

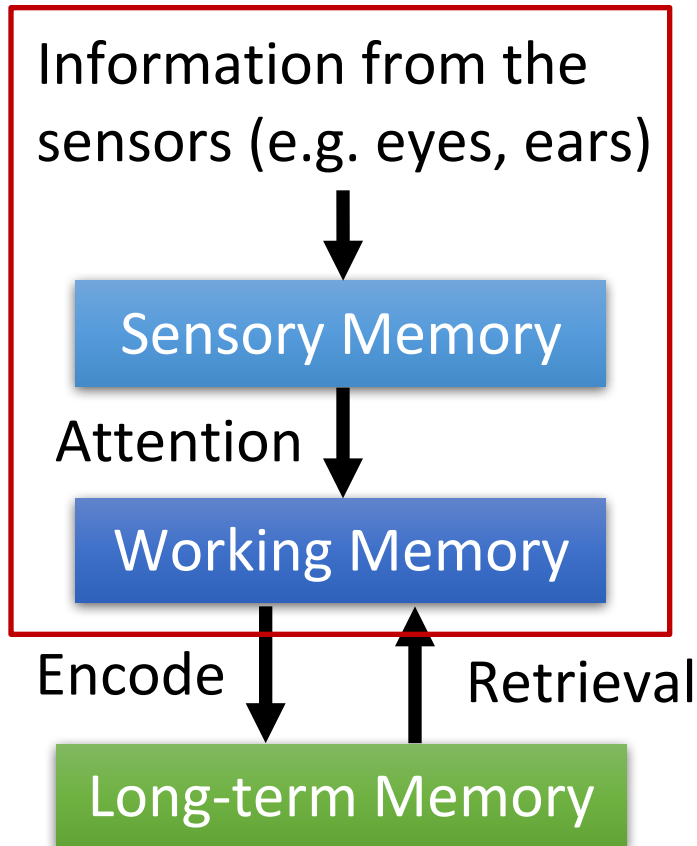
CS 291A: Deep Learning for NLP

Neural Networks: Attention and Memory

William Wang
UCSB Computer Science
william@cs.ucsb.edu

Slides adapted from V. Chen and H. Lee.

Attention and Memory

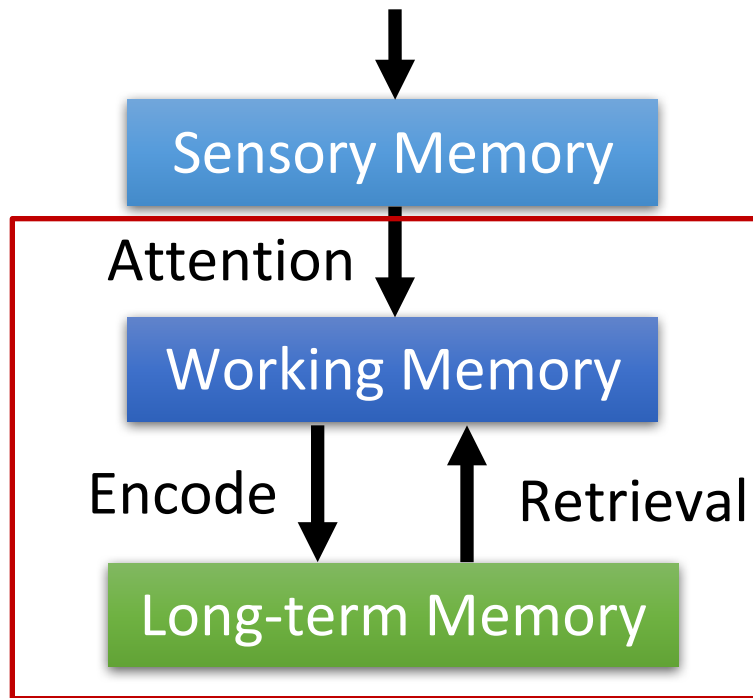


When the input is a very long sequence or an image

➡ Pay attention on partial of the input object each time

Attention and Memory

Information from the sensors (e.g. eyes, ears)



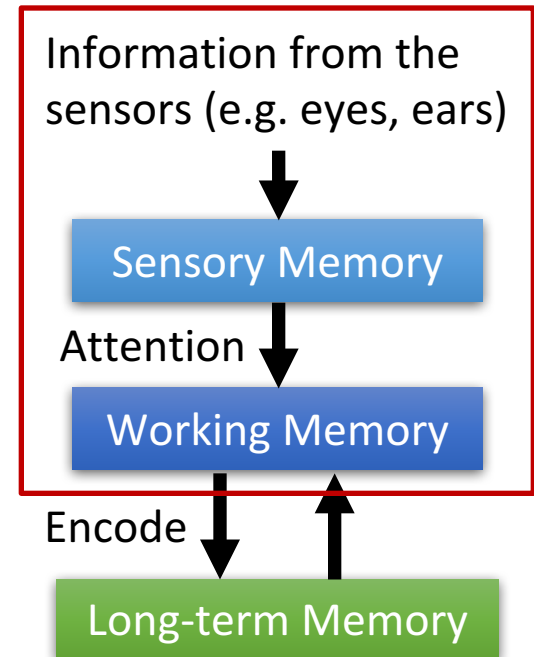
When the input is a very long sequence or an image

➡ Pay attention on partial of the input object each time

In RNN/LSTM, larger memory implies more parameters

➡ Increasing memory size will not increasing parameters

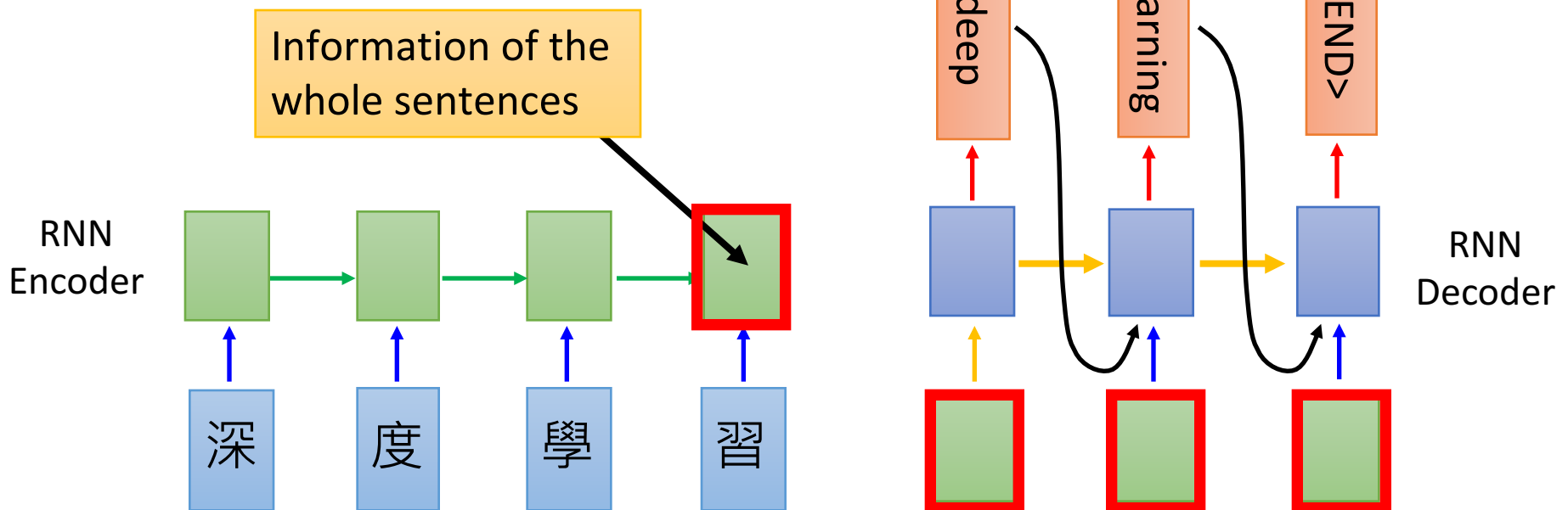
Attention on Sensory Info



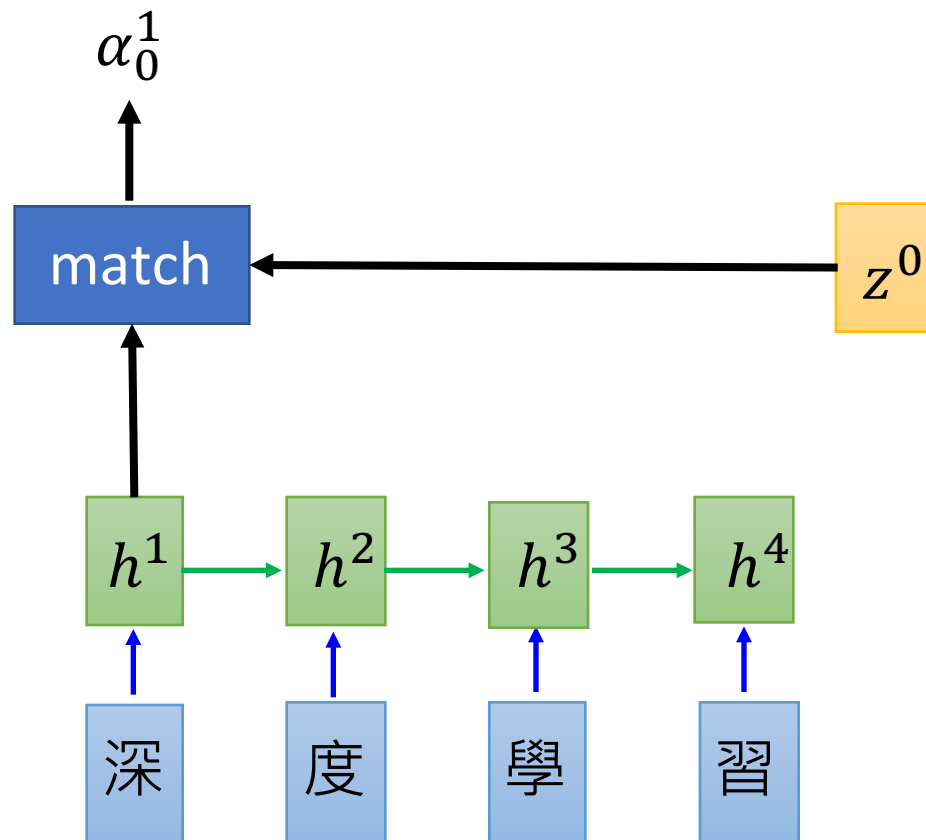
Machine Translation

Sequence-to-sequence learning: both input and output are both sequences *with different lengths*.

