

Ling 5801: Lecture Notes 15

From Probability Models to Sequence Models

Contents

15.1 Repeated trials	1
15.2 Interdependence	2
15.3 Example filtering (estimation of last hidden variable)	4
15.4 Most likely sequence estimation	4
15.5 Weighted finite-state automaton	5
15.6 Factored sequence model for speech recognition	7
15.7 Weighted PDA (bounded as hierarchic HMM)	7

15.1 Repeated trials

We can extend probability models to include unbounded number of trials.

This is well defined if models are re-used. This is called *stationarity*.

For example:

$$\langle R \times W_0 \times \dots \times W_T \times O_0 \times \dots \times O_T, 2^{R \times W_0 \times \dots \times W_T \times O_0 \times \dots \times O_T}, \mathbf{P} \rangle$$

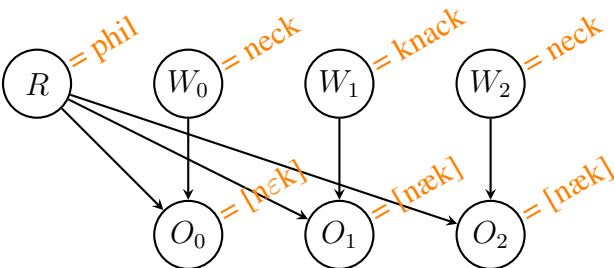
where $R = \{\text{ohio, phil}\}, W_t = \{\text{neck, knack}\}, O_t = \{[\text{nɛk}], [\text{næk}]\},$

$$P_{\theta_R}(R) \stackrel{\text{def}}{=} P(R)$$

$$P_{\theta_W}(W_t | R, W_0, \dots, W_{t-1}) \stackrel{\text{def}}{=} P(W_t)$$

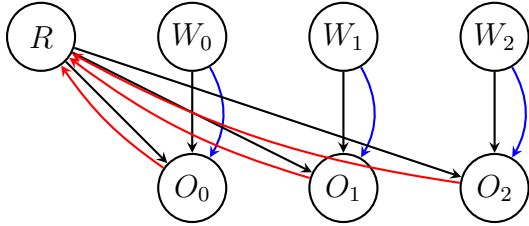
$$P_{\theta_O}(O_t | R, W_0, \dots, W_T, O_0, \dots, O_{t-1}) \stackrel{\text{def}}{=} P(O_t | R, W_t)$$

Graphically (with example values diagonally):



Inference:

$$\begin{aligned}
 P_{\theta_{Wag}}(r) &= \sum_{w_0, w_1, w_2, o_0, o_1, o_2} P_{\theta_{Wag}}(r, w_0, w_1, w_2, o_0, o_1, o_2) \\
 &\stackrel{\text{def}}{=} P(r) \cdot \left(\sum_{o_0, w_0} P(w_0) \cdot P(o_0 | r, w_0) \right) \cdot \left(\sum_{o_1, w_1} P(w_1) \cdot P(o_1 | r, w_1) \right) \cdot \left(\sum_{o_2, w_2} P(w_2) \cdot P(o_2 | r, w_2) \right)
 \end{aligned}$$



Complexity:

$$\mathcal{O}(n \cdot |\mathcal{D}_R| \cdot |\mathcal{D}_W| \cdot |\mathcal{D}_O|)$$

15.2 Interdependence

Models with unbounded number of variables can have variables be interdependent.

This kind of model is still stationary.

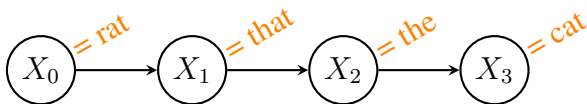
For example:

1. **Markov chain:** a simple sequence of successively dependent variables:

$$\begin{aligned}
 &\langle X_0 \times \dots \times X_T, 2^{X_0 \times \dots \times X_T}, P \rangle \\
 &P_{\theta_{X_0}}(X_0) \stackrel{\text{def}}{=} P(X_0) \\
 &P_{\theta_X}(X_t | X_0, \dots, X_{t-1}) \stackrel{\text{def}}{=} P(X_t | X_{t-1})
 \end{aligned}$$

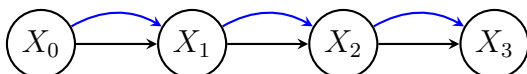
Note π_X for first variable, θ_X subsequently.

