# *Speech Recognition*

HUNG-YI LEE 李宏毅

# Two Points of Views



**Seq-to-seq**



Source of image:
李琳山老師
《數位語音處理概論》

**HMM**

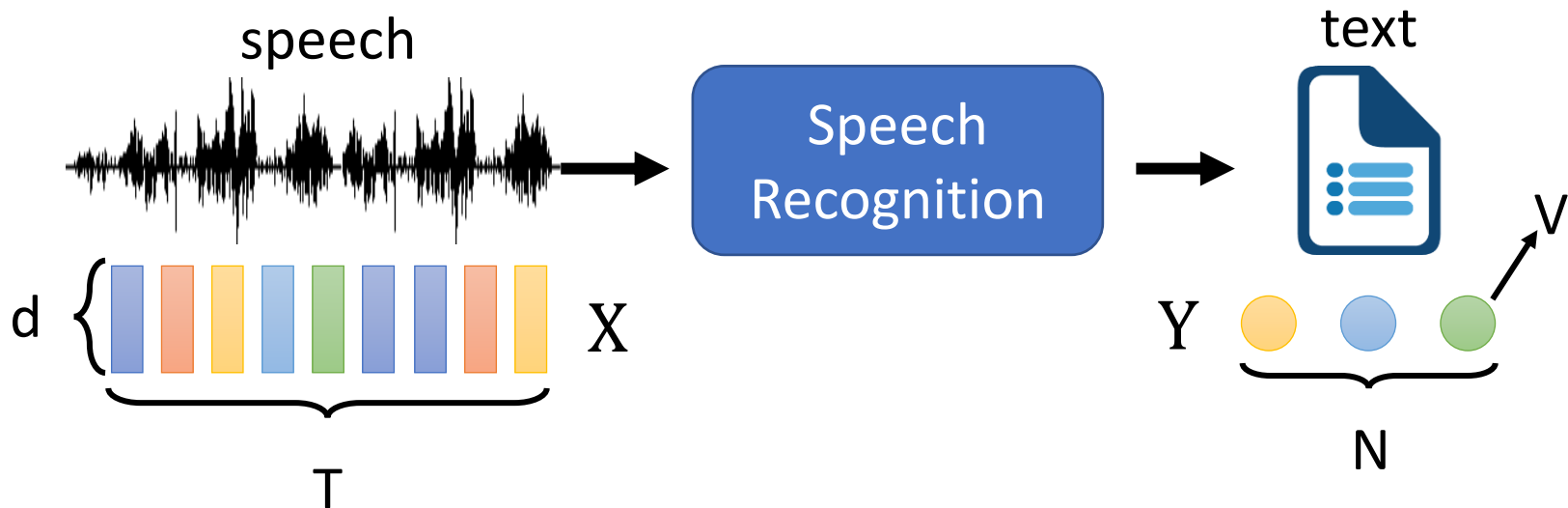# Hidden Markov Model (HMM)



speech

Speech Recognition

text

$X$  $Y$  $V$

d  $T$  $N$

$$P_\theta(X|\underline{Y}) = ?$$

The token here is small unit called state.

**_Training_**

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

**_Testing_**

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

# HMM

- A sentence corresponds to a sequence of **states**

what do you think

*Phoneme:*
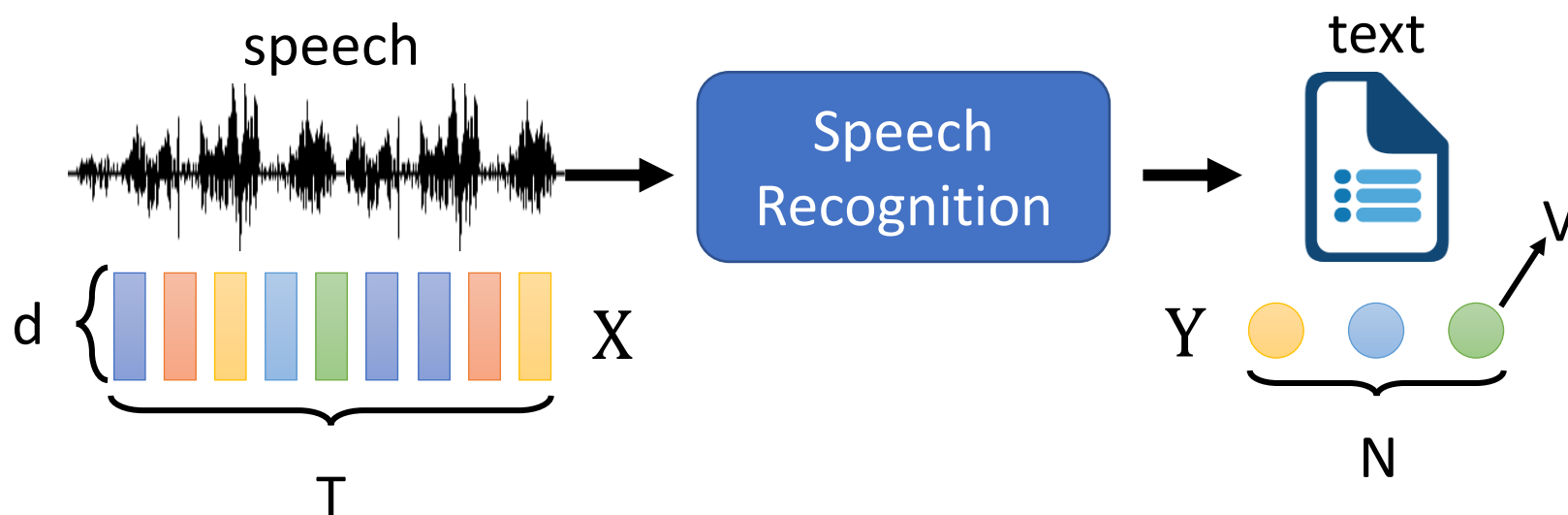
hh w aa t    d uw    y uw    th ih ng k