E.g. 深度學習 → deep learning



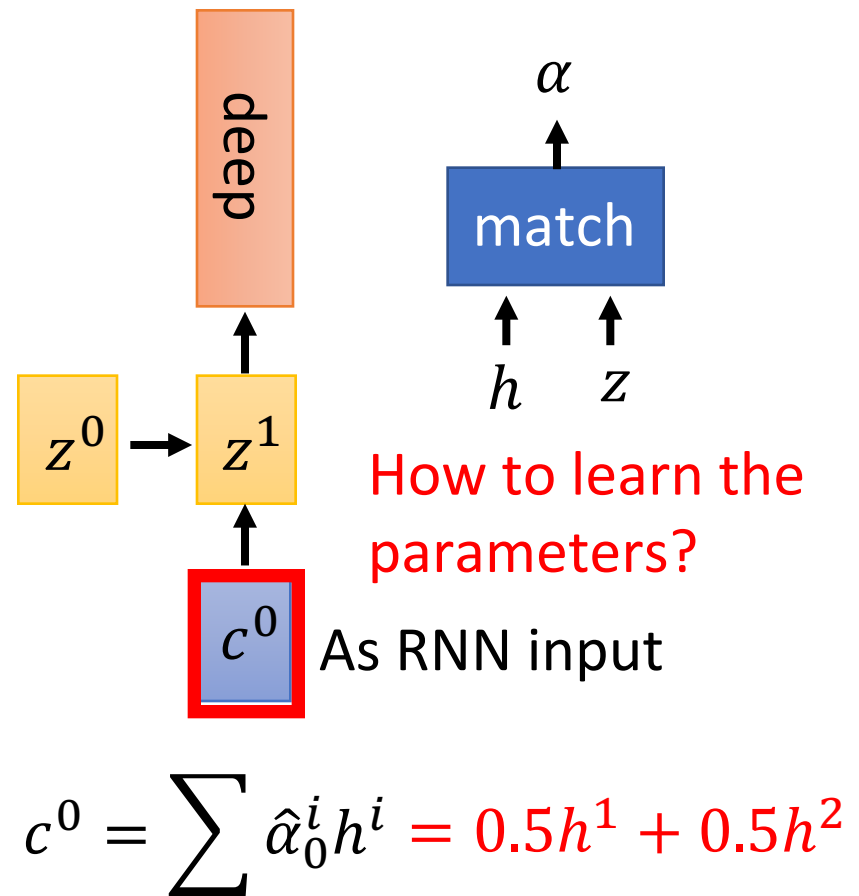
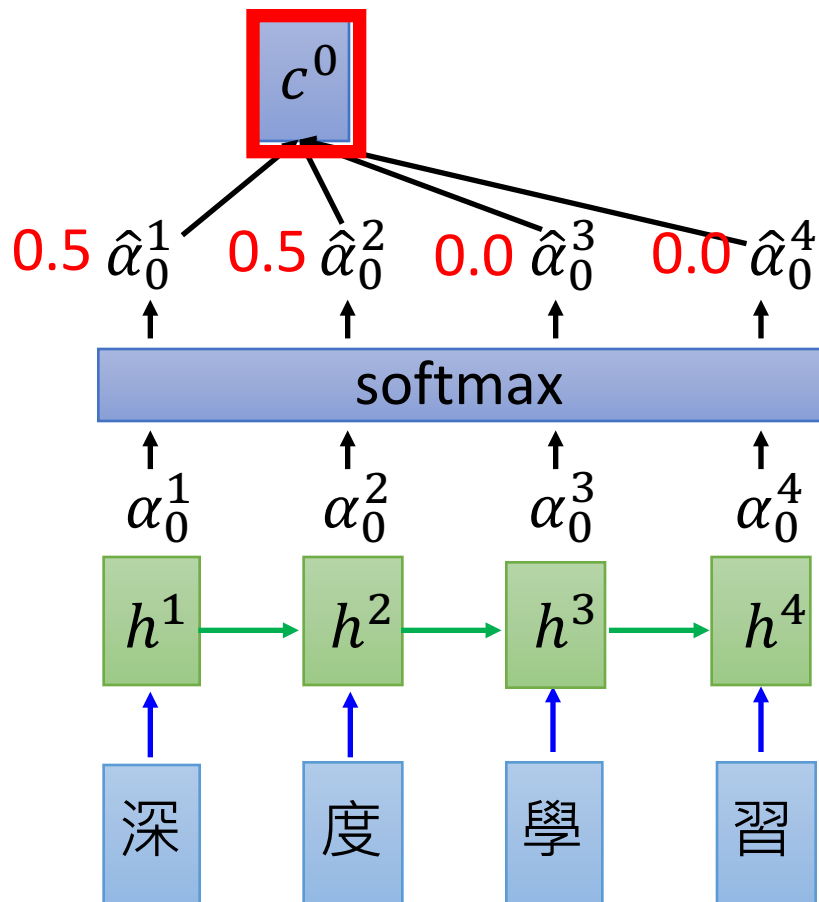
Machine Translation with Attention



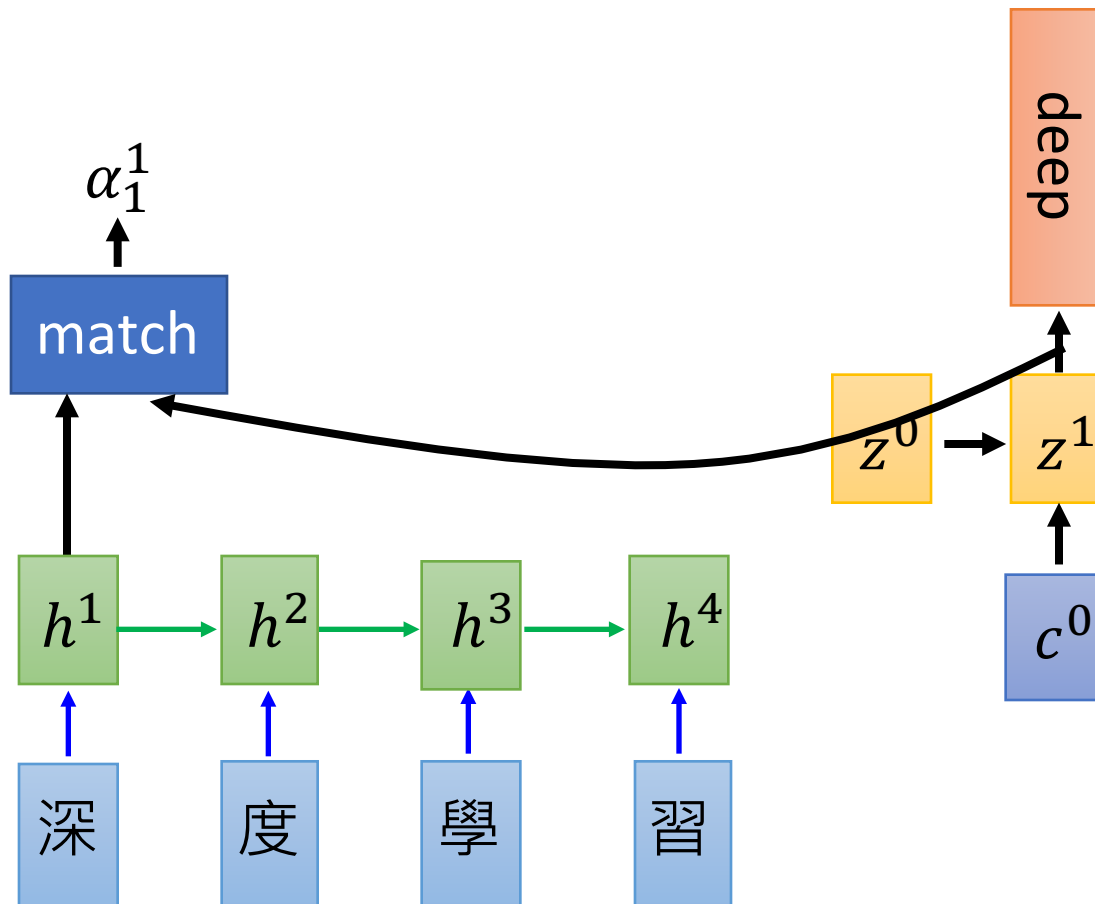
What is **match** ?

How to learn the parameters?

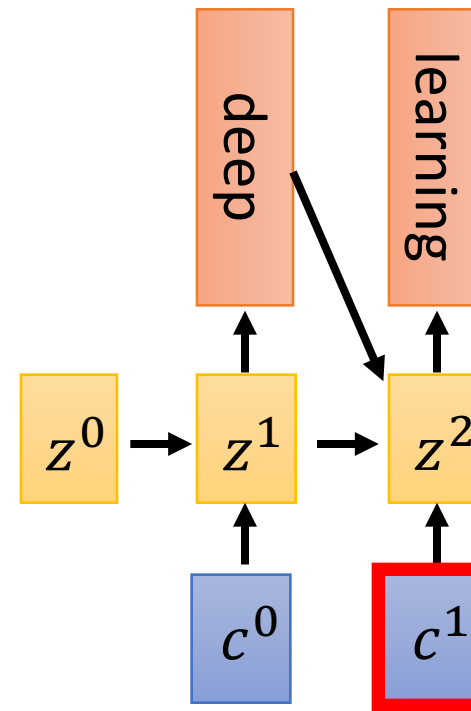
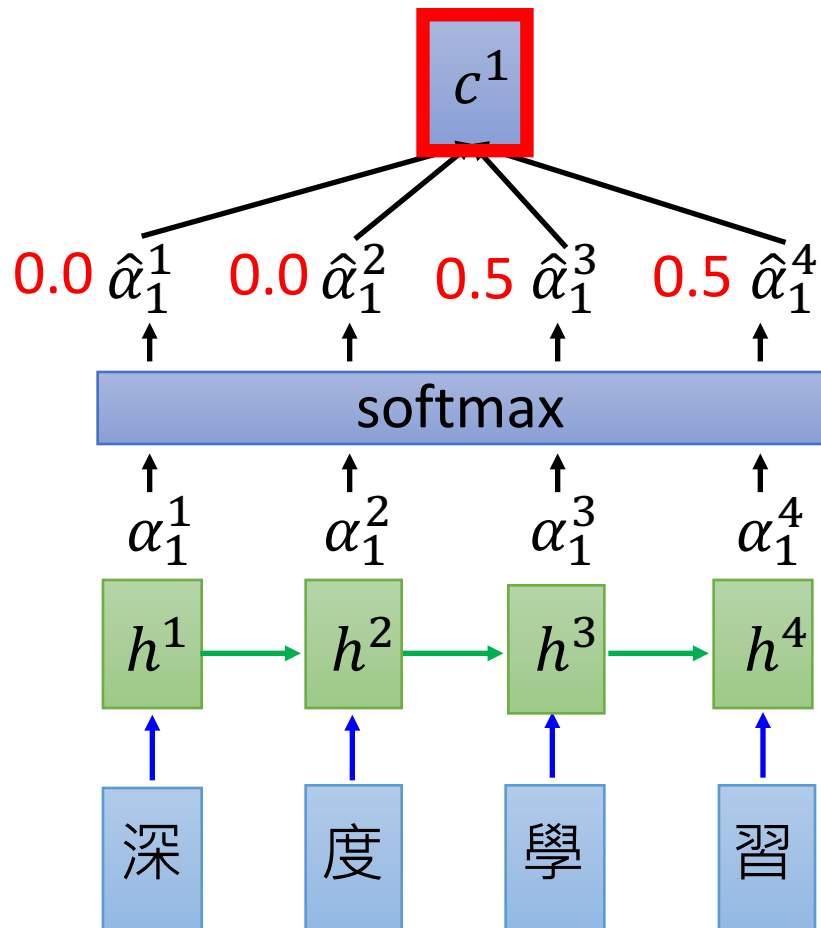
Machine Translation with Attention