Graphically:



Inference:

$$\begin{aligned}
 P_{\theta_{MC}}(x_3) &= \sum_{x_0, x_1, x_2} P_{\theta_{MC}}(x_0, x_1, x_2, x_3) \\
 &\stackrel{\text{def}}{=} \sum_{x_2} \left(\sum_{x_1} \left(\sum_{x_0} P(x_0) \cdot P(x_1 | x_0) \right) \cdot P(x_2 | x_1) \right) \cdot P(x_3 | x_2)
 \end{aligned}$$



Complexity:

$\mathcal{O}(n)$ (if X observed)

2. Hidden Markov model:

seq. of successively dependent hidden variables ea. w. dependent observation:

$$\langle Y_0 \times \dots \times Y_T \times X_0 \times \dots \times X_T, 2^{Y_0 \times \dots \times Y_T \times X_0 \times \dots \times X_T}, \mathbf{P} \rangle$$

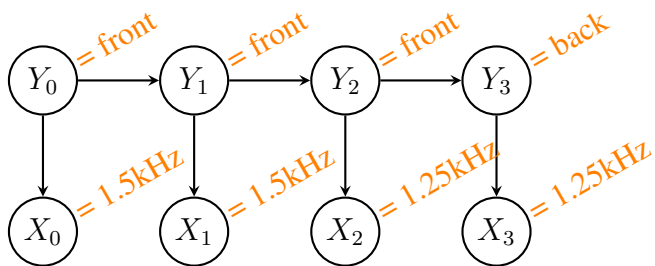
$$P_{\theta_{Y_0}}(Y_0) \stackrel{\text{def}}{=} P(Y_0)$$

$$P_{\theta_Y}(Y_t | Y_0, \dots, Y_{t-1}) \stackrel{\text{def}}{=} P(Y_t | Y_{t-1})$$

$$P_{\theta_X}(X_t | Y_0, \dots, Y_T, X_0, \dots, X_{t-1}) \stackrel{\text{def}}{=} P(X_t | Y_t)$$

Note π_Y for first variable, θ_Y subsequently.

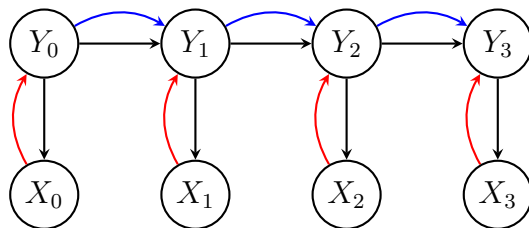
Graphically:



Inference:

$$P_{\theta_{HMM}}(y_3) = \sum_{x_0, x_1, x_2, x_3, y_0, y_1, y_2} P_{\theta_{HMM}}(x_0, x_1, x_2, x_3, y_0, y_1, y_2, y_3)$$

$$\stackrel{\text{def}}{=} \sum_{y_2} \left(\sum_{y_2} \left(\sum_{y_1} \left(\sum_{y_0} \left(P(y_0) \cdot \left(\sum_{x_0} P(x_0 | y_0) \right) \right) \cdot P(y_1 | y_0) \cdot \left(\sum_{x_1} P(x_1 | y_1) \right) \right) \right) \cdot P(y_2 | y_1) \cdot \left(\sum_{x_2} P(x_2 | y_2) \right) \right) \cdot P(y_3 | y_2) \cdot \left(\sum_{x_3} P(x_3 | y_3) \right)$$



Algorithm:

initialize table of possible states at each time step using start states

for y_0 in Y :

$$V[0, y_0] = P_{\pi_Y}(y_0) \cdot P_{\theta_X}(x_0 | y_0)$$

for each possible state y_{t-1} in V at time t , for each y_{t-1}, x_{t-1}, y_t in M , add y_t
for each t in $1..T$:

for each y_{t-1} in Y :

for each y_t in Y :

$$V[t, y_t] = V[t, y_t] + (V[t-1, y_{t-1}] \cdot P_{\theta_Y}(y_t | y_{t-1}) \cdot P_{\theta_X}(x_t | y_t))$$

Complexity:

$$\mathcal{O}(n \cdot |\mathcal{D}_Y| \cdot |\mathcal{D}_Y|) \quad (\text{if } X \text{ observed})$$

15.3 Example filtering (estimation of last hidden variable)

For example, in an HMM for estimating front/back tongue position with parameters:

$$\pi_Y = \begin{array}{|c|c|} \hline \text{front} & \text{back} \\ \hline .5 & .5 \\ \hline \end{array}$$

$$\theta_Y = \begin{array}{|c|c|c|} \hline Y_{t-1} & \text{front} & \text{back} \\ \hline \text{front} & .8 & .2 \\ \text{back} & .2 & .8 \\ \hline \end{array} \quad (\text{encodes inertia in tongue position})$$

$$\theta_X = \begin{array}{|c|c|c|c|c|c|c|c|} \hline Y_t & .5\text{kHz} & .75\text{kHz} & 1\text{kHz} & 1.25\text{kHz} & 1.5\text{kHz} & 1.75\text{kHz} & 2\text{kHz} \\ \hline \text{front} & 0 & 0 & .1 & .2 & .4 & .2 & .1 \\ \text{back} & .1 & .2 & .4 & .2 & .1 & 0 & 0 \\ \hline \end{array}$$