*Tri-phone:*

…… t-d+uw d-uw+y    uw-y+uw y-uw+th  ……

t-d+uw1 t-d+uw2 t-d+uw3  d-uw+y1 d-uw+y2 d-uw+y3

*State:*

# Hidden Markov Model (HMM)



speech

Speech Recognition

text

$d$ { X

$T$

Y

$V$

$N$

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

*__Training__*

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

emission

transition

$a \rightarrow b$
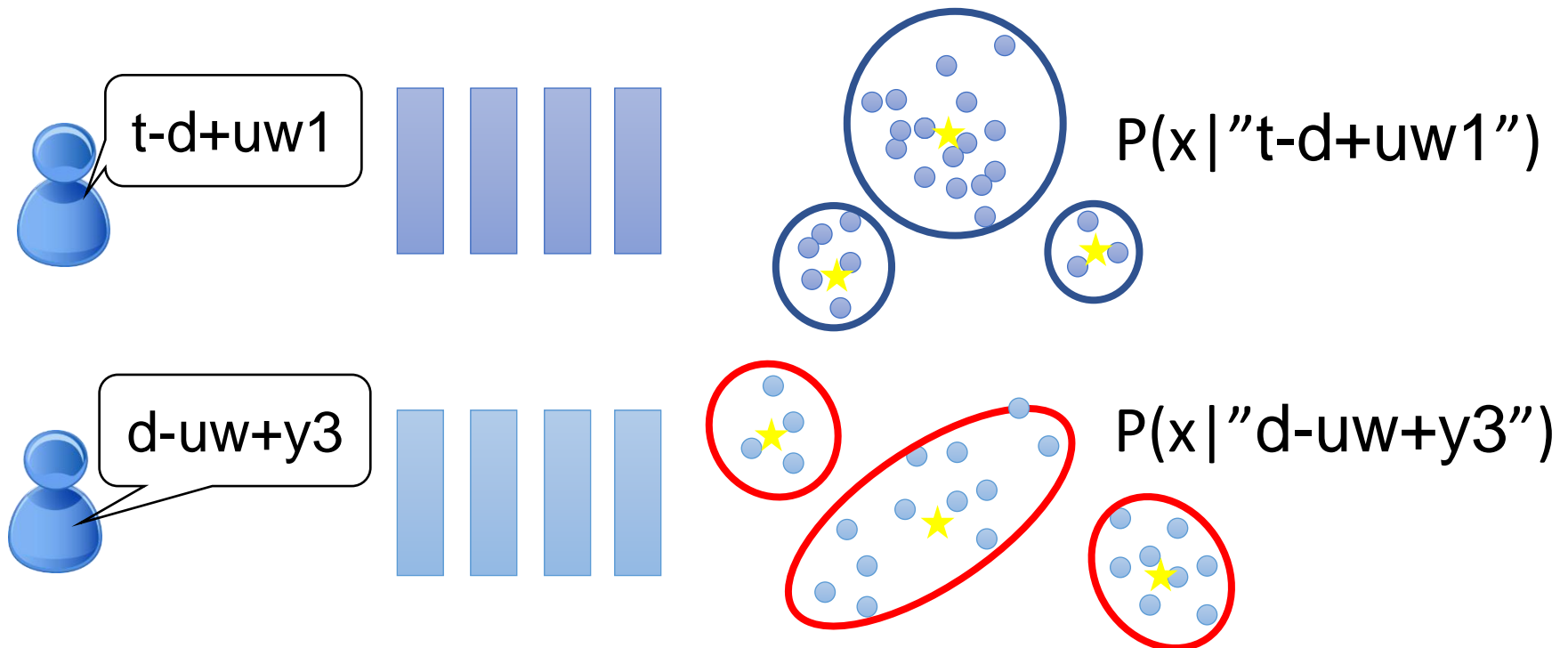
$p(b|a)$

*__Testing__*

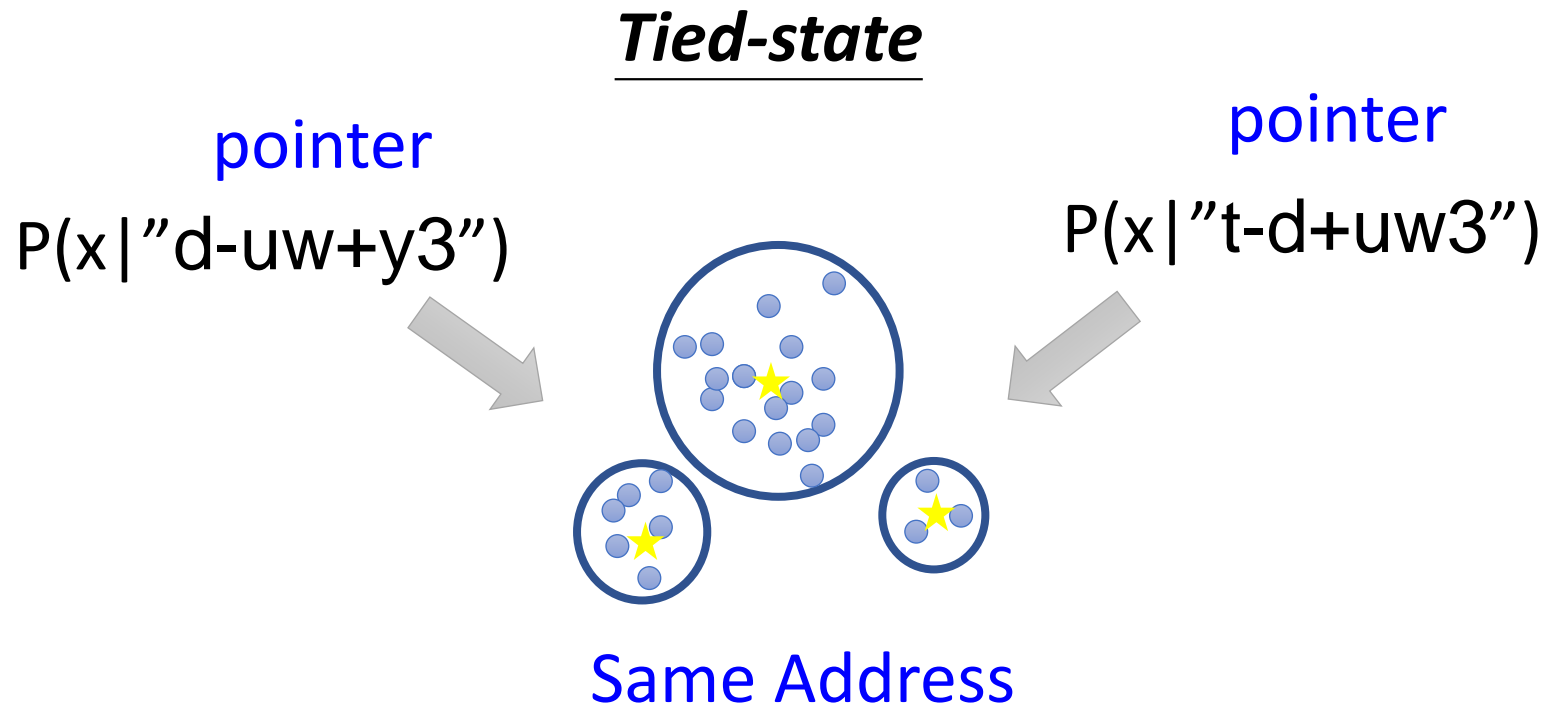$$Y^* = arg \max_{Y} log P_{\theta}(X|Y)$$

# HMM – Emission Probability

- Each state has a stationary distribution for acoustic features

Gaussian Mixture Model (GMM)



$P(x|\text{"t-d+uw1"})$

$P(x|\text{"d-uw+y3"})$
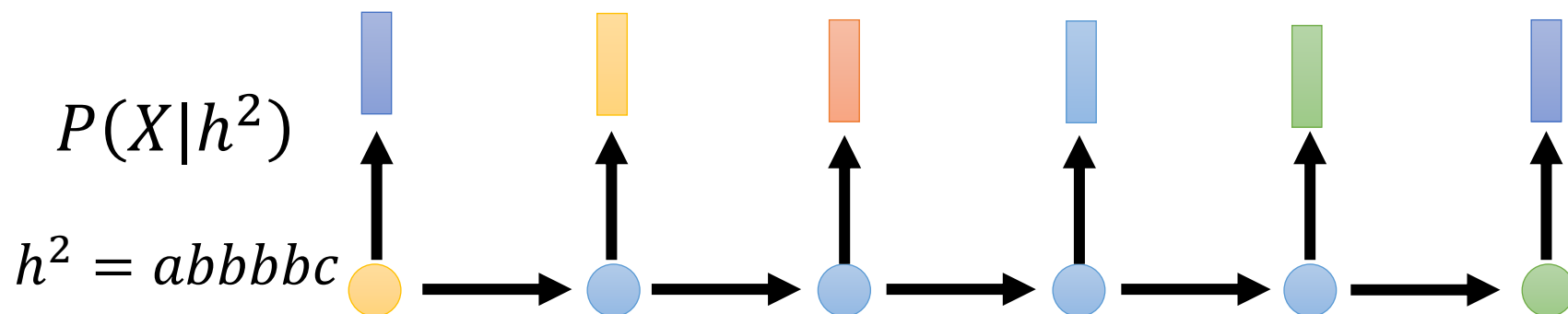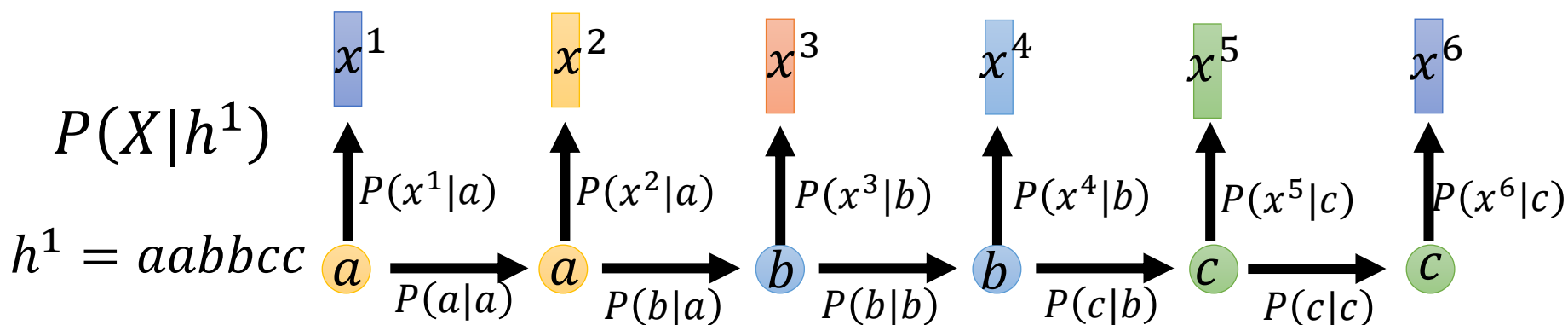
# HMM – Emission Probability

