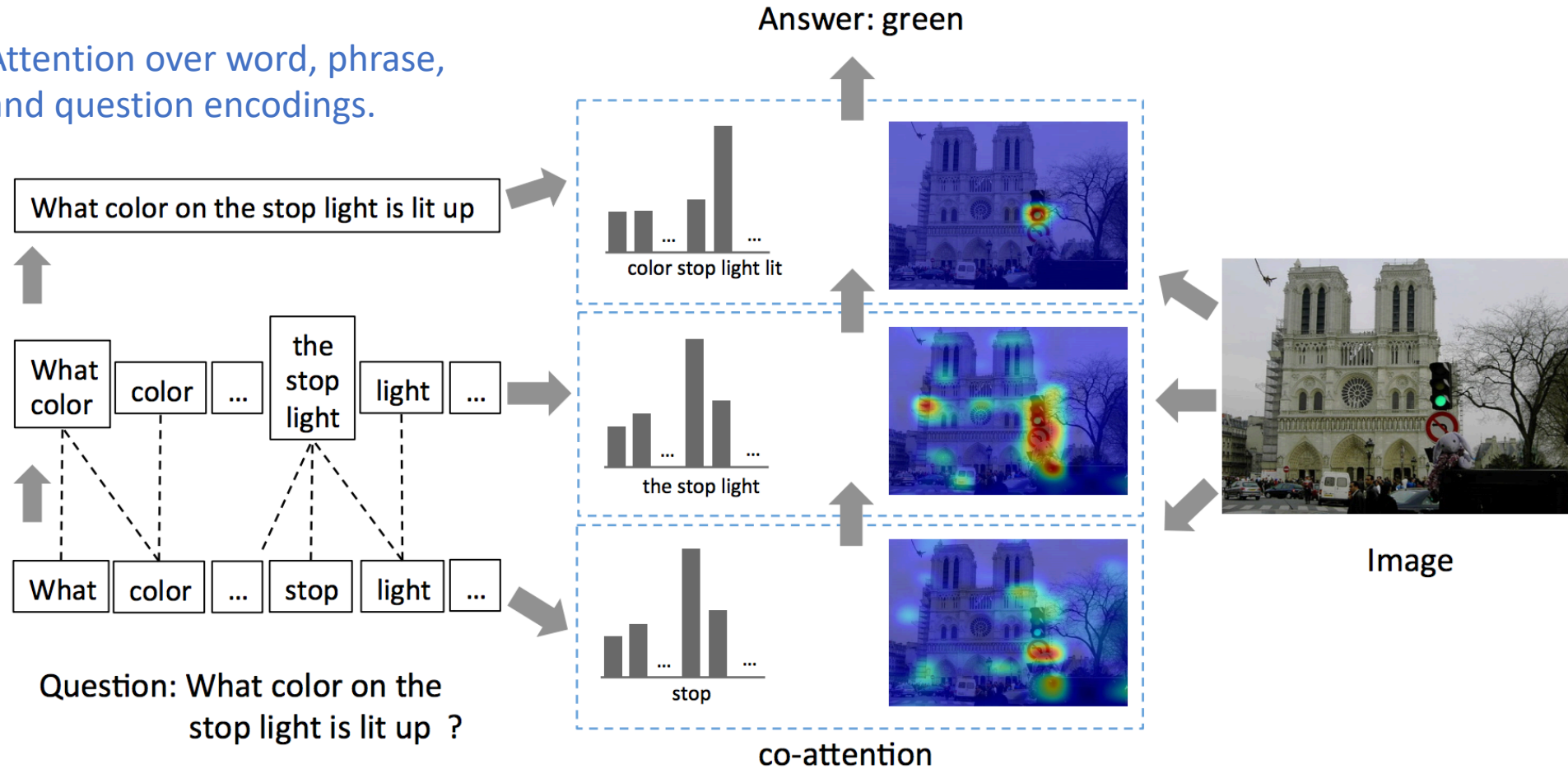


Stacked Attention Networks for Image Question Answering, Yang et al. CVPR 2016

Hierarchical question-image co-attention

- Attending jointly to both question and image in VQA

Attention over word, phrase, and question encodings.