we obtain forward messages calculating $P(Y_t, X_{1..t})$, marginalizing out all $Y_{1..t-1}$:

t	X_t	$P(Y_t = \text{front}, X_{1..t})$	$P(Y_t = \text{back}, X_{1..t})$
0	1.5kHz	$.5 \cdot .4 = .2$	$.5 \cdot .1 = .05$
1	1.5kHz	$.2 \cdot .8 \cdot .4 + .05 \cdot .2 \cdot .4 = .068$	$.2 \cdot .2 \cdot .1 + .05 \cdot .8 \cdot .1 = .008$
2	1.25kHz	$.068 \cdot .8 \cdot .2 + .008 \cdot .2 \cdot .2 = .0112$	$.068 \cdot .2 \cdot .2 + .008 \cdot .8 \cdot .2 = .004$

15.4 Most likely sequence estimation

Subst semiring: $\langle \mathbb{R}_0^\infty, +, \cdot, 0, 1 \rangle$ to $\langle \mathbb{R}_0^\infty \times (\mathcal{D}_Y \times \mathcal{D}_X)^*, \text{max-argmax}, \text{prod-pair}, \langle 0, \langle \rangle \rangle, \langle 1, \langle \rangle \rangle \rangle$

where:

$$p, x \text{ max-argmax } q, y = \begin{cases} \text{if } p > q & : p, x \\ \text{otherwise} & : q, y \end{cases}$$

$$p, x \text{ prod-pair } q, y = p \cdot q, \langle x, y \rangle$$

Algorithm:

initialize table of possible states at each time step using start states

for y_0 in Y :

$$V[0, y_0] = P_{\pi_Y}(y_0), y_0 \text{ prod-pair } P_{\theta_X}(x_0 | y_0), x_0$$

for each possible state y_{t-1} in V at time t , for each y_{t-1}, x_{t-1}, y_t in M , add y_t

for each t in $1..T$:

for each y_{t-1} in Y :

for each y_t in Y :

$$V[t, y_t] = V[t, y_t] \text{ max-argmax } (V[t-1, y_{t-1}] \text{ prod-pair } P_{\theta_Y}(y_t | y_{t-1}), y_t \text{ prod-pair } P_{\theta_X}(x_t | y_t), x_t)$$

Complexity:

$$\mathcal{O}(n \cdot |\mathcal{D}_Y| \cdot |\mathcal{D}_Y|) \quad (\text{if } X \text{ observed})$$

Example: obtain Viterbi MLS messages w. maximal $P(y_{1..t}, x_{1..t})$ ending at each y_t :

t	X_t	$Y_t = \text{front}$	$Y_t = \text{back}$
0	1.5	.5 · .4 = .2, ⟨front, 1.5⟩	.5 · .1 = .05, ⟨back, 1.5⟩
1	1.5	max(.2 · .8 · .4, .05 · .2 · .4) = .064, ⟨⟨front, 1.5⟩, ⟨front, 1.5⟩⟩	max(.2 · .2 · .1, .05 · .8 · .1) = .004, ⟨⟨back, 1.5⟩, ⟨back, 1.5⟩⟩
2	1.25	max(.064 · .8 · .2, .004 · .2 · .2) = .01024, ⟨⟨⟨front, 1.5⟩, ⟨front, 1.5⟩⟩, ⟨front, 1.25⟩⟩	max(.064 · .2 · .2, .004 · .8 · .2) = .00256, ⟨⟨⟨front, 1.5⟩, ⟨front, 1.5⟩⟩, ⟨back, 1.25⟩⟩

15.5 Weighted finite-state automaton

1. Semiring substitution from $\langle \{T, F\}, \vee, \wedge, F, T \rangle$ to $\langle \mathbb{R}_0^\infty, +, \cdot, 0, 1 \rangle$ allows weighted FSA:

initialize table of possible states at each time step using start states

for each q_t in Q :

$$V[0, q_t] = S[q_t]$$

for each possible state q_{t-1} in V at time t , for each q_{t-1}, x_{t-1}, q_t in M , add q_t

for each t in $1..T$:

for each q_{t-1} in Q :

for each q_t in Q :

$$V[t, q_t] = V[t, q_t] + (V[t-1, q_{t-1}] \cdot M[q_{t-1}, x_{t-1}, q_t])$$

2. Now replacing S with probabilities π_Q and factoring M into θ_X and θ_Q gives:

initialize table of possible states at each time step using start states

for q_t in Q :

$$V[0, q_t] = P_{\pi_Q}(q_t)$$

for each possible state q_{t-1} in V at time t , for each q_{t-1}, x_{t-1}, q_t in M , add q_t

for each t in $1..T$:

for each q_{t-1} in Q :

for each q_t in Q :

$$V[t, q_t] = V[t, q_t] + (V[t-1, q_{t-1}] \cdot P_{\theta_X}(x_{t-1} | q_{t-1}) \cdot P_{\theta_Q}(q_t | q_{t-1}, x_{t-1}))$$

where: $P_{\theta_X}(x_{t-1} | q_{t-1}) = \sum_{q_t} M[q_{t-1}, x_{t-1}, q_t]$

$$P_{\theta_Q}(q_t | q_{t-1}, x_{t-1}) = \frac{M[q_{t-1}, x_{t-1}, q_t]}{\sum_{q_t} M[q_{t-1}, x_{t-1}, q_t]}$$

