

Conditional Generation by RNN & Attention

Hung-yi Lee

李宏毅

Outline

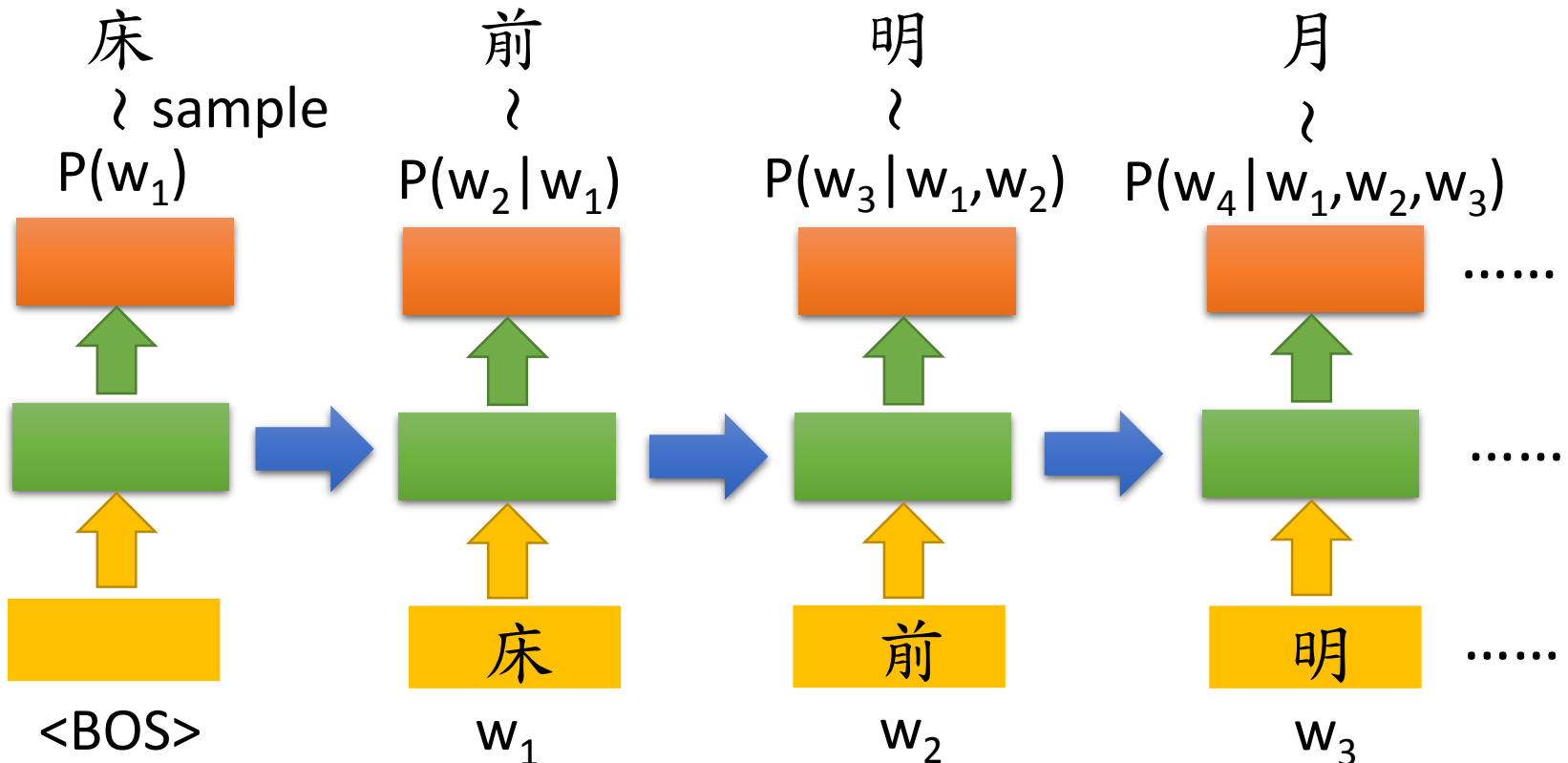
- Generation
- Attention
- Tips for Generation
- Pointer Network

Generation

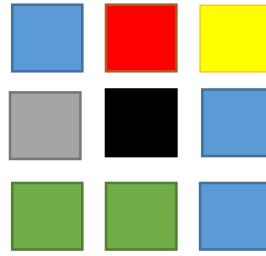
Generating a structured object component-by-component

Generation

- Sentences are composed of characters/words
 - Generating a character/word at each time by RNN



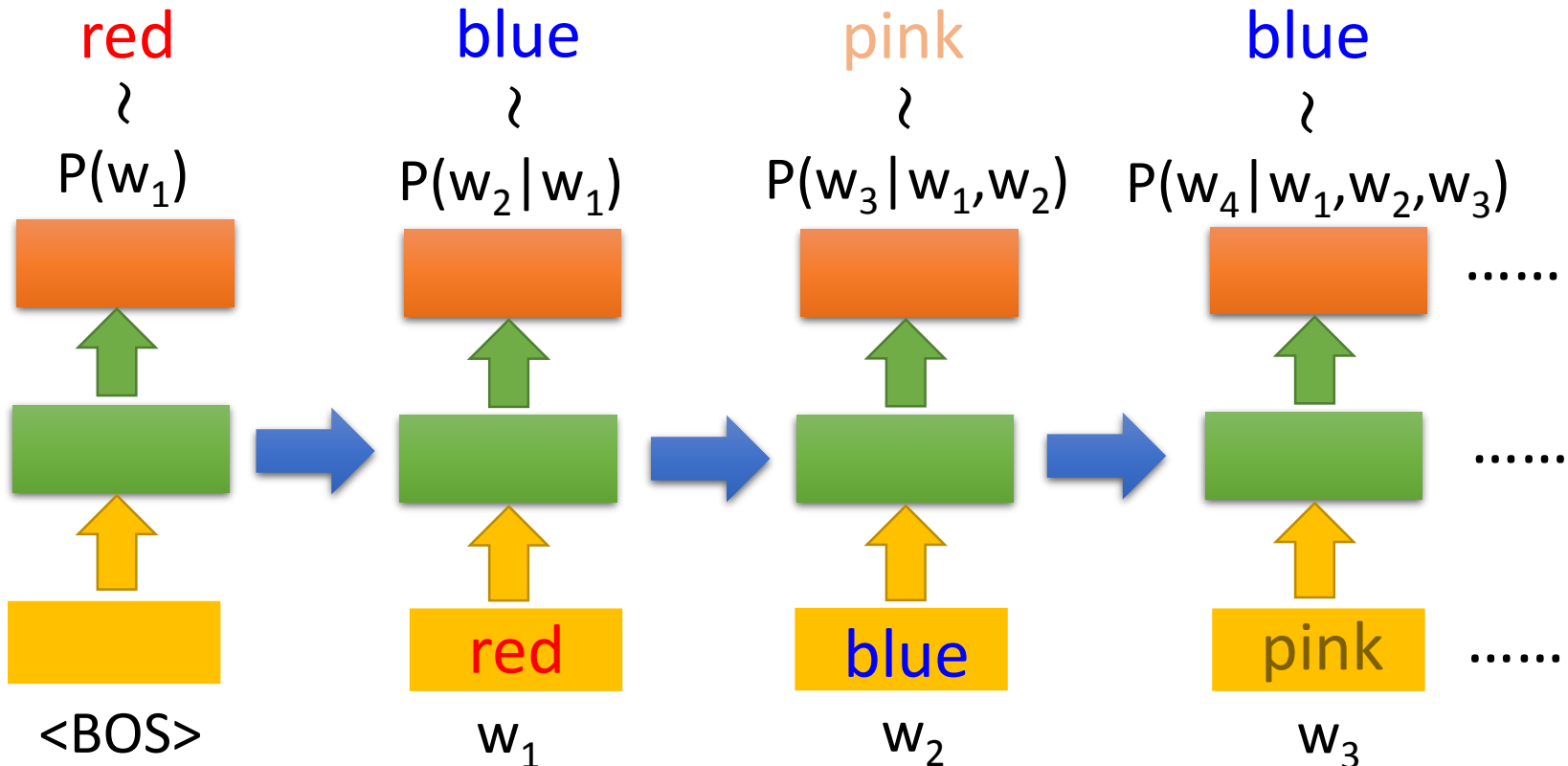
Generation



Consider as a sentence
blue red yellow gray

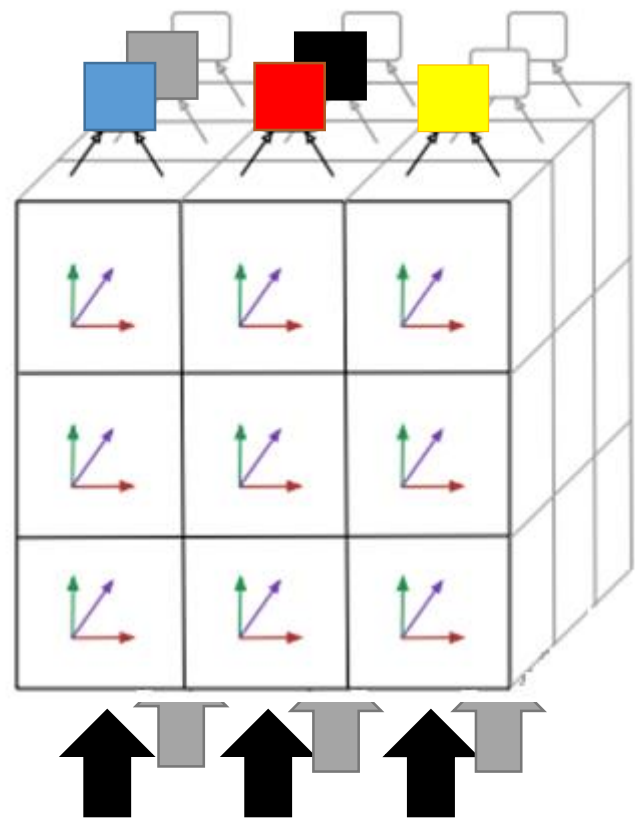
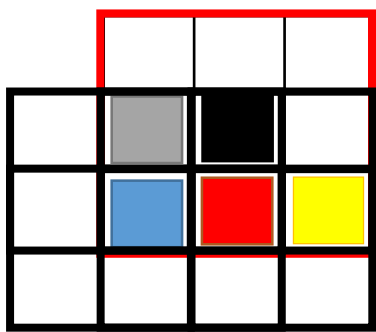
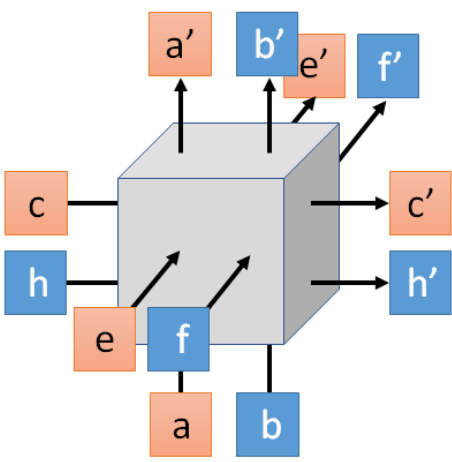
Train a language model
based on the “sentences”

- Images are composed of pixels
 - Generating a pixel at each time by RNN

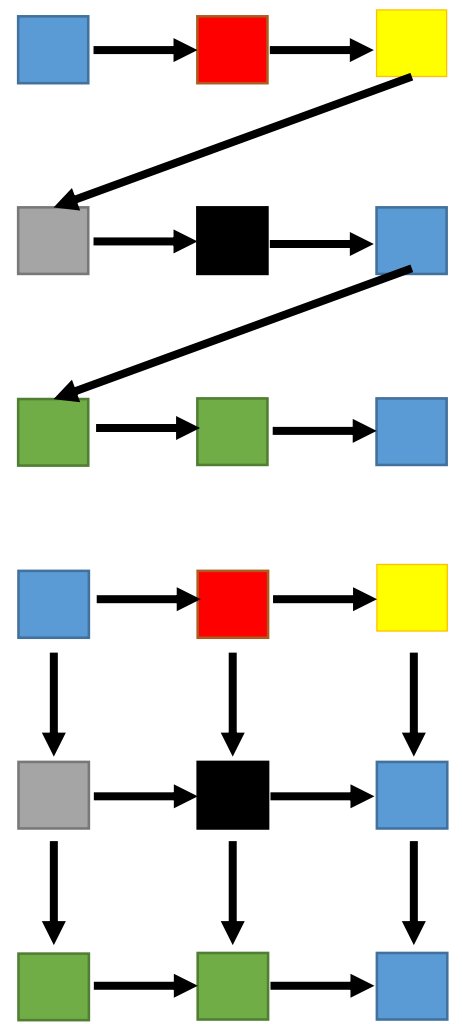


Generation

- Images are composed of pixels



3 x 3 images



Generation

- Image

- Aaron van den Oord, Nal Kalchbrenner, Koray Kavukcuoglu, Pixel Recurrent Neural Networks, arXiv preprint, 2016
- Aaron van den Oord, Nal Kalchbrenner, Oriol Vinyals, Lasse Espeholt, Alex Graves, Koray Kavukcuoglu, Conditional Image Generation with PixelCNN Decoders, arXiv preprint, 2016

- Video

- Aaron van den Oord, Nal Kalchbrenner, Koray Kavukcuoglu, Pixel Recurrent Neural Networks, arXiv preprint, 2016

- Handwriting

- Alex Graves, Generating Sequences With Recurrent Neural Networks, arXiv preprint, 2013

- Speech

- Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, Koray Kavukcuoglu, WaveNet: A Generative Model for Raw Audio, 2016

Conditional Generation

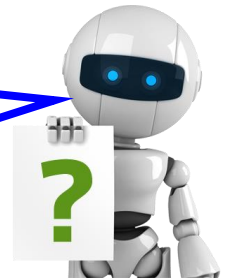
- We don't want to simply generate some random sentences.
- Generate sentences based on conditions:

Caption Generation

Given
condition:

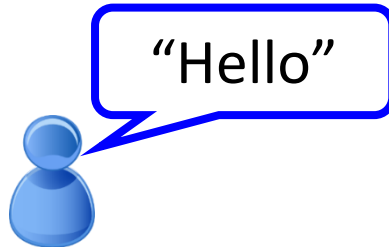


“A young girl
is dancing.”

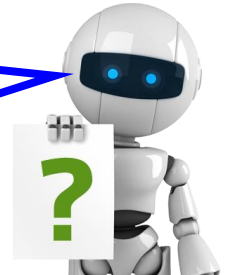


Chat-bot

Given
condition:



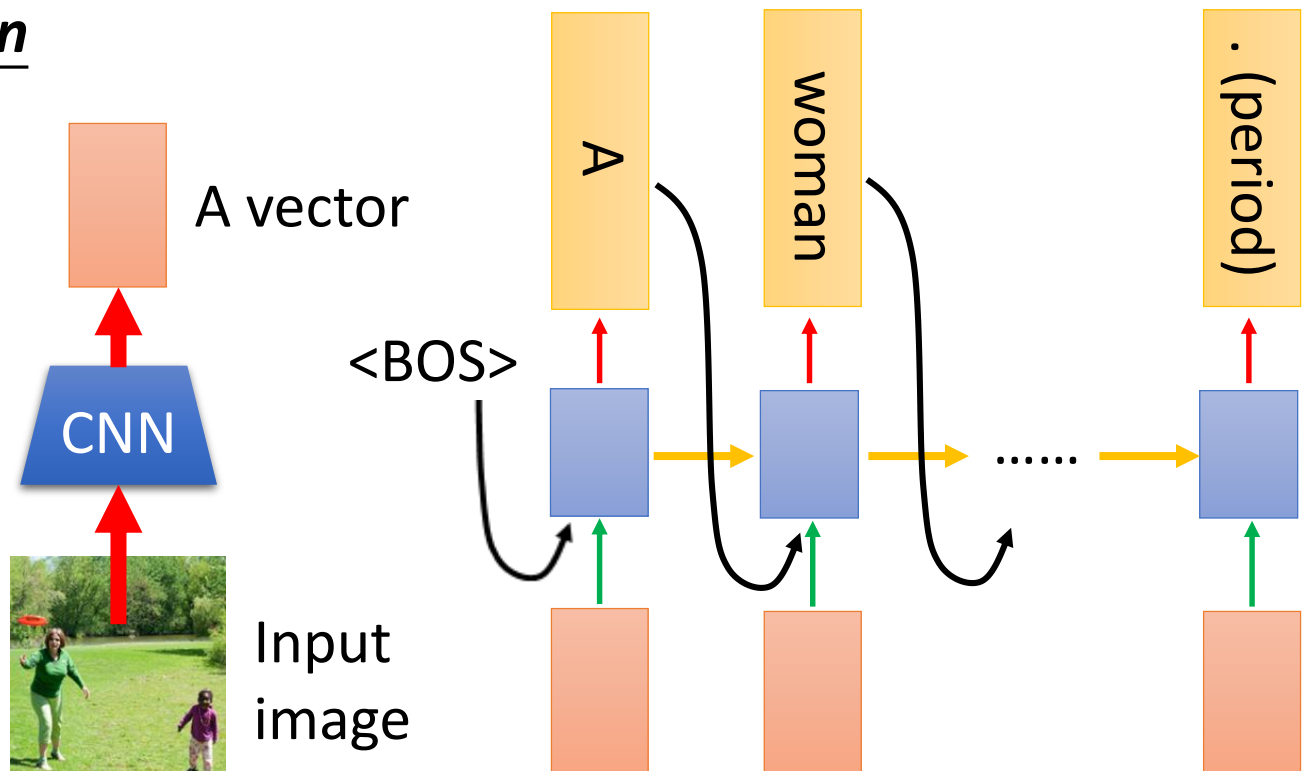
“Hello. Nice
to see you.”



