

# **Deep Learning**

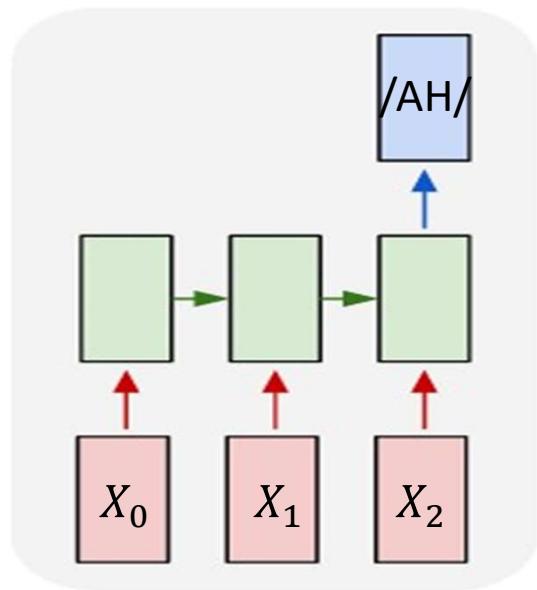
## **Sequence to Sequence models:**

### **Connectionist Temporal Classification**

# Sequence-to-sequence modelling

- Problem:
  - A sequence  $X_1 \dots X_N$  goes in
  - A different sequence  $Y_1 \dots Y_M$  comes out
- E.g.
  - Speech recognition: Speech goes in, a word sequence comes out
    - Alternately output may be phoneme or character sequence
  - Machine translation: Word sequence goes in, word sequence comes out
  - Dialog : User statement goes in, system response comes out
  - Question answering : Question comes in, answer goes out
- In general  $N \neq M$ 
  - No synchrony between  $X$  and  $Y$ .

# Basic model

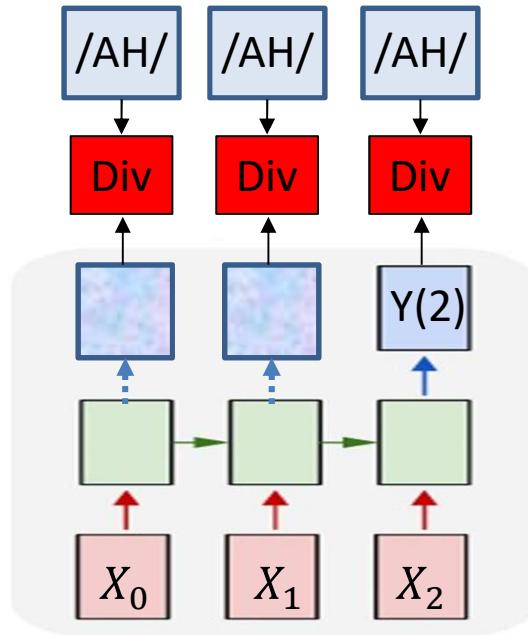


- Sequence of inputs produces a single output

# Training

Fix: Use these outputs too.

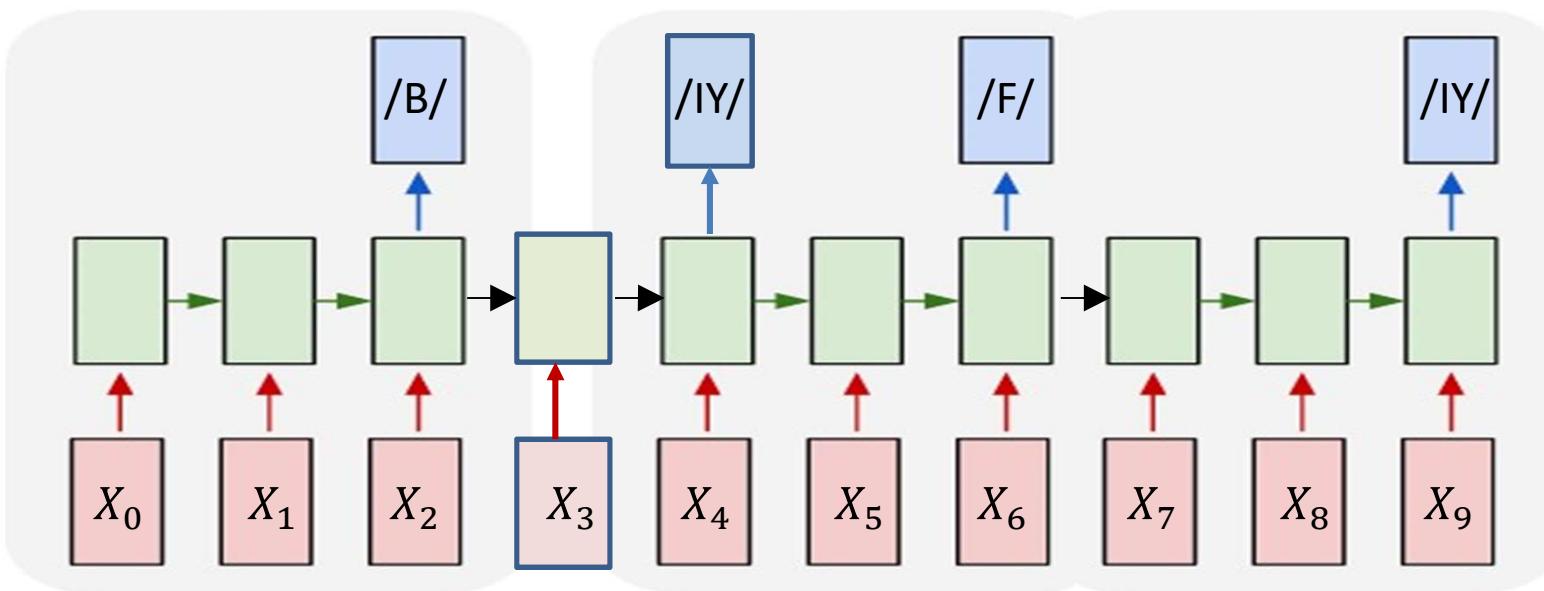
These too must ideally point to the correct phoneme



- Exploiting the untagged inputs: assume the same output for the entire input
- Define the divergence everywhere

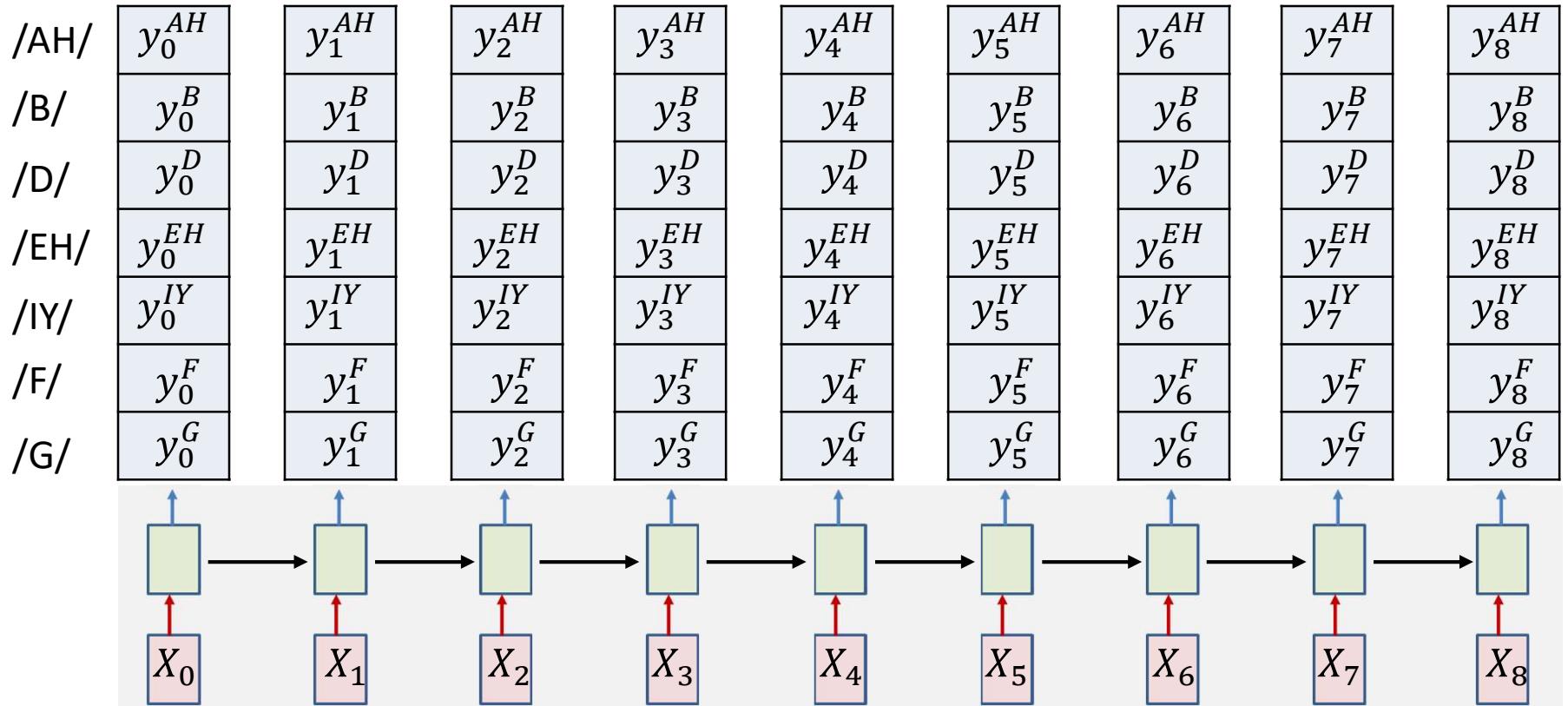
$$DIV(Y_{target}, Y) = \sum_t w_t Xent(Y(t), Phoneme)$$

# The more complex problem



- Objective: Given a sequence of inputs, asynchronously output a sequence of symbols
  - This is just a simple concatenation of many copies of the simple “output at the end of the input sequence” model we just saw
- But this simple extension complicates matters..
  - *When* do you output symbols for a continuous stream of input?

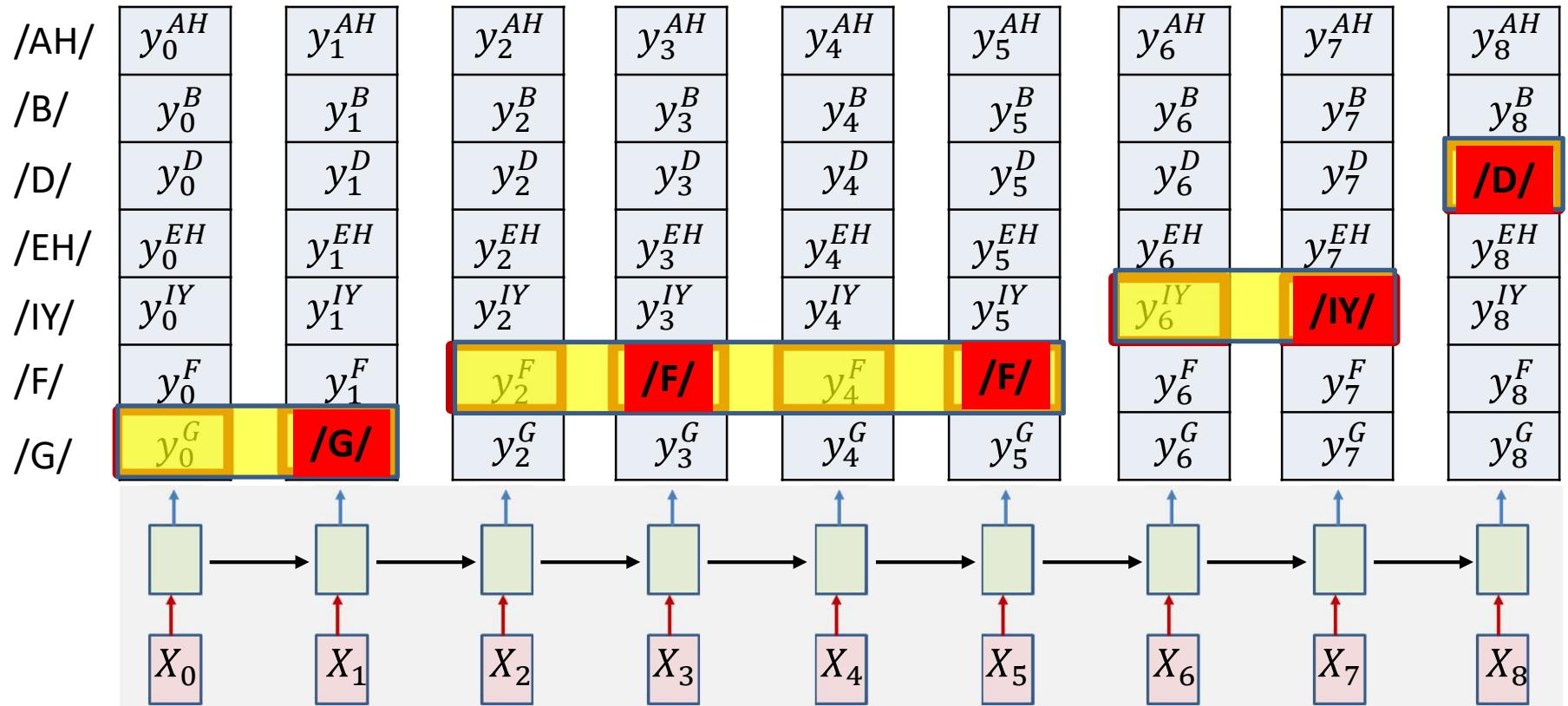
# Overall objective



- Find most likely symbol sequence given inputs

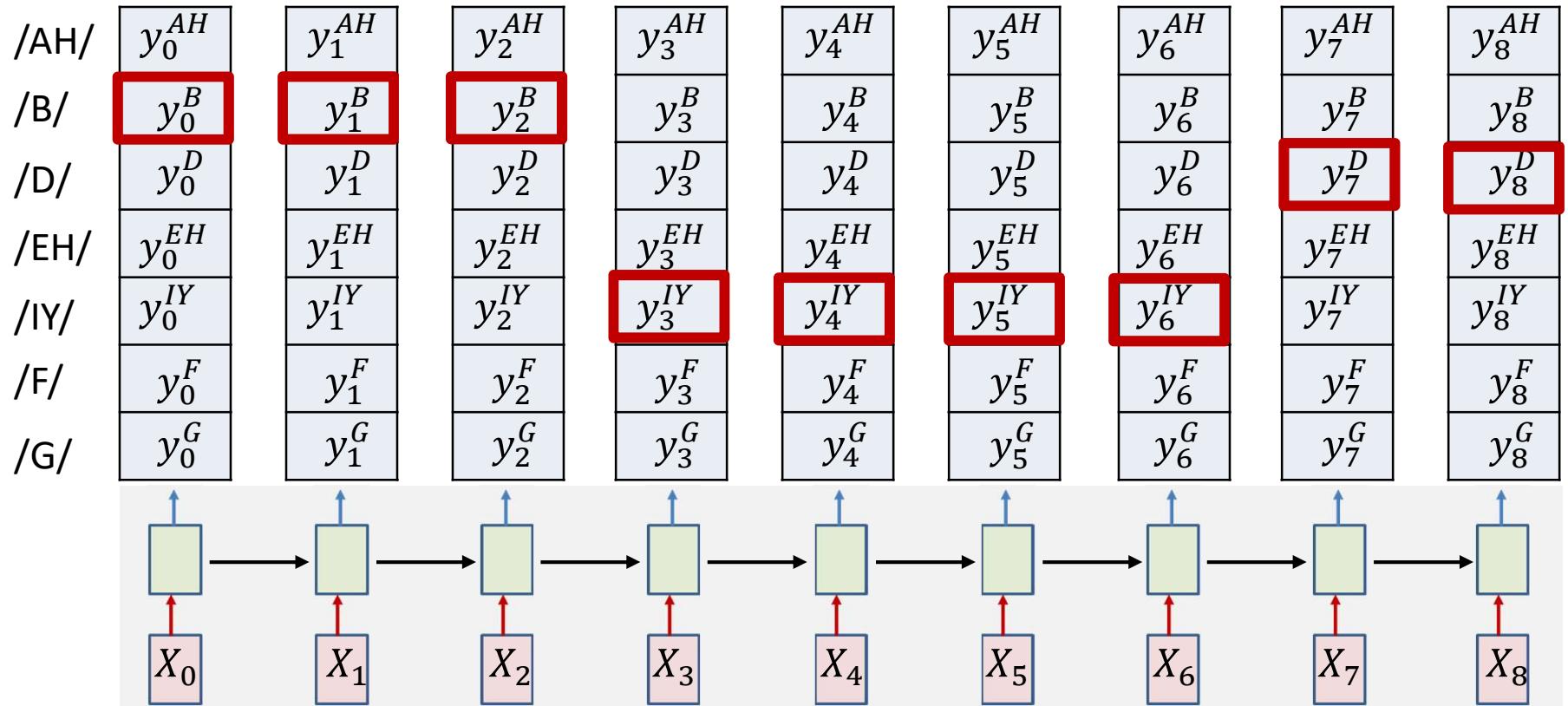
$$S_0 \dots S_{K-1} = \underset{S'_0 \dots S'_{K-1}}{\operatorname{argmax}} \text{prob}(S'_0 \dots S'_{K-1} | X_0 \dots X_{N-1})$$

# The actual output of the network



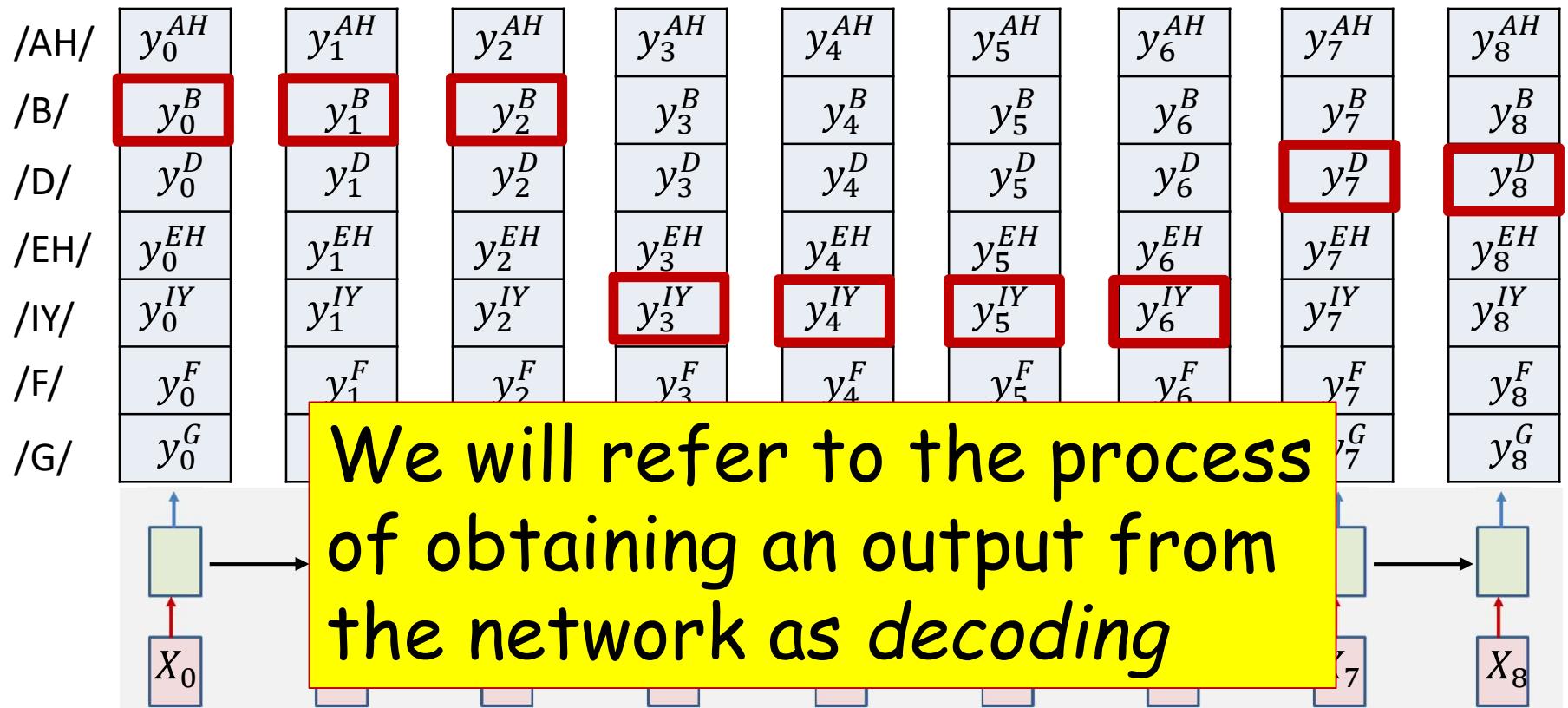
- Option 1: Simply select the most probable symbol at each time
  - Merge adjacent repeated symbols, and place the actual emission of the symbol in the final instant

