

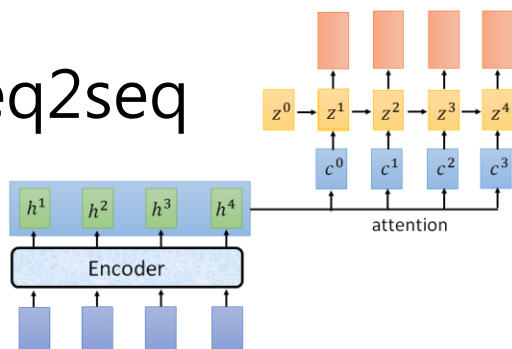


*Speech
Recognition*

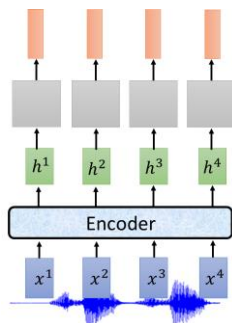
HUNG-YI LEE 李宏毅

Last Time

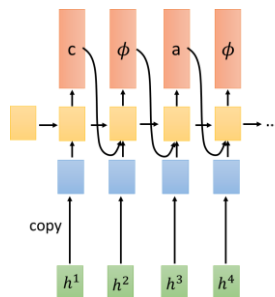
LAS: 就是 seq2seq



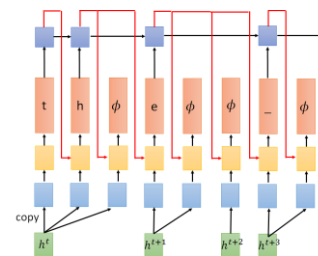
CTC: decoder 是 linear classifier 的 seq2seq



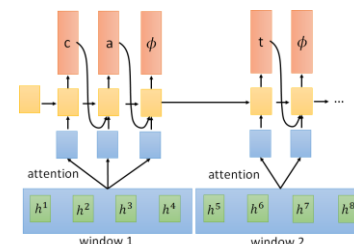
RNA: 輸入一個東西就要輸出一個東西的 seq2seq



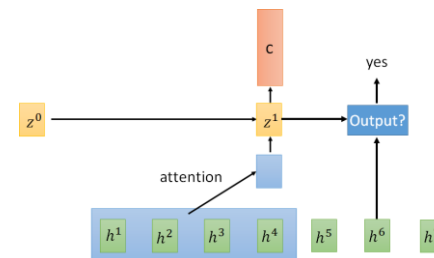
RNN-T: 輸入一個東西可以輸出多個東西的 seq2seq



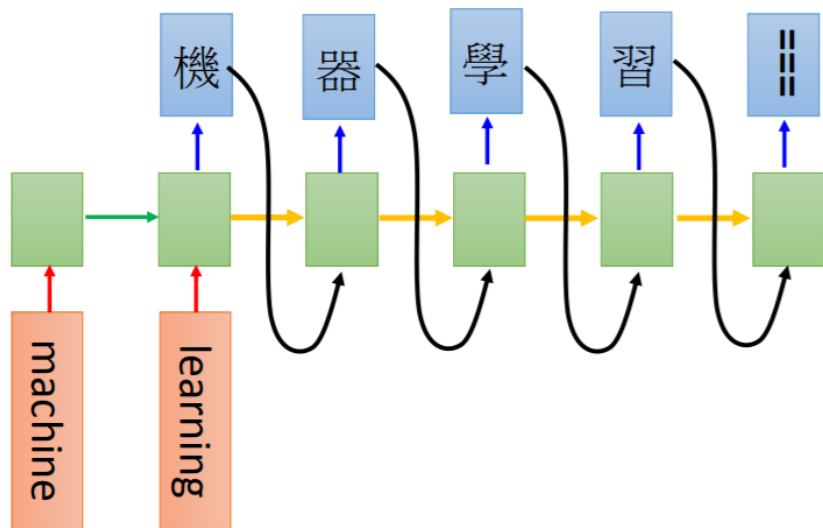
Neural Transducer: 每次輸入一個 window 的 RNN-T



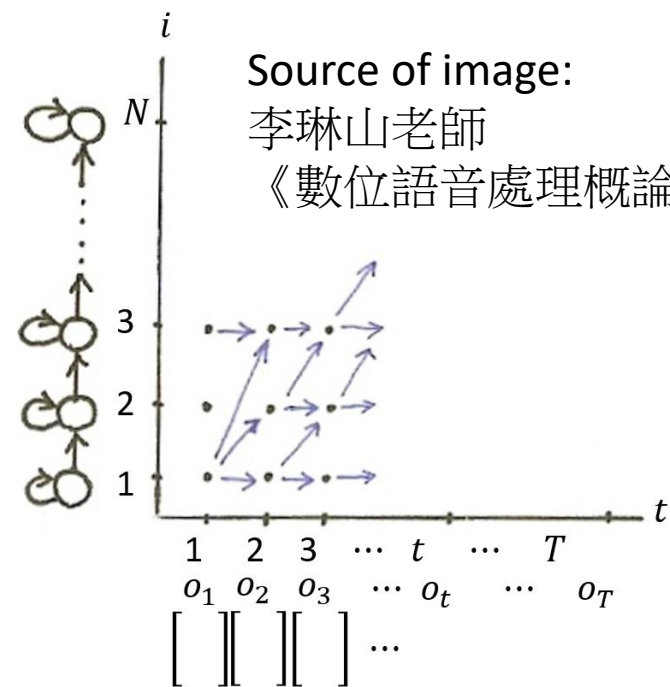
MoCha: window 移動伸縮自如的 Neural Transducer



Two Points of Views

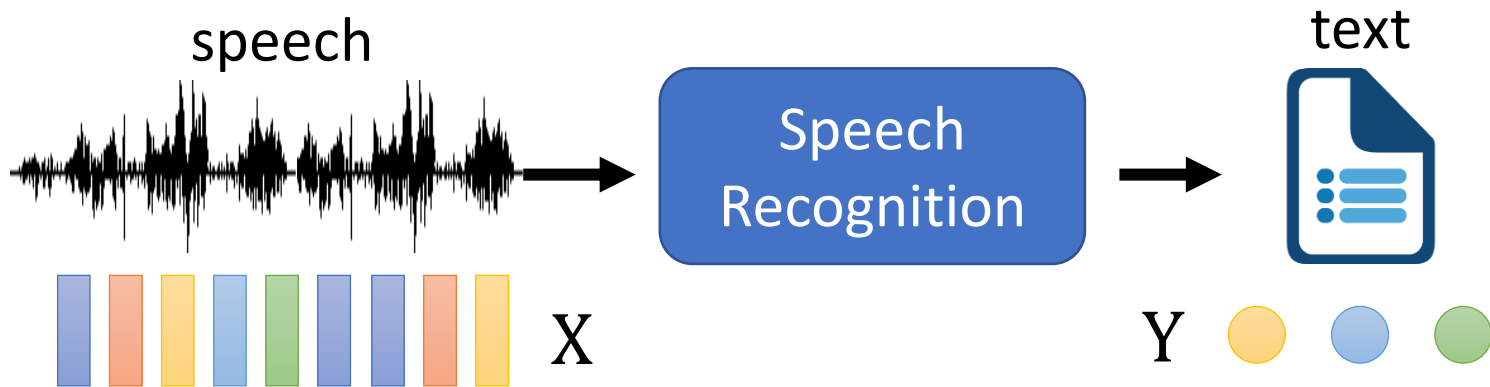


Seq-to-seq



HMM

Hidden Markov Model (HMM)



$$Y^* = \mathop{\text{arg max}}_Y P(Y|X)$$

Decode

$$= \mathop{\text{arg max}}_Y \frac{P(X|Y)P(Y)}{P(X)}$$

$$= \mathop{\text{arg max}}_Y P(X|Y)P(Y)$$

$P(X|Y)$: HMM

Acoustic Model

$P(Y)$:

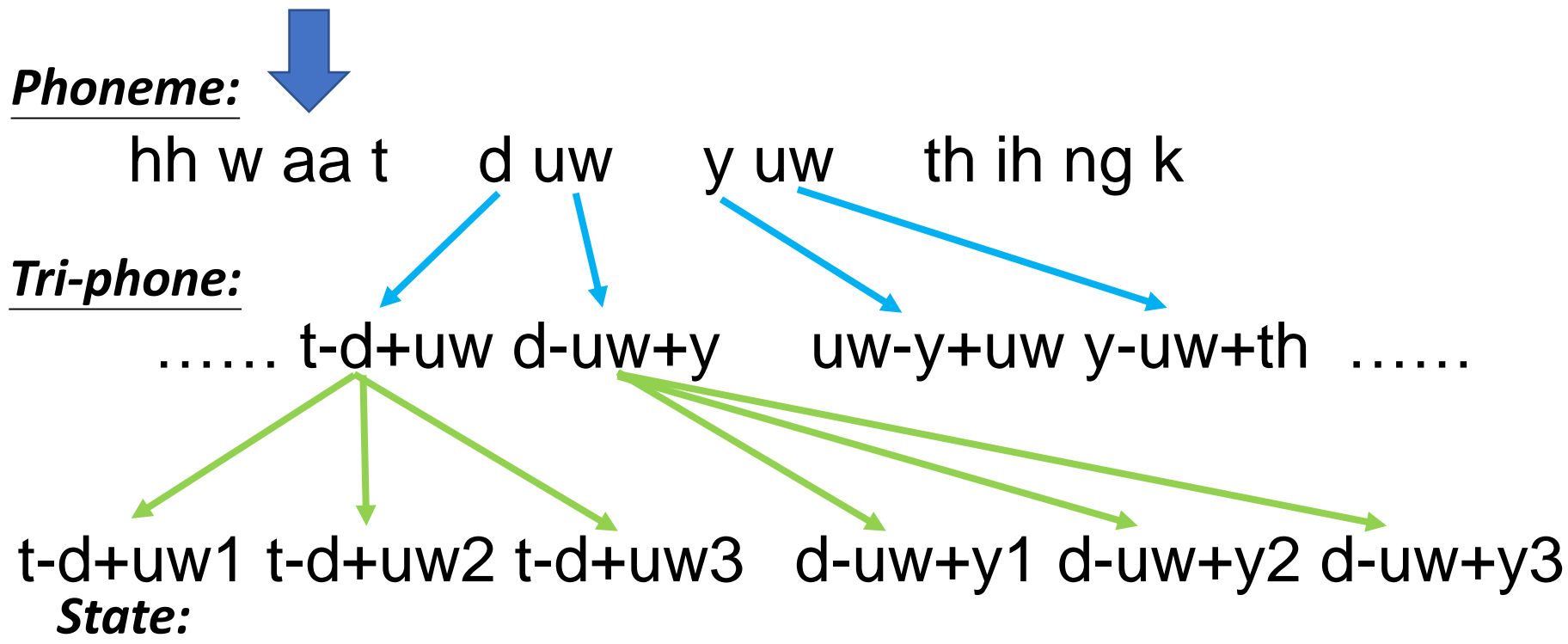
Language Model

HMM

$$P(X|Y) \longrightarrow P(X|S)$$

A token sequence Y corresponds to a sequence of **states** S

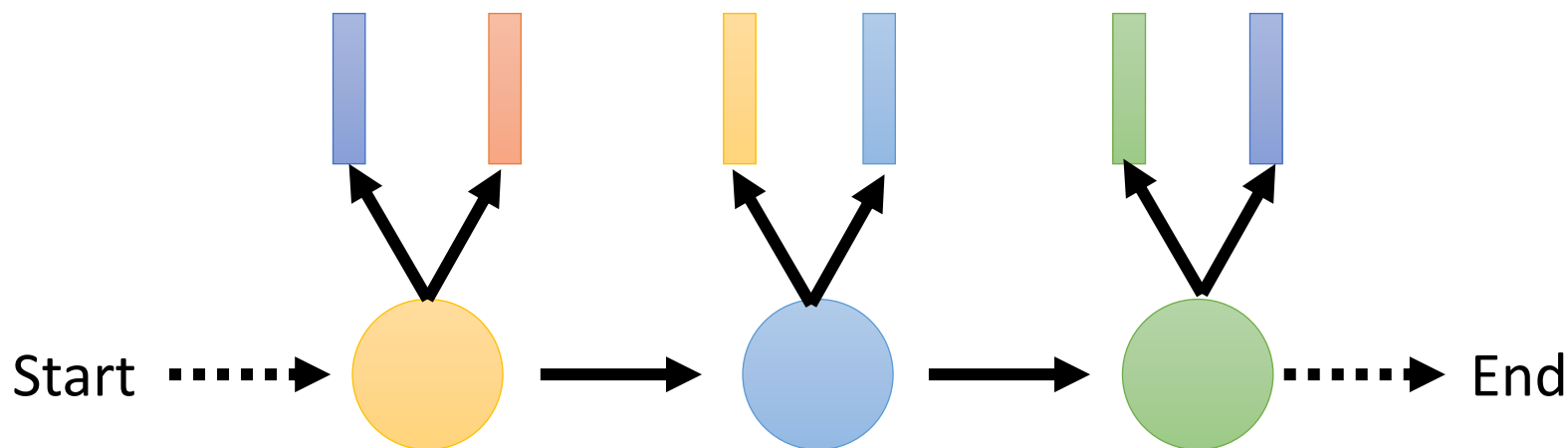
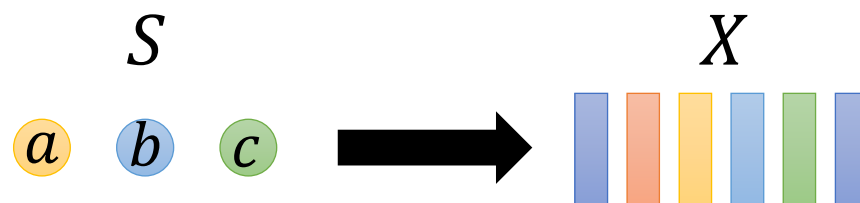
what do you think



HMM

$$P(X|Y) \longrightarrow P(X|S)$$

A sentence Y corresponds to a sequence of **states** S



HMM

$$P(X|Y) \longrightarrow P(X|S)$$

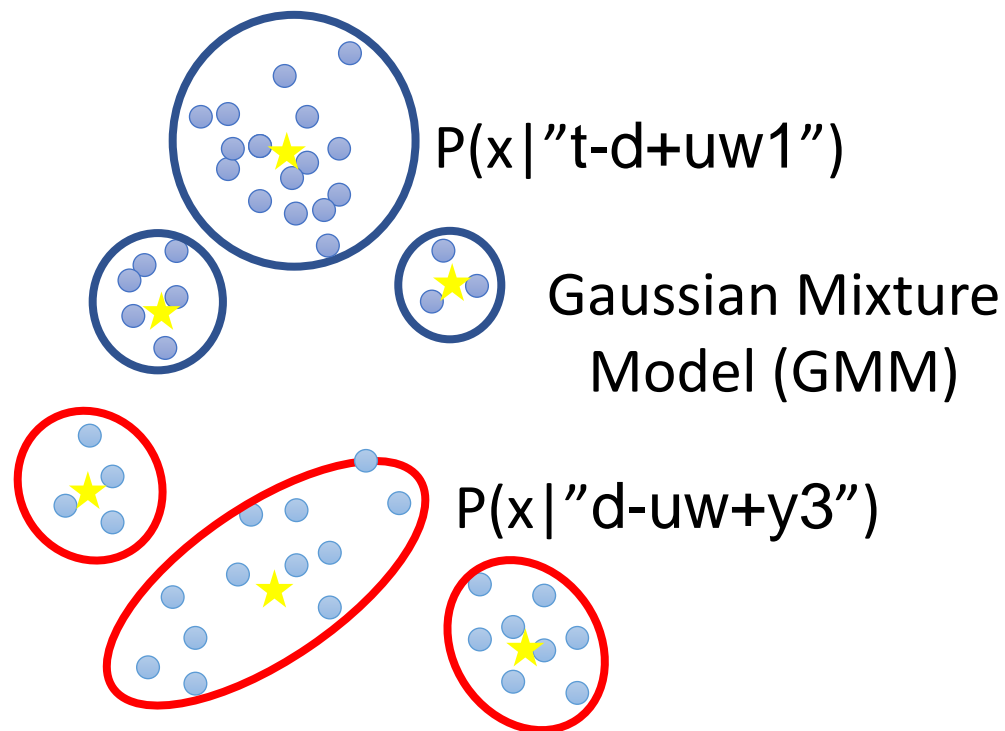
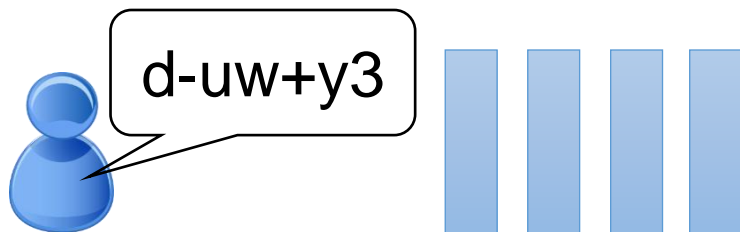
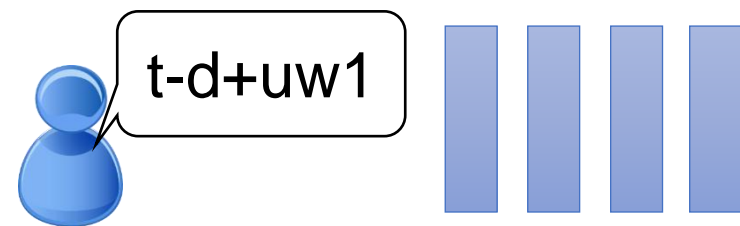
A sentence Y corresponds to a sequence of **states** S