- Each state has a stationary distribution for acoustic features

**_Tied-state_**

pointer

P(x|"d-uw+y3")

pointer

P(x|"t-d+uw3")

Same Address

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

$h = abccbc$ ✖

$h = abbbcccc$ ✖

emission

transition

$a \rightarrow b$

$p(b|a)$

$a$  $b$  $c$  alignment

which state generates which vector

$P(X|h^1)$

$x^1$  $x^2$  $x^3$  $x^4$  $x^5$  $x^6$

$h^1 = aabbcc$

$P(x^1|a)$  $P(x^2|a)$  $P(x^3|b)$  $P(x^4|b)$  $P(x^5|c)$  $P(x^6|c)$

$a \xrightarrow{P(a|a)} a \xrightarrow{P(b|a)} b \xrightarrow{P(b|b)} b \xrightarrow{P(c|b)} c \xrightarrow{P(c|c)} c$

$P(X|h^2)$

$h^2 = abbbbc$

# Before End-to-end – Tandem

$p(a|x^i)$  $p(b|x^i)$  $p(c|x^i)$ ......



New acoustic feature for HMM

Size of output layer = No. of states

DNN

$x^i$

Last hidden layer or bottleneck layer are also possible.

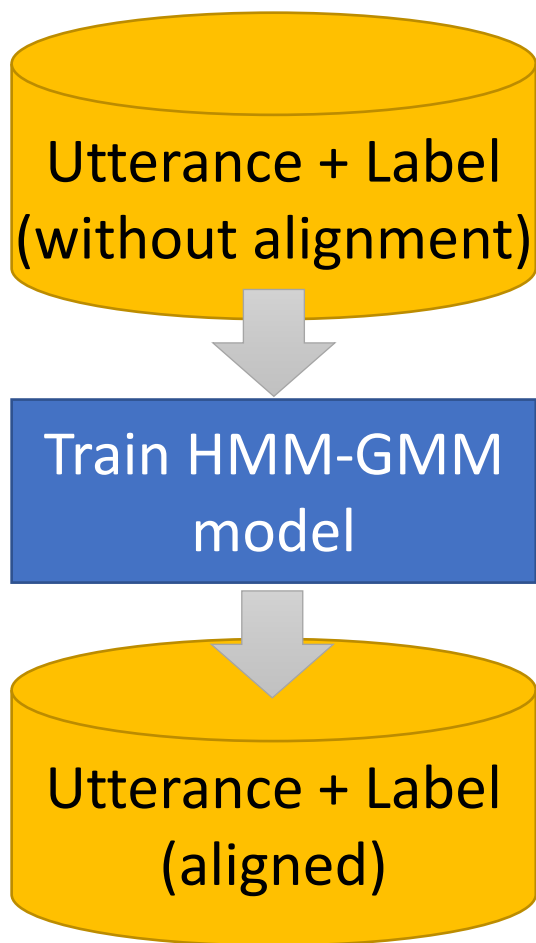# Before End-to-end – Tandem



Acoustic features:

{ a , b , c }
state sequence:

Acoustic features:

a   a   a   b   b   c   c

state sequence:

# Before End-to-end – Tandem

# Before End-to-end – Tandem

# Before End-to-end – Hybrid



$$P(x|a) = \frac{P(x,a)}{P(a)} = \frac{P(a|x)P(x)}{P(a)}$$

DNN output

Count from training data

# 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<br>conv 3x3, 64<br>conv 1x1, 256] | $5 - $ conv 3x3, 128 |
| Max pool | $4 - $ [conv 1x1, 128<br>conv 3x3, 128<br>conv 1x1, 512] | $5 - $ conv 3x3, 256 |
| $4 - $ conv 3x3, 192 | $6 - $ [conv 1x1, 256<br>conv 3x3, 256<br>conv 1x1, 1024] | $5 - $ conv 3x3, 512 |
| Max pool | $3 - $ [conv 1x1, 512<br>conv 3x3, 512<br>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]

# LAS

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

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

Beam Search

$P(Y|X) =?$



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

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

a    b

$p(a)$   $p(b)$   $p(EOS)$

Size V

$z^0 \rightarrow z^1 \rightarrow z^2 \rightarrow z^3$

$c^0$   $c^1$   $c^2$

# CTC, RNN-T

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

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

$P(Y|X) = ?$

$P(h|X)$     $h = a\ \phi\ b\ \phi \rightarrow a\ b$

$p(a)$     $p(\phi)$     $p(b)$     $p(\phi)$

$a$     $b$     $x^1$  $x^2$  $x^3$  $x^4$

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

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

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

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

$h^1$  $h^2$  $h^3$  $h^4$

Encoder

# HMM, CTC, RNN-T

**_HMM_**

**_CTC, RNN-T_**

$$\mathrm{P}_\theta(X|Y) = \boxed{\sum_{h \in align(Y)} P(X|h)}$$

$$\mathrm{P}_\theta(Y|X) = \boxed{\sum_{h \in align(Y)} P(h|X)}$$
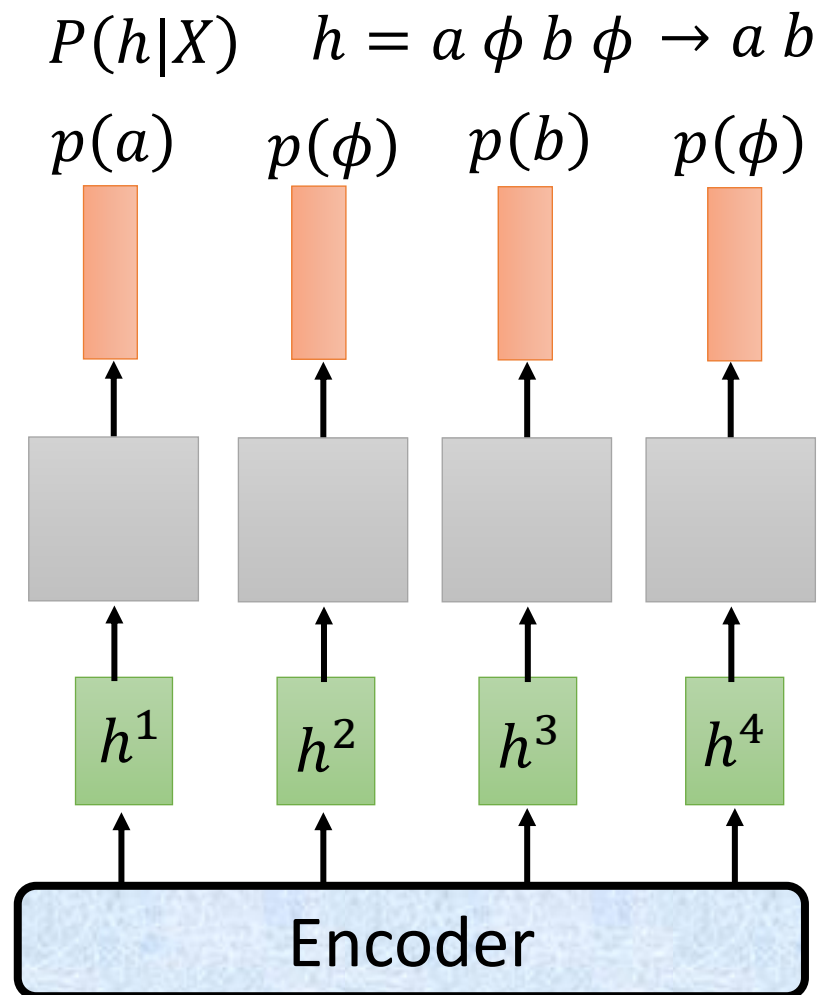
1. Enumerate all the possible alignments