Machine Translation with Attention

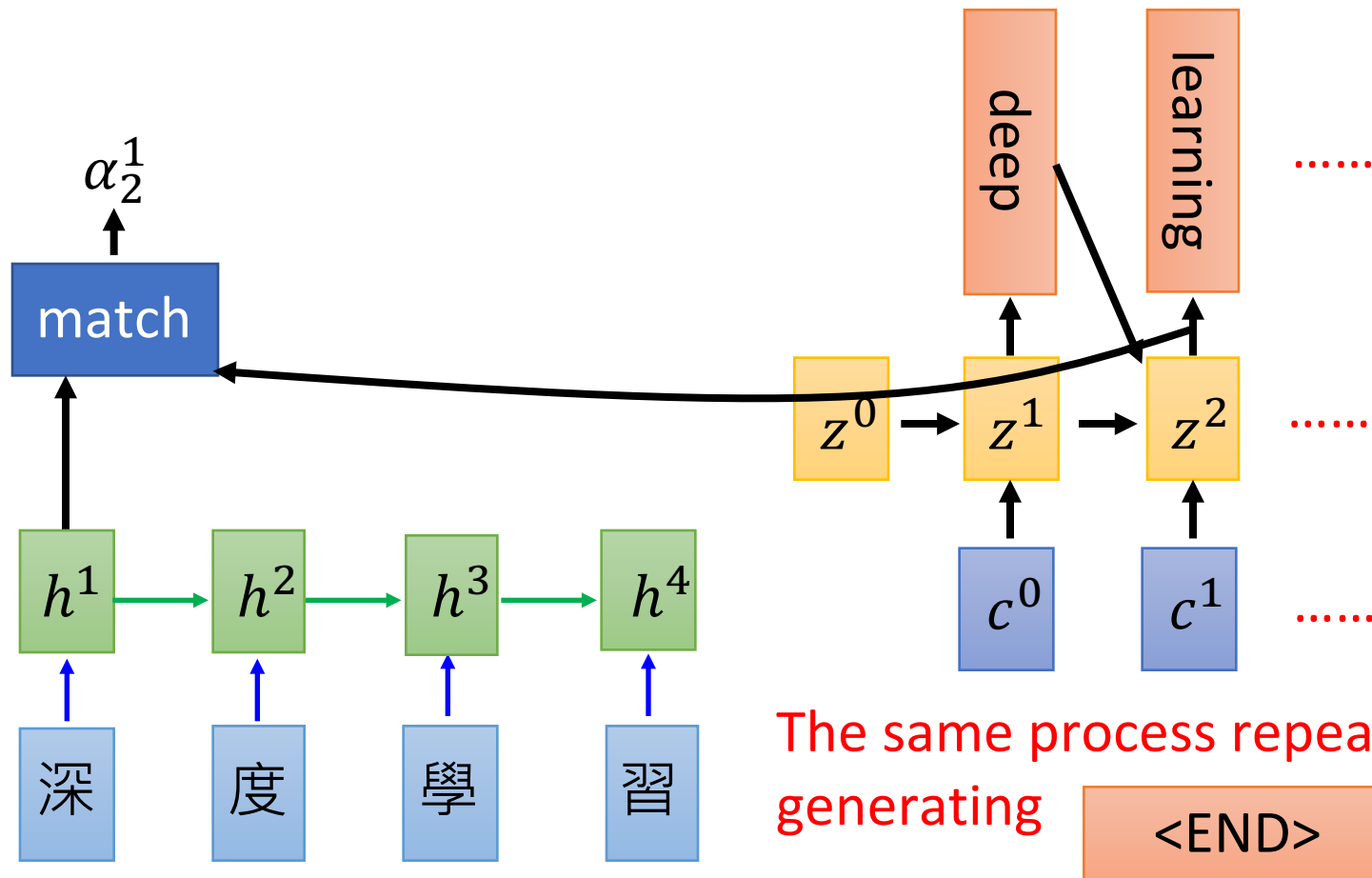


Machine Translation with Attention

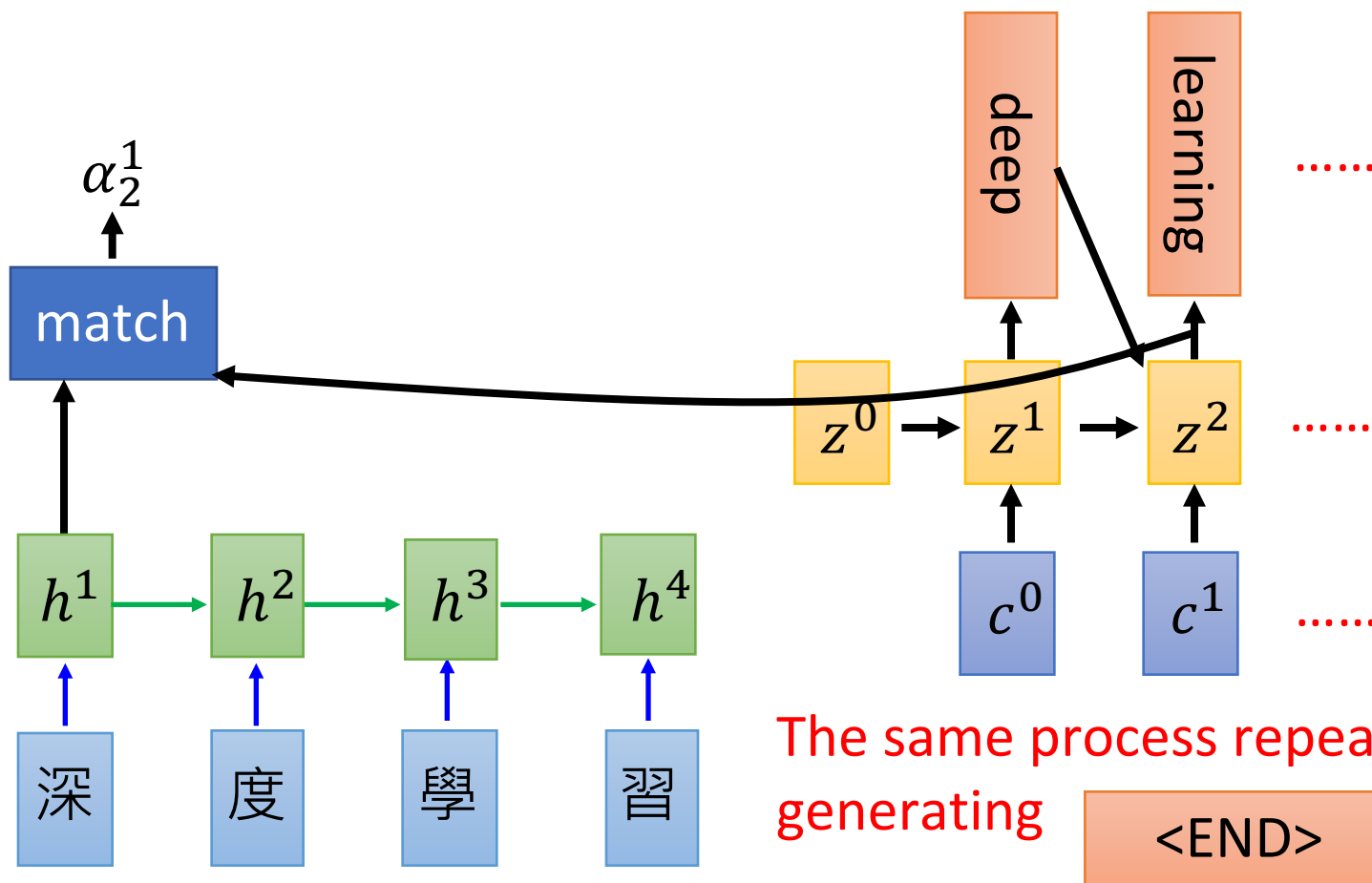


$$c^1 = \sum \hat{\alpha}_1^i h^i = 0.5h^3 + 0.5h^4$$

Machine Translation with Attention

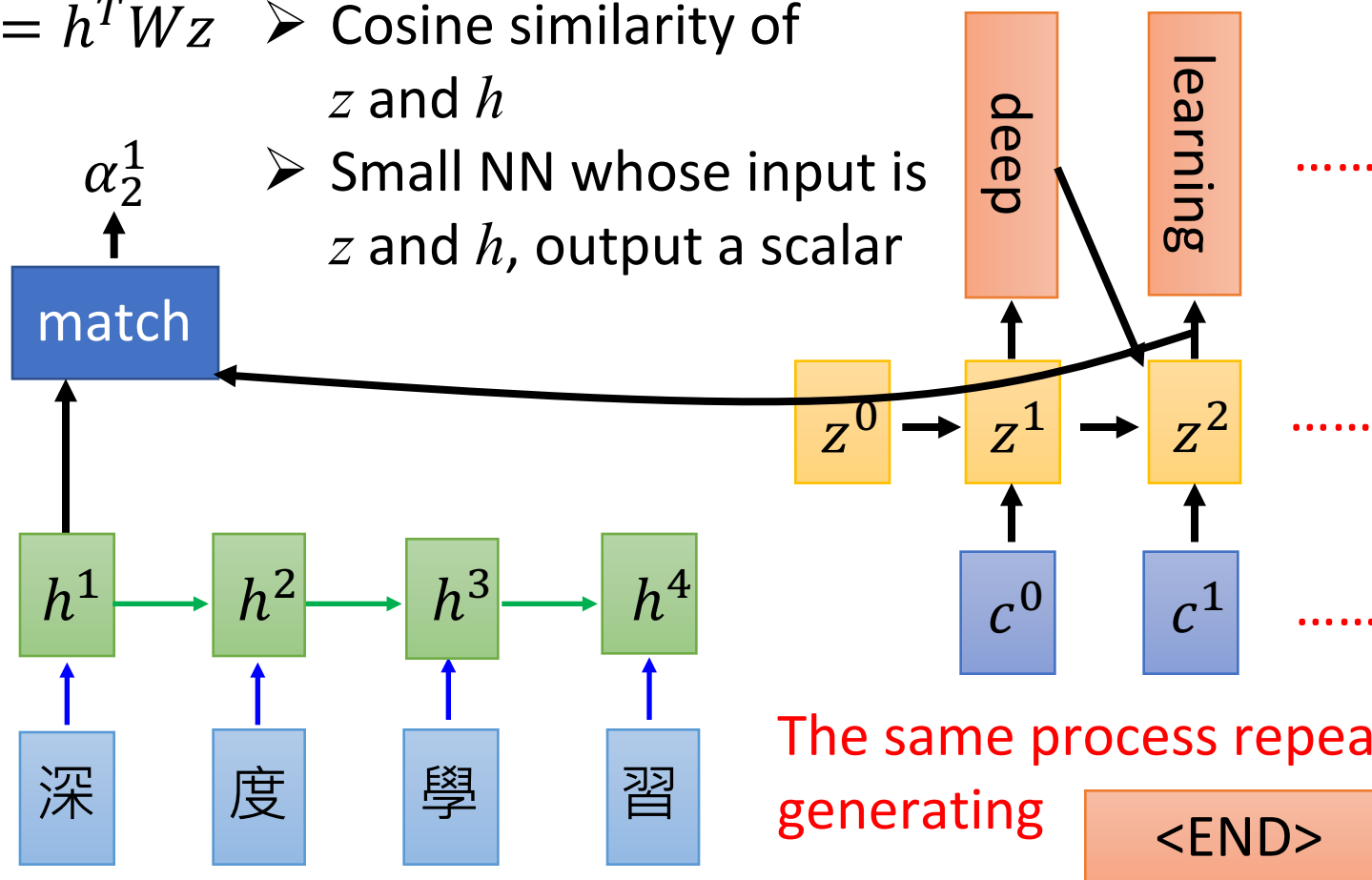


Group Discussion: how would you derive the attention weight α_2^1 ?



Group Discussion: how would you derive the attention weight α_2^1 ?

- $\alpha = h^T W z$
- Cosine similarity of z and h
- Small NN whose input is z and h , output a scalar