Transition Probability

Probability from one state to another

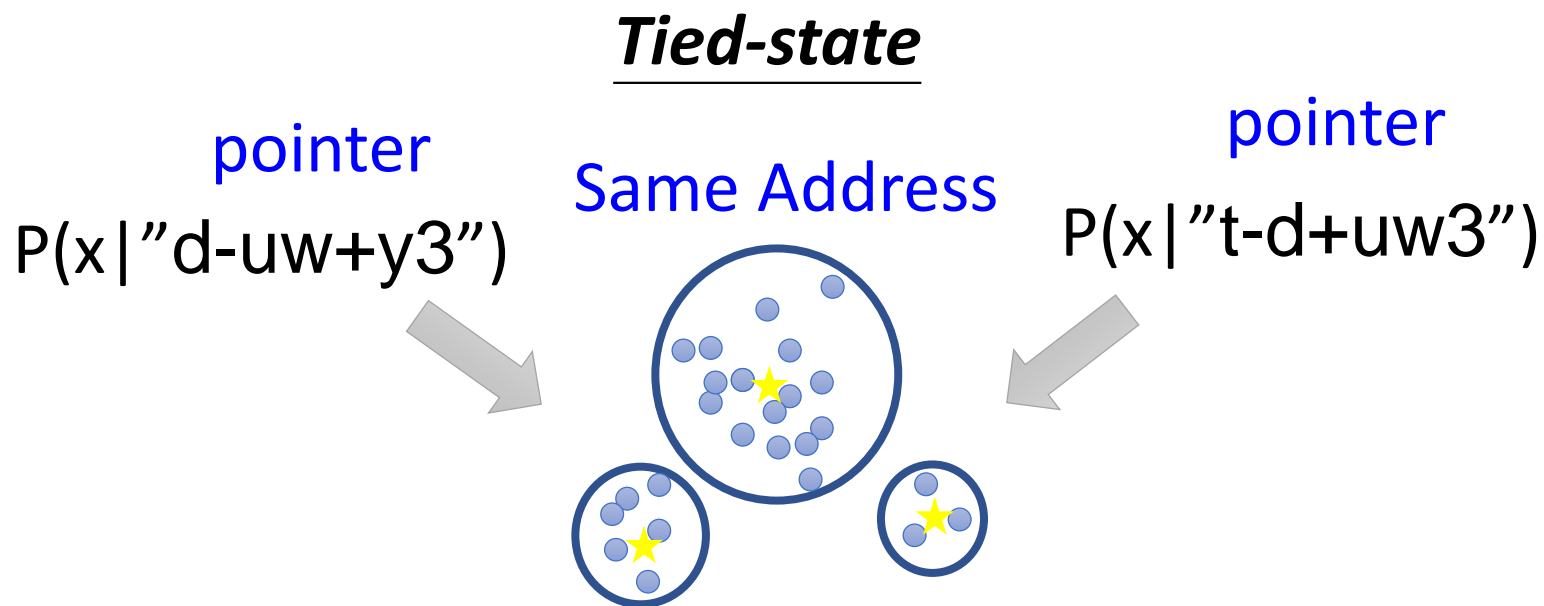
$$a \longrightarrow b \quad p(b|a)$$

Emission Probability



HMM – Emission Probability

- Too many states



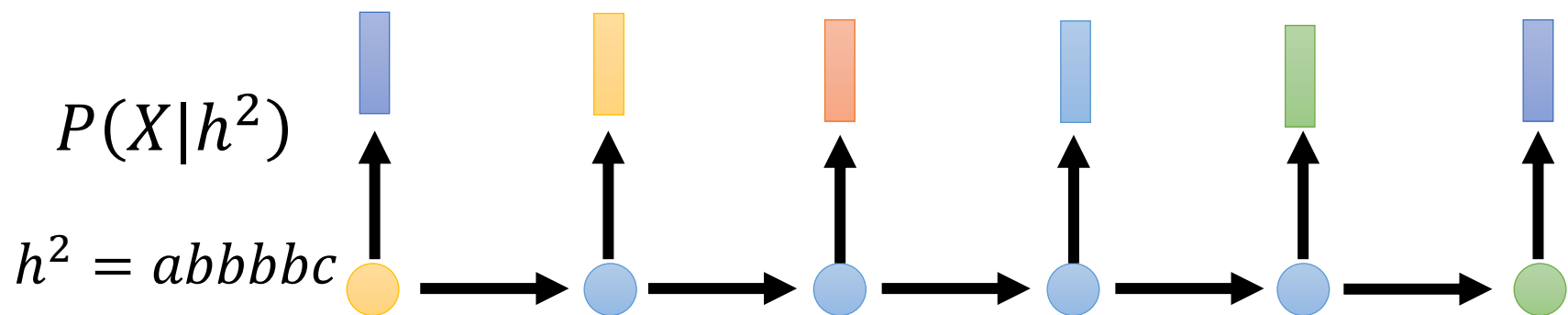
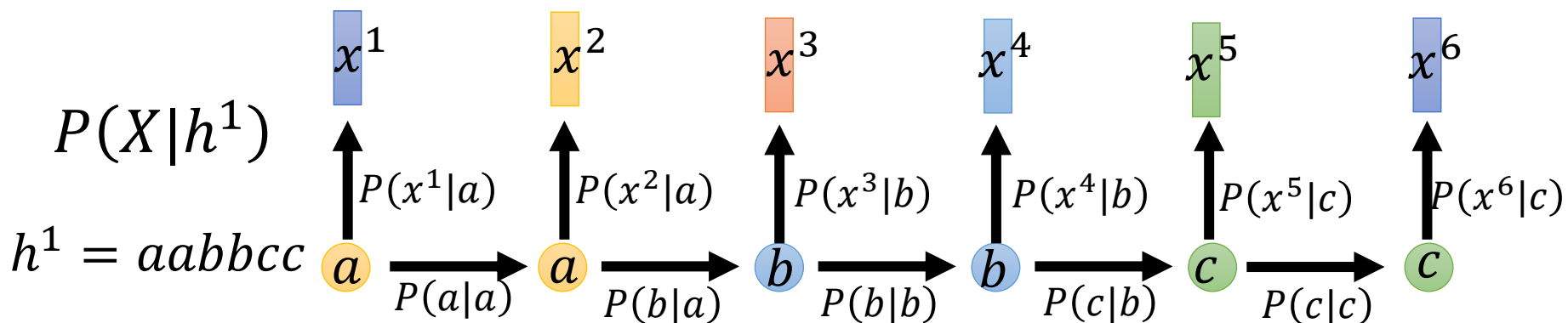
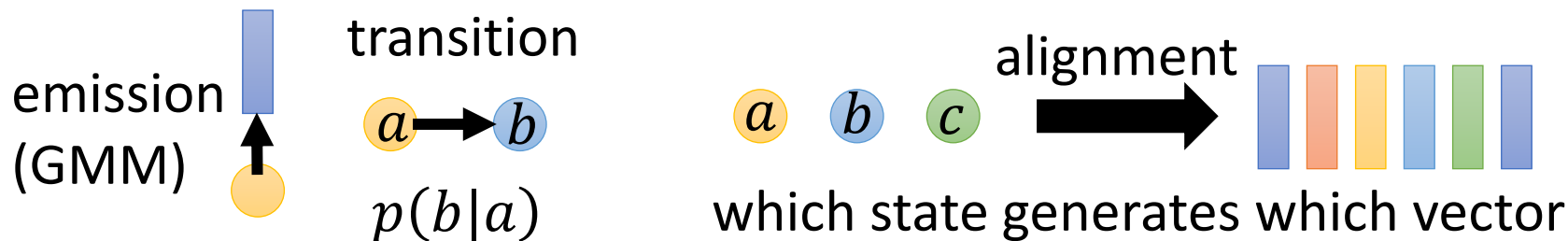
終極型態: Subspace GMM [Povey, et al., ICASSP'10]

(Geoffrey Hinton also published deep learning for ASR in the same conference)

[Mohamed , et al., ICASSP'10]

$$P_{\theta}(X|S) =? \sum_{h \in \text{align}(S)} P(X|h)$$

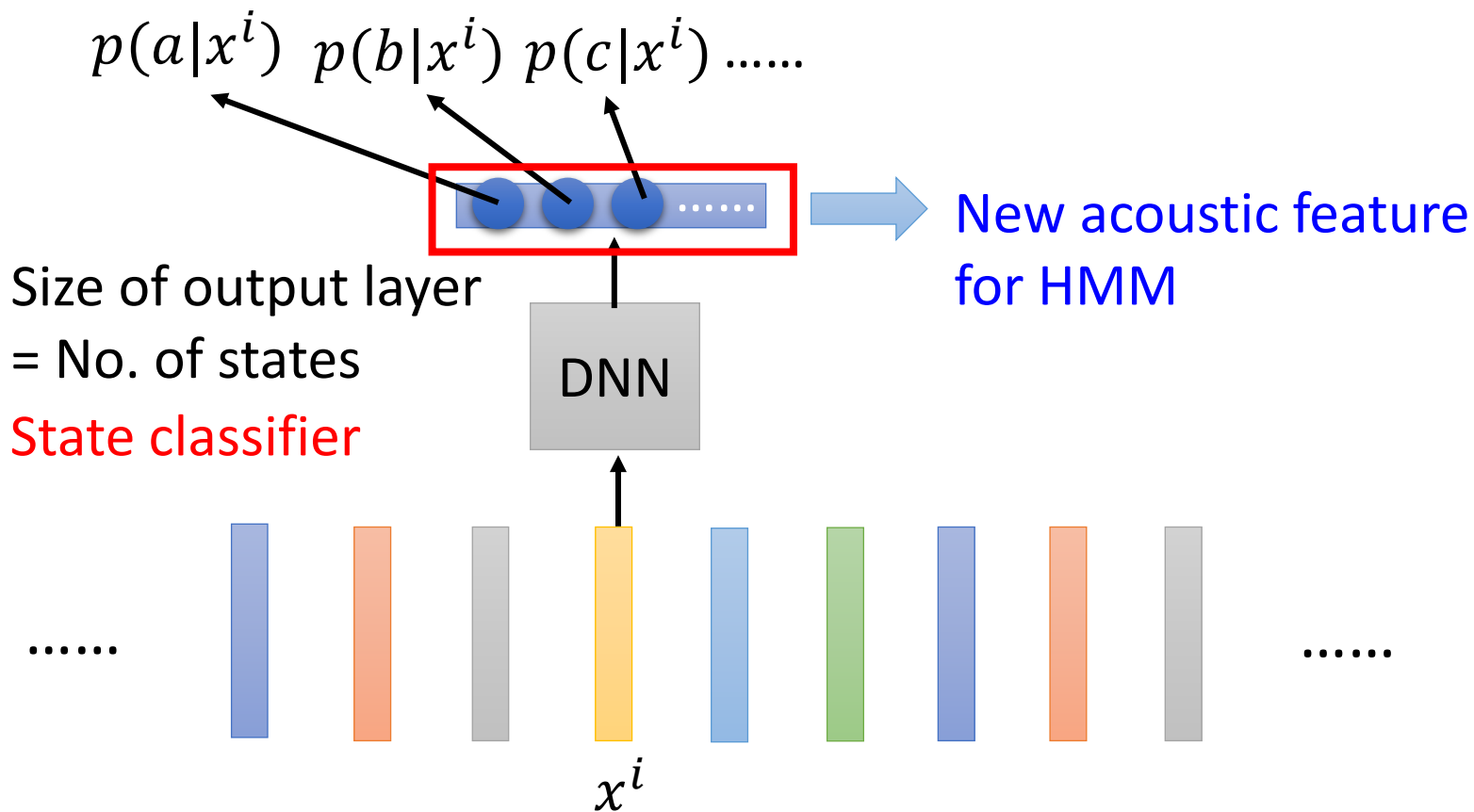
$h = abccbc$ ❌
 $h = abbbb$ ❌



A person is shown from behind, wearing a dark long-sleeved shirt, with their hands pressed against their head in a gesture of frustration or deep thought. The background is a light gray wall with several hand-drawn lightbulbs in various colors (green, blue, red, orange, purple) and the word "idea" written in cursive next to each. The lightbulbs have radiating lines around them, suggesting they are glowing. The overall scene is dimly lit, with the lightbulbs providing the primary source of illumination.

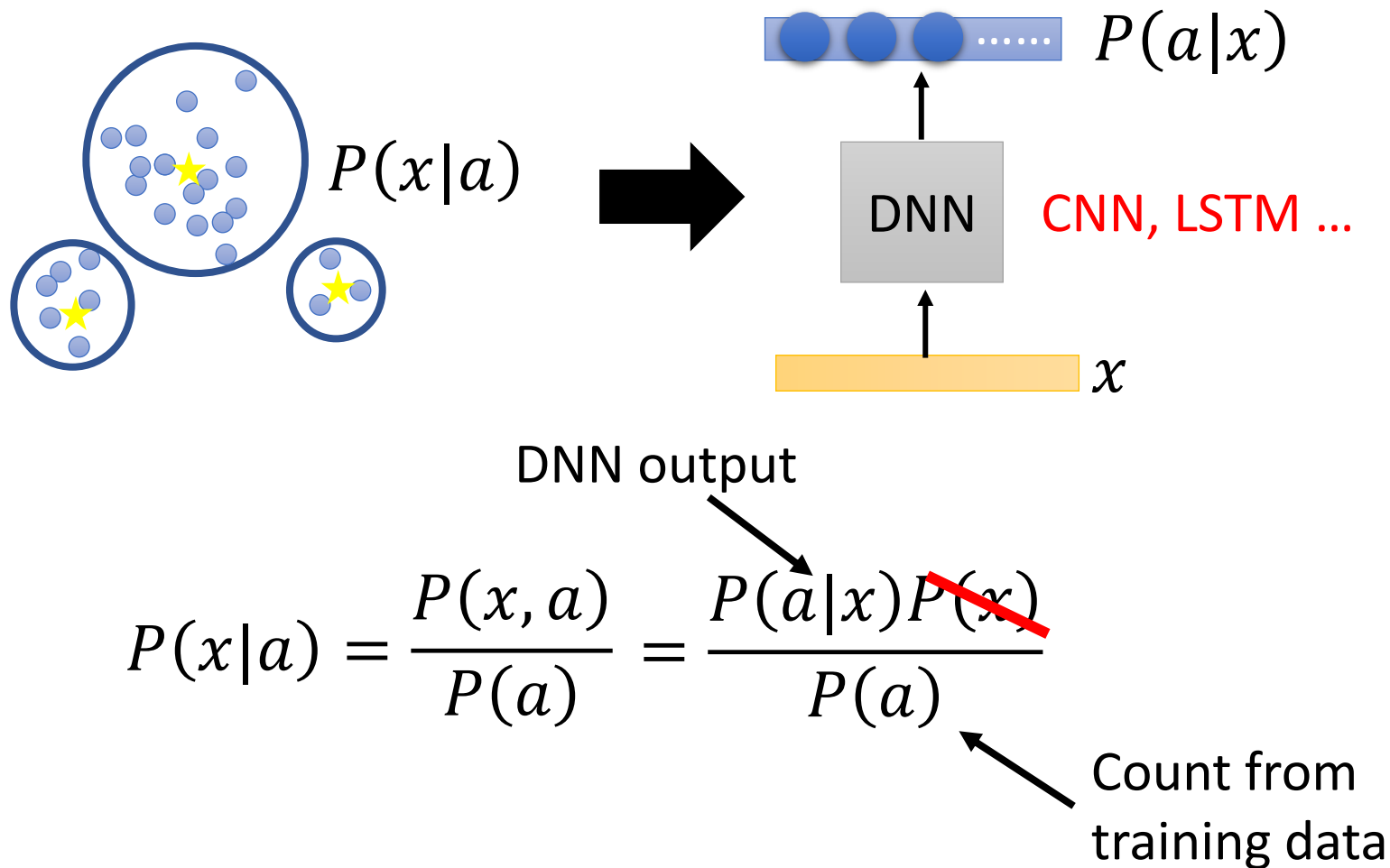
How to use Deep Learning?

Method 1: Tandem

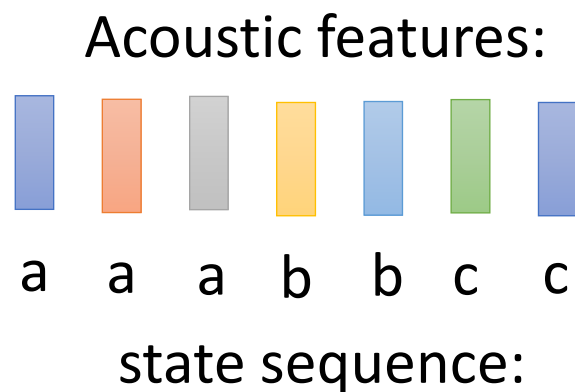
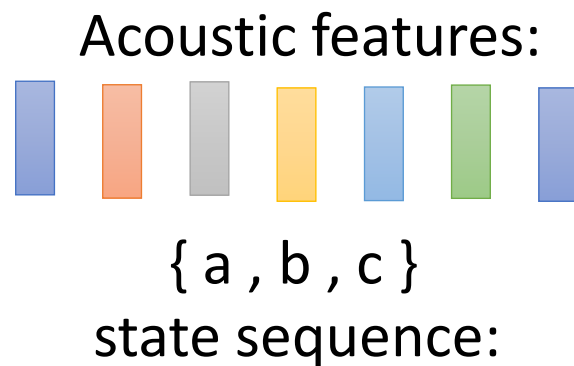
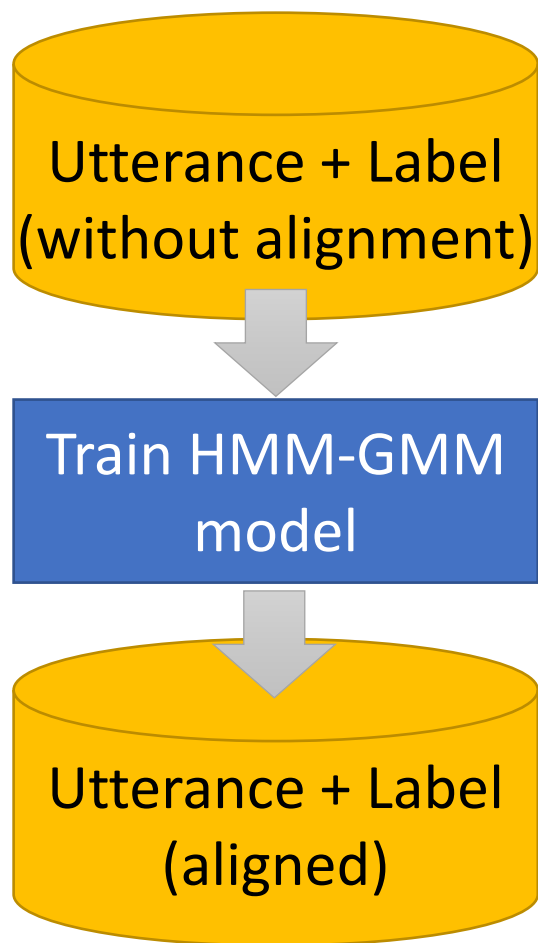


Last hidden layer or bottleneck layer are also possible.