# The actual output of the network



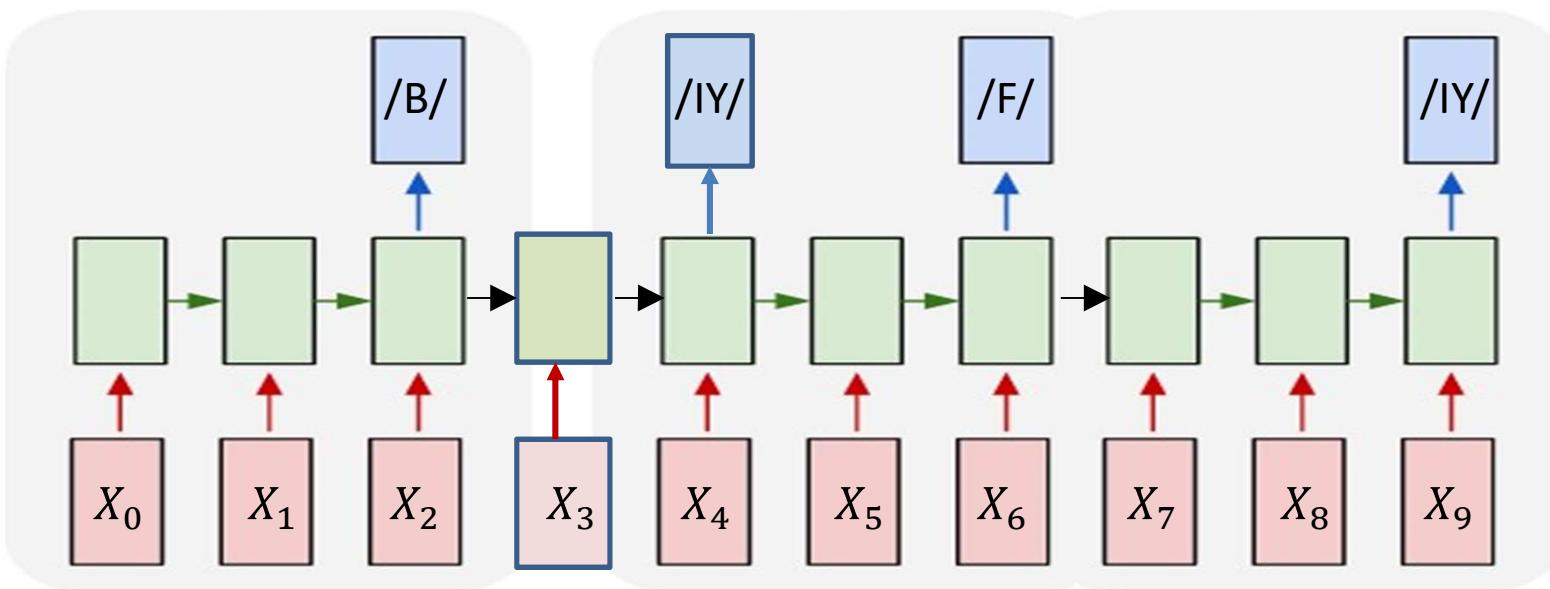
- Option 2: Impose external constraints on what sequences are allowed
  - E.g. only allow sequences corresponding to dictionary words

# The actual output of the network

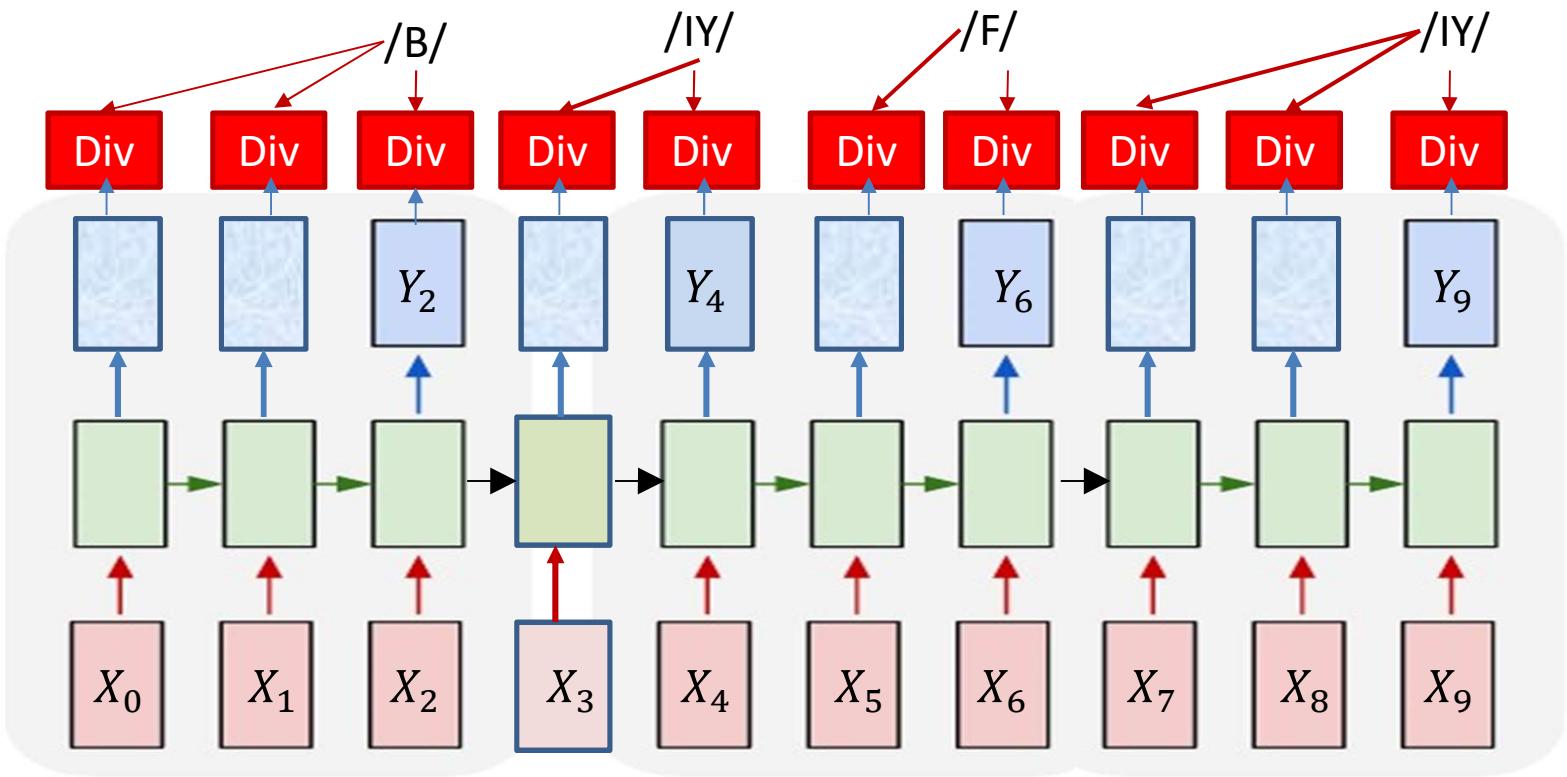


- Option 2: Impose external constraints on what sequences are allowed
  - E.g. only allow sequences corresponding to dictionary words

# Training



- Given output symbols *at the right locations*
  - The phoneme  $/B/$  ends at  $X_2$ ,  $/IY/$  at  $X_4$ ,  $/F/$  at  $X_6$ ,  $/IY/$  at  $X_9$



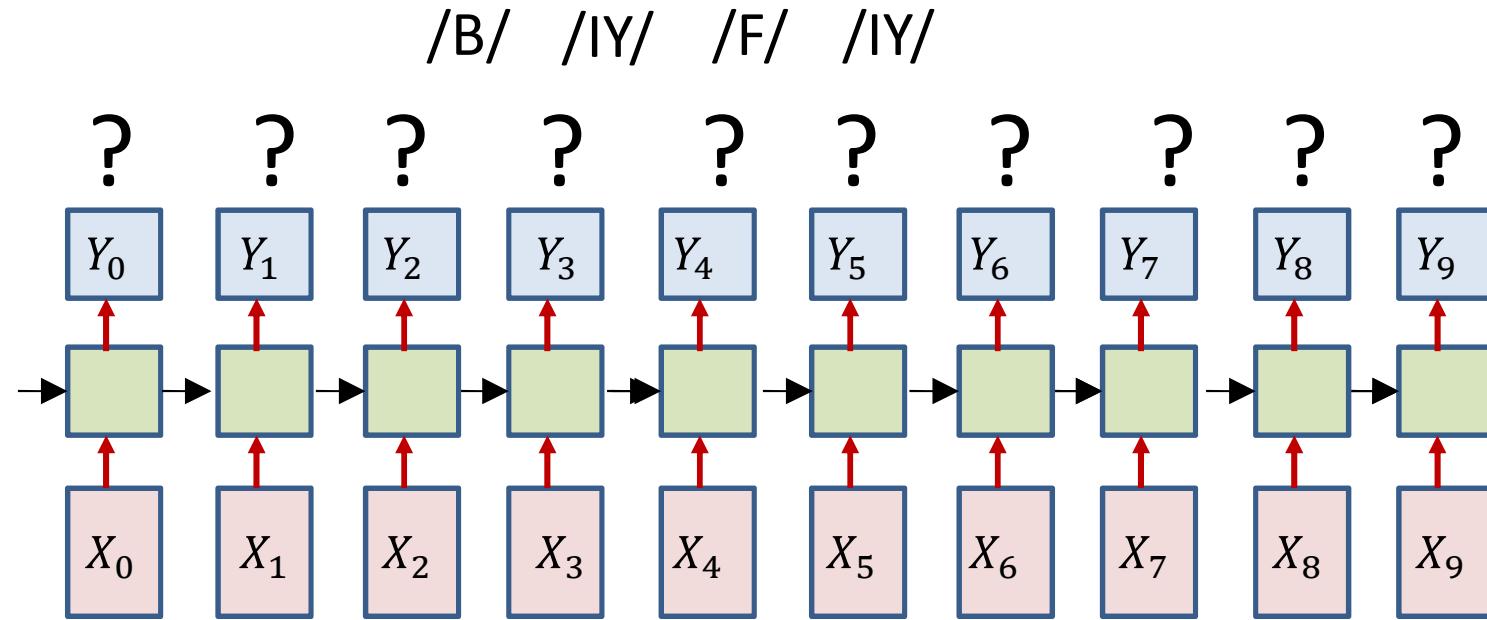
- Either just define Divergence as:

$$DIV = Xent(Y_2, B) + Xent(Y_4, IY) + Xent(Y_6, F) + Xent(Y_9, IY)$$

- Or *repeat the symbols over their duration*

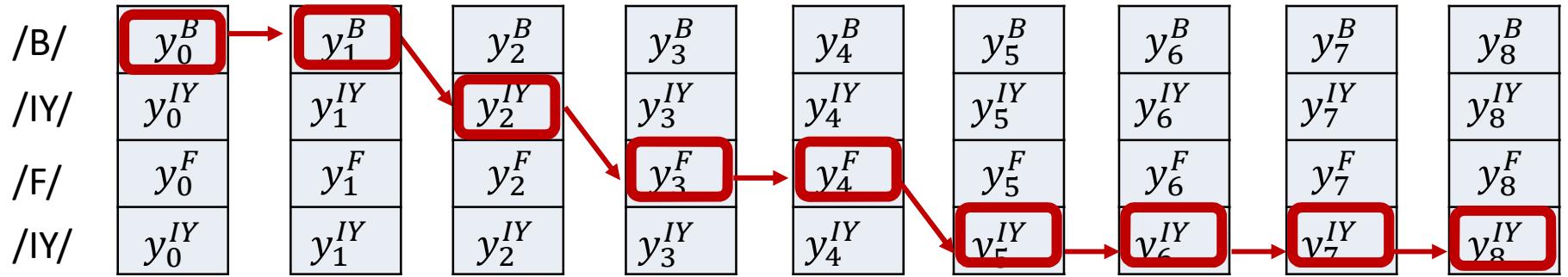
$$DIV = \sum_t Xent(Y_t, symbol_t) = - \sum_t \log Y(t, symbol_t)$$

# Problem: No timing information provided



- Only the sequence of output symbols is provided for the training data
  - But no indication of which one occurs where
- How do we compute the divergence?
  - And how do we compute its gradient w.r.t.  $Y_t$