Conditional Generation

- Represent the input condition as a vector, and consider the vector as the input of RNN generator

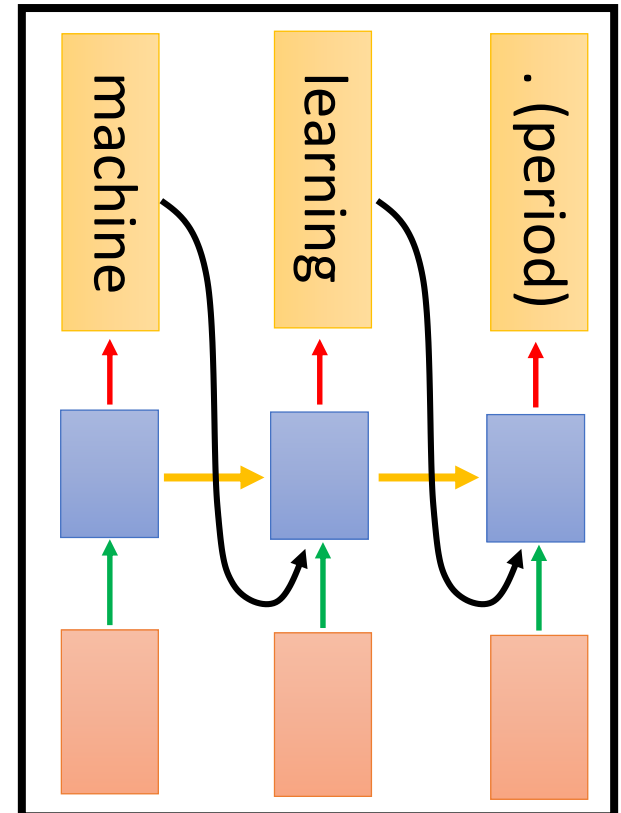
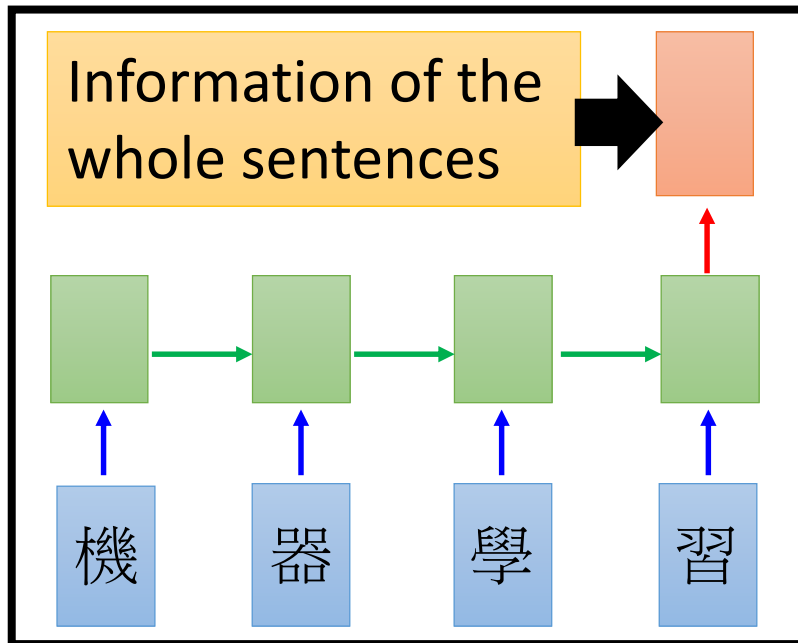
Image Caption Generation



Conditional Generation

Sequence-to-sequence learning

- Represent the input condition as a vector, and consider the vector as the input of RNN generator
- E.g. Machine translation / Chat-bot



Encoder ← Jointly train → Decoder

Conditional Generation

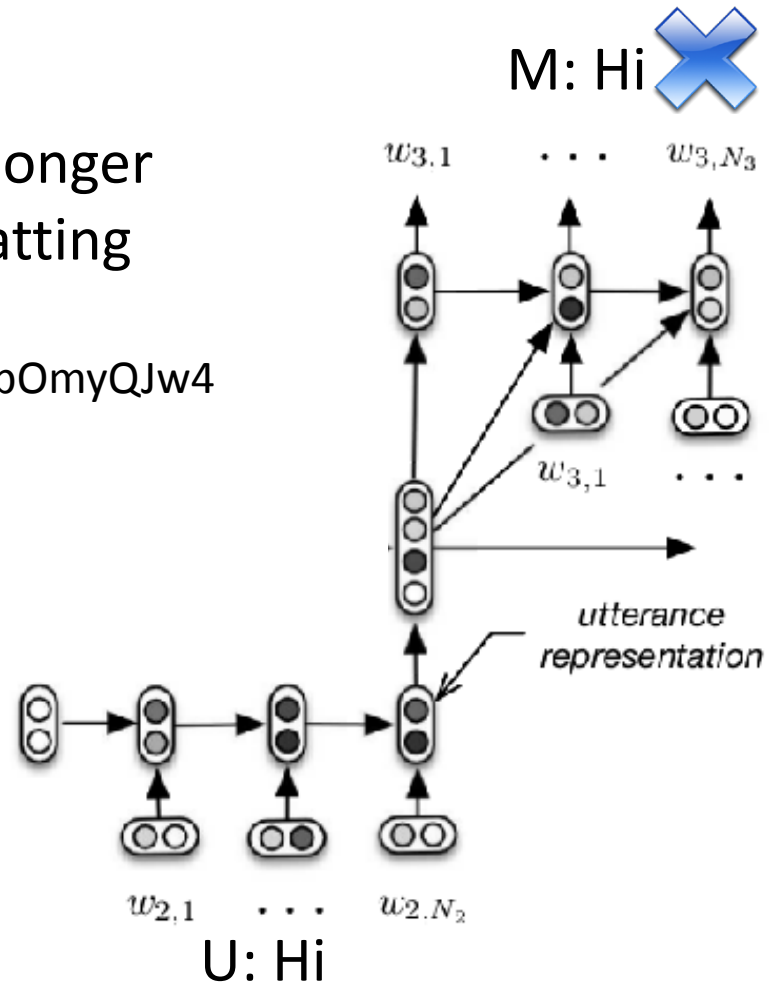
M: Hello

U: Hi

M: Hi

<https://www.youtube.com/watch?v=e2MpOmyQJw4>

Need to consider longer context during chatting



M: Hello

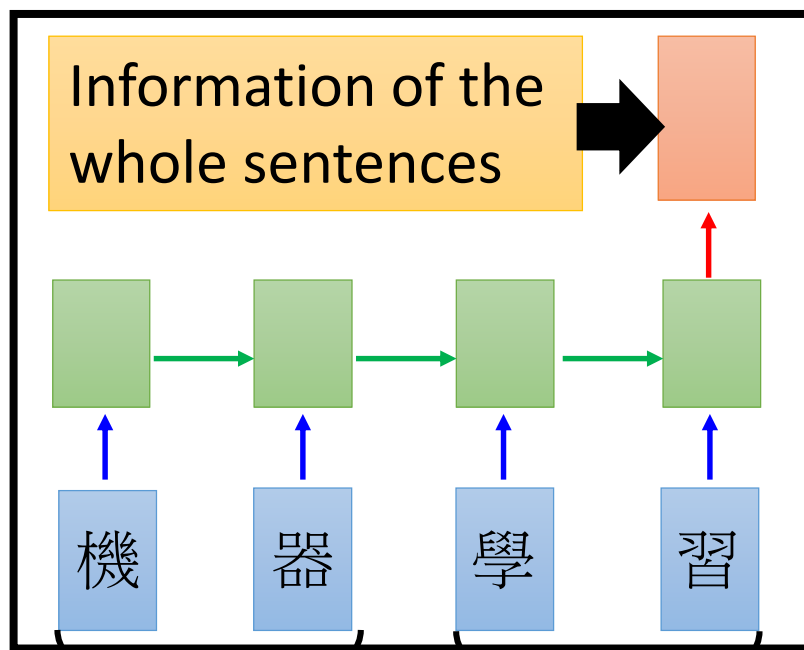
U: Hi

Attention

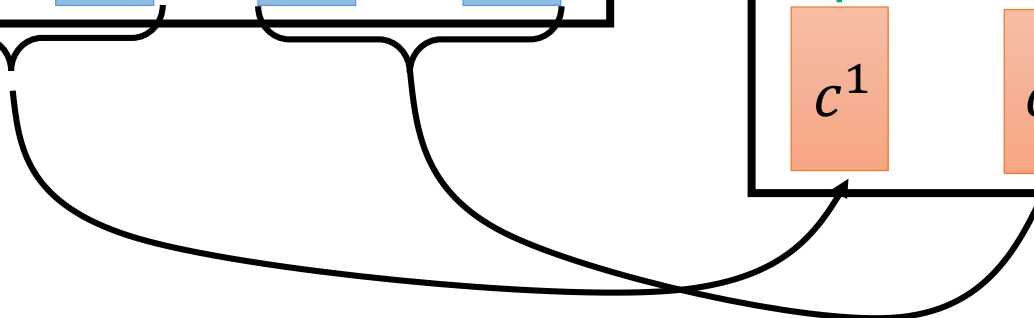
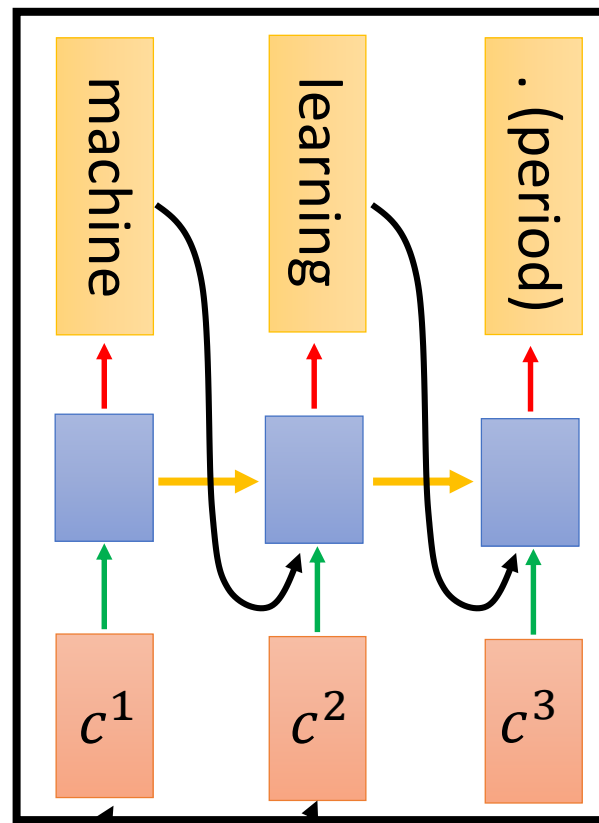
Dynamic Conditional Generation

Dynamic Conditional Generation

Encoder

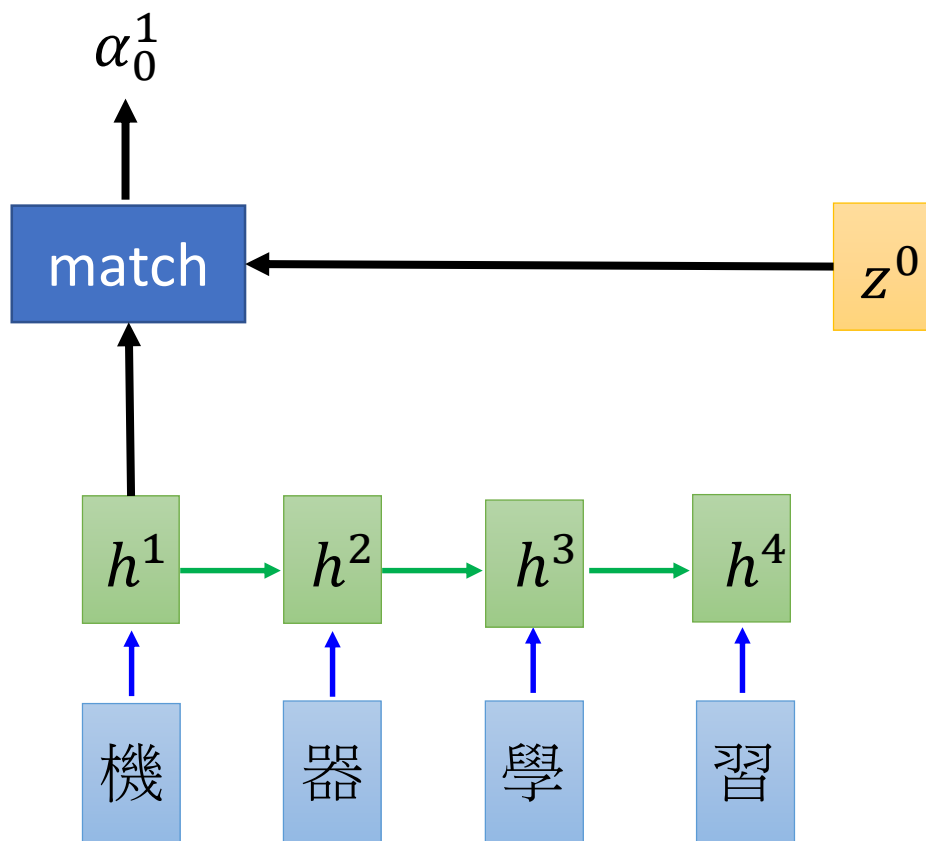


Decoder

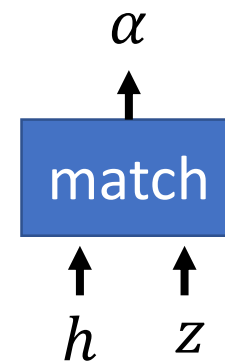


Machine Translation

- Attention-based model



Jointly learned
with other part
of the network



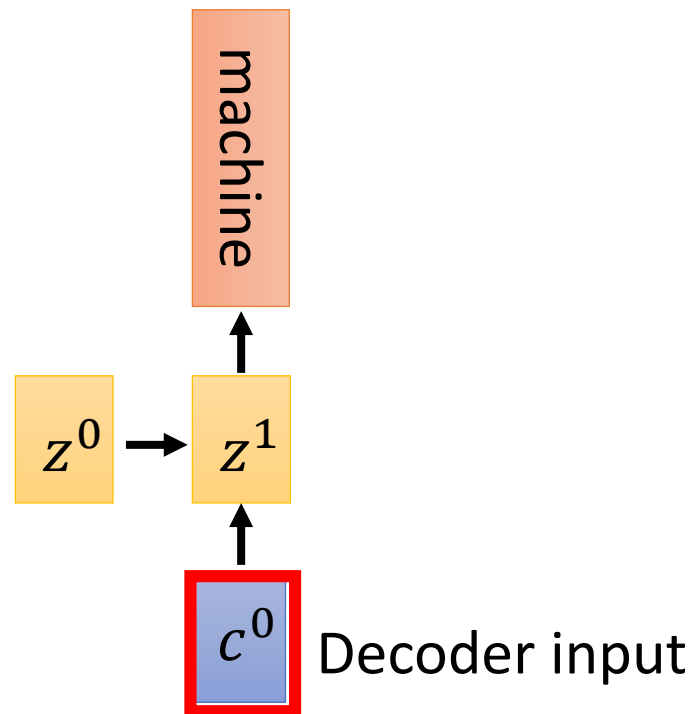
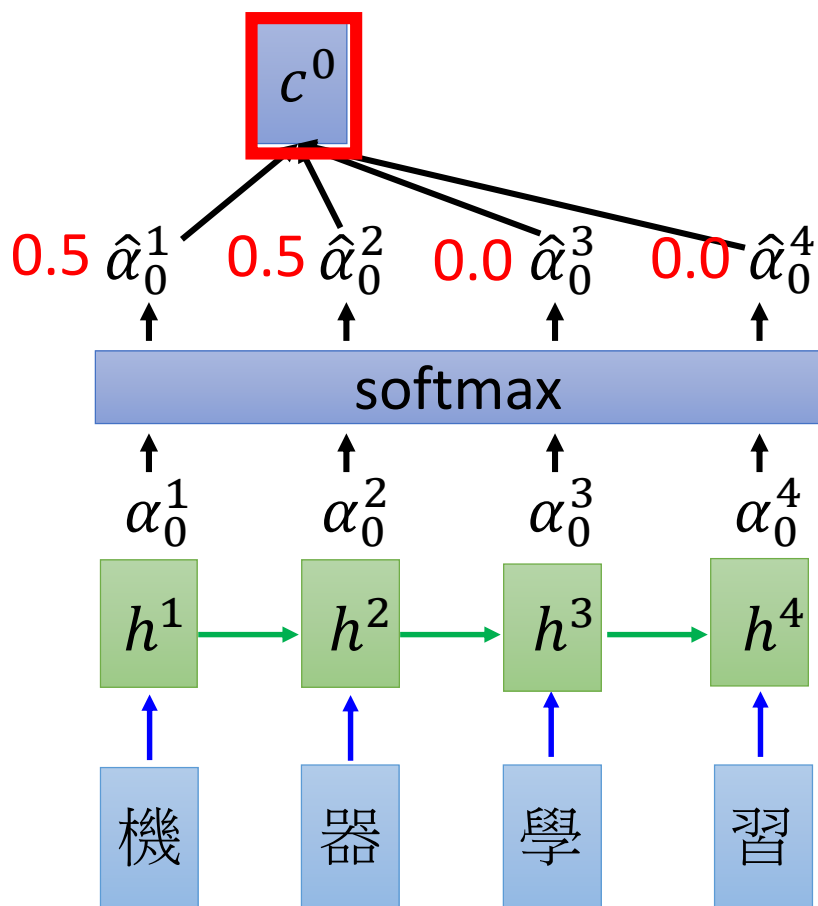
What is **match** ?

Design by yourself

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

Machine Translation

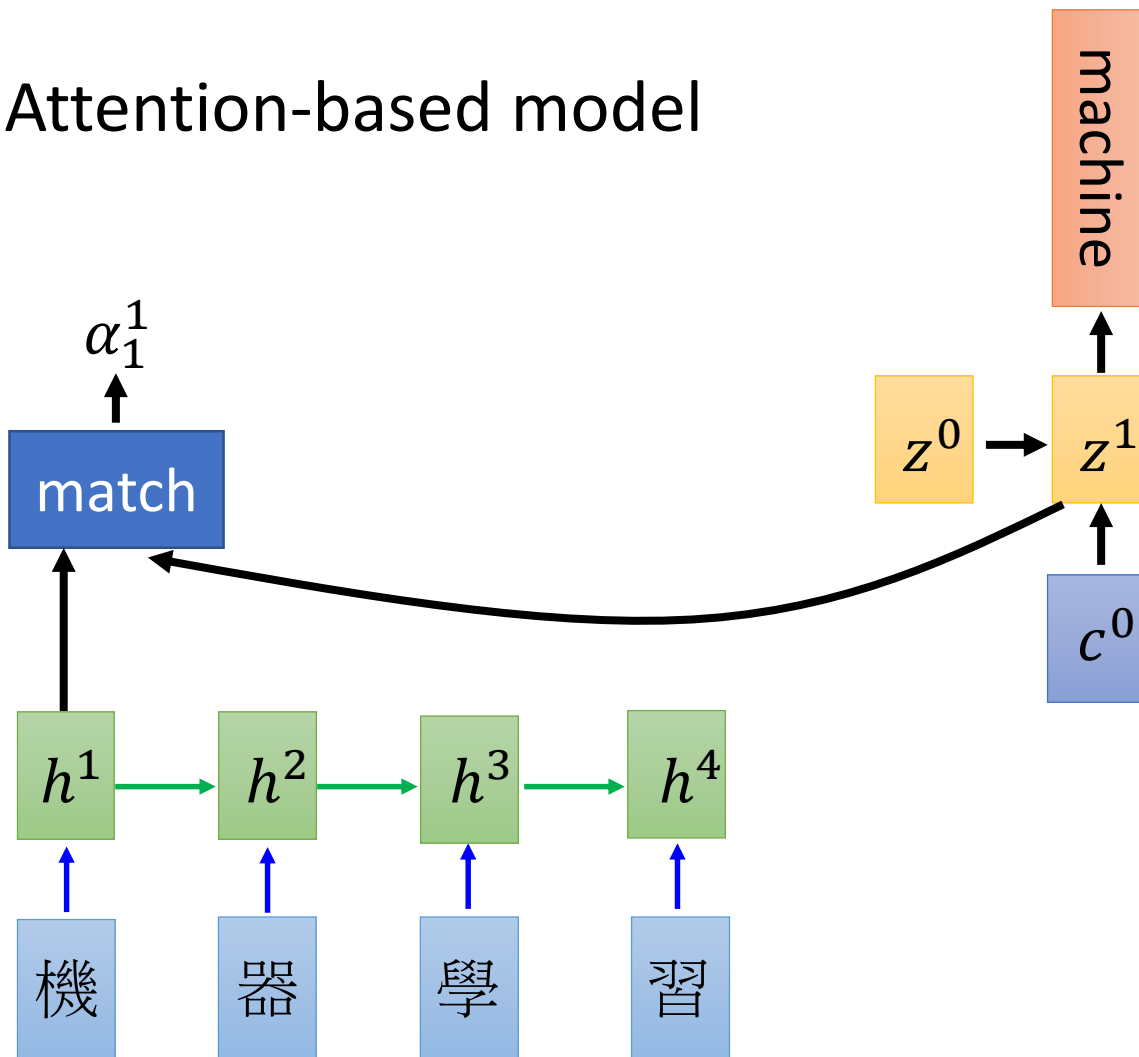
- Attention-based model



$$c^0 = \sum \hat{\alpha}_0^i h^i$$
$$= 0.5h^1 + 0.5h^2$$

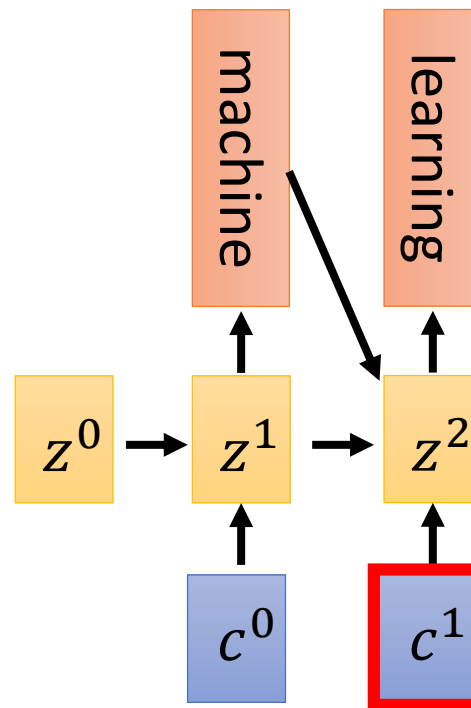
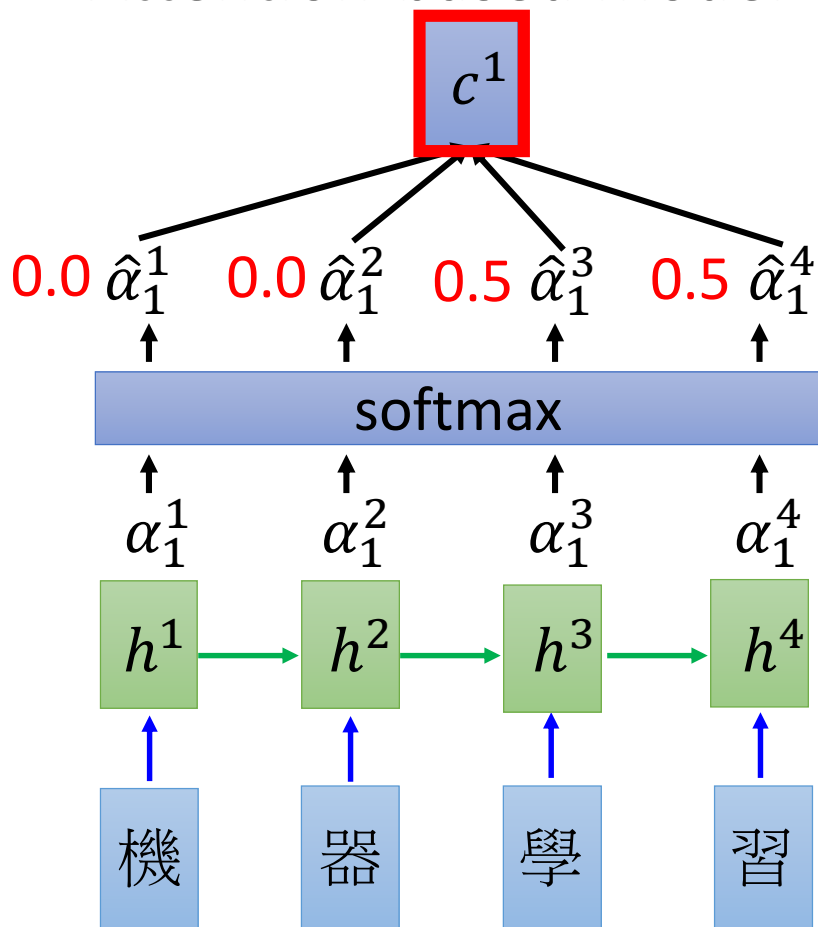
Machine Translation