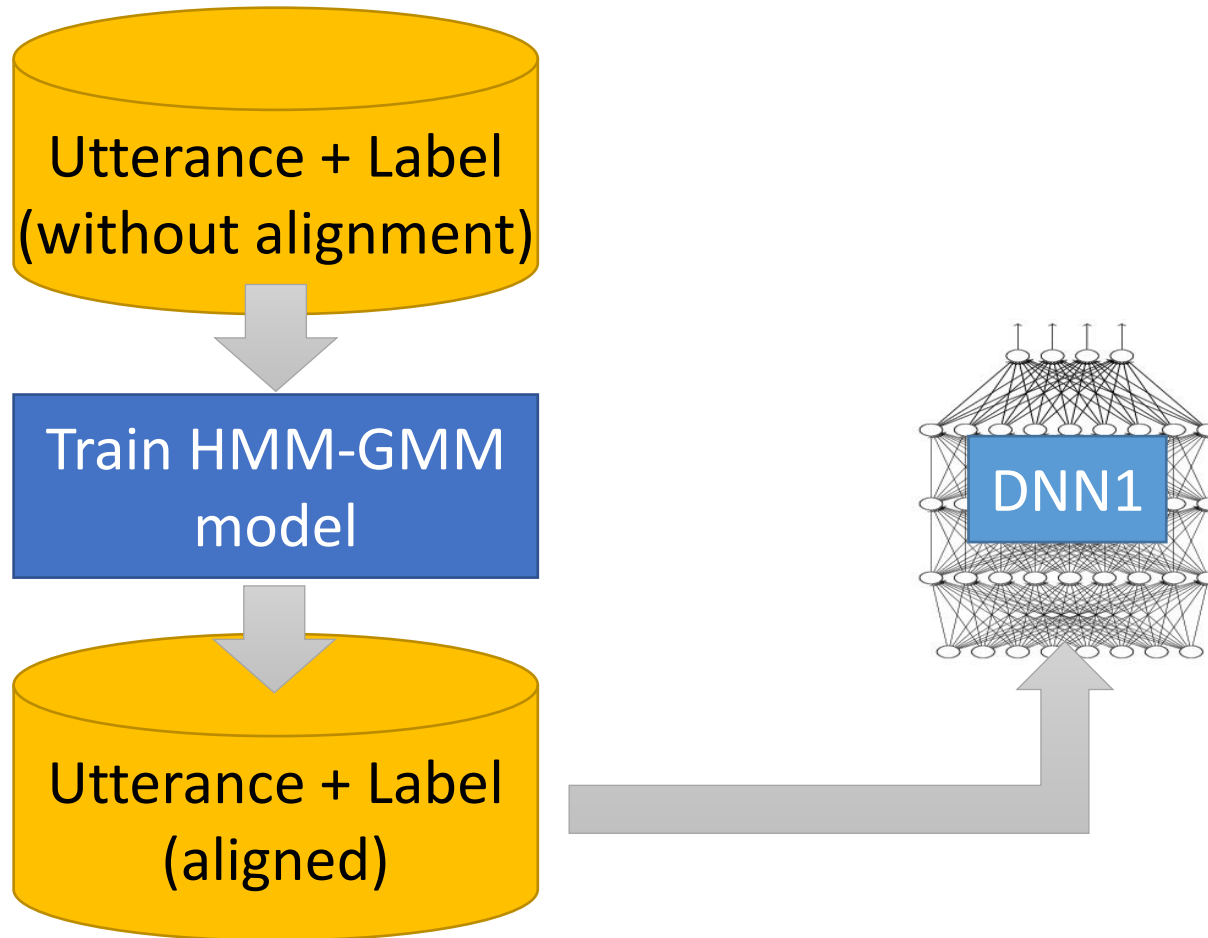
Method 2: DNN-HMM Hybrid



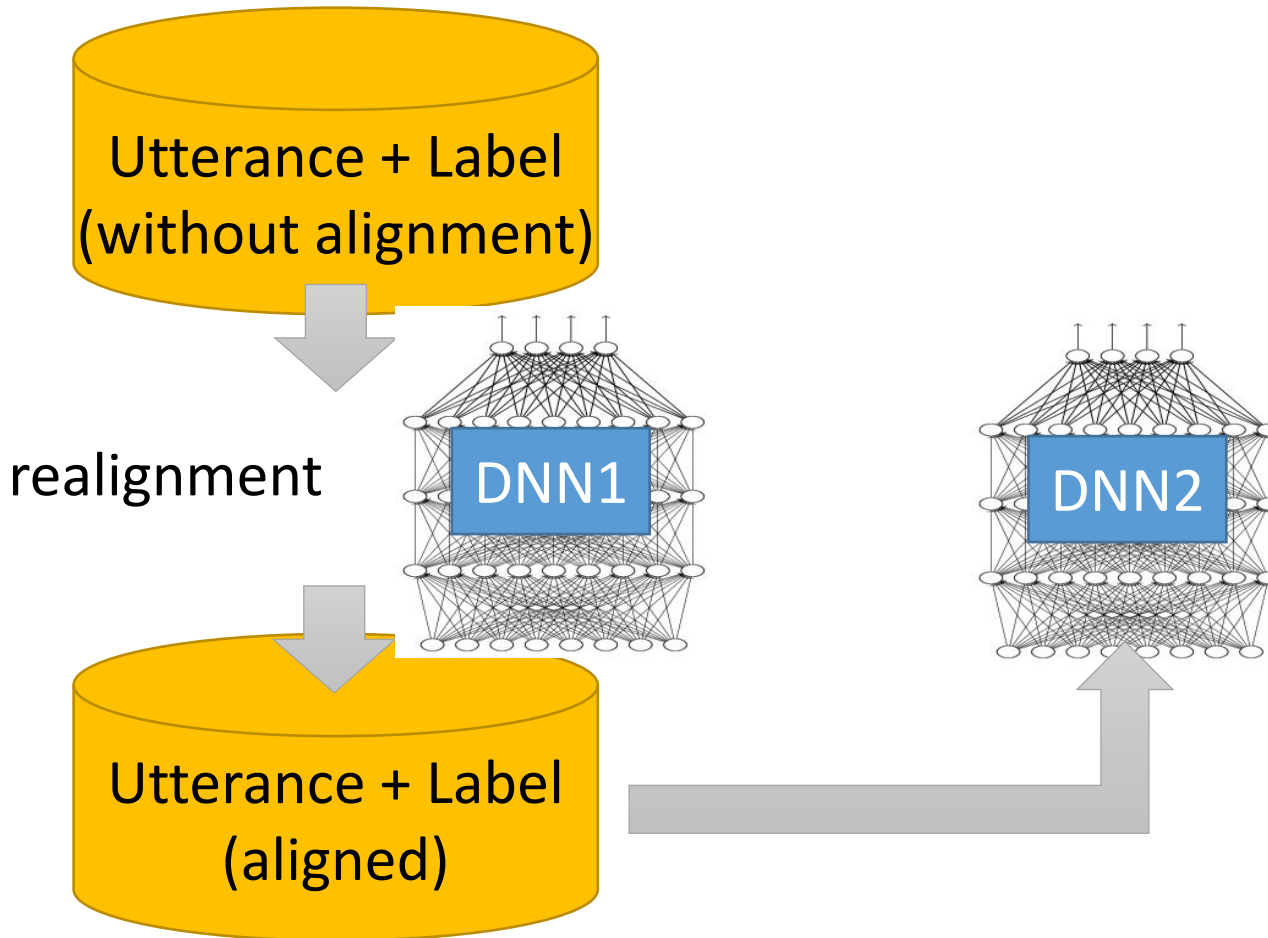
How to train a state classifier?



How to train a state classifier?



How to train a state classifier?



Human Parity!

- 微軟語音辨識技術突破重大里程碑：對話辨識能力達人類水準！(2016.10)

- <https://www.bnext.com.tw/article/41414/bn-2016-10-19-020437-216>

Machine 5.9% v.s. Human 5.9%

[Yu, et al., INTERSPEECH'16]

- IBM vs Microsoft: 'Human parity' speech recognition record changes hands again (2017.03)

- <http://www.zdnet.com/article/ibm-vs-microsoft-human-parity-speech-recognition-record-changes-hands-again/>

Machine 5.5% v.s. Human 5.1%

[Saon, et al., INTERSPEECH'17]

Very Deep

VGG Net (85M Parameters)	Residual-Net (38M Parameters)	LACE (65M Parameters)
14 weight layers	49 weight layers	22 weight layers
40x41 input	40x41 input	40x61 input
3 – conv 3x3, 96	3 – [conv 1x1, 64 conv 3x3, 64 conv 1x1, 256]	5 – conv 3x3, 128
Max pool	4 – [conv 1x1, 128 conv 3x3, 128 conv 1x1, 512]	5 – conv 3x3, 256
4 – conv 3x3, 192	6 – [conv 1x1, 256 conv 3x3, 256 conv 1x1, 1024]	5 – conv 3x3, 512
Max pool	3 – [conv 1x1, 512 conv 3x3, 512 conv 1x1, 2048]	5 – conv 3x3, 1024
4 – conv 3x3, 384	Average pool	1 – conv 3x4, 1
Max pool	Softmax (9000)	Softmax (9000)
2 – FC – 4096		
Softmax (9000)		

[Yu, et al., INTERSPEECH'16]

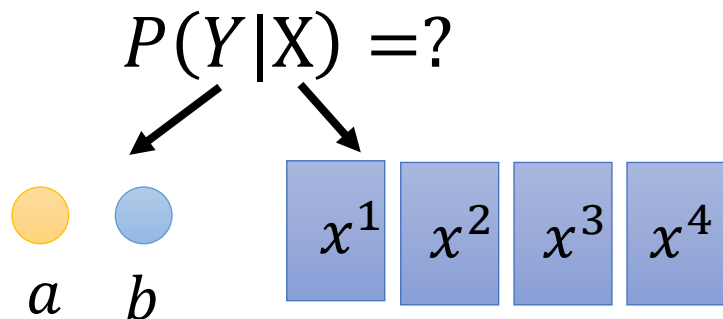


Back to End-to-end

LAS

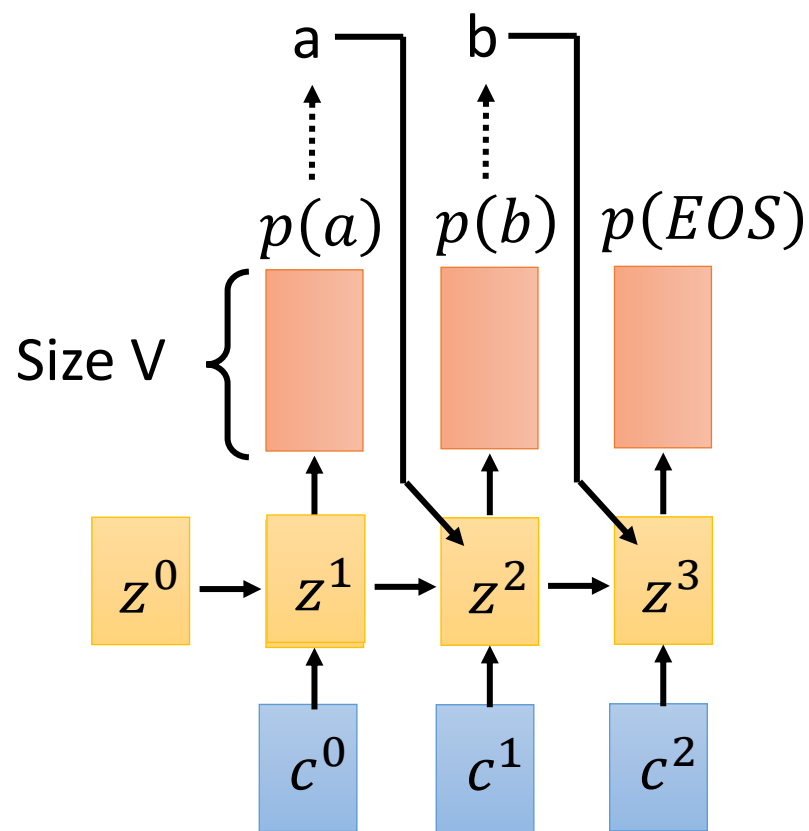
Decoding: $Y^* = \underset{Y}{\operatorname{arg\,max}} \log P(Y|X)$
Beam Search

Training: $\theta^* = \underset{\theta}{\operatorname{arg\,max}} \log P_{\theta}(\hat{Y}|X)$



- LAS directly computes $P(Y|X)$

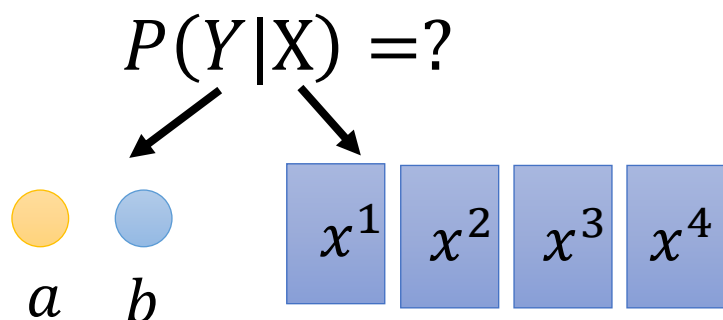
$$P(Y|X) = p(a|X)p(b|a, X)\dots$$



CTC, RNN-T

Decoding: $Y^* = \underset{Y}{\operatorname{arg\,max}} \log P(Y|X)$
Beam Search

Training: $\theta^* = \underset{\theta}{\operatorname{arg\,max}} \log P_{\theta}(\hat{Y}|X)$

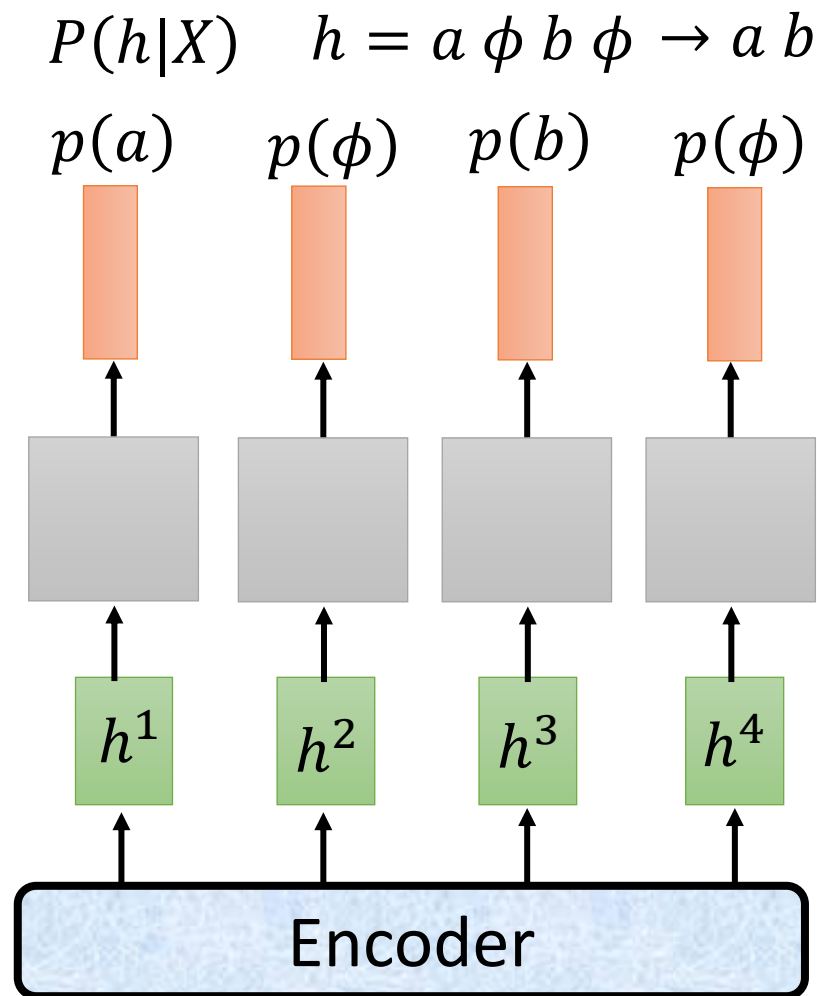


- LAS directly computes $P(Y|X)$

$$P(Y|X) = p(a|X)p(b|a, X)\dots$$

- CTC and RNN-T need **alignment**

$$P(Y|X) = \sum_{h \in \operatorname{align}(Y)} P(h|X)$$



HMM, CTC, RNN-T

HMM

$$P_{\theta}(X|S) = \sum_{h \in \text{align}(S)} P(X|h)$$

CTC, RNN-T

$$P_{\theta}(Y|X) = \sum_{h \in \text{align}(Y)} P(h|X)$$

1. Enumerate all the possible alignments
2. How to sum over all the alignments

3. Training: $\theta^* = \arg \max_{\theta} \log P_{\theta}(\hat{Y}|X)$ $\frac{\partial P_{\theta}(\hat{Y}|X)}{\partial \theta} = ?$

4. Testing (Inference, decoding):

$$Y^* = \arg \max_Y \log P(Y|X)$$

HMM, CTC, RNN-T

HMM

$$P(X|S) = \sum_{h \in \text{align}(S)} P(X|h)$$

CTC, RNN-T

$$P(Y|X) = \sum_{h \in \text{align}(Y)} P(h|X)$$

1. Enumerate all the possible alignments

2. How to sum over all the alignments

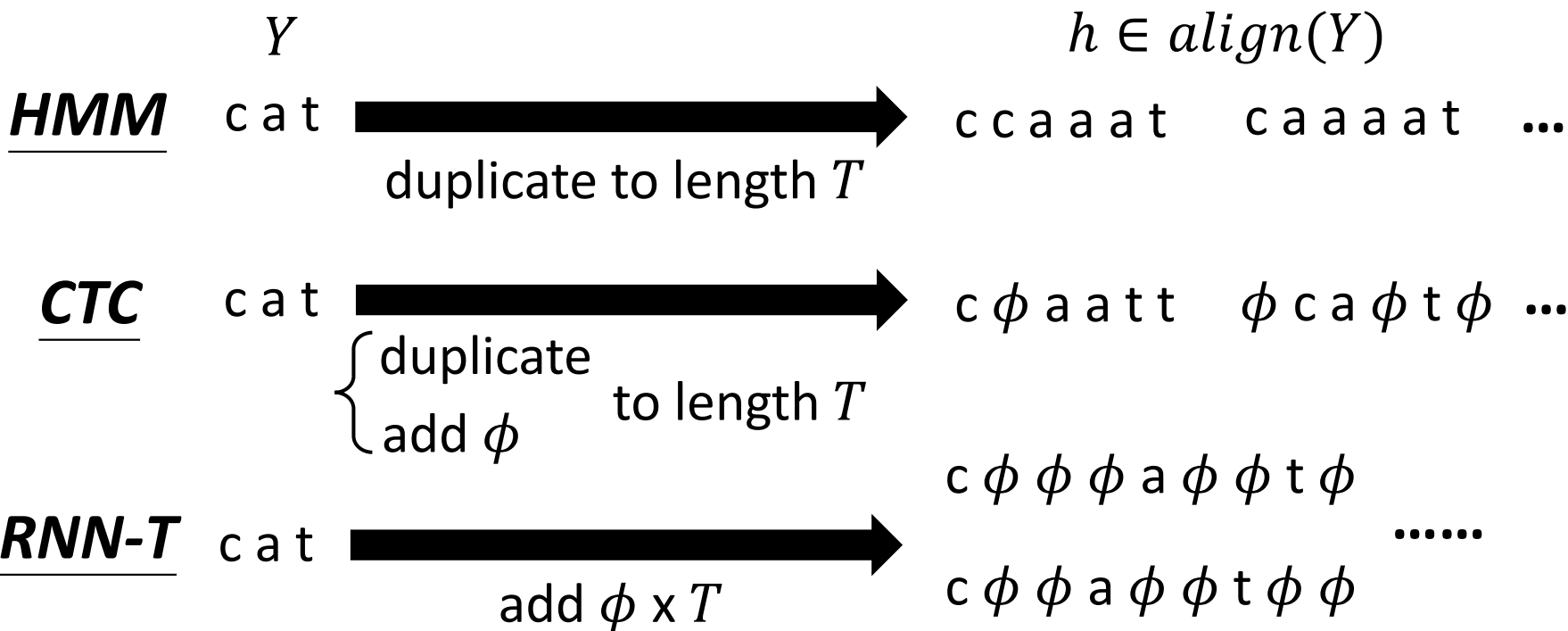
3. Training: $\theta^* = \arg \max_{\theta} \log P_{\theta}(\hat{Y}|X)$ $\frac{\partial P_{\theta}(\hat{Y}|X)}{\partial \theta} = ?$

4. Testing (Inference, decoding):


$$Y^* = \arg \max_Y \log P(Y|X)$$

All the alignments

你們在忙什麼 😊



HMM

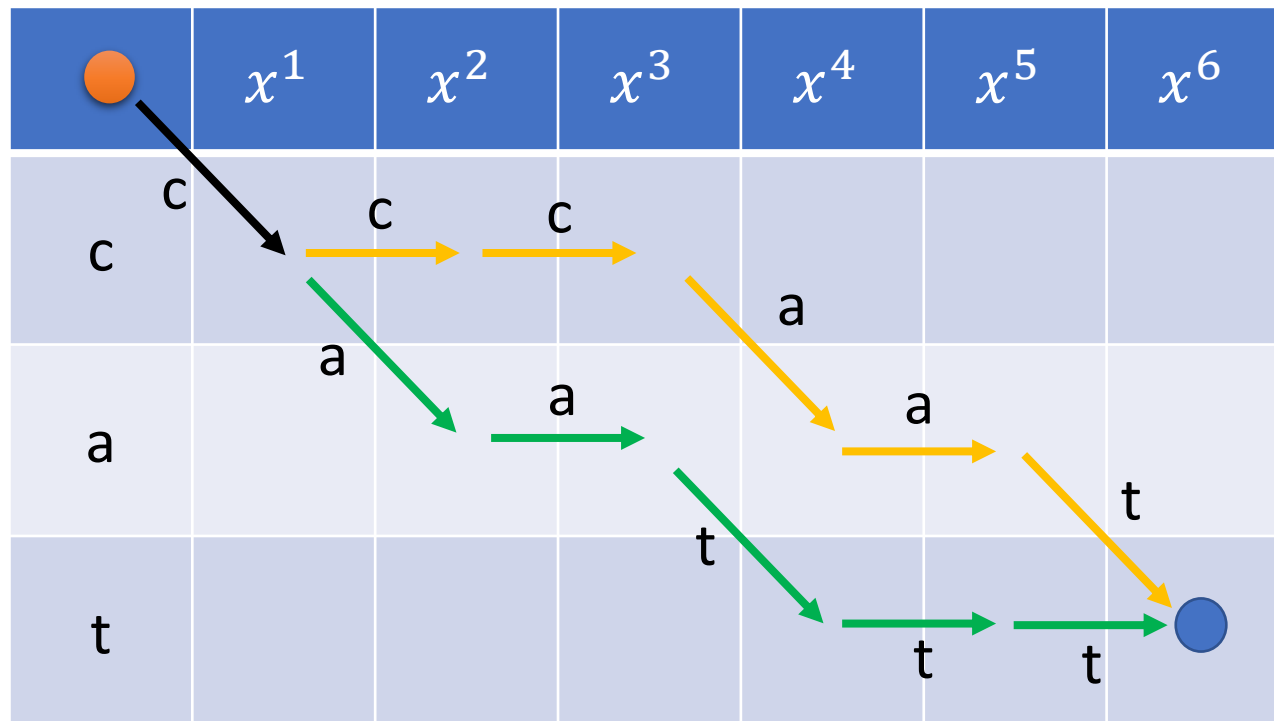
cat  ccaaat caaat ...
duplicate to length T