# Viterbi algorithm

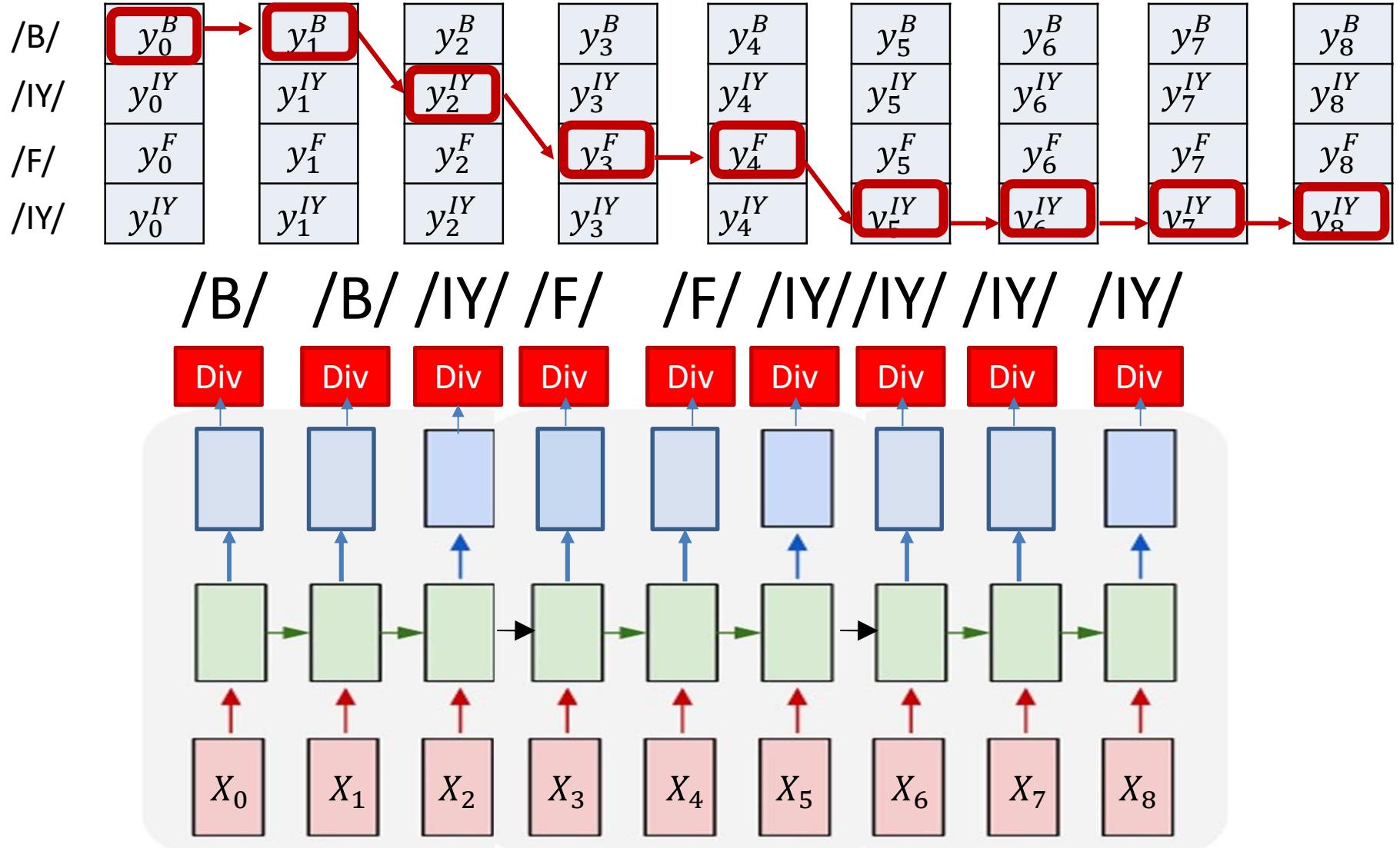


- $s(T - 1) = S(K - 1)$
- for  $t = T - 1$  down to 1

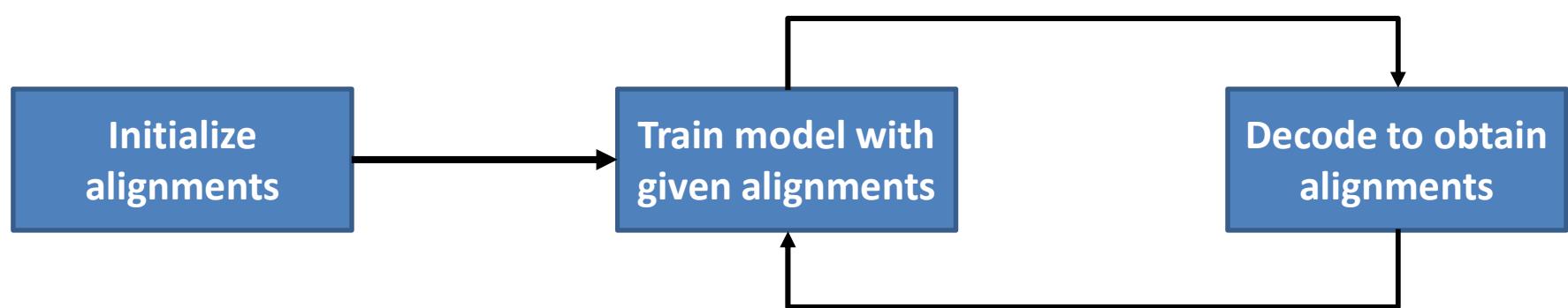
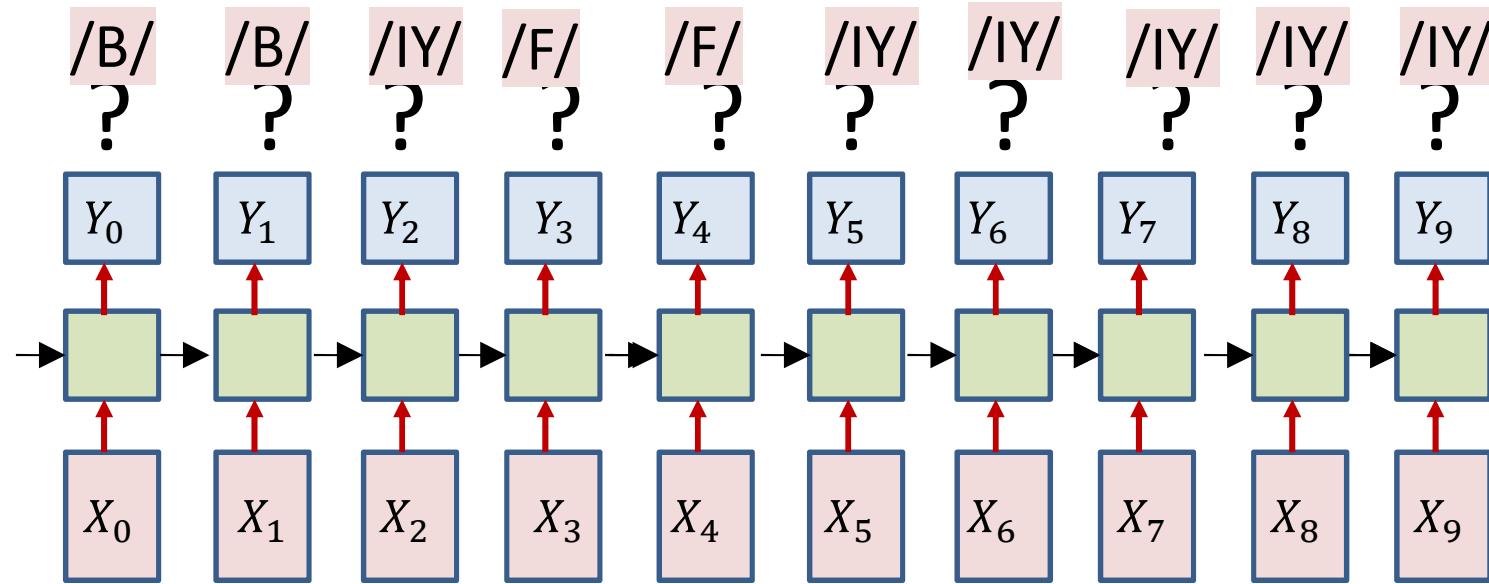
$$s(t - 1) = BP(s(t))$$

/B/ /B/ /IY/ /F/ /F/ /IY/ /IY/ /IY/ /IY/ /IY/

# Assumed targets for training with the Viterbi algorithm



# Iterative Estimate and Training

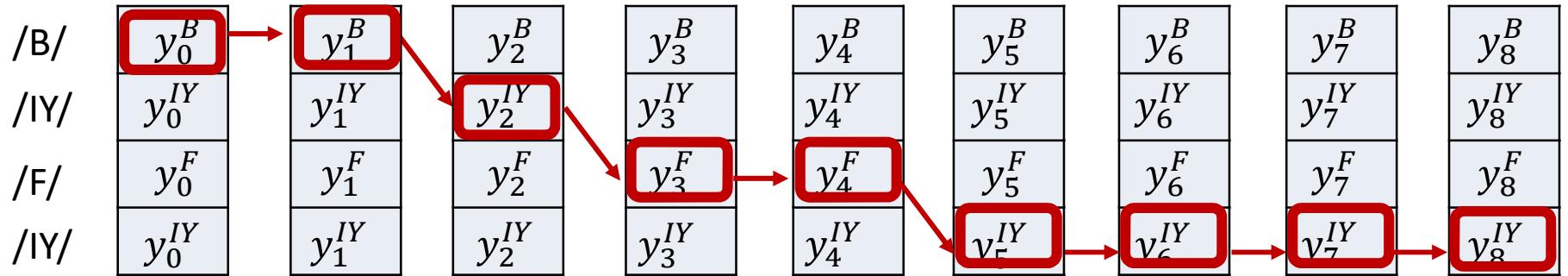


The “decode” and “train” steps may be combined into a single “decode, find alignment, compute derivatives” step for SGD and mini-batch updates

# Iterative update: Problem

- Approach heavily dependent on initial alignment
- Prone to poor local optima
- Alternate solution: Do not commit to an alignment during any pass..

# The reason for suboptimality

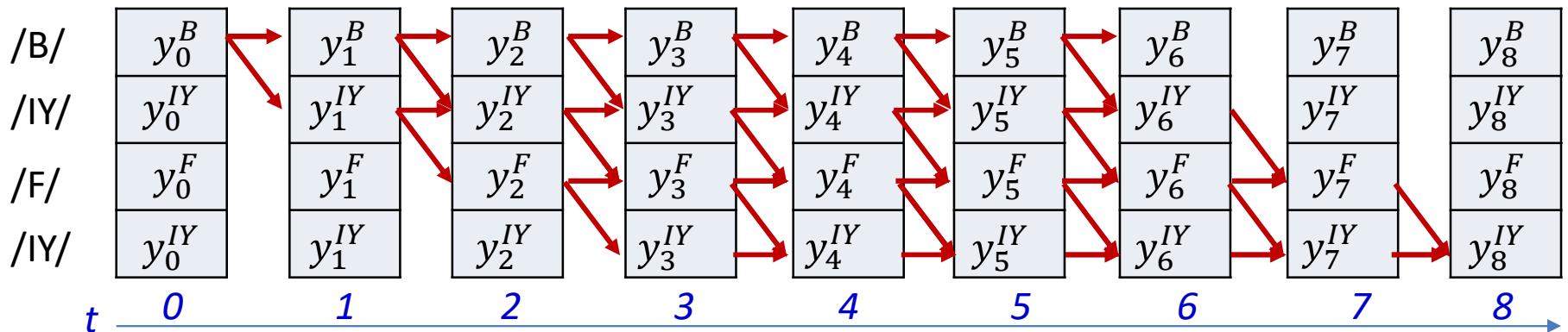


- We *commit* to the single “best” estimated alignment
  - The *most likely* alignment

$$DIV = - \sum_t \log Y(t, symbol_t^{bestpath})$$

- This can be way off, particularly in early iterations, or if the model is poorly initialized

# Averaging over *all* alignments

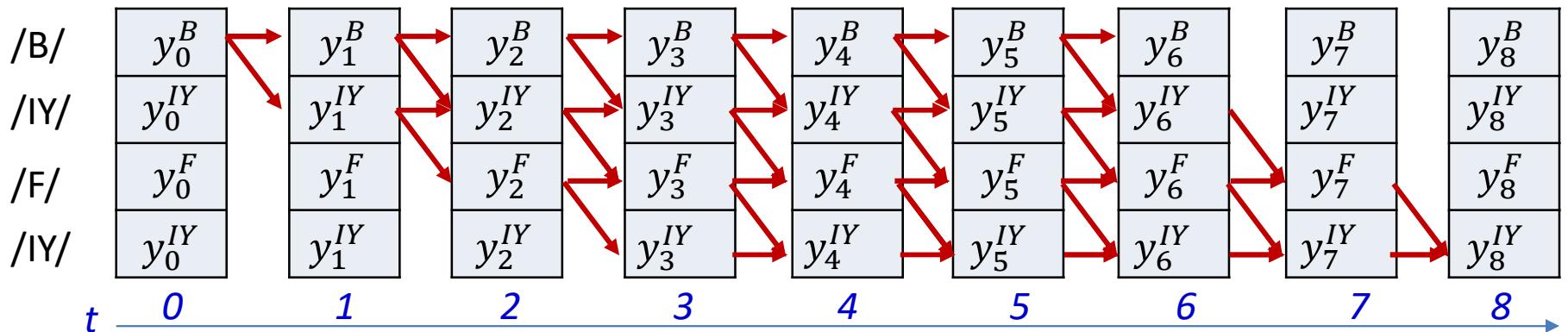


- Instead of only selecting the most likely alignment, use the statistical expectation over *all* possible alignments

$$DIV = E \left[ - \sum_t \log Y(t, s_t) \right]$$

- Use the *entire distribution of alignments*
- This will mitigate the issue of suboptimal selection of alignment

# The expectation over *all* alignments



$$DIV = E \left[ - \sum_t \log Y(t, s_t) \right]$$

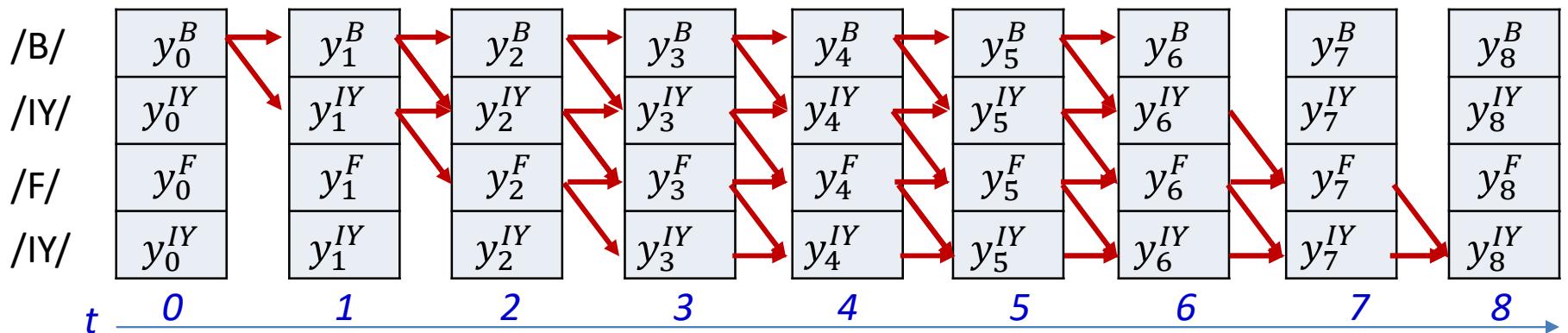
- Using the linearity of expectation

$$DIV = - \sum_t E[\log Y(t, s_t)]$$

- This reduces to finding the expected divergence *at each input*

$$DIV = - \sum_t \sum_{S \in S_1 \dots S_K} P(s_t = S | \mathbf{S}, \mathbf{X}) \log Y(t, s_t = S)$$

# The expectation over *all* alignments



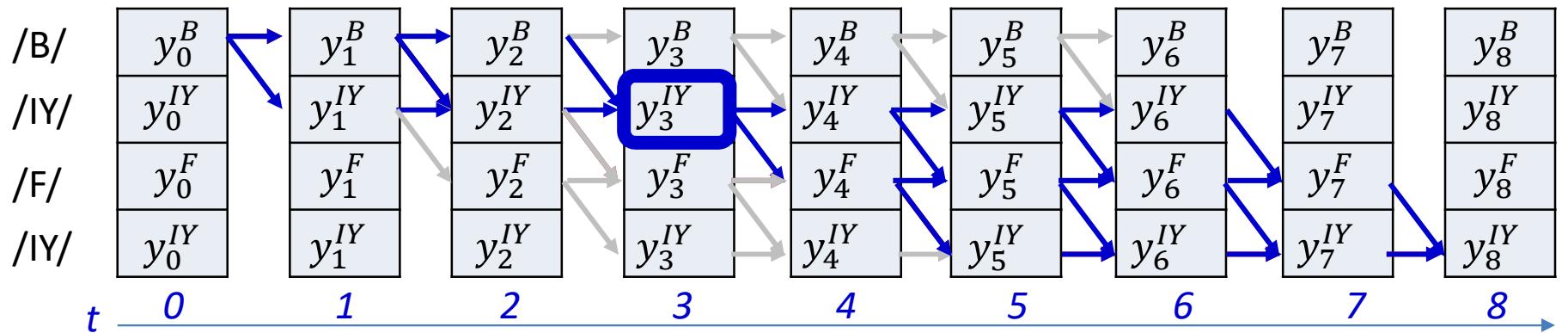
- The probability of seeing the specific symbol  $s$  at time  $t$ , given that the symbol sequence is an expansion of  $S = S_0 \dots S_{K-1}$  and given the input sequence  $X = X_0 \dots X_{N-1}$
- We need to be able to compute this

$$DIV = - \sum_t E[\log Y(t, s_t)]$$

– This reduces to finding the expected divergence at each input

$$DIV = - \sum_t \sum_{S \in S_1 \dots S_K} P(s_t = S | \mathbf{S}, \mathbf{X}) \log Y(t, s_t = S)$$

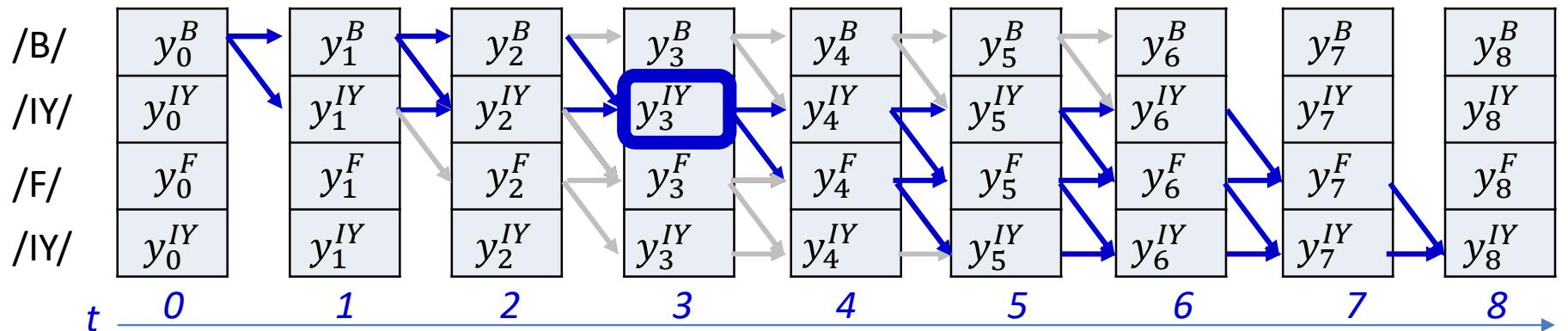
# A posteriori probabilities of symbols



$$P(s_t = S_r | \mathbf{S}, \mathbf{X}) \propto P(s_t = S_r, \mathbf{S} | \mathbf{X})$$

- $P(s_t = S_r, \mathbf{S} | \mathbf{X})$  is the total probability of all valid paths *in the graph for target sequence  $\mathbf{S}$*  that go through the symbol  $S_r$  (the  $r^{\text{th}}$  symbol in the sequence  $S_1 \dots S_K$ ) at time  $t$
- We will compute this using the “forward-backward” algorithm

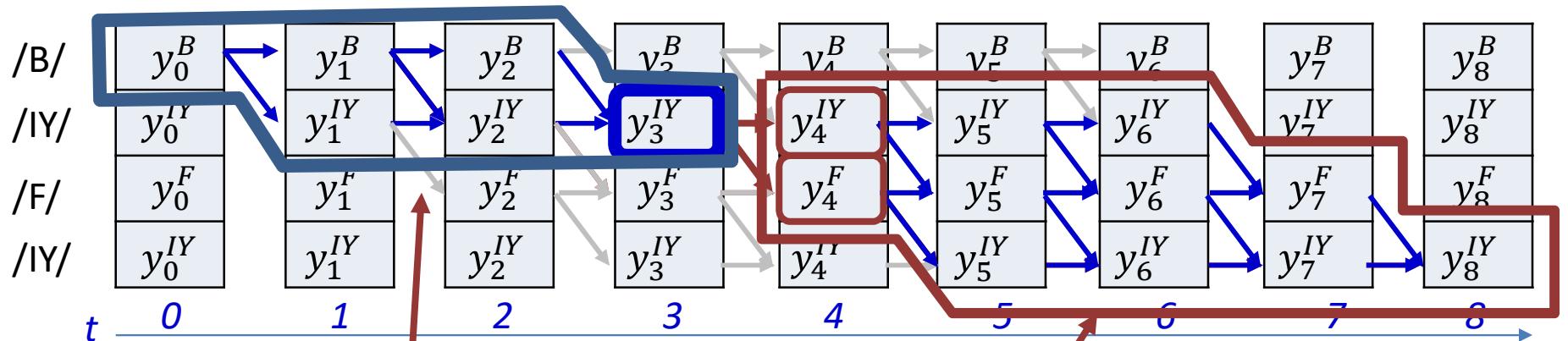
# A posteriori probabilities of symbols



- $P(s_t = S_r, \mathbf{S} | \mathbf{X})$  can be decomposed as

$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = P(S_1, \dots, S_K, s_t = S_r | \mathbf{X})$$

# A posteriori probabilities of symbols

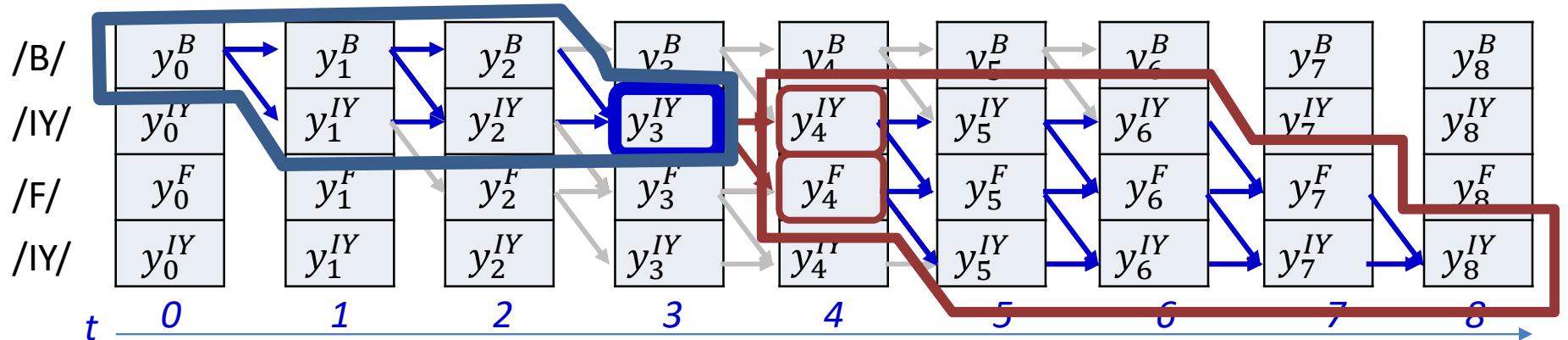


- $P(s_t = S_r, \mathbf{S} | \mathbf{X})$  can be decomposed as

$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = P(S_1, \dots, S_r, \dots, S_K, s_t = S_r | \mathbf{X}) \\ = P(S_1 \dots S_r, s_t = S_r, S_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K, | \mathbf{X})$$