- Attention-based model



Machine Translation

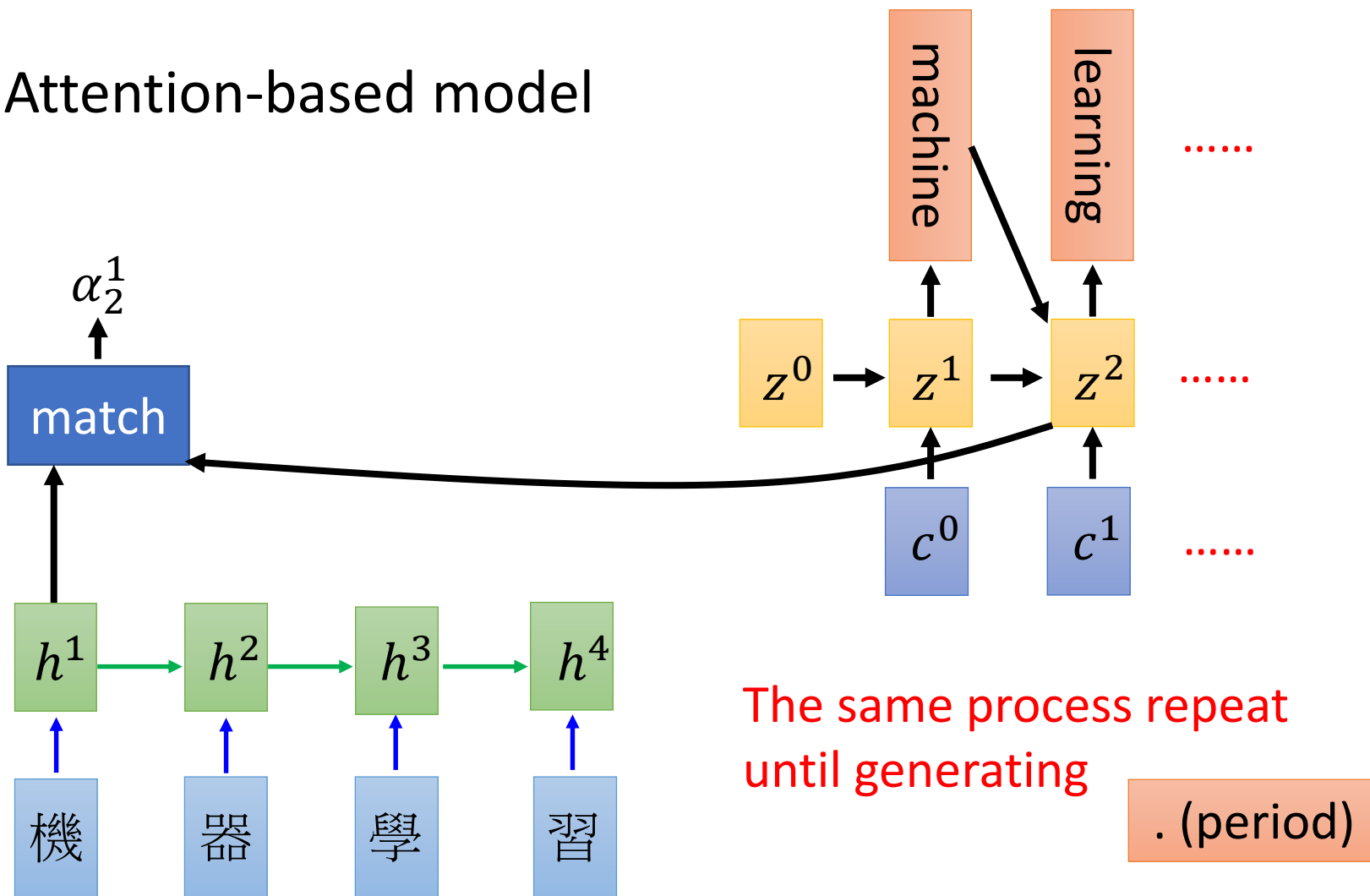
- Attention-based model



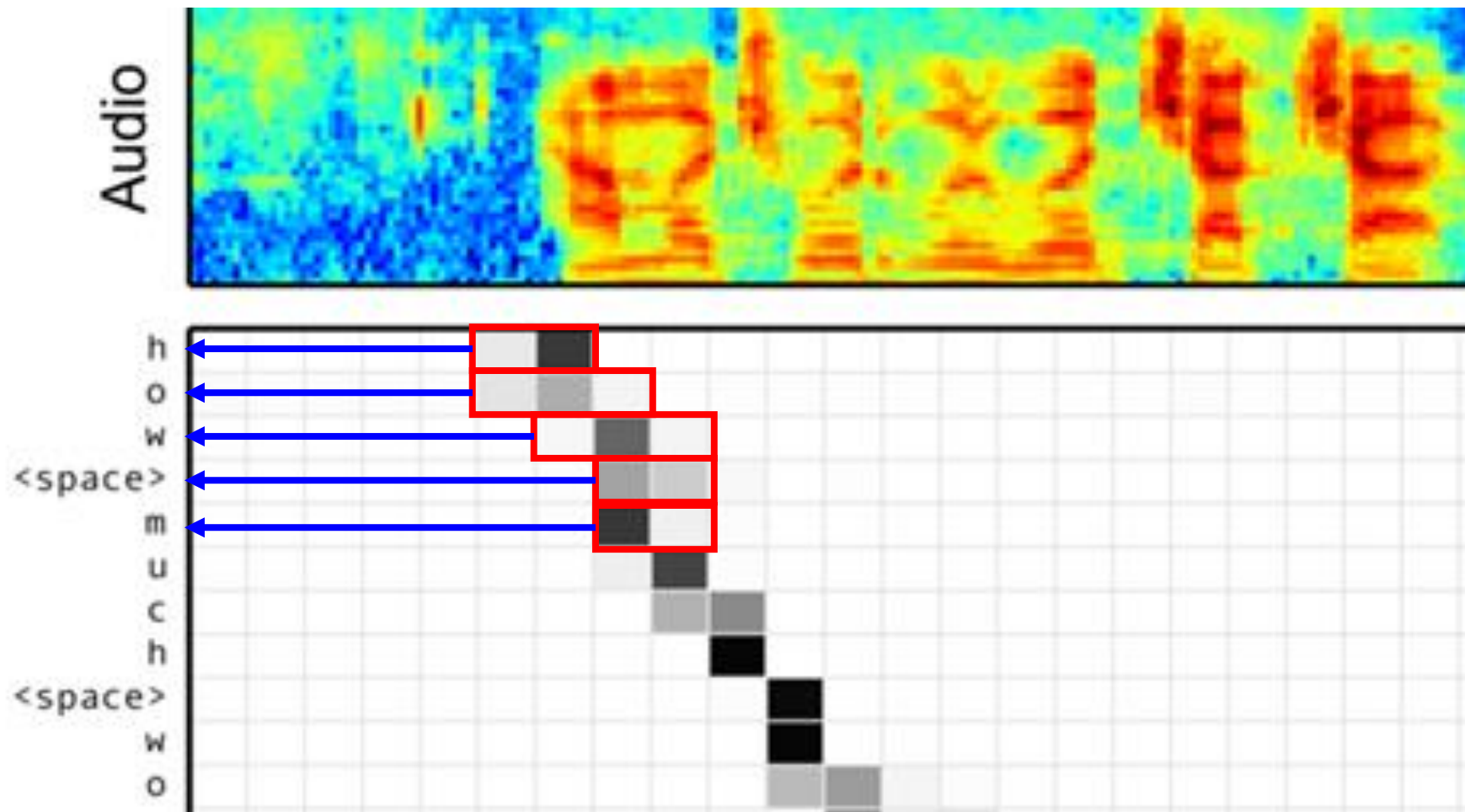
$$c^1 = \sum \hat{\alpha}_1^i h^i$$
$$= 0.5h^3 + 0.5h^4$$

Machine Translation

- Attention-based model



Speech Recognition



Model	Clean WER	Noisy WER
CLDNN-HMM [22]	8.0	8.9
LAS	14.1	16.5
LAS + LM Rescoring	10.3	12.0

William Chan, Navdeep Jaitly, Quoc V. Le, Oriol Vinyals, "Listen, Attend and Spell", ICASSP, 2016

Image Caption Generation

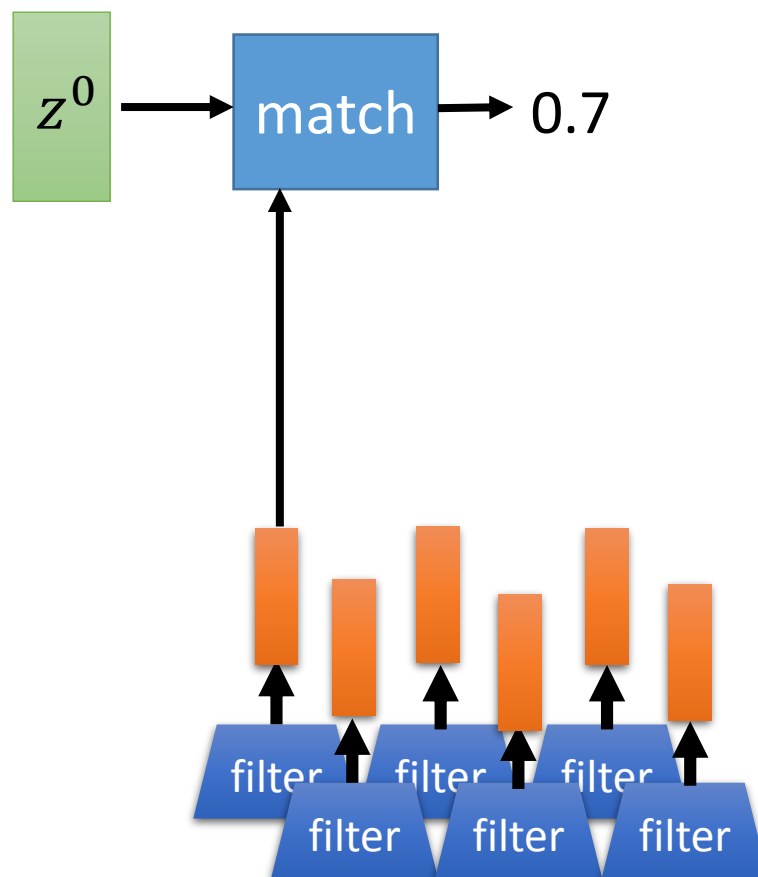
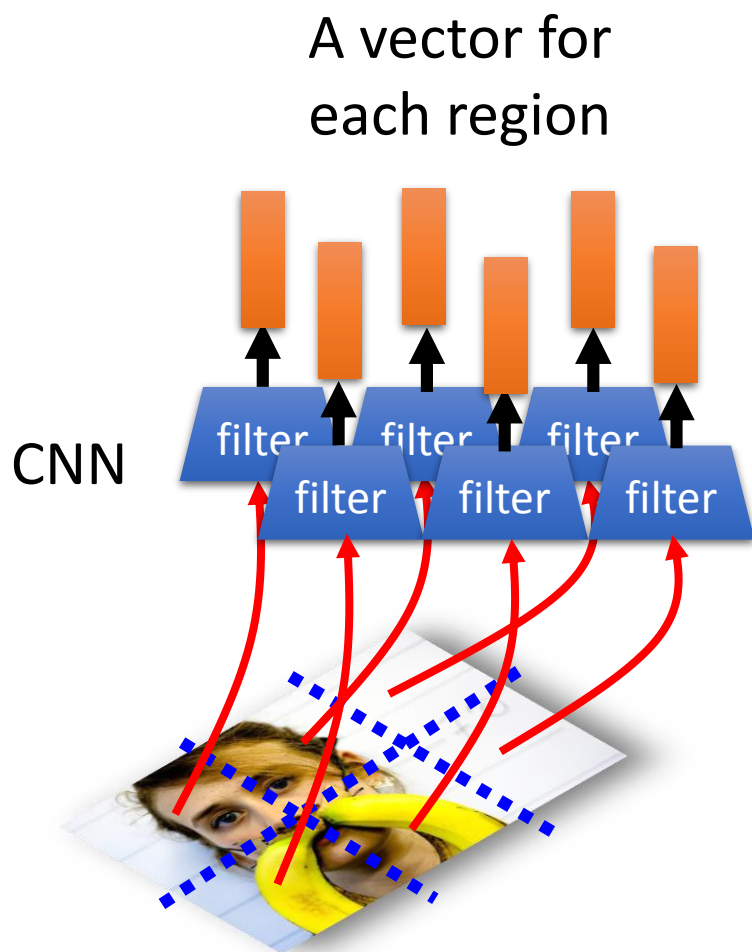


Image Caption Generation

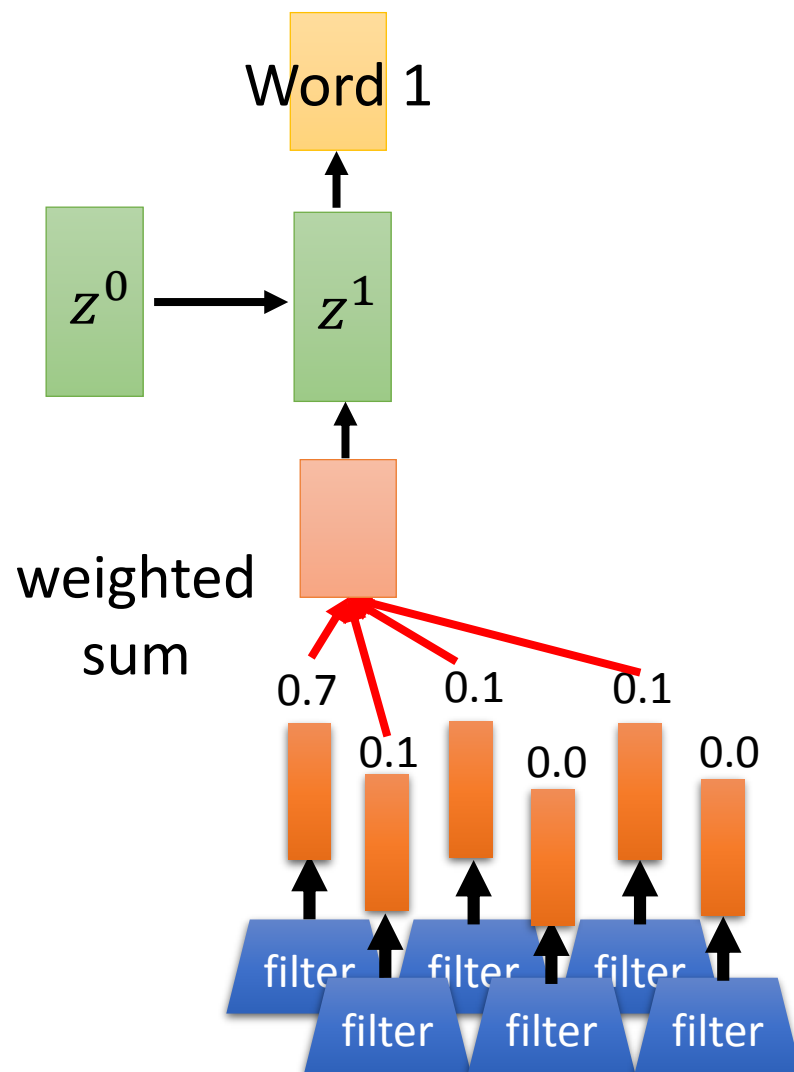
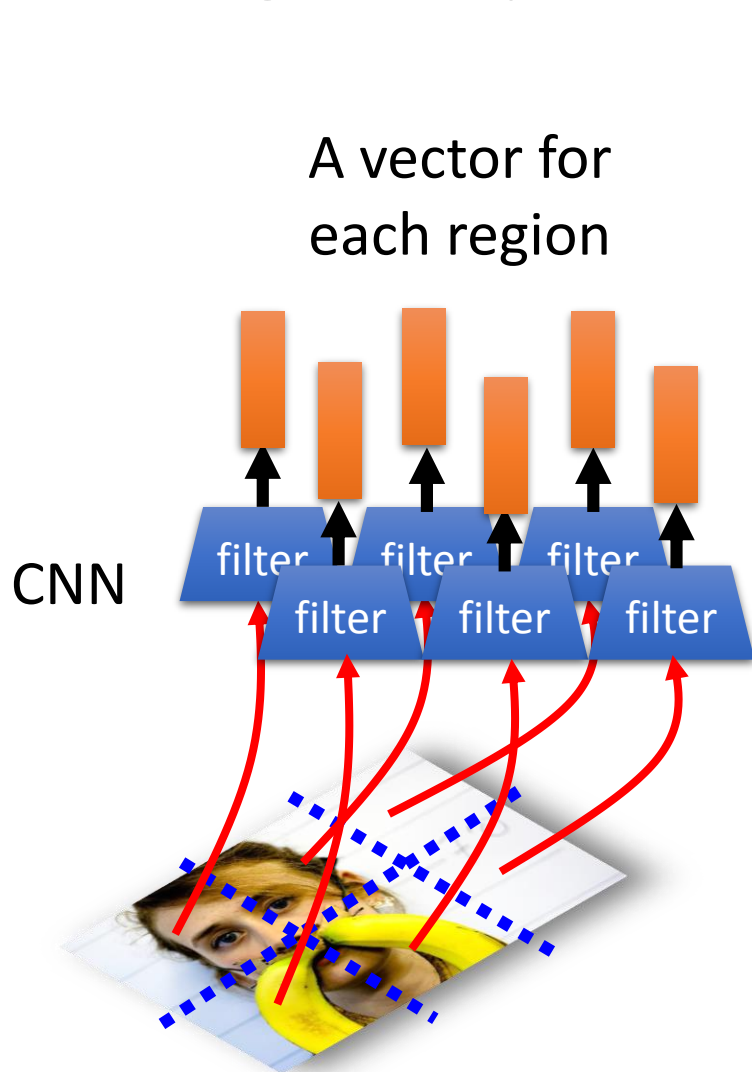