For $n = 1$ to N

output the n -th token t_n times

constraint: $t_1 + t_2 + \dots + t_N = T, t_n > 0$


Trellis Graph



 duplicate

 next token

HMM

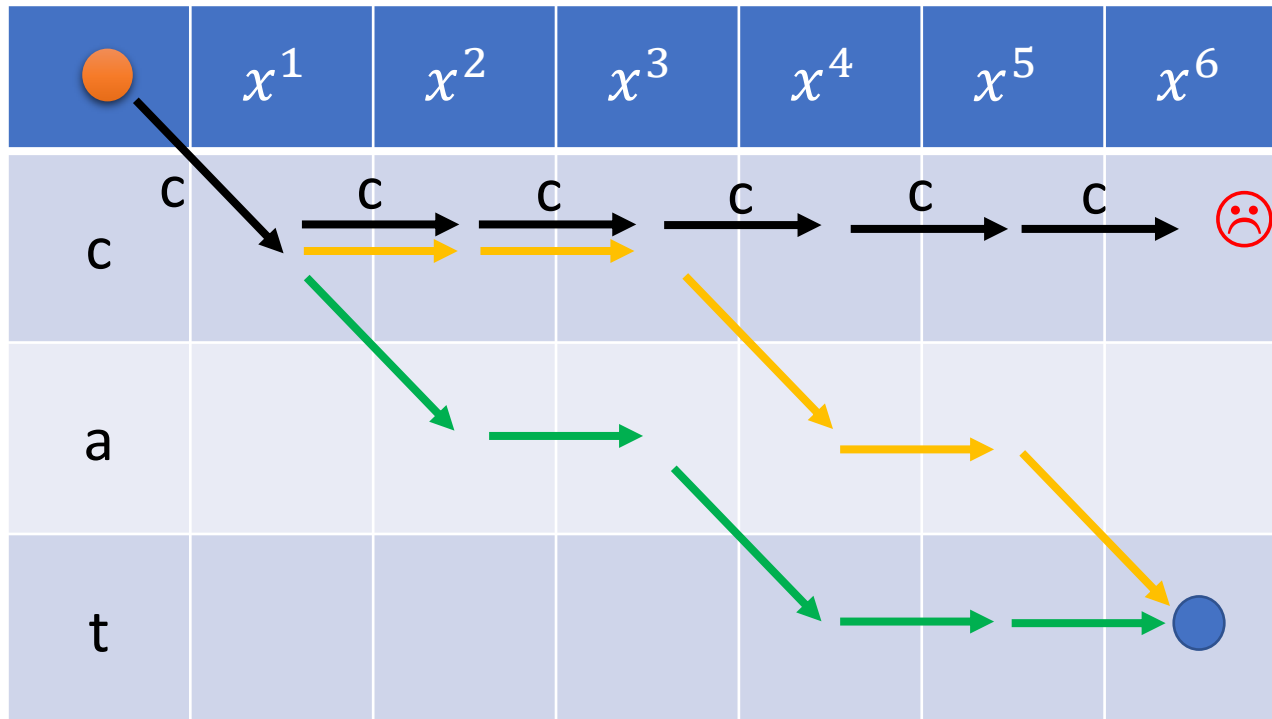
cat  ccaaat caaat ...
duplicate to length T

For $n = 1$ to N

output the n -th token t_n times

constraint: $t_1 + t_2 + \dots + t_N = T, t_n > 0$


Trellis Graph



 duplicate

 next token

CTC

c a t  c ϕ a a t t ϕ c a ϕ t ϕ ...
 { duplicate
 { add ϕ to length T

output " ϕ " c_0 times

For $n = 1$ to N

 output the n -th token t_n times

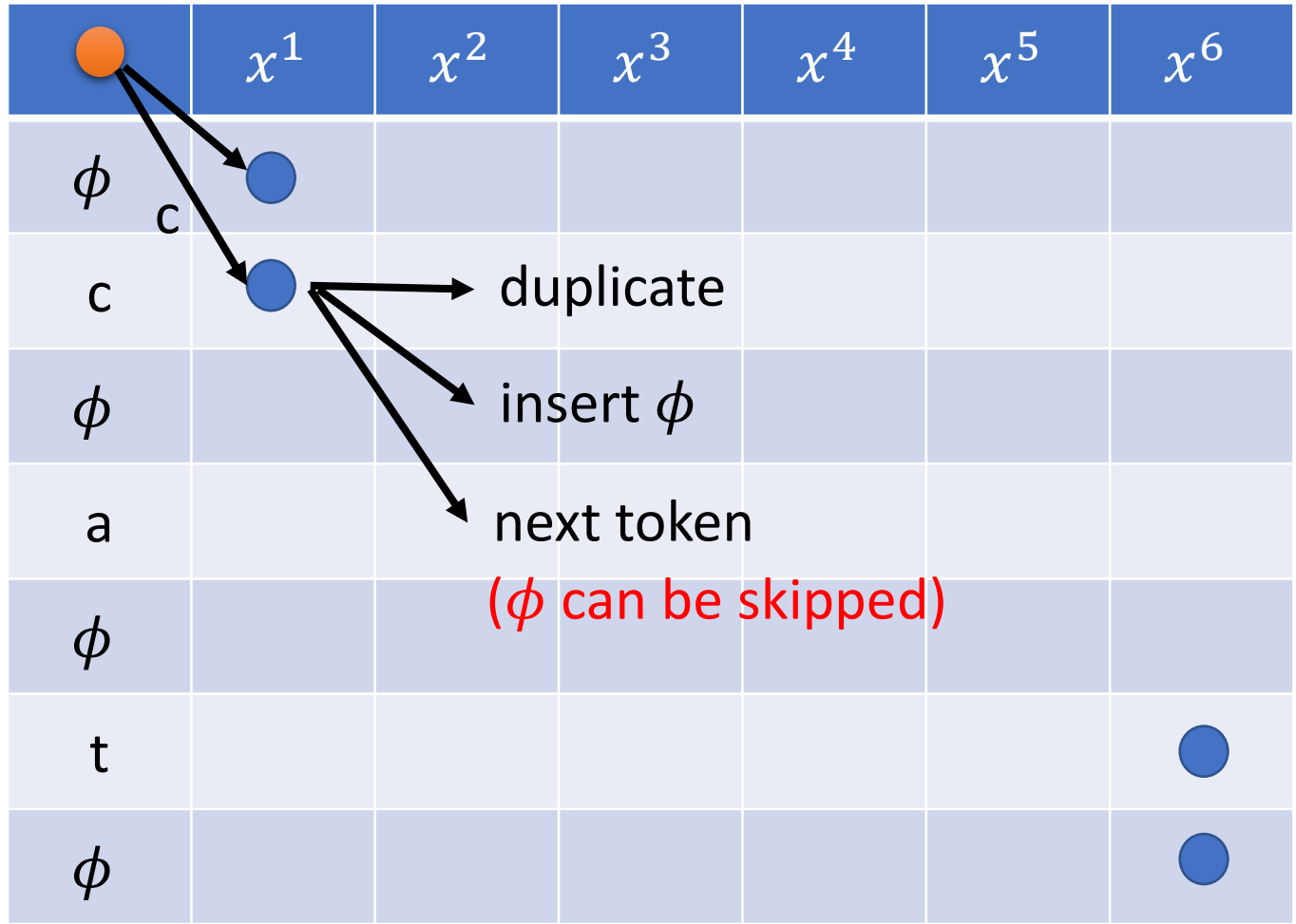
 output " ϕ " c_n times

constraint: $t_1 + t_2 + \dots + t_N +$
 $c_0 + c_1 + \dots + c_N = T$


$$t_n > 0 \quad c_n \geq 0$$

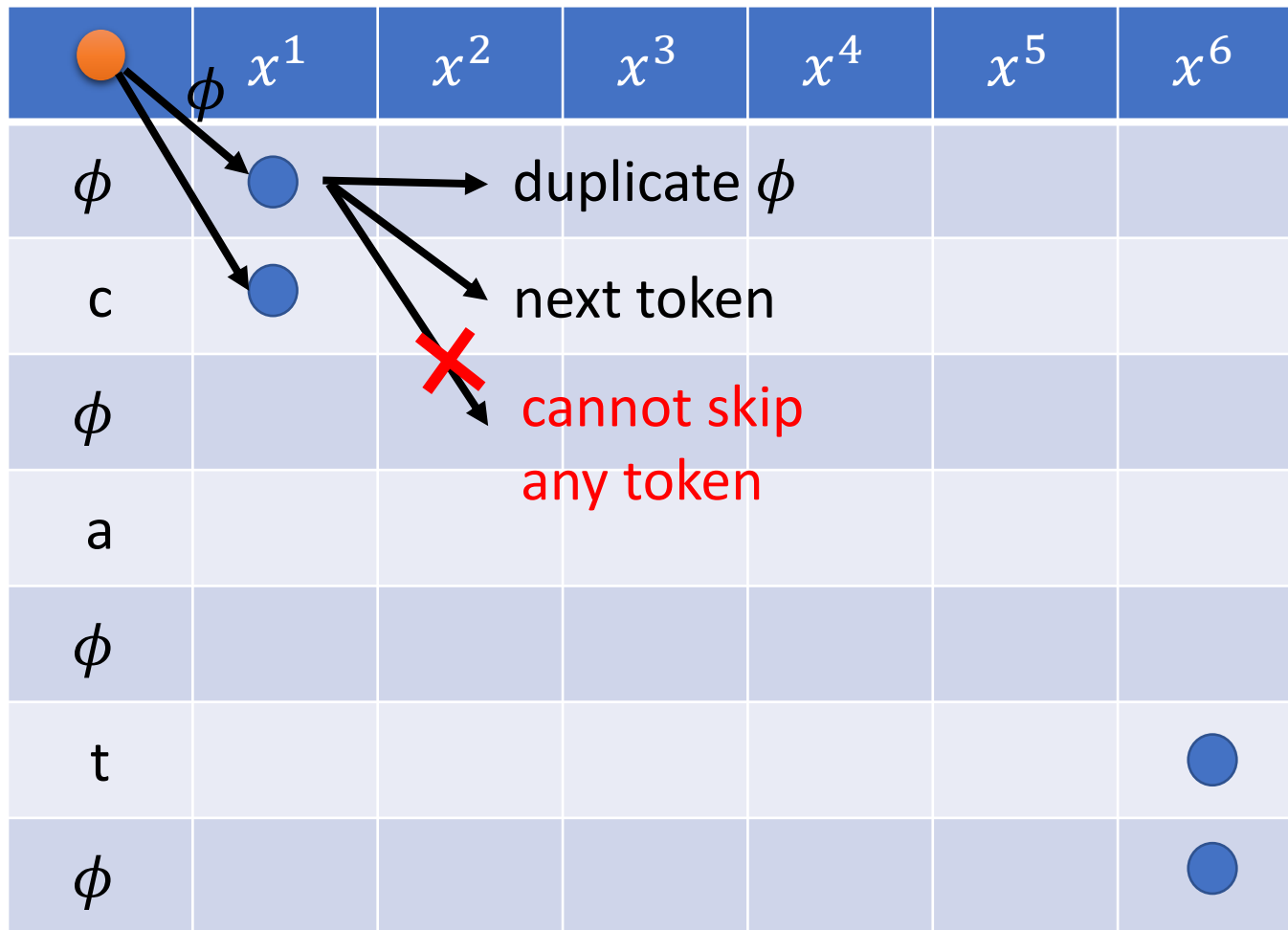
CTC

cat $\xrightarrow{\text{duplicate}} c\phi aatt \quad \phi ca\phi t\phi \dots$
 { add ϕ to length T



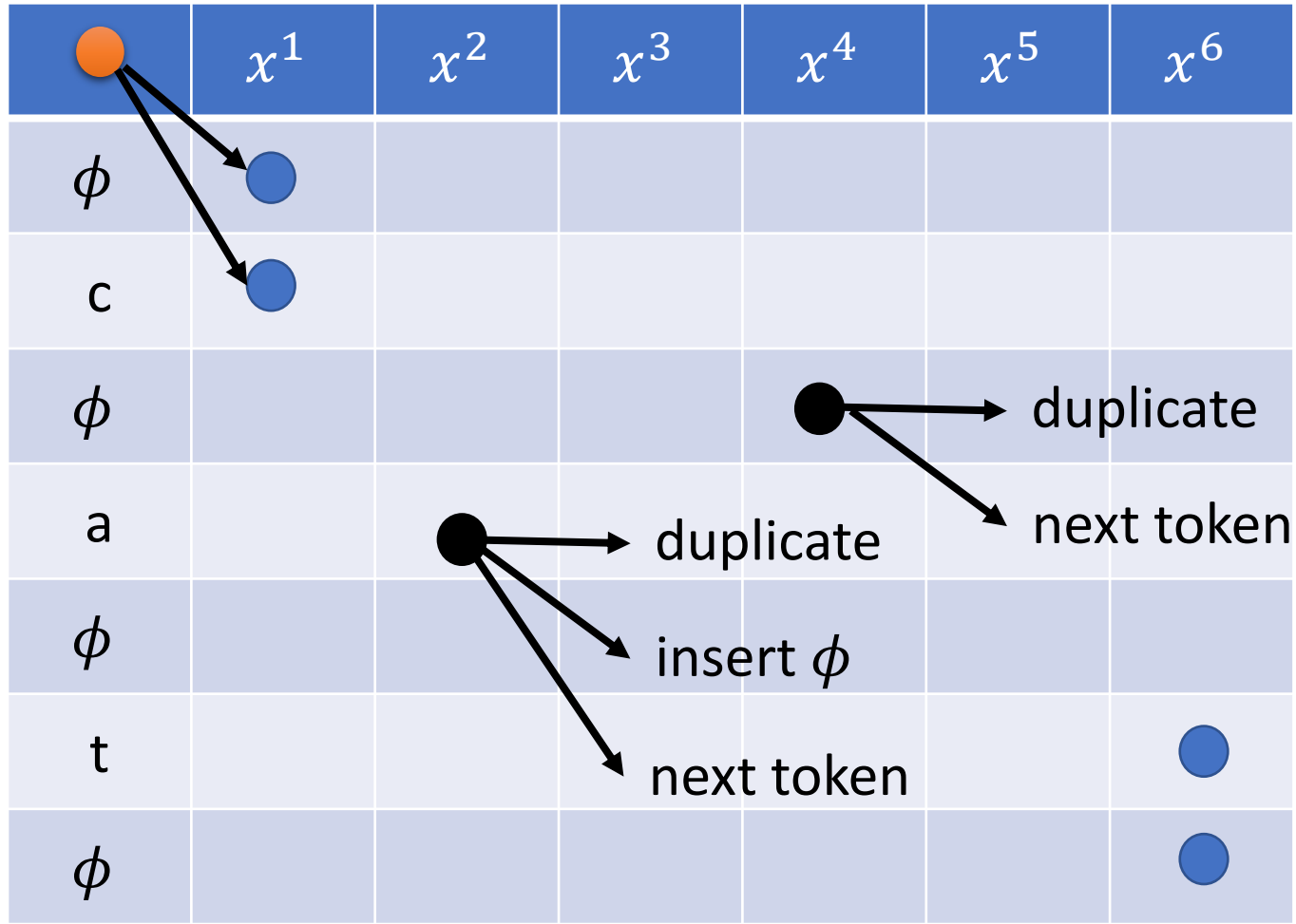
CTC

cat  c ϕ a a t t ϕ c a ϕ t ϕ ...
{ duplicate
add ϕ to length T



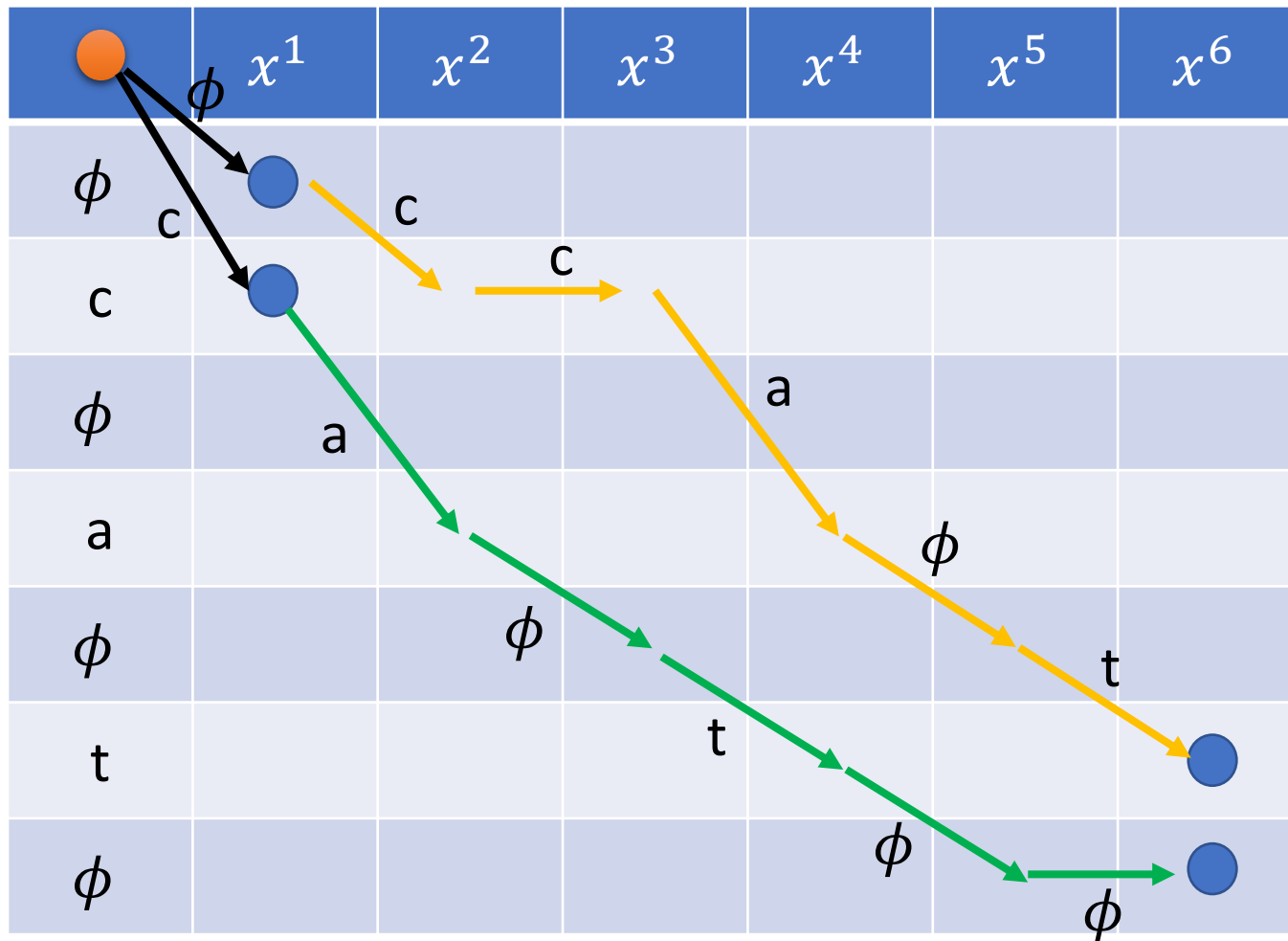
CTC

c a t $\xrightarrow{\text{duplicate}} c \phi a a t t \quad \phi c a \phi t \phi \dots$
{ duplicate
add ϕ to length T