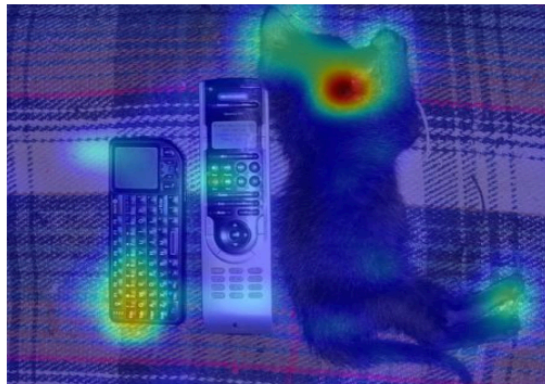


Hierarchical Question-Image Co-Attention for Visual Question Answering, Lu et al. NIPS 2016

Hierarchical question-image co-attention



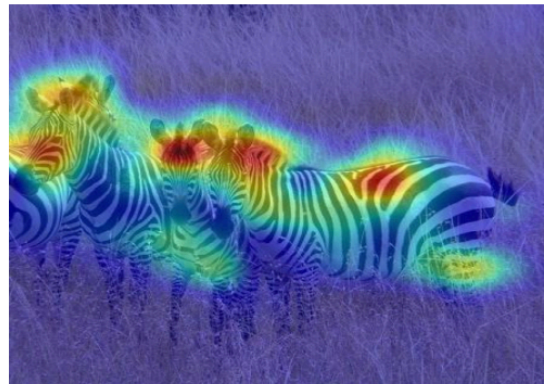
Q: what is the color of the kitten? A: **black**



what is the color of the kitten ?



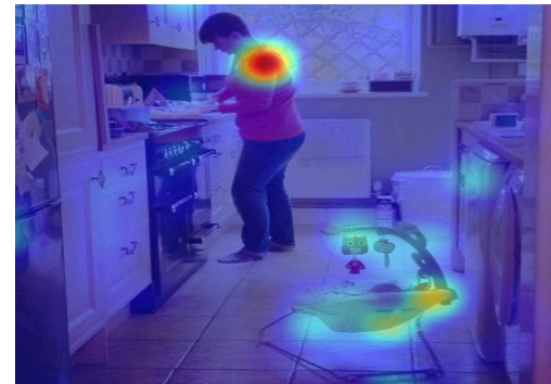
Q: what are standing in tall dry grass look at the tourists? A: **zebras**



what are standing in tall dry grass look at the tourists ?