- Where  $\text{succ}(S_r)$  is a symbol that can follow  $S_r$  in a sequence
  - Here it is either  $S_r$  or  $S_{r+1}$  (red blocks in figure)
  - The equation literally says that after the blue block, either of the two red arrows may be followed

# A posteriori probabilities of symbols

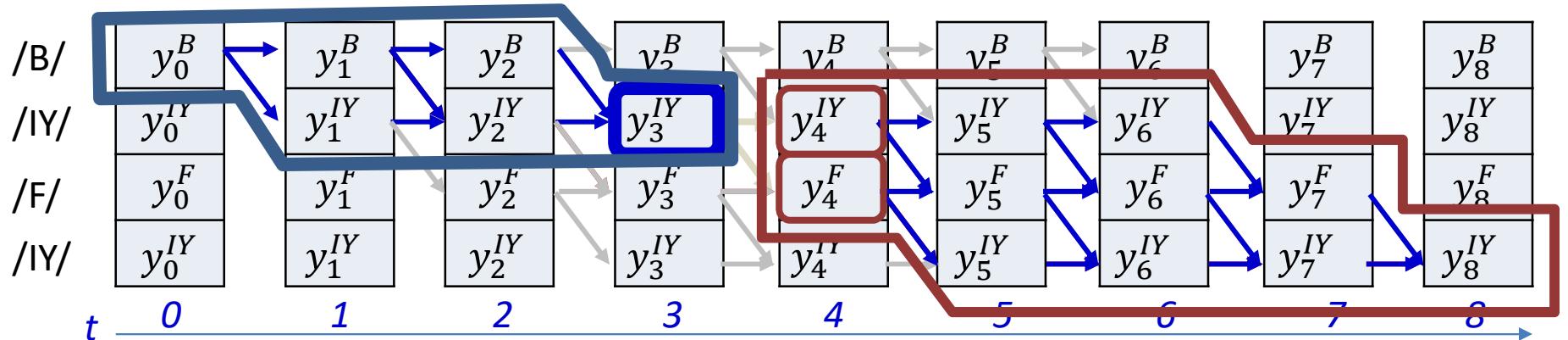


- $P(s_t = S_r, \mathbf{S} | \mathbf{X})$  can be decomposed as

$$\begin{aligned} P(s_t = S_r, \mathbf{S} | \mathbf{X}) &= P(S_1, \dots, S_r, \dots, S_K, s_t = S_r | \mathbf{X}) \\ &= P(S_1 \dots S_r, s_t = S_r, s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | S_1 \dots S_r, s_t = S_r | \mathbf{X}) \end{aligned}$$

- Using Bayes Rule  
 $= P(S_1 \dots S_r, s_t = S_r | \mathbf{X})P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | S_1 \dots S_r, s_t = S_r | \mathbf{X})$
- The probability of the subgraph in the blue outline, times the conditional probability of the red-encircled subgraph, given the blue subgraph

# A posteriori probabilities of symbols



- $P(s_t = S_r, \mathbf{S} | \mathbf{X})$  can be decomposed as

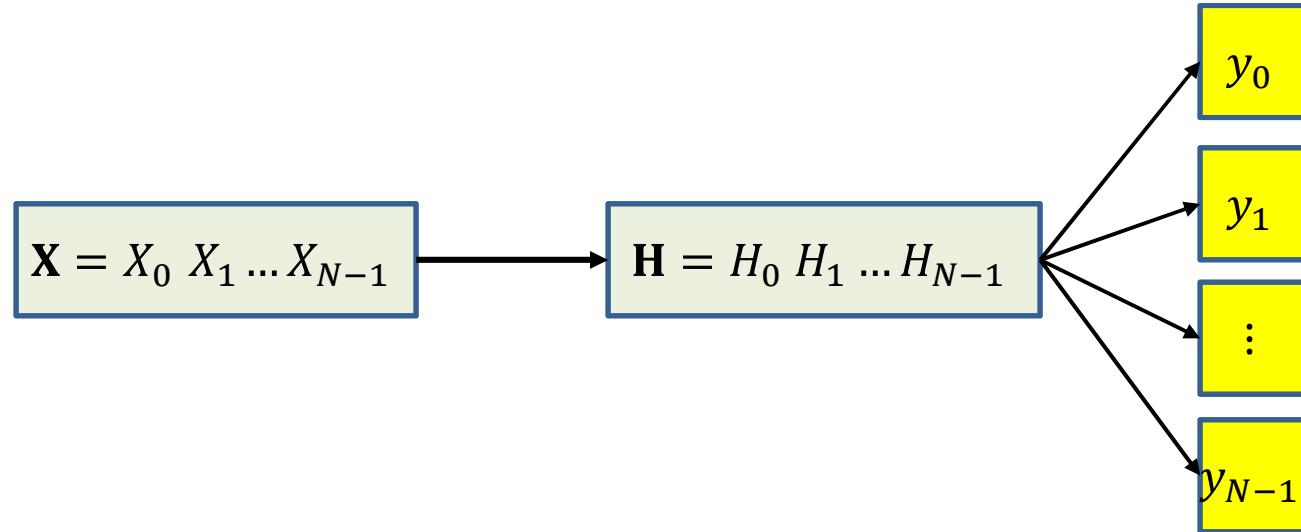
$$\begin{aligned} P(s_t = S_r, \mathbf{S} | \mathbf{X}) &= P(S_1, \dots, S_r, \dots, S_K, s_t = S_r | \mathbf{X}) \\ &= P(S_1 \dots S_r, s_t = S_r, s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X}) \end{aligned}$$

- Using Bayes Rule
- $= P(S_1 \dots S_r, s_t = S_r | \mathbf{X})P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | S_1 \dots S_r, s_t = S_r | \mathbf{X})$
- For a recurrent network without feedback from the output we can make the conditional independence assumption:

$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = P(S_1 \dots S_r, s_t = S_r | \mathbf{X})P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})$$

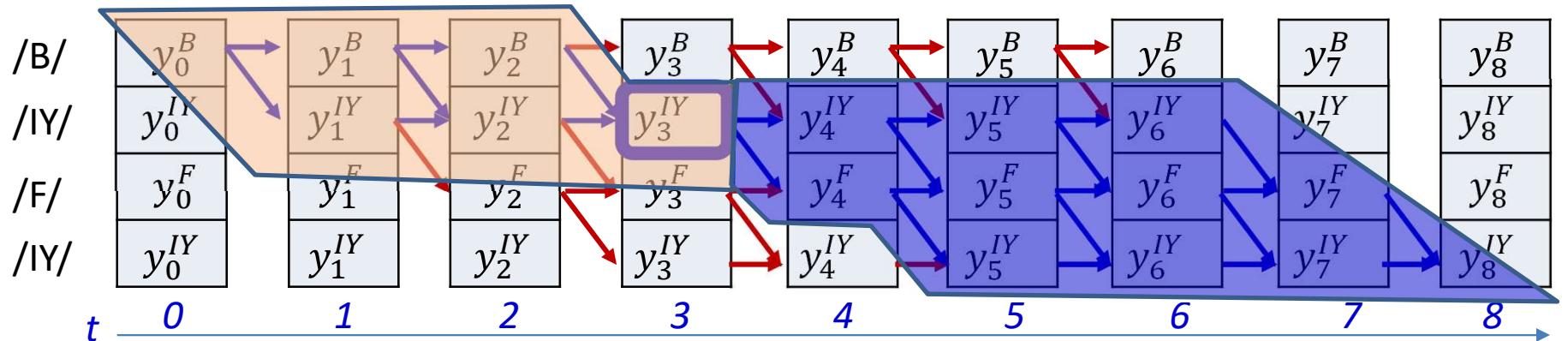
Assuming past output symbols do not directly feed back into the net

# Conditional independence



- **Dependency graph:** Input sequence  $\mathbf{X} = X_0 \ X_1 \dots X_{N-1}$  governs hidden variables  $\mathbf{H} = H_0 \ H_1 \dots H_{N-1}$
- Hidden variables govern output predictions  $y_0, y_1, \dots y_{N-1}$  individually
- $y_0, y_1, \dots y_{N-1}$  are conditionally independent given  $\mathbf{H}$
- Since  $\mathbf{H}$  is deterministically derived from  $\mathbf{X}$ ,  $y_0, y_1, \dots y_{N-1}$  are also conditionally independent given  $\mathbf{X}$ 
  - This wouldn't be true if the relation between  $\mathbf{X}$  and  $\mathbf{H}$  were not deterministic or if  $\mathbf{X}$  is unknown, or if there were direct connections between the  $y$ s

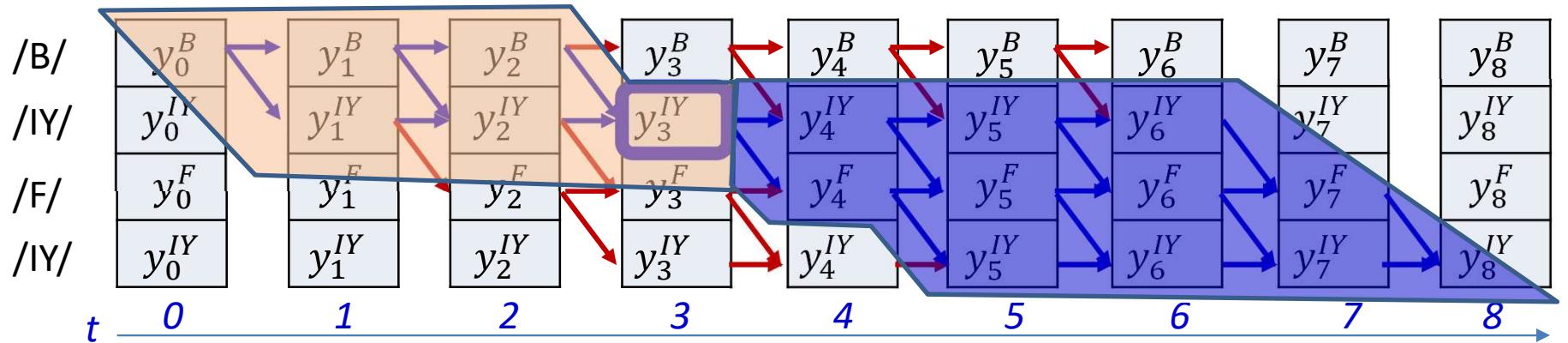
# A posteriori symbol probability



$$\begin{aligned}
 & P(s_t = S_r, \mathbf{S} | \mathbf{X}) \\
 &= \underline{P(S_1 \dots S_r, s_t = S_r | \mathbf{X})} \underline{P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})}
 \end{aligned}$$

- We will call the first term the *forward probability*  $\alpha(t, r)$
- We will call the second term the *backward probability*  $\beta(t, r)$

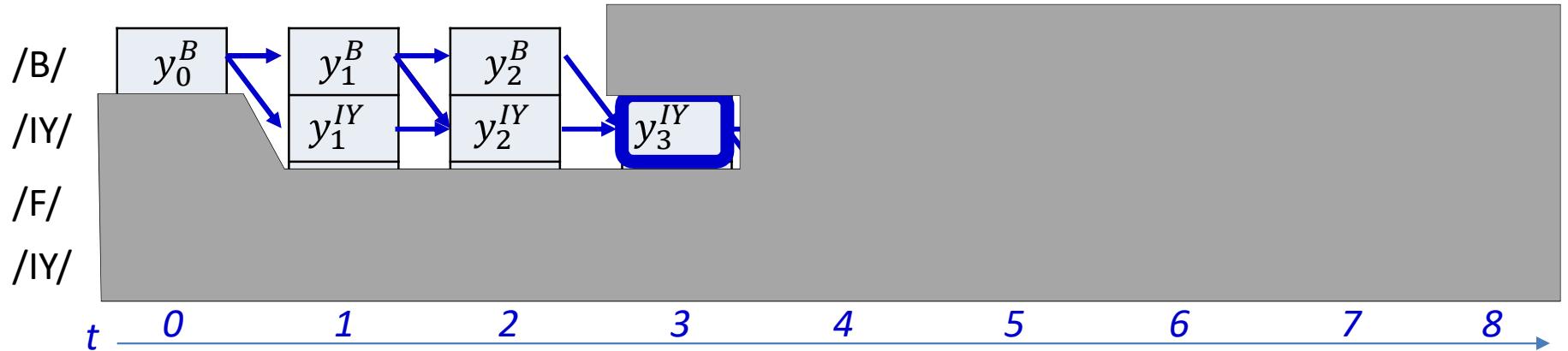
# A posteriori symbol probability



$$\begin{aligned}
 & P(s_t = S_r, \mathbf{S} | \mathbf{X}) \\
 &= \boxed{P(S_1 \dots S_r, s_t = S_r | \mathbf{X})} \underbrace{P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})}
 \end{aligned}$$

- We will call the first term the *forward probability*  $\alpha(t, r)$
- We will call the second term the *backward probability*  $\beta(t, r)$

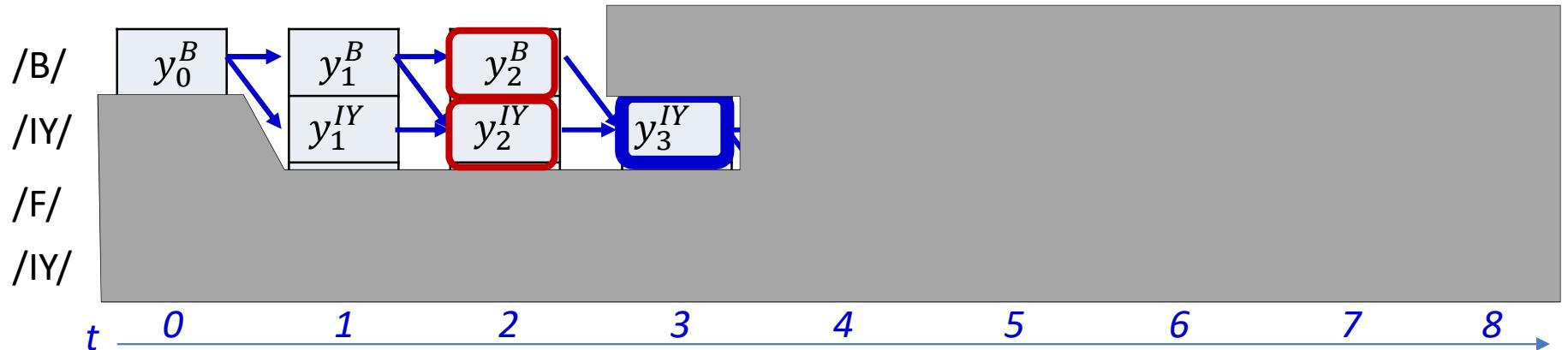
# Computing $\alpha(t, r)$ : Forward algorithm



$$\alpha(t, r) = P(S_1..S_r, s_t = S_r | \mathbf{X})$$

- The  $\alpha(t, r)$  is the total probability of the subgraph shown

# Computing $\alpha(t, r)$ : Forward algorithm



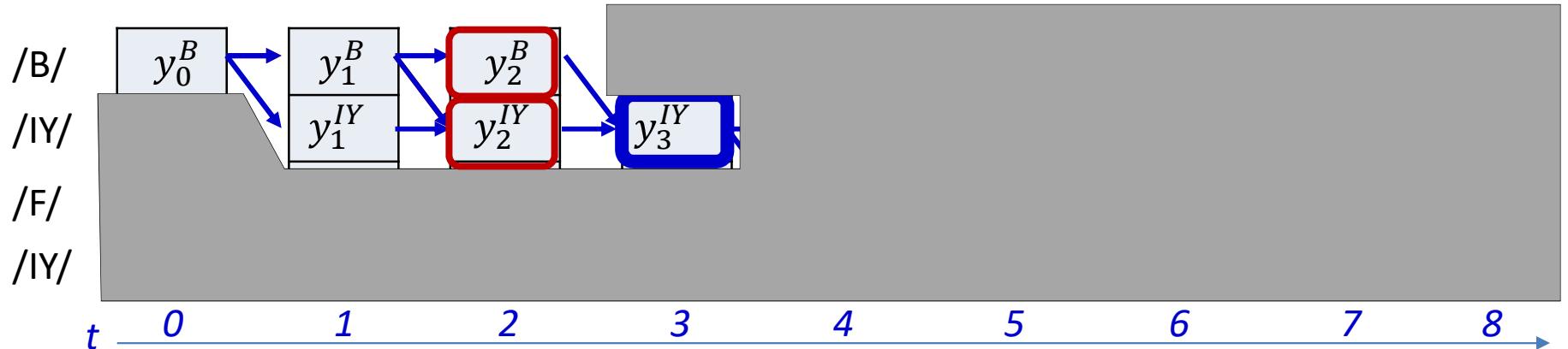
$$\alpha(t, r) = P(S_1 \dots S_r, s_t = S_r | \mathbf{X})$$

- The  $\alpha(t, r)$  is the total probability of the subgraph shown
- We can marginalize the symbol at time  $t-1$

$$\alpha(t, r) = \sum_{q: S_q \in \text{pred}(S_r)} P(S_1 \dots S_q, s_{t-1} = S_q, s_t = S_r | \mathbf{X})$$

- Where  $\text{pred}(S_r)$  is any symbol that is permitted to come before an  $S_r$  and may include  $S_r$
- $q$  is its row index, and can take values  $r$  and  $r - 1$  in this example

# Computing $\alpha(t, r)$ : Forward algorithm



$$\alpha(t, r) = P(S_1 \dots S_r, s_t = S_r | \mathbf{X})$$

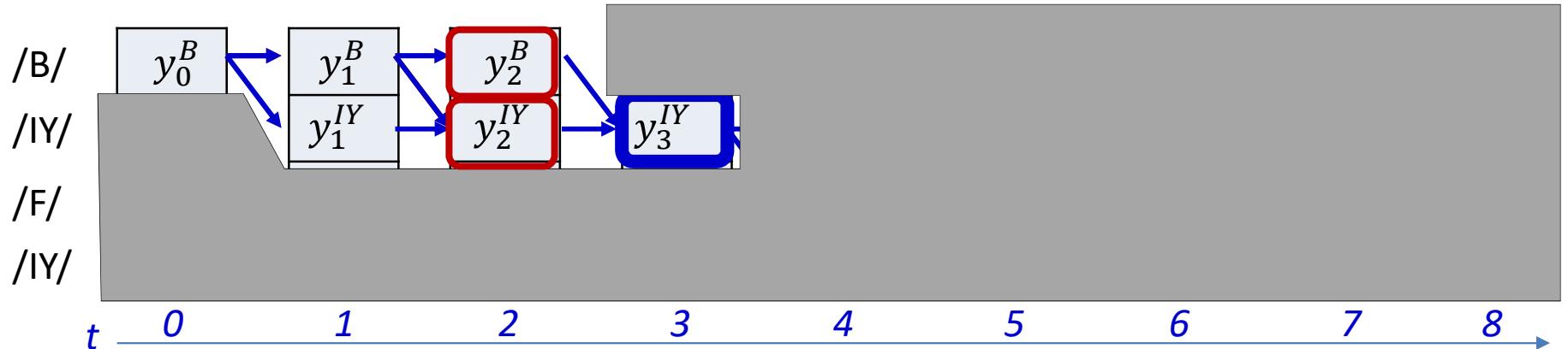
- The  $\alpha(t, r)$  is the total probability of the subgraph shown
- We can marginalize out the symbol at time  $t-1$

$$\alpha(t, r) = \sum_{q: S_q \in \text{pred}(S_r)} P(S_1 \dots, \underline{S_q}, \underline{s_{t-1} = S_q}, \underline{s_t = S_r} | \mathbf{X})$$