3. This recognizer can be expressed as a probability model:

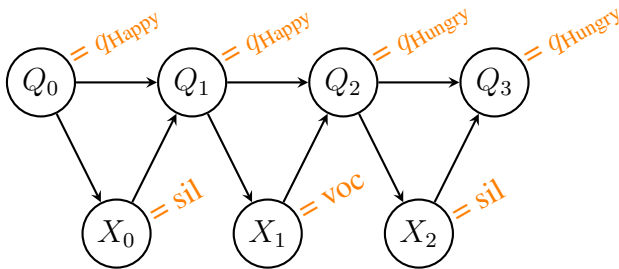
$$\langle Q_0 \times X_0 \times Q_1 \times \dots \times X_{T-1} \times Q_T, 2^{Q_0 \times X_0 \times Q_1 \times \dots \times X_{T-1} \times Q_T}, P \rangle$$

$$P_{\theta_{Q_0}}(Q_0) \stackrel{\text{def}}{=} P(Q_0)$$

$$P_{\theta_X}(X_t | Q_0, X_0, Q_1, \dots, X_{t-1}, Q_t) \stackrel{\text{def}}{=} P(X_t | Q_t)$$

$$P_{\theta_Q}(Q_t | Q_0, X_0, Q_1, \dots, X_{t-1}, Q_{t-1}, X_{t-1}) \stackrel{\text{def}}{=} P(Q_t | Q_{t-1}, X_{t-1})$$

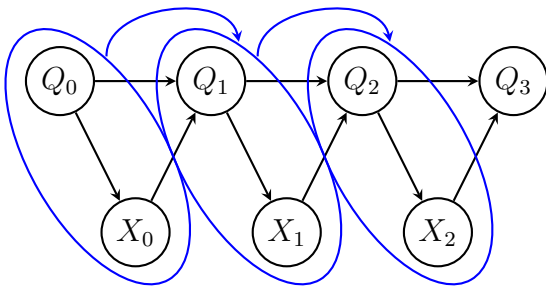
Graphically:



Inference:

$$P_{\theta_{WFSA}}(q_3) = \sum_{x_0, x_1, x_2, q_0, q_1, q_2} P_{\theta_{WFSA}}(x_0, x_1, x_2, q_0, q_1, q_2, q_3)$$

$$\stackrel{\text{def}}{=} \sum_{q_2} \left(\sum_{q_2, x_2} \left(\sum_{q_1, x_1} \left(\sum_{q_0, x_0} \left(P(q_0) \cdot P(x_0 | q_0) \right) \cdot P(q_1 | q_0, x_0) \cdot P(x_1 | q_1) \right) \cdot P(q_2 | q_1, x_1) \cdot P(x_2 | q_2) \right) \cdot P(q_3 | q_2, x_2) \right)$$

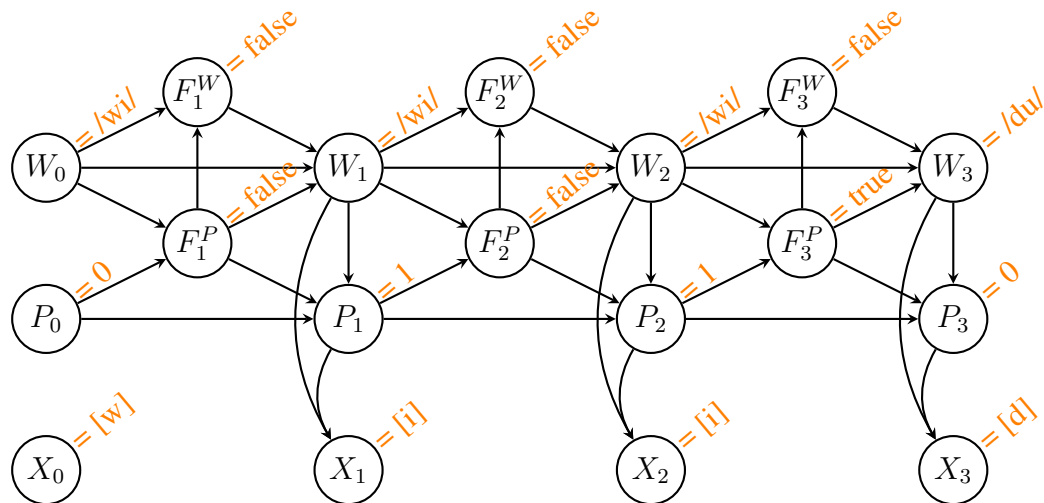


Complexity:

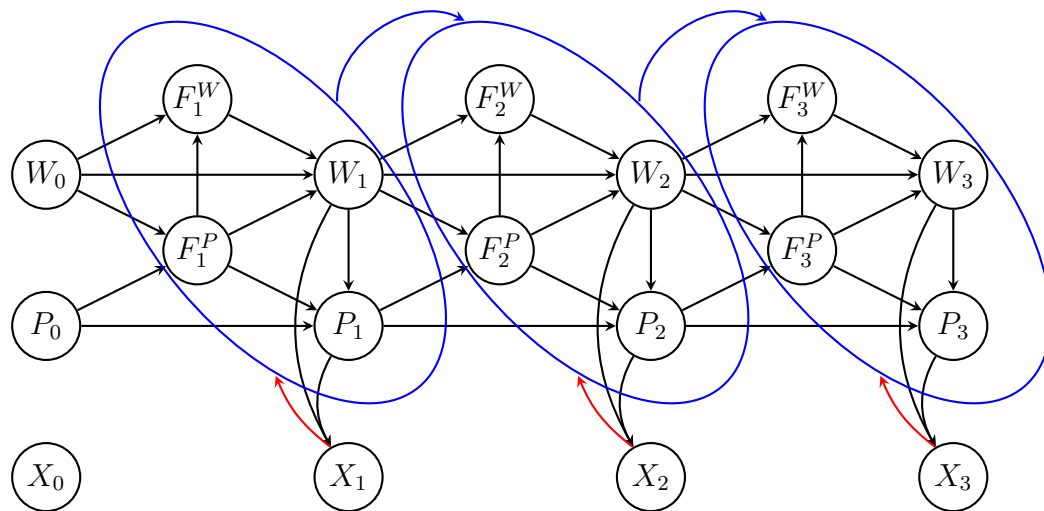
$$\mathcal{O}(n \cdot |D_Q| \cdot |D_Q|) \quad (\text{if } X \text{ observed})$$

15.6 Factored sequence model for speech recognition

Synchronize high-/low-level sequences with true/false 'final state' variables:



Inference:



Complexity:

$$\mathcal{O}(n \cdot |\mathcal{D}_W| \cdot |\mathcal{D}_P| \cdot |\mathcal{D}_F| \cdot |\mathcal{D}_F| \cdot |\mathcal{D}_W| \cdot |\mathcal{D}_P|) \quad (\text{if } X \text{ observed})$$

15.7 Weighted PDA (bounded as hierarchic HMM)

1. Semiring substitution from $\langle \{T, F\}, \vee, \wedge, F, T \rangle$ to $\langle \mathbb{R}_0^\infty, +, \cdot, 0, 1 \rangle$ allows weighted PDA:

for each time step t :

for each previous state q_{t-1}^d and store $q_{t-1}^{1..d-1}$ (where $d-1$ is the depth of the store):

i) for each final state f_t^d and current state q_t^d : (expand +) trans

$$V[t, q_t^d q_{t-1}^{1..d-1}] = V[t, q_t^d q_{t-1}^{1..d-1}] + (V[t-1, q_{t-1}^d q_{t-1}^{1..d-1}] \cdot M[q_{t-1}^d, \epsilon, x_{t-1}, f_t^d, \epsilon] \cdot M[f_t^d, q_{t-1}^{d-1}, \epsilon, q_t^d, q_{t-1}^d])$$

ii) for each previous state at deeper level q_{t-1}^{d+1} , final state f_t^{d+1} and current state q_t^{d+1} :

