# Fundamentals of
# Hidden Markov Model (HMM)
# (II)  CHMM with applications to speech recognition

Ren-yuan Lyu

Dept of Computer Science & Information Engineering,

Chang Gung University,

Guei-shan, Taoyuan, Taiwan

rylyu@mail.cgu.edu.tw

Reference:
1.    X. Huang, "Spoken Language Processing", Chap 8
2.    L. Rabiner, "Fundamentals of Speech Recognition", Chap 6
3.    HTK Book, http://htk.eng.cam.ac.uk/

# Continuous HMM (CHMM)

- The (state-dependent) observation probability distribution, $\underline{\mathbf{B}} = \{b_i(o)\}$

  By assuming $O(t)$ be a state-dependent random process,

  it is enough to specify $P(O(t) = o \mid S(t) = i)$ to completely describe $O(t)$,
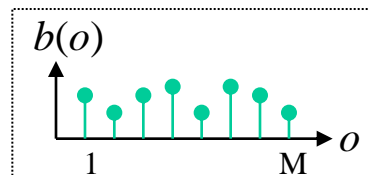
  as long as $S(t)$ is given.

  Let

  $$b_i(o) = P(O(t) = o \mid S(t) = i), \quad i \in \Omega_S = \{1,2,3,...,N\},$$

  $$o \in \Omega_O = R,$$

  $$\underline{B} = \{b_i(o)\},$$

  $$b_i(o) = \sum_{k=1}^{M} c_{ik} \cdot f\left(o; \mu_{ik}, \sigma_{ik}^2\right)$$

  $$\text{where } f\left(o; \mu_{ik}, \sigma_{ik}^2\right) = \frac{1}{\sqrt{2\pi\sigma_{ik}^2}} e^{-\frac{1}{2}\frac{(o-\mu_{ik})^2}{\sigma_{ik}^2}}, \quad \forall i \in \Omega_S$$
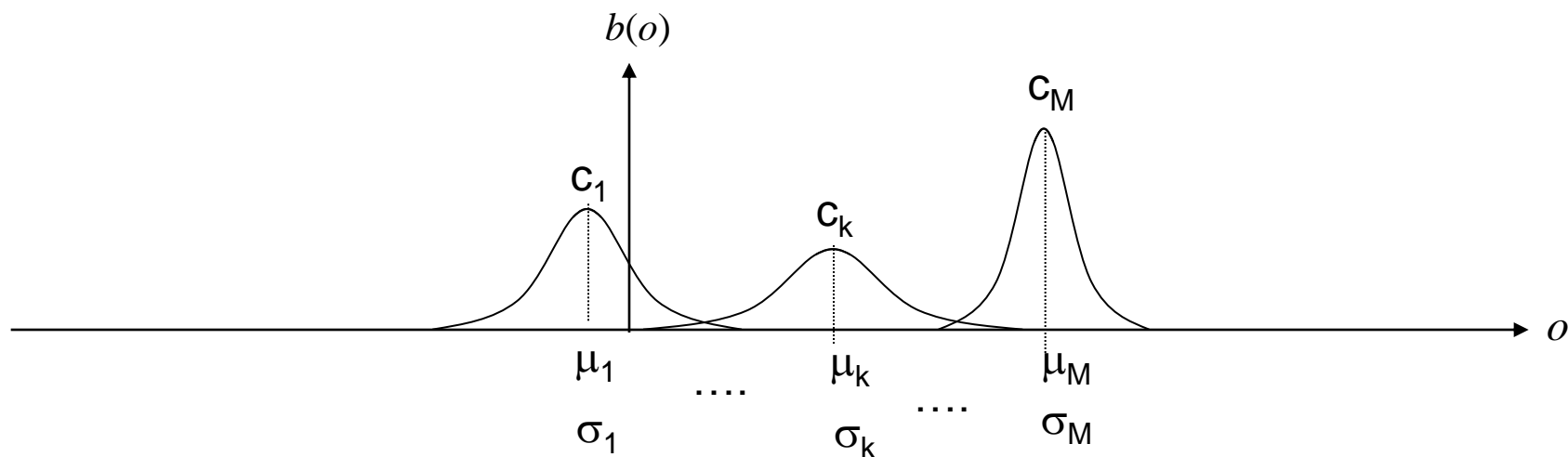
  $$\sum_{k=1}^{M} c_{ik} = 1, \quad \forall i \in \Omega_S$$