2. How to sum over all the alignments

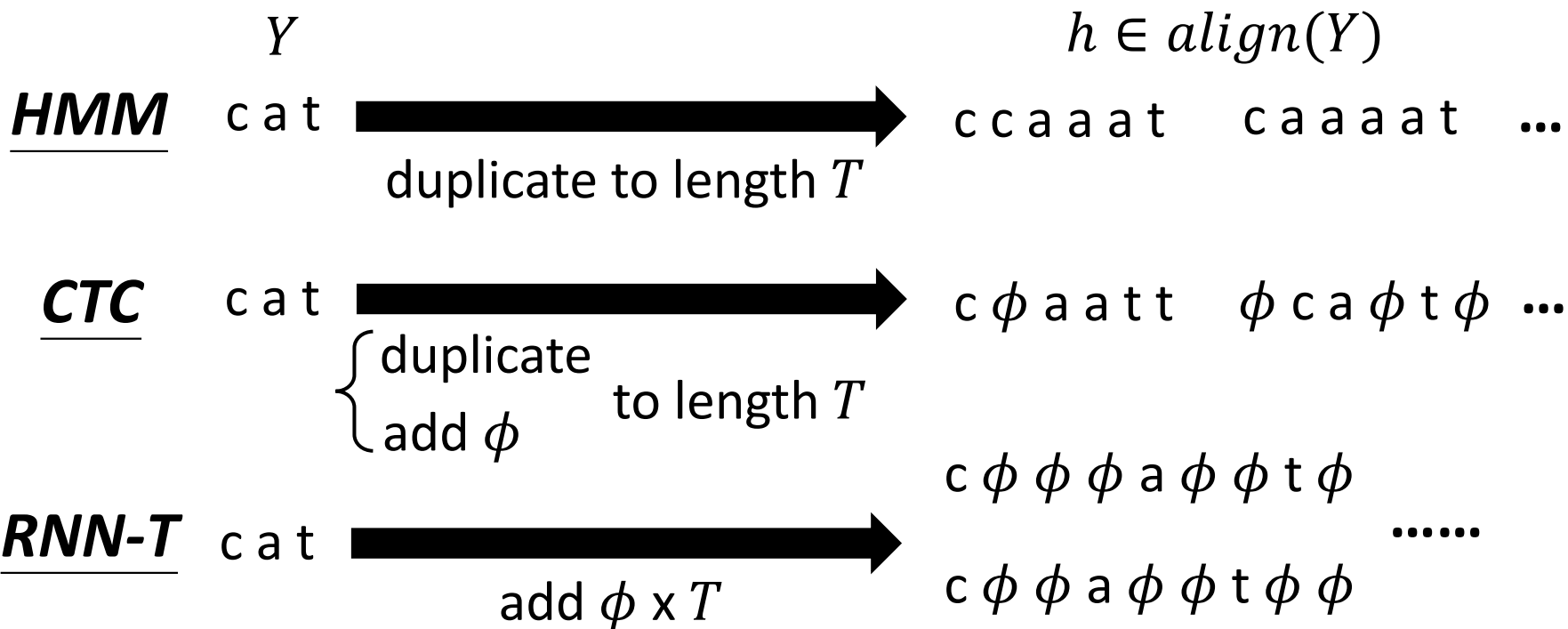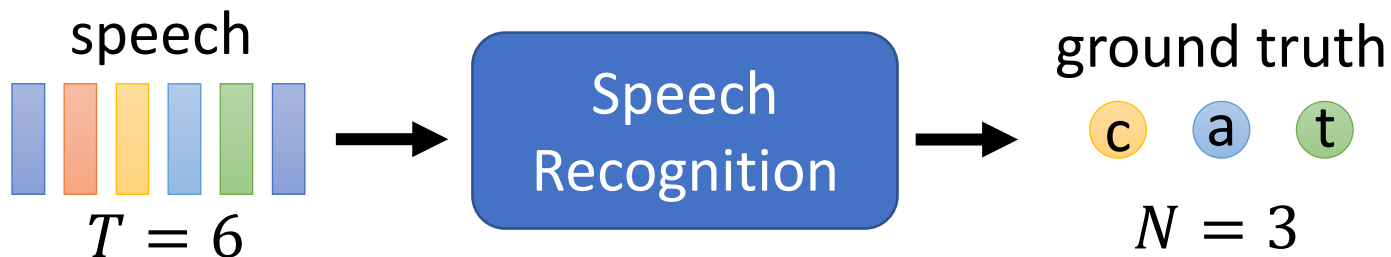3. Training: $\boxed{\theta^* = arg \max_\theta log\mathrm{P}_\theta(\hat{Y}|X)}$ $\quad \dfrac{\partial P(\hat{Y}|X)}{\partial \theta} = ?$