Speech Recognition with Attention

Alignment between the Characters and Audio

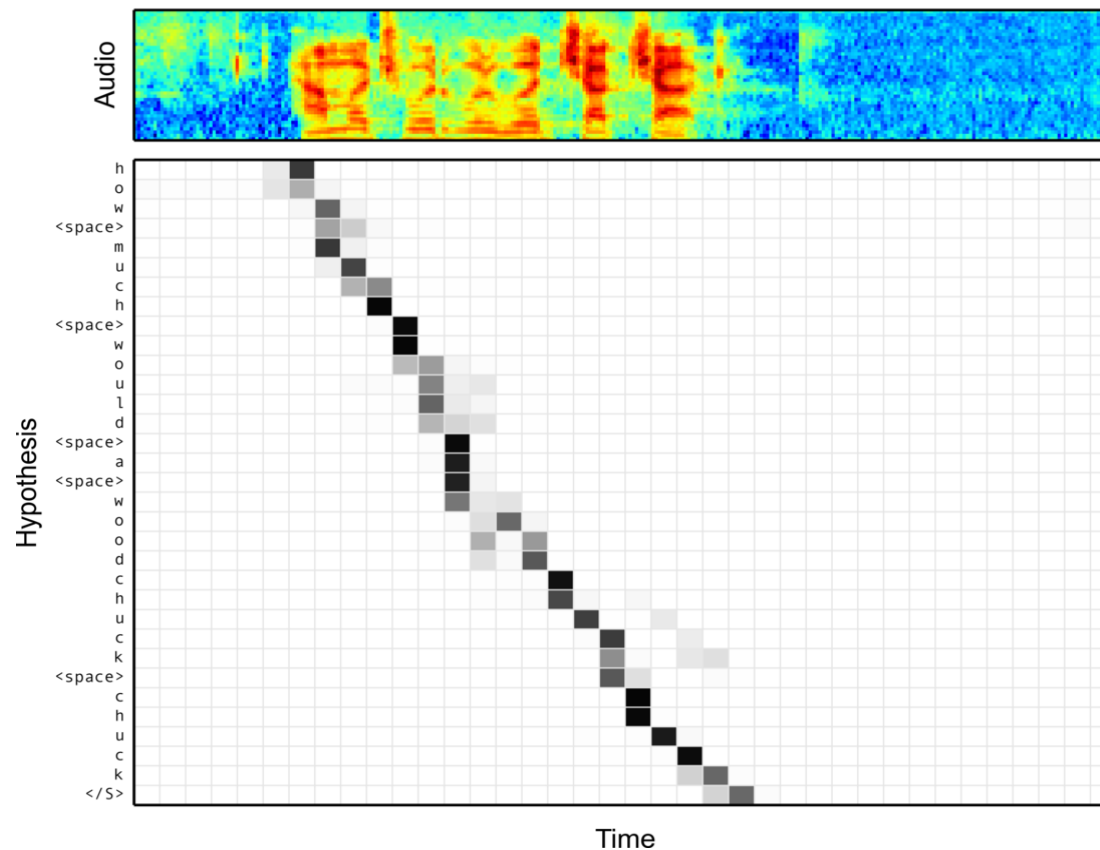


Image Captioning

Input: image

Output: word sequence

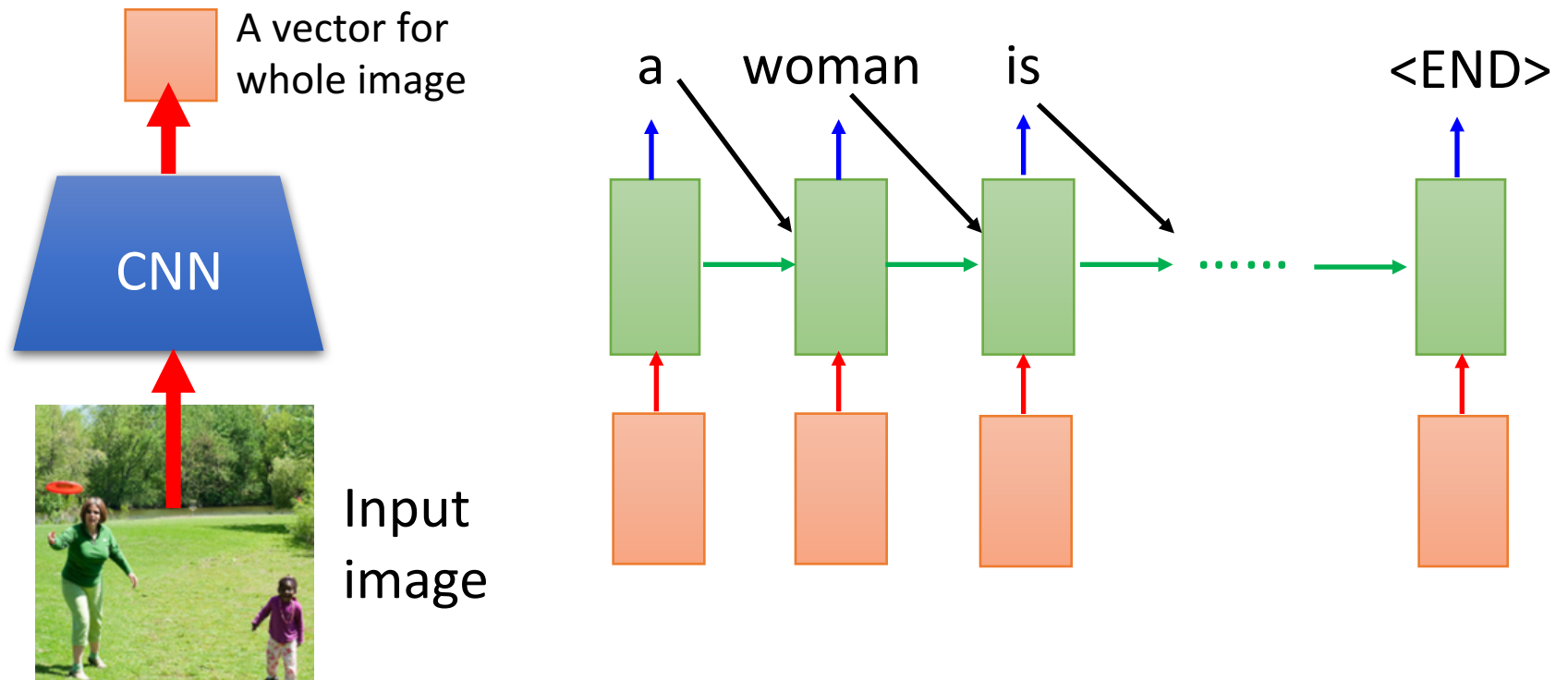


Image Captioning with Attention

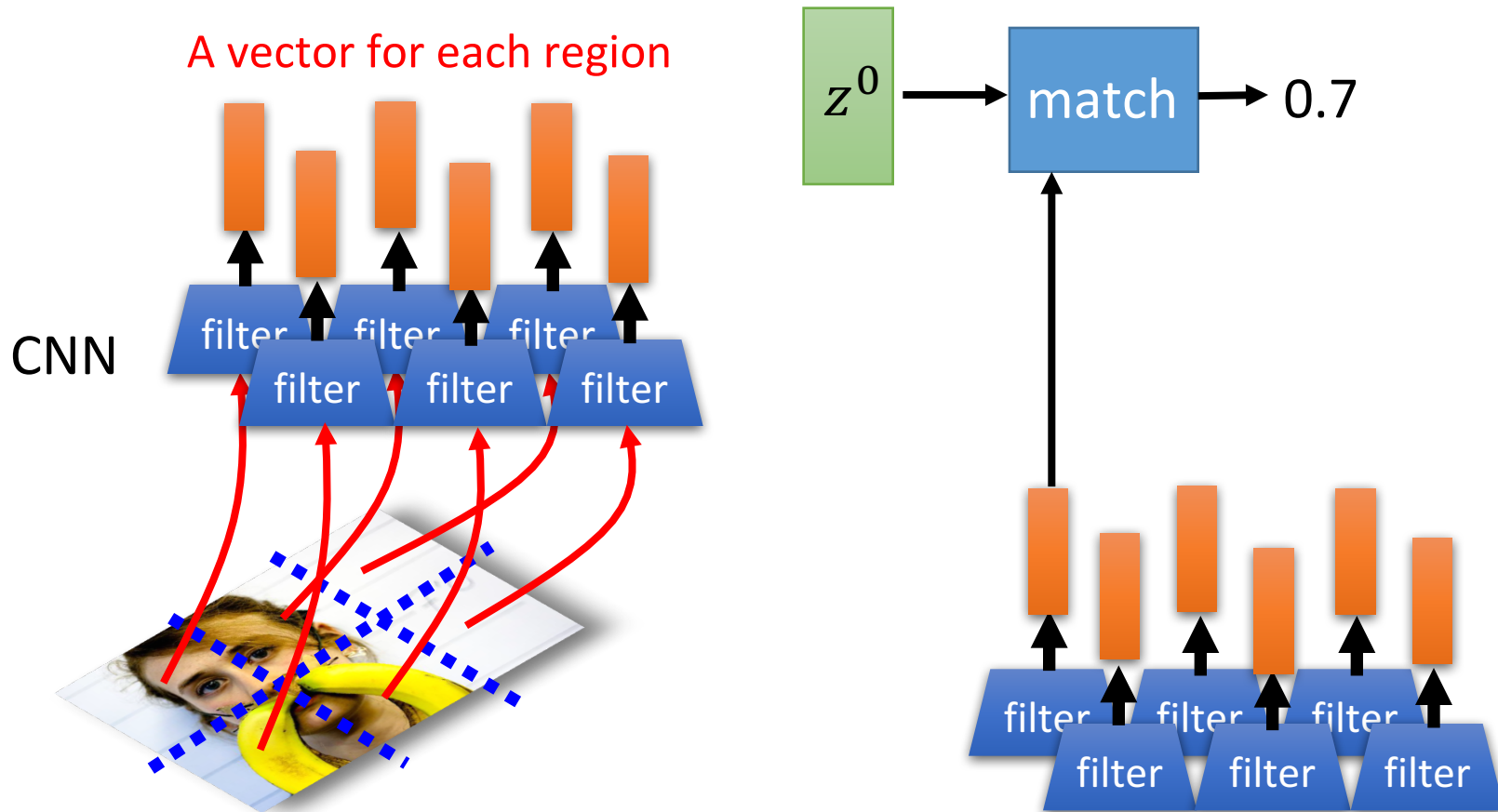


Image Captioning with Attention

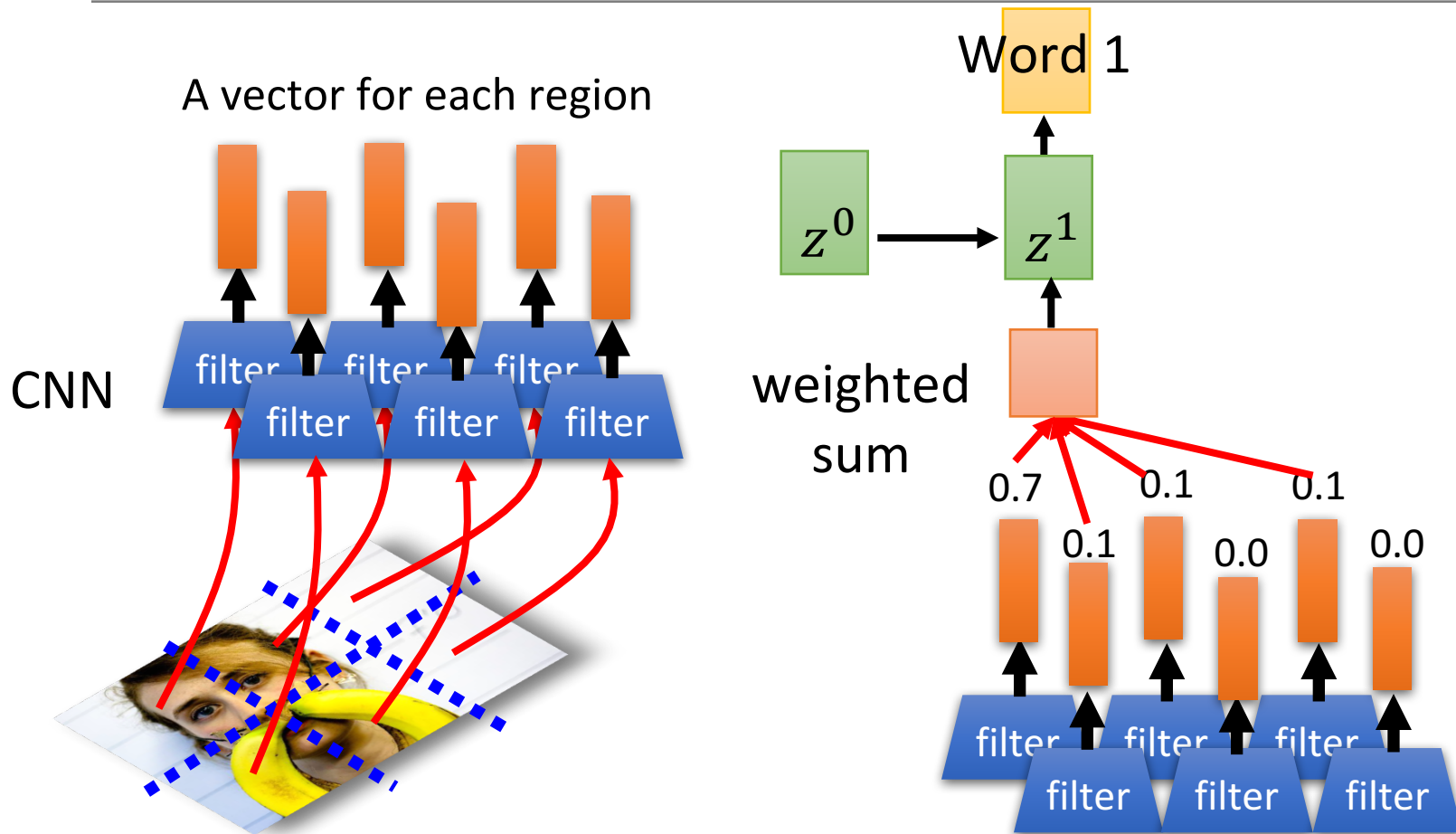


Image Captioning with Attention

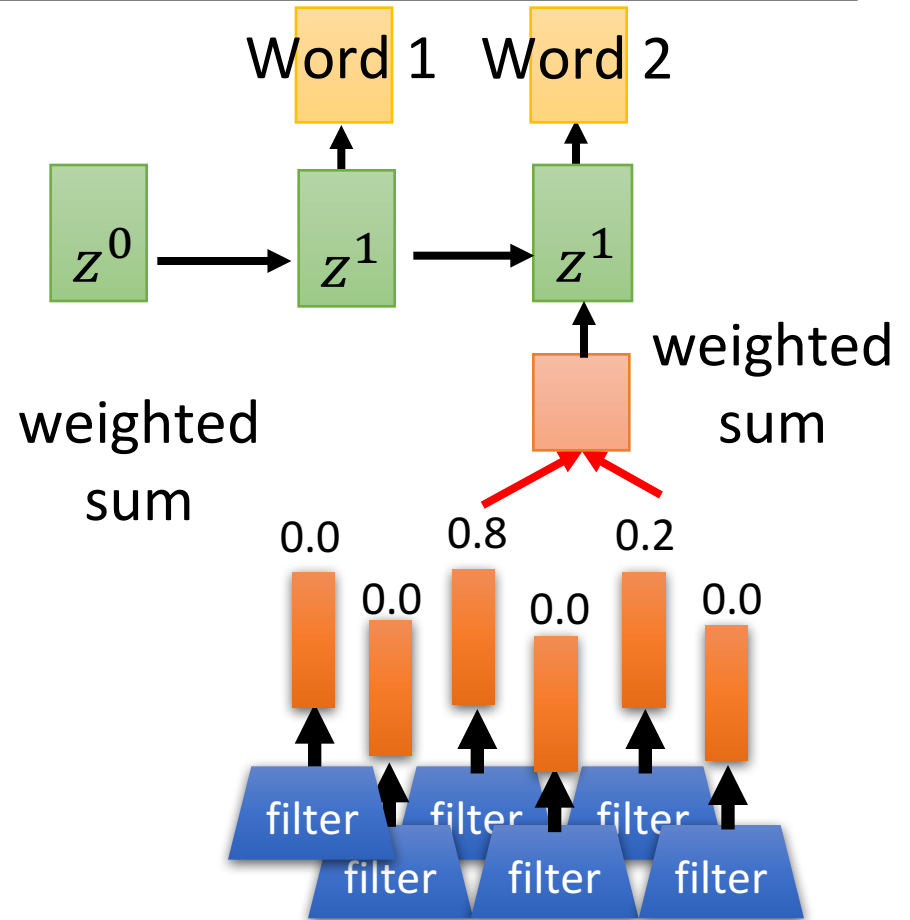
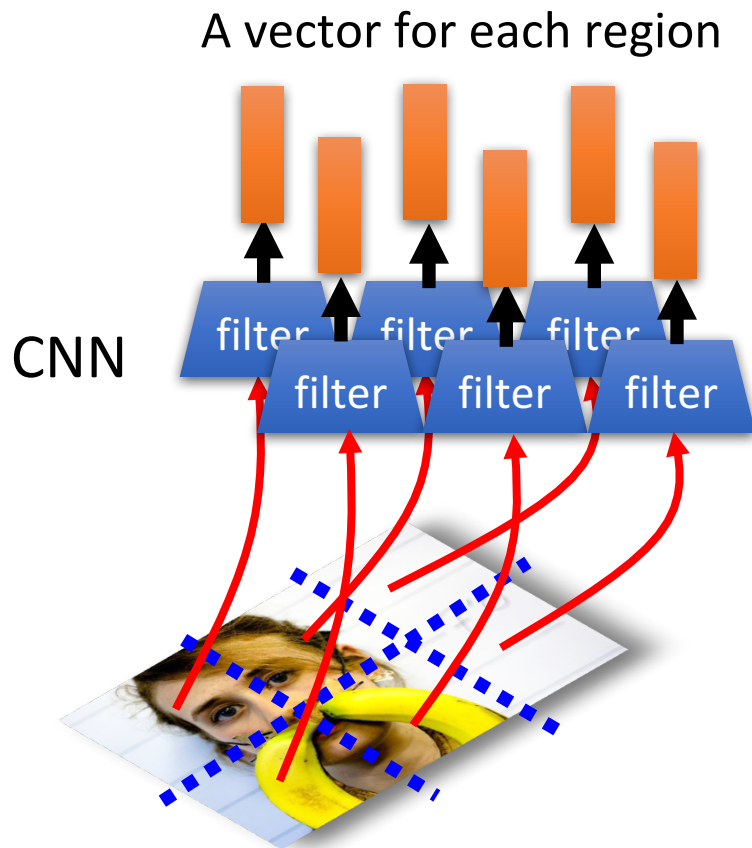
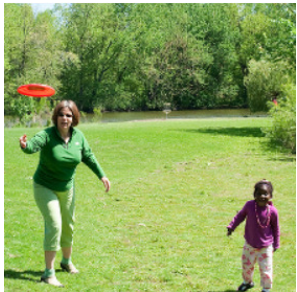


Image Captioning

Good examples



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



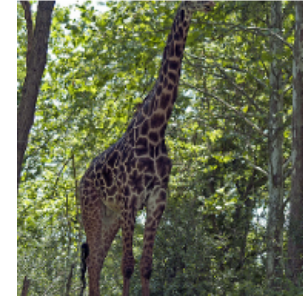
A stop sign is on a road with a mountain in the background.



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



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

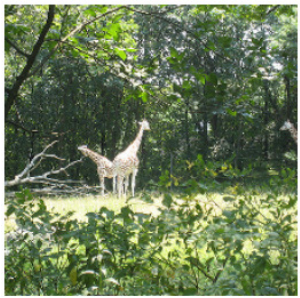