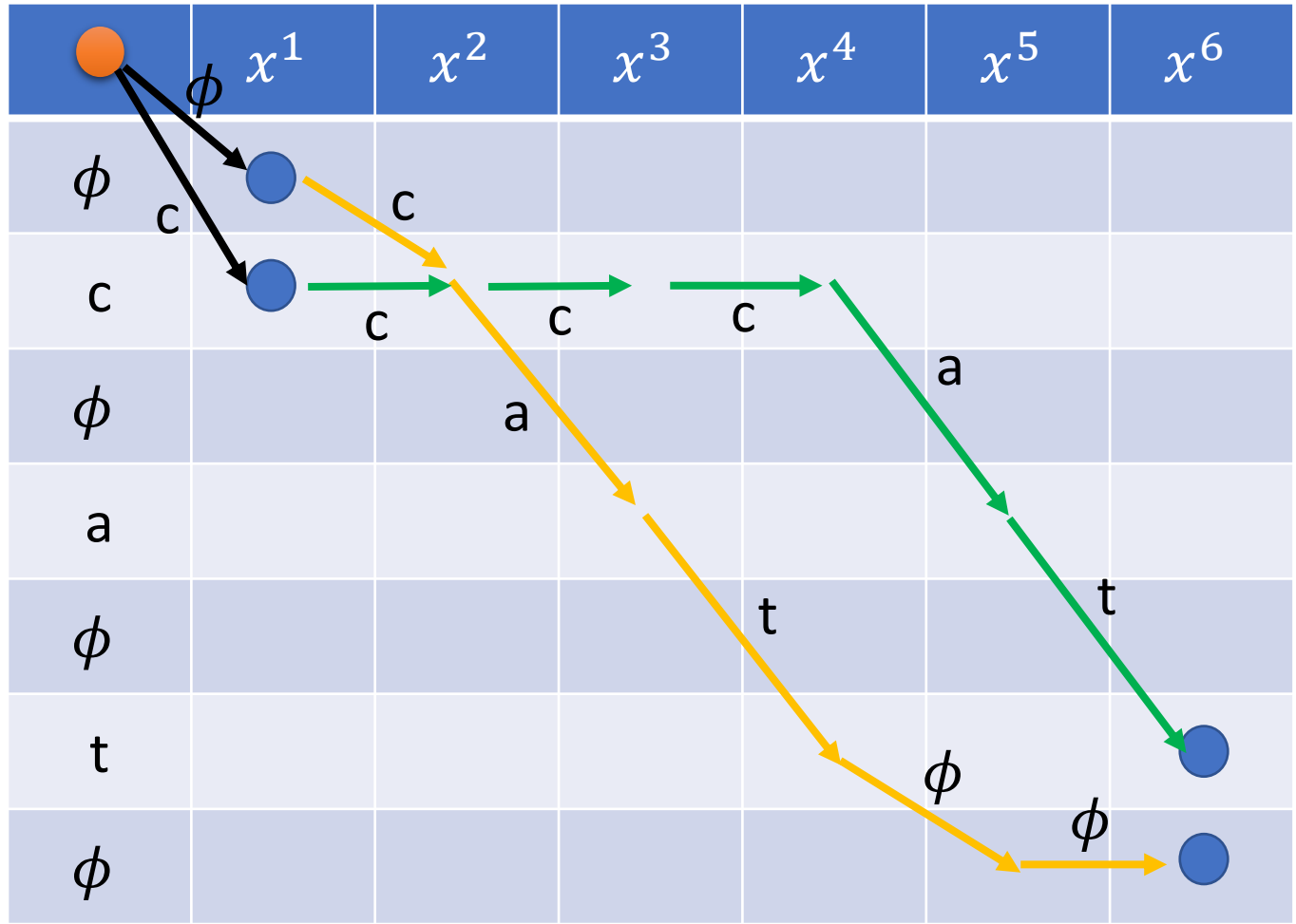
CTC

cat $\xrightarrow{\text{duplicate to length } T}$ $c\phi aatt \quad \phi ca\phi t\phi \dots$
duplicate to length T
add ϕ



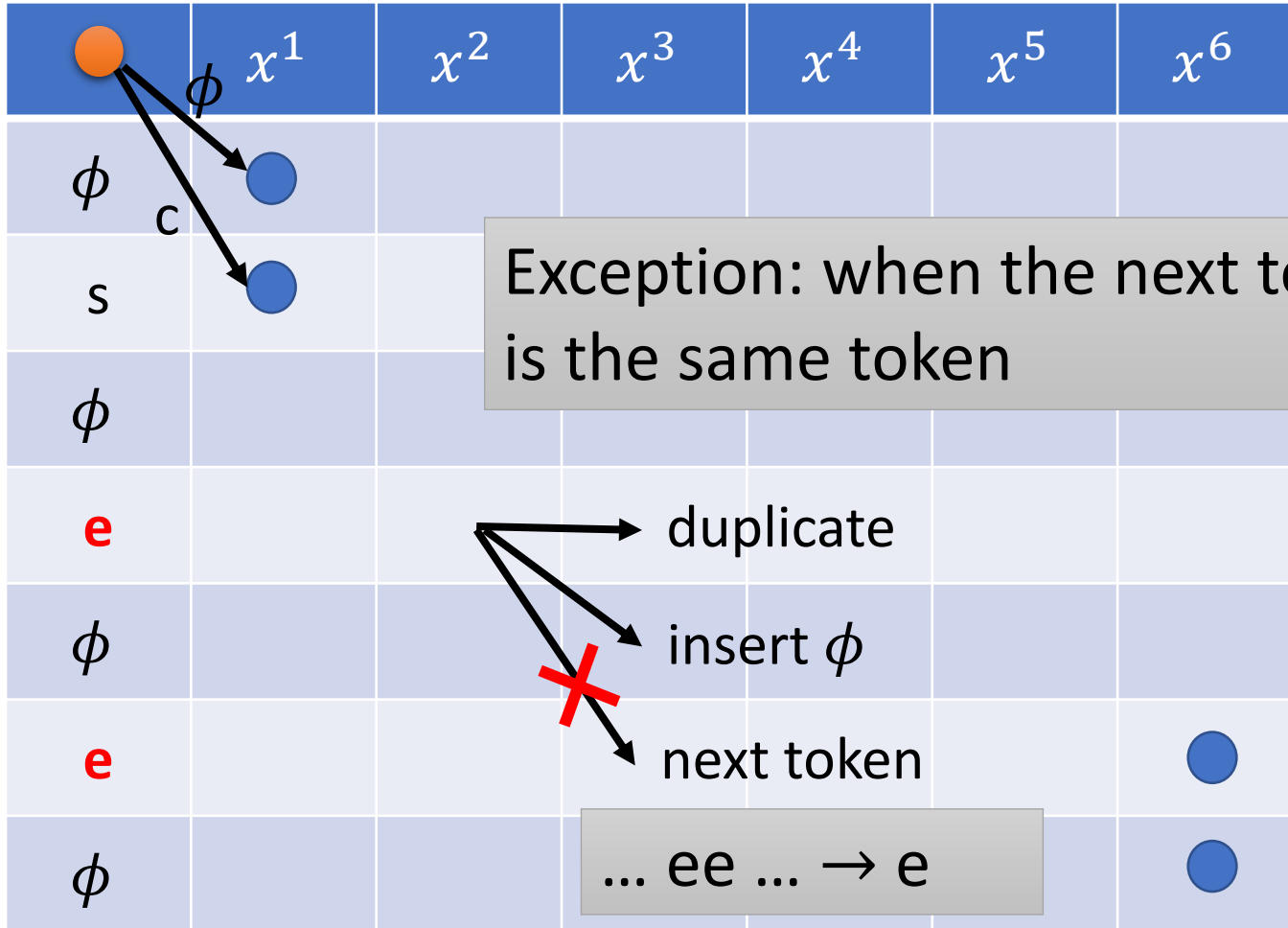
CTC

cat $\xrightarrow{\hspace{10em}}$ $c\phi aatt \quad \phi ca\phi t\phi \dots$
duplicate
add ϕ to length T



CTC

cat $\xrightarrow{\text{duplicate}} c\phi aatt \quad \phi ca\phi t\phi \dots$
{ duplicate
add ϕ to length T



RNN-T

cat

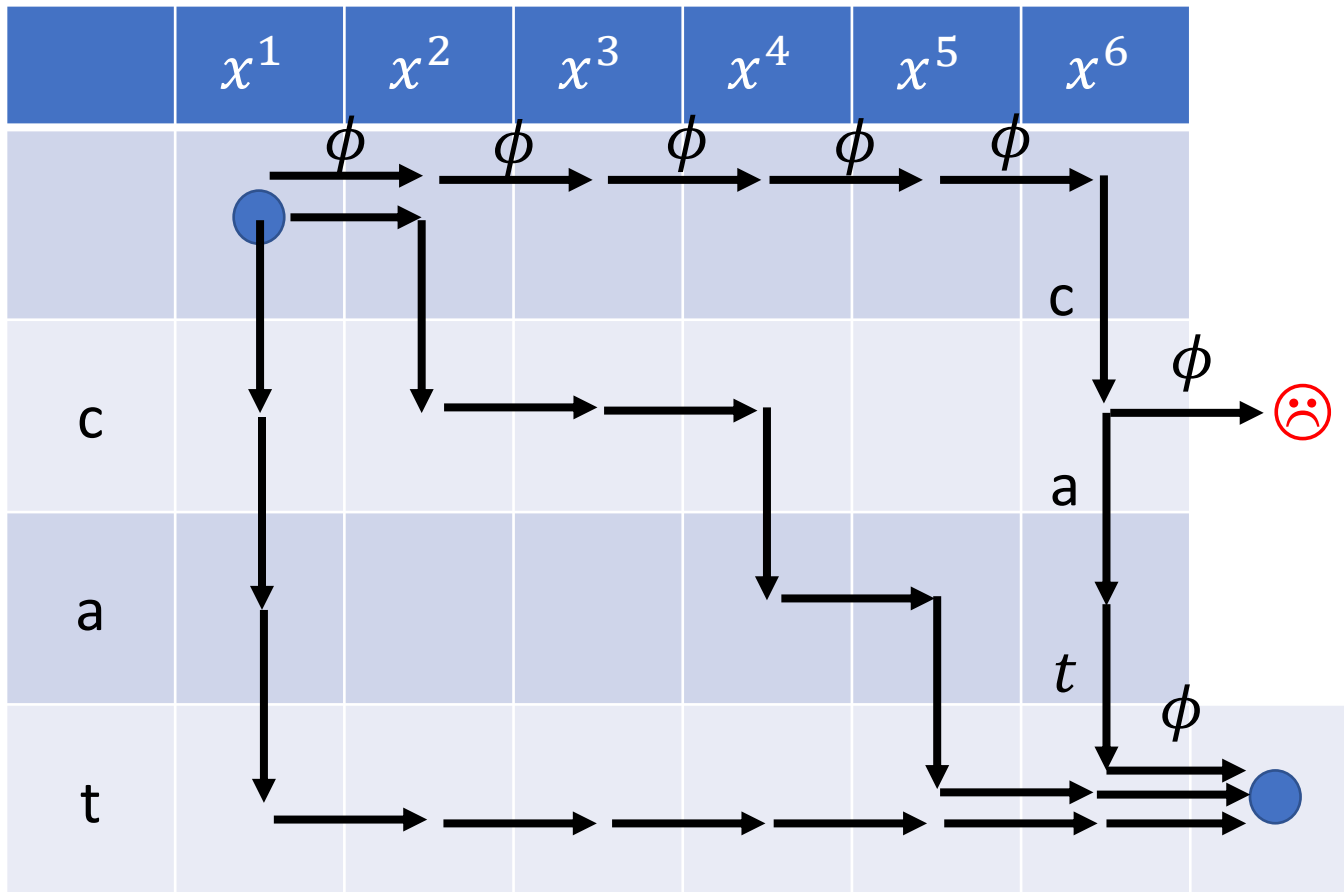


add $\phi \times T$

c ϕ ϕ ϕ a ϕ ϕ t ϕ

.....

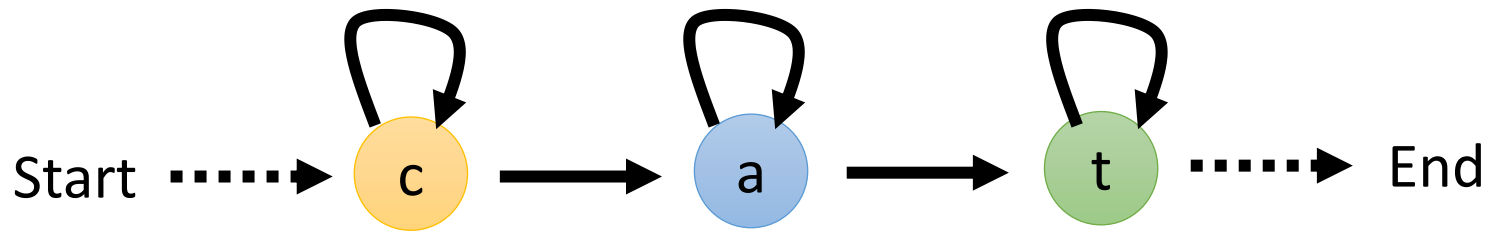
c ϕ ϕ a ϕ ϕ t ϕ ϕ



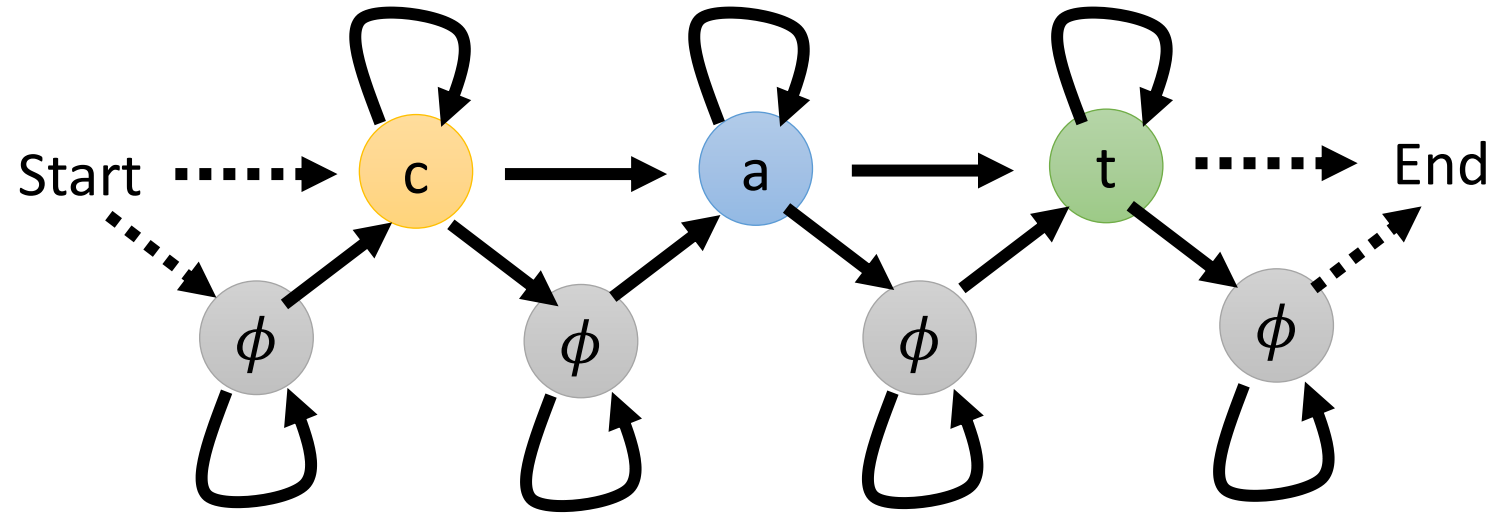
→ Insert ϕ

↓ output token

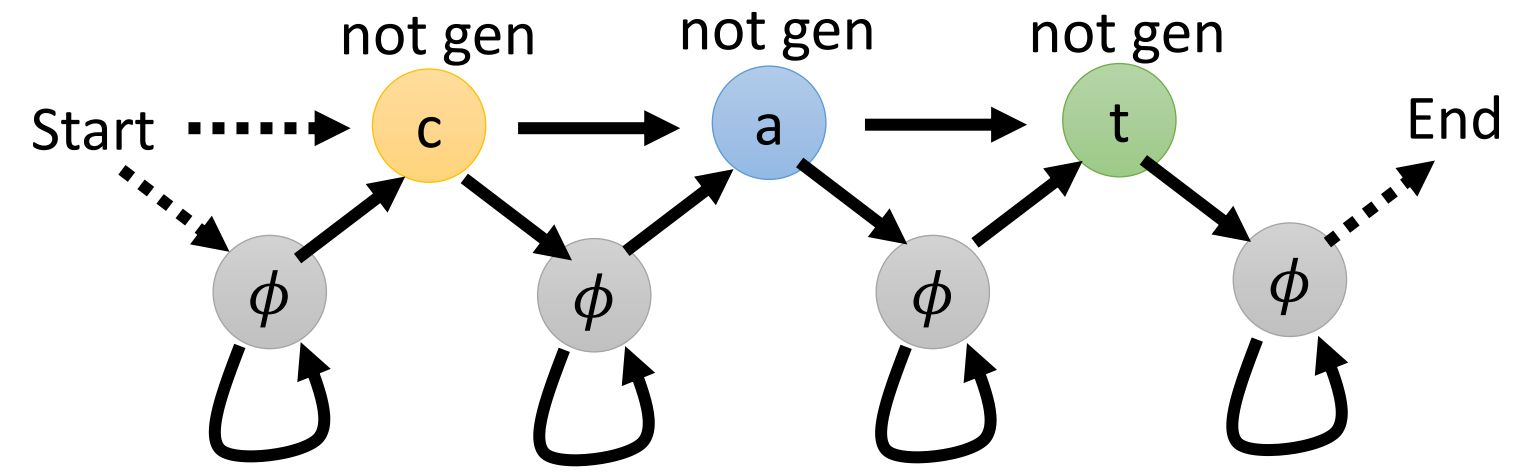
HMM



CTC



RNN-T



HMM, CTC, RNN-T

HMM

$$P(X|Y) = \sum_{h \in \text{align}(Y)} P(X|h)$$

CTC, RNN-T

$$P(Y|X) = \sum_{h \in \text{align}(Y)} P(h|X)$$

1. Enumerate all the possible alignments

2. How to sum over all the alignments

3. Training:

$$\theta^* = \arg \max_{\theta} \log P_{\theta}(\hat{Y}|X) \quad \frac{\partial P_{\theta}(\hat{Y}|X)}{\partial \theta} = ?$$

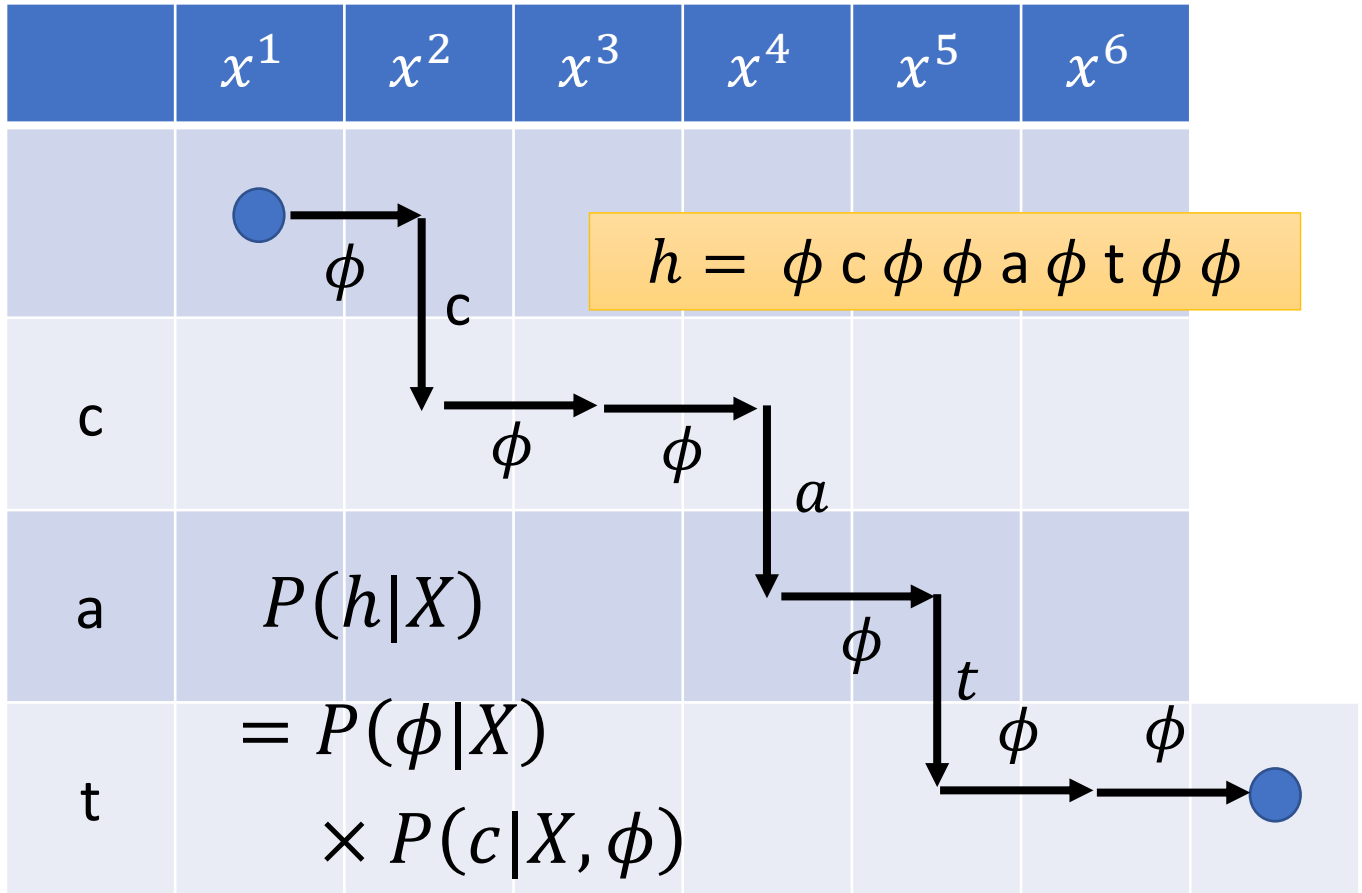
4. Testing (Inference, decoding):

$$Y^* = \arg \max_Y \log P(Y|X)$$



This part is challenging.