$b(o)$

1          M          $o$

for each $b(o), o \in \Omega_O$

we need

$\{ b_k \mid k \in [1, 2, .., M] \}$

$b(o)$

$c_1$          $c_k$          $c_M$

$\mu_1$   ....   $\mu_k$   ....   $\mu_M$          $o$

$\sigma_1$          $\sigma_k$          $\sigma_M$

$$b(o) = \sum_{k=1}^{M} c_k \cdot f\left(o; \mu_k, \sigma_k^2\right)$$

$$\text{where } f\left(o; \mu_k, \sigma_k^2\right) = \frac{1}{\sqrt{2\pi\sigma_k^2}} e^{-\frac{1}{2}\frac{(o-\mu_k)^2}{\sigma_k^2}}$$

for each $b(o), o \in \Omega_O$

we need

$\{ c_k, \mu_k, \sigma_k^2 \mid k \in [1, 2, .., M] \}$

3

$$\alpha_{s_t}(t) \equiv P(o_1, o_2...., o_t, s_t) = \sum_{s_{t-1}=1}^{N} \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t)$$

$$\beta_{s_t}(t) \equiv P(o_{t+1}, o_{t+2},..., o_{T-1}, o_T \mid s_t) = \sum_{s_{t+1}=1}^{N} a_{s_t s_{t+1}} \cdot b_{s_{t+1}}(o_{t+1}) \cdot \beta_{s_{t+1}}(t+1)$$

$$\eta_{s_{t-1}s_t}(t) \equiv P(s_{t-1}, s_t, \underline{O}) = \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot b_{s_t}(o_t) \cdot \beta_{s_t}(t)$$

$$\eta_{s_t}(t) \equiv P(s_t, \underline{O}) = \sum_{s_{t-1}=1}^{N} \eta_{s_{t-1}s_t}(t) = \alpha_{s_t}(t) \cdot \beta_{s_t}(t)$$

$$P(\underline{O}) = \sum_{s_T=1}^{N} \alpha_{s_T}(T) = \beta_{s_0}(0) = \sum_{s_t=1}^{N} \eta_{s_t}(t)$$

$$\gamma_{s_{t-1}s_t}(t) \equiv P(s_{t-1}, s_t \mid \underline{O}) = \eta_{s_{t-1}s_t}(t)/P(\underline{O})$$

$$\gamma_{s_t}(t) \equiv P(s_t \mid \underline{O}) = \eta_{s_t}(t)/P(\underline{O})$$

## The training formula for

c, μ, σ

$$\xi_{s_{t-1}s_t k_t}(t) \equiv P(s_{t-1}, s_t, k_t, \underline{O}) = \alpha_{s_{t-1}}(t-1) \cdot a_{s_{t-1}s_t} \cdot c_{s_t} \cdot f_{s_t k_t}(o_t) \cdot \beta_{s_t}(t)$$

$$\xi_{s_t k_t}(t) \equiv P(s_t, k_t, \underline{O}) = \sum_{s_{t-1}=1}^{N} \xi_{s_{t-1}s_t k_t}(t)$$

$$\zeta_{s_t k_t}(t) \equiv P(s_t, k_t \mid \underline{O}) = \xi_{s_t k_t}(t)/P(\underline{O})$$

$$\zeta_{s_t}(t) \equiv P(s_t \mid \underline{O}) = \sum_{k_t=1}^{M} \zeta_{s_t k_t}(t) \overset{尚未證明}{========} \gamma_{s_t}(t)$$

$$f_{s_t k_t}(o_t) \equiv f\left(o_t; \mu_{s_t k_t}, \sigma^2_{s_t k_t}\right)$$

$k_t$ : the mixture index at time $t$

$f\left(o; \mu, \sigma^2\right)$ : the Gaussian distribution

$$\hat{c}_{jk} = \frac{<\zeta_{jk}(t)>}{<\zeta_j(t)>} = \frac{\frac{1}{T}\sum_{t=1}^{T}\zeta_{jk}(t)}{\frac{1}{T}\sum_{t=1}^{T}\zeta_j(t)} = \frac{\sum_{t=1}^{T}\zeta_{jk}(t)}{\sum_{t=1}^{T}\zeta_j(t)} ========== \frac{Z_{jk}}{\Gamma_j}$$