Image Caption Generation

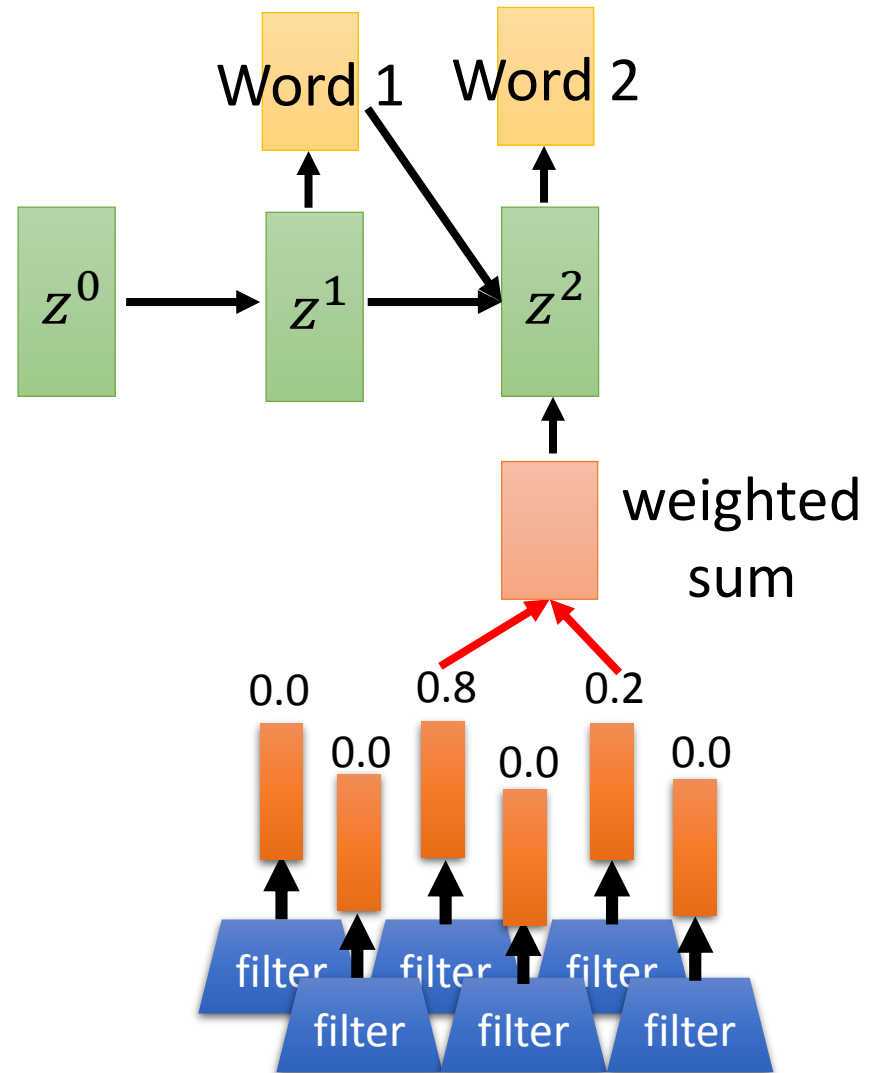
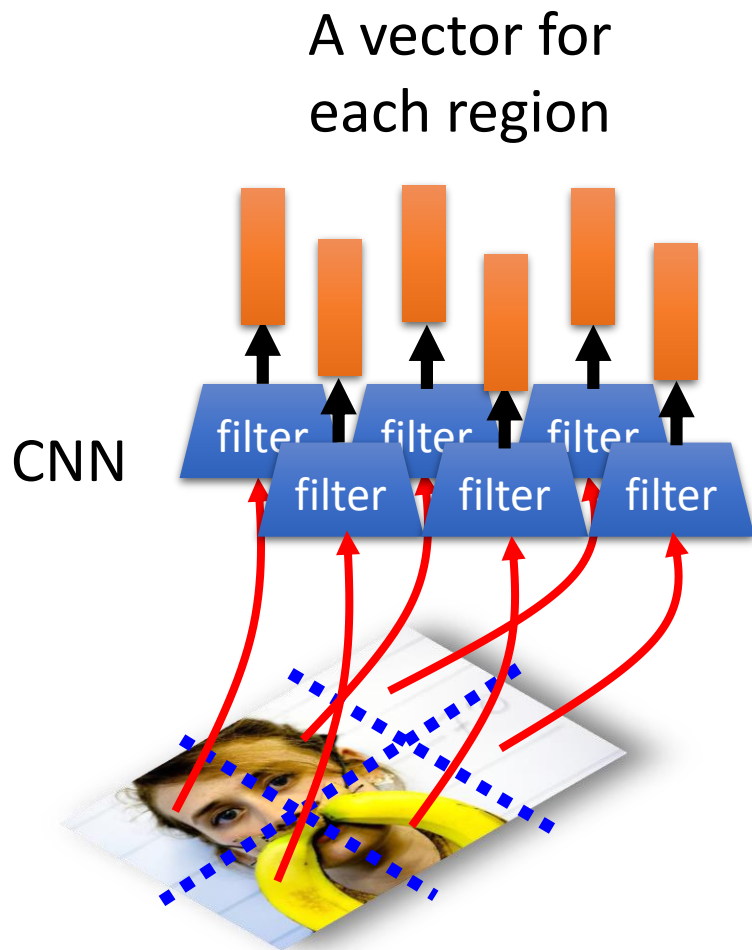
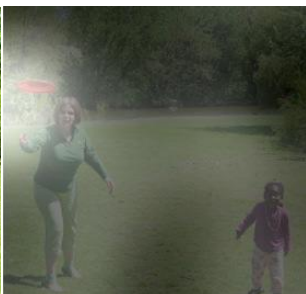


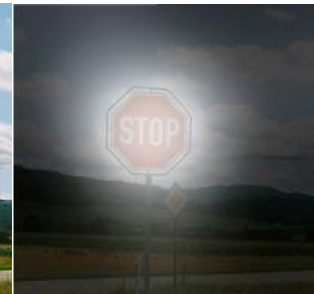
Image Caption Generation



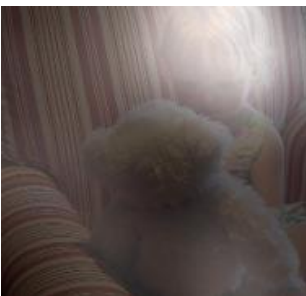
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



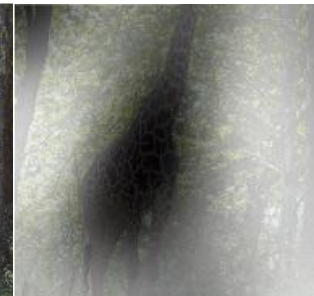
A stop sign is on a road with a mountain in the background.



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



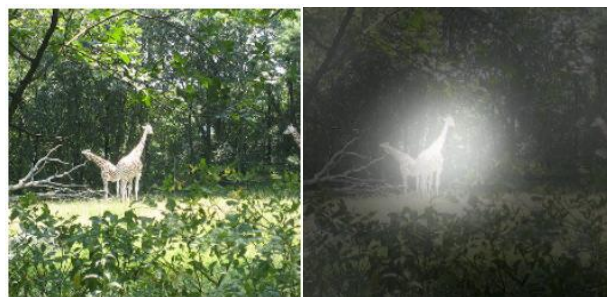
A group of people sitting on a boat in the water.



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

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML, 2015

Image Caption Generation



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.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML, 2015



Ref: A man and a woman ride a motorcycle

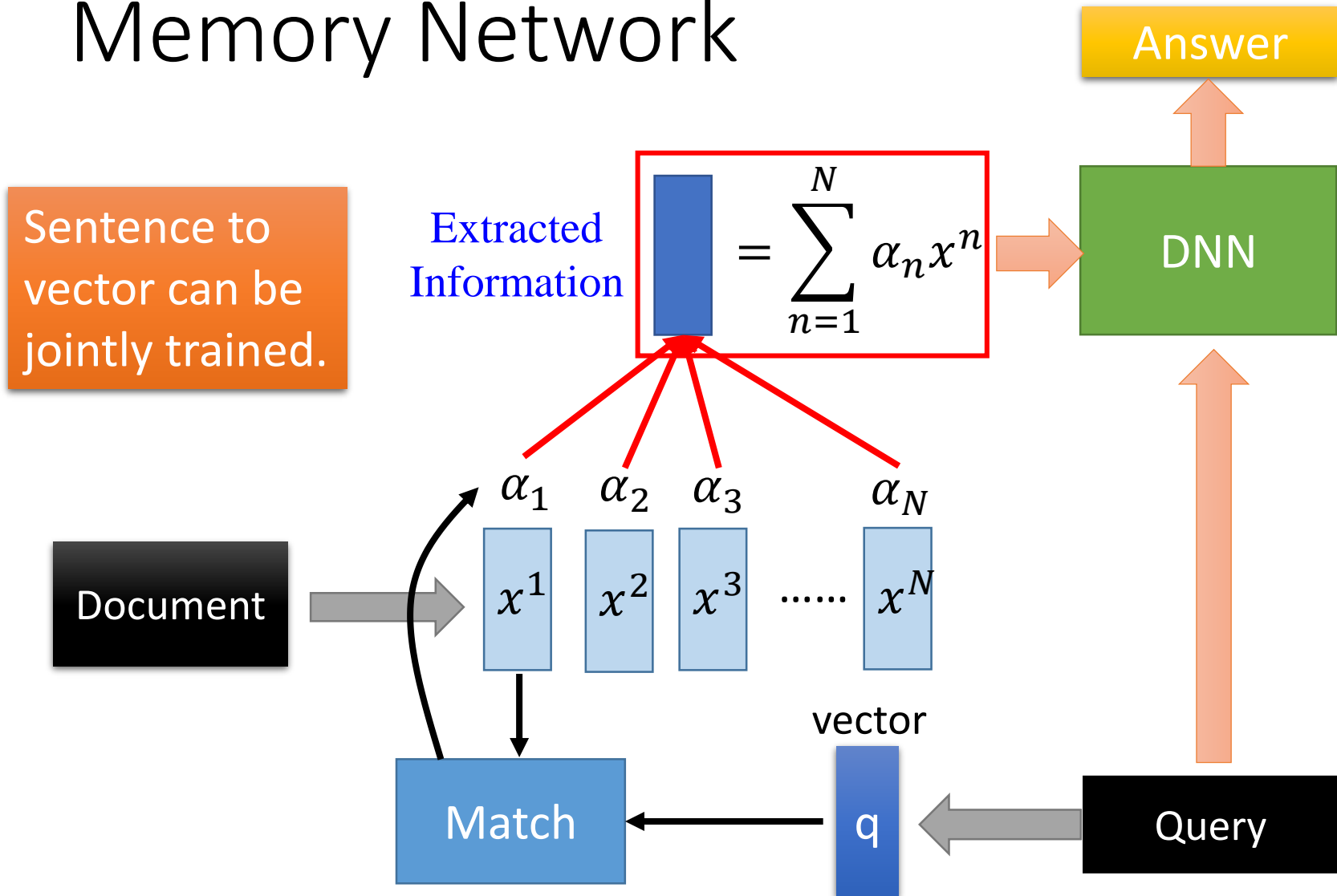
A **man** and a **woman** are **talking** on the **road**



Ref: A woman is frying food

Someone is **frying** a **fish** in a **pot**

Memory Network



Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, Rob Fergus, "End-To-End Memory Networks", NIPS, 2015

Memory Network

Jointly learned

Document

Extracted Information

$$\begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} = \sum_{n=1}^N \alpha_n h^n$$