4. Testing (Inference, decoding):

$$Y^* = arg \max_Y log\mathrm{P}_\theta(Y|X)$$

**HMM**  c a t  ⟶  c c a a a t  c a a a a t  ...

duplicate to length $T$

For n = 1 to $N$

output the n-th token $t_n$ times

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

Trellis Graph



⟶ duplicate

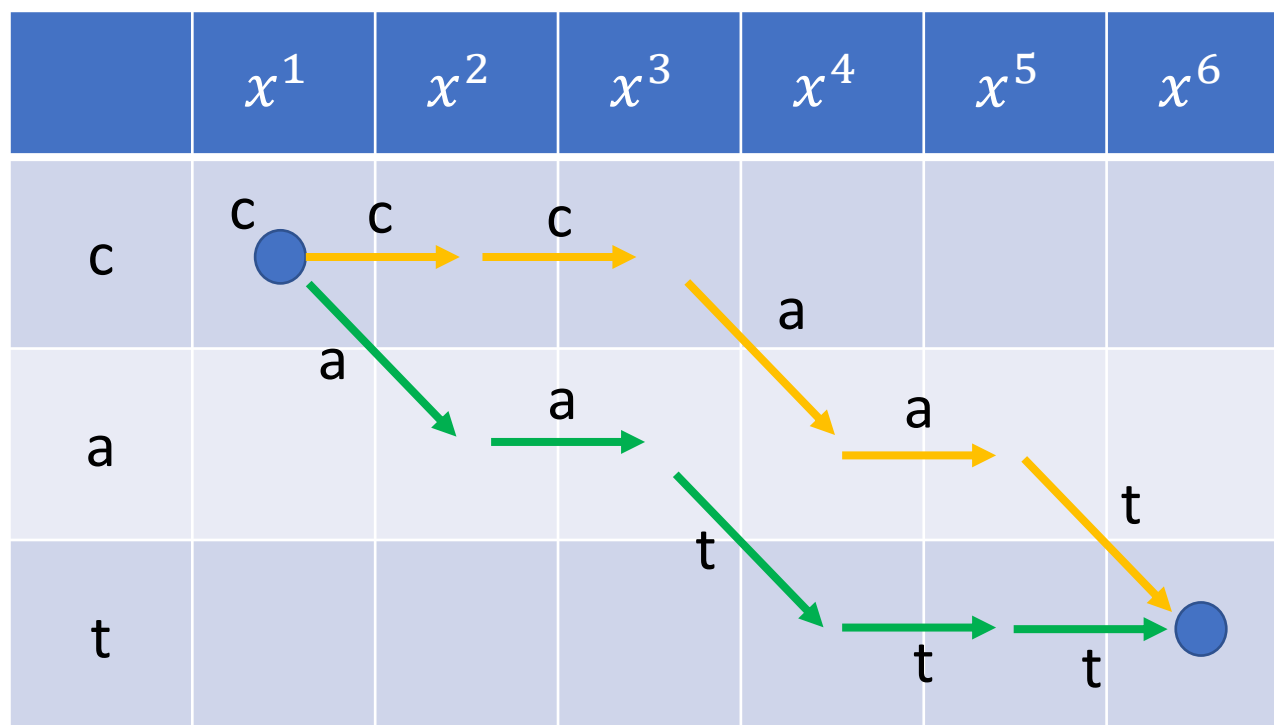↘ next token

**_HMM_**   c a t   ➡   c c a a a t   c a a a a t   ...

duplicate to length $T$

For n = 1 to $N$

   output the n-th token $t_n$ times

**_constraint_**: $t_1 + t_2 + \cdots t_N = T$, $t_n > 0$

Trellis Graph

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| c | c | c | c | c | c | c ☹ |
| a | | | | | | |
| t | | | | | | |

⟶ duplicate

↘ next token

**_CTC_**   c a t ➡️   c $\phi$ a a t t   $\phi$ c a $\phi$ t $\phi$ ...

$\begin{cases} \text{duplicate} \\ \text{add } \phi \end{cases}$ to length $T$

output "$\phi$" $c_0$ times

For n = 1 to $N$

    output the n-th token $t_n$ times

    output "$\phi$" $c_n$ times

**_constraint_**:   $t_1 + t_2 + \cdots t_N +$

$$c_0 + c_1 + \cdots c_N = T$$

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

**_CTC_**   c a t   →   c $\phi$ a a t t   $\phi$ c a $\phi$ t $\phi$ ...

$\begin{cases} \text{duplicate} \\ \text{add } \phi \end{cases}$ to length $T$

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| $\phi$ | ● | | | | | |
| c | ● | | | | | |
| $\phi$ | | | | | | |
| a | | | | | | |
| $\phi$ | | | | | | |
| t | | | | | | ● |
| $\phi$ | | | | | | ● |

duplicate

insert $\phi$

next token

($\phi$ can be skipped)

**_CTC_**　c a t　→　c $\phi$ a a t t　$\phi$ c a $\phi$ t $\phi$ ...

$\begin{cases} \text{duplicate} \\ \text{add } \phi \end{cases}$ to length $T$

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| $\phi$ | ● | | | | | |
| c | ● | | | | | |
| $\phi$ | | | | | | |
| a | | | | | | |
| $\phi$ | | | | | | |
| t | | | | | | ● |
| $\phi$ | | | | | | ● |

duplicate $\phi$

next token

✖ cannot skip any token

**_CTC_**  c a t  ⟶  c $\phi$ a a t t   $\phi$ c a $\phi$ t $\phi$ ...

$\begin{cases} \text{duplicate} \\ \text{add } \phi \end{cases}$ to length $T$

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| $\phi$ | ● | | | | | |
| c | ● | | | | | |
| $\phi$ | | | | ● duplicate / insert $\phi$ | | |
| a | | ● duplicate / insert $\phi$ / next token | | | | |
| $\phi$ | | | | | | |
| t | | | | | | ● |
| $\phi$ | | | | | | ● |

**_CTC_** c a t →(duplicate / add $\phi$ — to length $T$) → c $\phi$ a a t t $\quad$ $\phi$ c a $\phi$ t $\phi$ ...



| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| $\phi$ | $\phi$ ● | | | | | |
| c | c ● | | | | | |
| $\phi$ | | | | | | |
| a | | | | | | |
| $\phi$ | | | | | | |
| t | | | | | | ● |
| $\phi$ | | | | | | ● |

**_CTC_**  c a t  $\longrightarrow$  c $\phi$ a a t t   $\phi$ c a $\phi$ t $\phi$ ...

$\begin{cases} \text{duplicate} \\ \text{add } \phi \end{cases}$ to length $T$

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| $\phi$ | $\phi$ ● | | | | | |
| c | c ● | c | c | c | | |
| $\phi$ | | a | | | a | |
| a | | | | | | |
| $\phi$ | | | | t | | |
| t | | | | | $\phi$ | t |
| $\phi$ | | | | | $\phi$ | ● |

**_CTC_**   c a t $\longrightarrow$   c $\phi$ a a t t   $\phi$ c a $\phi$ t $\phi$ **...**

$\begin{cases} \text{duplicate} \\ \text{add } \phi \end{cases}$ to length $T$

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| $\phi$ | ● | | | | | |
| s | ● | | | | | |
| $\phi$ | | | | | | |
| **e** | | | duplicate | | | |
| $\phi$ | | | insert $\phi$ | | | |
| **e** | | | next token | | | ● |
| $\phi$ | | | … ee … → e | | | ● |

Exception: when the next token is the same token

**_RNN-T_** c a t $\xrightarrow[\text{add } \phi \text{ x } T]{}$ c $\phi$ $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$

c $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$ $\phi$ ......

Put some $\phi$ (option)

Put some $\phi$ (option)

c     a     t

Put some $\phi$ (option)

Put some $\phi$

at least once

output "$\phi$" $c_0$ times

For n = 1 to $N$

    output the n-th token 1 times

    output "$\phi$" $c_n$ times

**_constraint_**: $c_0 + c_1 + \cdots c_N = T$

$c_N > 0$

$c_n \geq 0$ for n = 1 to $N-1$

**_RNN-T_** c a t $\xrightarrow{\text{add } \phi \text{ x } T}$ c $\phi$ $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$ ......

c $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$ $\phi$

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| | | | | | | |
| c | | | | | | |
| a | | | | | | |
| t | | | | | | |

$\phi$ c

$\longrightarrow$ Insert $\phi$

$\downarrow$ output token

$\phi$

**RNN-T** c a t →(add $\phi$ × $T$)→ c $\phi$ $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$

c $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$ $\phi$   ……

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| | | | | | | |
| c | | | | | | |
| a | | | | | | |
| t | | | | | | |

→ Insert $\phi$

↓ output token

**_RNN-T_** c a t $\xrightarrow{\text{add } \phi \text{ x } T}$ 

c $\phi$ $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$ ......

c $\phi$ $\phi$ a $\phi$ $\phi$ t $\phi$ $\phi$



|     | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|-----|-------|-------|-------|-------|-------|-------|
| c   |       |       |       |       |       |       |
| a   |       |       |       |       |       |       |
| t   |       |       |       |       |       |       |

→ Insert $\phi$

↓ output token

# HMM, CTC, RNN-T

**HMM**

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

**CTC, RNN-T**

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

1. Enumerate all the possible alignments

2. How to sum over all the alignments

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

4. Testing (Inference, decoding):

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

# Score Computation



| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| | | | | | | |
| c | | | | | | |
| a | | | | | | |
| t | | | | | | |

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

$P(h|X)$

$\phi$  c  $\phi$  $\phi$  $a$  $\phi$  $t$  $\phi$  $\phi$

→ Insert $\phi$

↓ output token

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

# Score Computation

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|

$p_{1,0}(\phi)$

$p_{2,0}(c)$

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

c

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

a

$p_{4,2}(\phi)$

$P(h|X)$

$p_{5,2}(t)$

t

$p_{5,3}(\phi)$ $p_{6,3}(\phi)$

# Score Computation



$p_{4,2}(\phi)$

$p_{4,2}(t)$

Because $\phi$ is not considered!