$$\hat{\mu}_{jk} = \frac{<\zeta_{jk}(t) \cdot o_t>}{<\zeta_{jk}(t)>} = \frac{\sum_{t=1}^{T}\zeta_{jk}(t) \cdot o_t}{\sum_{t=1}^{T}\zeta_{jk}(t)} =============== \frac{M_{jk}}{Z_{jk}}$$

$$\hat{\sigma}^2_{jk} = \frac{<\zeta_{jk}(t) \cdot \left(o_t - \hat{\mu}_{jk}\right)^2>}{<\zeta_{jk}(t)>} = \frac{\sum_{t=1}^{T}\zeta_{jk}(t) \cdot \left(o_t - \hat{\mu}_{jk}\right)^2}{\sum_{t=1}^{T}\zeta_{jk}(t)} ==== \frac{U_{jk}}{Z_{jk}}$$

4

# The training algorithm for c, μ, σ

$$\Gamma_i = \sum_{t=0}^{T-1} \gamma_i(t)$$

$$\Gamma_{ij} = \sum_{t=1}^{T} \gamma_{ij}(t)$$

$$\Gamma_0 = \gamma_0(0)$$

$$\Gamma_{0j} = \gamma_{0j}(1)$$

$$\Gamma_j = \sum_{t=1}^{T} \gamma_j(t)$$

$$\Delta_{jo} = \sum_{t=1}^{T} \gamma_j(t) \cdot \delta(o_t - o)$$

- - - - - - - - - - - - - - - - - - - - -

$$\hat{\pi}_j = \hat{a}_{0j} \mid_{\underline{o}} = \frac{\Gamma_{0j}}{\Gamma_0}$$

$$\hat{a}_{ij} \mid_{\underline{o}} = \frac{\Gamma_{ij}}{\Gamma_i}$$

$$\hat{b}_j(o) \mid_{\underline{o}} = \frac{\Delta_{jo}}{\Gamma_j}$$

$$Z_{jk} = \sum_{t=1}^{T} \zeta_{jk}(t)$$

$$\Rightarrow \hat{c}_{jk} = \frac{Z_{jk}}{\Gamma_j}$$

$$M_{jk} = \sum_{t=1}^{T} \zeta_{jk}(t) \cdot o_t$$

$$\Rightarrow \hat{\mu}_{jk} = \frac{M_{jk}}{Z_{jk}}$$

$$V_{jk} = \sum_{t=1}^{T} \zeta_{jk}(t) \cdot \left(o_t - \hat{\mu}_{jk}\right)^2$$

$$\Rightarrow \hat{\sigma}_{jk}^2 = \frac{V_{jk}}{Z_{jk}}$$

For multi-dimensional observation vector

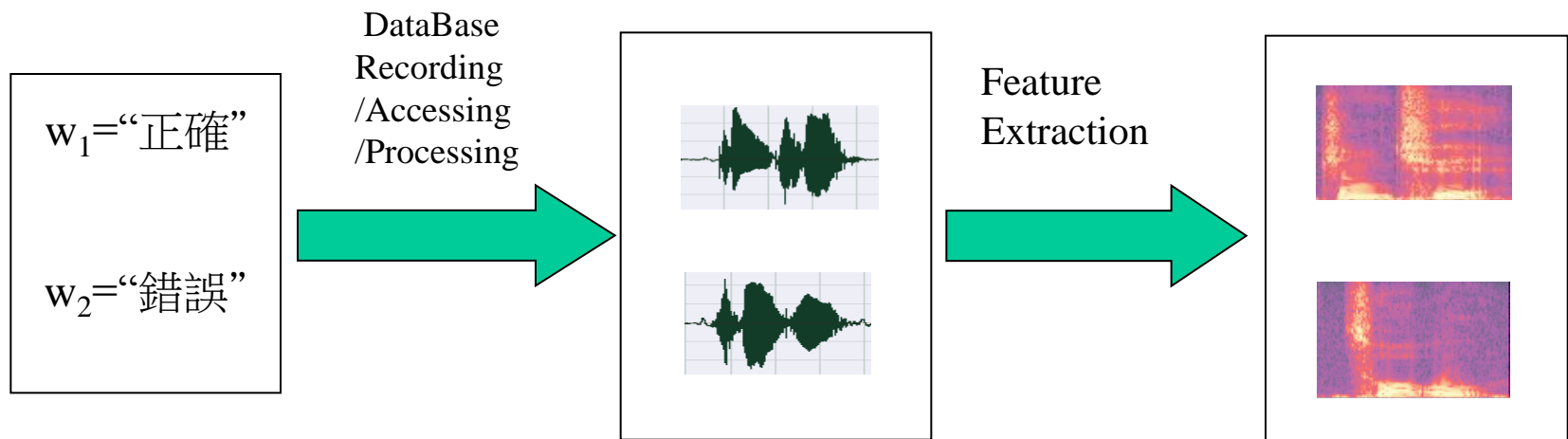$$\vec{M}_{jk} = \sum_{t=1}^{T} \zeta_{jk}(t) \cdot \vec{o}_t$$