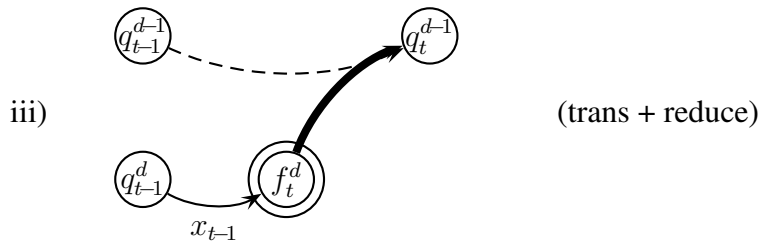
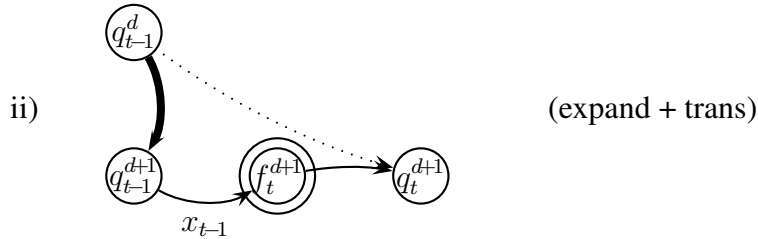
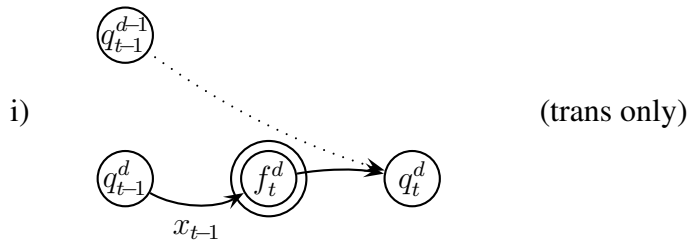
$$V[t, q_t^{d+1} q_{t-1}^{1..d-1}] = V[t, q_t^{d+1} q_{t-1}^{1..d-1}] + (V[t-1, q_{t-1}^d q_{t-1}^{1..d-1}] \cdot M[q_{t-1}^d, \epsilon, \epsilon, q_{t-1}^{d+1}, q_{t-1}^d] \cdot M[q_{t-1}^{d+1}, \epsilon, x_{t-1}, f_t^{d+1}, \epsilon] \cdot M[f_t^{d+1}, q_{t-1}^d, \epsilon, q_t^{d+1}, q_{t-1}^d])$$

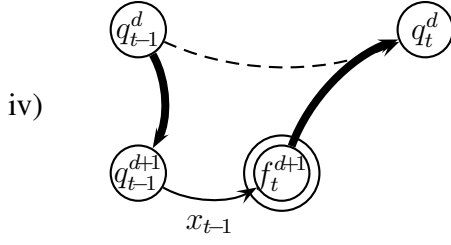
iii) for each final state f_t^d and current state q_t^{d-1} : (expand +) trans + reduce

$$V[t, q_t^{d-1} q_{t-1}^{1..d-2}] = V[t, q_t^{d-1} q_{t-1}^{1..d-2}] + (V[t-1, q_{t-1}^d q_{t-1}^{1..d-2}] \cdot M[q_{t-1}^d, \epsilon, x_{t-1}, f_t^d, \epsilon] \cdot M[f_t^d, q_{t-1}^{d-1}, \epsilon, q_t^{d-1}, \epsilon])$$

iv) for each previous state at deeper level q_{t-1}^{d+1} and final state f_t^{d+1} and current state q_t^d :

$$V[t, q_t^d q_{t-1}^{1..d-2}] = V[t, q_t^d q_{t-1}^{1..d-2}] + (V[t-1, q_{t-1}^d q_{t-1}^{1..d-2}] \cdot M[q_{t-1}^d, \epsilon, \epsilon, q_{t-1}^{d+1}, q_{t-1}^d] \cdot M[q_{t-1}^{d+1}, \epsilon, x_{t-1}, f_t^{d+1}, \epsilon] \cdot M[f_t^{d+1}, q_{t-1}^d, \epsilon, q_t^d, \epsilon])$$





(expand + trans + reduce)

2. Now replacing S with probabilities π_Q and factoring M into θ_X , θ_F , and θ_Q gives:

for each time step t :

for each previous state q_{t-1}^d and store $q_{t-1}^{1..d-1}$ (where $d-1$ is the depth of the store):

i) for each final state f_t^d and current state q_t^d : (expand +) trans

$$V[t, q_t^d q_{t-1}^{1..d-1}] = V[t, q_t^d q_{t-1}^{1..d-1}] + (V[t-1, q_{t-1}^d q_{t-1}^{1..d-1}] \cdot P_{\theta_X}(x_{t-1} | q_{t-1}^d) \cdot P_{\theta_F}(f_t^d | -, -, q_{t-1}^d, x_{t-1}) \cdot P_{\theta_Q}(q_t^d | f_t^d, -, q_{t-1}^d, -))$$

ii) for each previous state at deeper level q_{t-1}^{d+1} , final state f_t^{d+1} and current state q_t^{d+1} :

$$V[t, q_t^{d+1} q_t^d q_{t-1}^{1..d-1}] = V[t, q_t^{d+1} q_t^d q_{t-1}^{1..d-1}] + (V[t-1, q_{t-1}^d q_{t-1}^{1..d-1}] \cdot P_{\theta_X}(x_{t-1} | q_{t-1}^d) \cdot P_{\theta_F}(f_t^{d+1} | -, q_{t-1}^d, -, x_{t-1}) \cdot P_{\theta_Q}(q_t^{d+1} | f_t^{d+1}, -, q_{t-1}^d, -))$$

iii) for each final state f_t^d and current state q_t^{d-1} : (expand +) trans + reduce

$$V[t, q_t^{d-1} q_{t-1}^{1..d-2}] = V[t, q_t^{d-1} q_{t-1}^{1..d-2}] + (V[t-1, q_{t-1}^d q_{t-1}^{d-1} q_{t-1}^{1..d-2}] \cdot P_{\theta_X}(x_{t-1} | q_{t-1}^d) \cdot P_{\theta_F}(f_t^d | -, -, q_{t-1}^d, x_{t-1}) \cdot P_{\theta_Q}(q_t^{d-1} | -, f_t^d, -, q_{t-1}^d))$$

