

# Sequence Generation

Hung-yi Lee

李宏毅

# Outline

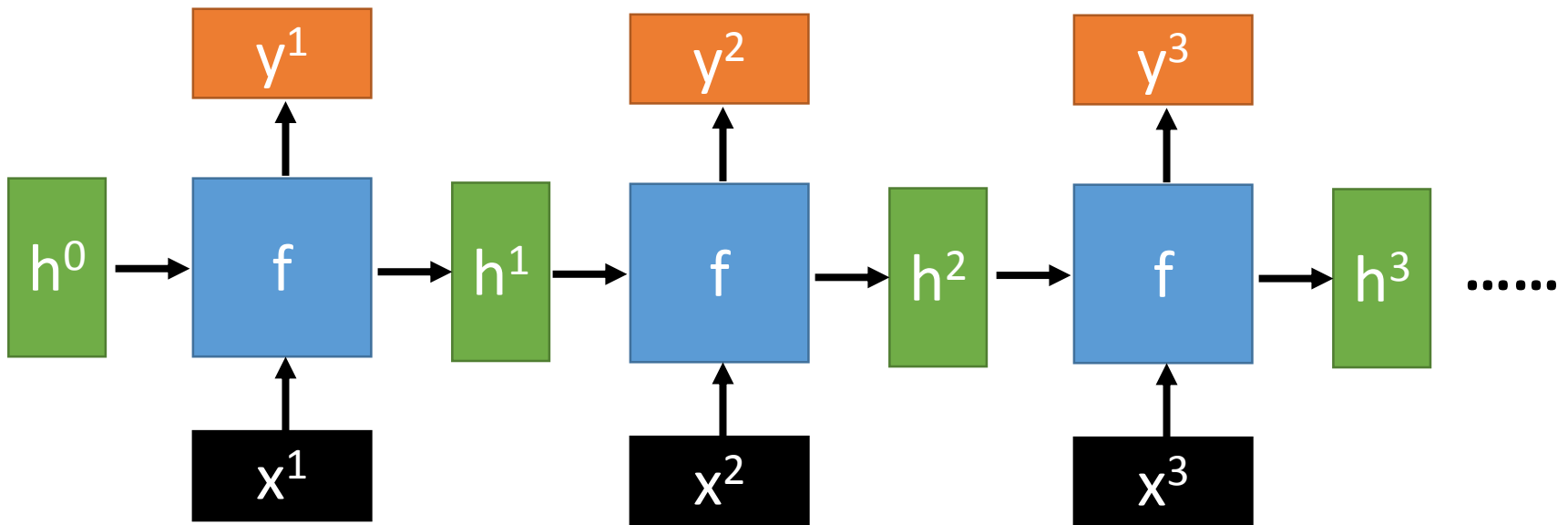
- RNN with Gated Mechanism
- Sequence Generation
- Conditional Sequence Generation
- Tips for Generation

# RNN with Gated Mechanism

# Recurrent Neural Network

- Given function  $f: h', y = f(h, x)$

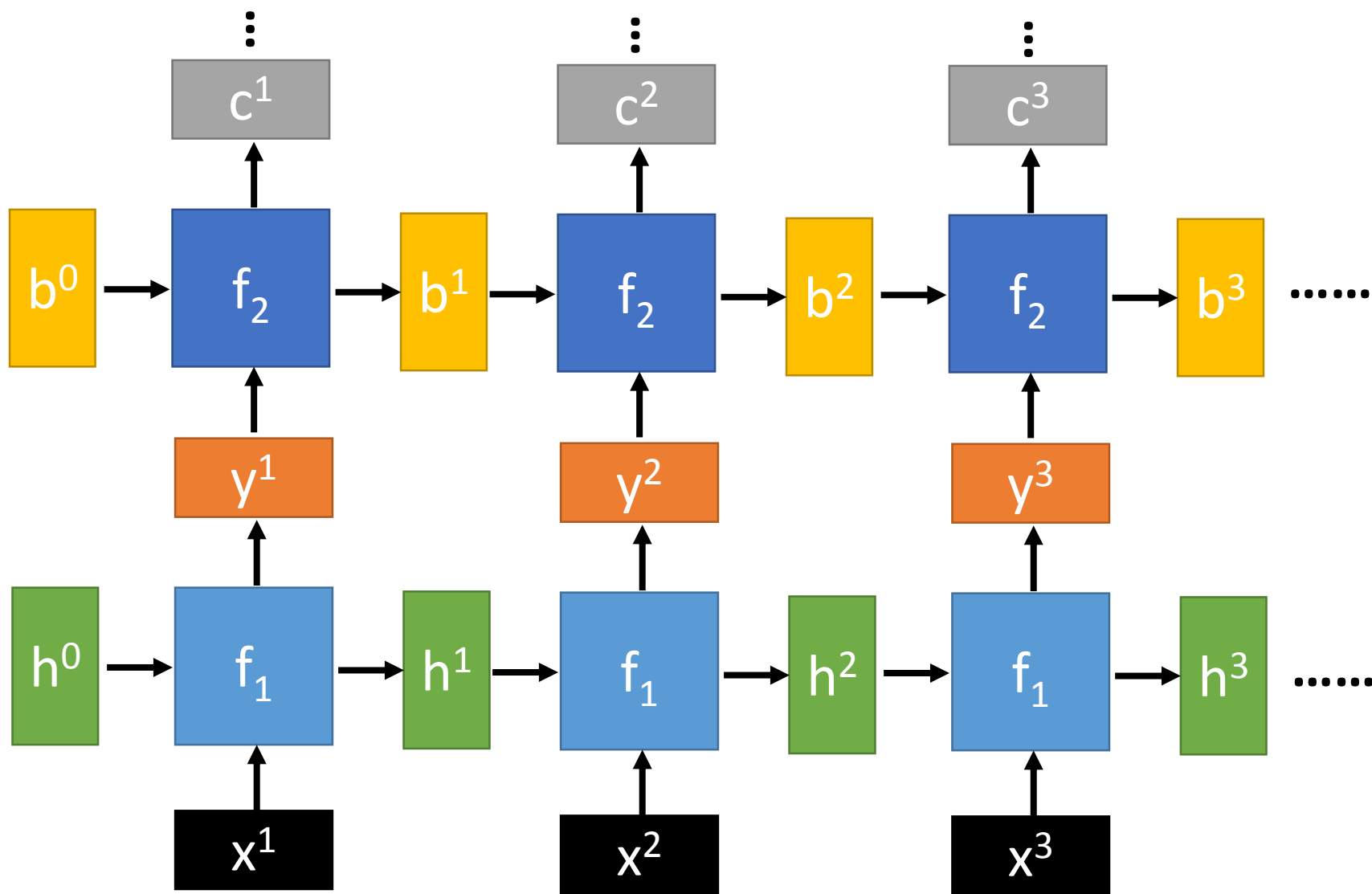
$h$  and  $h'$  are vectors with the same dimension



No matter how long the input/output sequence is, we only need one function  $f$

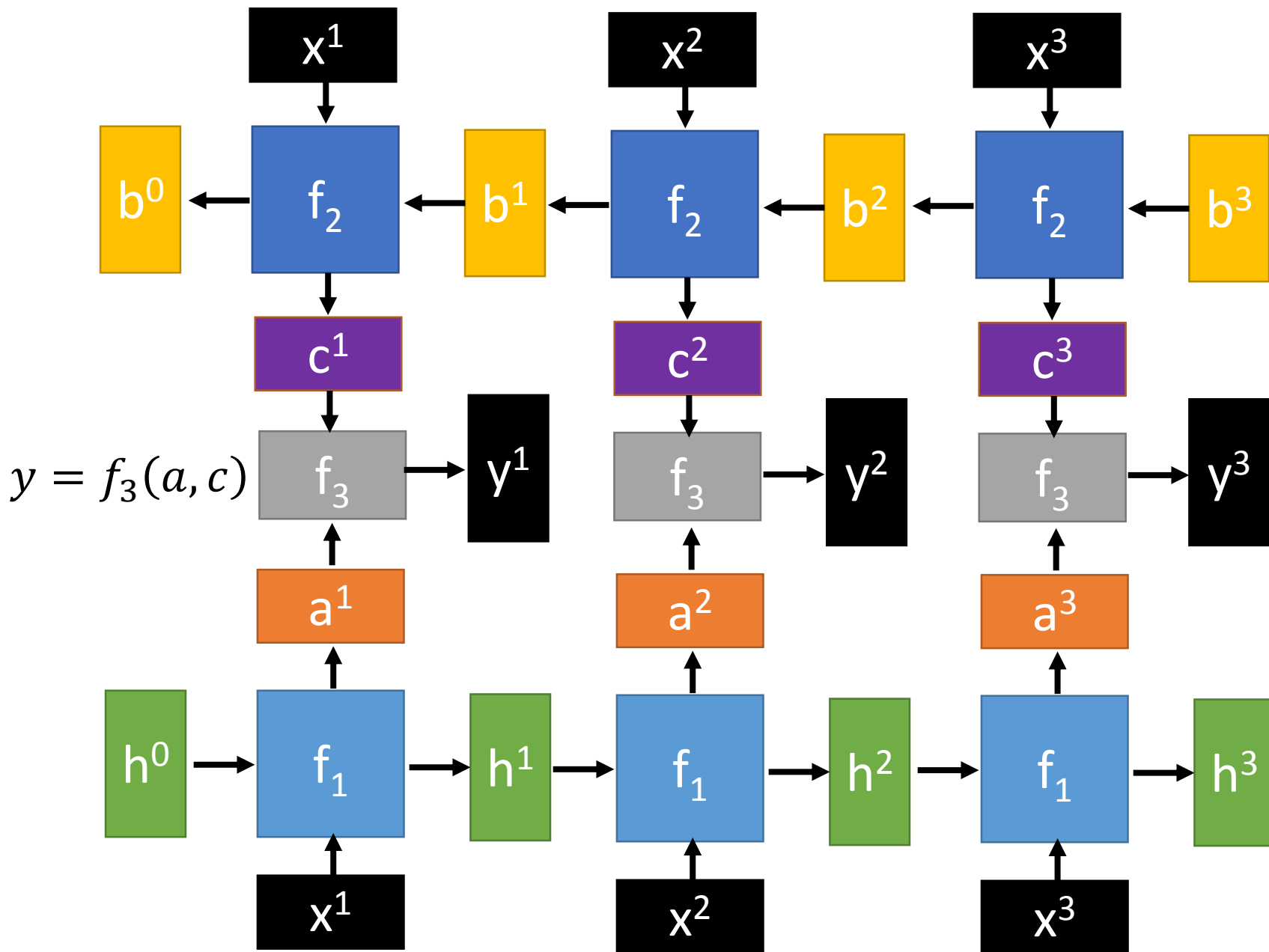
# Deep RNN

$$h', y = f_1(h, x) \quad b', c = f_2(b, y) \quad \dots$$



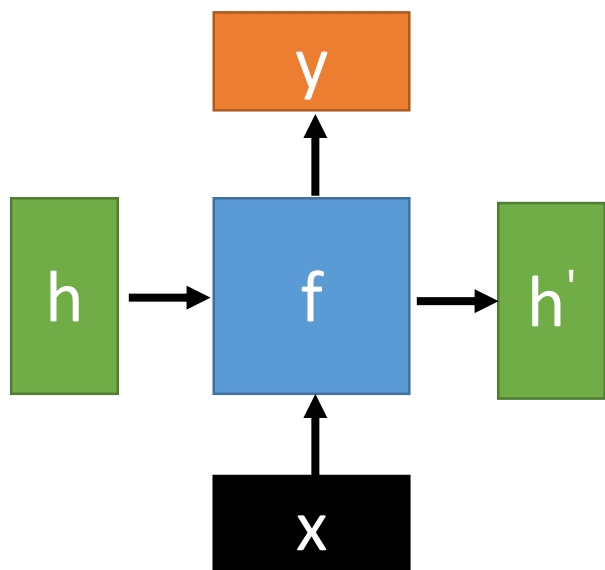
# Bidirectional RNN

$$h', a = f_1(h, x) \quad b', c = f_2(b, x)$$



# Naïve RNN

- Given function  $f: h', y = f(h, x)$



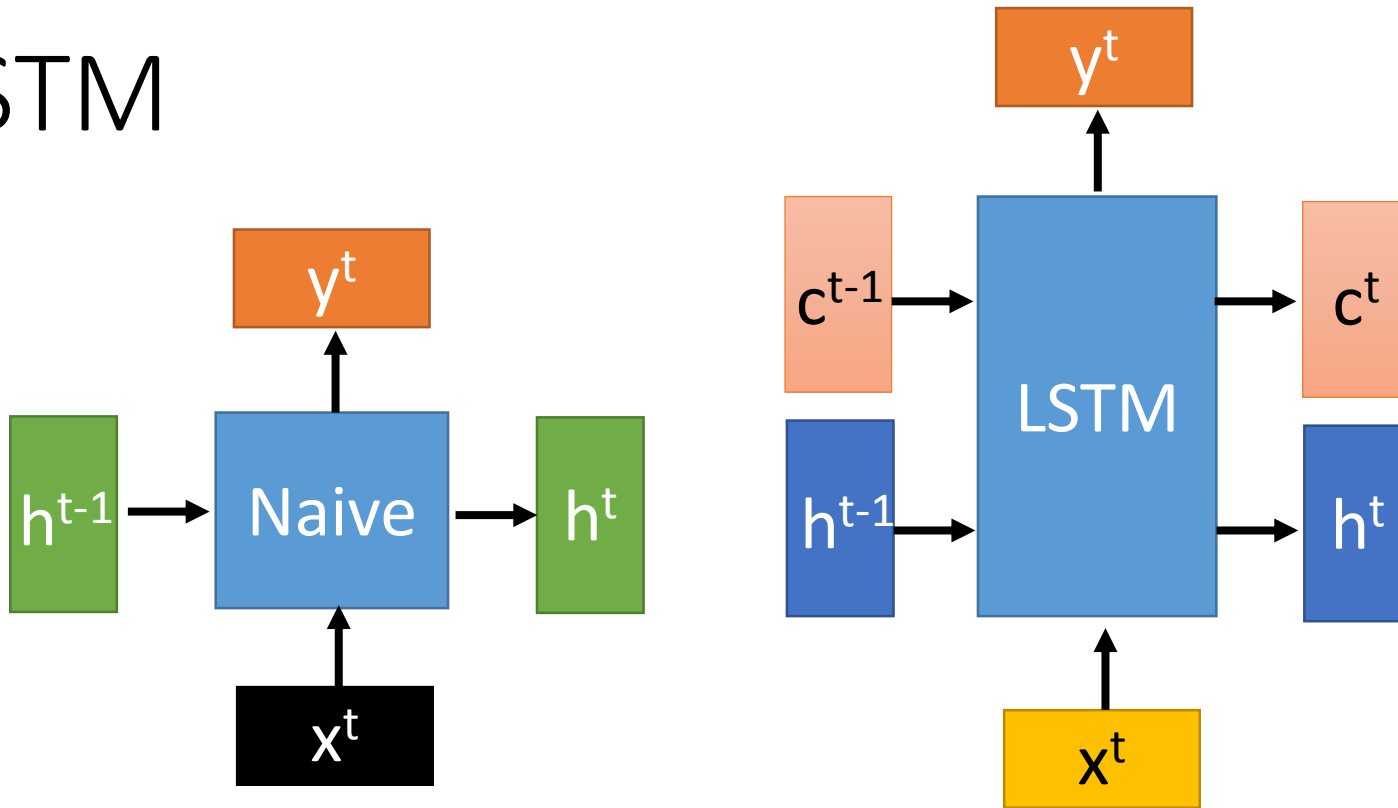
$$h' = \sigma( W^h h + W^i x )$$

$$y = \sigma( W^o h' )$$

softmax

Ignore bias here

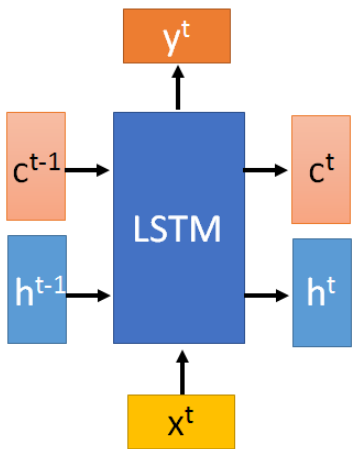
# LSTM



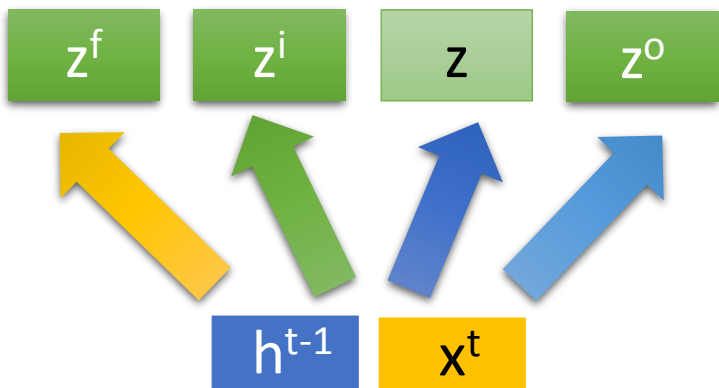
$c$  changes slowly  $\Rightarrow c^t$  is  $c^{t-1}$  added by something

$h$  changes faster  $\Rightarrow h^t$  and  $h^{t-1}$  can be very different





$c^{t-1}$



$$z = \tanh\left( W \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix} \right)$$

$$z^i = \sigma\left( W^i \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix} \right)$$