$$\Rightarrow \hat{\vec{\mu}}_{jk} = \frac{\vec{M}_{jk}}{Z_{jk}}$$

$$\underline{\underline{V}}_{jk} = \sum_{t=1}^{T} \zeta_{jk}(t) \cdot \left(\vec{o}_t - \hat{\vec{\mu}}_{jk}\right) \cdot \left(\vec{o}_t - \hat{\vec{\mu}}_{jk}\right)^T$$

$$\Rightarrow \underline{\underline{U}}_{jk} = \frac{\underline{\underline{V}}_{jk}}{Z_{jk}}$$

# Isolated Word Recognition using CHMM

- Data Preparation

$w_1$="正確"

$w_2$="錯誤"

DataBase
Recording
/Accessing
/Processing

Feature
Extraction

- Training

$\underline{O}$ for $w_1$

HMM$_1$ for $w_1$

Training

Make $P(O|\pi_1,A_1,B_1)$
As large as possible

$\underline{O}$ for $w_2$

HMM$_2$ for $w_2$

Training

Make $P(O|\pi_2,A_2,B_2)$
As large as possible

7

- Recognition

$$HMM_1 \text{ for } w_1$$



$$P(\underline{O}|\pi_1, A_1, B_1)$$

$$\underline{O}$$



$$i* = \underset{i=1,2}{\text{Argmax}} \ P(\underline{O}|\pi_i, A_i, B_i)$$

$$W = w_{i*}$$

$$HMM_2 \text{ for } w_2$$



$$P(\underline{O}|\pi_2, A_2, B_2)$$

# Continuous Speech Recognition

# Large Vocabulary Speech Recognition

- Sentence = "今天,天氣,不錯"
- Sentence = $\underline{W}$ = $W_1W_2\ldots.W_{(Nw)}$
  - "今天,天氣,不錯" = "今天", "天氣", "不錯"
- Word, W=$\underline{C}$= $C_1C_2\ldots.C_{(Nc)}$
  - "今天" = "今", "天"
- Character, C = $\underline{S}$ = $S_1S_2\ldots.S_{(Ns)}$
  - "今" = "zin"
- Syllable, S = $\underline{P}$ =$P_1P_2\ldots.P_{(Np)}$
  - "zin" = "z", "i", "n"
- Phone, P, has some variations
  - mono-phone,
    » "**z**", "**i**", "**n**"
  - bi-phone,
    » "**z**+i", "**i**+n", "**n**+sil"
  - tri-phone,
    » "sil-**z**+i", "z-**i**+n", "i-**n**+sil"
  - Initial/Final
    » "**z**", "**in**"
    » "**z**+i", "**in**"
    » "**z**+in", "**in**"

要用1個HMM來代表何層次的語言單位？

| | |
|---|---|
| N(Sentence) | = ∞ |
| N(Word) | = 100K |
| N(Character) | = 10K |
| N(Syllable) | = 1K |
| N(Phone) | = |
| N(Mono-phone) | = .1K |
| N(bi-phone) | = .5K |
| N(tri-phone) | = 1 K |
| N(Initial/Final) | = .5K |