A giraffe standing in a forest with trees in the background.



Image Captioning

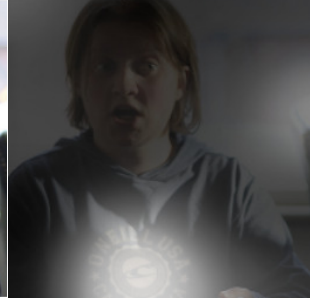
Bad examples



A large white bird standing in a forest.



A woman holding a clock in her hand.



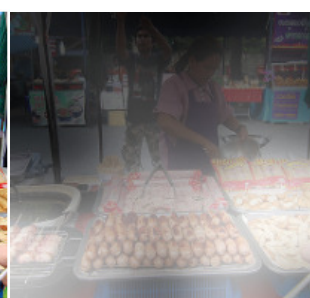
A man wearing a hat and a hat on a skateboard.



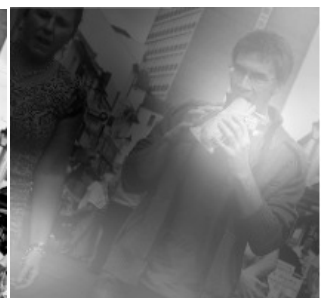
A person is standing on a beach with a surfboard.



A woman is sitting at a table with a large pizza.



A man is talking on his cell phone while another man watches.

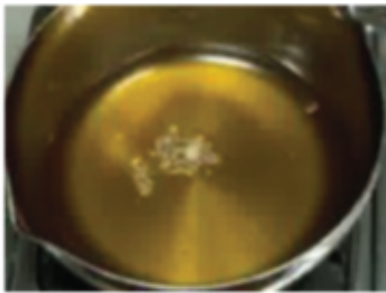


Video Captioning



Ref: A man and a woman ride a motorcycle
A **man** and a **woman** are **talking** on the **road**

Video Captioning



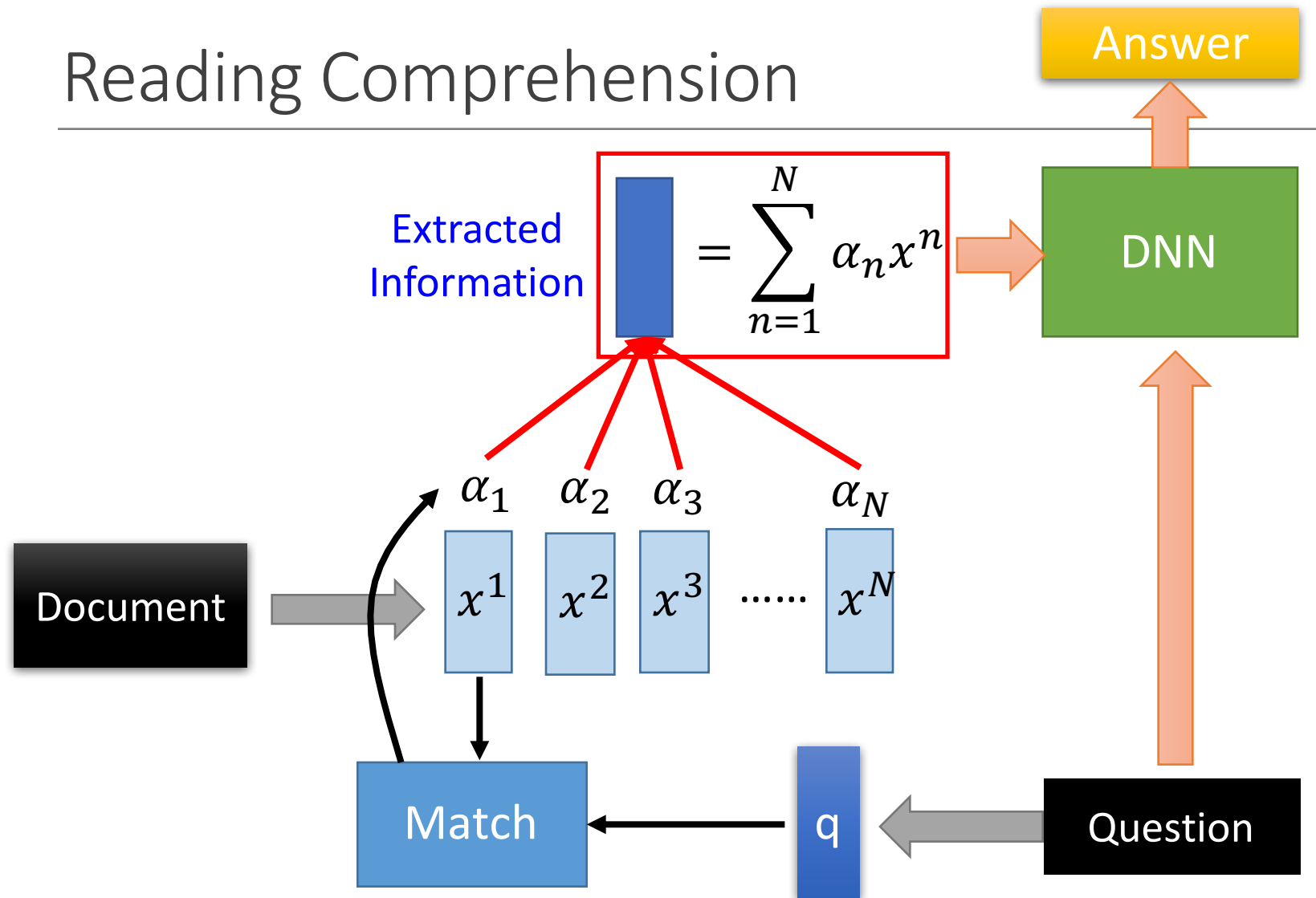
Ref: A woman is frying food
Someone is **frying** a **fish** in a **pot**

Group Discussion:

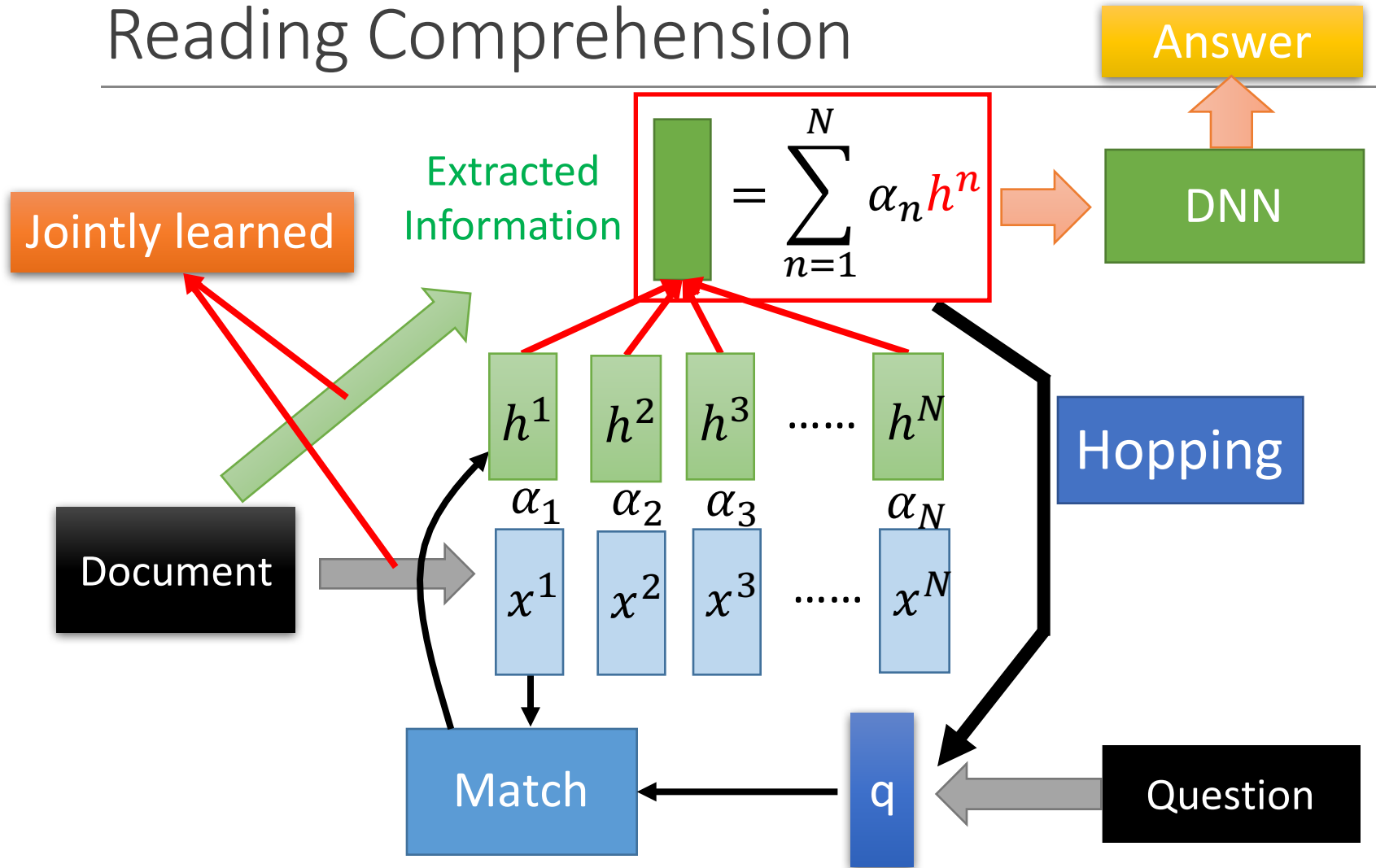
Learning attention weights vs using cosine similarity, which one is better? Why?