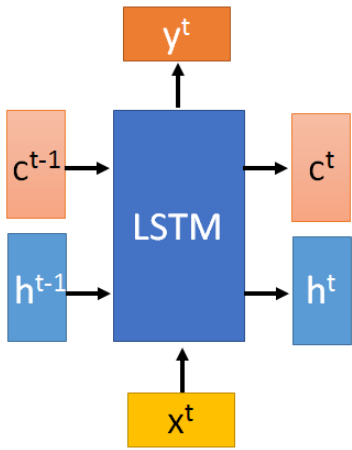
Input gate

$$z^f = \sigma\left( W^f \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix} \right)$$

forget gate

$$z^o = \sigma\left( W^o \begin{bmatrix} x^t \\ h^{t-1} \end{bmatrix} \right)$$

output gate

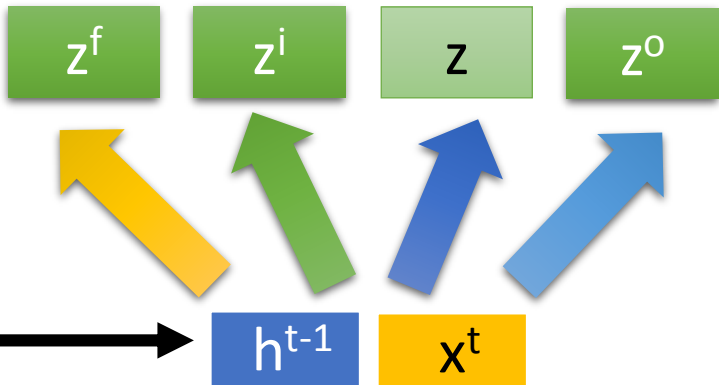


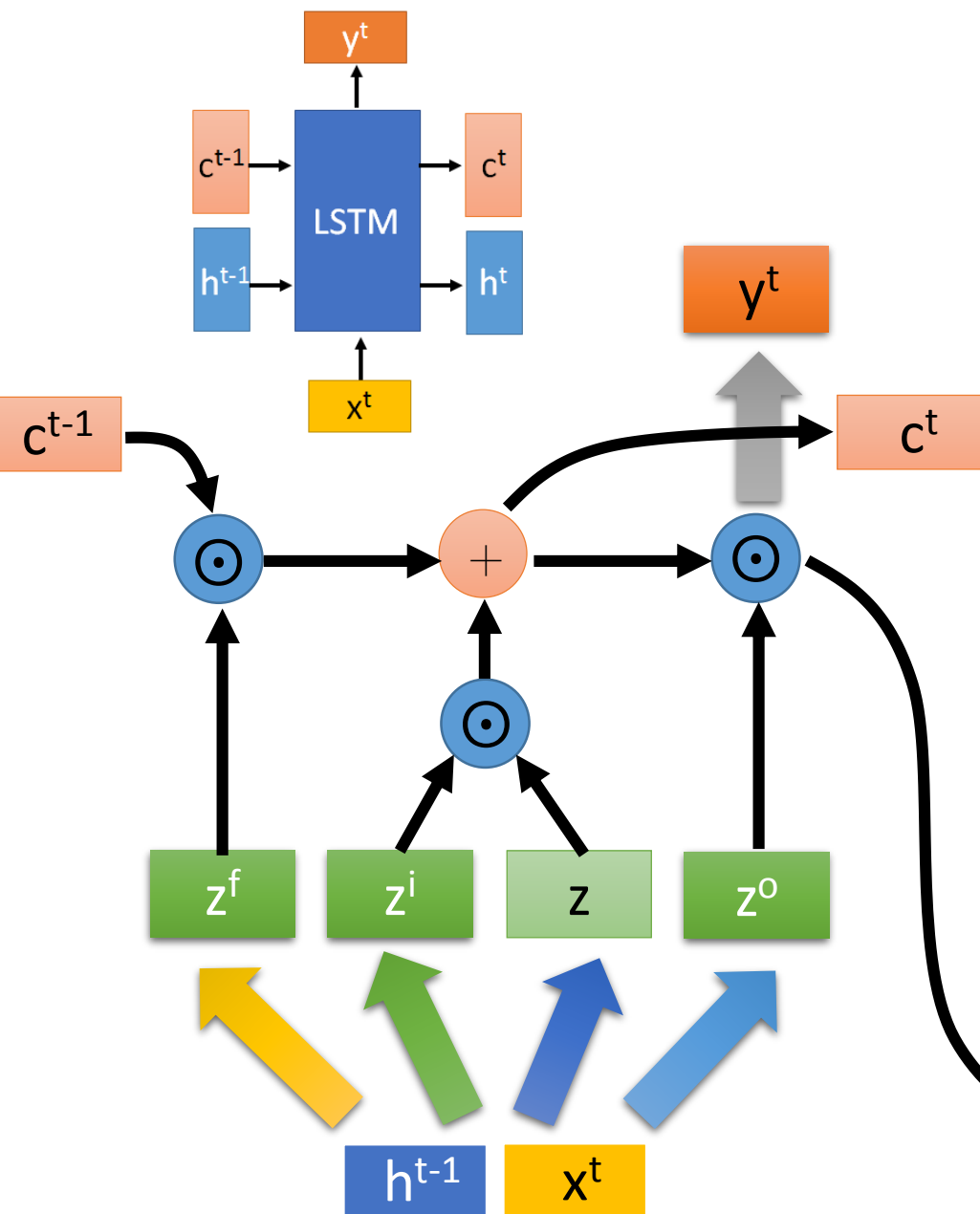
$$z = \tanh\left( \begin{matrix} \text{blue box } W & \text{orange box} \\ \text{orange box } c^{t-1} & \text{blue box } h^{t-1} \end{matrix} \right)$$

↑  
diagonal

$z^o$   $z^f$   $z^i$  obtained by the same way

“peephole”



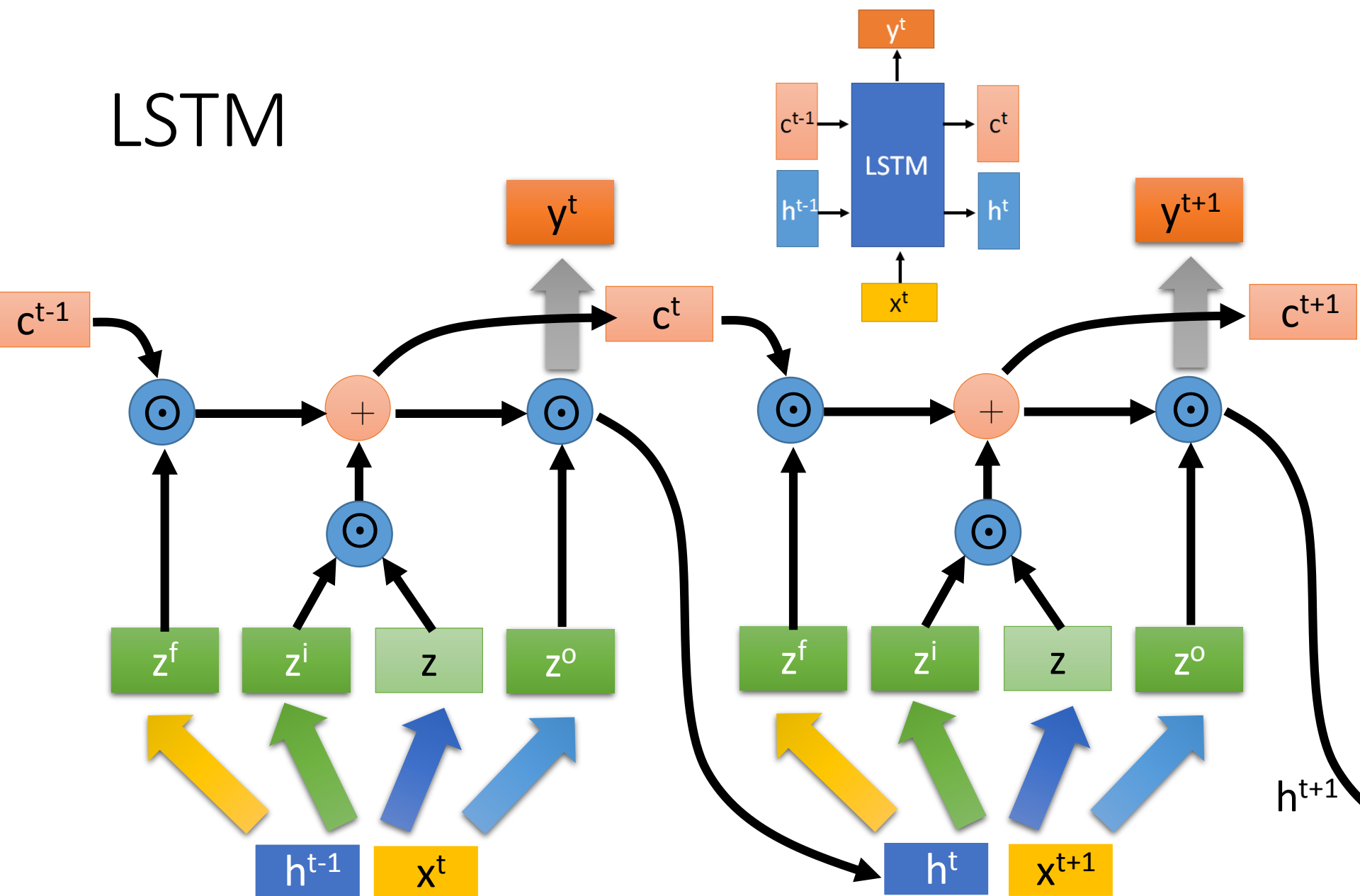


$$c^t = z^f \odot c^{t-1} + z^i \odot z$$

$$h^t = z^o \odot \tanh(c^t)$$

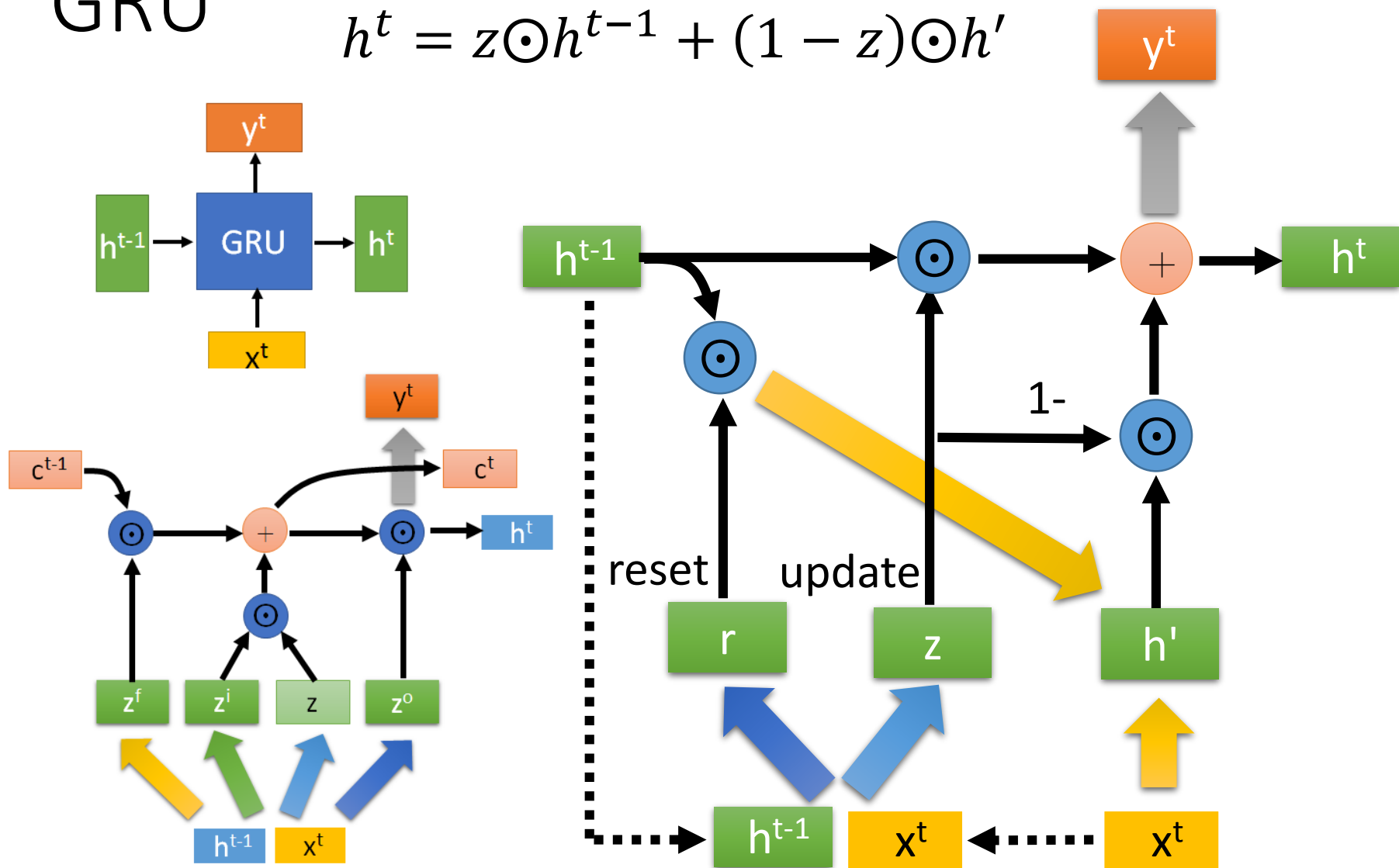
$$y^t = \sigma(W' h^t)$$

# LSTM

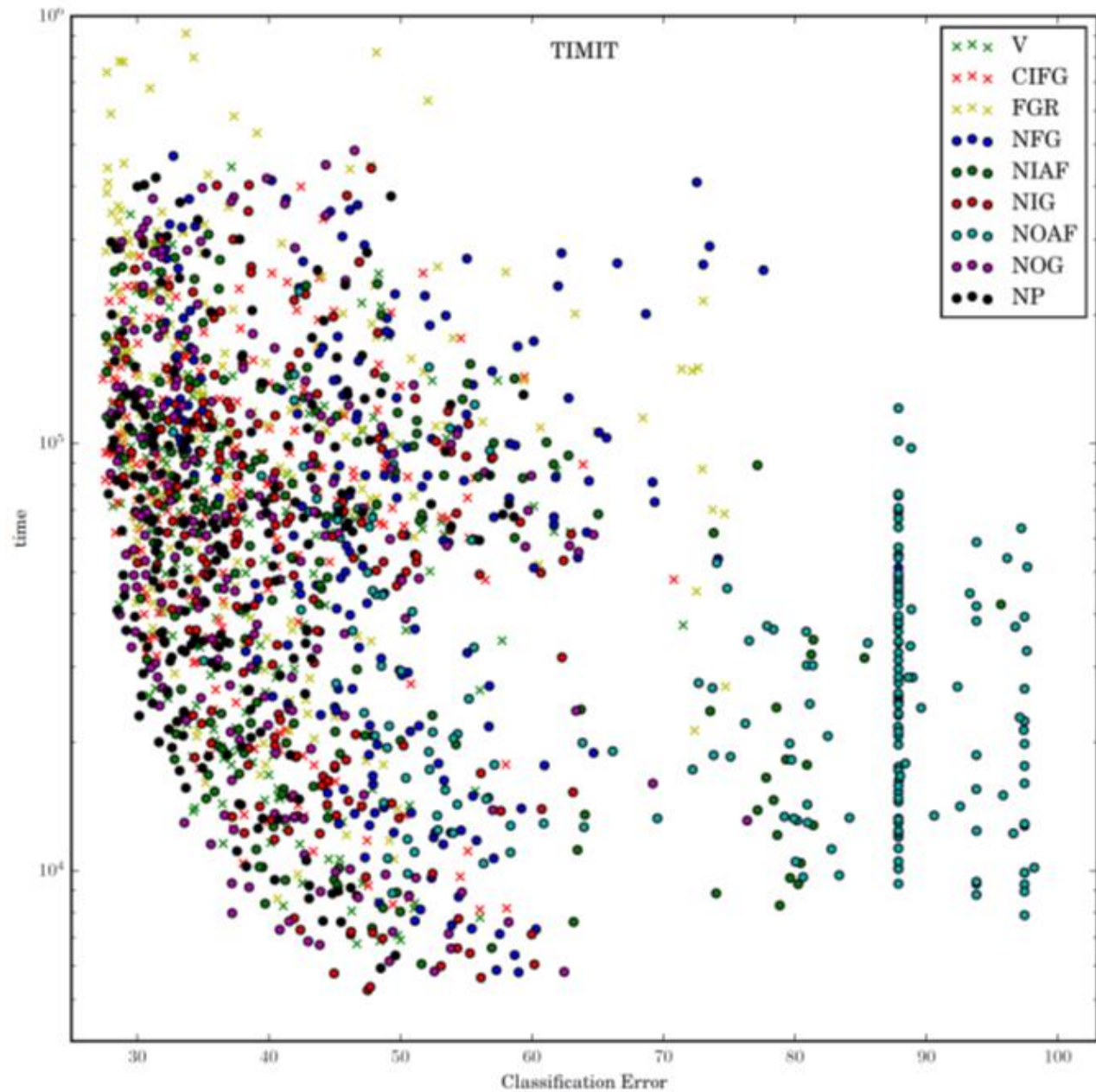


# GRU

$$h^t = z \odot h^{t-1} + (1 - z) \odot h'$$

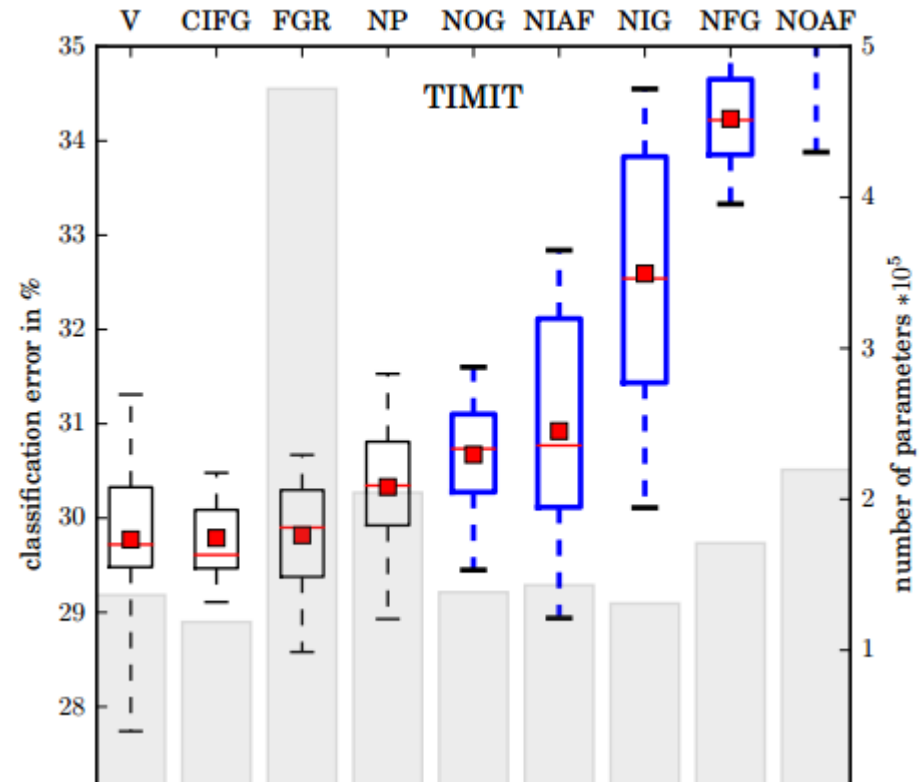


# LSTM: A Search Space Odyssey



# LSTM: A Search Space Odyssey

1. No Input Gate (NIG)
2. No Forget Gate (NFG)
3. No Output Gate (NOG)
4. No Input Activation Function (NIAF)
5. No Output Activation Function (NOAF)
6. No Peepholes (NP)
7. Coupled Input and Forget Gate (CIFG)
8. Full Gate Recurrence (FGR)