- Using the conditional independence assumed

$$\alpha(t, r) = \sum_{q: S_q \in \text{pred}(S_r)} P(S_1 \dots, S_q, s_{t-1} = S_q | \mathbf{X}) P(s_t = S_r | \mathbf{X})$$

# Computing $\alpha(t, r)$ : Forward algorithm

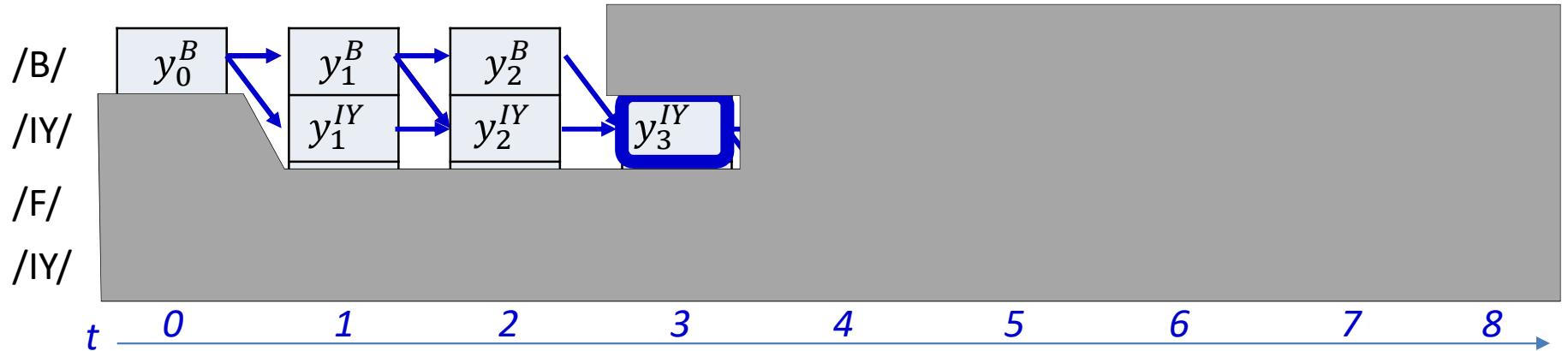


$$\alpha(t, r) = P(S_1 \dots S_r, s_t = S_r | \mathbf{X})$$

- The  $\alpha(t, r)$  is the total probability of the subgraph shown

$$\begin{aligned} \alpha(t, r) &= \sum_{q: S_q \in pred(S_r)} P(S_1 \dots, S_q, s_{t-1} = S_q | \mathbf{X}) P(s_t = S_r | \mathbf{X}) \\ &= \sum_{q: S_q \in pred(S_r)} \alpha(t-1, q) y_t^{S_r} \end{aligned}$$

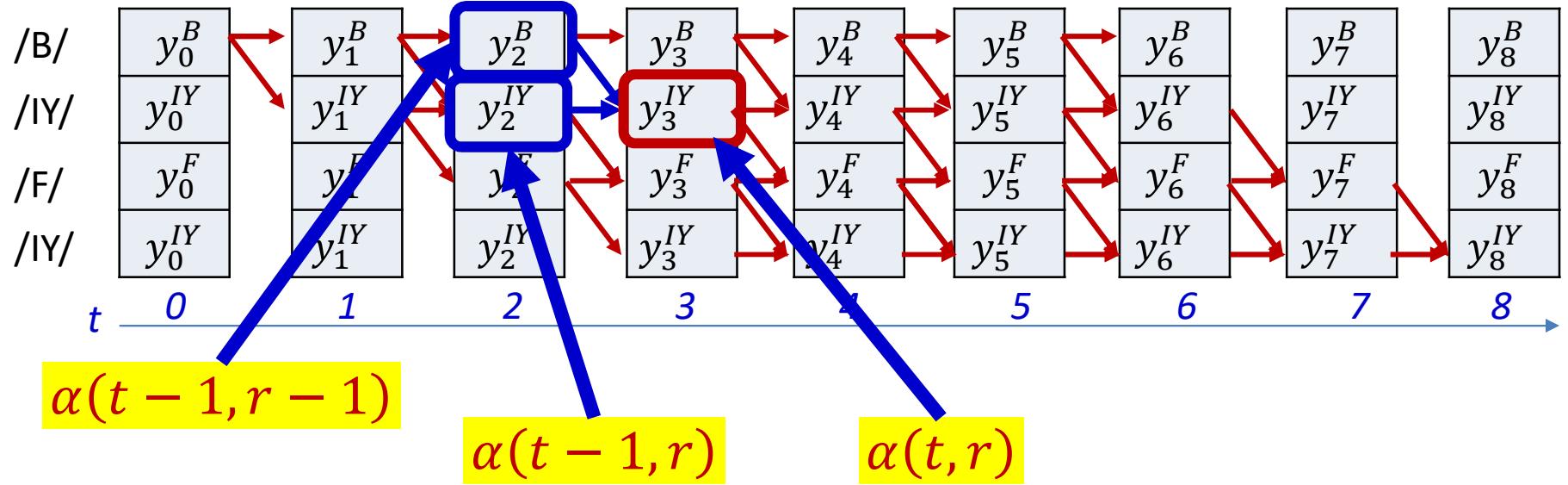
# Forward algorithm



$$\alpha(t, r) = \sum_{q : S_q \in pred(S_r)} \alpha(t - 1, q) y_t^{S_r}$$

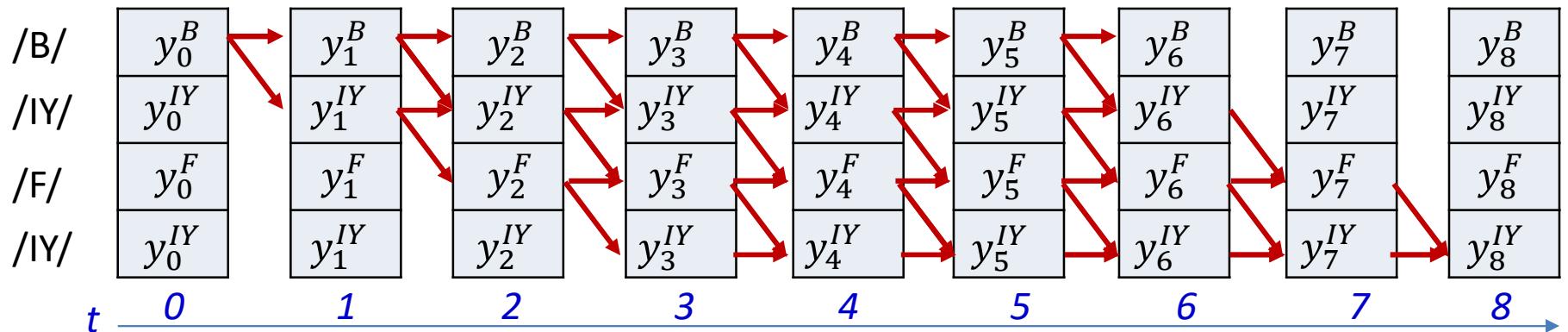
- The  $\alpha(t, r)$  is the total probability of the subgraph shown

# Forward algorithm



$$\alpha(t, r) = (\alpha(t-1, r) + \alpha(t-1, r-1)) y_t^{S(r)}$$

# Forward algorithm



- Initialization:

$$\alpha(0,1) = y_0^{S(1)}, \quad \alpha(0,r) = 0, \quad r > 1$$

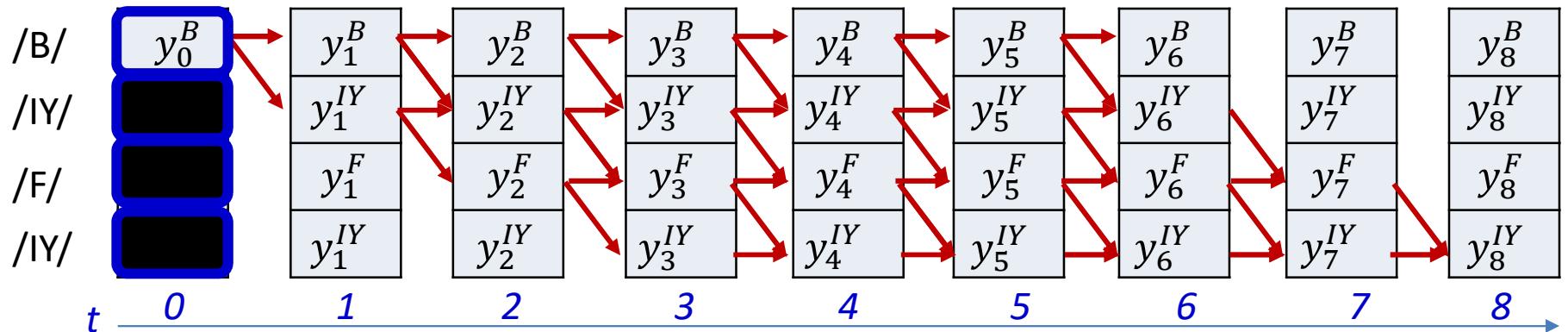
- for  $t = 1 \dots T - 1$

$$\alpha(t,1) = \alpha(t-1,1)y_t^{S(1)}$$

for  $l = 2 \dots K$

- $\alpha(t,l) = (\alpha(t-1,l) + \alpha(t-1,l-1))y_t^{S(l)}$

# Forward algorithm



- Initialization:

$$\alpha(0,1) = y_0^{S(1)}, \quad \alpha(0,r) = 0, \quad r > 1 \quad \leftarrow$$

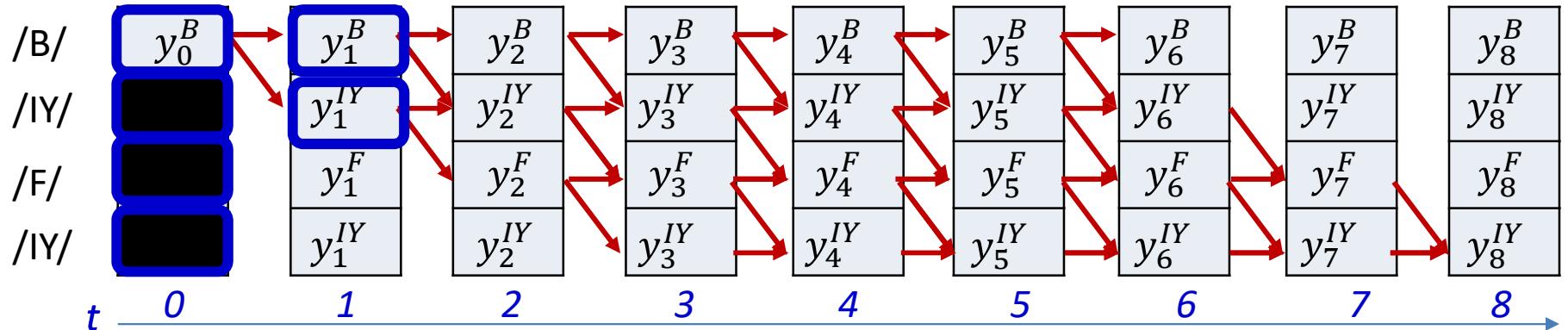
- for  $t = 1 \dots T - 1$

$$\alpha(t,1) = \alpha(t-1,1)y_t^{S(1)}$$

for  $l = 2 \dots K$

- $\alpha(t,l) = (\alpha(t-1,l) + \alpha(t-1,l-1))y_t^{S(l)}$

# Forward algorithm



- Initialization:

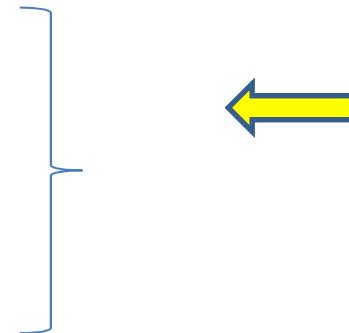
$$\alpha(0,1) = y_0^{S(1)}, \quad \alpha(0,r) = 0, \quad r > 1$$

- for  $t = 1 \dots T - 1$

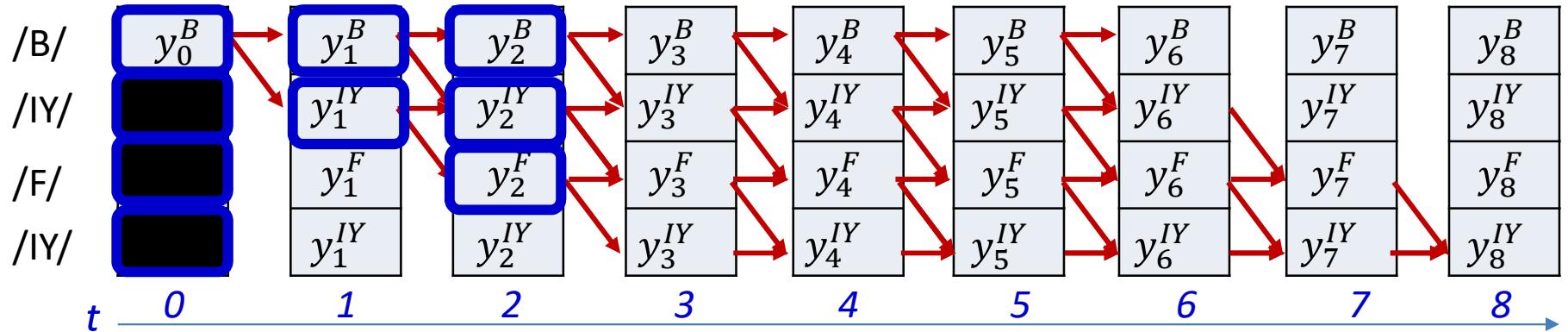
$$\alpha(t,1) = \alpha(t-1,1)y_t^{S(1)}$$

for  $l = 2 \dots K$

- $\alpha(t,l) = (\alpha(t-1,l) + \alpha(t-1,l-1))y_t^{S(l)}$



# Forward algorithm



- Initialization:

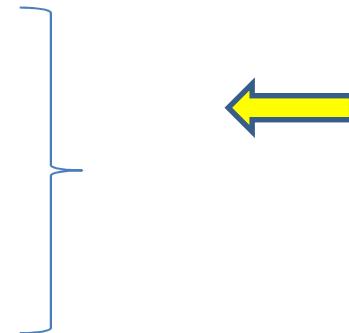
$$\alpha(0,1) = y_0^{S(1)}, \quad \alpha(0,r) = 0, \quad r > 1$$

- for  $t = 1 \dots T - 1$

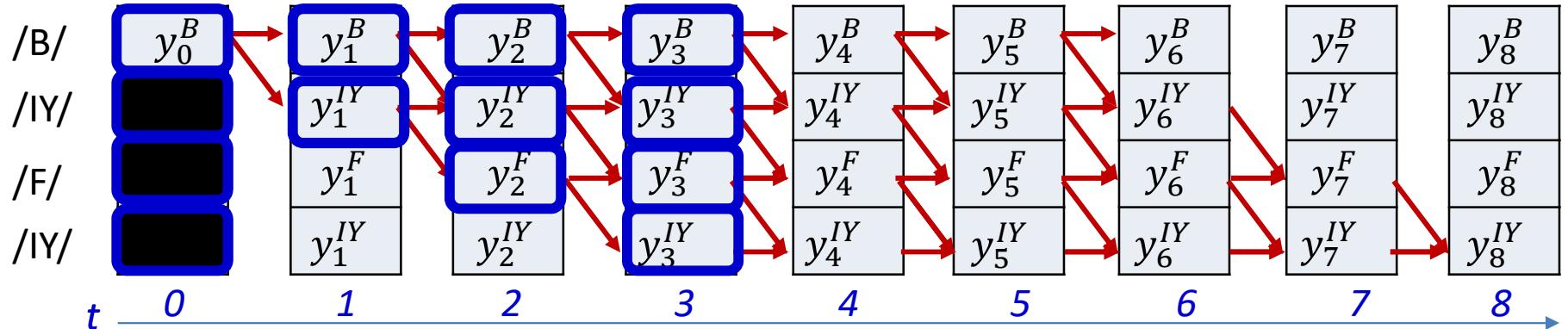
$$\alpha(t,1) = \alpha(t-1,1)y_t^{S(1)}$$

for  $l = 2 \dots K$

- $\alpha(t,l) = (\alpha(t-1,l) + \alpha(t-1,l-1))y_t^{S(l)}$



# Forward algorithm



- Initialization:

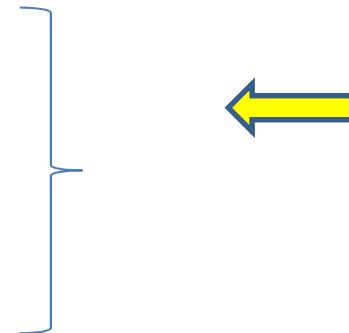
$$\alpha(0,1) = y_0^{S(1)}, \quad \alpha(0,r) = 0, \quad r > 1$$

- for  $t = 1 \dots T - 1$

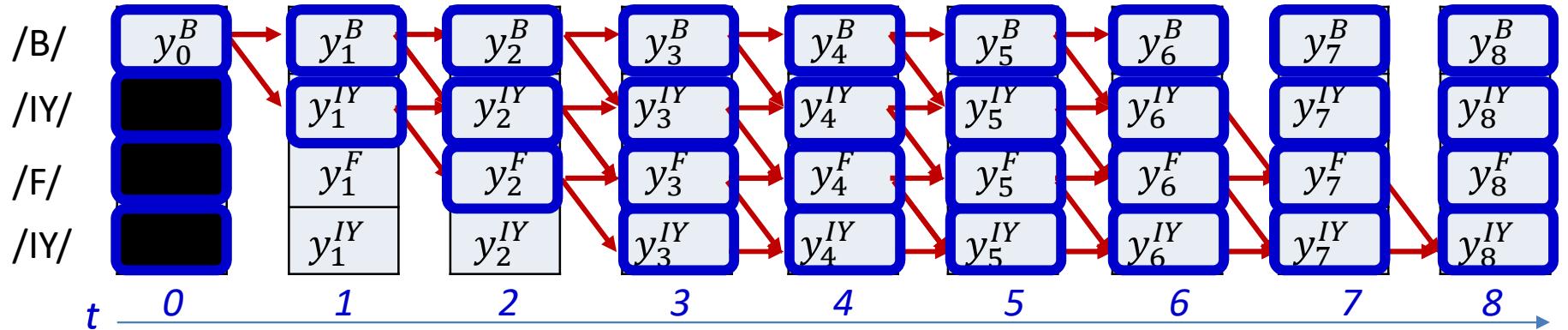
$$\alpha(t,1) = \alpha(t-1,1)y_t^{S(1)}$$

for  $l = 2 \dots K$

$$\cdot \quad \alpha(t,l) = (\alpha(t-1,l) + \alpha(t-1,l-1))y_t^{S(l)}$$



# Forward algorithm



- Initialization:

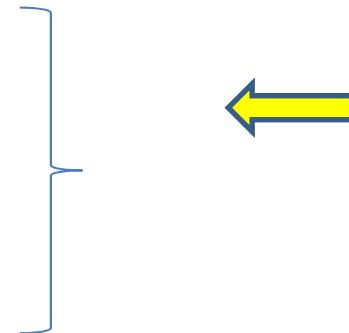
$$\alpha(0,1) = y_0^{S(1)}, \quad \alpha(0,r) = 0, \quad r > 1$$

- for  $t = 1 \dots T - 1$

$$\alpha(t,1) = \alpha(t-1,1)y_t^{S(1)}$$

for  $l = 2 \dots K$

$$\cdot \quad \alpha(t,l) = (\alpha(t-1,l) + \alpha(t-1,l-1))y_t^{S(l)}$$



# In practice..

- The recursion

$$\alpha(t, l) = (\alpha(t - 1, l) + \alpha(t - 1, l - 1))y_t^{S(l)}$$

will generally underflow

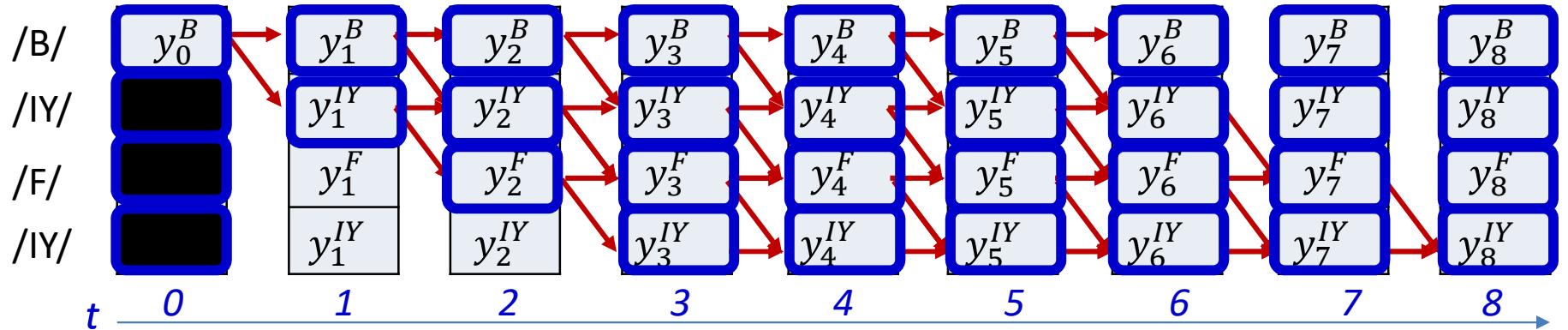
- Instead we can do it in the *log* domain

$$\log \alpha(t, l)$$

$$= \log(e^{\log \alpha(t-1, l)} + e^{\log \alpha(t-1, l-1)}) + \log y_t^{S(l)}$$

- This can be computed entirely without underflow

# Forward algorithm: Alternate statement



- The algorithm can also be stated as follows which separates the graph probability from the observation probability. This is needed to compute derivatives
- Initialization:

$$\hat{\alpha}(0,1) = 1, \quad \hat{\alpha}(0,r) = 0, \quad r > 1$$

$$\alpha(0,r) = \hat{\alpha}(0,r)y_0^{S(r)}, \quad 1 \leq r \leq K$$

- for  $t = 1 \dots T - 1$

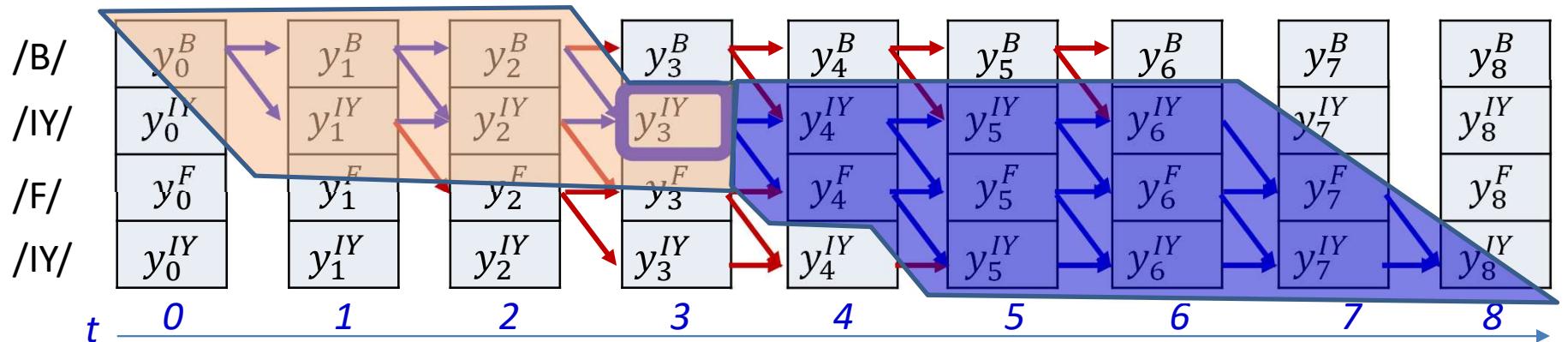
$$\hat{\alpha}(t,1) = \alpha(t-1,1)$$

for  $l = 2 \dots K$

- $\hat{\alpha}(t,l) = \alpha(t-1,l) + \alpha(t-1,l-1)$

$$\alpha(t,r) = \hat{\alpha}(t,r)y_t^{S(r)}, \quad 1 \leq r \leq K$$

# A posteriori symbol probability



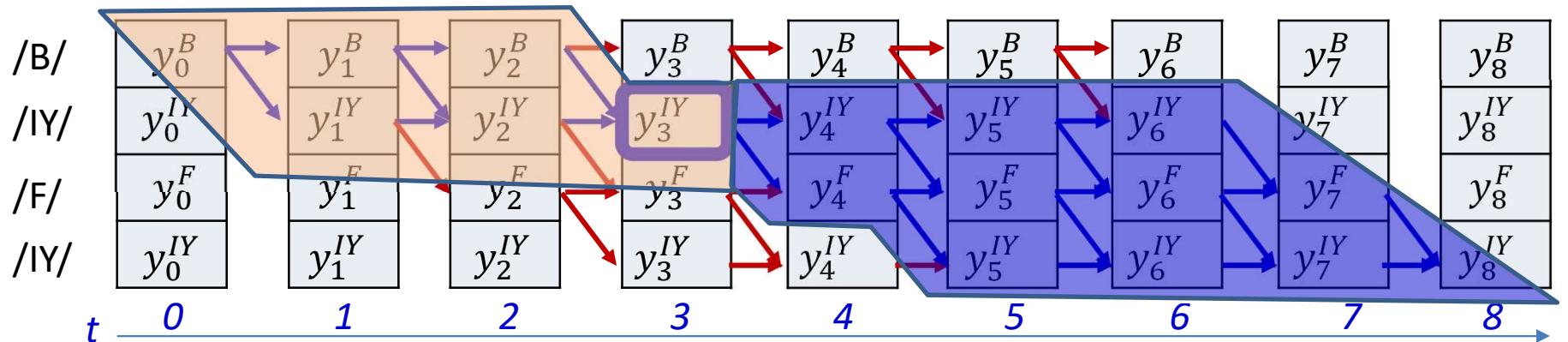
$$\begin{aligned}
 & P(s_t = S_r, \mathbf{S} | \mathbf{X}) \\
 &= \boxed{P(S_1 \dots S_r, s_t = S_r | \mathbf{X})} \boxed{P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})}
 \end{aligned}$$

- We will call the first term the *forward probability*  $\alpha(t, r)$
- We will call the second term the *backward probability*  $\beta(t, r)$



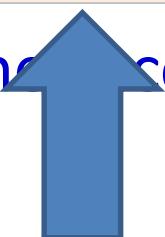
We have seen how to compute this

# A posteriori symbol probability



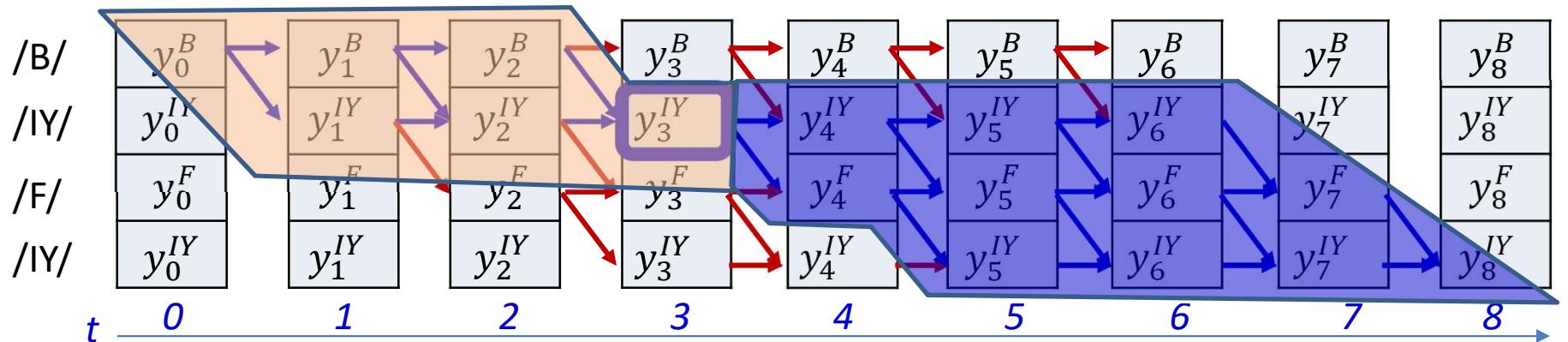
$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = \alpha(t, r) P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})$$

- We will call the first term the *forward probability*  $\alpha(t, r)$
- We will call the second term the *backward probability*  $\beta(t, r)$



We have seen how to compute this

# A posteriori symbol probability



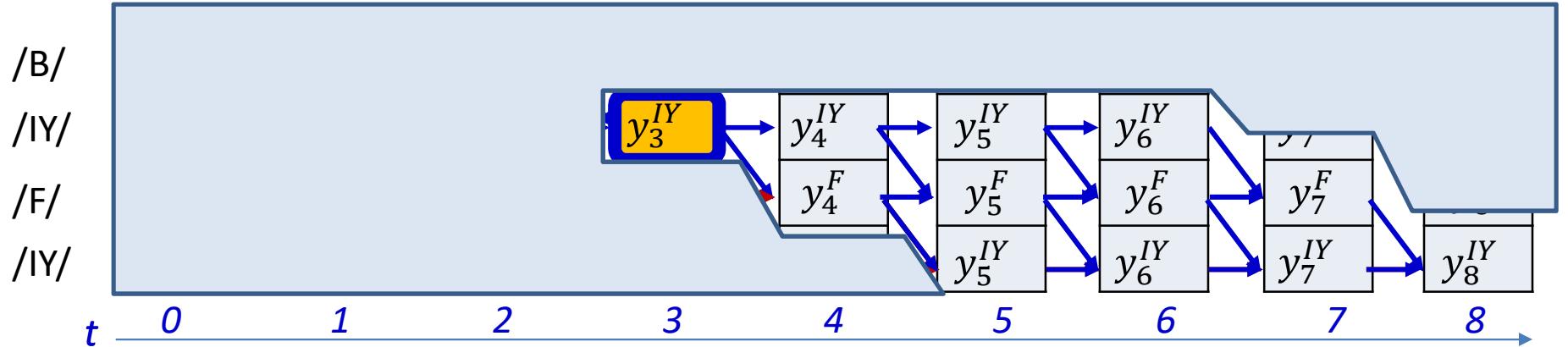
$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = \alpha(t, r) P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})$$

- We will call the first term the *forward probability*  $\alpha(t, r)$
- We will call the second term the *backward probability*  $\beta(t, r)$



Lets look at this

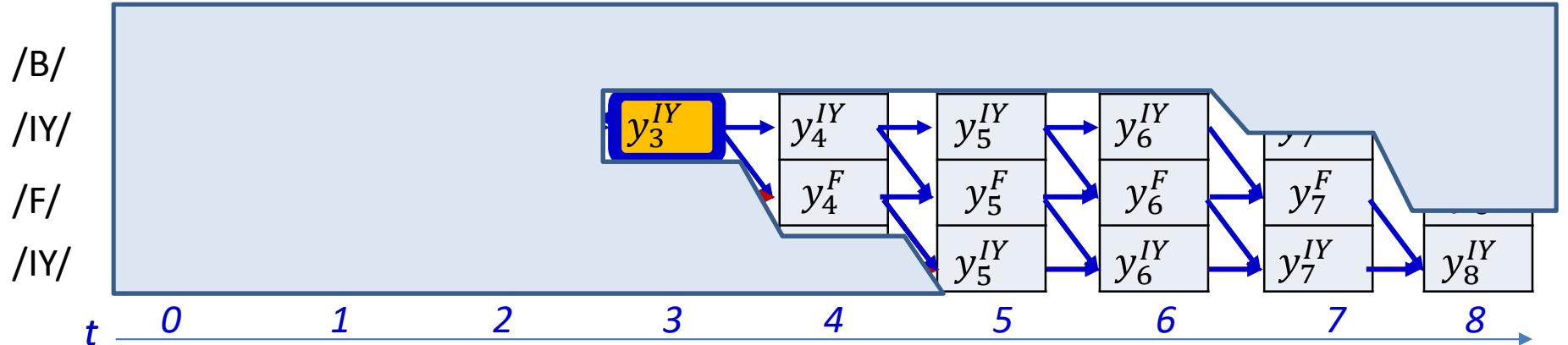
# The backward probability



$$\beta(t, r) = P(S_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})$$

- $\beta(t, r)$  is the probability of the exposed subgraph, not including the orange shaded box

# The backward probability

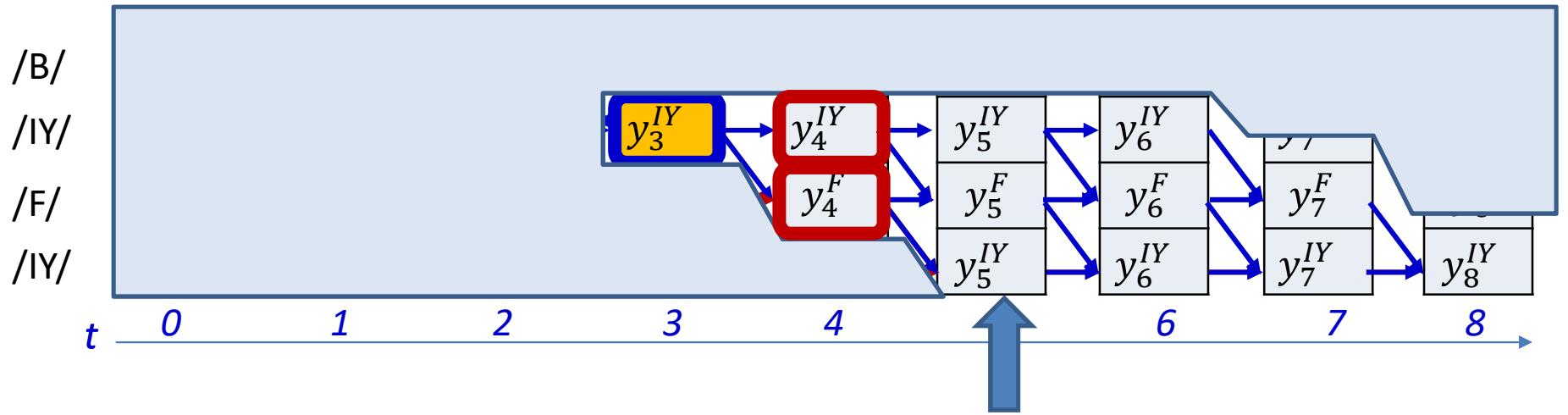


$$\beta(t, r) = P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})$$

- $\beta(t, r)$  is the probability of the exposed subgraph, not including the orange shaded box
- Note that RHS includes *all* potential successors of  $S_r$ . Lets expand this out by explicitly summing over all potential successors

$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(s_{t+1} = S_q, S_q, \dots, S_K | \mathbf{X})$$

# The backward probability

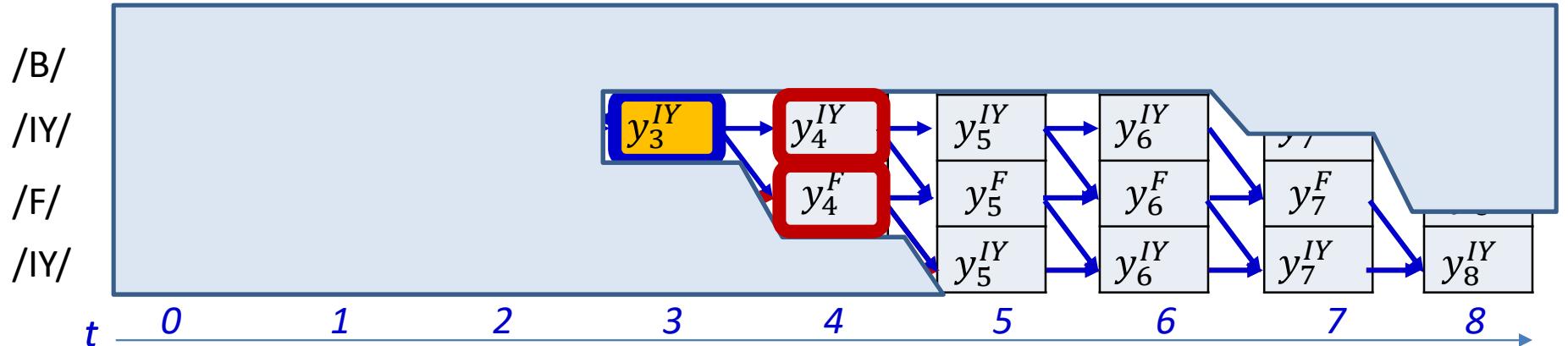


$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(S_{t+1} = S_q, S_q, \dots, S_K | \mathbf{X})$$

- Lets expand this out in terms of the *successors of  $S_q$*  at  $t+2$ 
  - Explicitly consider all possible successors, to cover all possibilities (the two red boxes)

$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(S_{t+1} = S_q, S_{t+2} \in \text{succ}(S_q), \text{succ}(S_q), \dots, S_K | \mathbf{X})$$