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

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

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| | | | | | | |
| c | | | | | | |
| a | | | | | | |
| t | | | | | | |

$\alpha_{4,1}$

generate "a"

$\alpha_{4,2}$

$\alpha_{3,2}$

read $x^4$

generate "$\phi$",

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

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

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| | | | | | | |
| c | | | | | | |
| a | | | | | | |
| t | | | | | | |

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

# HMM, CTC, RNN-T

*HMM*                               *CTC, RNN-T*

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

1. Enumerate all the possible alignments

2. How to sum over all the alignments

3. Training:    $\boxed{\theta^* = arg\,\max_\theta\, log P_\theta(\hat{Y}|X)}$    $\dfrac{\partial P(\hat{Y}|X)}{\partial \theta} = ?$
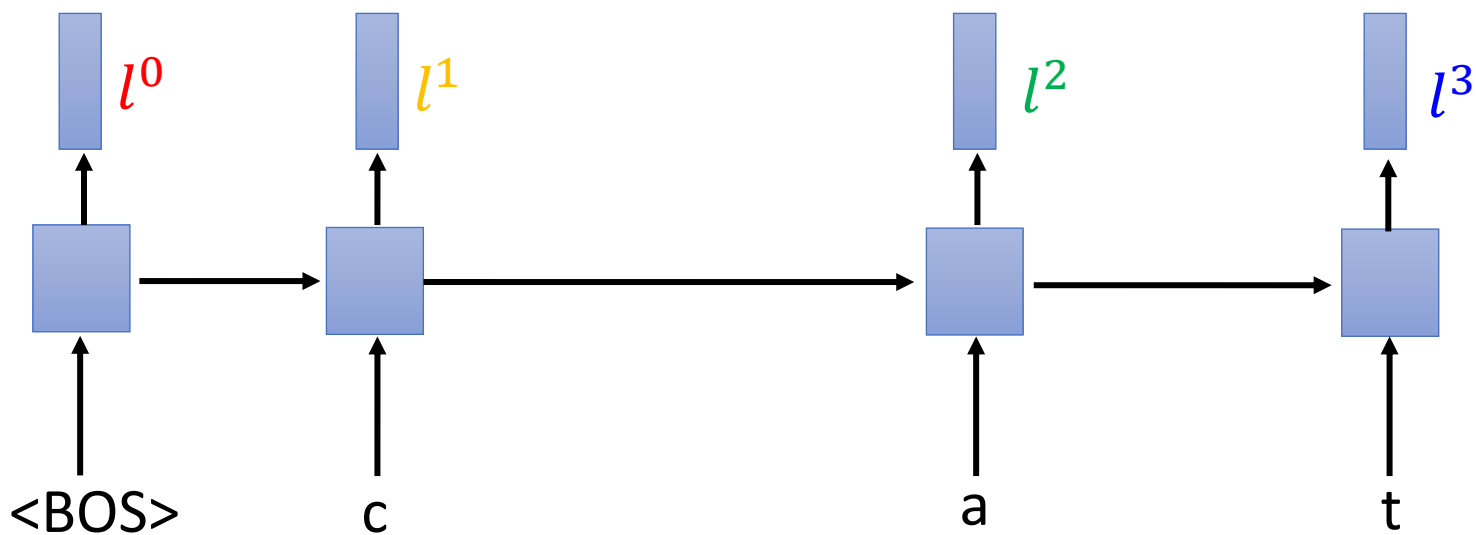
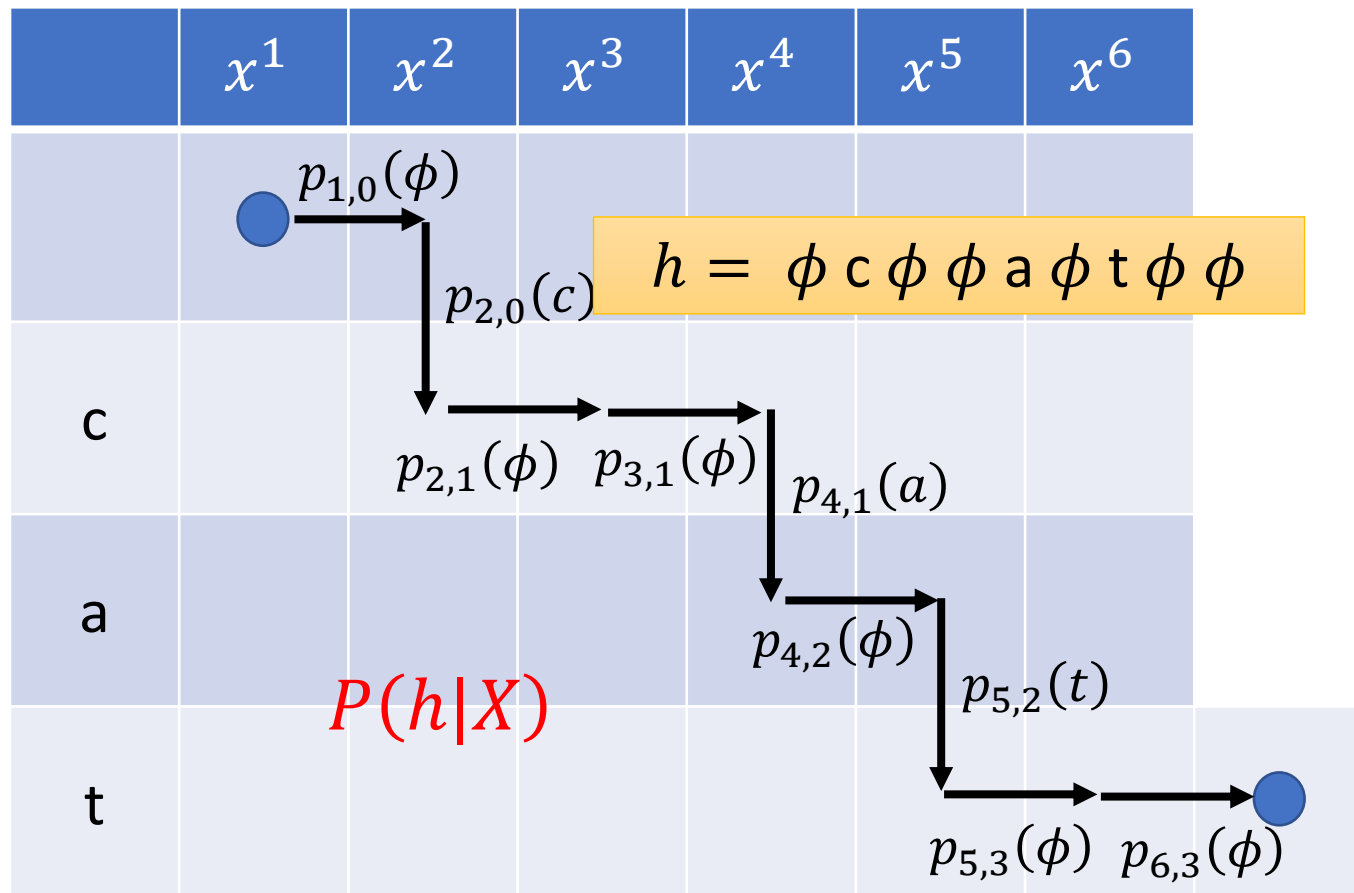4. Testing (Inference, decoding):
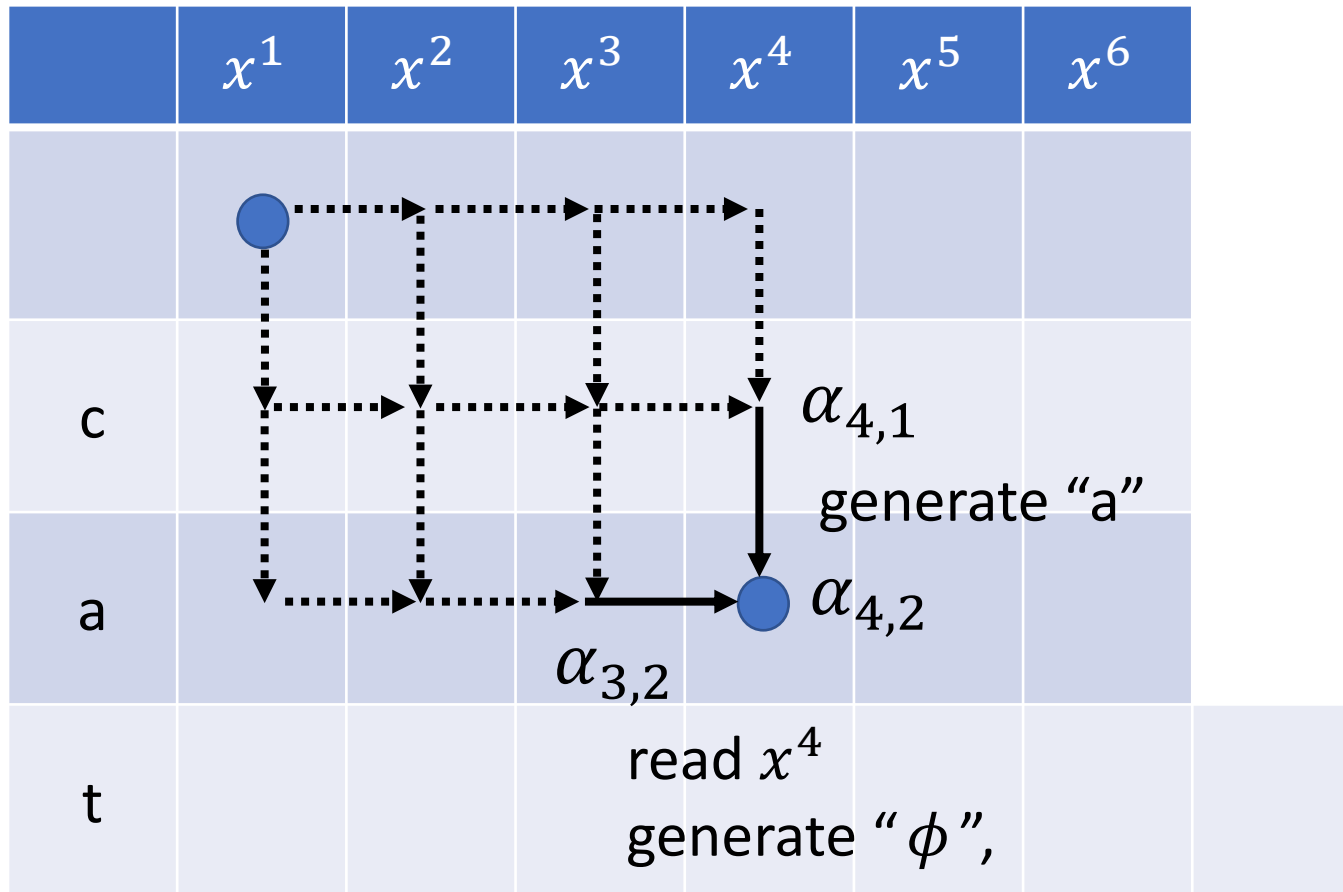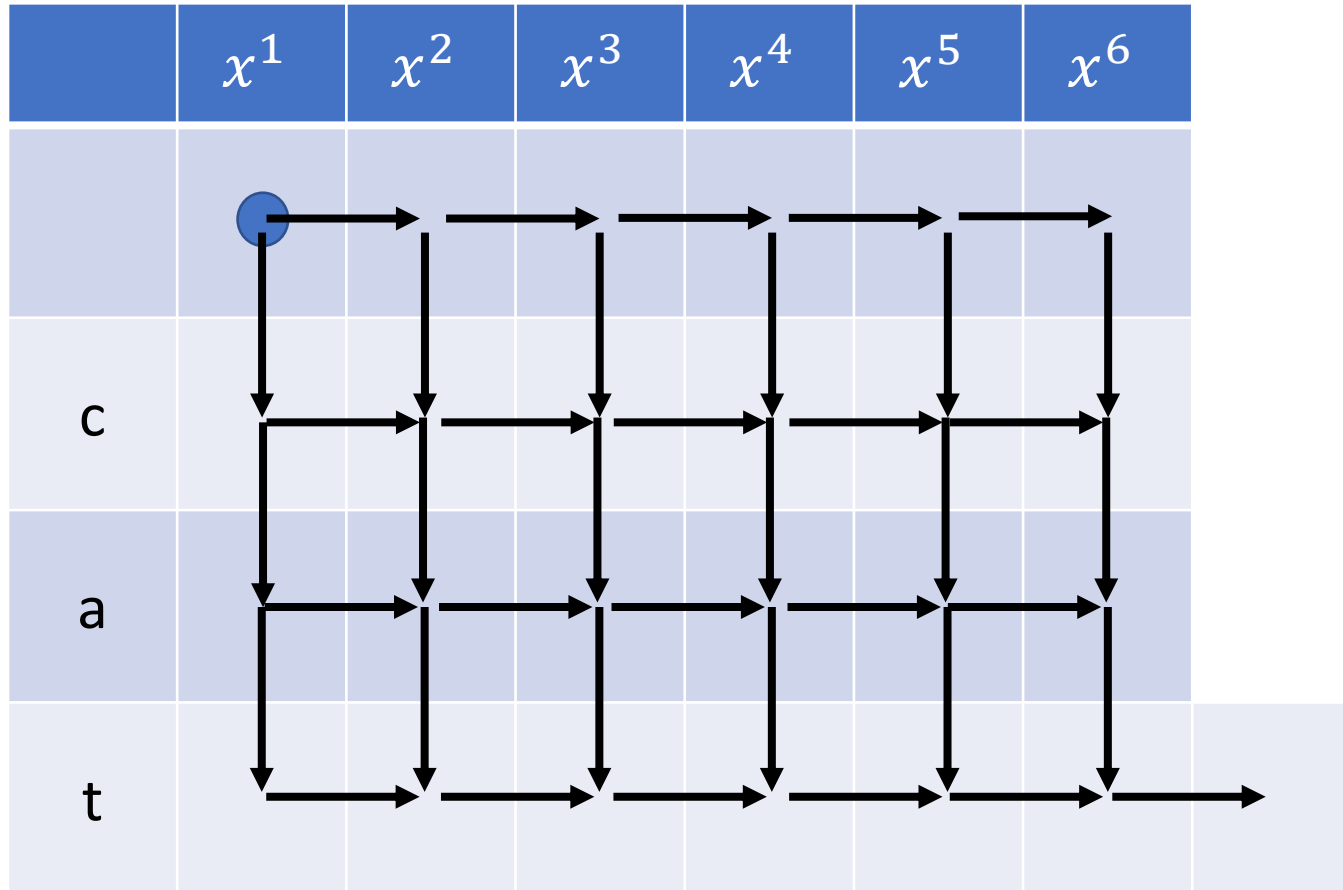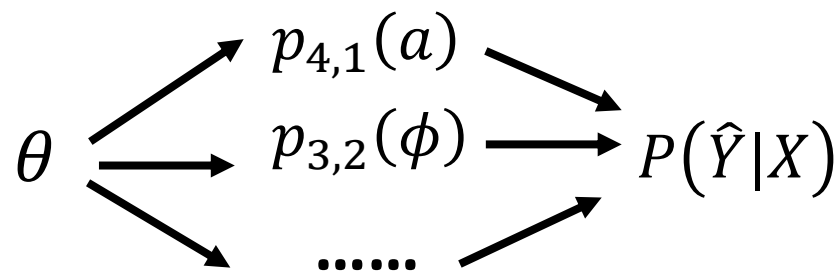
$$Y^* = arg\,\max_Y\, log P_\theta(Y|X)$$