- $\alpha = h^T W z$
- Cosine similarity of z and h

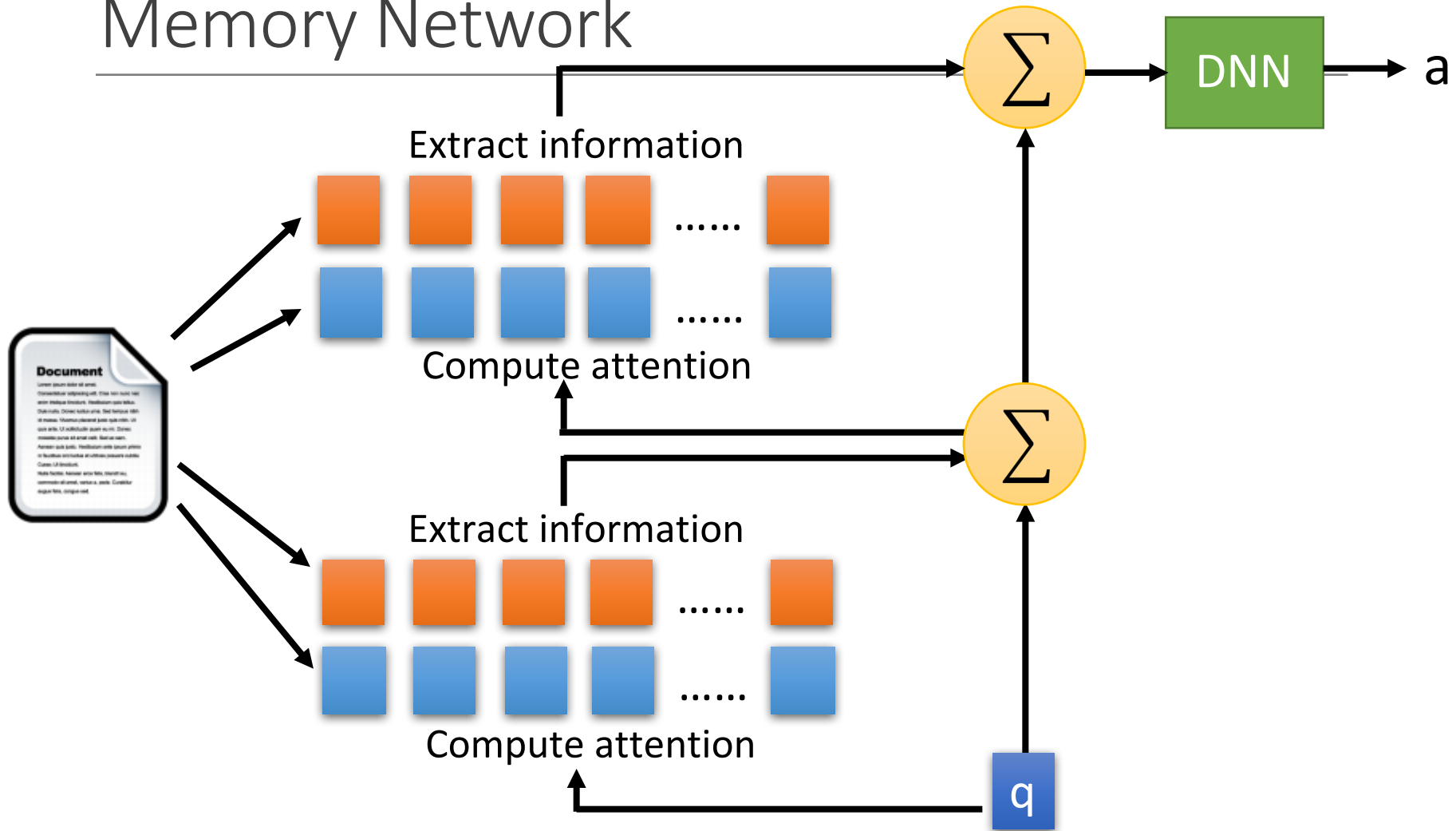
Reading Comprehension



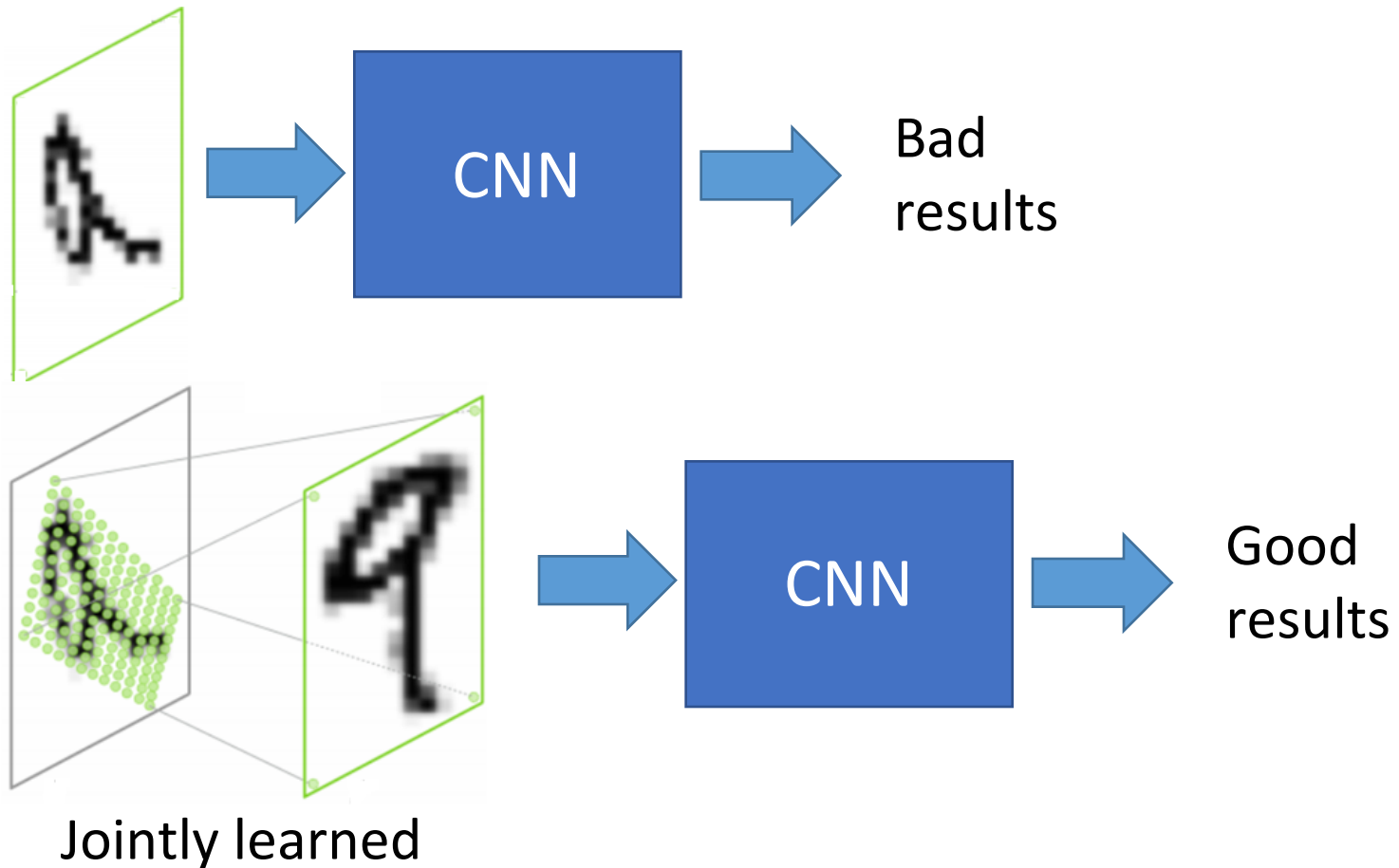
Reading Comprehension



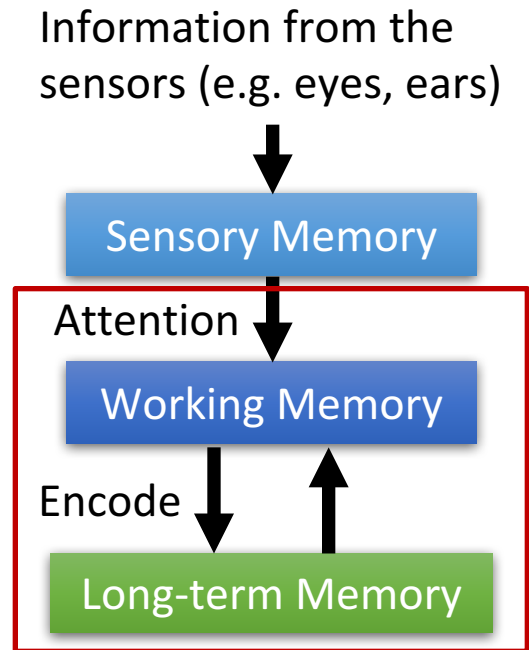
Memory Network



Special Attention: Spatial Transformers



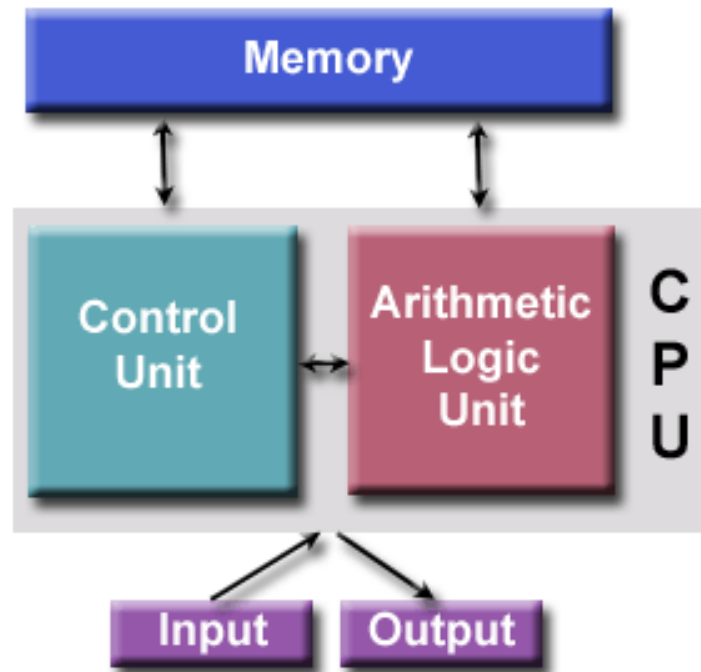
Attention on Memory



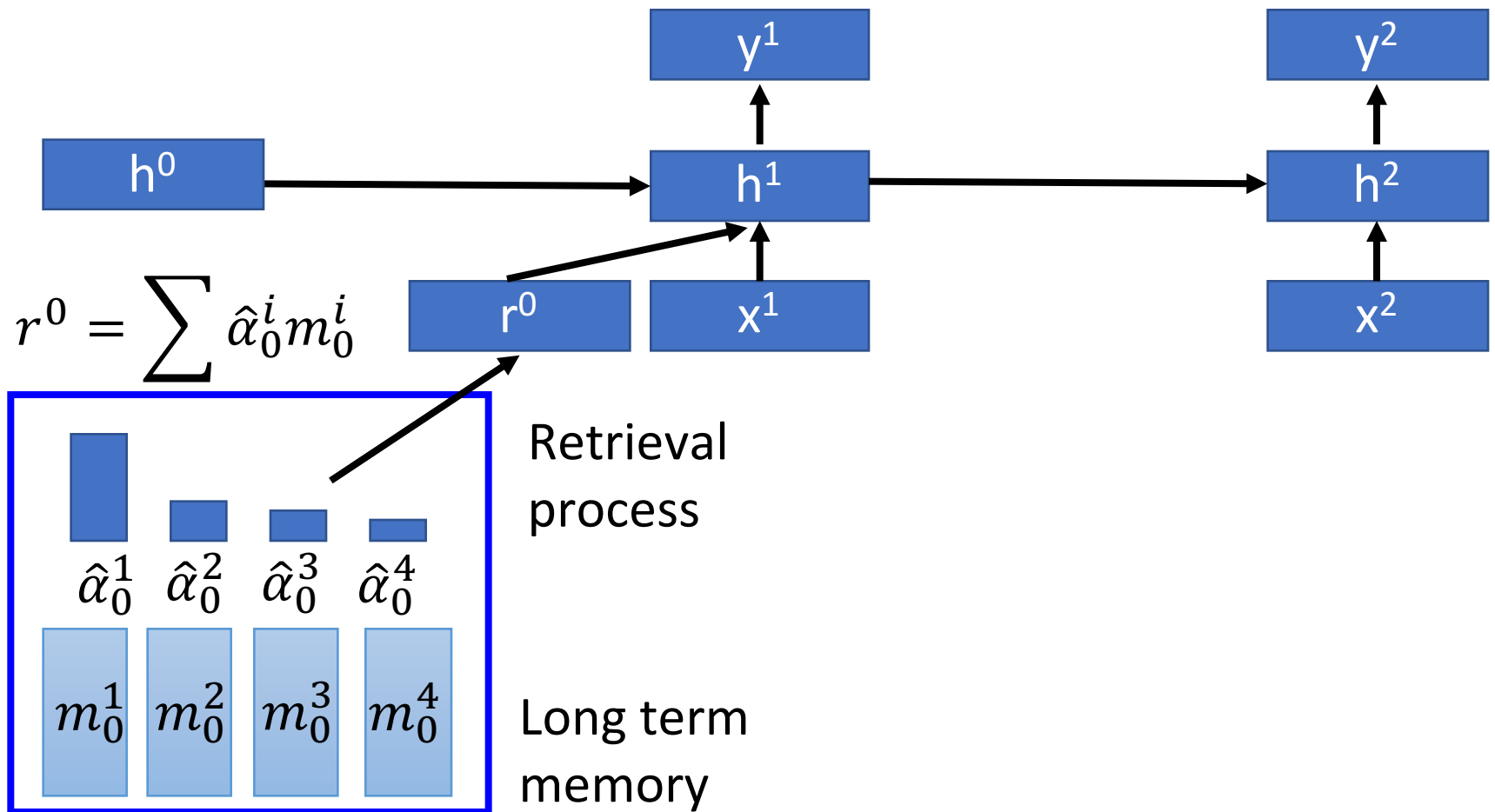
Neural Turing Machine

Von Neumann architecture

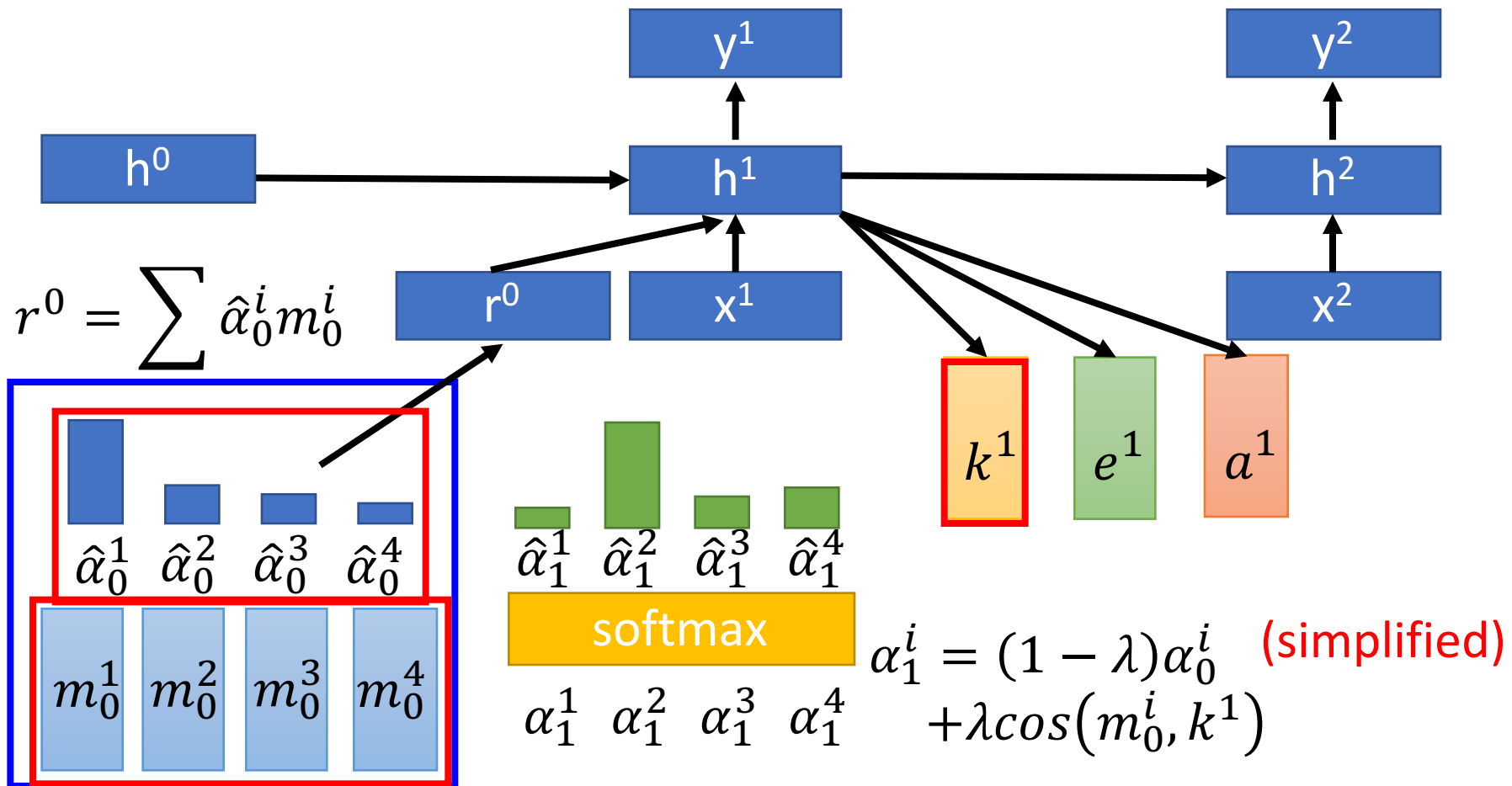
Neural Turing Machine is an advanced RNN/LSTM.