# The backward probability



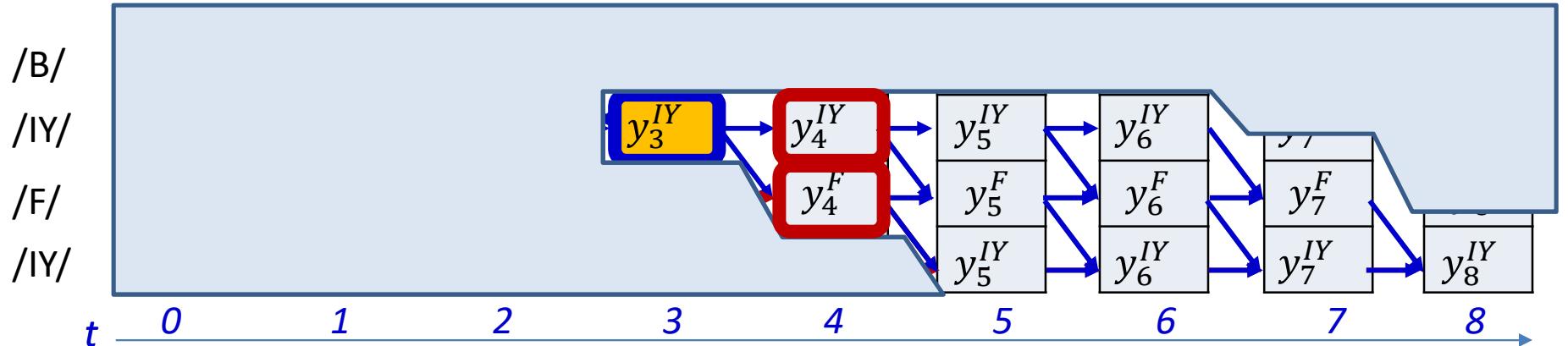
- Expressing  $\beta(t, r)$  in terms of the successors of  $S_q$

$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(s_{t+1} = S_q, s_{t+2} \in \text{succ}(S_q), \dots, s_K | \mathbf{X})$$

- Using our assumption of conditional independence

$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(s_{t+1} = S_q | \mathbf{X}) P(s_{t+2} \in \text{succ}(S_q), \dots, s_K | \mathbf{X})$$

# The backward probability



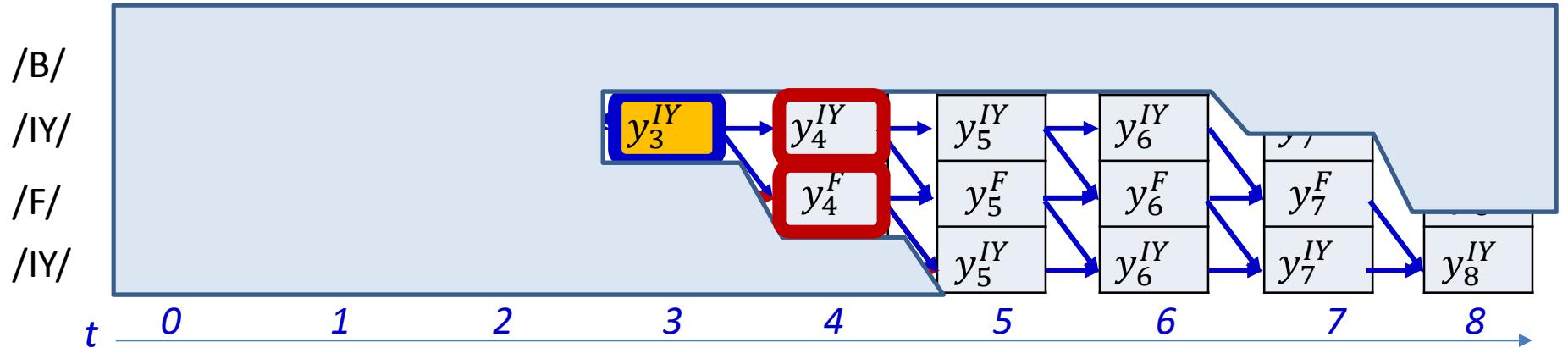
- Expressing  $\beta(t, r)$  in terms of the successors of  $S_q$

$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(s_{t+1} = S_q, s_{t+2} \in \text{succ}(S_q), \dots, s_K | \mathbf{X})$$

- Using our assumption of copy conditional independence  $\beta(t + 1, q)$

$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(s_{t+1} = S_q | \mathbf{X}) P(s_{t+2} \in \text{succ}(S_q), \dots, s_K | \mathbf{X})$$

# The backward probability



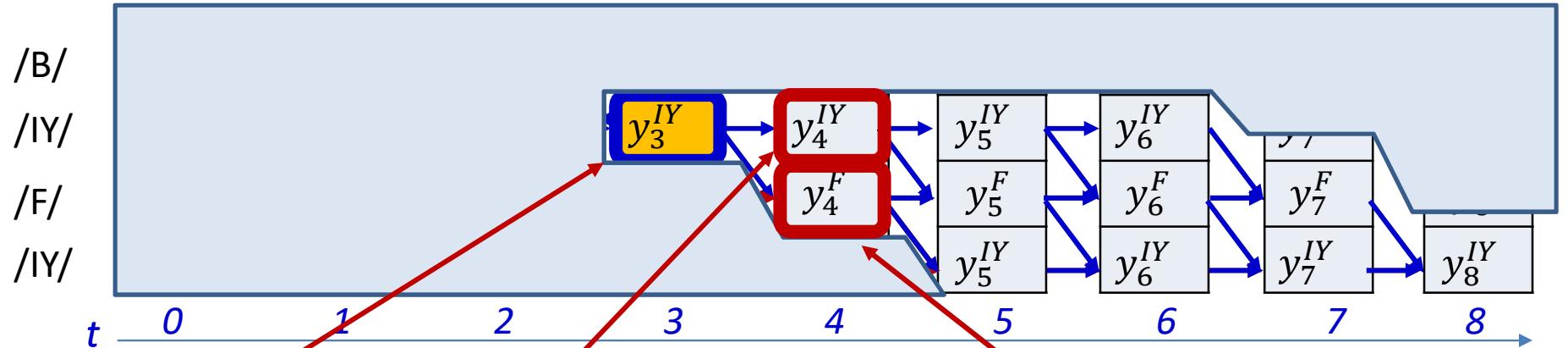
- Expressing  $\beta(t, r)$  in terms of the successors of  $S_q$

$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} P(s_{t+1} = S_q, s_{t+2} \in \text{succ}(S_q), \dots, s_K | \mathbf{X})$$

- Using our assumption of conditional independence

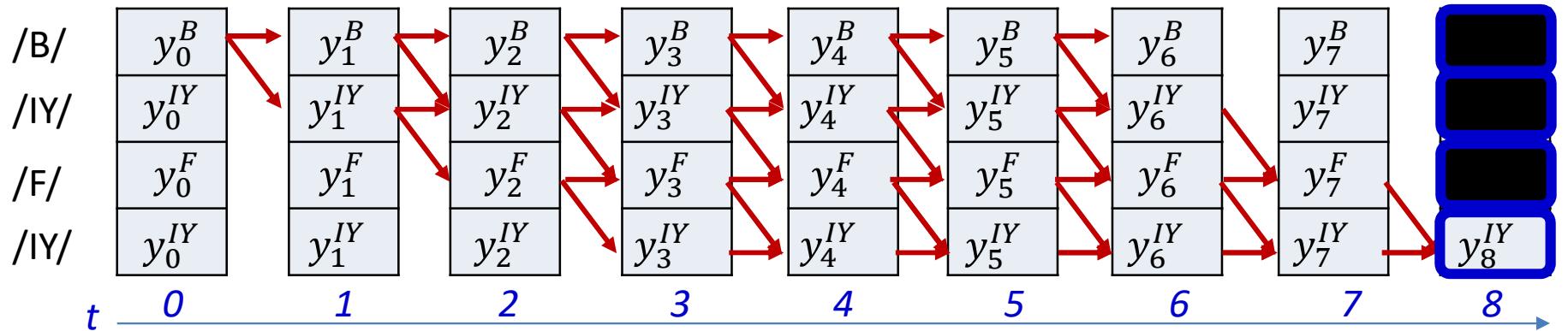
$$\beta(t, r) = \sum_{q: S_q \in \text{succ}(S_r)} y_{t+1}^{S_q} \beta(t + 1, q)$$

# Backward algorithm



$$\beta(t, r) = y_{t+1}^{S(r)} \beta(t + 1, r) + y_{t+1}^{S(r+1)} \beta(t + 1, r + 1)$$

# Backward algorithm



- Initialization:

$$\beta(T-1, K) = 1, \quad \beta(T-1, r) = 0, \quad r < K$$



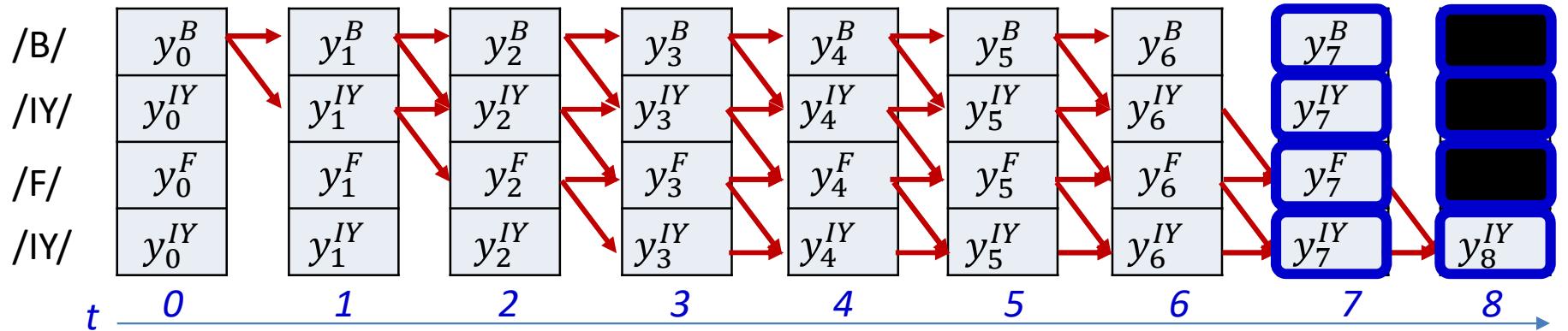
- for  $t = T-2$  down to 0

$$\beta(t, K) = \beta(t+1, K) y_{t+1}^{S(K)}$$

for  $l = K-1 \dots 1$

- $\beta(t, r) = y_{t+1}^{S(l)} \beta(t+1, r) + y_{t+1}^{S(r+1)} \beta(t+1, r+1)$

# Backward algorithm



- Initialization:

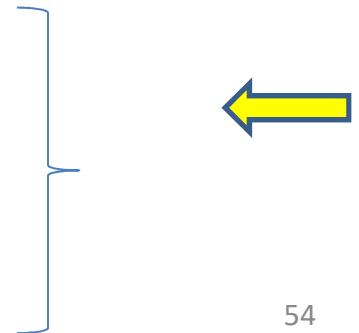
$$\beta(T-1, K) = 1, \quad \beta(T-1, r) = 0, \quad r < K$$

- for  $t = T-2$  down to 0

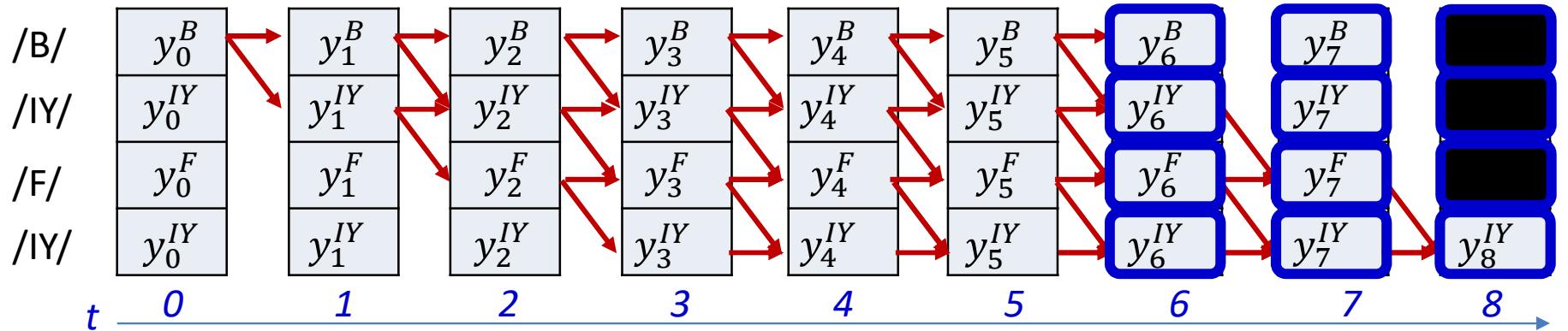
$$\beta(t, K) = \beta(t+1, K) y_{t+1}^{S(K)}$$

for  $l = K-1 \dots 1$

- $$\beta(t, r) = y_{t+1}^{S(l)} \beta(t+1, r) + y_{t+1}^{S(r+1)} \beta(t+1, r+1)$$



# Backward algorithm



- Initialization:

$$\beta(T-1, K) = 1, \quad \beta(T-1, r) = 0, \quad r < K$$

- for  $t = T - 2$  down to 0

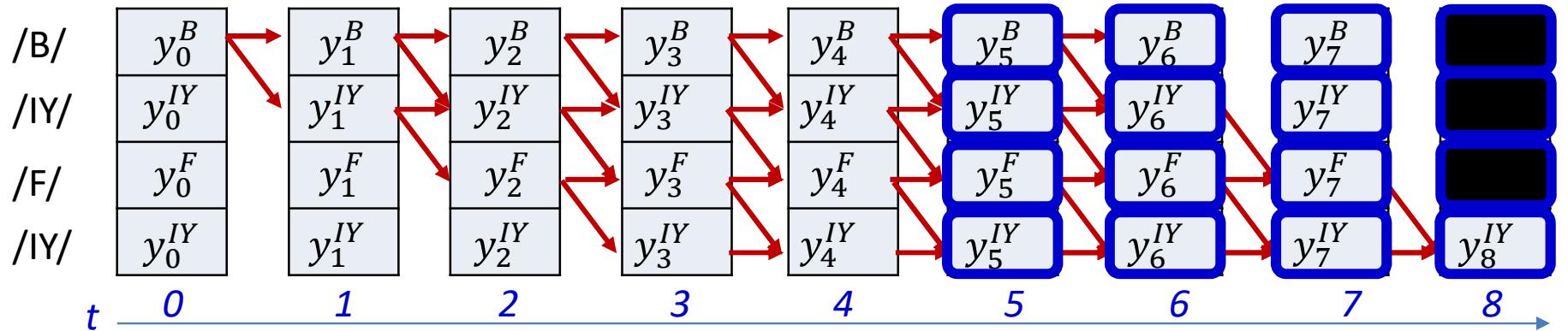
$$\beta(t, K) = \beta(t+1, K) y_{t+1}^{S(K)}$$

for  $l = K - 1 \dots 1$

- $\beta(t, r) = y_{t+1}^{S(l)} \beta(t+1, r) + y_{t+1}^{S(r+1)} \beta(t+1, r+1)$



# Backward algorithm



- Initialization:

$$\beta(T-1, K) = 1, \quad \beta(T-1, r) = 0, \quad r < K$$

- for  $t = T-2$  down to 0

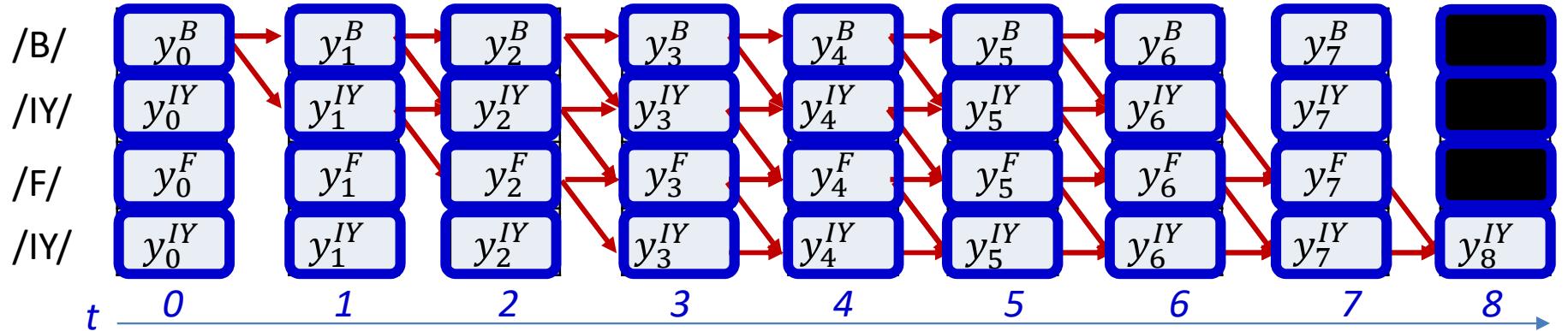
$$\beta(t, K) = \beta(t+1, K) y_{t+1}^{S(K)}$$

for  $l = K-1 \dots 1$

- $$\beta(t, r) = y_{t+1}^{S(l)} \beta(t+1, r) + y_{t+1}^{S(r+1)} \beta(t+1, r+1)$$



# Backward algorithm



- Initialization:

$$\beta(T-1, K) = 1, \quad \beta(T-1, r) = 0, \quad r < K$$

- for  $t = T - 2$  down to 0

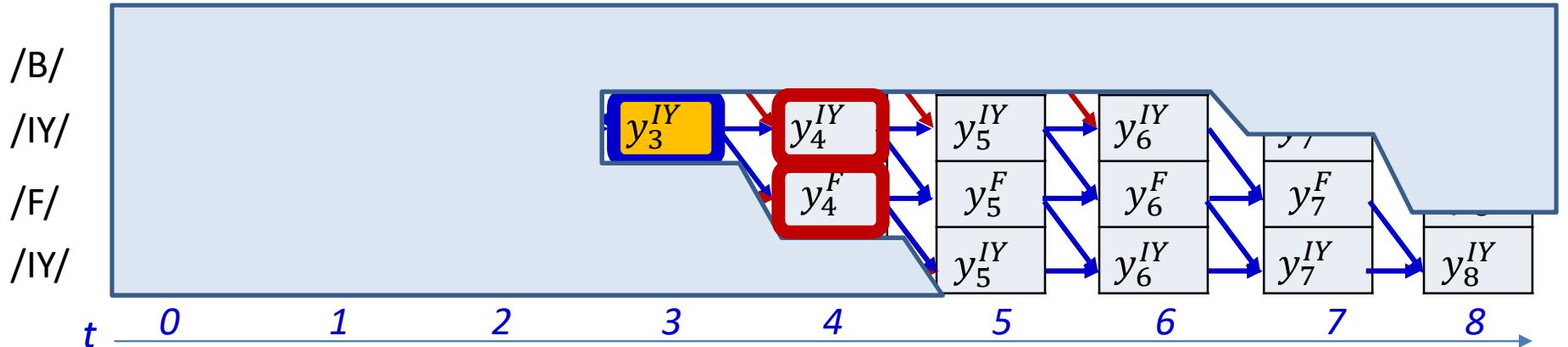
$$\beta(t, K) = \beta(t+1, K) y_{t+1}^{S(K)}$$

for  $l = K - 1 \dots 1$

- $\beta(t, r) = y_{t+1}^{S(l)} \beta(t+1, r) + y_{t+1}^{S(r+1)} \beta(t+1, r+1)$