# Continuous Syllable Recognition

# Review of Viterbi Algorithm in HMM



$$\alpha_j(0) = \begin{cases} 1, j = 1 \\ 0, j \neq 1, j \in [1..N] \end{cases}$$
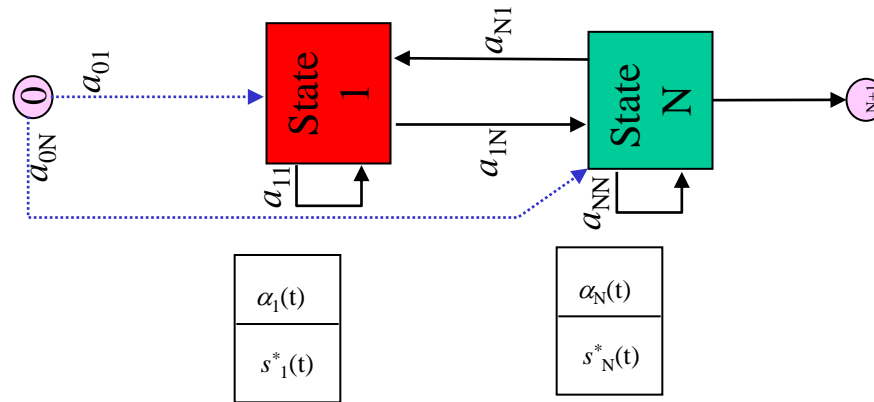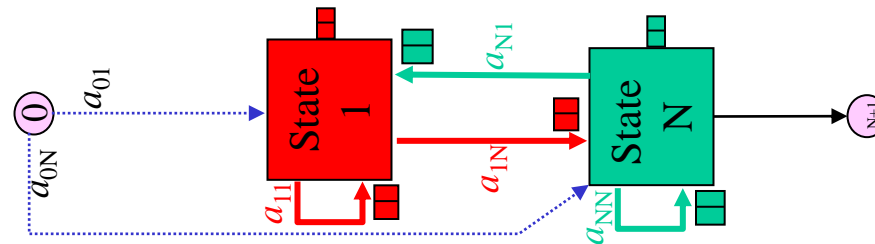
$$s_j^*(t) = i^* = \underset{i \in [1..N]}{ArgMax}\left(\alpha_i(t-1) \cdot a_{ij}\right) \quad , j \in [1..N]$$

$$\alpha_j(t) = \left(\alpha_{i^*}(t-1) \cdot a_{i^*j}\right) \cdot b_j(o_t) \quad , j \in [1..N]$$

$$s_{N+1}^*(t) = N$$

$$\alpha_{N+1}(t) = \left(\alpha_N(t) \cdot a_{N(N+1)}\right)$$

12

# Token Passing



$$s_j^*(t) = i^* = \underset{i \in [1..N]}{ArgMax}\left(\alpha_i(t-1) \cdot a_{ij}\right) \quad , j \in [1..N]$$

$$\alpha_j(t) = \left(\alpha_{i^*}(t-1) \cdot a_{i^* j}\right) \cdot b_j(o_t) \quad , j \in [1..N]$$
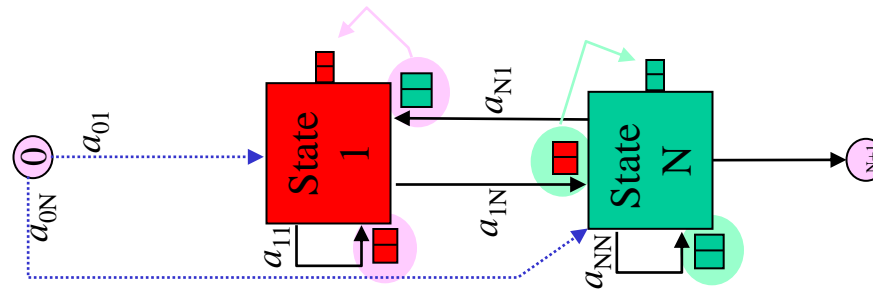
13