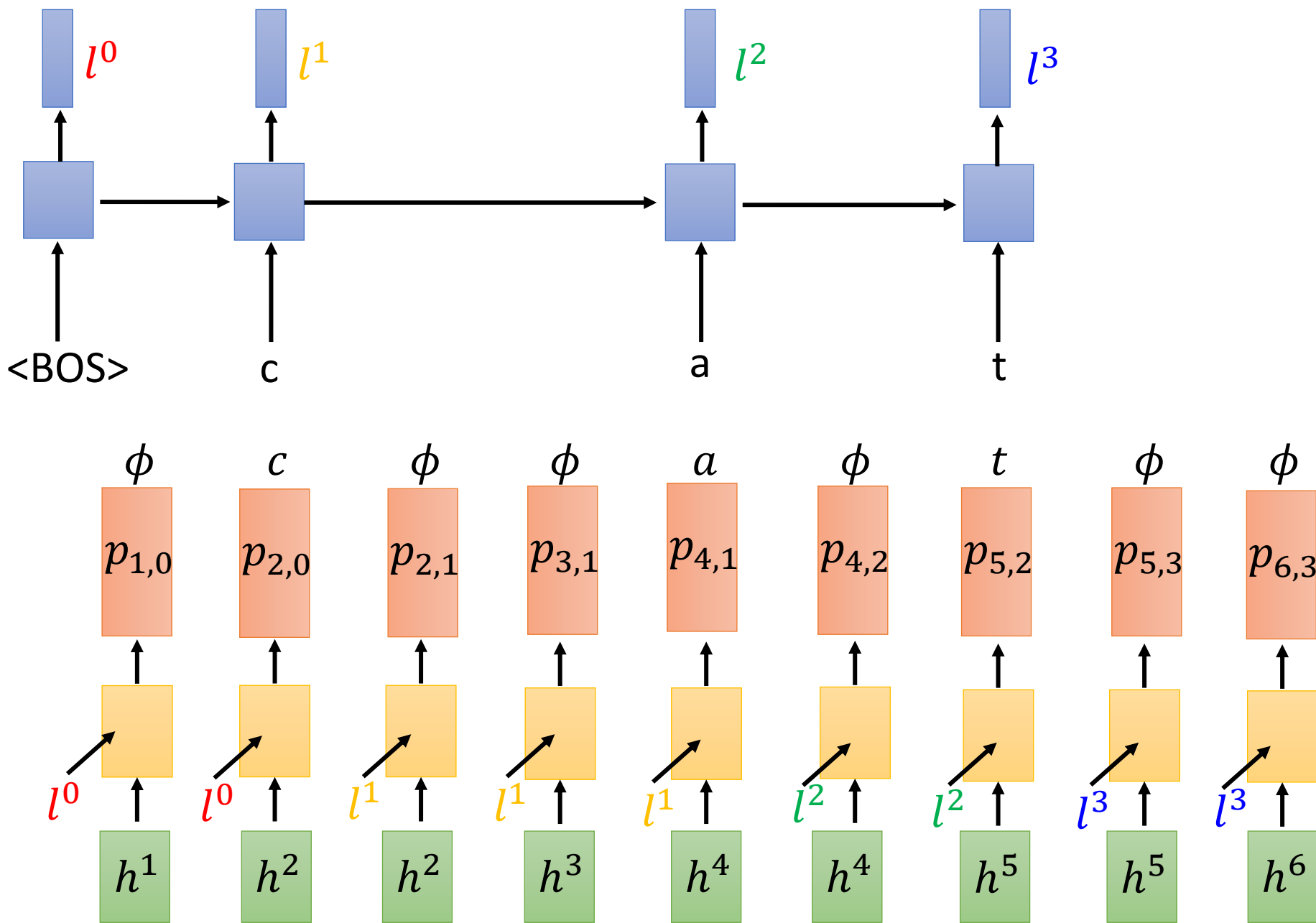
Score Computation



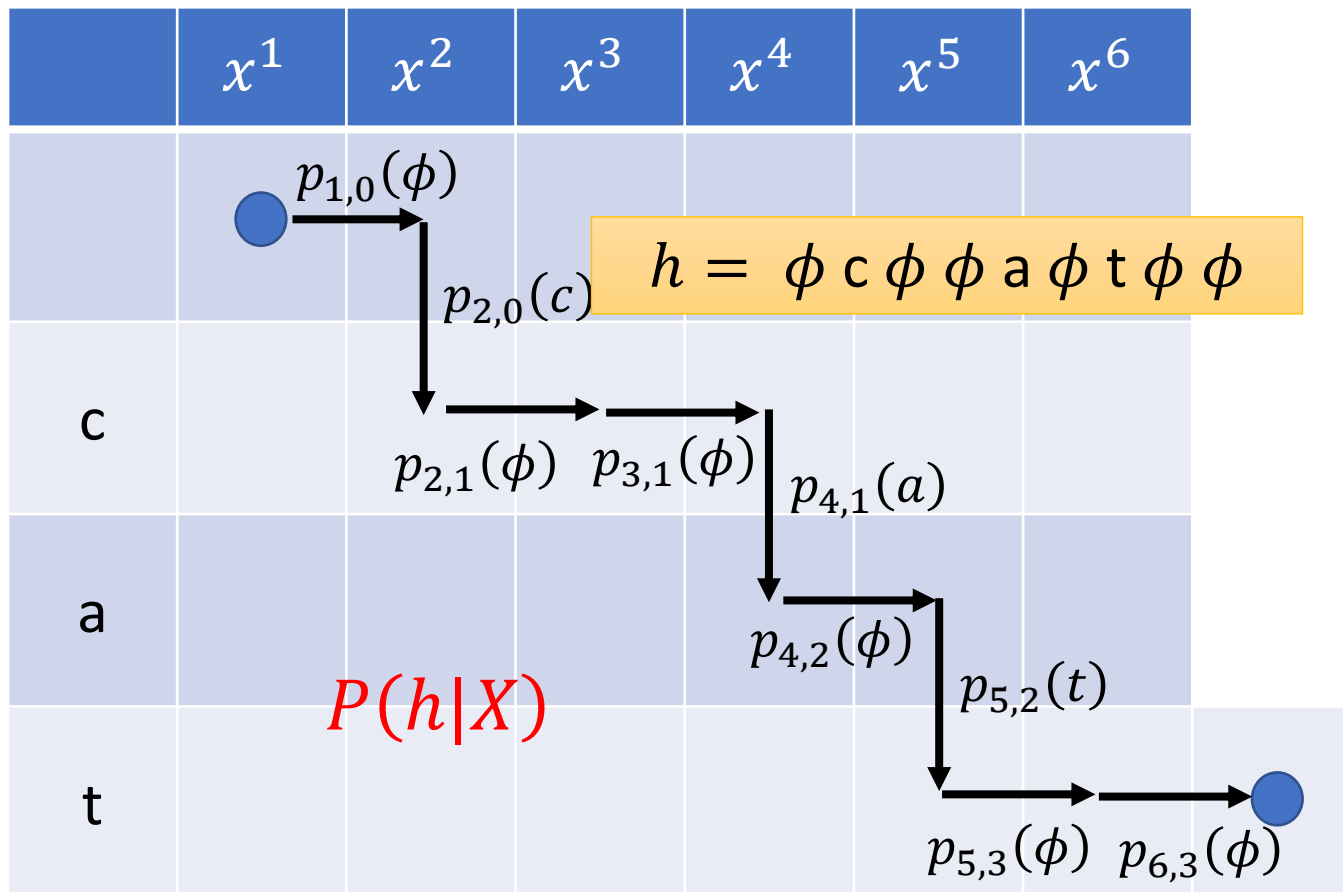
Insert ϕ
 output token

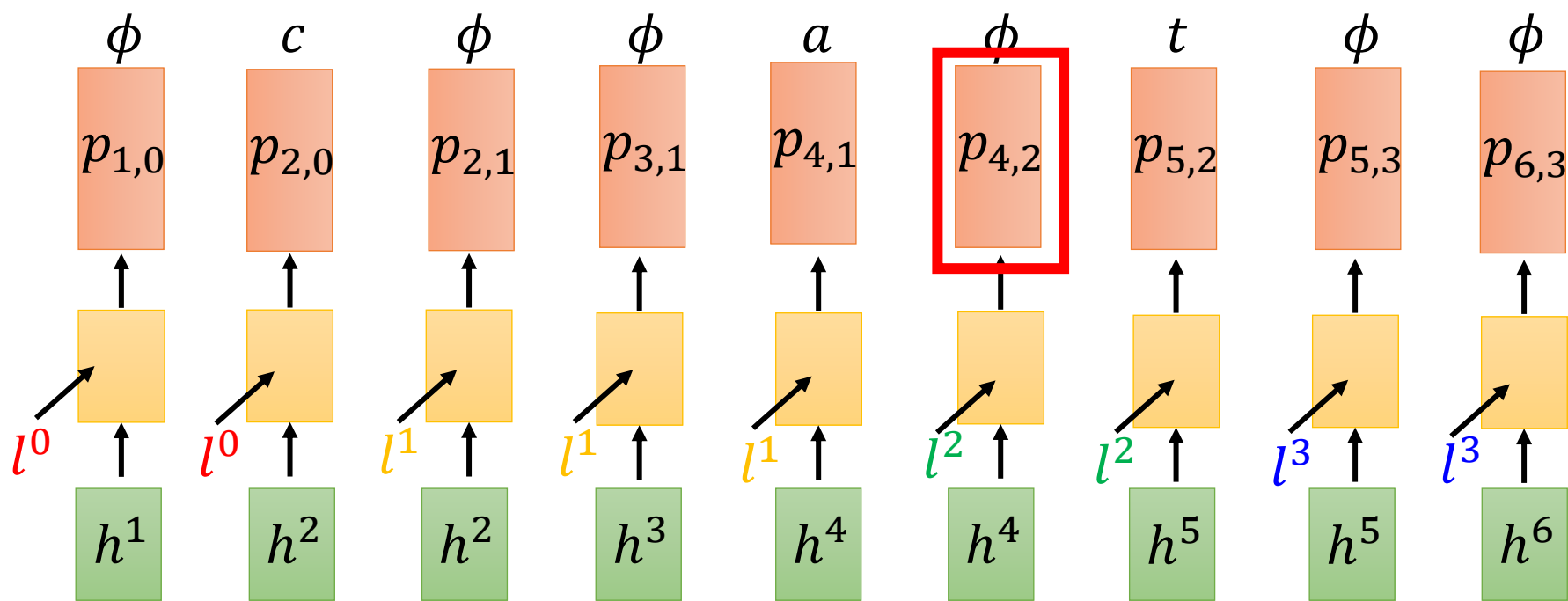
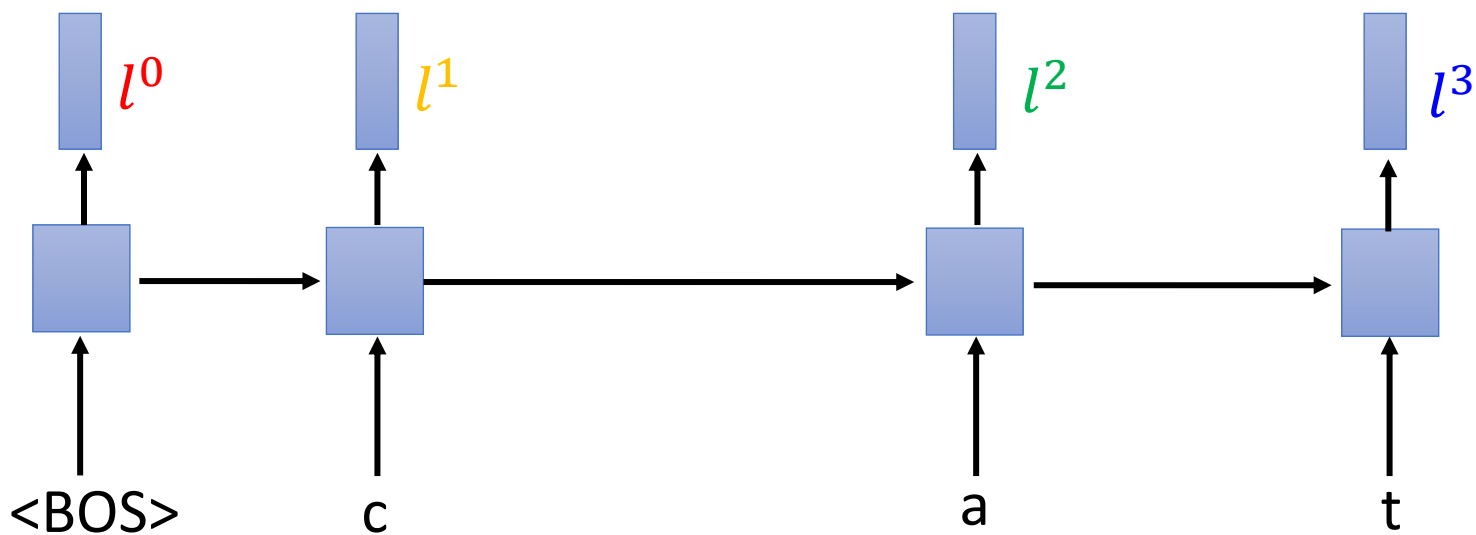
$$\begin{aligned}
 &P(h|X) \\
 &= P(\phi|X) \\
 &\quad \times P(c|X, \phi) \\
 &\quad \times P(\phi|X, \phi c) \dots \dots
 \end{aligned}$$

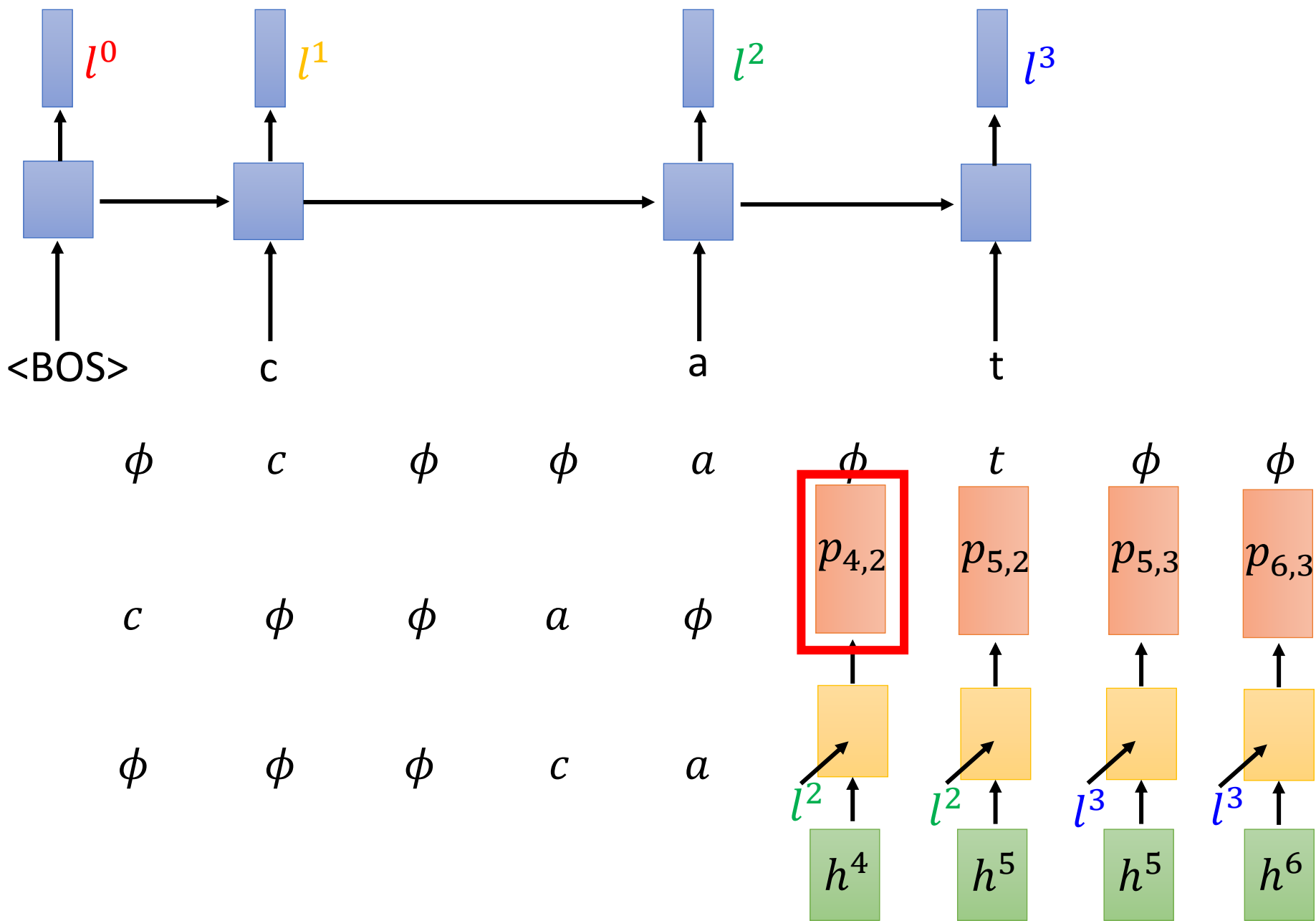
$$h = \phi c \phi \phi a \phi t \phi \phi$$