Standard LSTM works well

Simply LSTM: coupling input and forget gate, removing peephole

Forget gate is critical for performance

Output gate activation function is critical

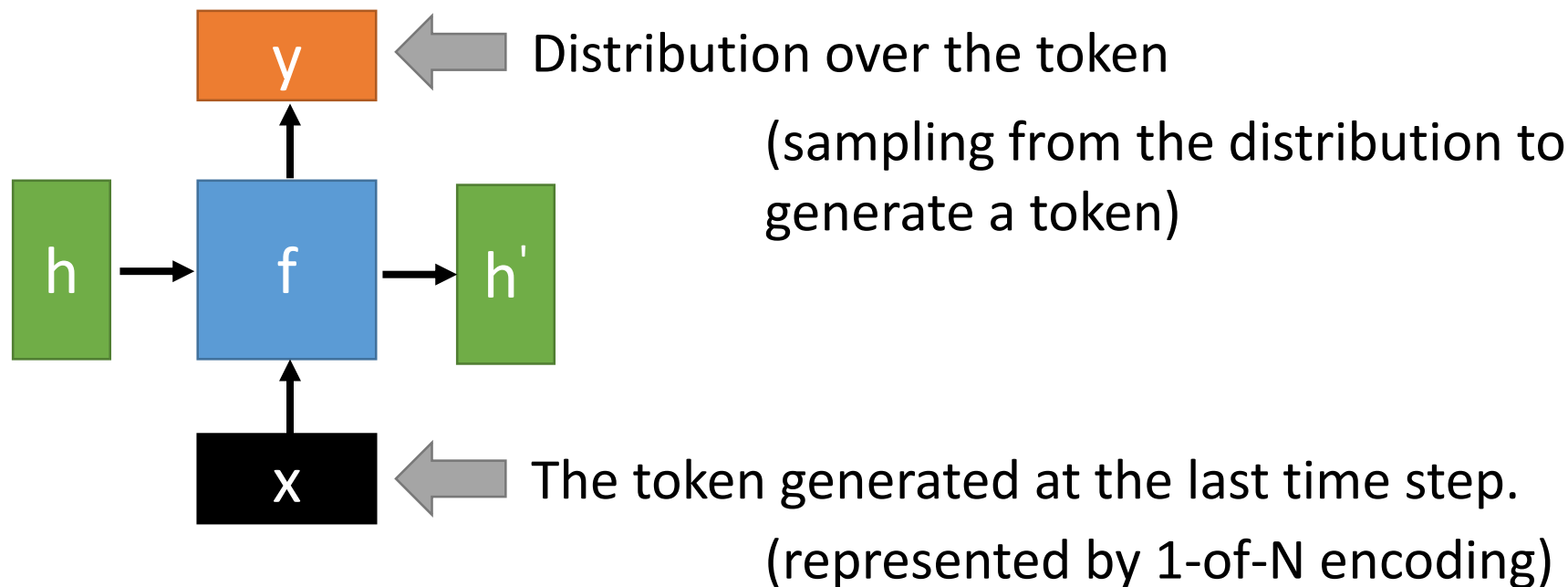
# Sequence Generation



# Generation



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



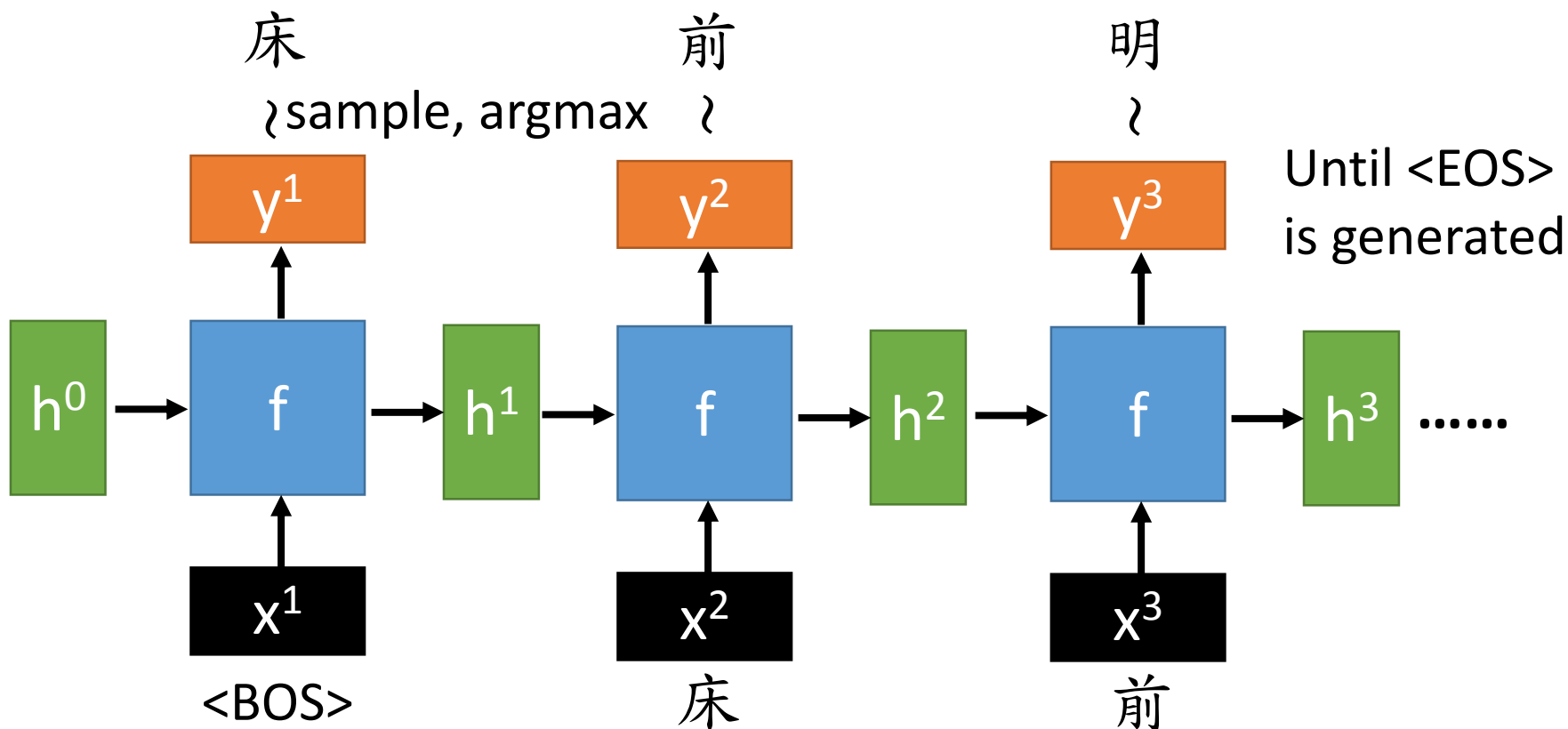
# Generation

$$y^1: P(w | \langle \text{BOS} \rangle)$$

$$y^2: P(w | \langle \text{BOS} \rangle, \text{床})$$

$$y^3: P(w | \langle \text{BOS} \rangle, \text{床}, \text{前})$$

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

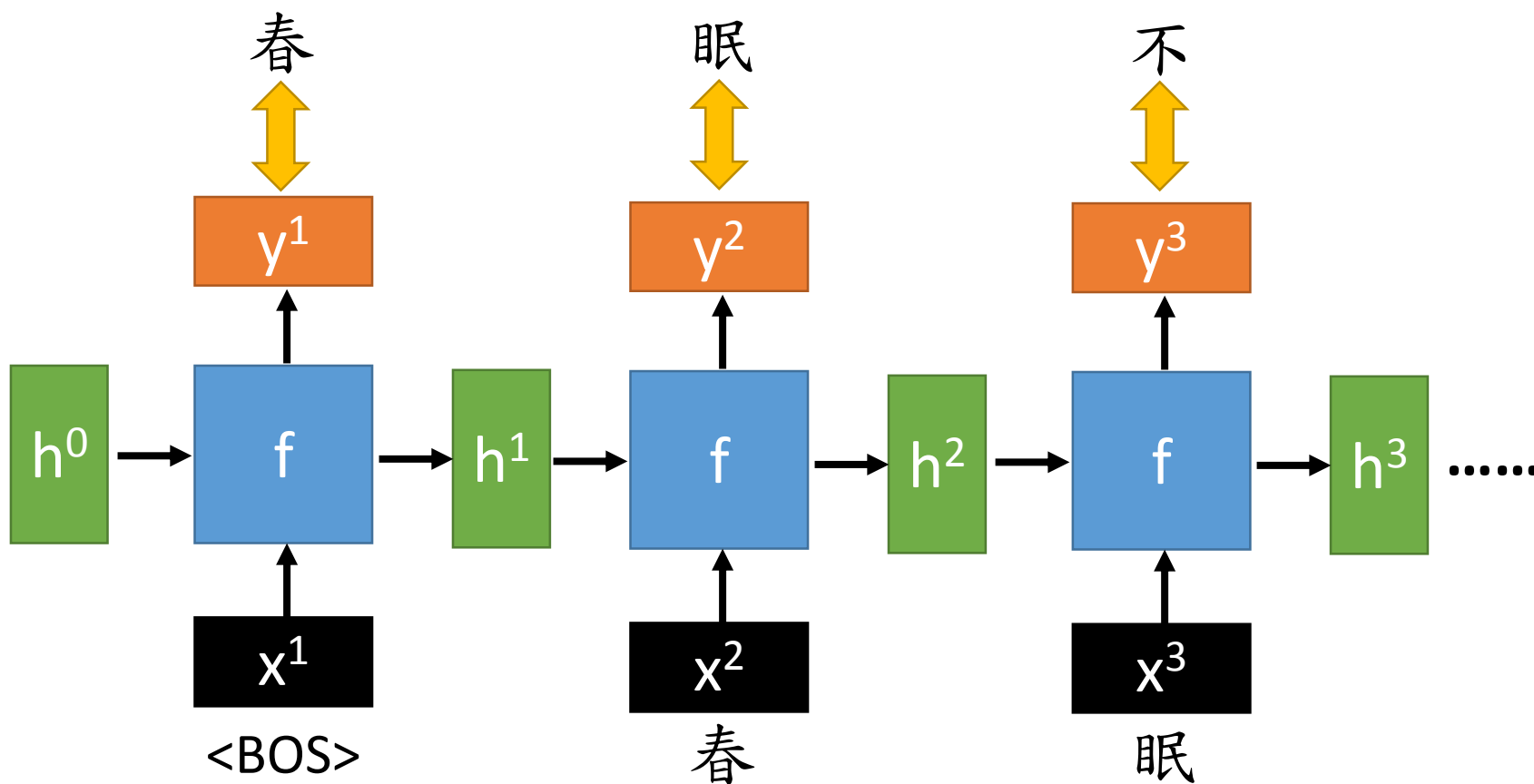


# Generation

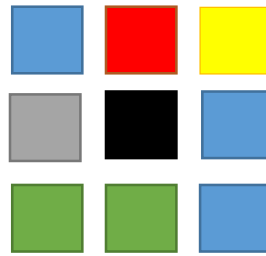
↕ : minimizing cross-entropy

- Training

Training data: 春 眠 不 覺 曉



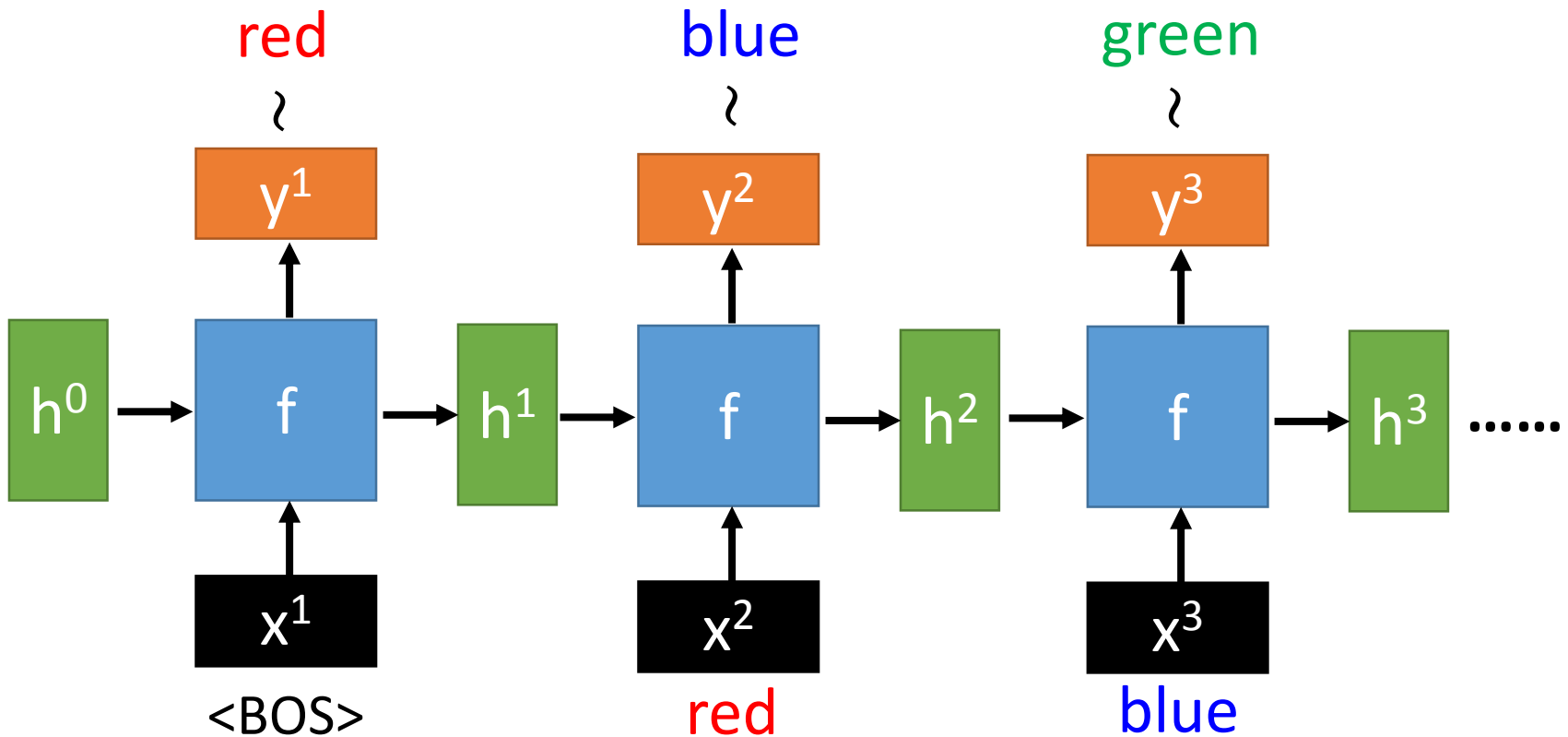
# Generation



Consider as a sentence  
blue red yellow gray .....