For each state i,
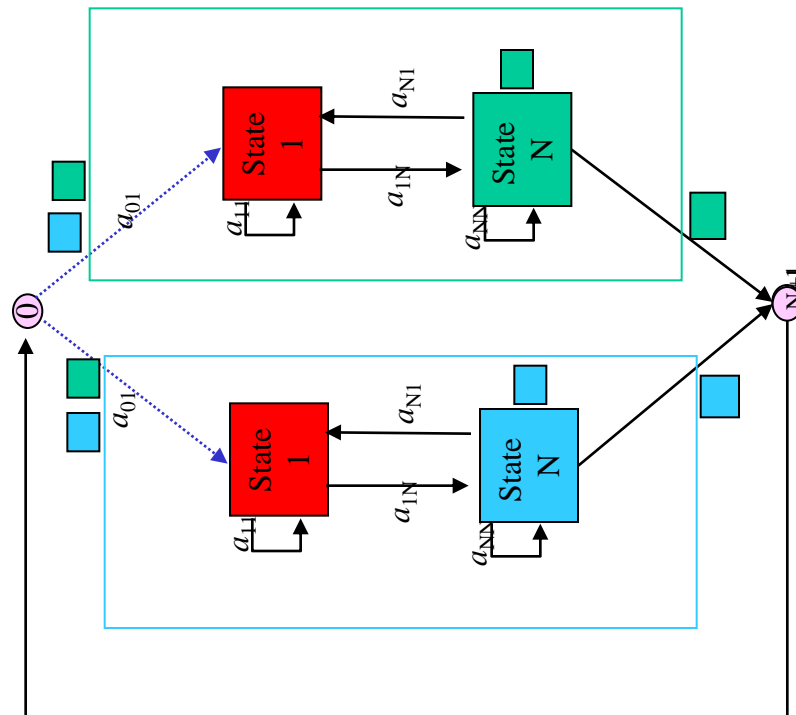
Passing the token of state i to all its connecting state j