h^1 h^2 h^3 h^N

α_1 α_2 α_3 α_N
 x^1 x^2 x^3 x^N

Match

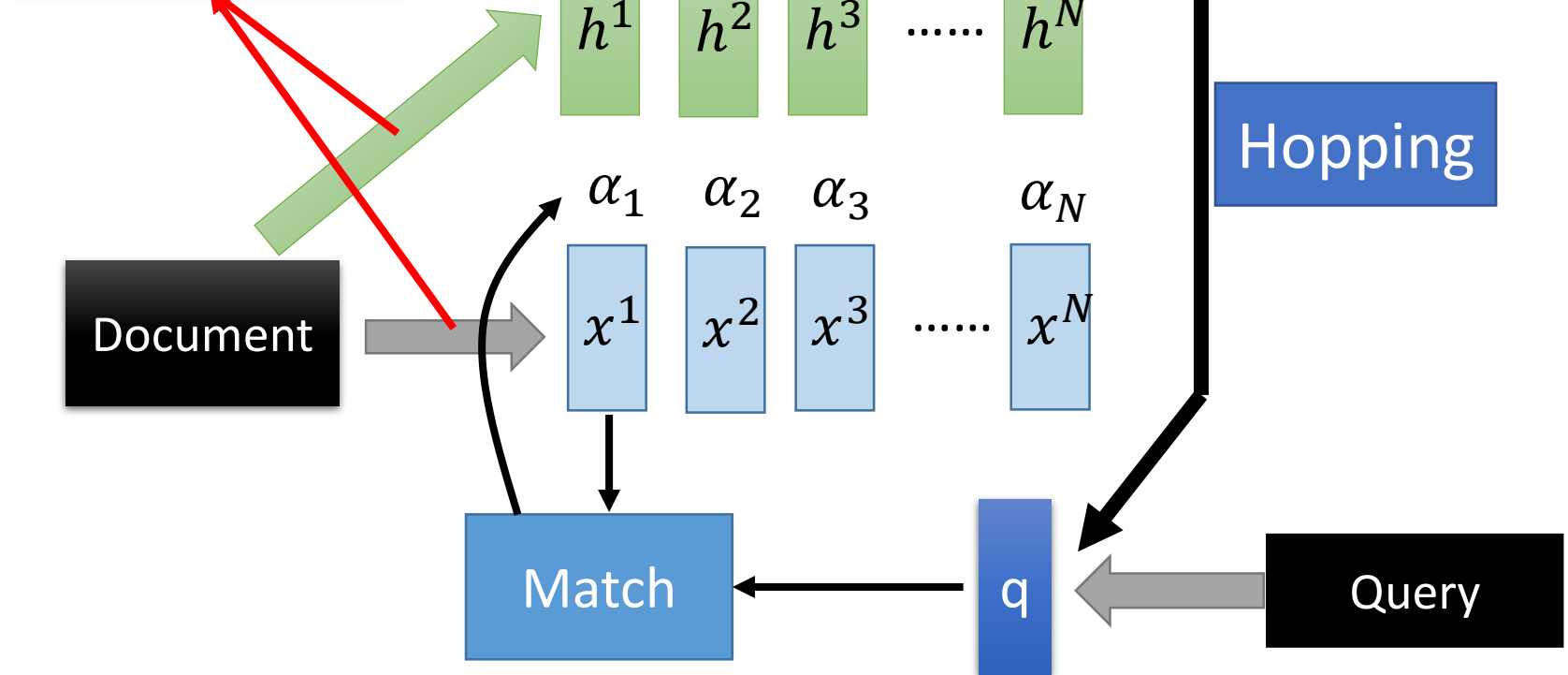
q

Hopping

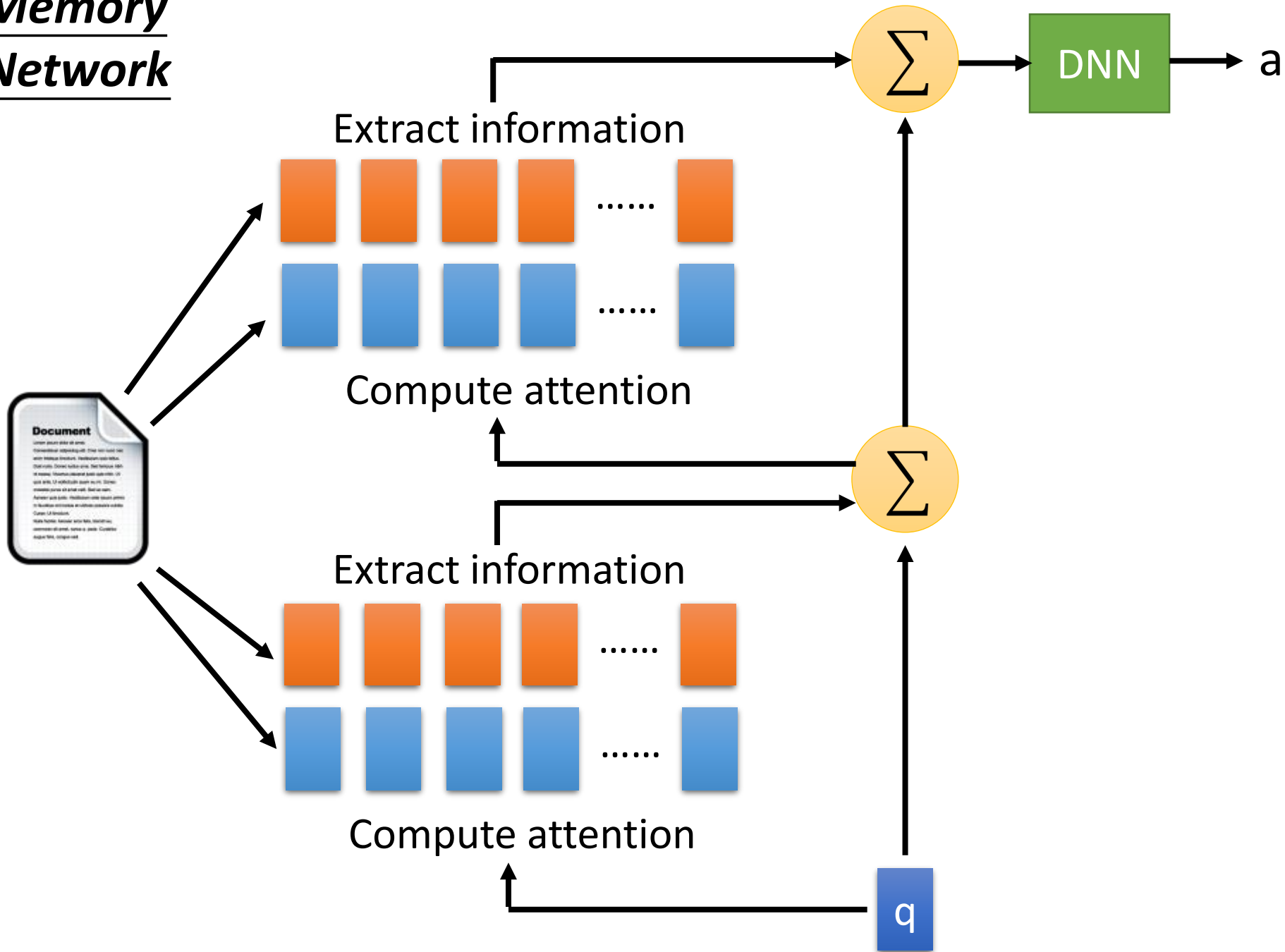
Query

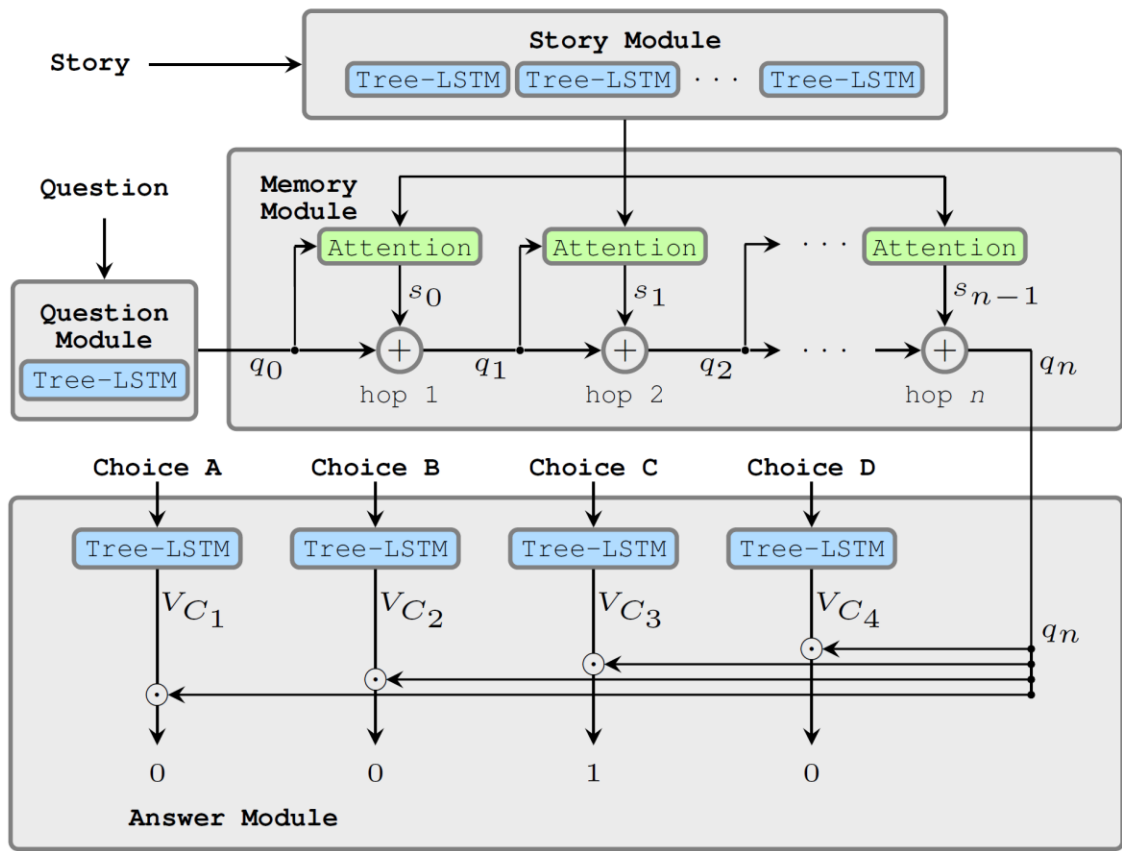
DNN

Answer

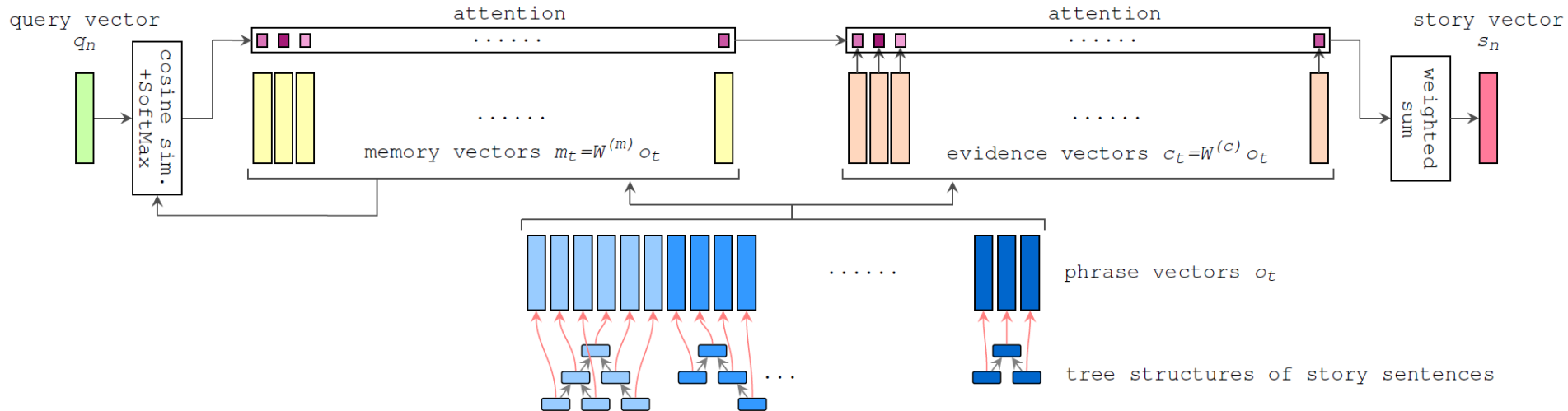
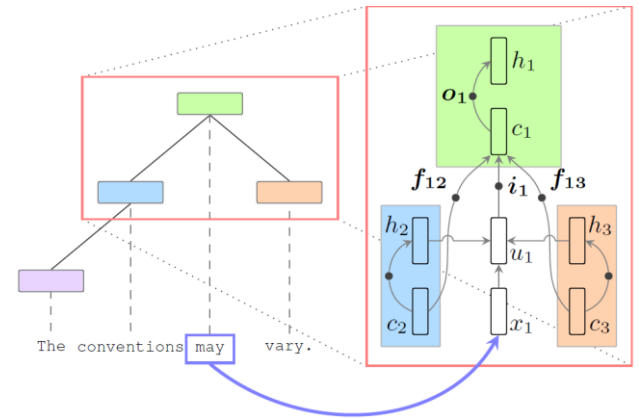


Memory Network





Wei Fang, Juei-Yang Hsu, Hung-yi Lee, Lin-Shan Lee, "Hierarchical Attention Model for Improved Machine Comprehension of Spoken Content", SLT, 2016

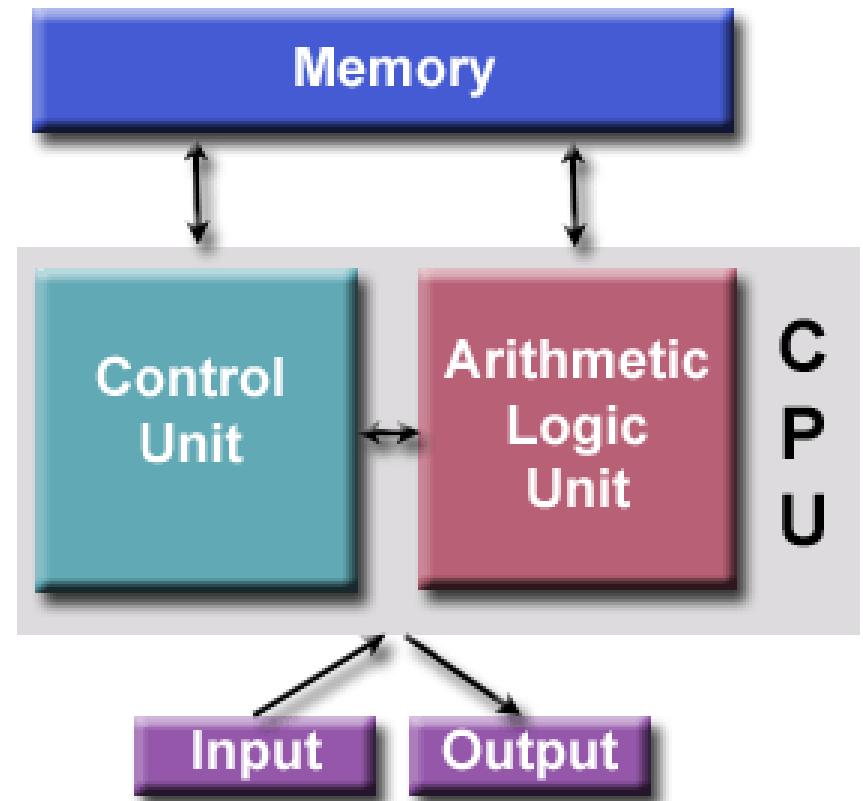


Neural Turing Machine

- von Neumann architecture

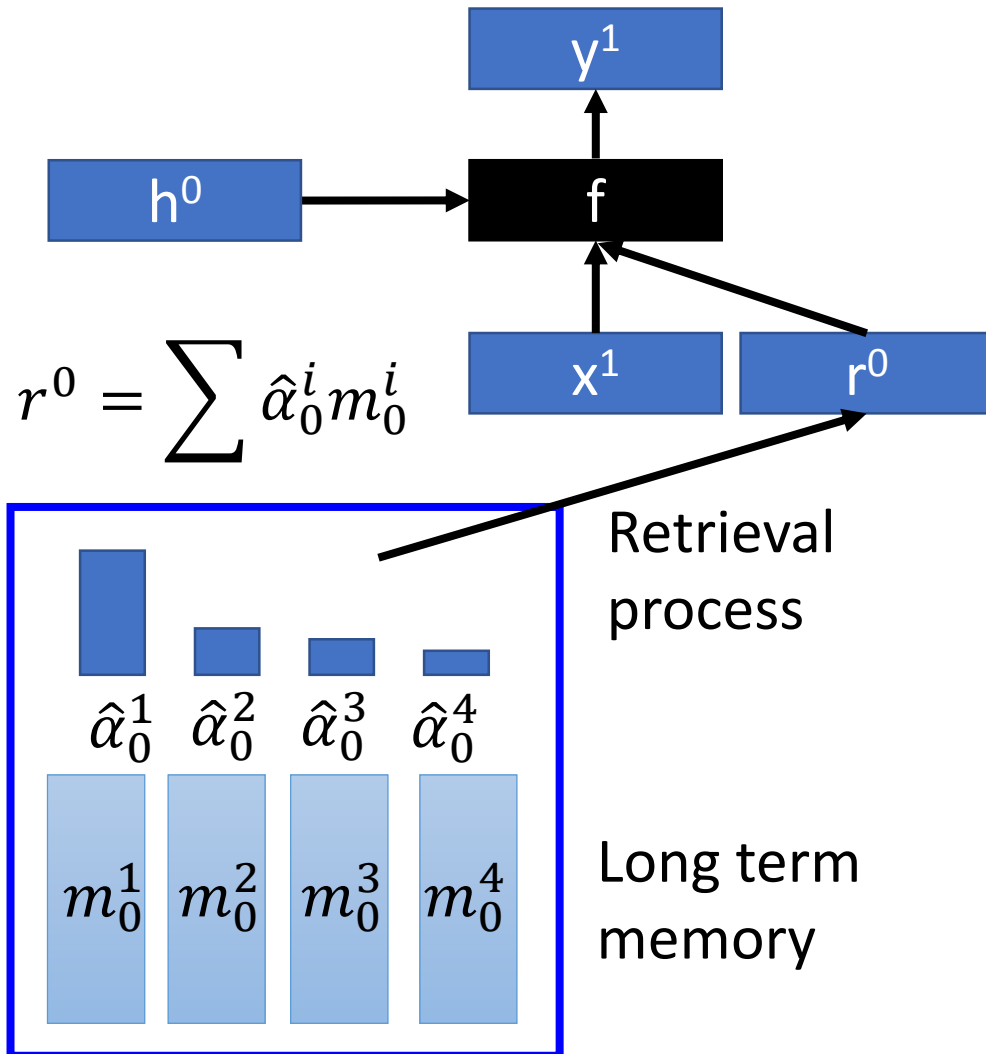
Neural Turing Machine
not only read from
memory

Also modify the memory
through attention

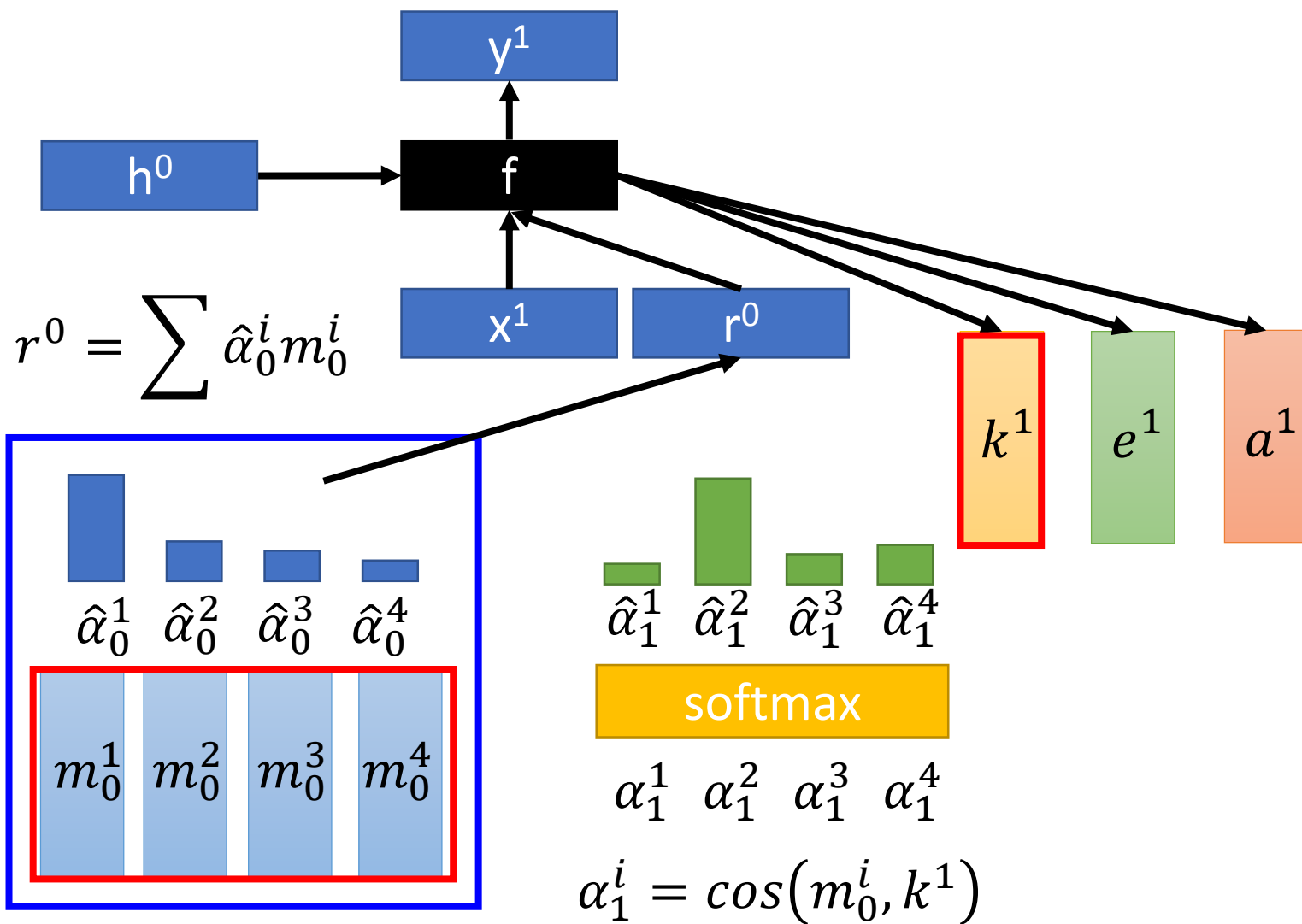


<https://www.quora.com/How-does-the-Von-Neumann-architecture-provide-flexibility-for-program-development>

Neural Turing Machine

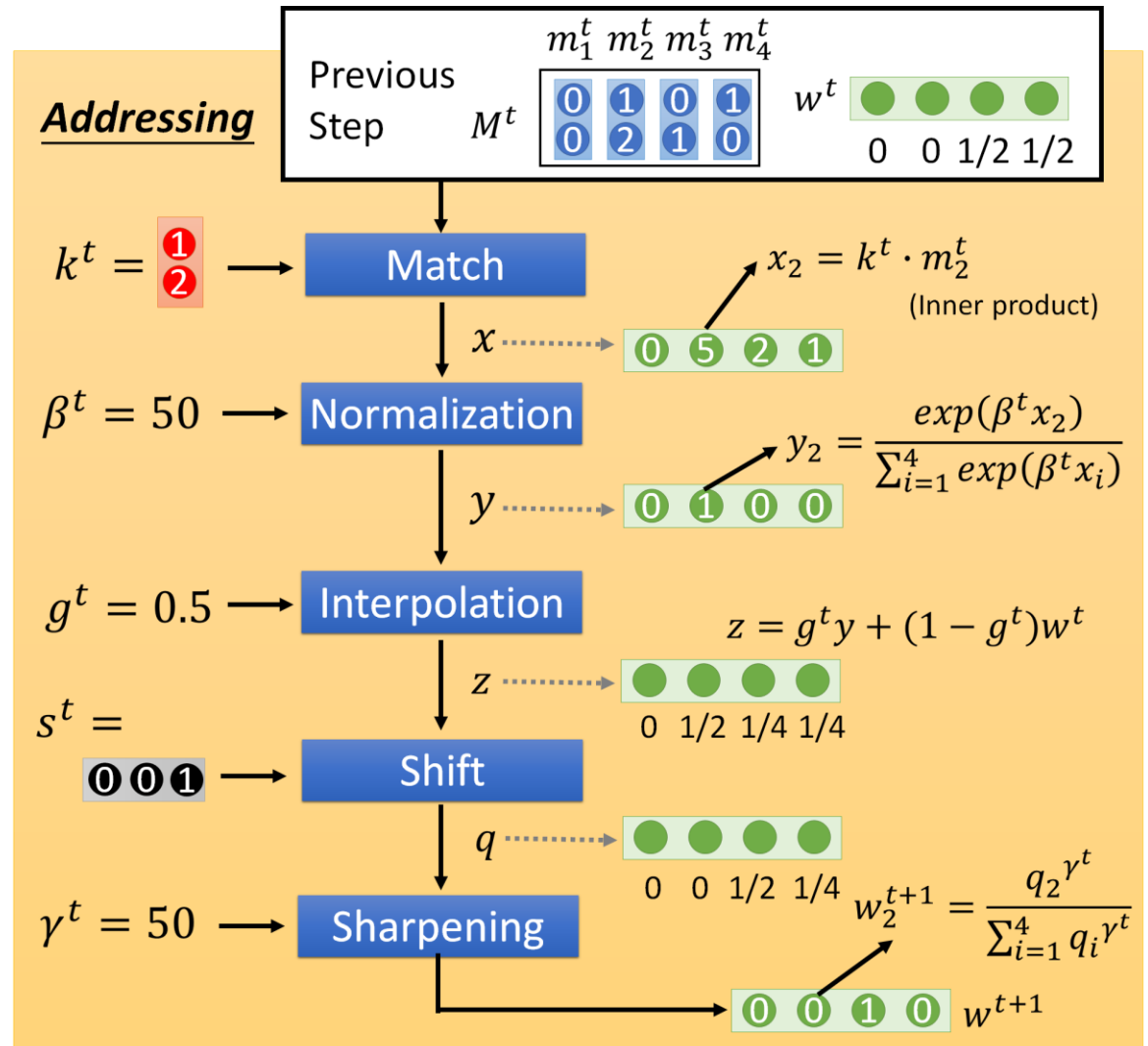


Neural Turing Machine



Neural Turing Machine

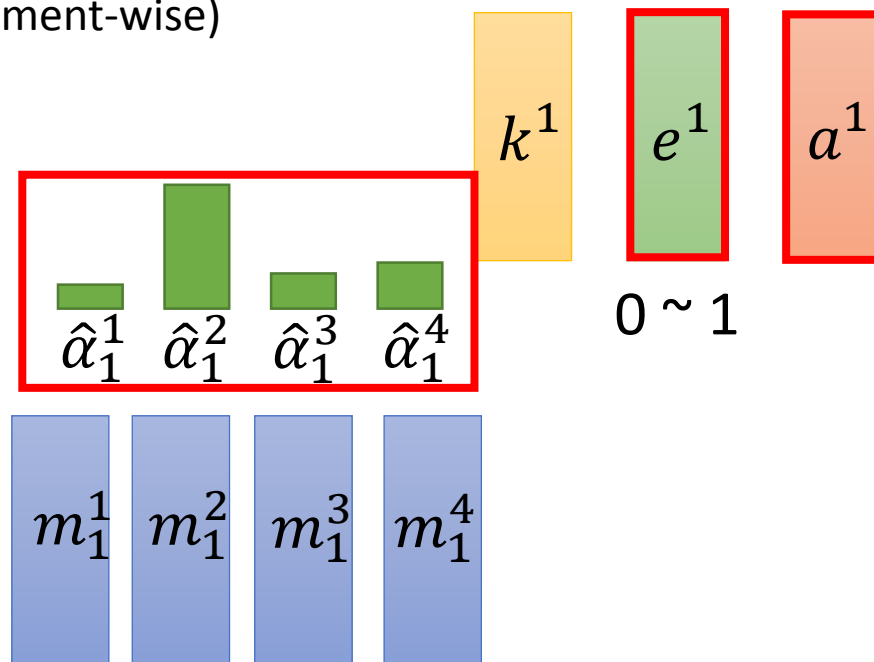
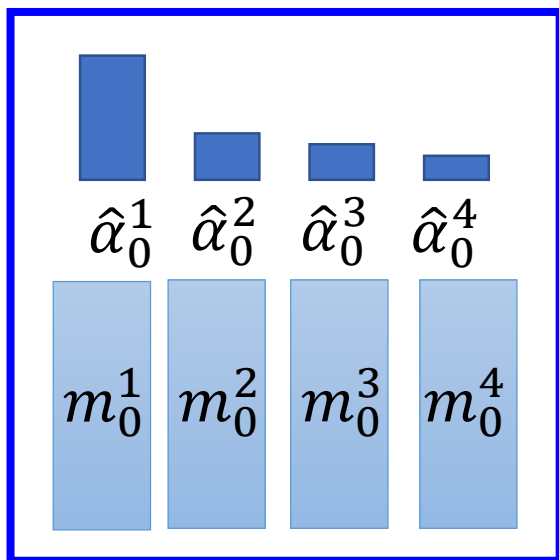
- Real version