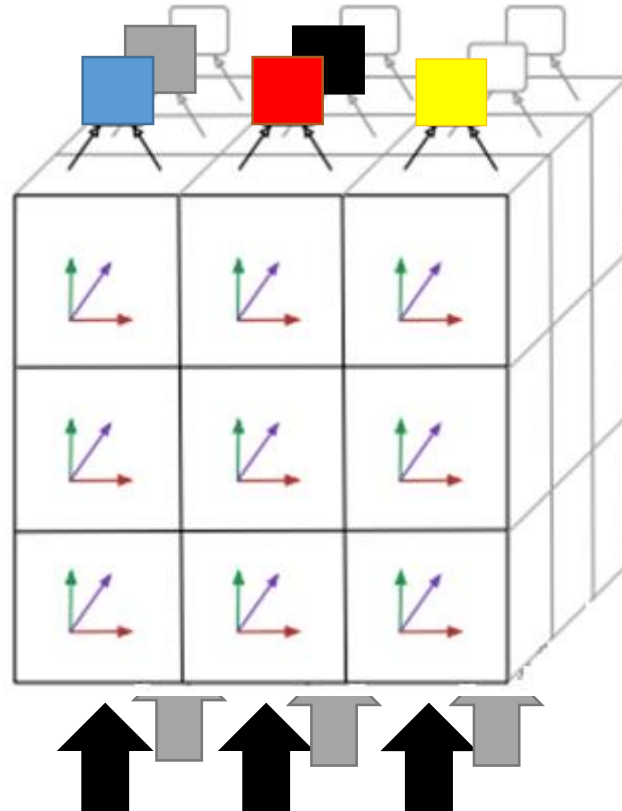
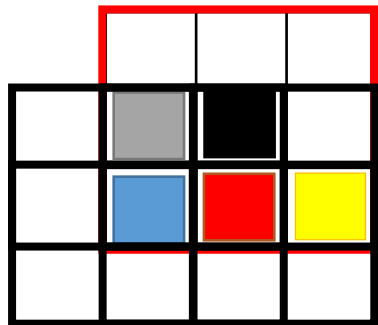
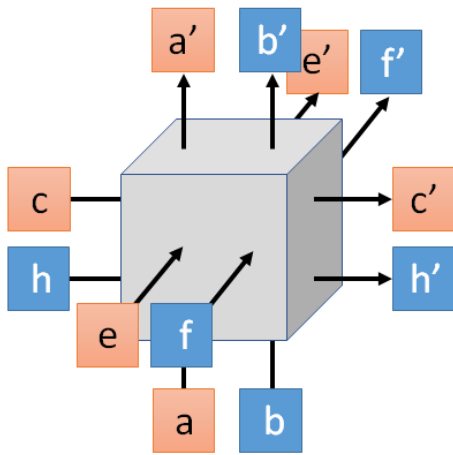
Train a RNN based on the  
“sentences”

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

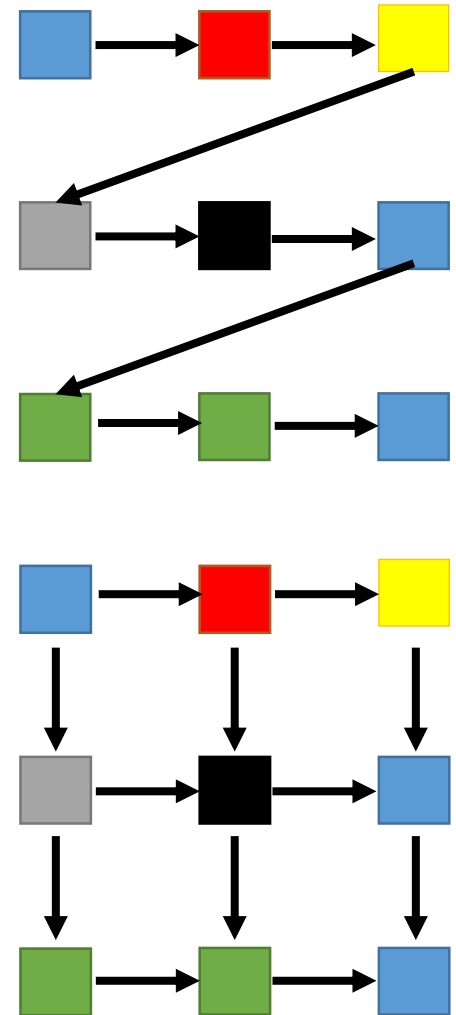


# Generation - PixelRNN

- Images are composed of pixels



3 x 3 images



# Conditional Sequence Generation

# Conditional Generation

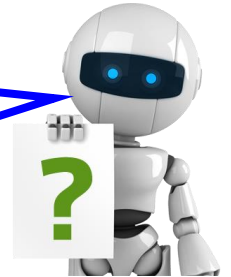
- We don't want to simply generate some random sentences.
- Generate sentences based on conditions:

## Caption Generation

Given  
condition:

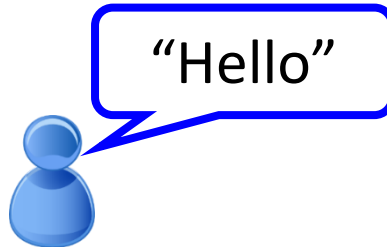


“A young girl  
is dancing.”

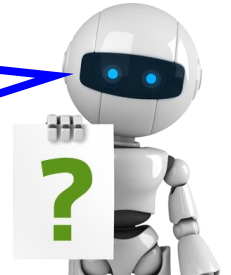


## Chat-bot

Given  
condition:



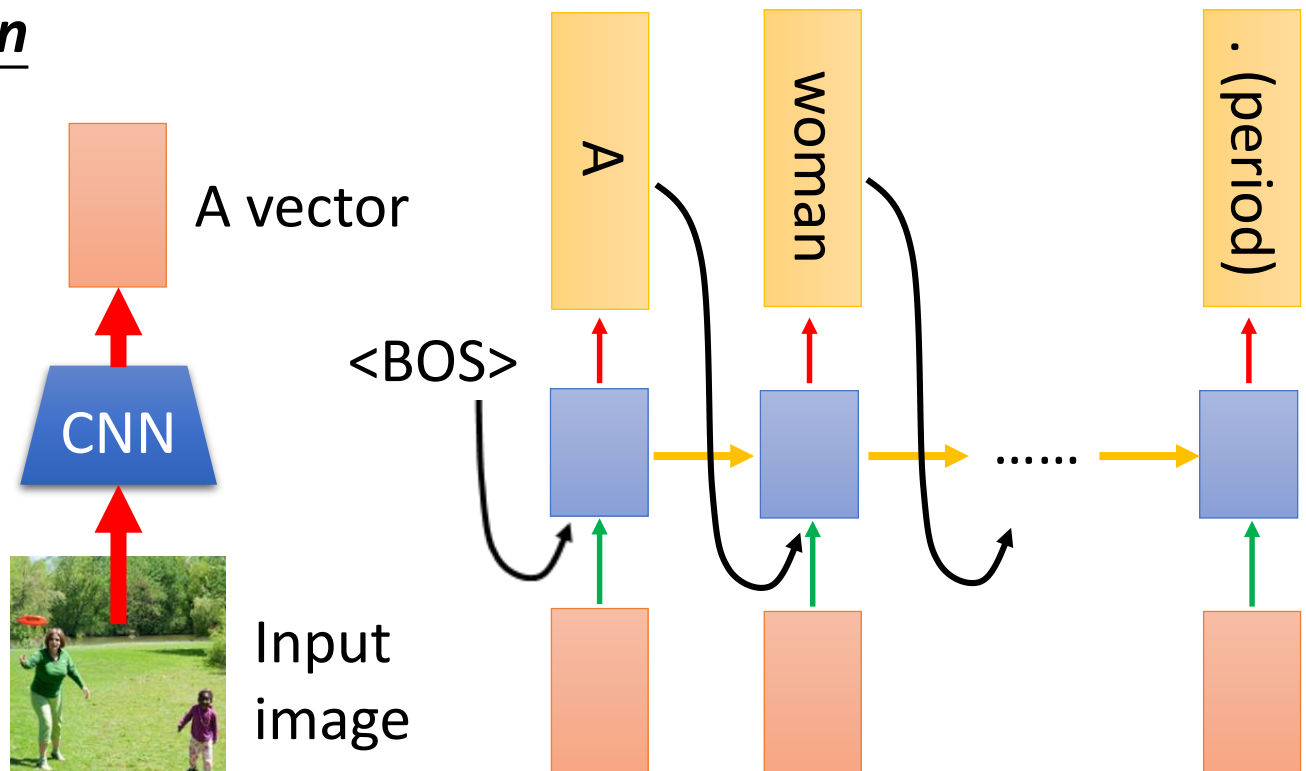
“Hello. Nice  
to see you.”



# Conditional Generation

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

## Image Caption Generation

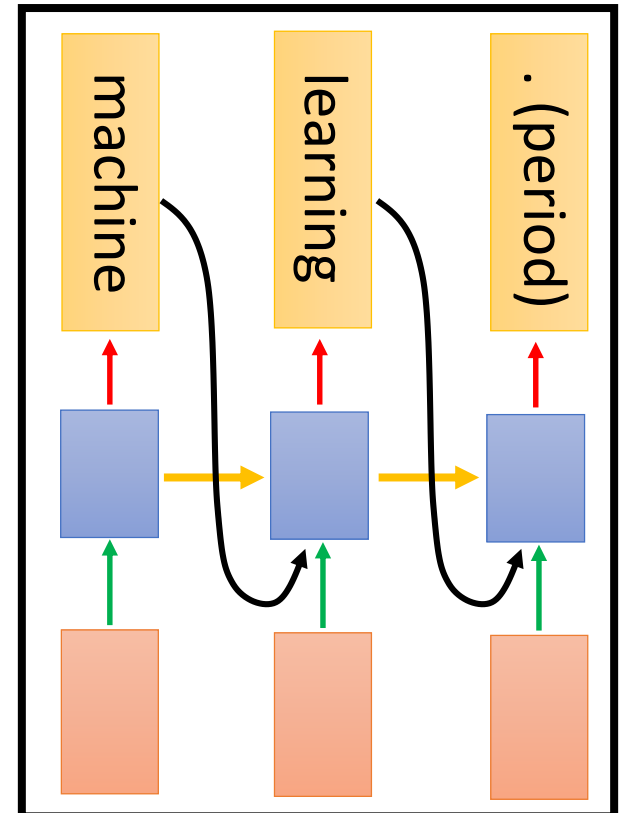
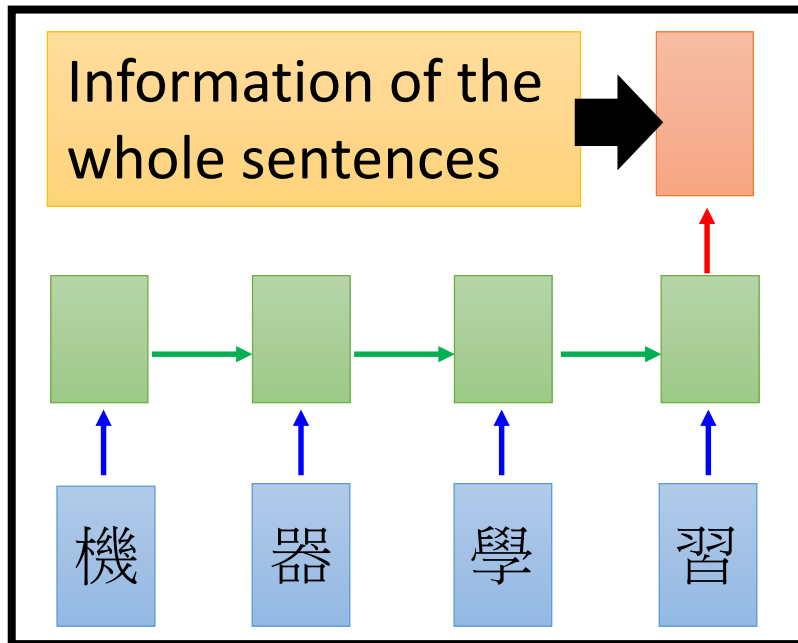




# Conditional Generation

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

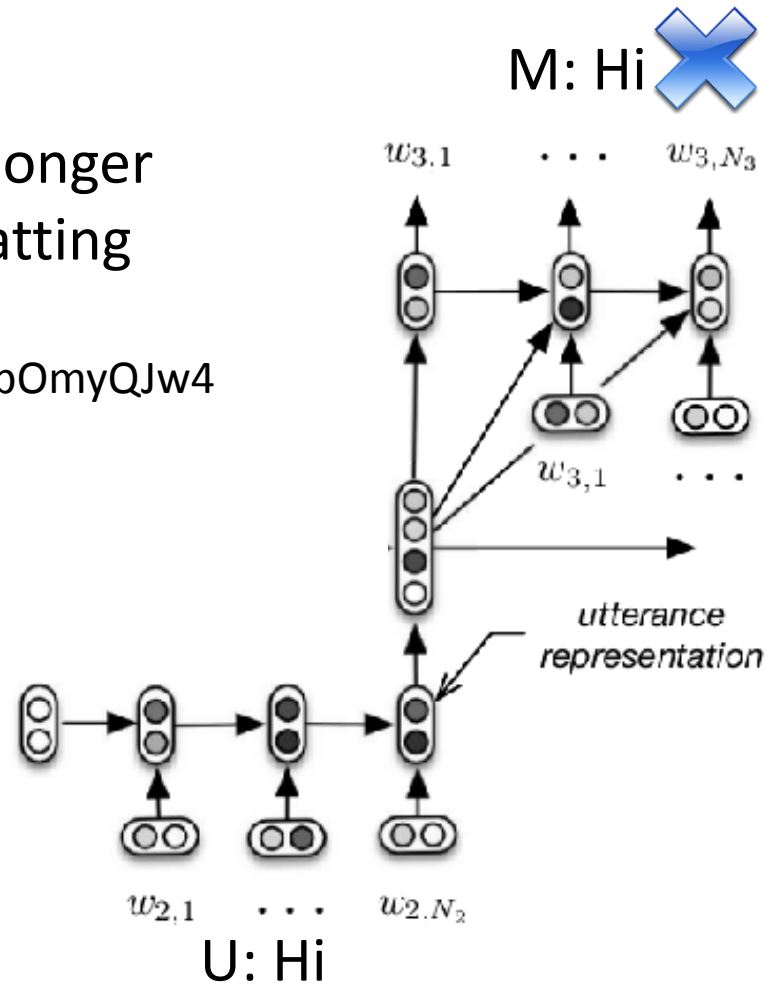
M: Hello

U: Hi

M: Hi

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

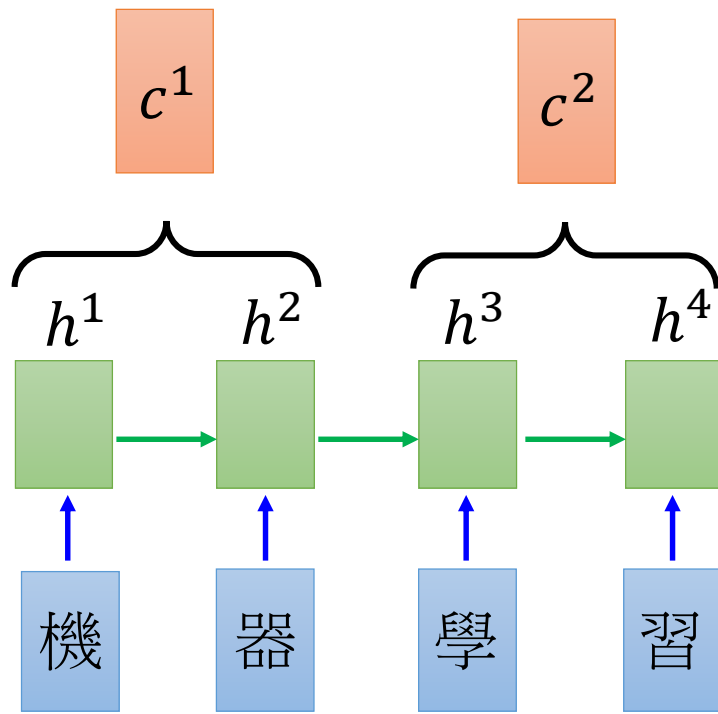
Need to consider longer context during chatting