Q: where is the woman while her baby is sleeping? A: **kitchen**



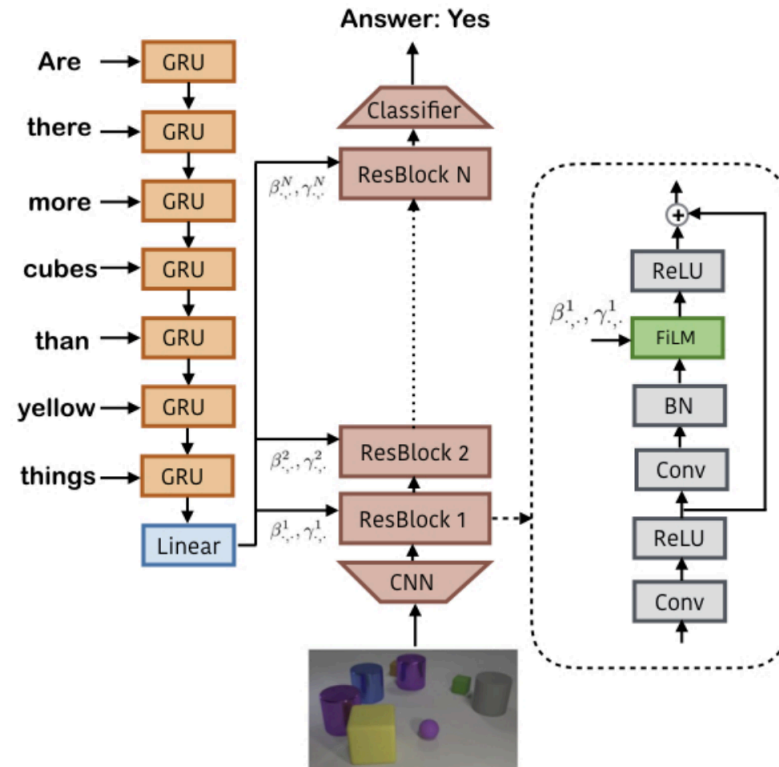
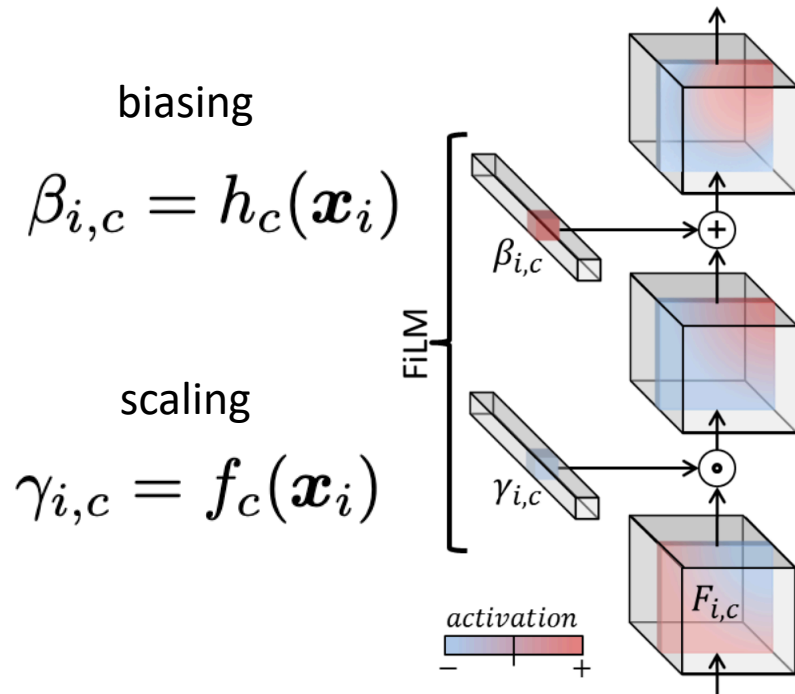
where is the woman while her baby is sleeping ?

Hierarchical Question-Image Co-Attention for Visual Question Answering, Lu et al. NIPS 2016

FiLM: Feature-wise Linear Modulation

- Applying attention by scaling and biasing CNN layers

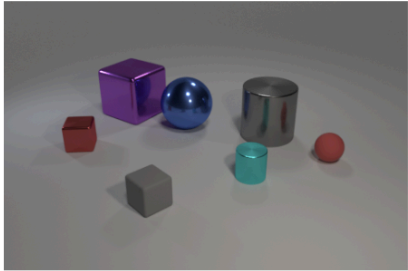
Feature-wise transformations are learned functions of input



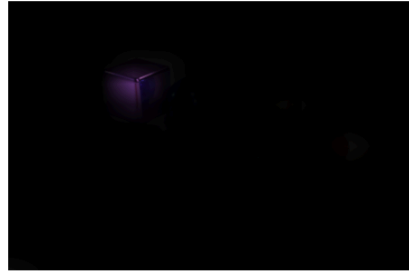
FiLM: Visual Reasoning with a General Conditioning Layer, Perez et al. AAAI 2018

FiLM: Feature-wise Linear Modulation

Q: What shape is the...