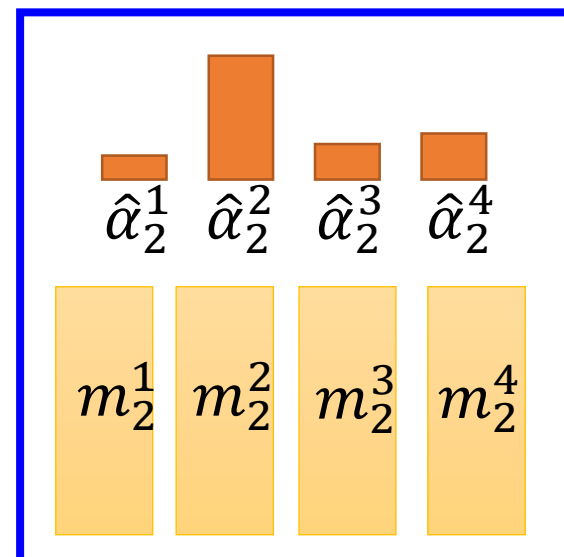
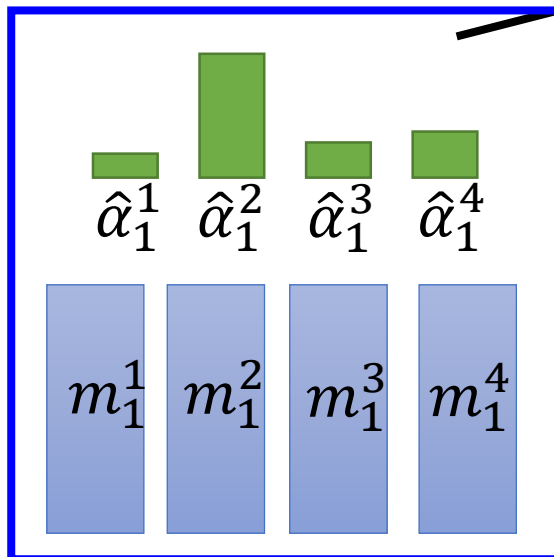
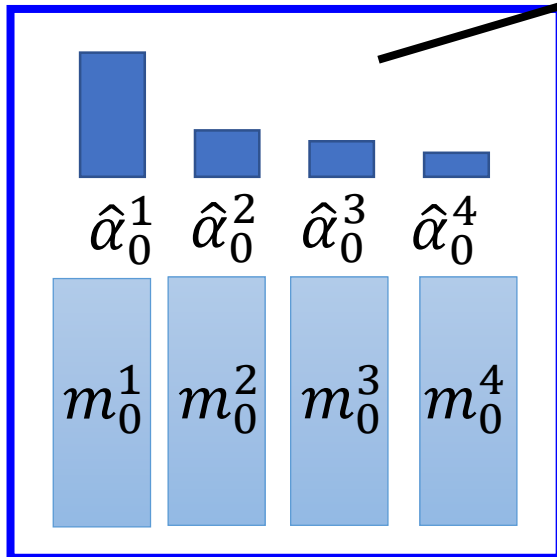
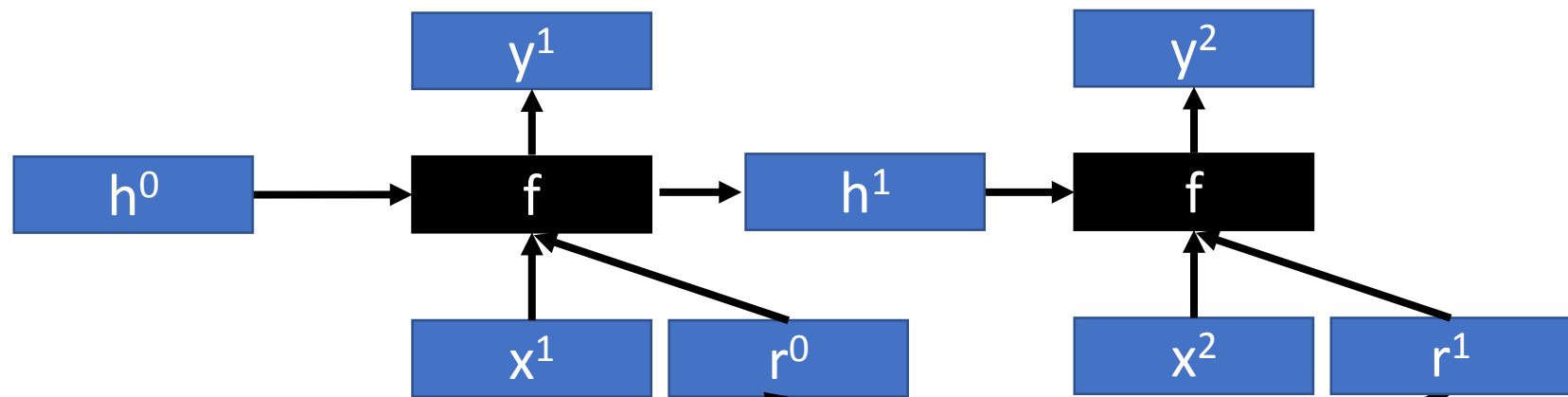
Neural Turing Machine

$$m_1^i = m_0^i - \hat{\alpha}_1^i e^1 \odot m_0^i + \hat{\alpha}_1^i a^1$$

(element-wise)



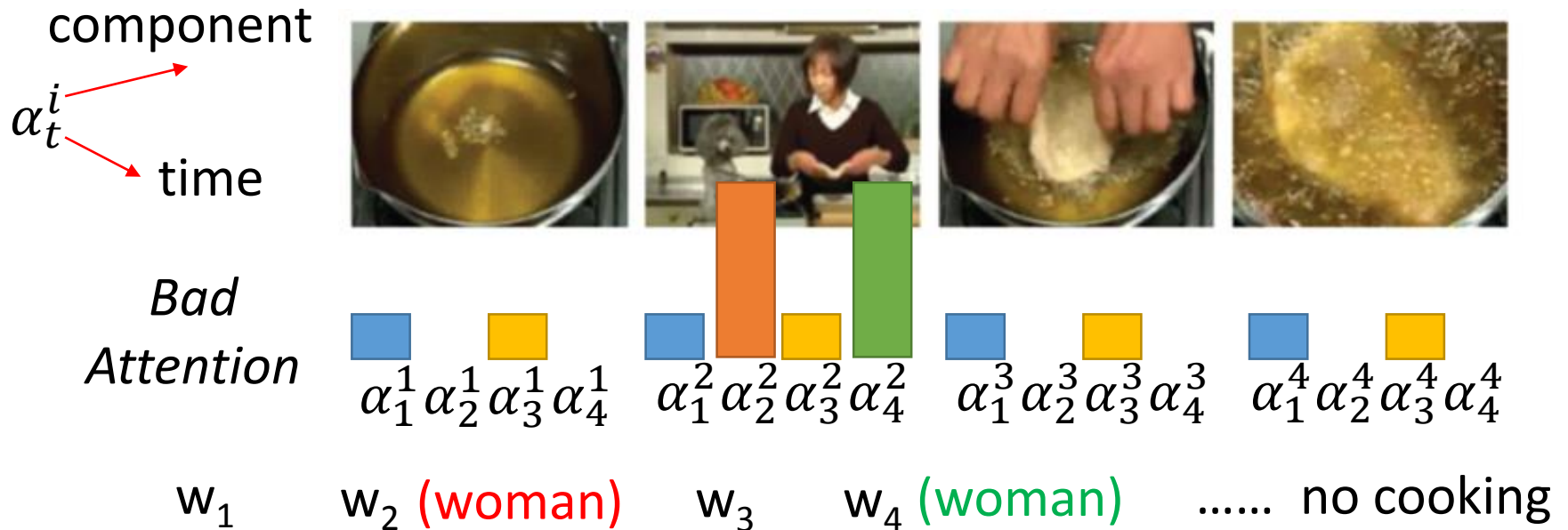
Neural Turing Machine



Tips for Generation

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML, 2015

Attention



Good Attention: each input component has approximately the same attention weight

E.g. Regularization term: $\sum_i \left(\tau - \sum_t \alpha_t^i \right)^2$

For each component

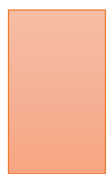
Over the generation

Mismatch between Train and Test

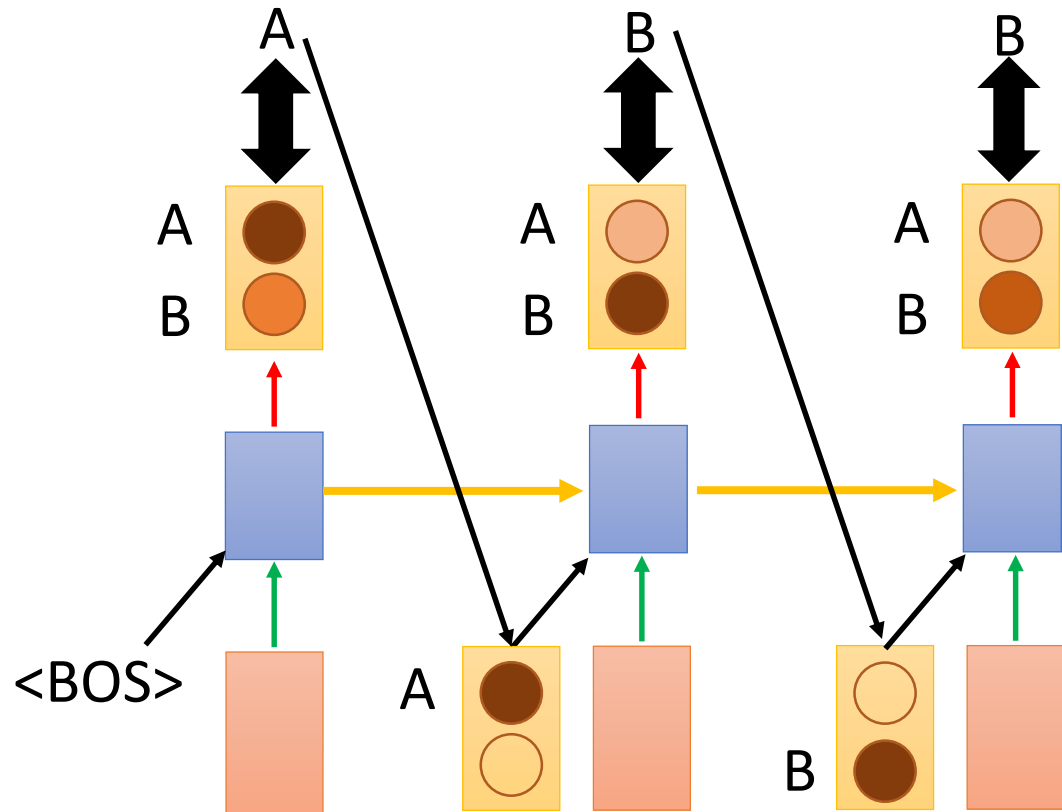
- Training

$$C = \sum_t C_t$$

Minimizing cross-entropy of each component

 : condition

Reference:



Mismatch between Train and Test

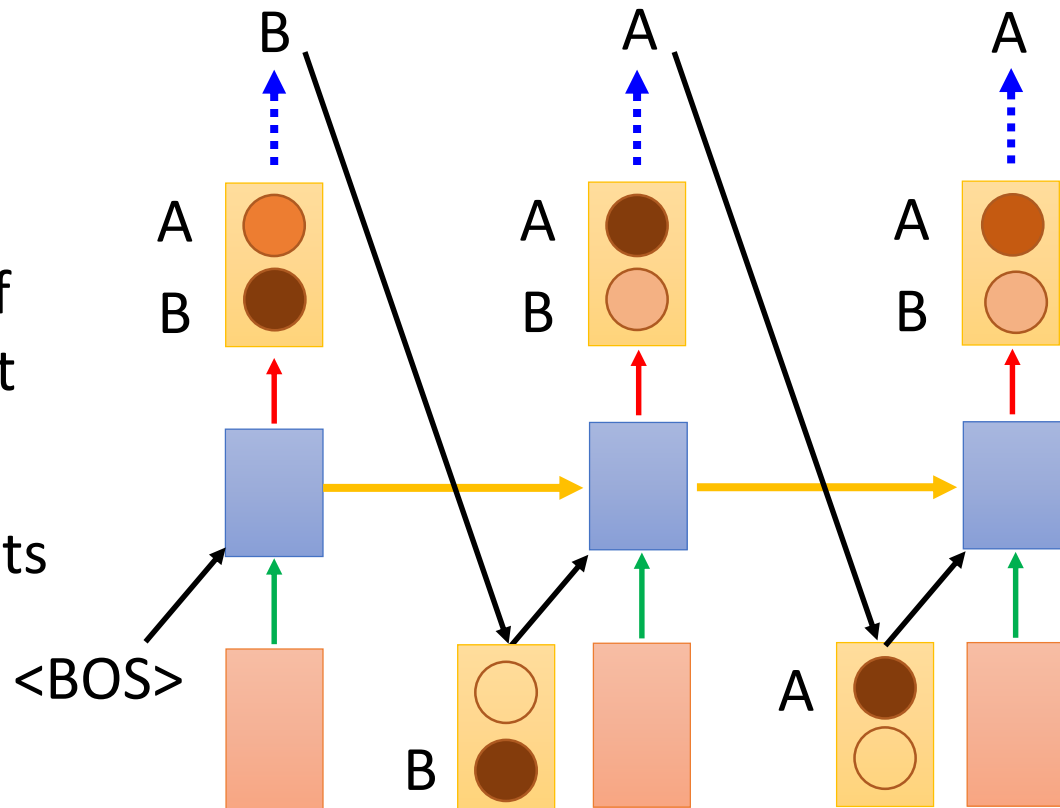
- Generation

We do not know the reference

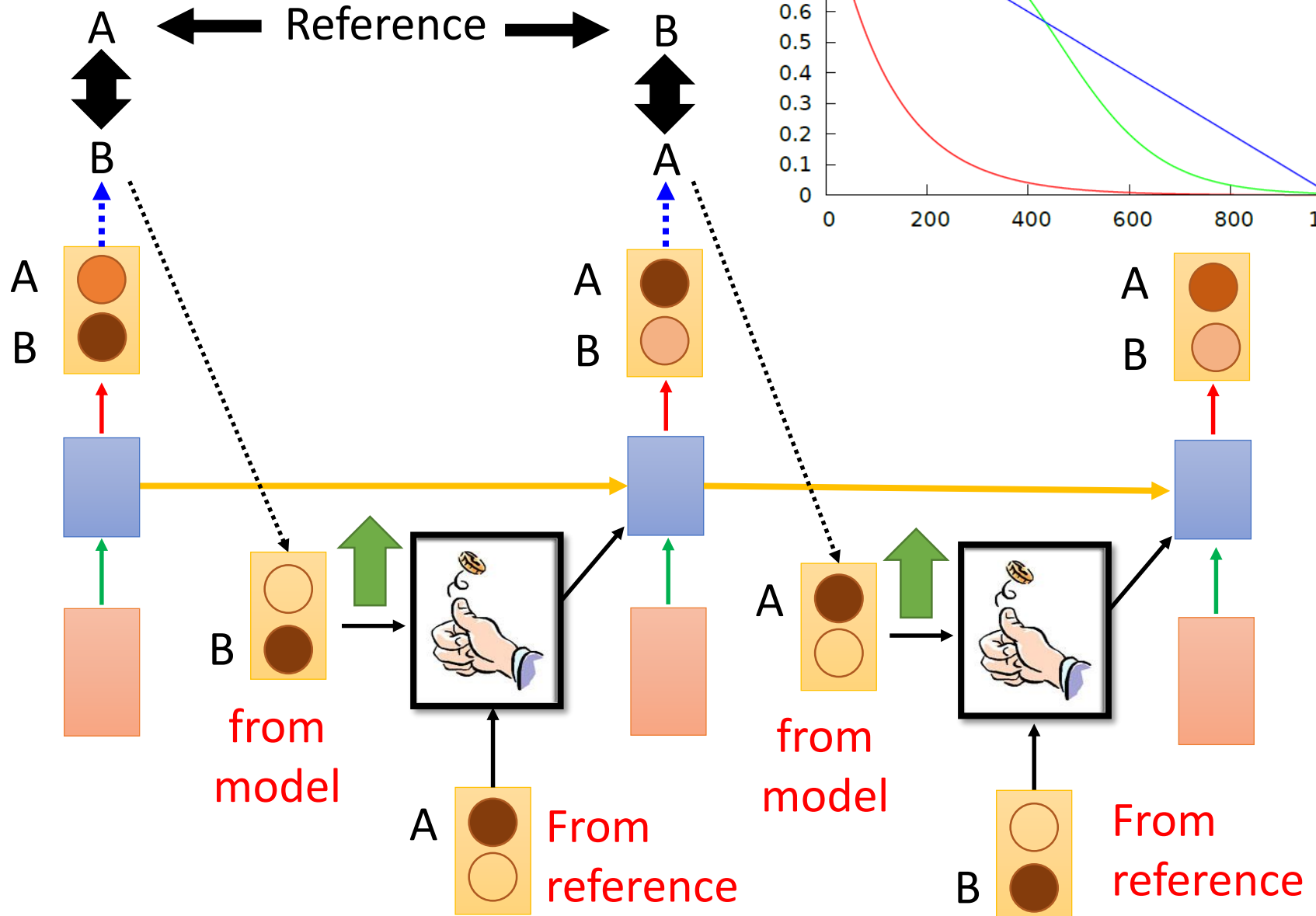
Testing: Output of model is the input of the next step.

Training: the inputs are reference.

Exposure Bias



Scheduled Sampling



Scheduled Sampling

- Caption generation on MSCOCO

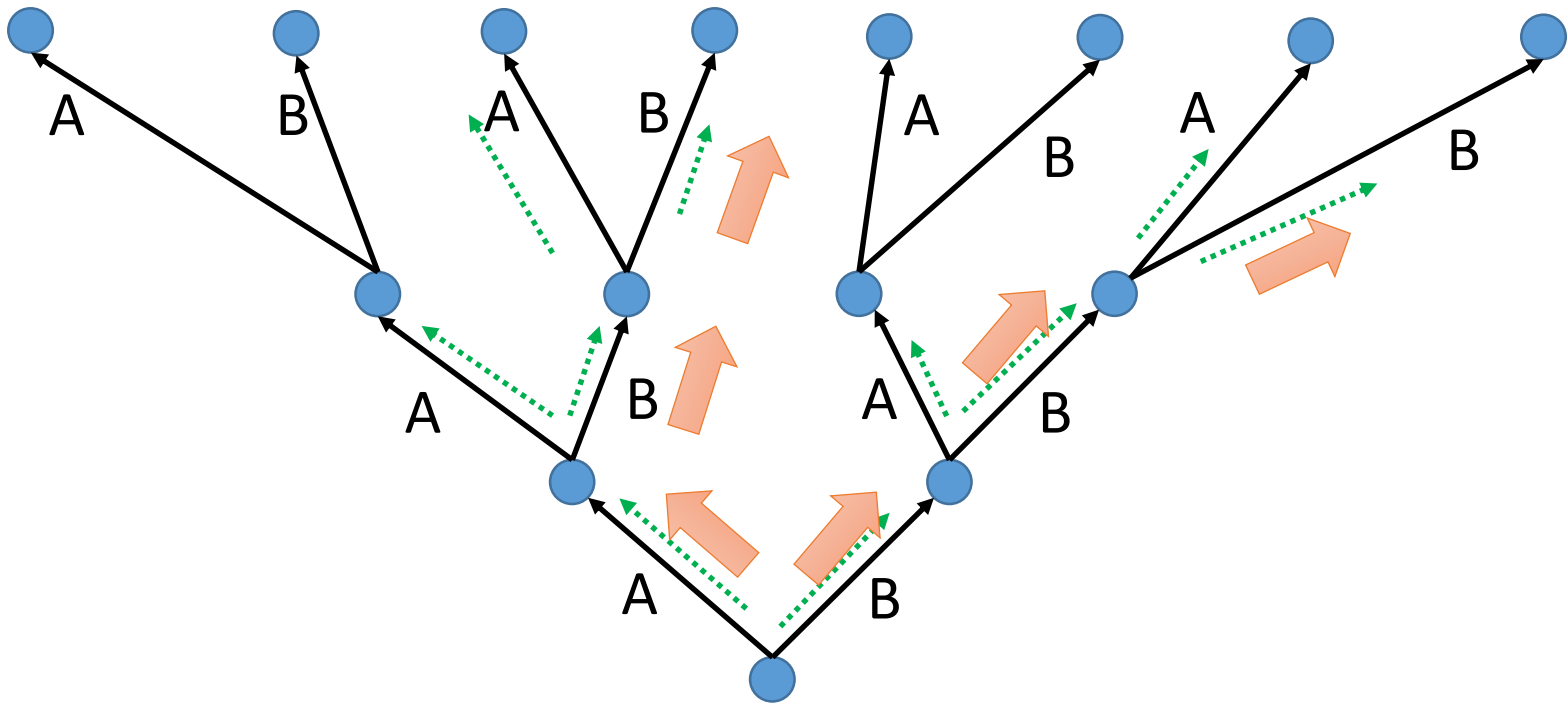
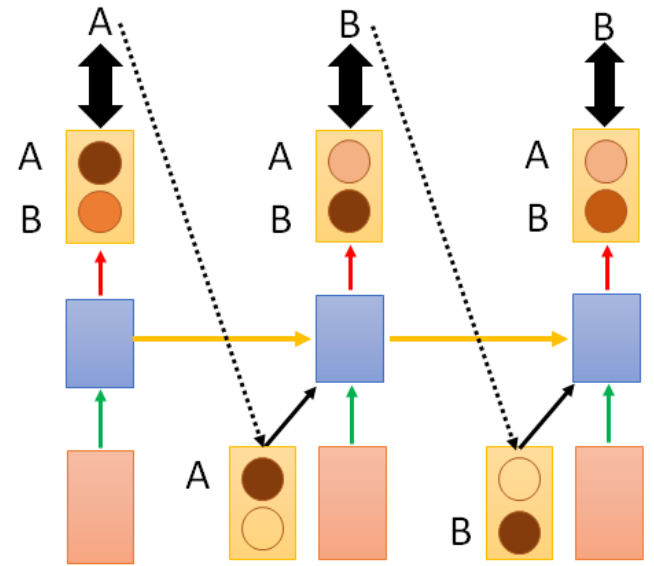
	BLEU-4	METEOR	CIDER
Always from reference	28.8	24.2	89.5
Always from model	11.2	15.7	49.7
Scheduled Sampling	30.6	24.3	92.1