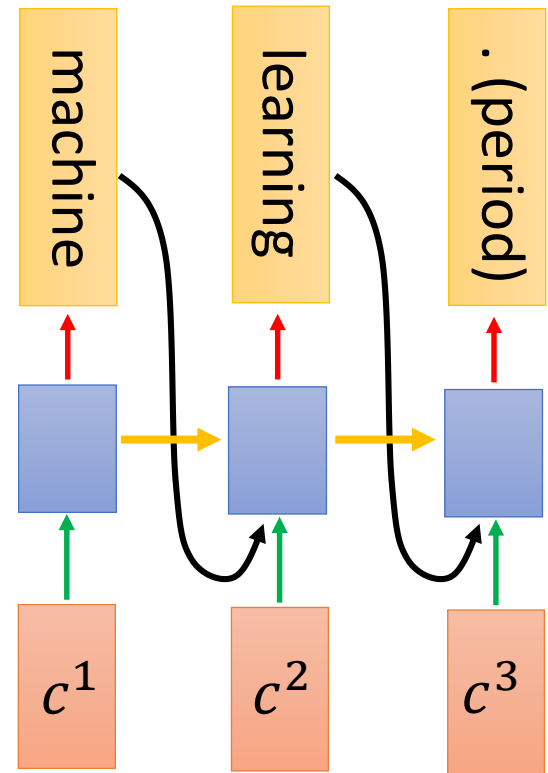
M: Hello

U: Hi

# Dynamic Conditional Generation



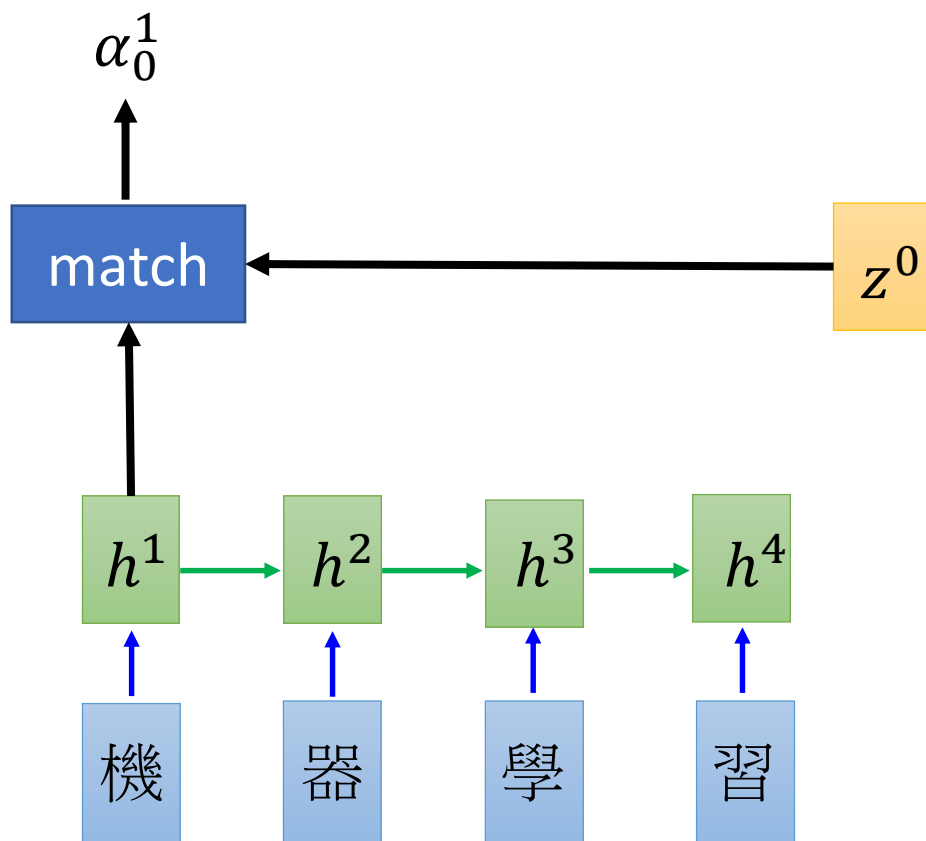
Encoder



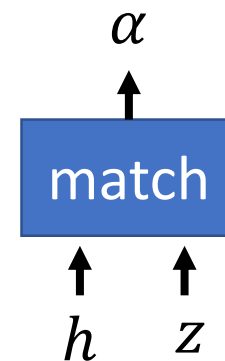
Decoder

# Machine Translation

- Attention-based model



Jointly learned  
with other part  
of the network



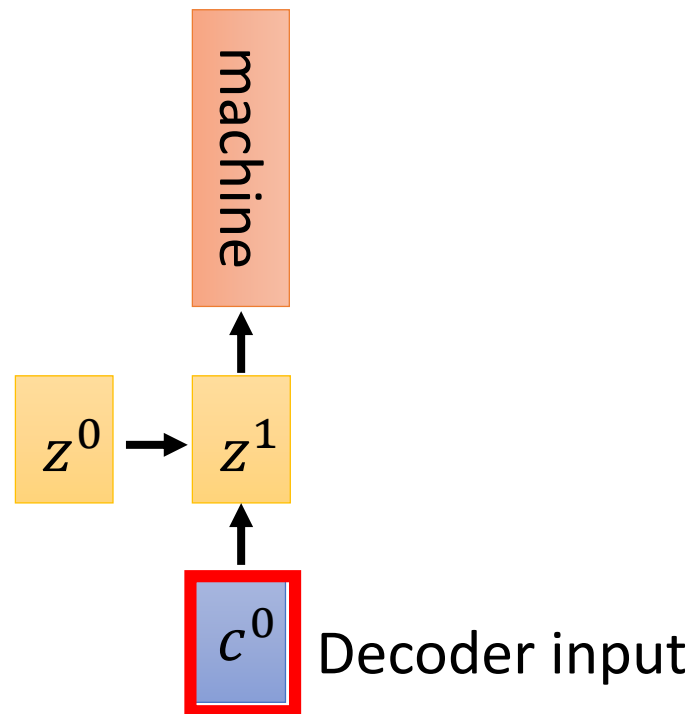
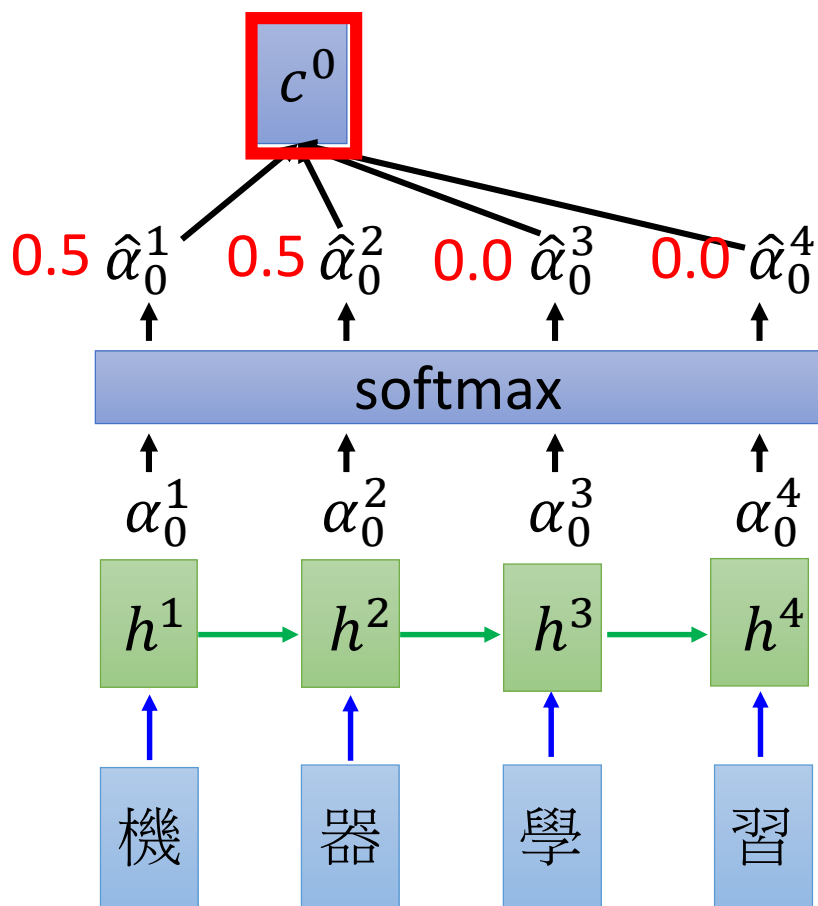
What is **match** ?

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

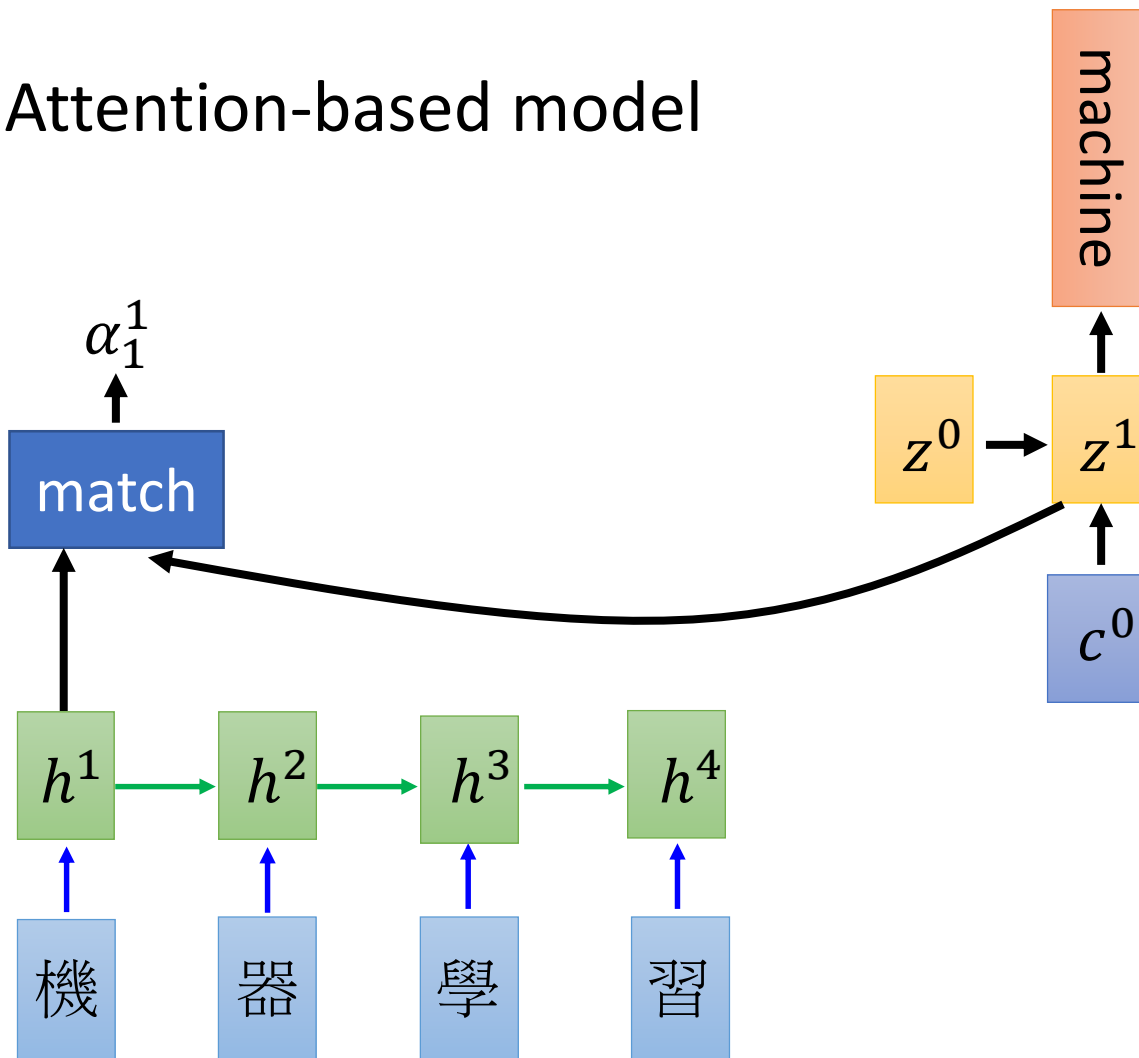
- Attention-based model



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

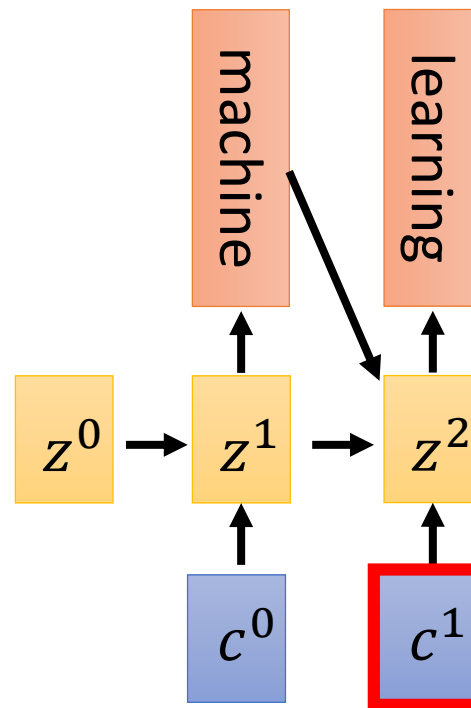
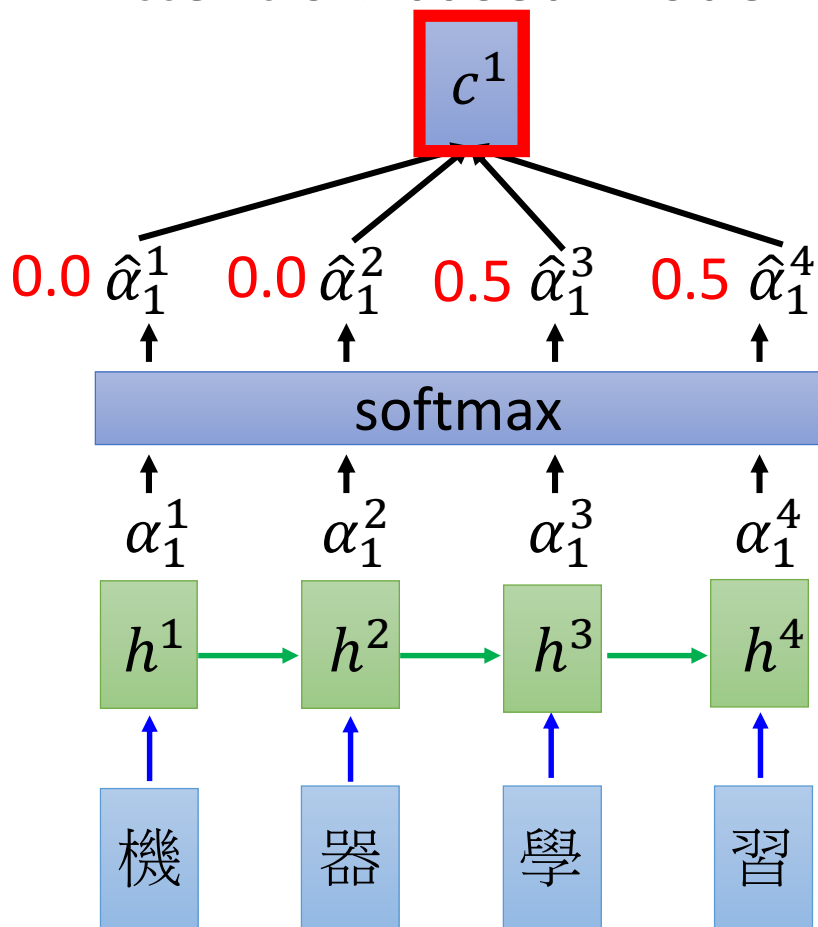
# Machine Translation