# Training

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



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

$$\phi \; c \; \phi \; \phi \; a \; \phi \; t \; \phi \; \phi$$

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

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

Each arrow is
a component



$$\theta \longrightarrow \begin{matrix} p_{4,1}(a) \\ p_{3,2}(\phi) \\ \cdots\cdots \end{matrix} \longrightarrow P(\hat{Y}|X)$$

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |
|---|---|---|---|---|---|---|
| | | | | | | |
| c | | | | | | |
| a | | | | | | |
| t | | | | | | |

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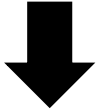
$$\frac{\partial p_{4,1}(a)}{\partial \theta} = ?$$

Backpropagation
(through time)

$$\frac{\partial P(\hat{Y}|X)}{\partial \theta} =? \quad \frac{\partial p_{4,1}(a)}{\partial \theta} \boxed{\frac{\partial P(\hat{Y}|X)}{\partial p_{4,1}(a)}} + \frac{\partial p_{3,2}(\phi)}{\partial \theta} \frac{\partial P(\hat{Y}|X)}{\partial p_{3,2}(\phi)} + \cdots$$
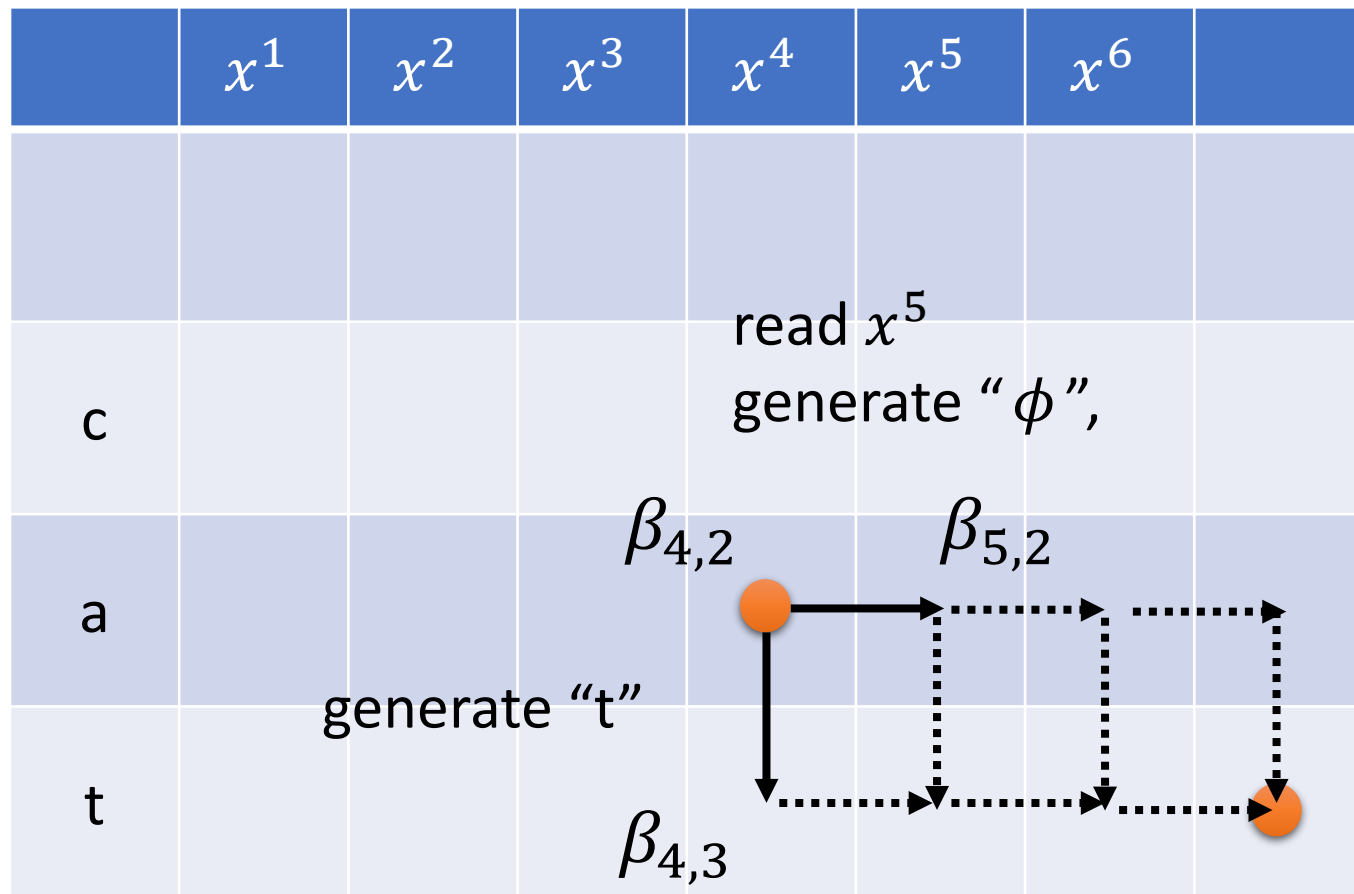
$$P(\hat{Y}|X) = \sum_{h \ with \ p_{4,1}(a)} \boxed{P(h|X)} + \sum_{h \ without \ p_{4,1}(a)} P(h|X)$$

$$\Downarrow$$

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

$$\frac{\partial P(\hat{Y}|X)}{\partial p_{4,1}(a)} = \sum_{h \ with \ p_{4,1}(a)} other \quad = \sum_{h \ with \ p_{4,1}(a)} \frac{P(h|X)}{p_{4,1}(a)}$$

$$= \frac{1}{p_{4,1}(a)} \sum_{h \ with \ p_{4,1}(a)} P(h|X)$$

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

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

| | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |
| c | | | | read $x^5$ generate "$\phi$", | | | |
| a | | | | $\beta_{4,2}$ | $\beta_{5,2}$ | | |
| t | | | generate "t" | $\beta_{4,3}$ | | | |

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



|  | $x^1$ | $x^2$ | $x^3$ | $x^4$ | $x^5$ | $x^6$ |  |
|---|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  |  |
| c |  |  |  | $\alpha_{4,1}$ |  |  |  |
| a |  |  |  | $\beta_{4,2}$ |  |  |  |
| t |  |  |  |  |  |  |  |

$p_{4,1}(a)$

# HMM, CTC, RNN-T

*HMM*                                                   *CTC, RNN-T*

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

1. Enumerate all the possible alignments

2. How to sum over all the alignments

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

4. Testing (Inference, decoding):

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

# Summary

| | LAS | CTC | RNN-T |
|---|---|---|---|
| Decoder | Not independent | independent | Not independent |
| Alignment | Not explicit (Soft alignment) | Yes | Yes |
| Training | Just train it | Sum over alignment | Sum over alignment |
| Streaming | 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