Score Computation

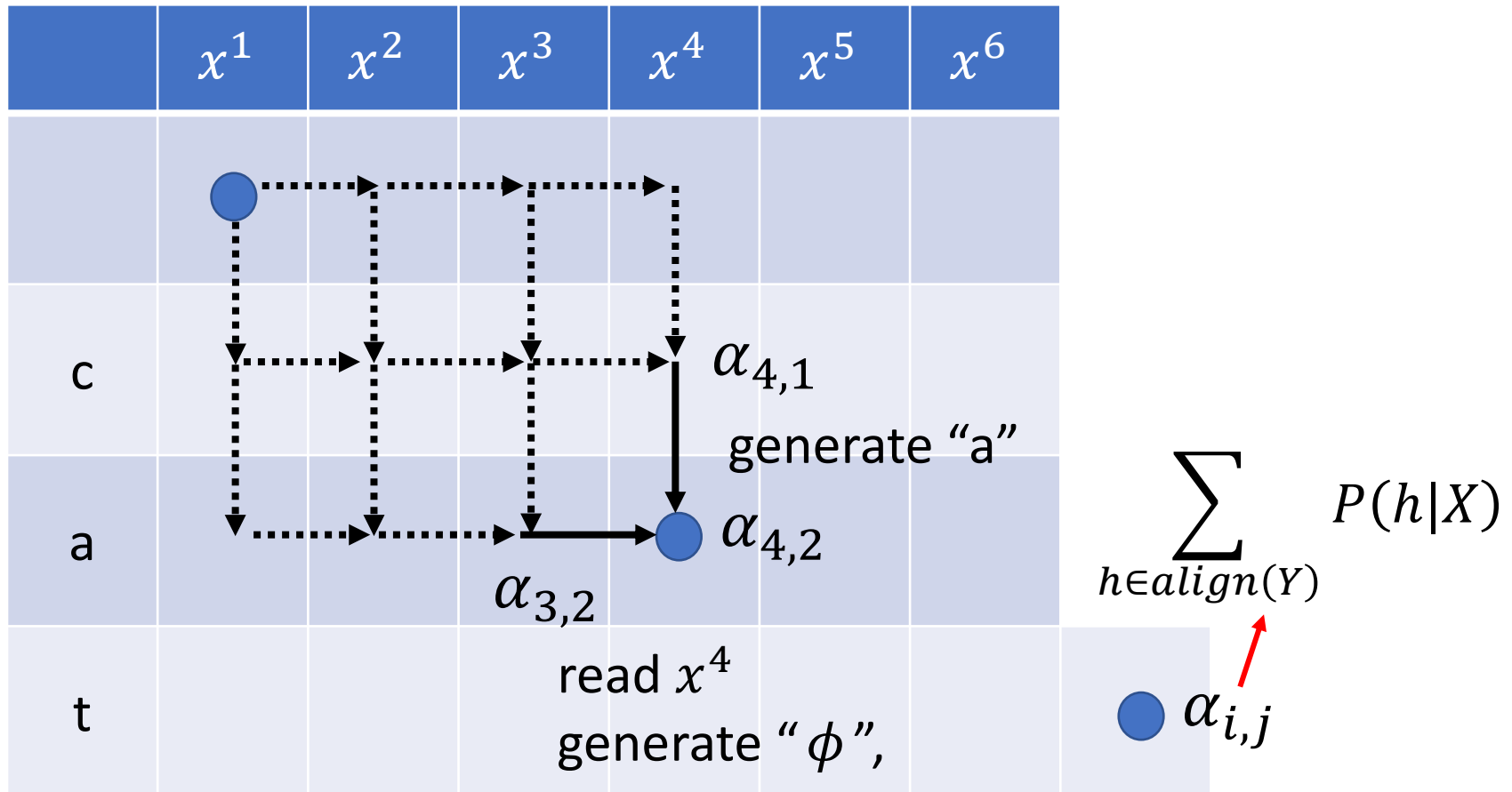






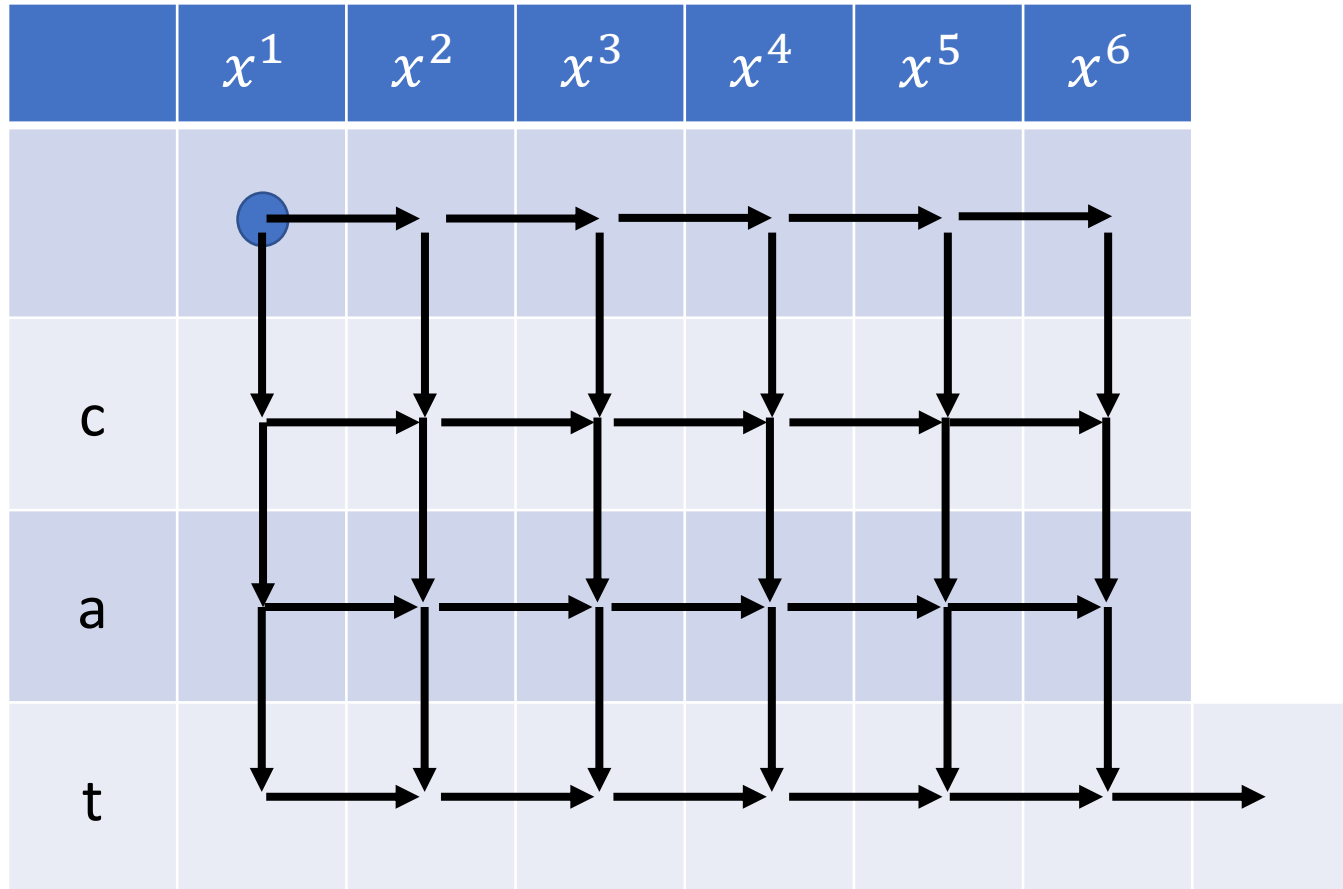
$\alpha_{i,j}$: the summation of the scores of all the alignments that read i -th acoustic features and output j -th tokens

$$\alpha_{4,2} = \alpha_{4,1}p_{4,1}(a) + \alpha_{3,2}p_{3,2}(\phi)$$



$\alpha_{i,j}$: the summation of the scores of all the alignments that read i-th acoustic features and output j-th tokens

$$\alpha_{4,2} = \alpha_{4,1}p_{4,1}(a) + \alpha_{3,2}p_{3,2}(\phi)$$



You can compute summation of the scores of all the alignments.

HMM, CTC, RNN-T

HMM

$$P_{\theta}(X|Y) = \sum_{h \in \text{align}(Y)} P(X|h)$$

CTC, RNN-T

$$P_{\theta}(Y|X) = \sum_{h \in \text{align}(Y)} P(h|X)$$

1. Enumerate all the possible alignments
2. How to sum over all the alignments

3. Training:

$$\theta^* = \arg \max_{\theta} \log P_{\theta}(\hat{Y}|X)$$

$$\frac{\partial P_{\theta}(\hat{Y}|X)}{\partial \theta} = ?$$

4. Testing (Inference, decoding):

$$Y^* = \arg \max_Y \log P(Y|X)$$

Training

$$\theta^* = \arg \max_{\theta} \log P(\hat{Y}|X)$$

$$P(\hat{Y}|X) = \sum_h P(h|X)$$

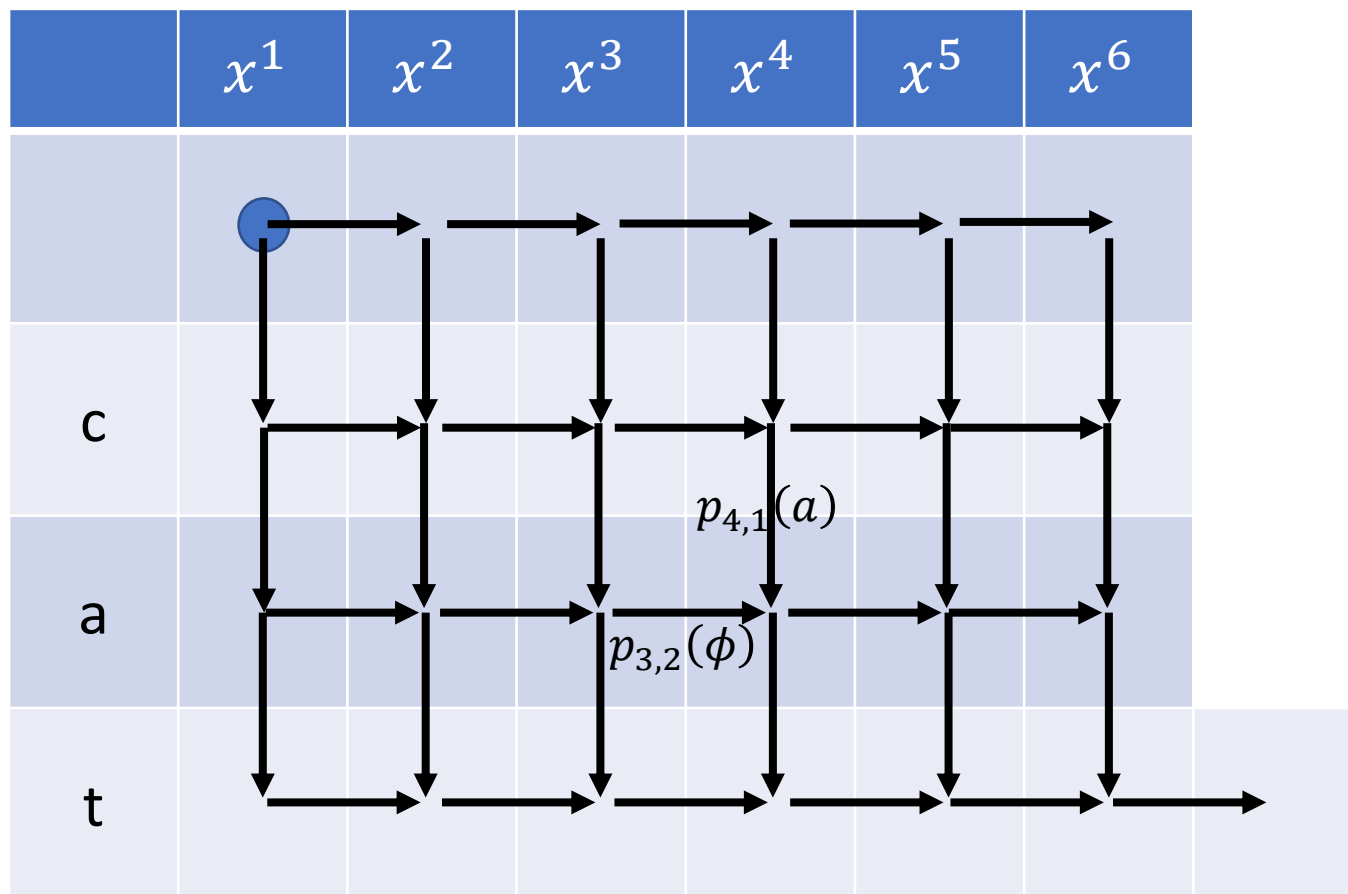
ϕ c ϕ ϕ a ϕ t ϕ ϕ

$p_{1,0}(\phi)$ $p_{2,0}(c)$ $p_{2,1}(\phi)$ $p_{3,1}(\phi)$ $p_{4,1}(a)$ $p_{4,2}(\phi)$ $p_{5,2}(t)$ $p_{5,3}(\phi)$ $p_{6,3}(\phi)$

$$\frac{\partial P(\hat{Y}|X)}{\partial \theta} = ?$$

$$P(\hat{Y}|X) = \sum_h P(h|X)$$

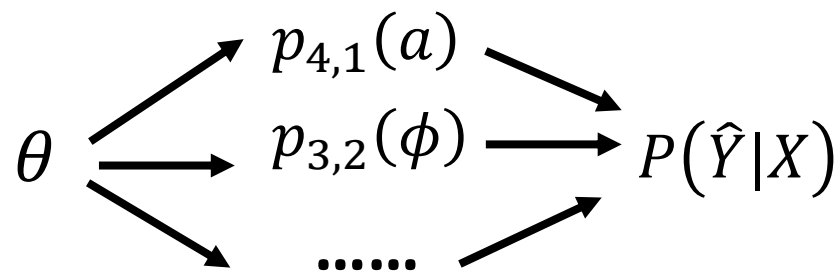
$p_{1,0}(\phi)$ $p_{2,0}(c)$ $p_{2,1}(\phi)$ $p_{3,1}(\phi)$ $p_{4,1}(a)$ $p_{4,2}(\phi)$ $p_{5,2}(t)$ $p_{5,3}(\phi)$ $p_{6,3}(\phi)$



Each arrow is a component in $P(\hat{Y}|X) = \sum_h P(h|X)$

Training

$$\theta^* = \arg \max_{\theta} \log P(\hat{Y}|X)$$



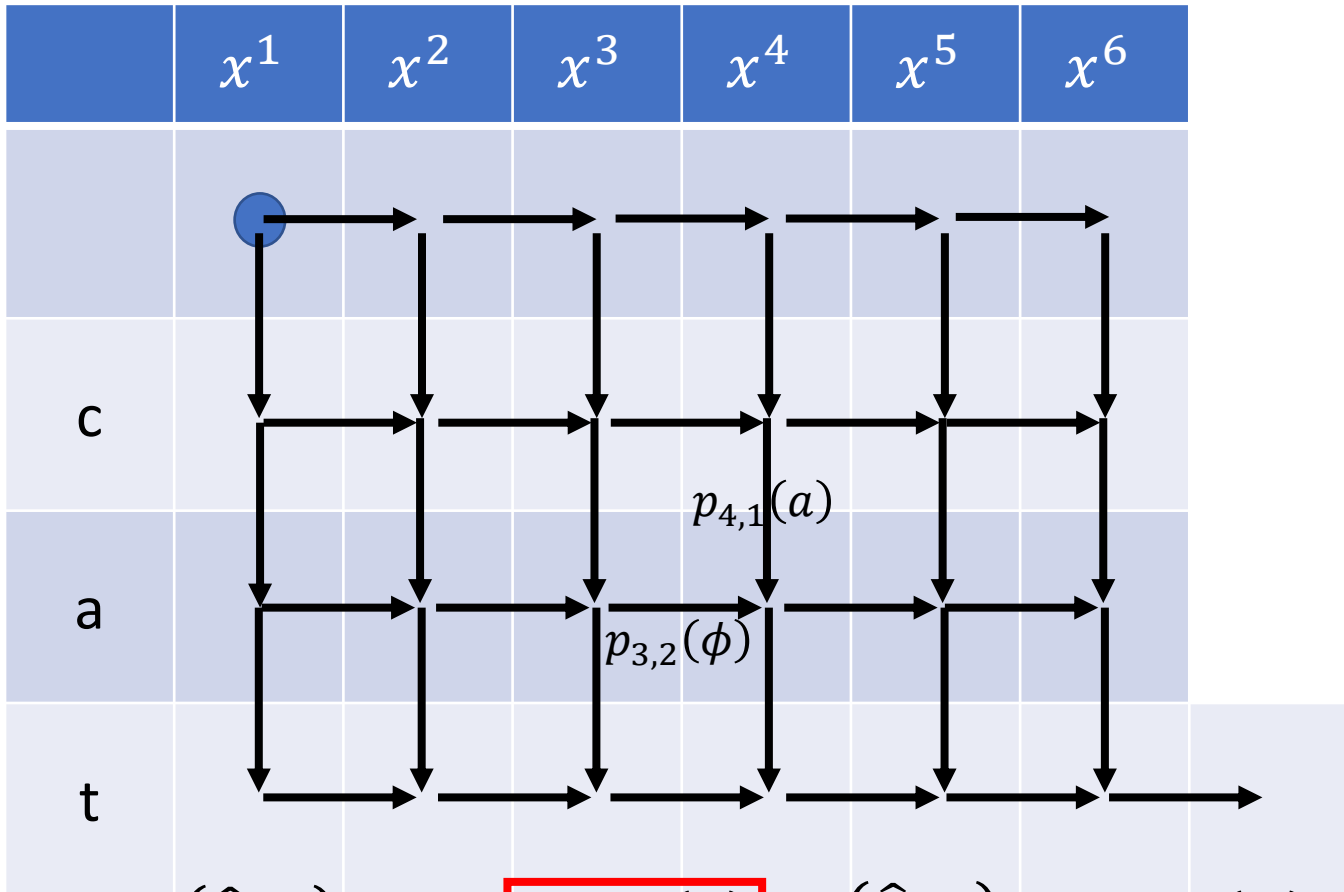
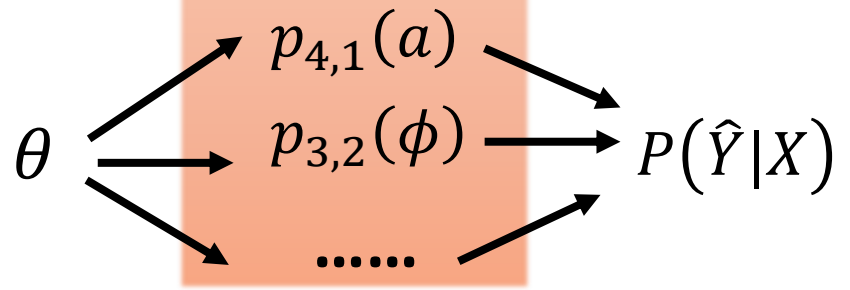
$$P(\hat{Y}|X) = \sum_h P(h|X)$$