- Attention-based model



# Machine Translation

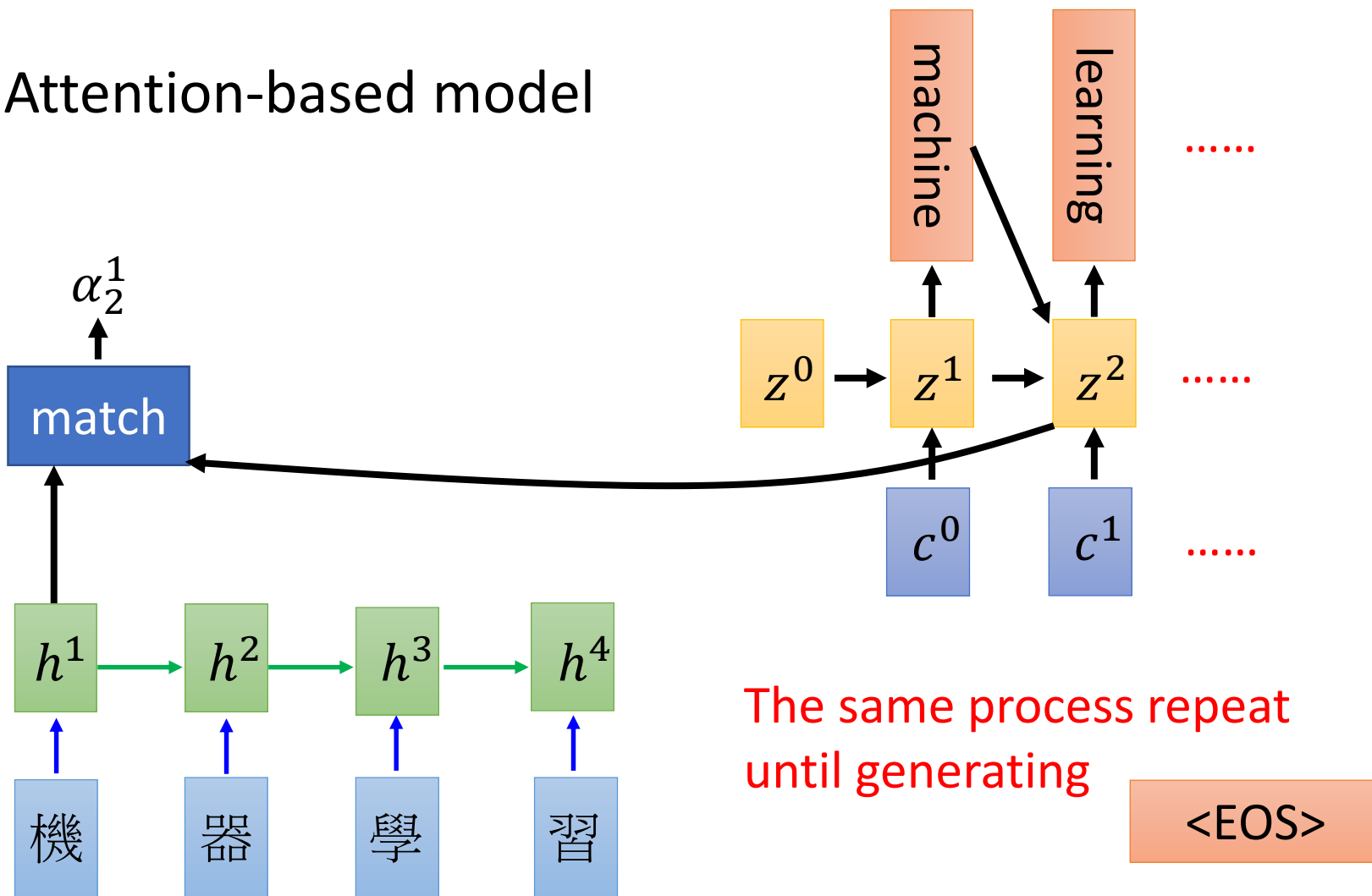
- Attention-based model



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

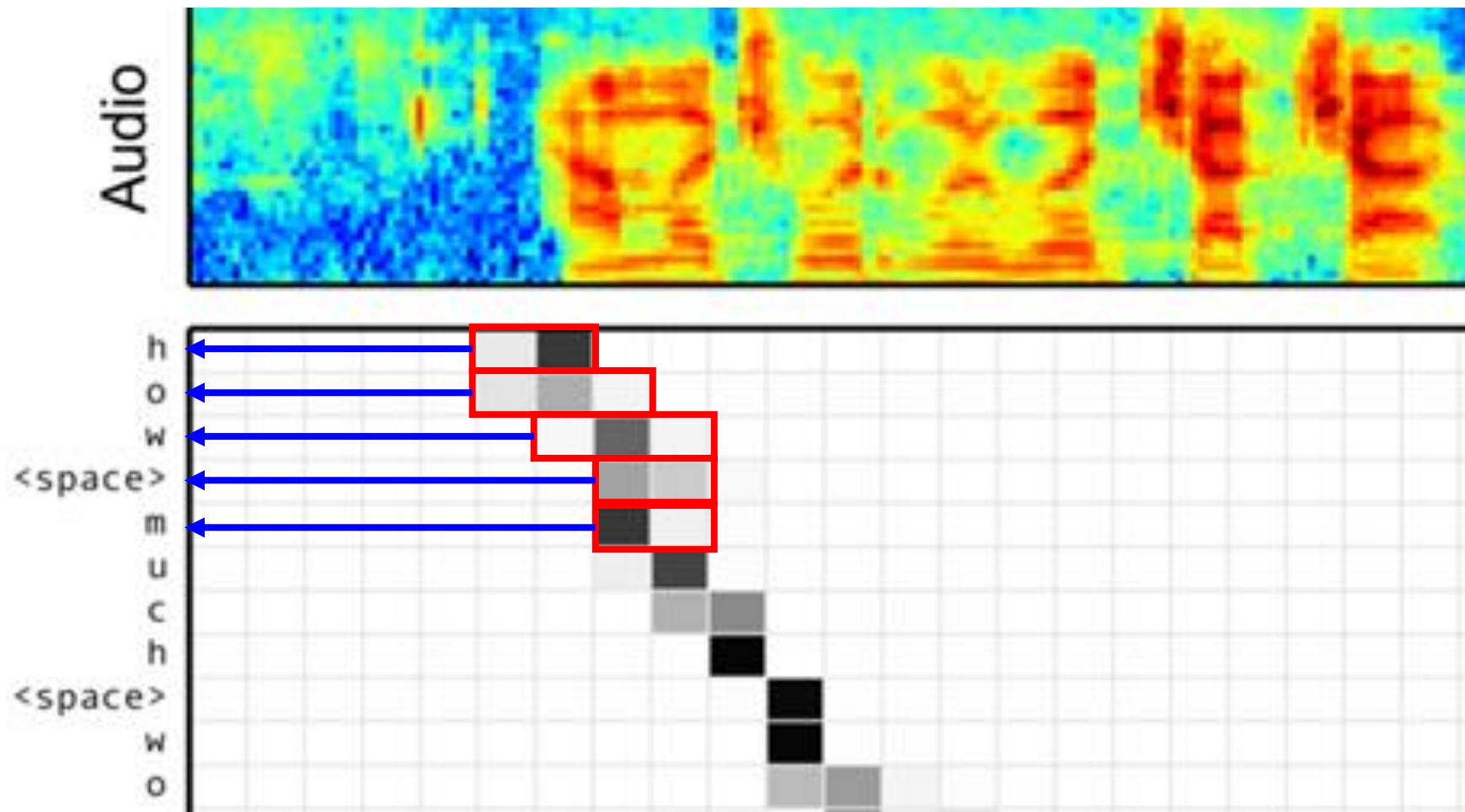
# Machine Translation

- Attention-based model





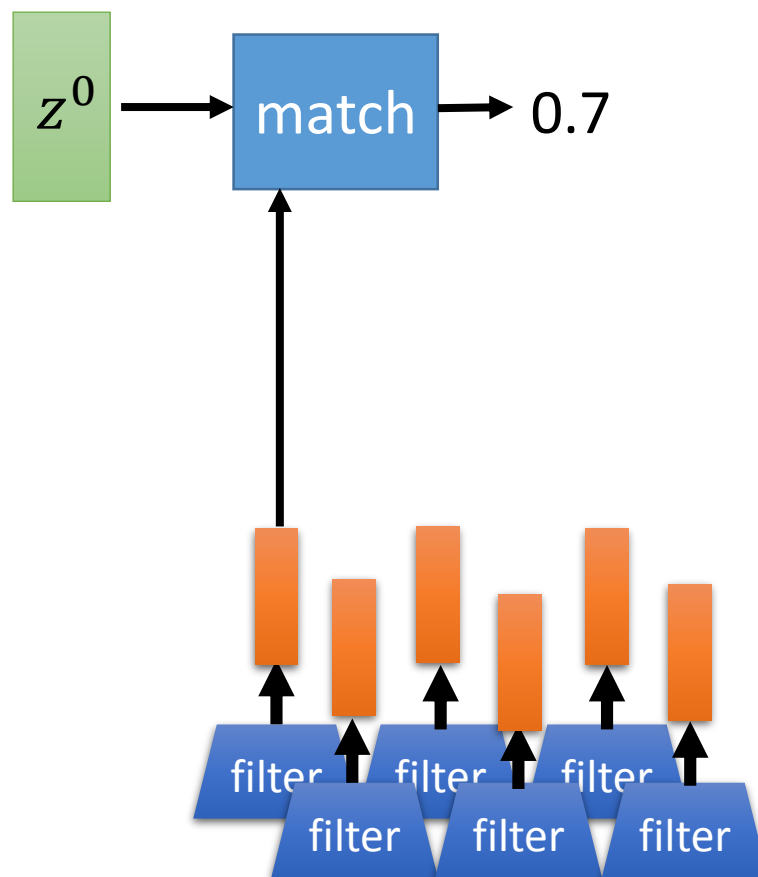
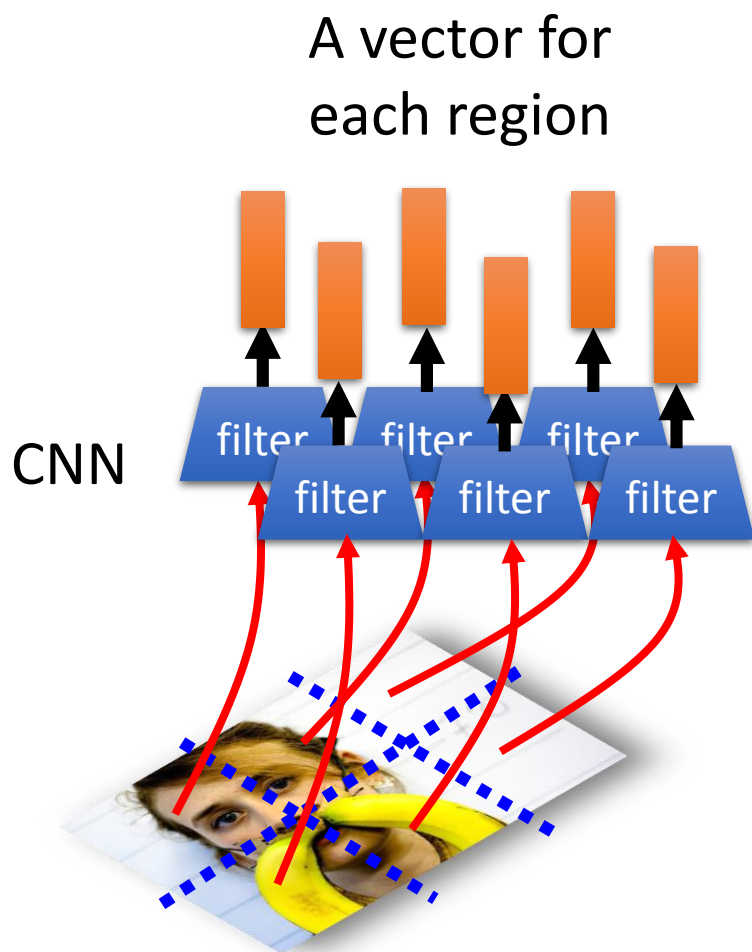
# Speech Recognition



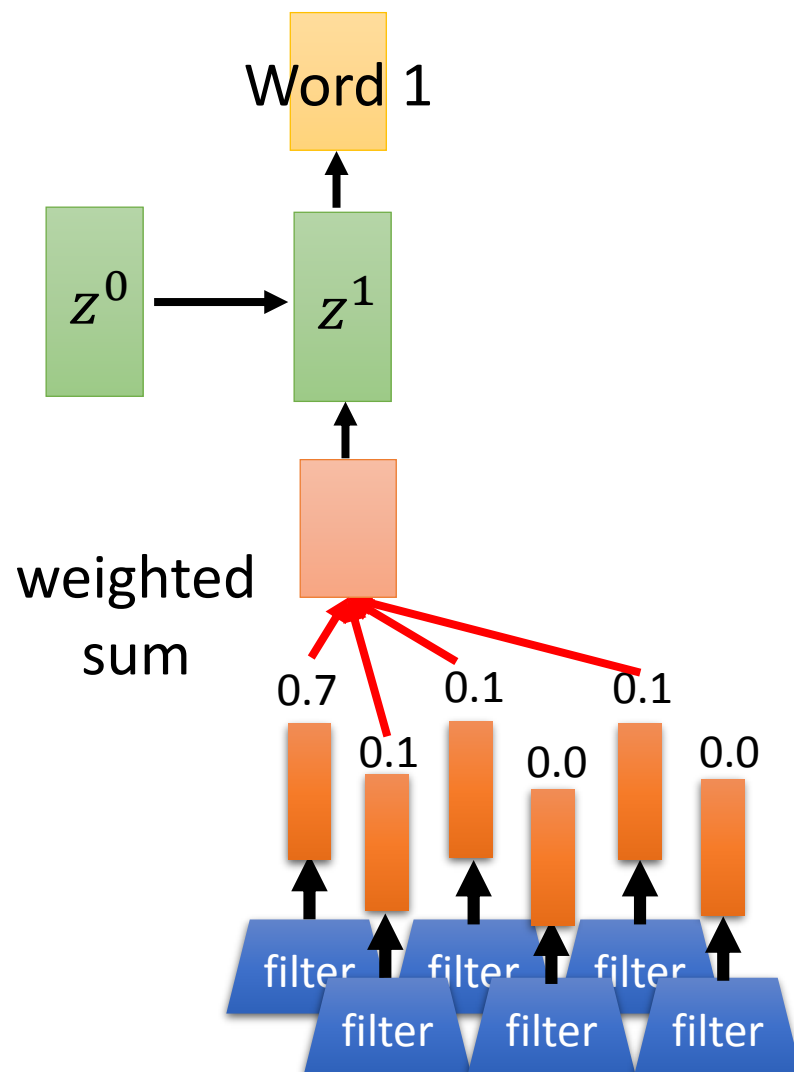
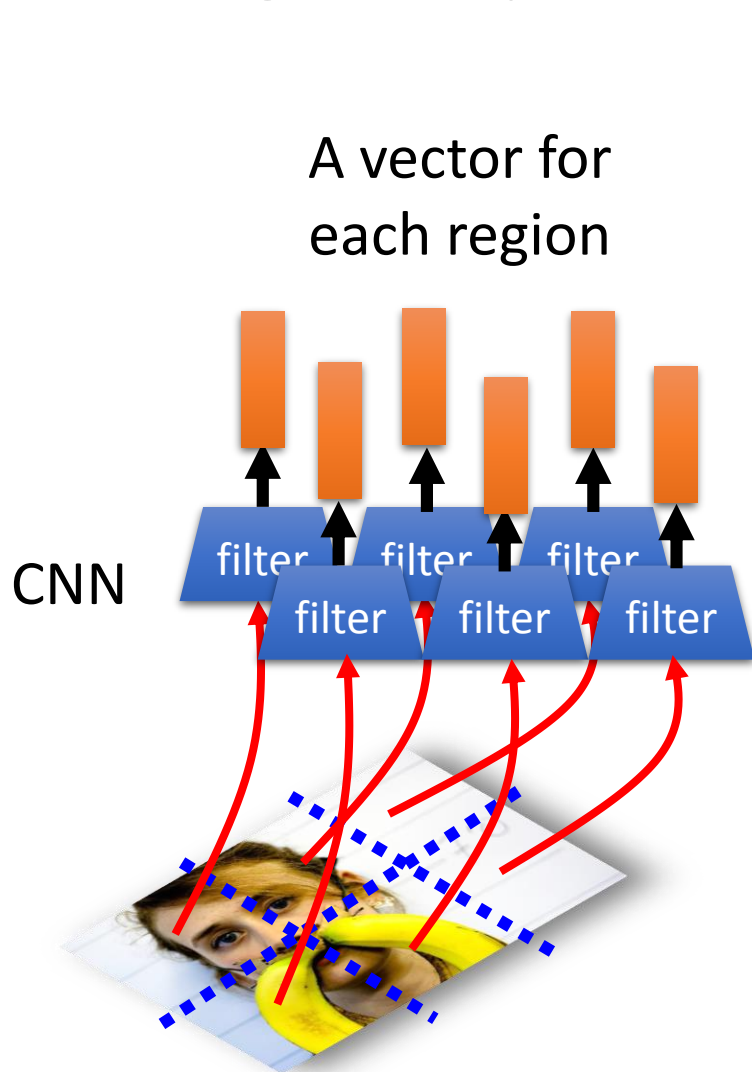
Model	Clean WER	Noisy WER
CLDNN-HMM [22]	8.0	8.9
LAS	14.1	16.5
LAS + LM Rescoring	10.3	12.0

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

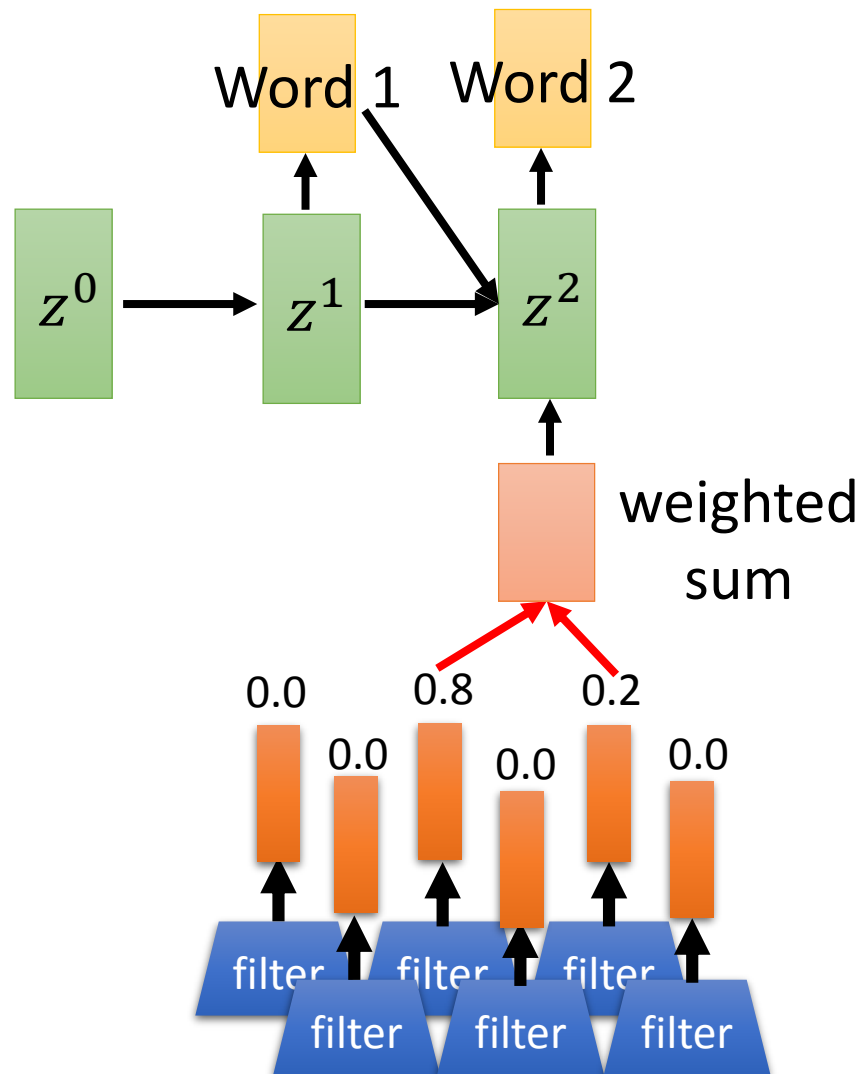
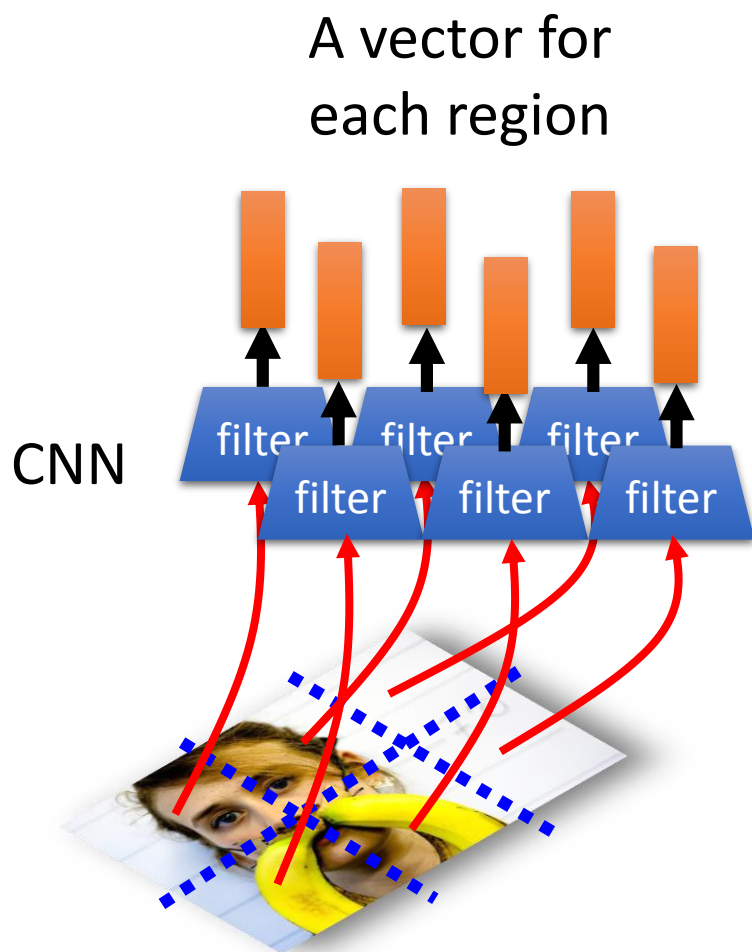
# Image Caption Generation



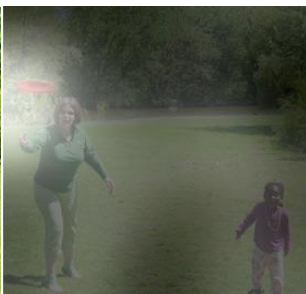
# Image Caption Generation



# Image Caption Generation



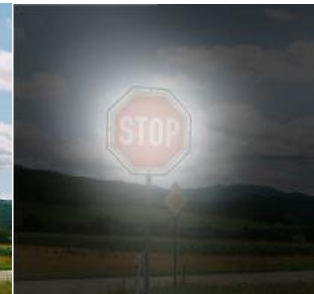
# Image Caption Generation



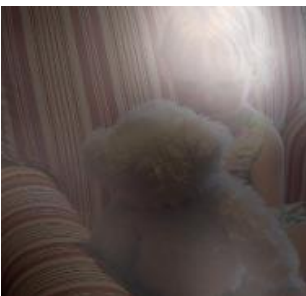
A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



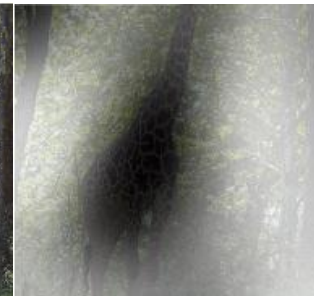
A stop sign is on a road with a mountain in the background.