For each state j,

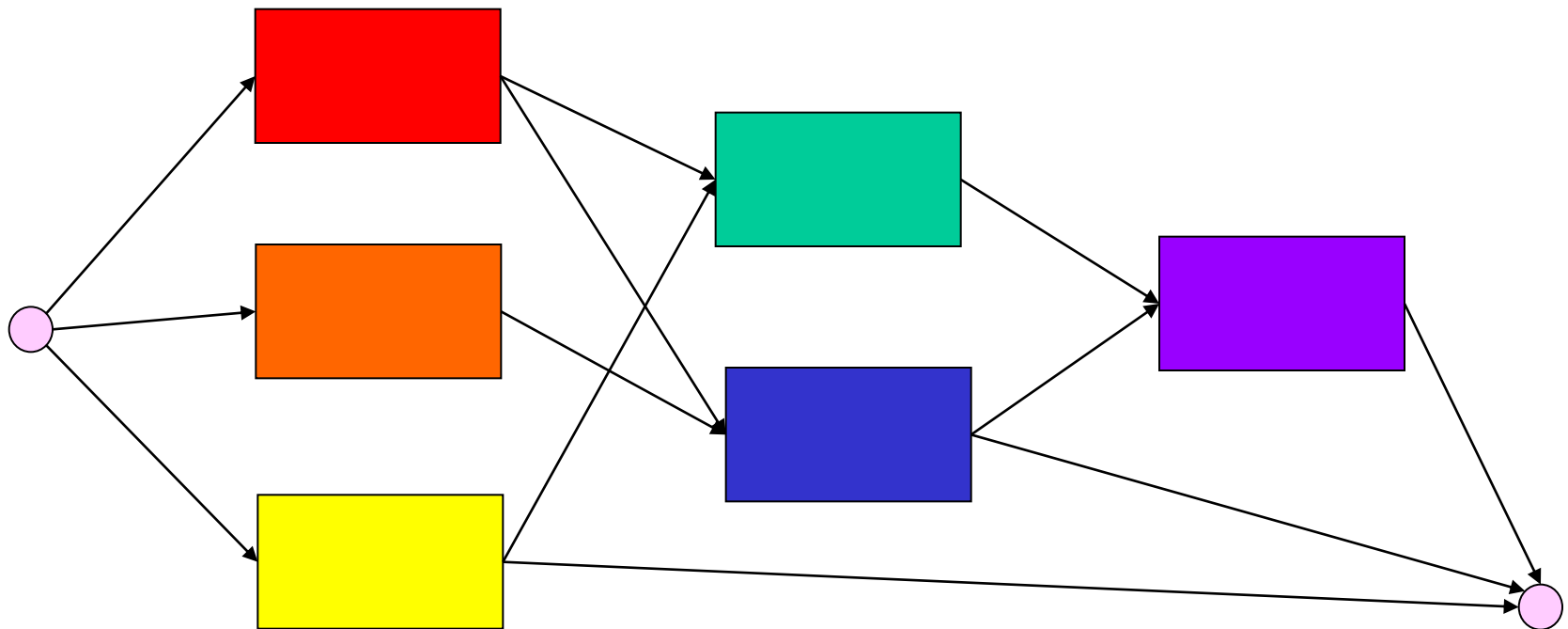Find the best token of all tokens which are passed to state j,
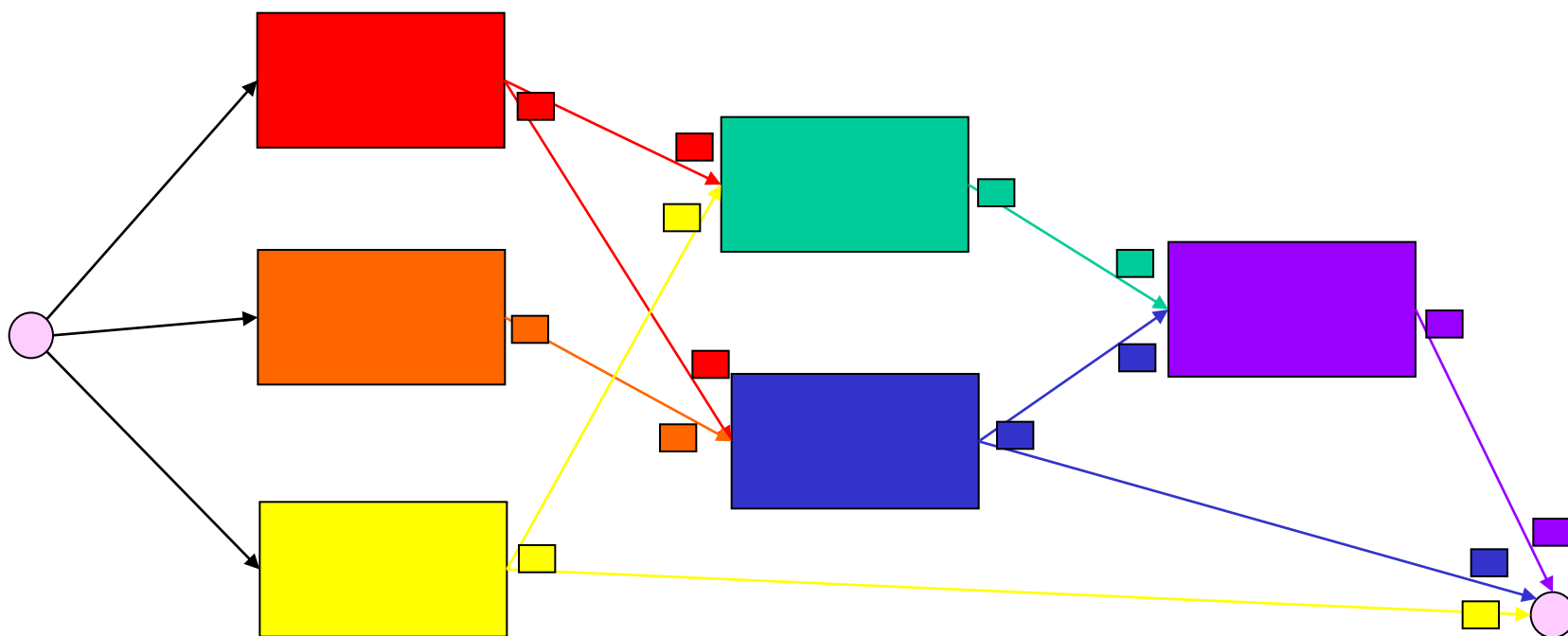
Then update it as the new token of state j

這是一個動態的圖，隨著t的前進，
token 會不斷從各個model框框流出，
再進入各個model，
歷經最佳分數選取的機制，再Update目前token後，
繼續流出…

15

假設根據「文法規則」，model之間可以如下圖連接。



16

隨著$o_1,o_2,\ldots o_t,\ldots$不斷進來，會不斷有token在model間流來流去。



在任意時間點t,這些token會不斷從每個model流出，並流入連接的下一個model，經歷選擇最佳token，並經過model的「消化」後，產生新的token，再流出。

每一個token (at time t, for model w)記載著到時間t為止，

17