Samy Bengio, Oriol Vinyals, Navdeep Jaitly, Noam Shazeer, Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks, arXiv preprint, 2015

Beam Search

Keep several best path at each step

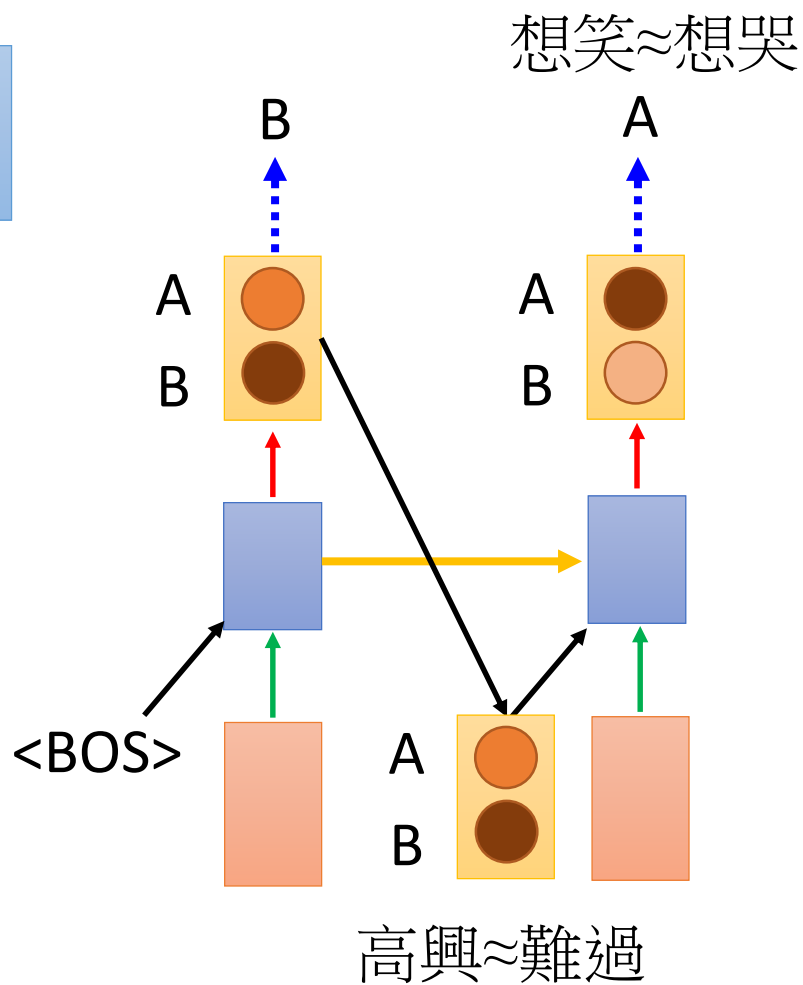
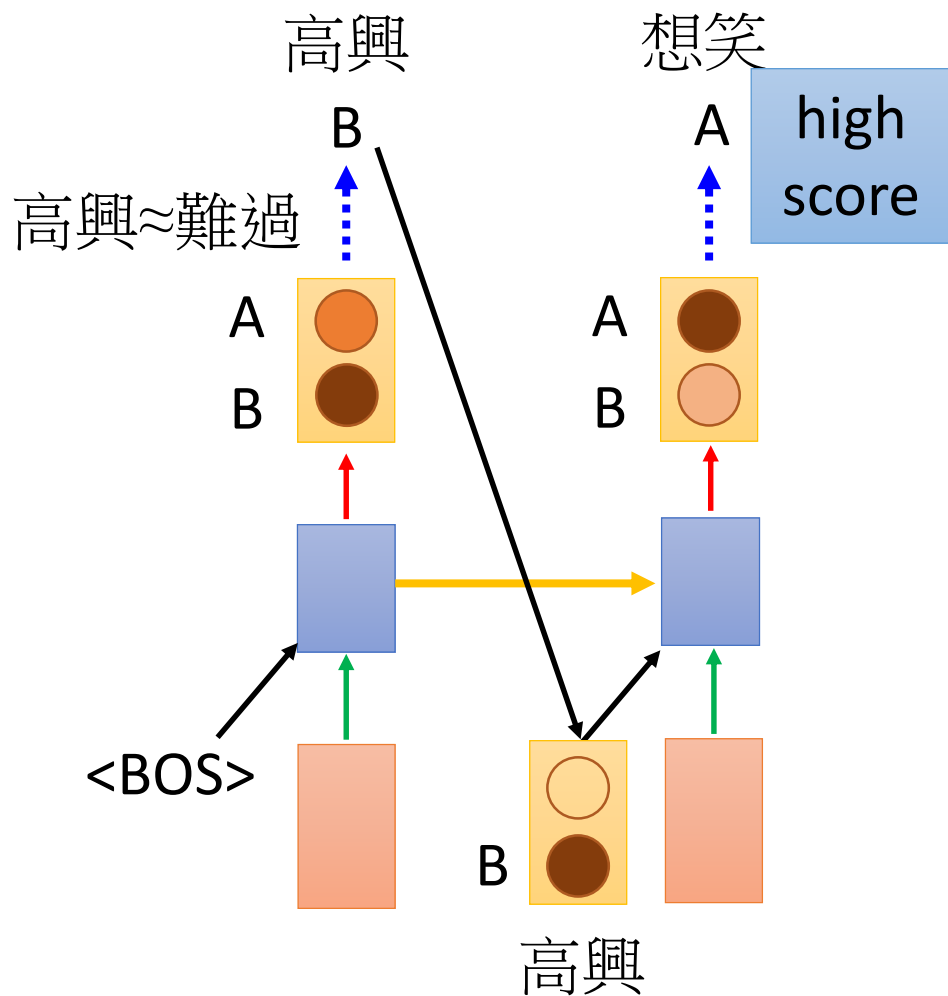
Beam size = 2



Better Idea?

U: 你覺得如何?

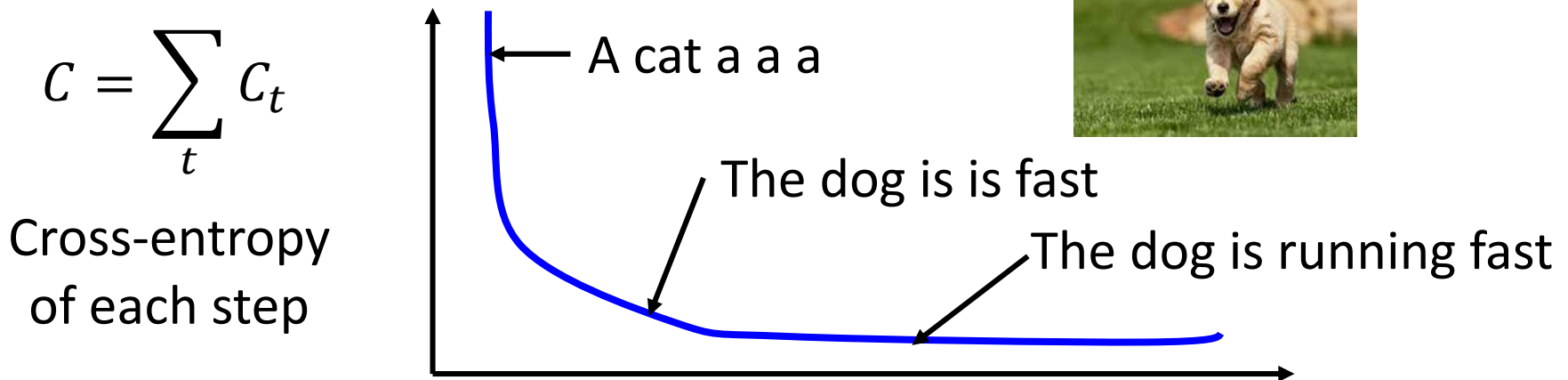
M: 高興想笑 or 難過想哭



Object level v.s. Component level

- Minimizing the error defined on component level is not equivalent to improving the generated objects

Ref: The dog is running fast



Optimize object-level criterion instead of component-level cross-entropy. object-level criterion: $R(y, \hat{y})$

Gradient Descent?

y : generated utterance, \hat{y} : ground truth

Reinforcement learning?

Start with
observation s_1

Observation s_2

Observation s_3



Obtain reward
 $r_1 = 0$

Action a_1 : "right"

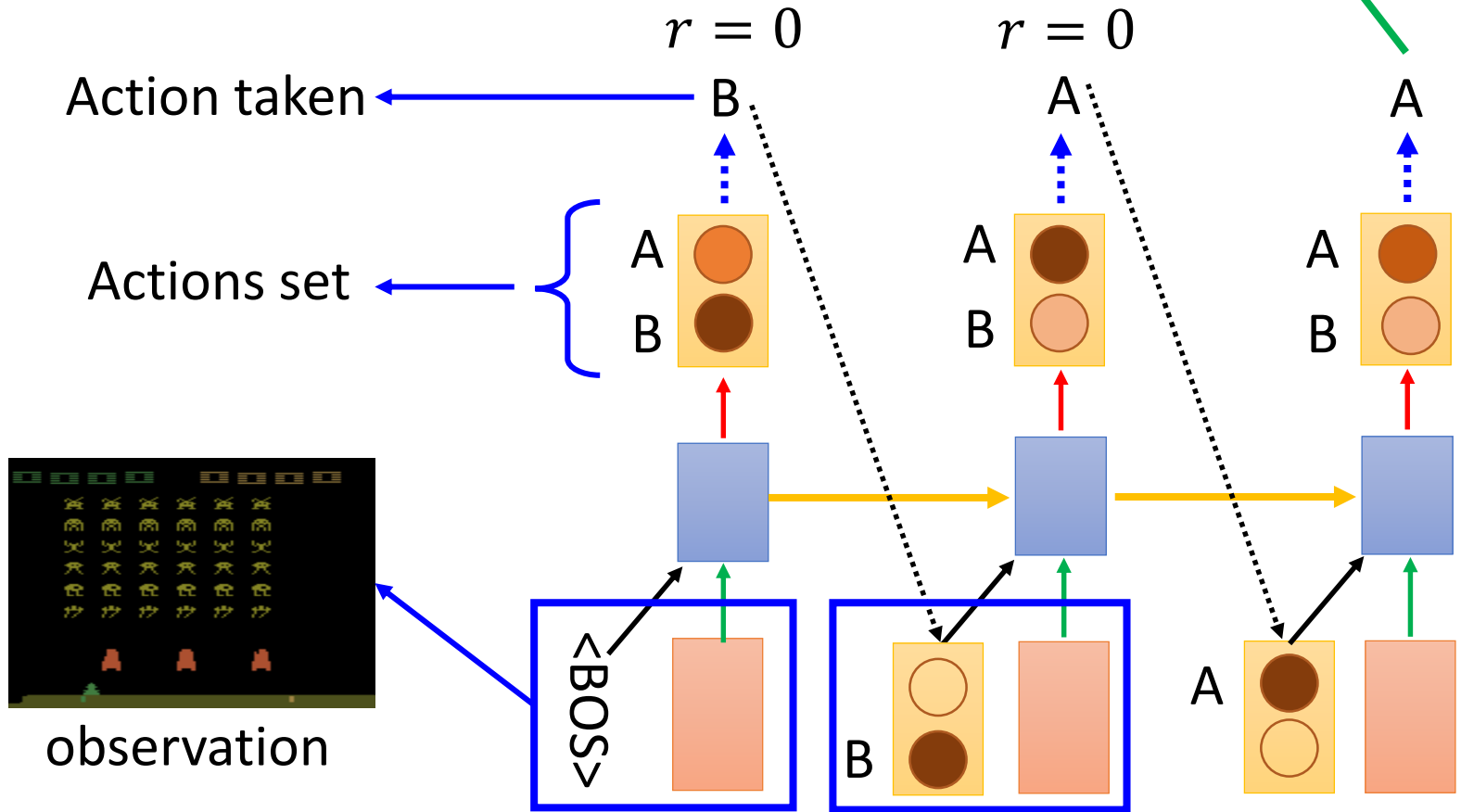


Obtain reward
 $r_2 = 5$

Action a_2 : "fire"
(kill an alien)

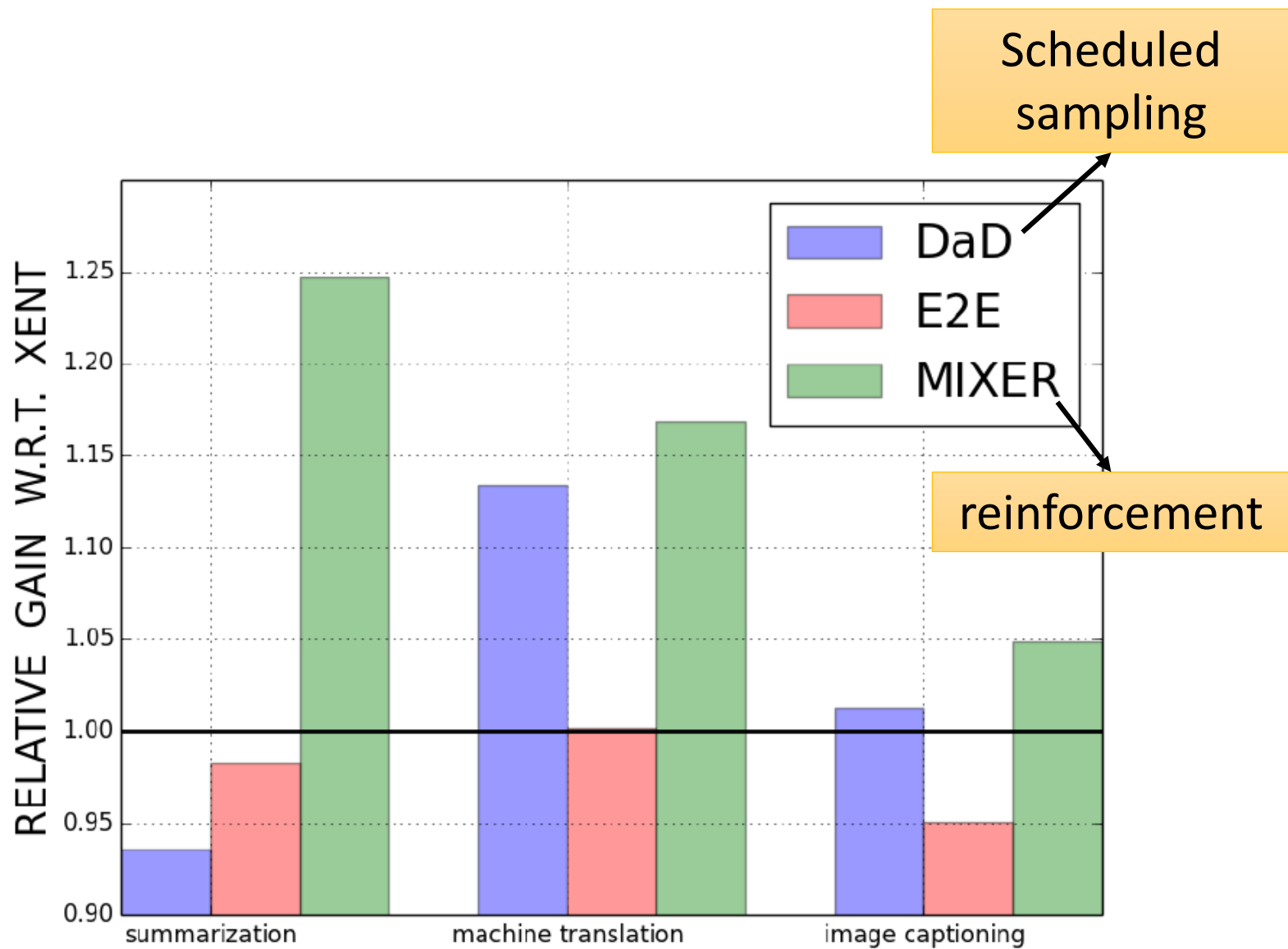
Reinforcement learning?

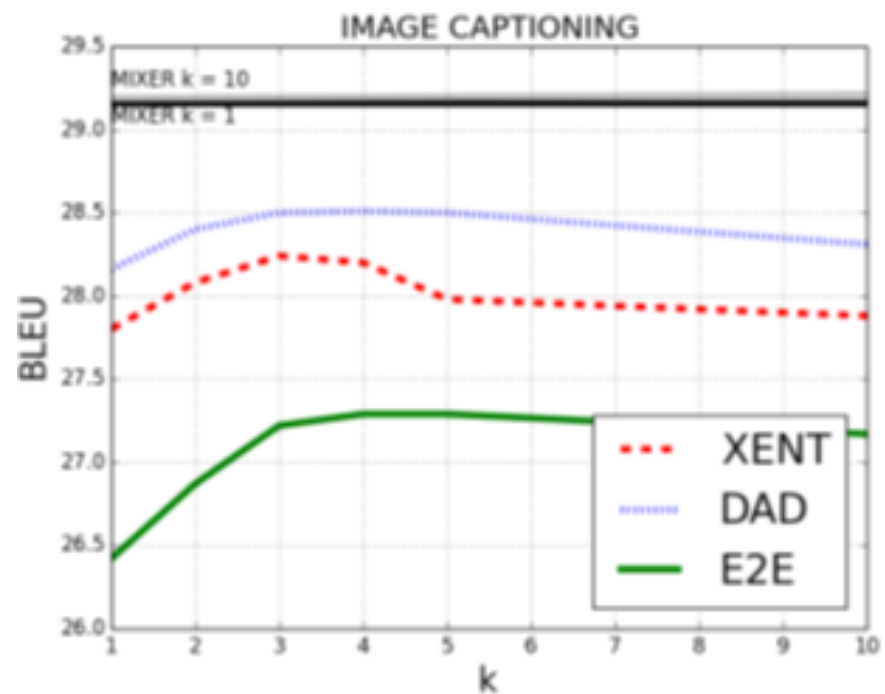
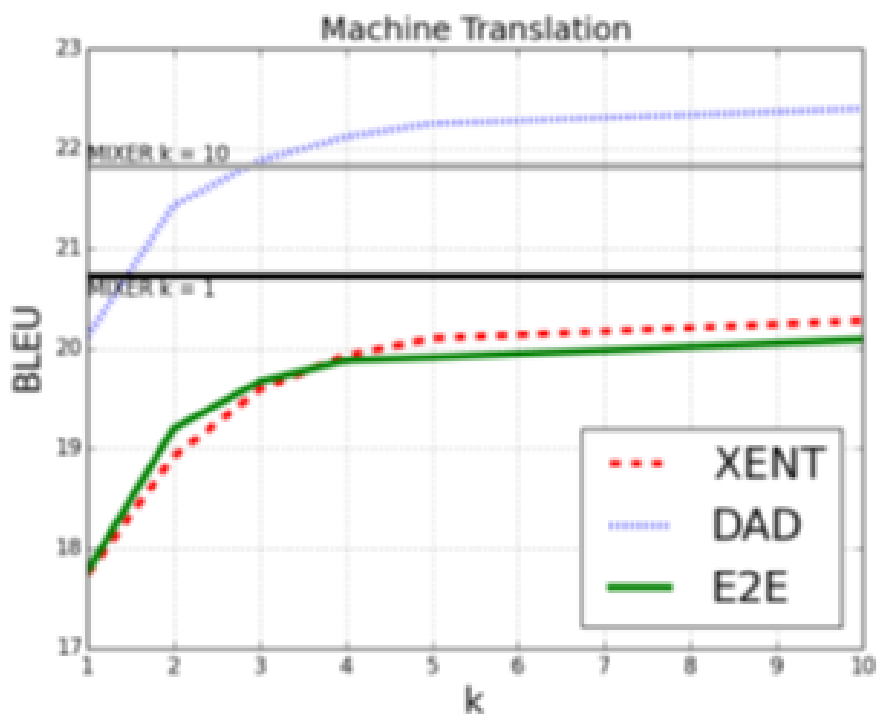
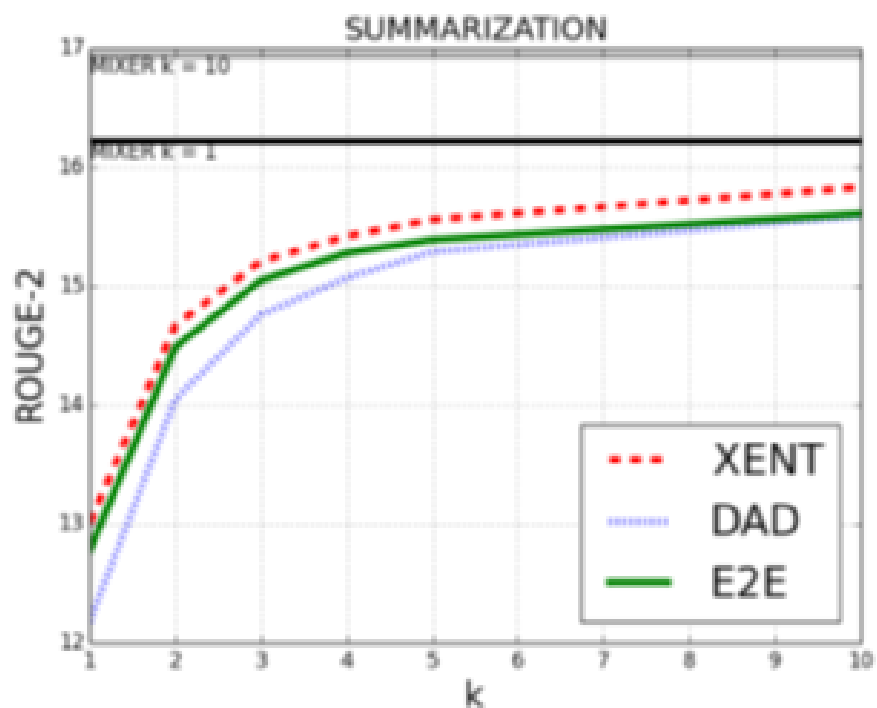
reward:
 $R(\text{"BAA"}, \text{reference})$