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



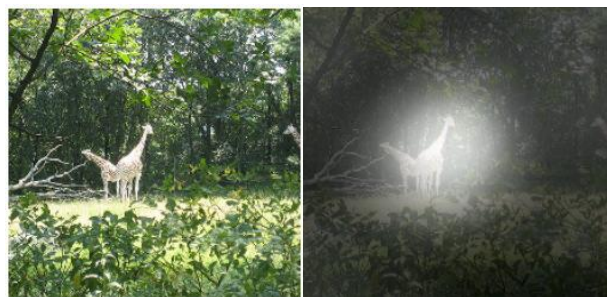
A group of people sitting on a boat in the water.



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

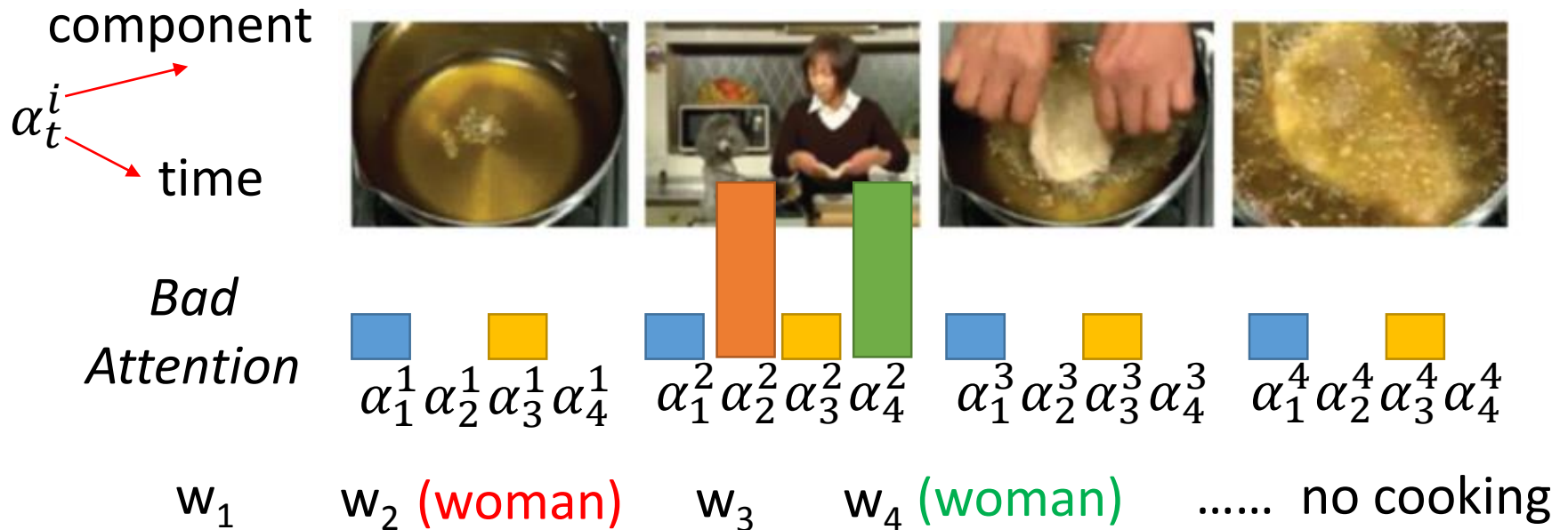
**Someone** is **frying** a **fish** in a **pot**

# 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

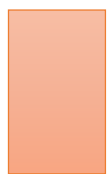
Over the generation

# Mismatch between Train and Test

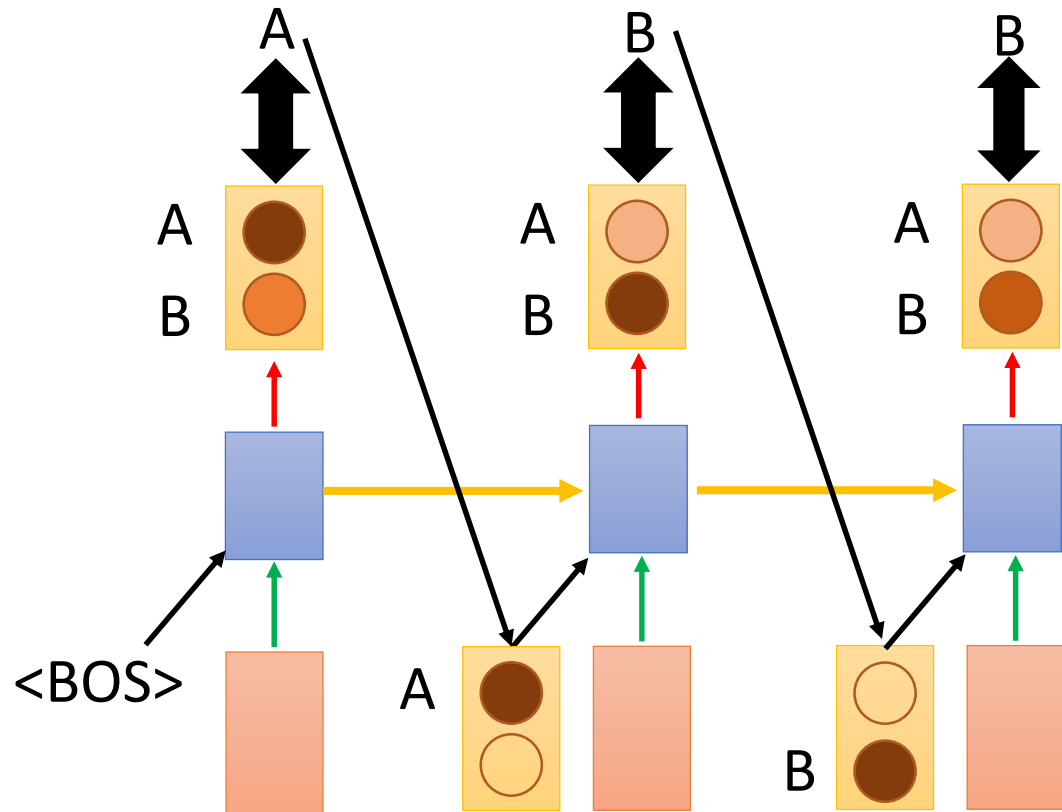
- Training

$$C = \sum_t C_t$$

Minimizing cross-entropy of each component

 : condition

Reference:



# Mismatch between Train and Test

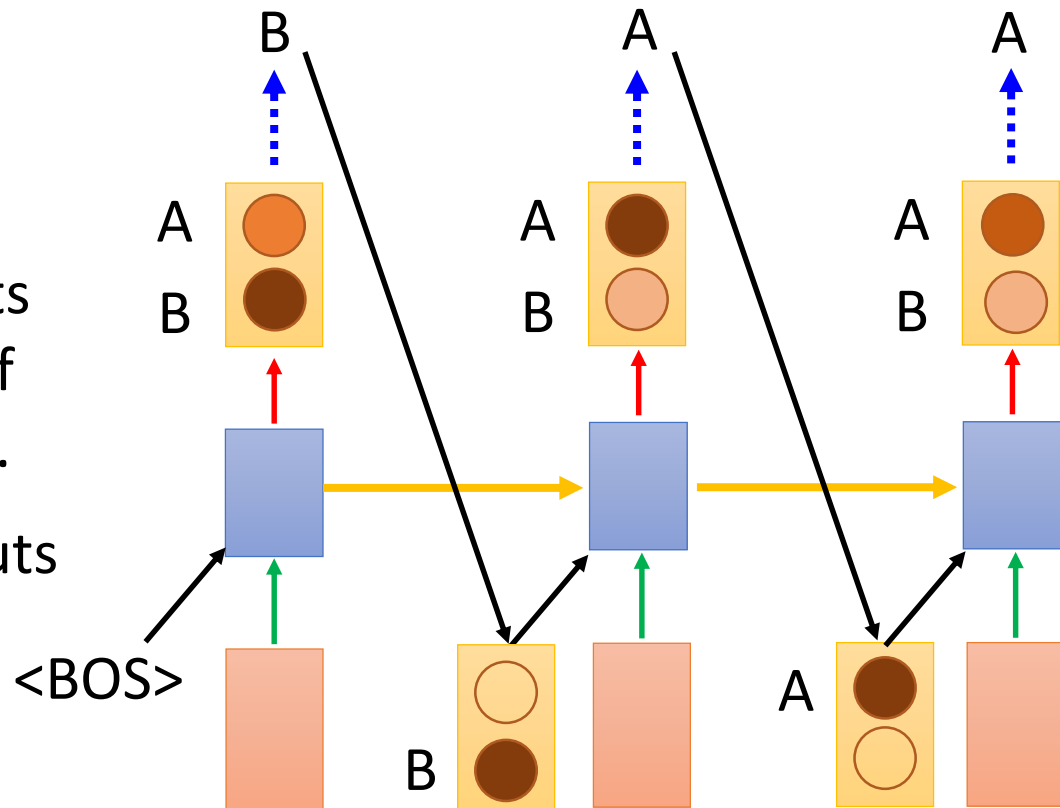
- Generation

We do not know the reference

Testing: The inputs are the outputs of the last time step.

Training: The inputs are reference.

*Exposure Bias*



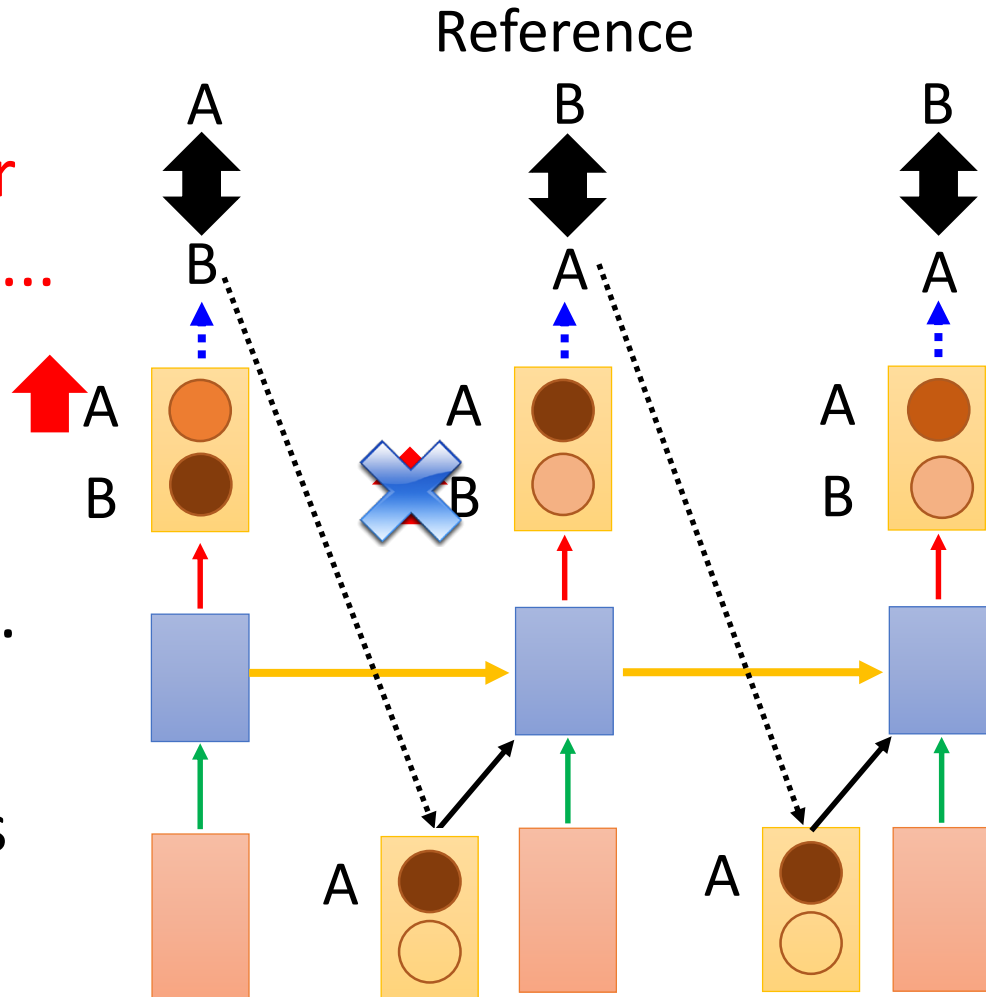


# Modifying Training Process?

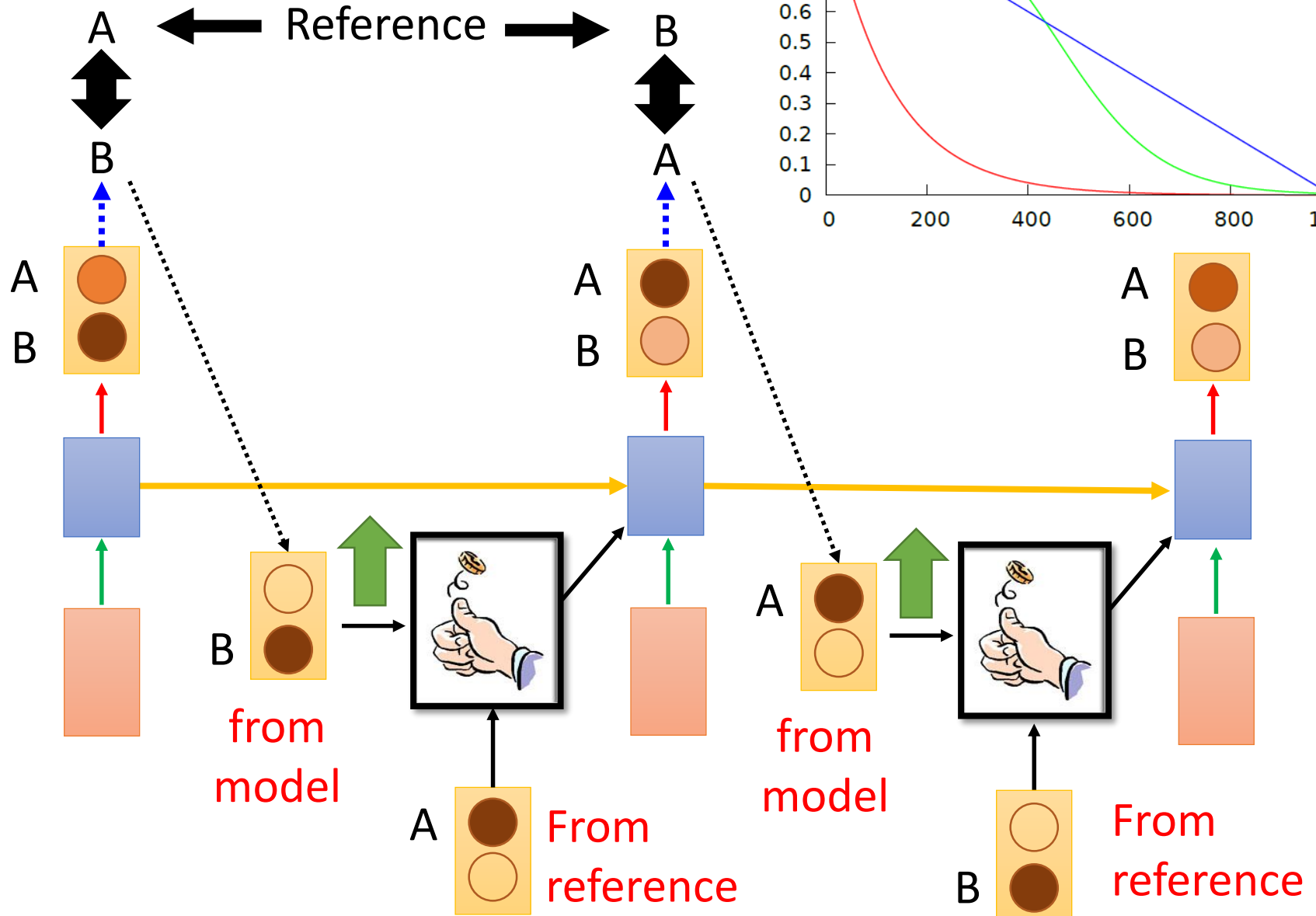
When we try to decrease the loss for both steps 1 and 2 .....

Training is matched to testing.

In practice, it is hard to train in this way.