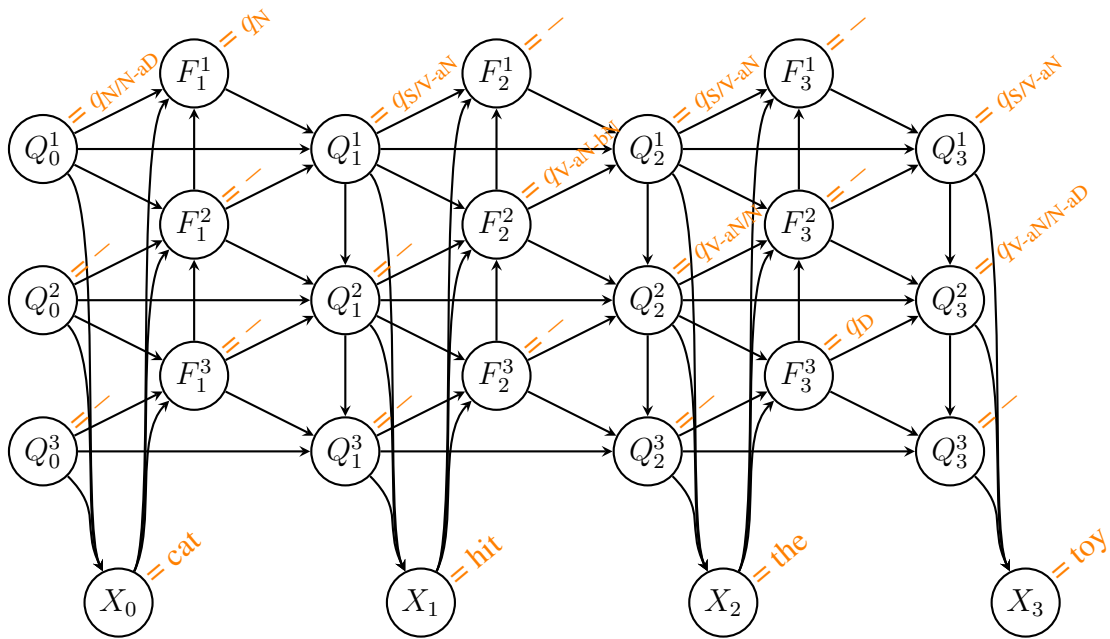
iv) for each previous state at deeper level q_{t-1}^{d+1} and final state f_t^{d+1} and current state q_t^d :

$$V[t, q_t^d q_t^{d-1} q_{t-1}^{1..d-2}] = V[t, q_t^d q_t^{d-1} q_{t-1}^{1..d-2}] + (V[t-1, q_{t-1}^d q_{t-1}^{d-1} q_{t-1}^{1..d-2}] \cdot P_{\theta_X}(x_{t-1} | q_{t-1}^d) \cdot P_{\theta_F}(f_t^{d+1} | -, q_{t-1}^d, -, x_{t-1}) \cdot P_{\theta_Q}(q_t^d | -, f_t^{d+1}, -, q_{t-1}^d))$$