Neural Turing Machine



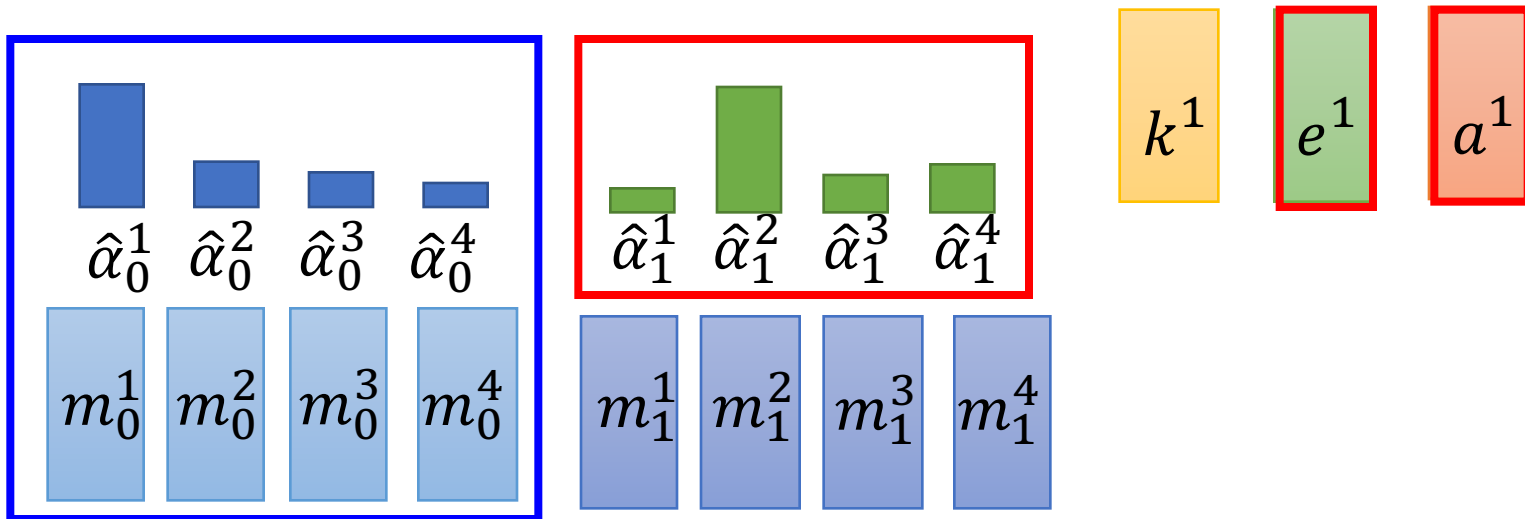
Neural Turing Machine



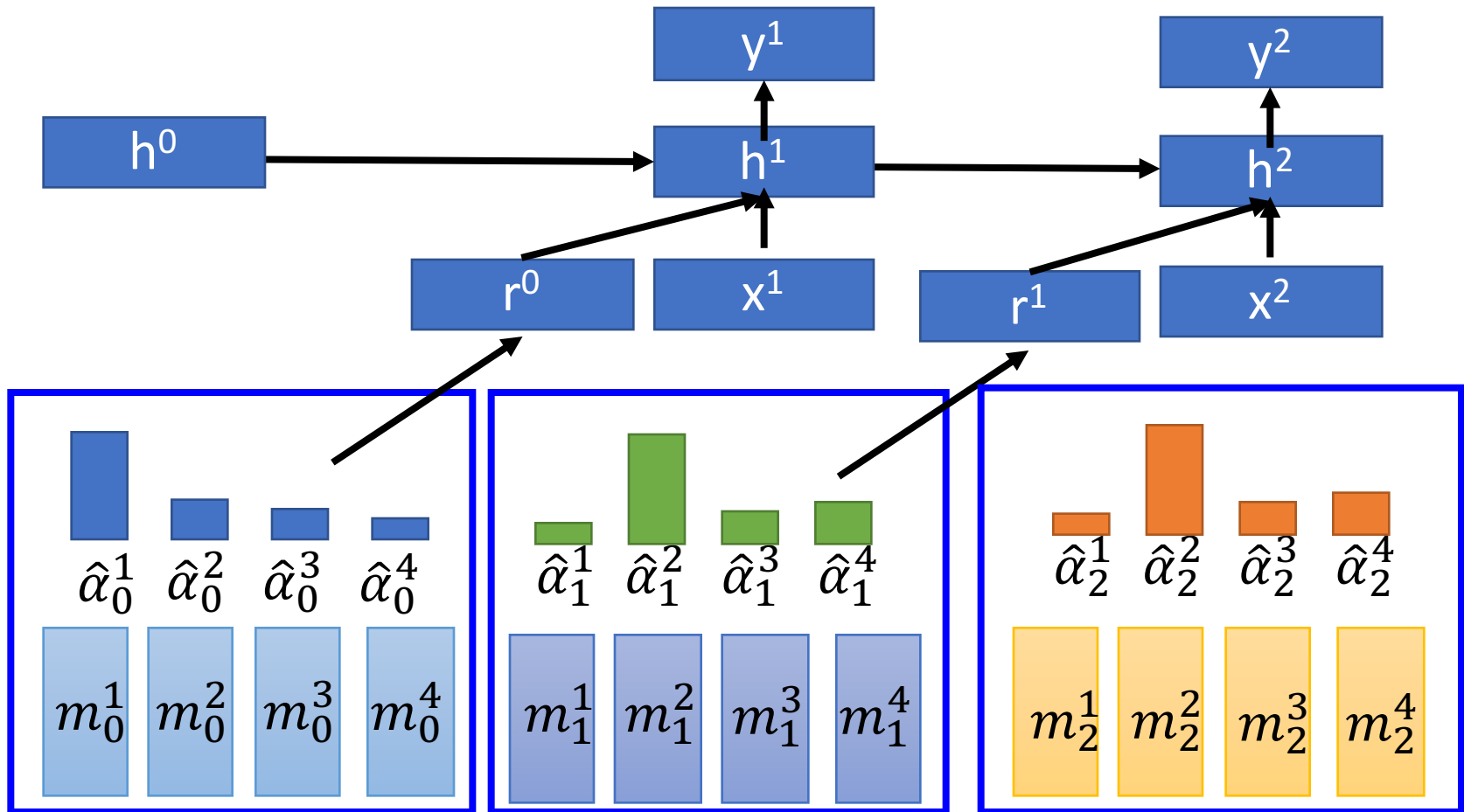
Neural Turing Machine

$$m_1^i = m_0^i * \begin{bmatrix} 1 & -\hat{\alpha}_1^i e^1 \end{bmatrix} + \hat{\alpha}_1^i a^1 \quad \rightarrow \text{Encode process}$$

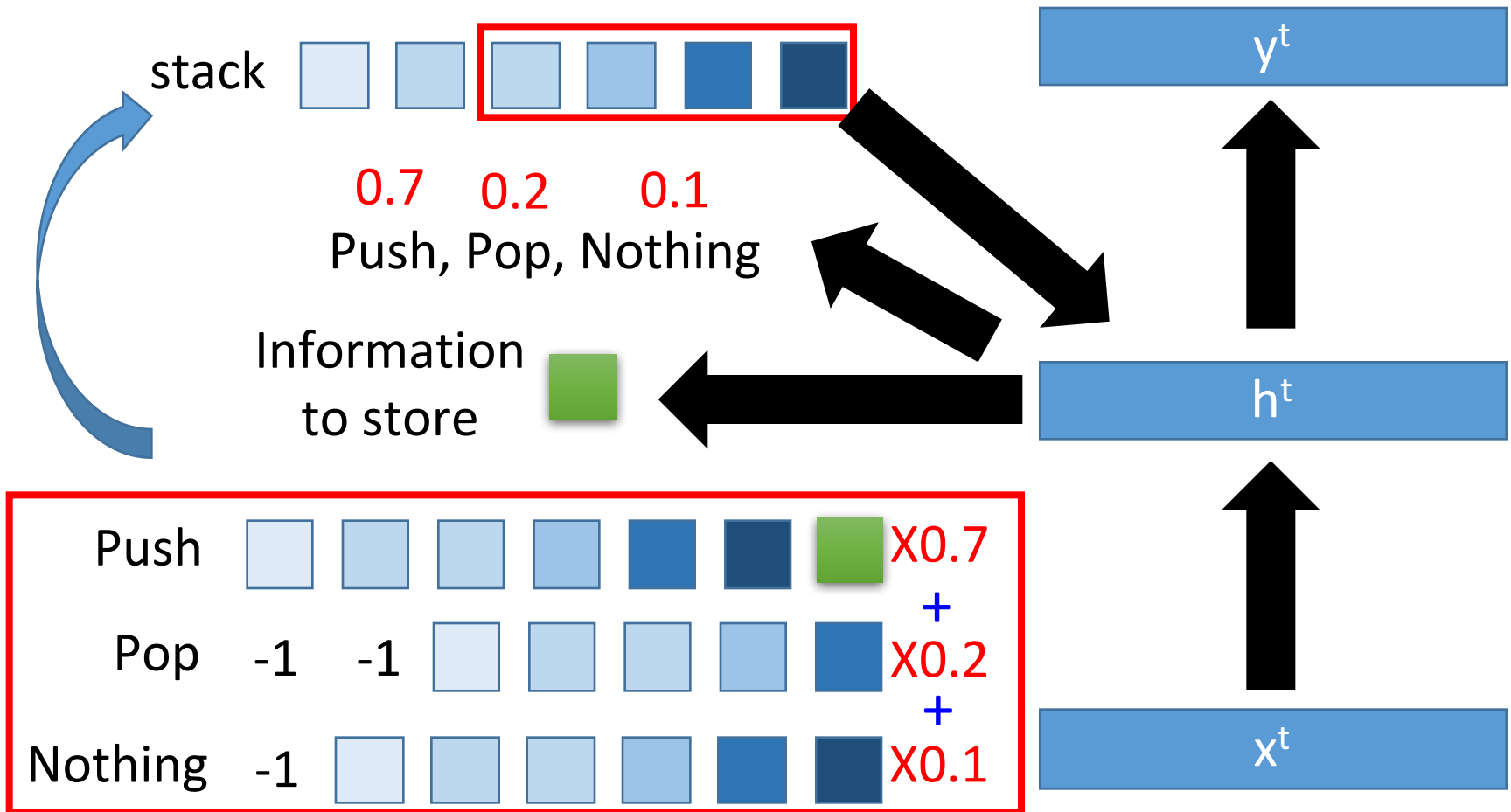
(element-wise)



Neural Turing Machine



Stack RNN

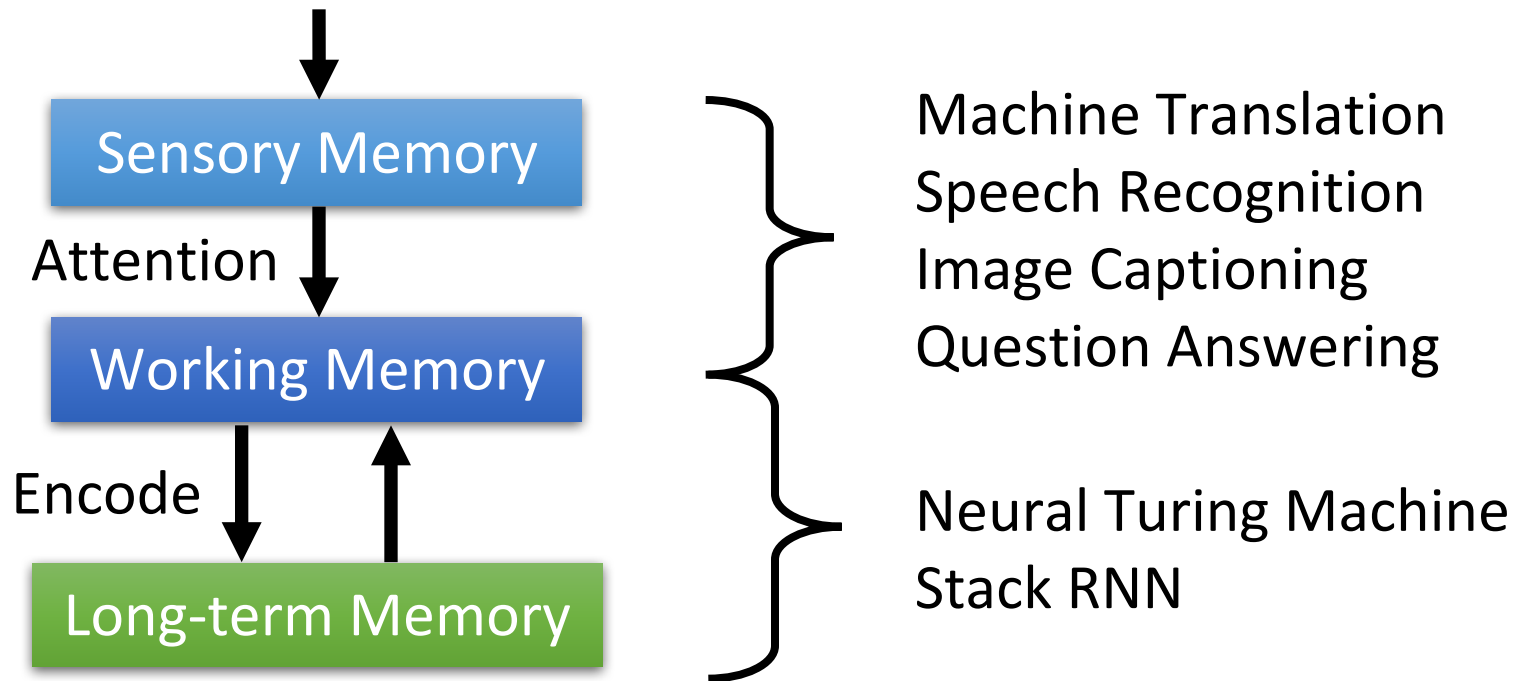


Quick Question: why stack but not queue?



Concluding Remarks

Information from the
sensors (e.g. eyes, ears)



Reference

End-To-End Memory Networks. S. Sukhbaatar, A. Szlam, J. Weston, R. Fergus. arXiv Pre-Print, 2015.

Neural Turing Machines. Alex Graves, Greg Wayne, Ivo Danihelka. arXiv Pre-Print, 2014

Ask Me Anything: Dynamic Memory Networks for Natural Language Processing. Kumar et al. arXiv Pre-Print, 2015

Neural Machine Translation by Jointly Learning to Align and Translate. D. Bahdanau, K. Cho, Y. Bengio; International Conference on Representation Learning 2015.

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. Kelvin Xu et. al.. arXiv Pre-Print, 2015.

Attention-Based Models for Speech Recognition. Jan Chorowski, Dzmitry Bahdanau, Dmitriy Serdyuk, Kyunghyun Cho, Yoshua Bengio. arXiv Pre-Print, 2015.

A Neural Attention Model for Abstractive Sentence Summarization. A. M. Rush, S. Chopra and J. Weston. EMNLP 2015.