# Scheduled Sampling

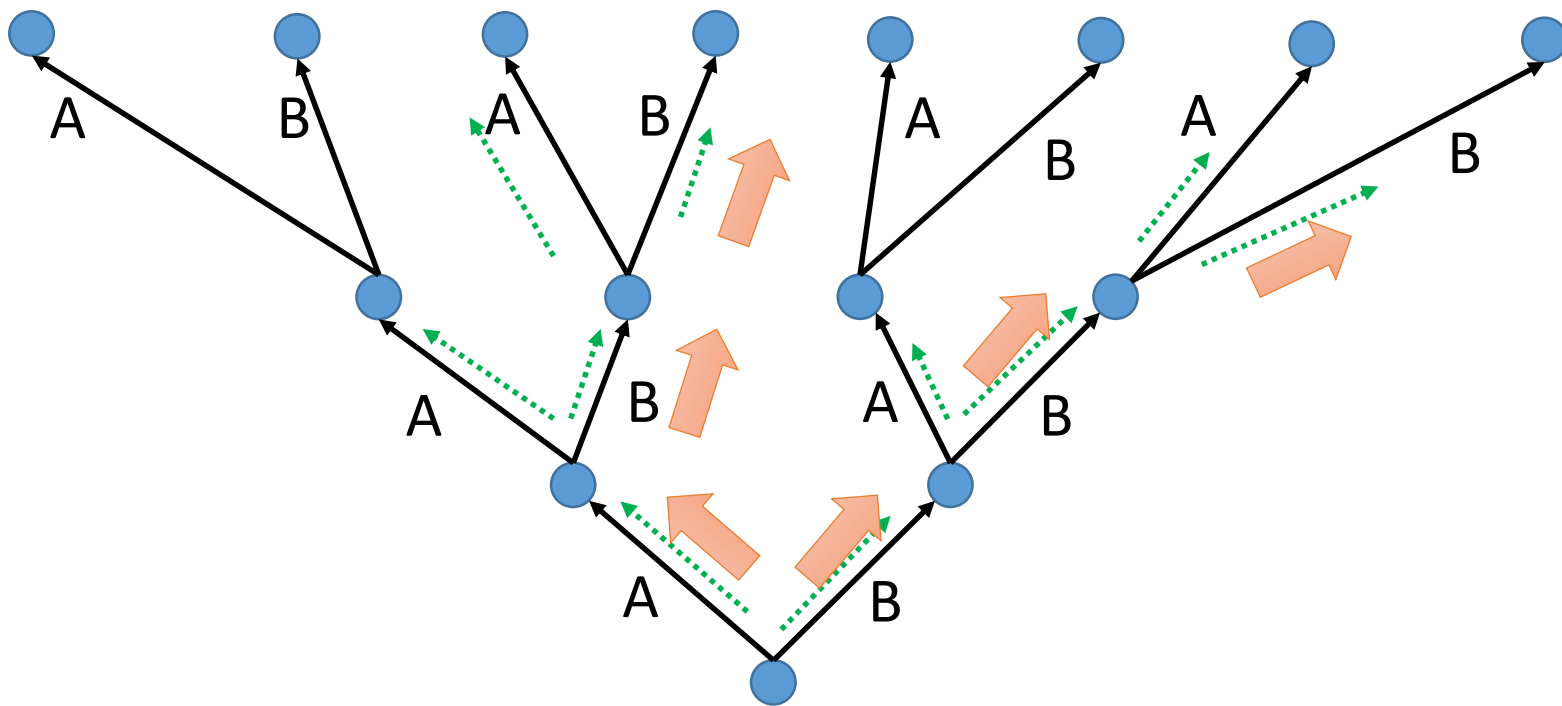
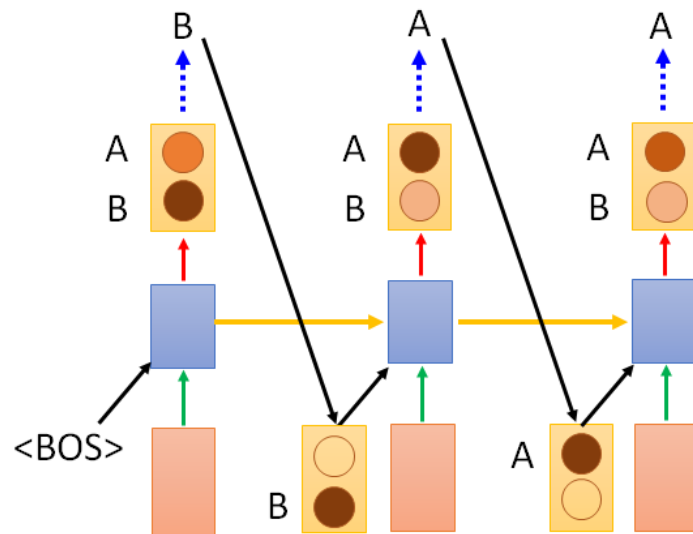




# Beam Search

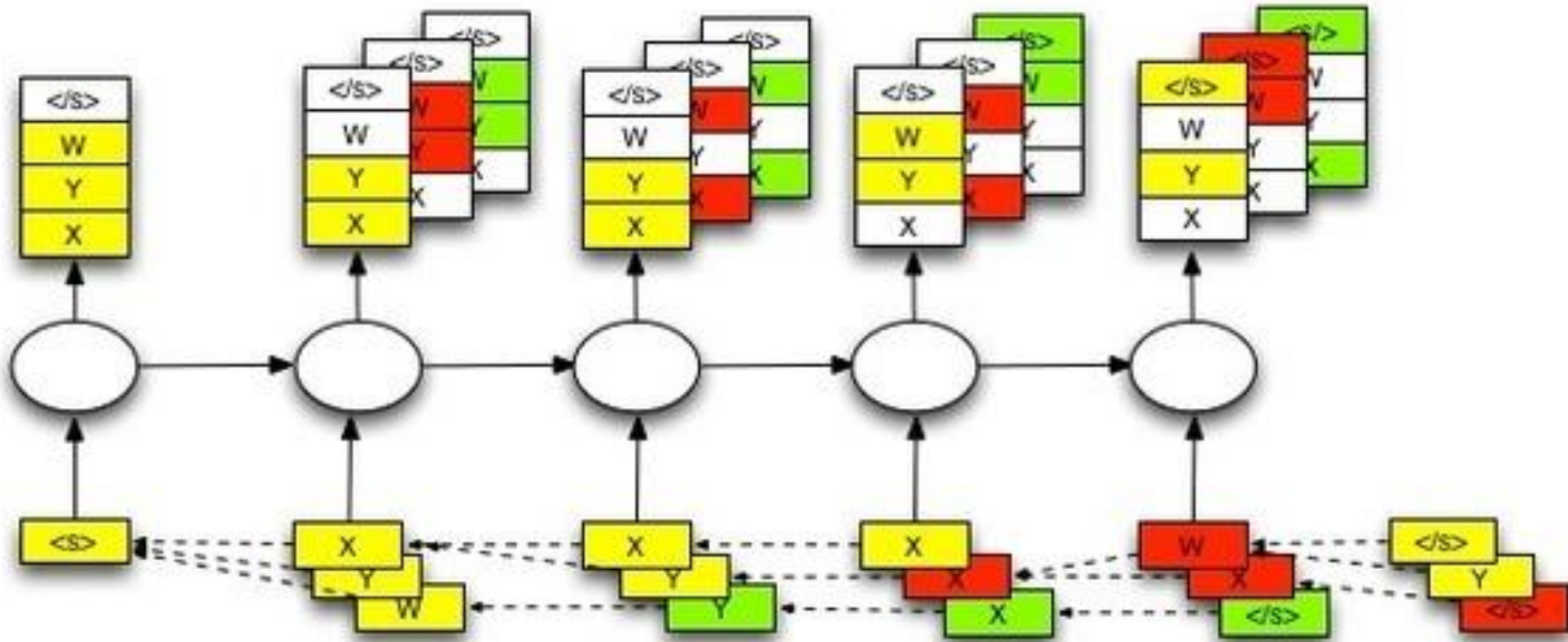
Keep several best path at each step

Beam size = 2





# Beam Search



The size of beam is 3 in this example.

# Better Idea?

I am .....



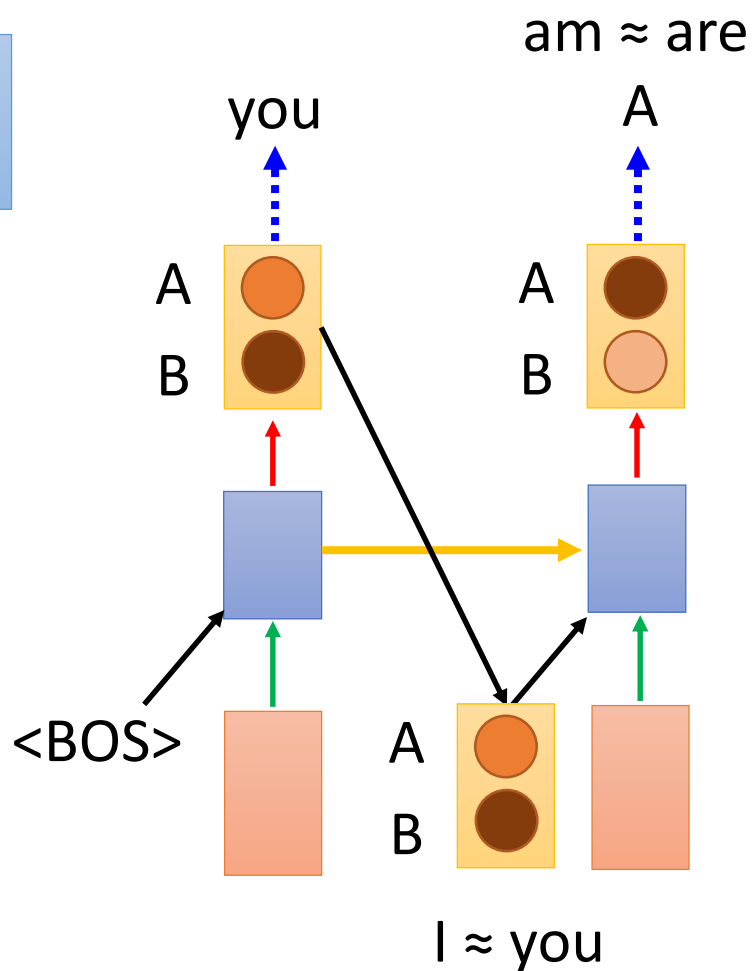
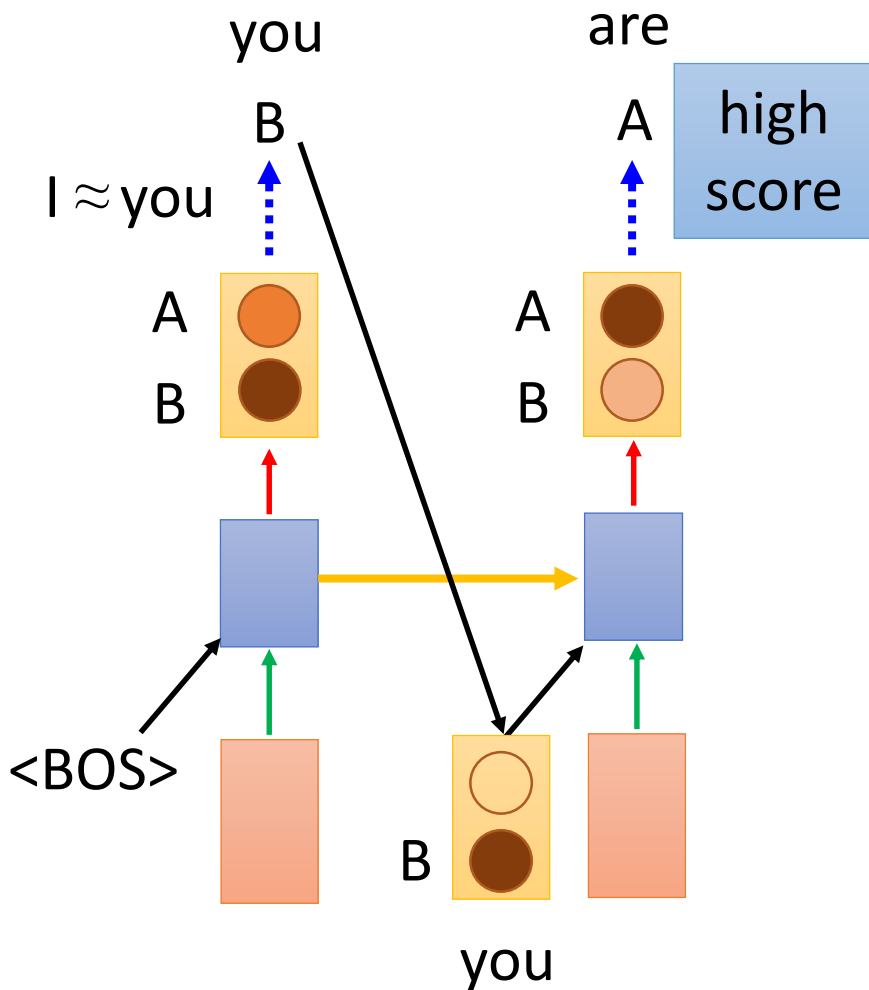
You are .....



I are .....



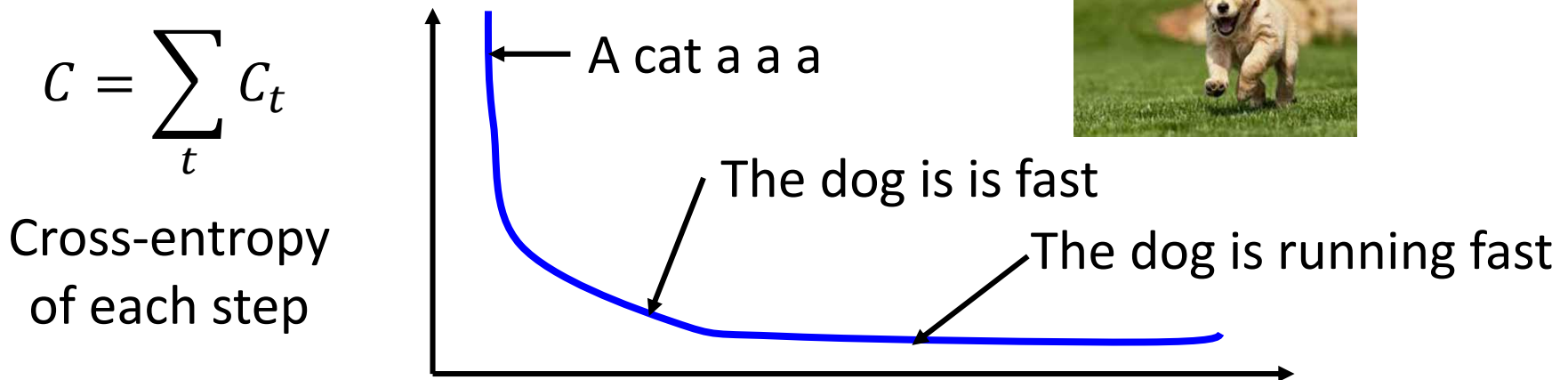
You am .....



# 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"

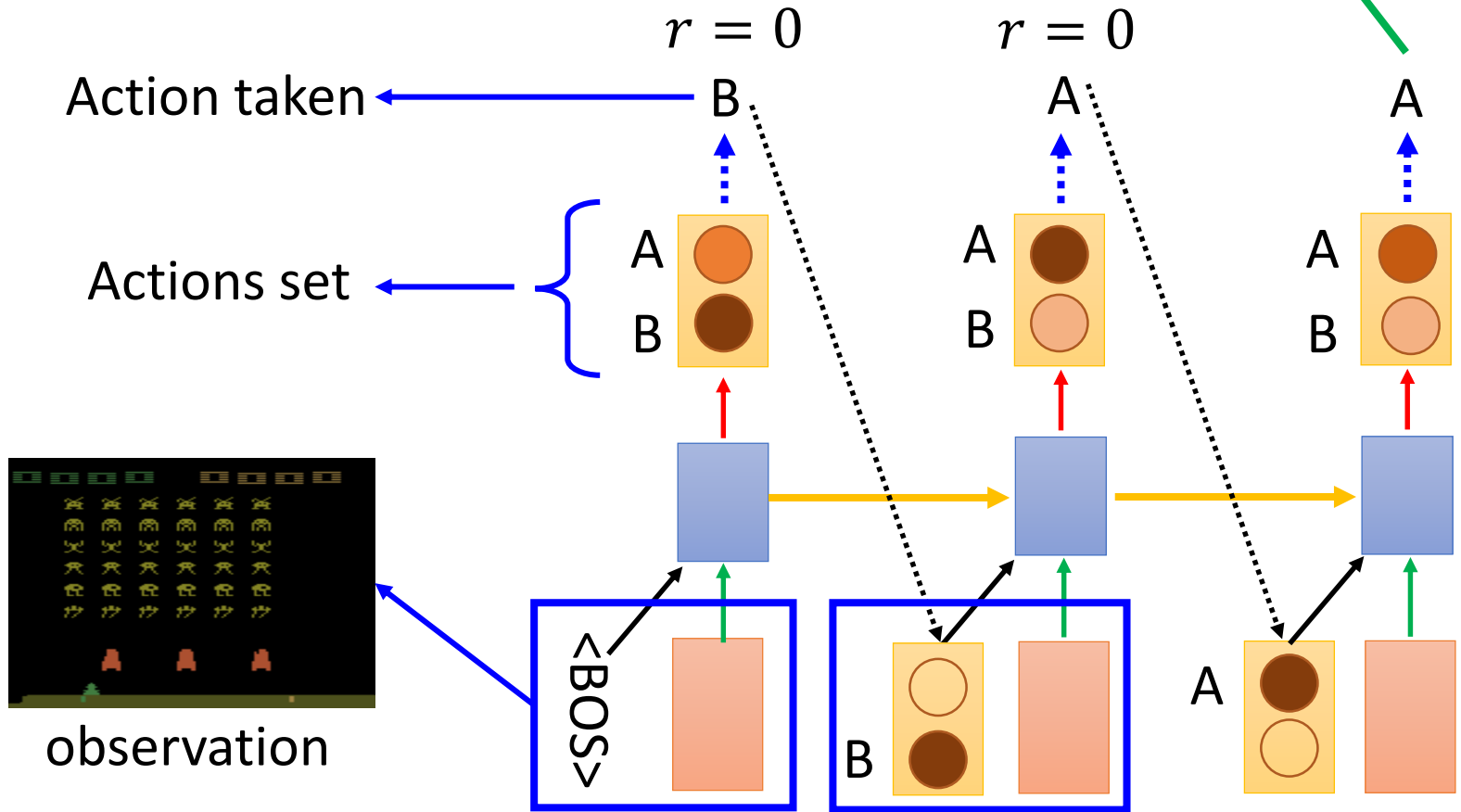


Obtain reward  
 $r_2 = 5$

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