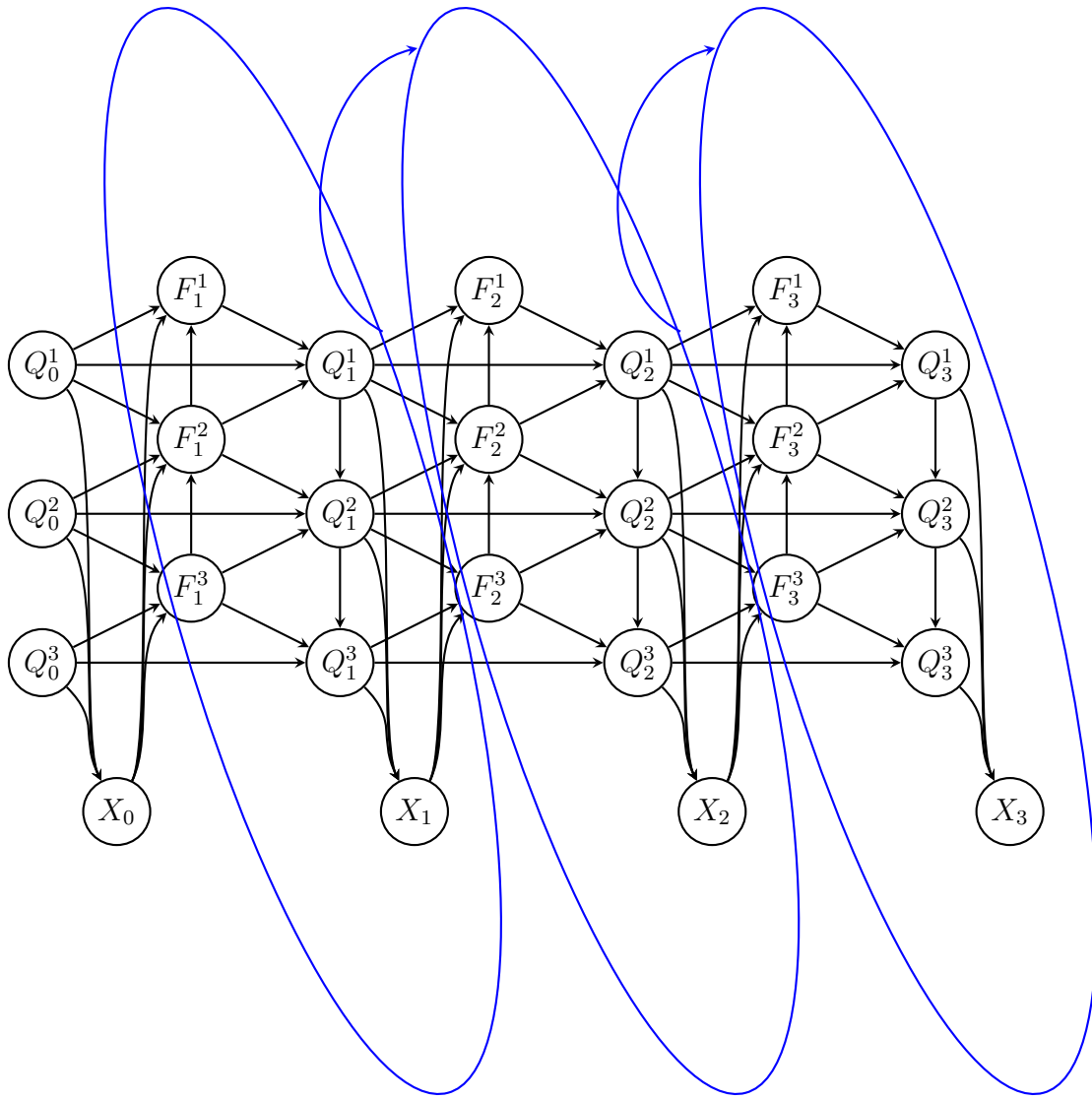
where:

$$\begin{aligned} P_{\theta_X}(x_{t-1} | q_{t-1}^d) &= \sum_{f_t^d} M[q_{t-1}^d, \epsilon, x_{t-1}, f_t^d, \epsilon] \\ &\quad + \sum_{f_t^{d+1}, q_{t-1}^{d+1}} M[q_{t-1}^d, \epsilon, \epsilon, q_{t-1}^{d+1}, q_{t-1}^d] \cdot M[q_{t-1}^{d+1}, \epsilon, x_{t-1}, f_t^{d+1}, \epsilon] \\ P_{\theta_F}(f_t^d | -, -, q_{t-1}^d, x_{t-1}) &= M[q_{t-1}^d, \epsilon, x_{t-1}, f_t^d, \epsilon] \\ P_{\theta_F}(f_t^{d+1} | -, q_{t-1}^d, -, x_{t-1}) &= \sum_{q_{t-1}^{d+1}} M[q_{t-1}^d, \epsilon, \epsilon, q_{t-1}^{d+1}, q_{t-1}^d] \cdot M[q_{t-1}^{d+1}, \epsilon, x_{t-1}, f_t^{d+1}, \epsilon] \\ P_{\theta_Q}(q_t^d | -, f_t^{d+1}, -, q_{t-1}^d) &= M[f_t^{d+1}, q_{t-1}^d, \epsilon, q_{t-1}^d, \epsilon] \\ P_{\theta_Q}(q_t^d | f_t^d, -, q_{t-1}^d, -) &= M[f_t^d, q_{t-1}^d, \epsilon, q_t^d, q_{t-1}^d] \text{ where } q_{t-1}^{d-1} = q_{t-1}^d \end{aligned}$$

3. This recognizer can be expressed as a probability model:



Inference:



Complexity:

$$\mathcal{O}(n \cdot |\mathcal{D}_Q|^D \cdot |\mathcal{D}_F|^D \cdot |\mathcal{D}_Q|^D) \text{ (if } X \text{ observed)}$$

Some observations:

- generates incremental probabilities (good for predicting reading times)
- connects one rule per word (good for efficient/human-like semantic composition)