ϕ c ϕ ϕ a ϕ t ϕ ϕ

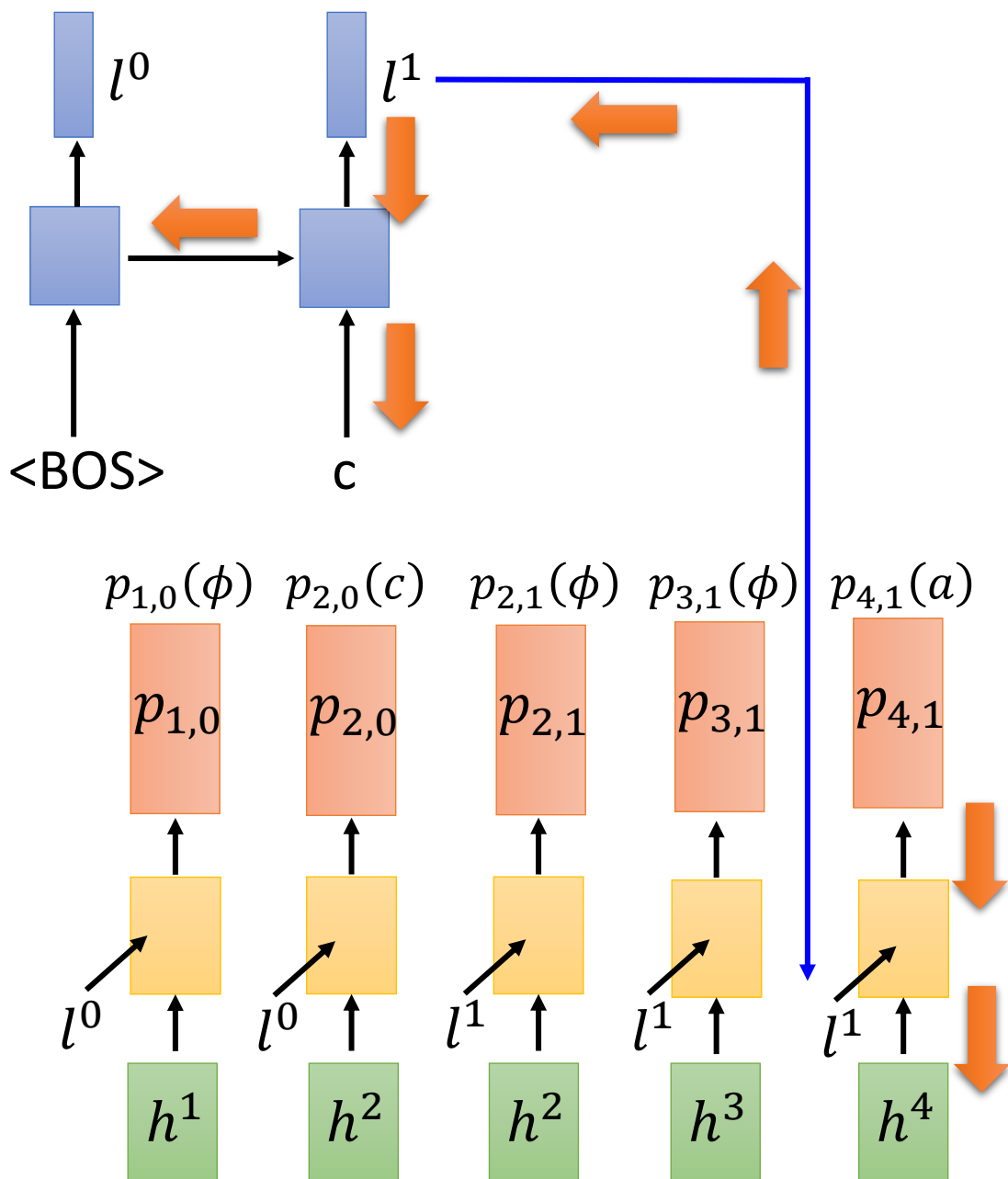
$p_{1,0}(\phi)$ $p_{2,0}(c)$ $p_{2,1}(\phi)$ $p_{3,1}(\phi)$ $p_{4,1}(a)$ $p_{4,2}(\phi)$ $p_{5,2}(t)$ $p_{5,3}(\phi)$ $p_{6,3}(\phi)$

$$\frac{\partial P(\hat{Y}|X)}{\partial \theta} = ? \quad \frac{\partial p_{4,1}(a)}{\partial \theta} \frac{\partial P(\hat{Y}|X)}{\partial p_{4,1}(a)} + \frac{\partial p_{3,2}(\phi)}{\partial \theta} \frac{\partial P(\hat{Y}|X)}{\partial p_{3,2}(\phi)} + \dots$$

Each arrow is a component



$$\frac{\partial P(\hat{Y}|X)}{\partial \theta} = ? \quad \boxed{\frac{\partial p_{4,1}(a)}{\partial \theta}} \frac{\partial P(\hat{Y}|X)}{\partial p_{4,1}(a)} + \frac{\partial p_{3,2}(\phi)}{\partial \theta} \frac{\partial P(\hat{Y}|X)}{\partial p_{3,2}(\phi)} + \dots$$



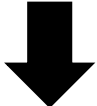
$$\frac{\partial p_{4,1}(a)}{\partial \theta} = ?$$

Backpropagation
(through time)

To encoder

$$\frac{\partial P(\hat{Y}|X)}{\partial \theta} \stackrel{=?}{=} \frac{\partial p_{4,1}(a)}{\partial \theta} \frac{\partial P(\hat{Y}|X)}{\partial p_{4,1}(a)} + \frac{\partial p_{3,2}(\phi)}{\partial \theta} \frac{\partial P(\hat{Y}|X)}{\partial p_{3,2}(\phi)} + \dots$$

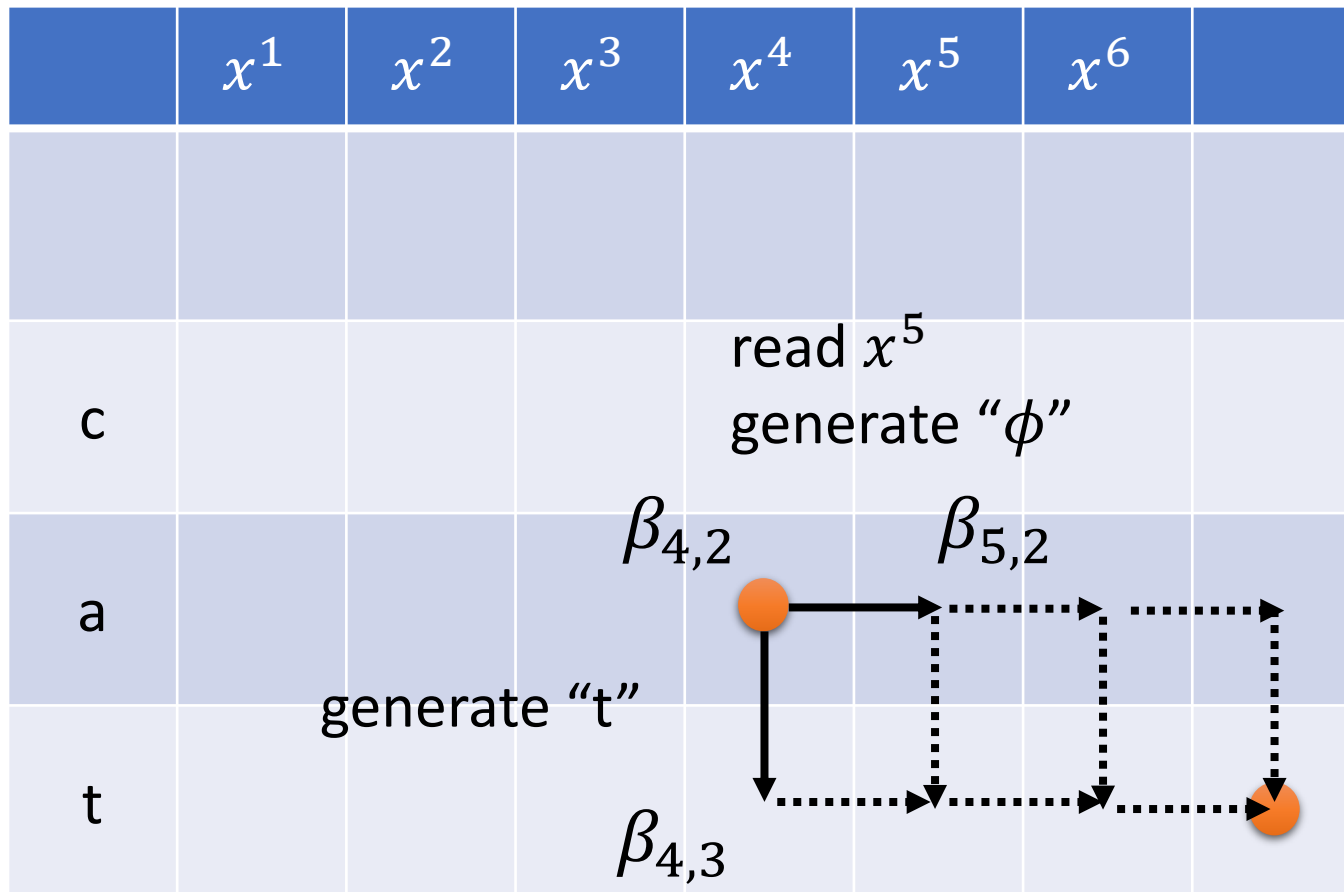
$$P(\hat{Y}|X) = \sum_{h \text{ with } p_{4,1}(a)} P(h|X) + \sum_{h \text{ without } p_{4,1}(a)} P(h|X)$$


 $p_{4,1}(a) \times \text{other}$

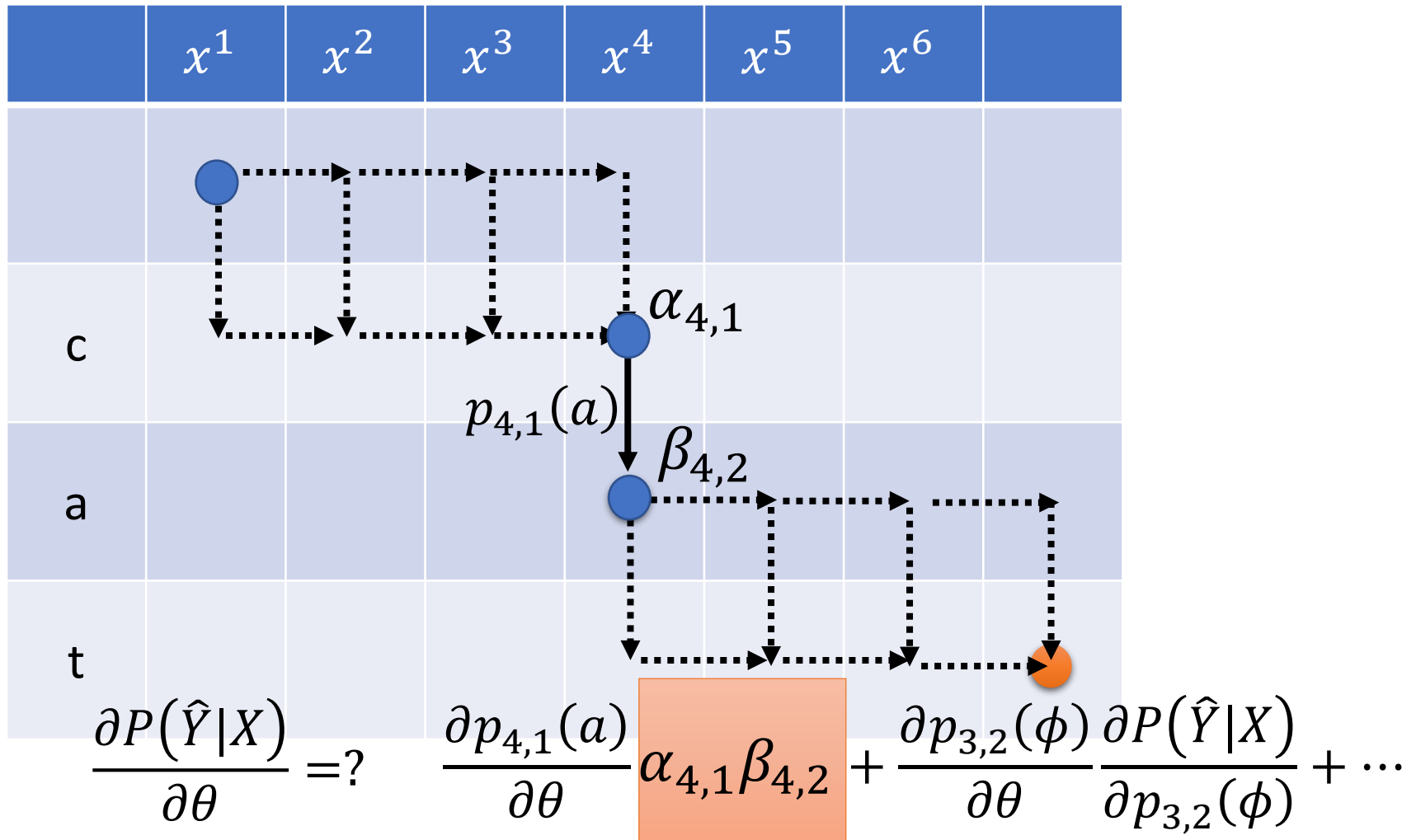
$$\begin{aligned} \frac{\partial P(\hat{Y}|X)}{\partial p_{4,1}(a)} &= \sum_{h \text{ with } p_{4,1}(a)} \text{other} = \sum_{h \text{ with } p_{4,1}(a)} \frac{P(h|X)}{p_{4,1}(a)} \\ &= \frac{1}{p_{4,1}(a)} \sum_{h \text{ with } p_{4,1}(a)} P(h|X) \end{aligned}$$

$\beta_{i,j}$: the summation of the score of all the alignments starting from i-th acoustic features and j-th tokens

$$\beta_{4,2} = \beta_{4,3}p_{4,2}(t) + \beta_{5,2}p_{4,2}(\phi)$$



$$\frac{\partial P(\hat{Y}|X)}{\partial p_{4,1}(a)} = \frac{1}{p_{4,1}(a)} \sum_{a \text{ with } p_{4,1}(a)} P(h|X) \alpha_{4,1} p_{4,1}(a) \beta_{4,2}$$



HMM, CTC, RNN-T

HMM

$$P_{\theta}(X|Y) = \sum_{h \in \text{align}(Y)} P(X|h)$$

CTC, RNN-T

$$P_{\theta}(Y|X) = \sum_{h \in \text{align}(Y)} P(h|X)$$

1. Enumerate all the possible alignments
2. How to sum over all the alignments

3. Training: $\theta^* = \arg \max_{\theta} \log P_{\theta}(\hat{Y}|X)$ $\frac{\partial P_{\theta}(\hat{Y}|X)}{\partial \theta} = ?$

4. Testing (Inference, decoding):

$$Y^* = \arg \max_Y \log P(Y|X)$$

Decoding

$$Y^* = \arg \max_Y \log P(Y|X)$$

理想

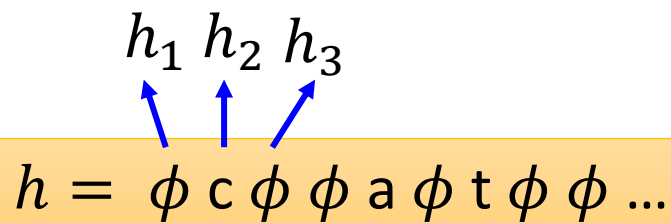
$$= \arg \max_Y \log \sum_{h \in \text{align}(Y)} P(h|X) \quad \max_{h \in \text{align}(Y)} P(h|X)$$

現實

$$\approx \arg \max_Y \max_{h \in \text{align}(Y)} \log P(h|X)$$

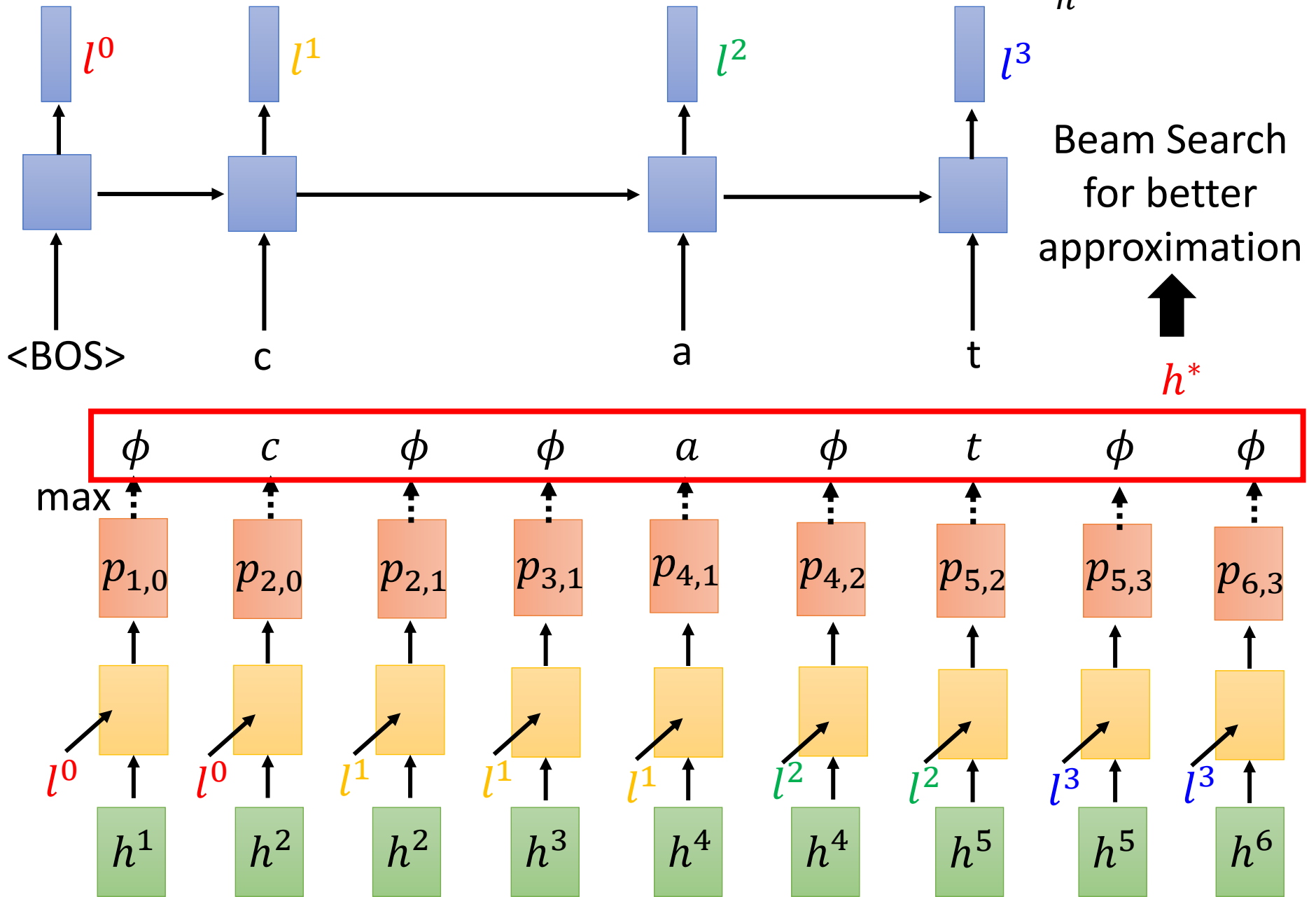
$$Y^* = \text{align}^{-1}(h^*)$$

$$h^* = \arg \max_h \log P(h|X)$$



$$P(h|X) = P(h_1|X)P(h_2|X, h_1)P(h_3|X, h_1, h_2) \dots$$

$$h^* = \arg \max_h \log P(h|X)$$



Summary

	LAS	CTC	RNN-T
Decoder	dependent	independent	dependent
Alignment	not explicit (soft alignment)	Yes	Yes
Training	just train it	sum over alignment	sum over alignment
On-line	No	Yes	Yes

Reference

- [Yu, et al., INTERSPEECH'16] Dong Yu, Wayne Xiong, Jasha Droppo, Andreas Stolcke , Guoli Ye, Jinyu Li , Geoffrey Zweig, Deep Convolutional Neural Networks with Layer-wise Context Expansion and Attention, INTERSPEECH, 2016
- [Saon, et al., INTERSPEECH'17] George Saon, Gakuto Kurata, Tom Sercu, Kartik Audhkhasi, Samuel Thomas, Dimitrios Dimitriadis, Xiaodong Cui, Bhuvana Ramabhadran, Michael Picheny, Lynn-Li Lim, Bergul Roomi, Phil Hall, English Conversational Telephone Speech Recognition by Humans and Machines, INTERSPEECH, 2017
- [Povey, et al., ICASSP'10] Daniel Povey, Lukas Burget, Mohit Agarwal, Pinar Akyazi, Kai Feng, Arnab Ghoshal, Ondrej Glembek, Nagendra Kumar Goel, Martin Karafiat, Ariya Rastrow, Richard C. Rose, Petr Schwarz, Samuel Thomas, Subspace Gaussian Mixture Models for speech recognition, ICASSP, 2010
- [Mohamed , et al., ICASSP'10] Abdel-rahman Mohamed and Geoffrey Hinton, Phone recognition using Restricted Boltzmann Machines, ICASSP, 2010