...purple thing? **A:** cube



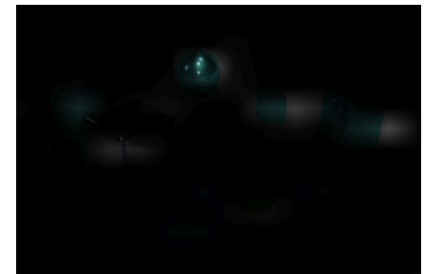
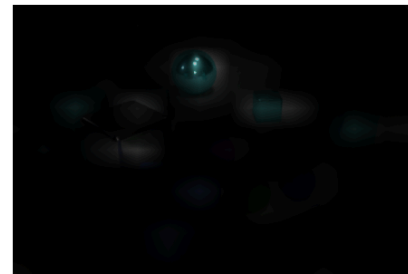
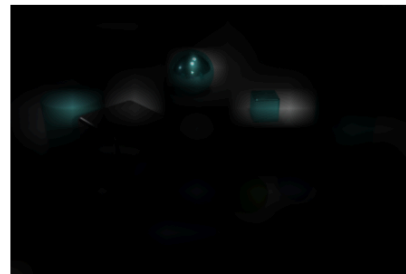
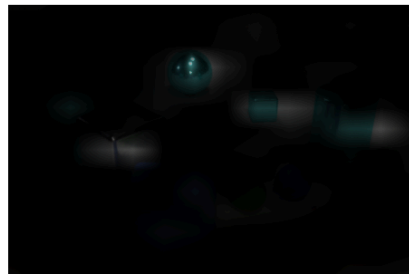
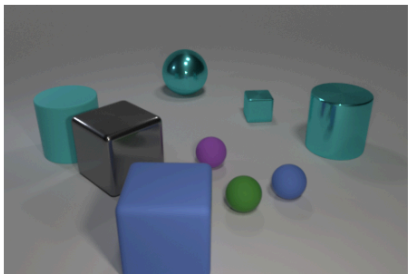
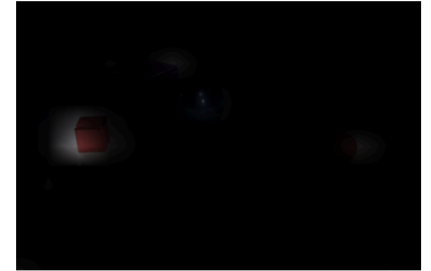
...blue thing? **A:** sphere



...red thing right of the blue thing? **A:** sphere



...red thing left of the blue thing? **A:** cube



Q: How many cyan things are...

...right of the gray cube?
A: 3

...left of the small cube?
A: 2

...right of the gray cube and left of the small cube? **A:** 1

...right of the gray cube or left of the small cube? **A:** 4 (**P:** 3)

'Soft' vs 'hard' attention

- **Soft:** Each attention candidate is weighted by α_i

$$\hat{v} = \sum_{i=1}^k \alpha_i \mathbf{v}_i$$

- Easy to train (smooth and differentiable)
- But can be expensive over large input

- **Hard:** Use α_i as a sample probability to pick *one* attention candidate as input to subsequent layers
 - Trainable with REINFORCE approaches (Xu et al. ICML 2015), or Gumbel-Softmax (Jang et al. ICLR 2017)