Marc'Aurelio Ranzato, Sumit Chopra, Michael Auli, Wojciech Zaremba, "Sequence Level Training with Recurrent Neural Networks", ICLR, 2016

The action we take influence the observation in the next step

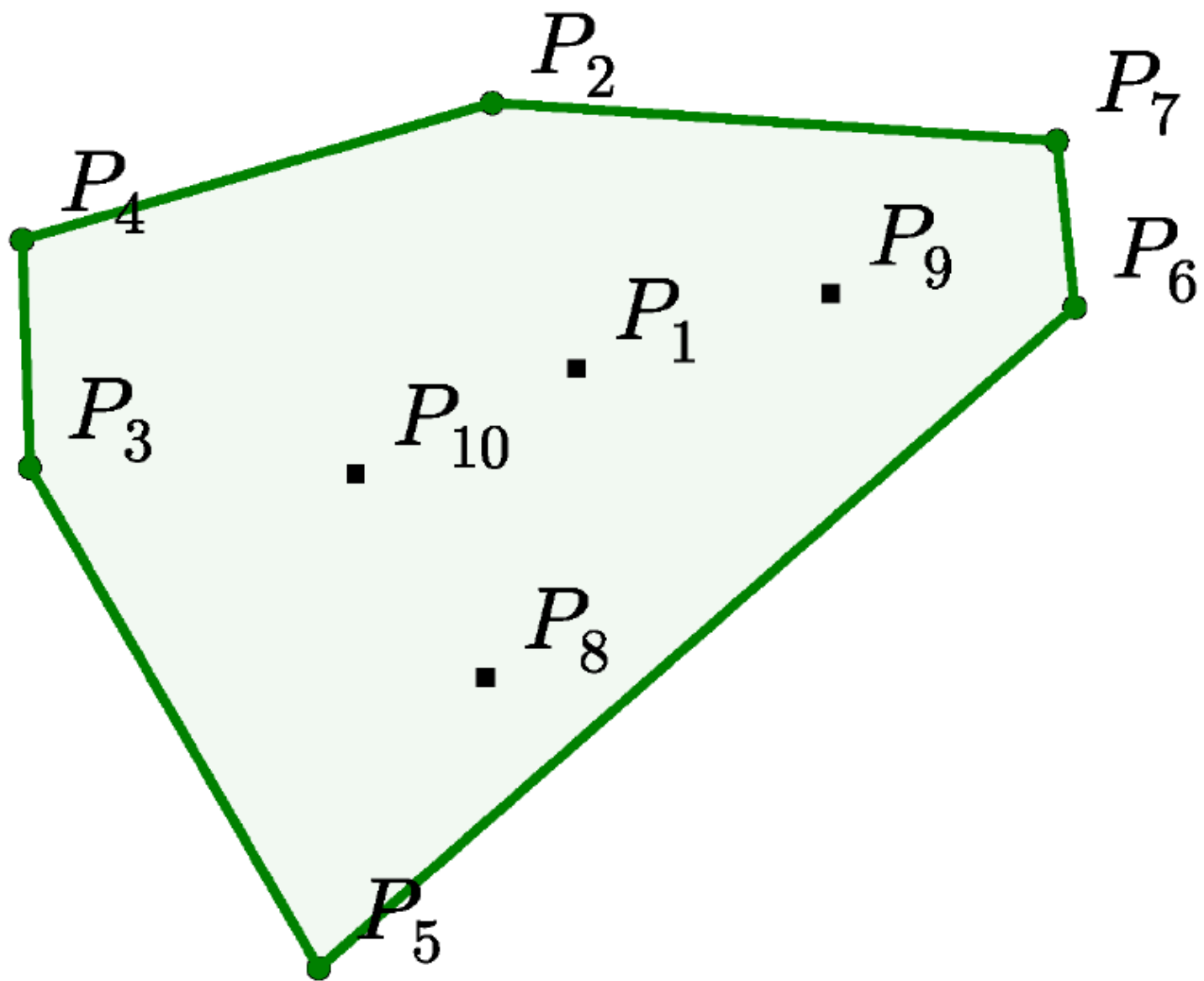


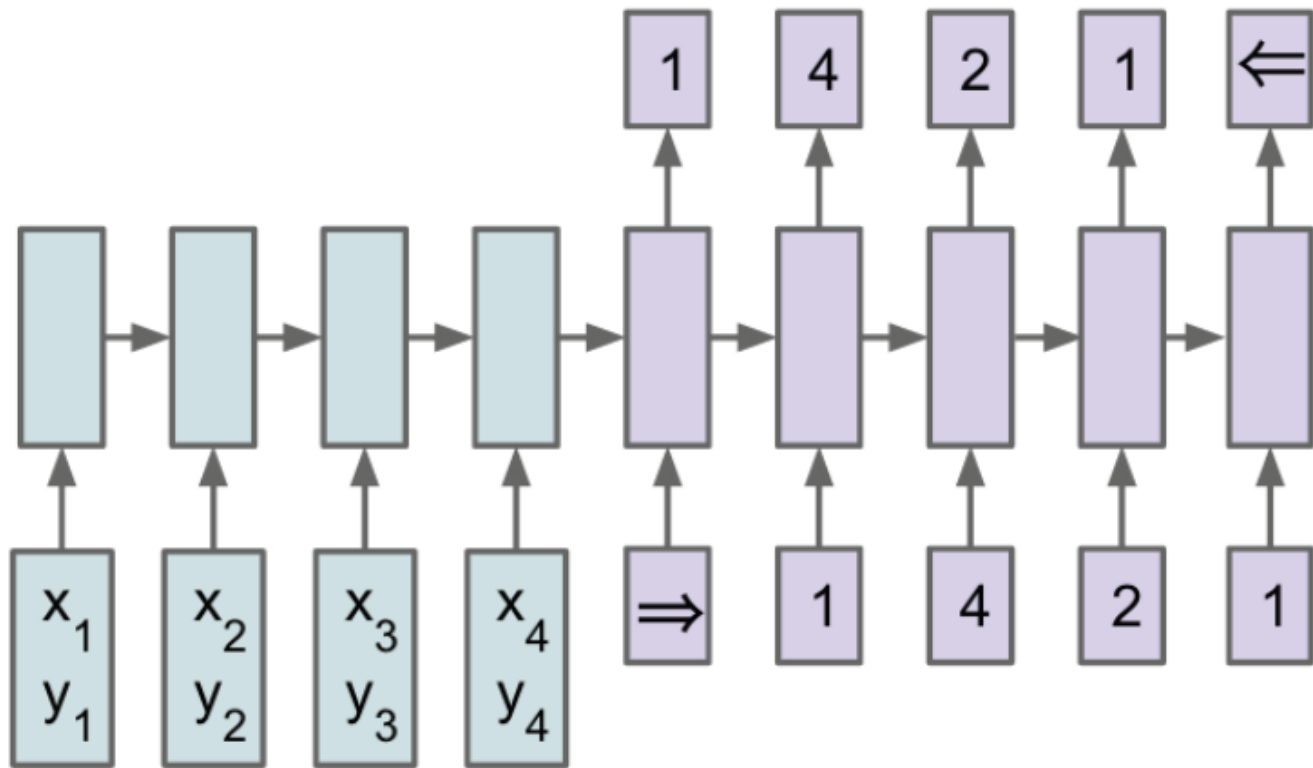
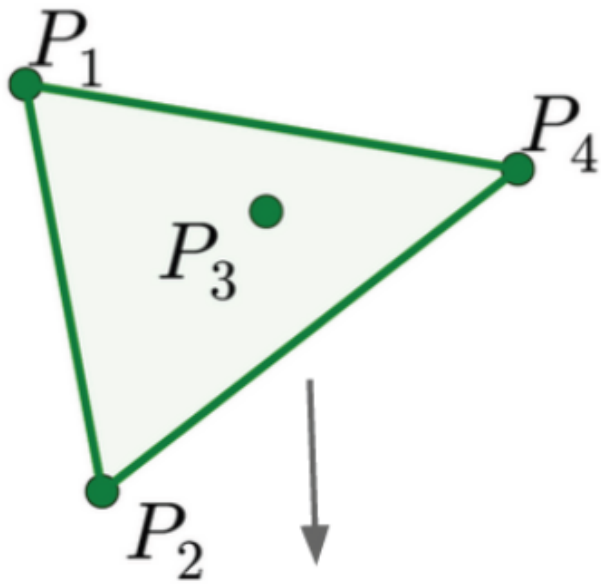


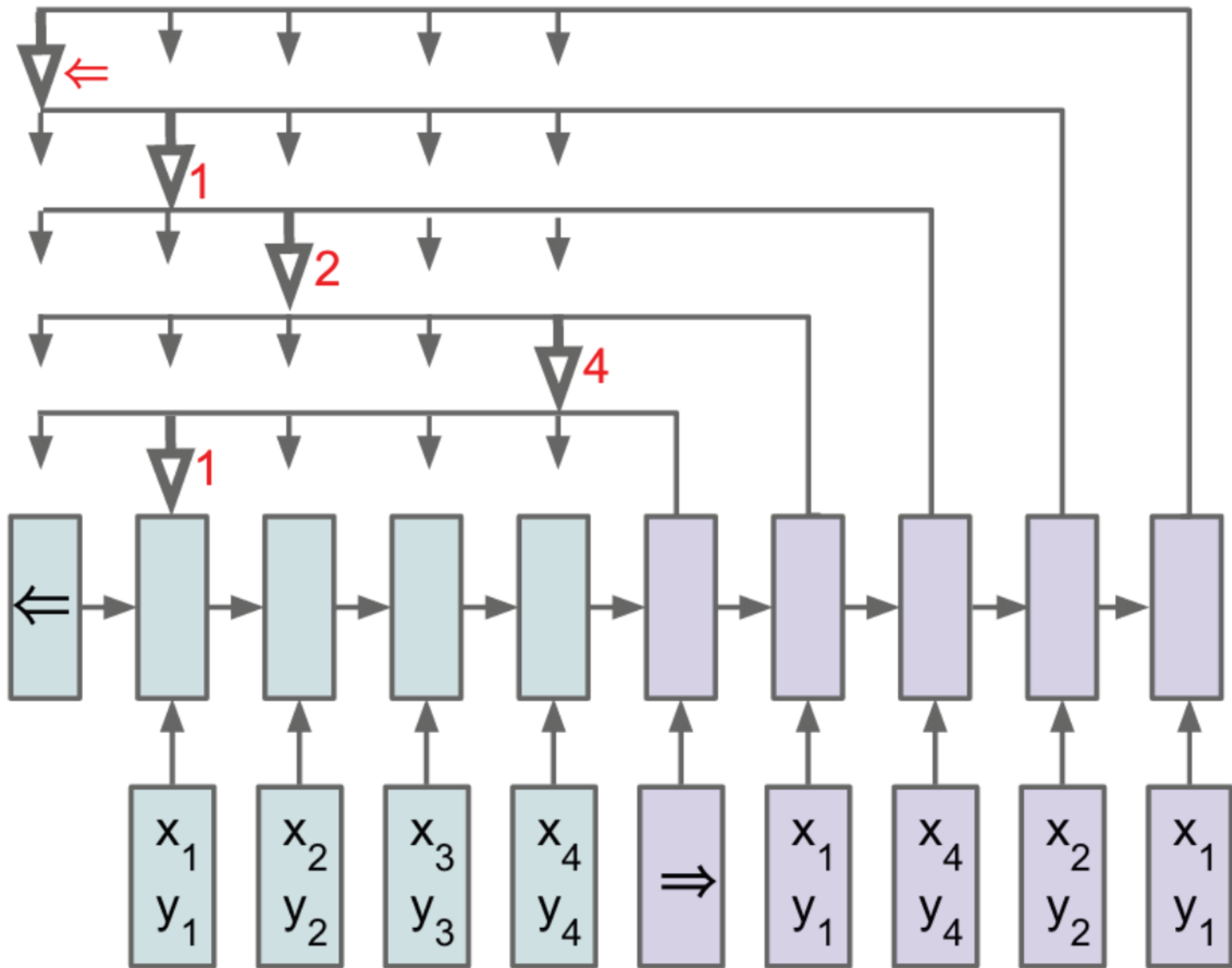
DAD: Scheduled Sampling
 MIXER: reinforcement

Pointer Network

Oriol Vinyals, Meire Fortunato, Navdeep Jaitly, Pointer Network, NIPS, 2015

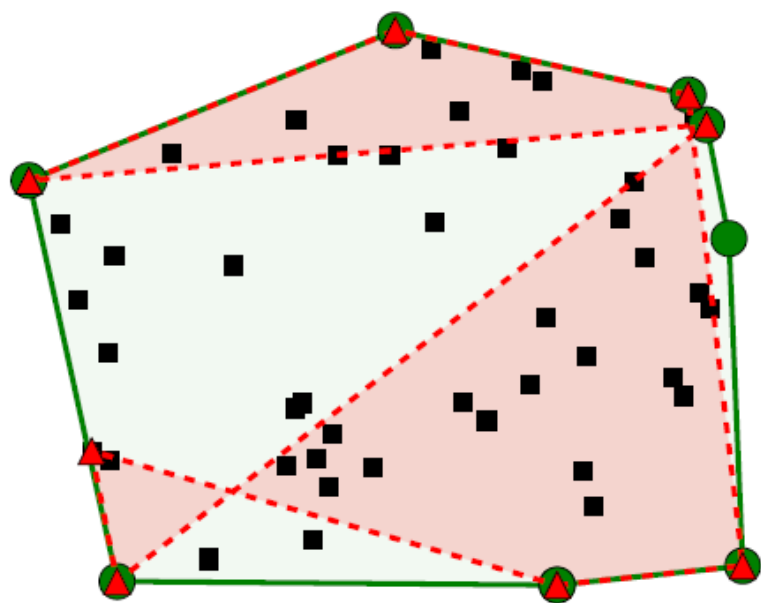






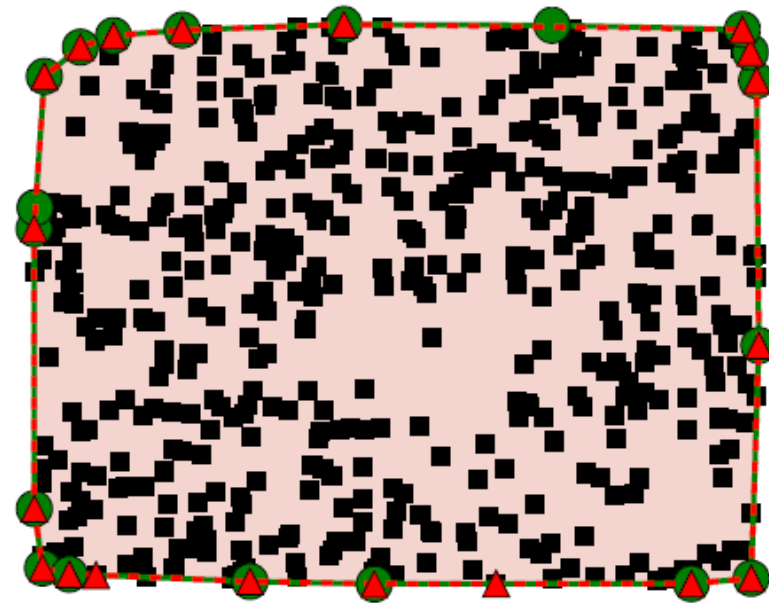
METHOD	TRAINED n	n	ACCURACY	AREA
LSTM [1]	50	50	1.9%	FAIL
+ATTENTION [5]	50	50	38.9%	99.7%
PTR-NET	50	50	72.6%	99.9%
LSTM [1]	5	5	87.7%	99.6%
PTR-NET	5-50	5	92.0%	99.6%
LSTM [1]	10	10	29.9%	FAIL
PTR-NET	5-50	10	87.0%	99.8%
PTR-NET	5-50	50	69.6%	99.9%
PTR-NET	5-50	100	50.3%	99.9%
PTR-NET	5-50	200	22.1%	99.9%
PTR-NET	5-50	500	1.3%	99.2%

● Ground Truth ▲ Predictions



(a) LSTM, $m=50$, $n=50$

● Ground Truth ▲ Predictions

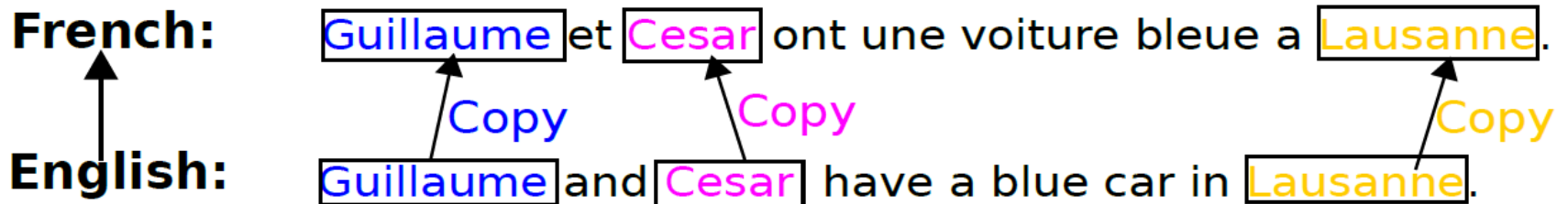


(d) Ptr-Net, $m=5-50$, $n=500$

Applications

Jiatao Gu, Zhengdong Lu, Hang Li, Victor O.K. Li,
“Incorporating Copying Mechanism in Sequence-to-Sequence Learning”, ACL, 2016
Caglar Gulcehre, Sungjin Ahn, Ramesh
Nallapati, Bowen Zhou, Yoshua Bengio, “Pointing the
Unknown Words”, ACL, 2016

Machine Translation



Chat-bot

User: X寶你好，我是宏毅

Machine: 宏毅你好，很高興認識你