Soft



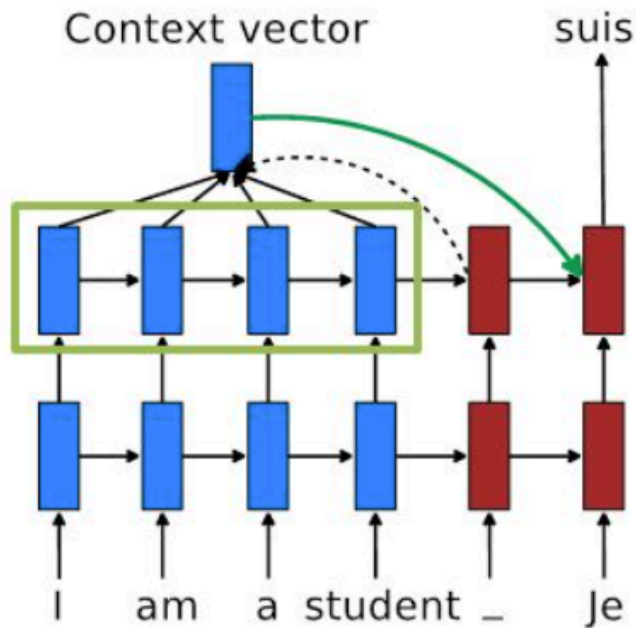
Hard

bird

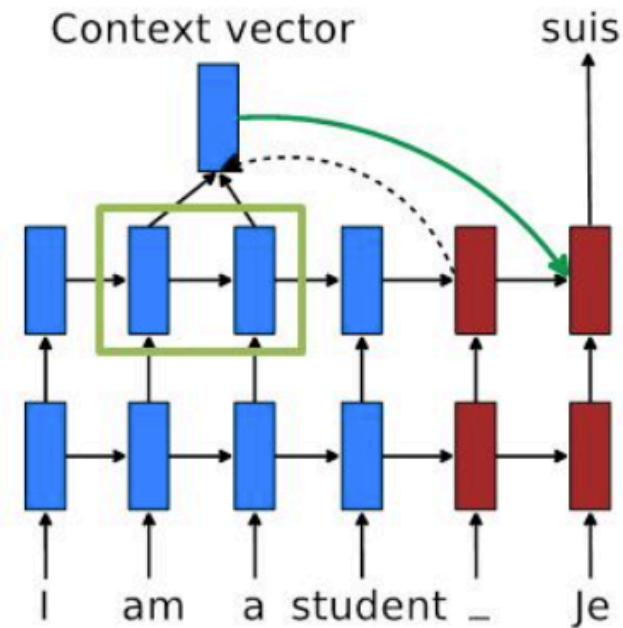
Xu et al. ICML 2015

'Global' vs 'local' attention

- **Global**: attention over the entire input
- **Local**: attention over a window (or subset) of the input



Global: ***all*** source states.



Local: ***subset*** of source states.

Luong et al, 2015

Attention

Attention According to Cognitive Psychology / Neuroscience

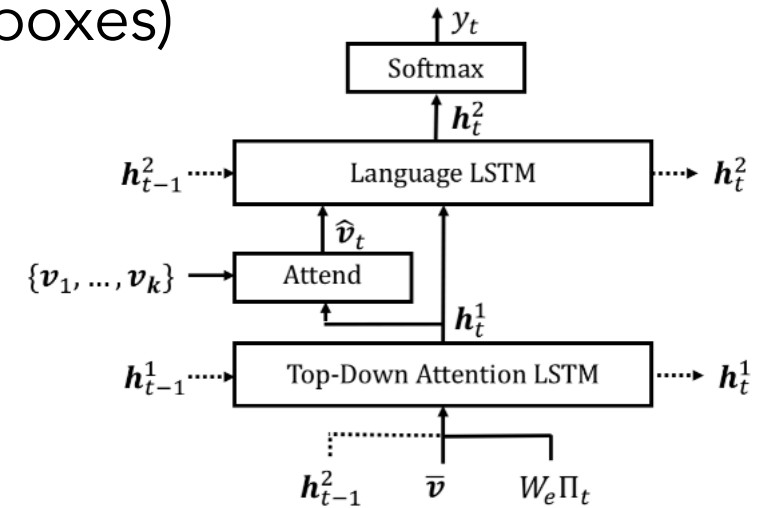
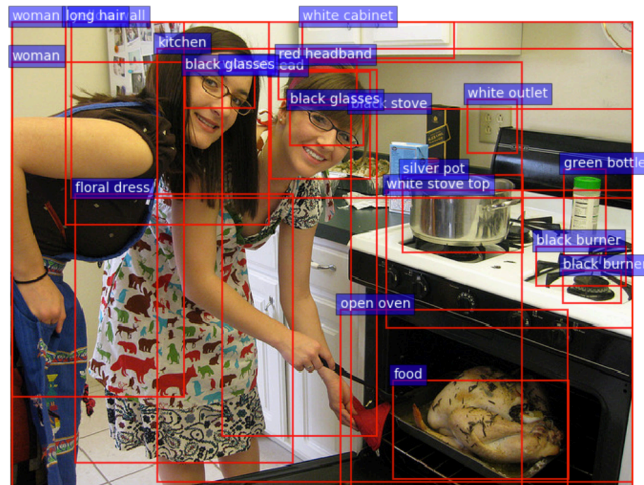
- A set of mechanisms that limit some processing to a subset of incoming stimuli (reducing computational demands)
- Can be driven 'top-down' by task demands (i.e. volitionally)
- Can be driven 'bottom-up' by salient stimuli (i.e. involuntarily)
- Visual attention can be applied to features, objects and spatial regions, as well as temporal cues (anticipating events)

Buschman and Miller 2007, Scholl 2001

Bottom-up and top-down attention

Combines 'top-down' and 'bottom-up' attention

- 'top-down' attention = *soft* attention over image conditioned on the task
 - VQA: question
 - Image captioning: what has been output before
- 'bottom-up' attention = *hard* attention using Faster-RCNN to identify image regions (object bounding boxes)



Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering, Anderson et al. CVPR 2018

Bottom-up and top-down attention

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



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



Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering, Anderson et al. CVPR 2018

Bottom-up and top-down attention

Q: What color is illuminated on the traffic light?



Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering, Anderson et al. CVPR 2018

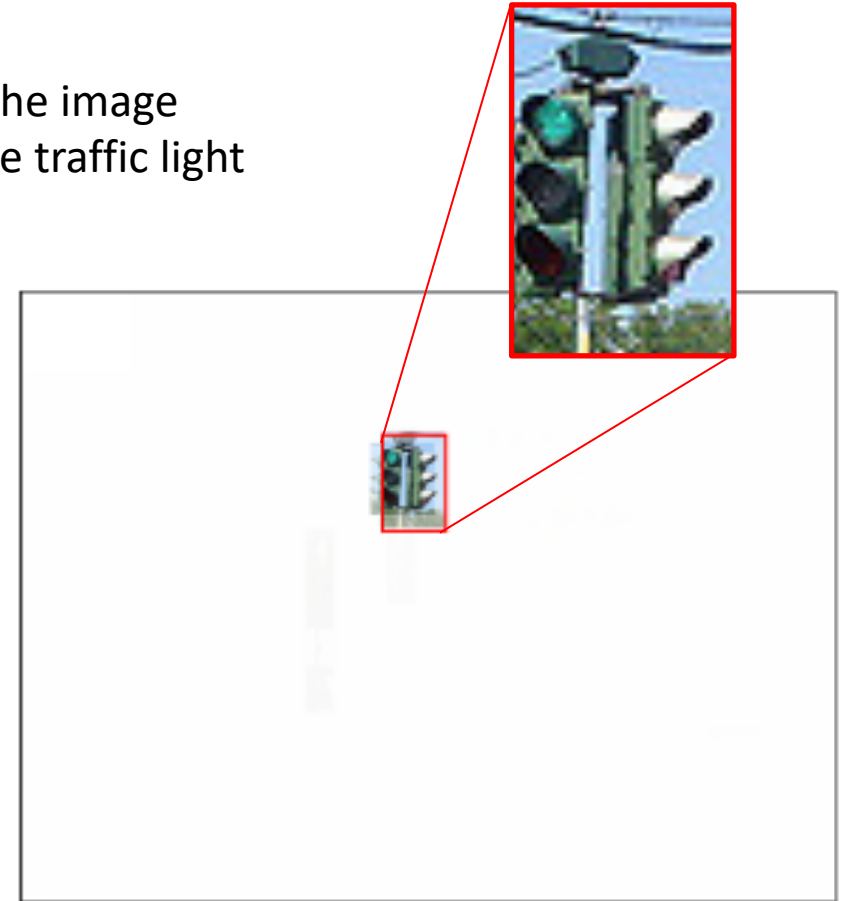
Bottom-up and top-down attention

Q: What color is illuminated on the traffic light?

A: **green**



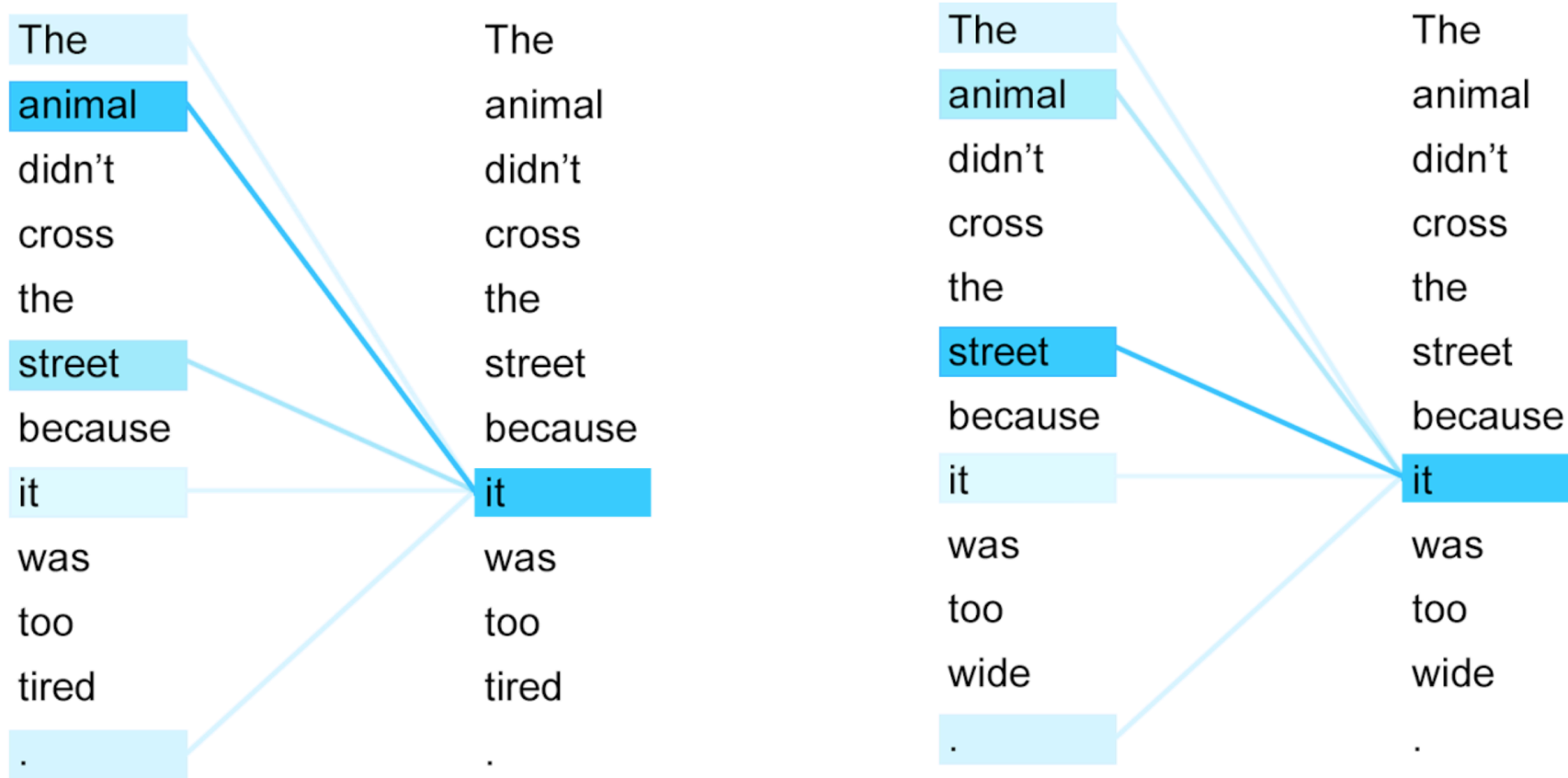
Focus on region of the image corresponding to the traffic light



Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering, Anderson et al. CVPR 2018

Self-attention

- Attention (correlation) with different parts of itself



- Transformers: modules with scaled dot-product self-attention

Types of attention scores

Attention function, f

$$a_i = g(\mathbf{c}_i, \mathbf{z})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$

- Dot-product attention:

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{c}_i$$

- Scaled dot-product attention:

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{c}_i / \sqrt{d}$$

- Bilinear / multiplicative attention:

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{z}^\top \mathbf{W} \mathbf{c}_i \in \mathbb{R}$$

where \mathbf{W} is a weight matrix

- Additive attention (essentially MLP):

$$g(\mathbf{c}_i, \mathbf{z}) = \mathbf{v}^\top \tanh(\mathbf{W}_1 \mathbf{c}_i + \mathbf{W}_2 \mathbf{z})$$

where $\mathbf{W}_1, \mathbf{W}_2$ are weight matrices and \mathbf{v} is a weight vector

Query-key-value view of attention

Attention function, f

$$a_i = g(\mathbf{c}_i, \mathbf{z})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{c}_i$$



Attention function, f

$$a_i = g(\mathbf{k}_i, \mathbf{q})$$

$$\boldsymbol{\alpha} = \text{softmax}(\mathbf{a})$$

$$\hat{\mathbf{c}} = \sum_{i=1}^k \alpha_i \mathbf{v}_i$$

Projected query, key, value



$$\mathbf{q} = W_Q \mathbf{z}$$

$$\mathbf{k}_i = W_K \mathbf{c}_i$$

$$\mathbf{v}_i = W_V \mathbf{c}_i$$



Matrix form

$$\mathbf{q} = W_Q \mathbf{z}$$

$$K = W_K C^T$$

$$V = W_V C^T$$

More next Thursday when we discuss transformers

Next week

- Monday: Paper presentations and discussions
 - Show, attend, and tell (captioning)
 - Guest presenter: Ali Gholami
 - MattNet (referring expressions)
- Thursday: Pretraining with Transformers