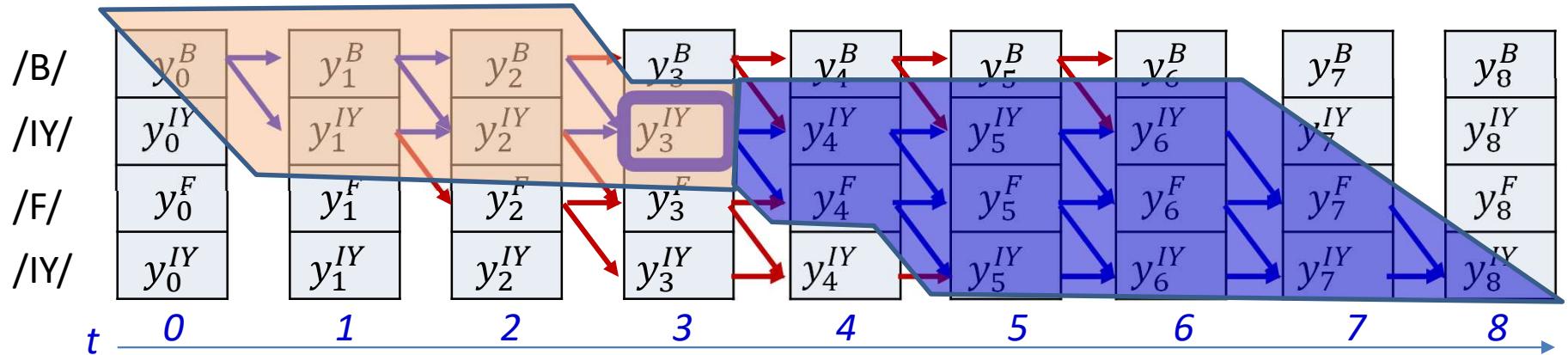
# Alternate Backward algorithm



$$\hat{\beta}(t, r) = y_t^{S(r)} (\hat{\beta}(t + 1, r) + \hat{\beta}(t + 1, r + 1))$$

- Some implementations of the backward algorithm will use the above formula
- Note that here the probability of the observation at  $t$  is also factored into beta
- It will have to be unfactored later (we'll see how)

# The joint probability

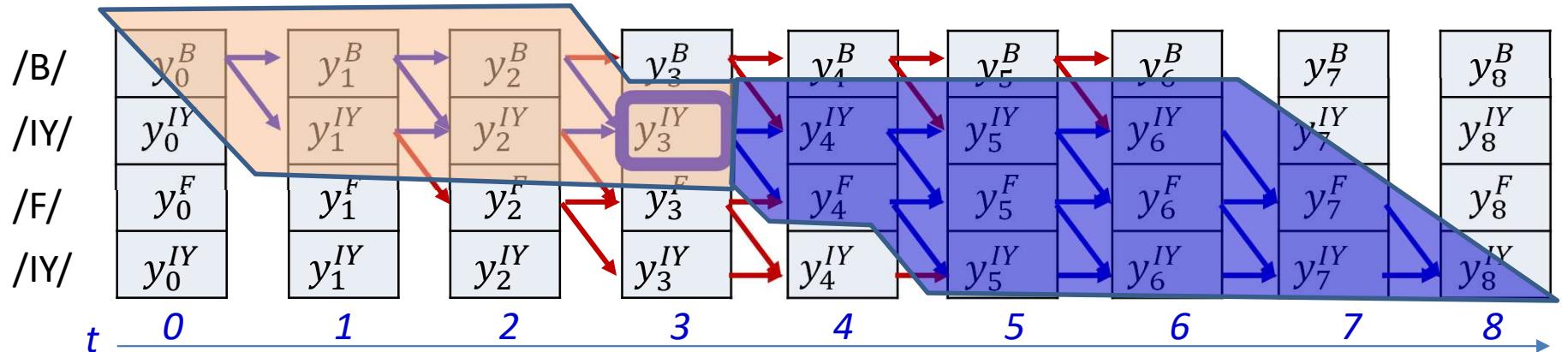


$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = \alpha(t, r) P(s_{t+1} \in \text{succ}(S_r), \text{succ}(S_r), \dots, S_K | \mathbf{X})$$

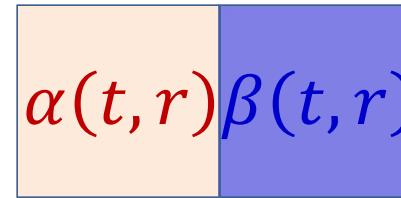
- We will call the first term the *forward probability*  
 $\alpha(t, r)$
- We will call the second term the *backward probability*  
 $\beta(t, r)$

We now can compute this

# The joint probability



$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = \alpha(t, r) \beta(t, r)$$

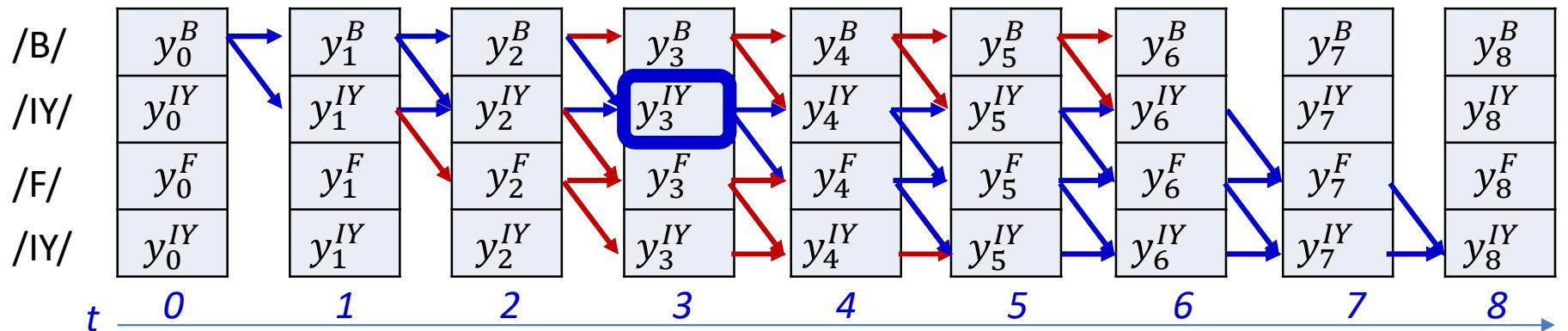


- We will call the first term the *forward probability*  $\alpha(t, r)$
- We will call the second term the *backward probability*  $\beta(t, r)$

Forward algo

Backward algo

# The posterior probability

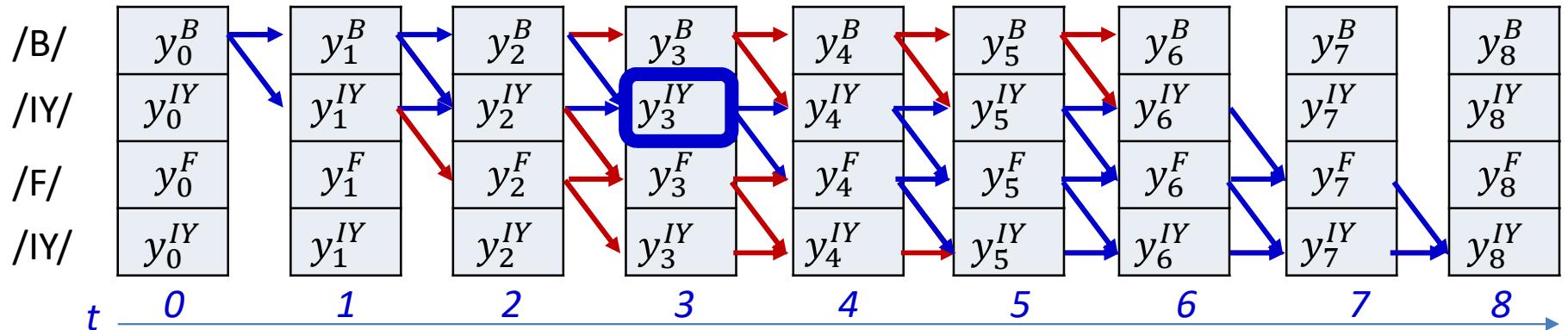


$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = \alpha(t, r) \beta(t, r)$$

- The *posterior* is given by

$$P(s_t = S_r | \mathbf{S}, \mathbf{X}) = \frac{P(s_t = S_r, \mathbf{S} | \mathbf{X})}{\sum_{S'_r} P(s_t = S'_r, \mathbf{S} | \mathbf{X})} = \frac{\alpha(t, r) \beta(t, r)}{\sum_{r'} \alpha(t, r') \beta(t, r')}$$

# The posterior probability

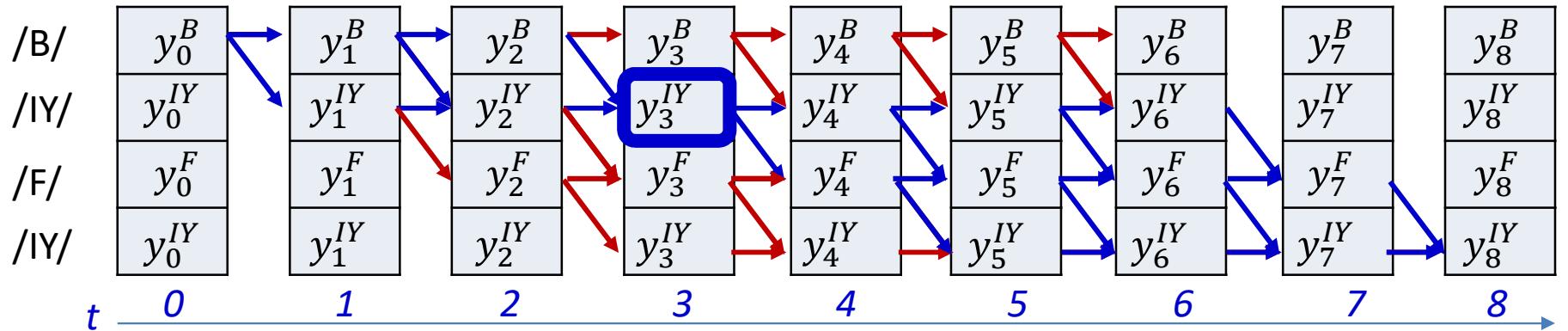


- Let the posterior  $P(s_t = S_r | \mathbf{S}, \mathbf{X})$  be represented by  $\gamma(t, r)$

$$\gamma(t, r) = \frac{\alpha(t, r)\beta(t, r)}{\sum_{r'} \alpha(t, r')\beta(t, r')}$$

Note that it is a function of  $y_t^{S_r}$  through  $\alpha(t, r)$

# The posterior probability



$$P(s_t = S_r, \mathbf{S} | \mathbf{X}) = \alpha(t, r) \beta(t, r)$$

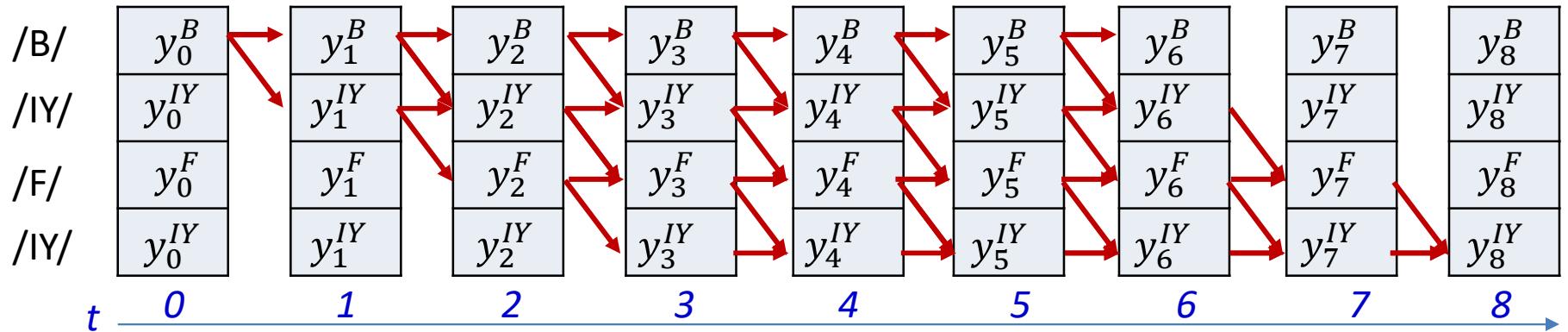
- The *posterior* is given by

$$\gamma(t, r) = \frac{\alpha(t, r) \beta(t, r)}{\sum_{r'} \alpha(t, r') \beta(t, r')}$$

- We can also write this using the modified beta formula as (you will see this in papers)

$$\gamma(t, r) = \frac{\frac{1}{y_t^{S(r)}} \alpha(t, r) \hat{\beta}(t, r)}{\sum_{r'} \frac{1}{y_t^{S(r)}} \alpha(t, r) \hat{\beta}(t, r)}$$

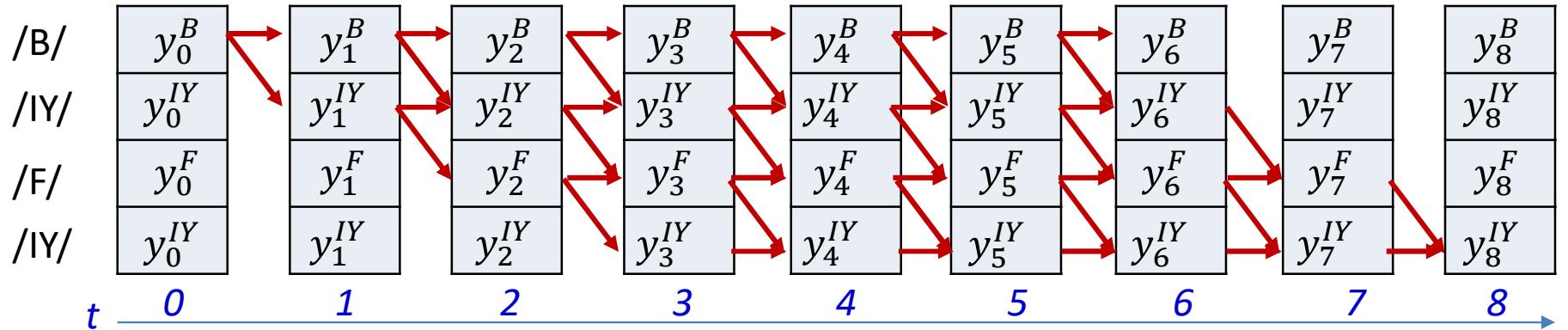
# The expected divergence



$$DIV = - \sum_t \sum_{s \in S_1 \dots S_K} P(s_t = s | \mathbf{S}, \mathbf{X}) \log Y(t, s_t = s)$$

$$DIV = - \sum_t \sum_r \gamma(t, r) \log \textcolor{red}{y}_t^{S(r)}$$

# The expected divergence



$$DIV = - \sum_t \sum_{s \in S_1 \dots S_K} P(s_t = s | \mathbf{S}, \mathbf{X}) \log Y(t, s_t = s)$$

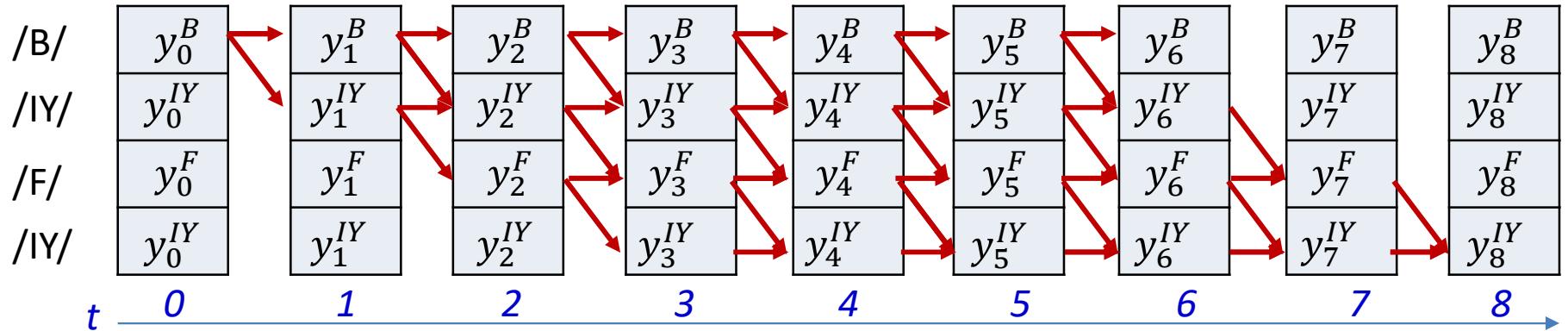
$$DIV = - \sum_t \sum_r \gamma(t, r) \log \mathbf{y}_t^{S(r)}$$

- The derivative of the divergence w.r.t the output  $Y_t$  of the net at any time:

$$\nabla_{Y_t} DIV = \left[ \frac{dDIV}{dy_t^1} \quad \frac{dDIV}{dy_t^2} \quad \dots \quad \frac{dDIV}{dy_t^L} \right]$$

- Components will be non-zero only for symbols that occur in the training instance

# The expected divergence



$$DIV = - \sum_t \sum_{s \in S_1 \dots S_K} P(s_t = s | \mathbf{S}, \mathbf{X}) \log Y(t, s_t = s)$$

$$DIV = - \sum_t \sum_r \gamma(t, r) \log y_t^{s(r)}$$

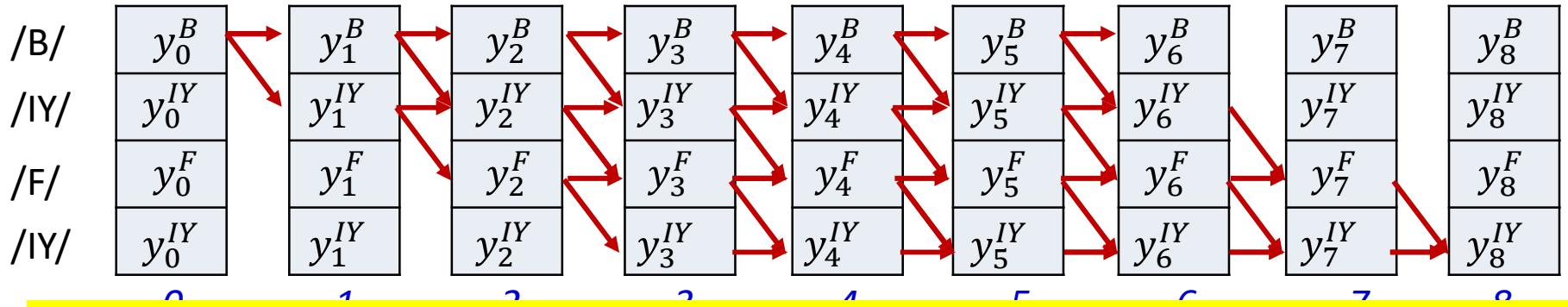
- The derivative of the divergence w.r.t the output  $Y_t$  of the net at any time:

$$\nabla_{Y_t} DIV = \left[ \frac{dDIV}{dy_t^1} \right] \left[ \frac{dDIV}{dy_t^2} \right] \dots$$

Must compute these terms from here, keeping in mind that  $\gamma(t, r)$  is also a function of  $y_t^{s(r)}$

- Components will be non-zero only for symbols that occur

# The expected divergence



$$\frac{dDIV}{dy_t^l} = - \sum_{r : S(r)=l} \frac{d}{dy_t^{S(r)}} \left( \frac{\alpha(t, r) \beta(t, r)}{\sum_{r'} \alpha(t, r') \beta(t, r')} \log y_t^{S(r)} \right)$$

$$DIV = - \sum_t \sum_r \gamma(t, r) \log y_t^{S(r)}$$

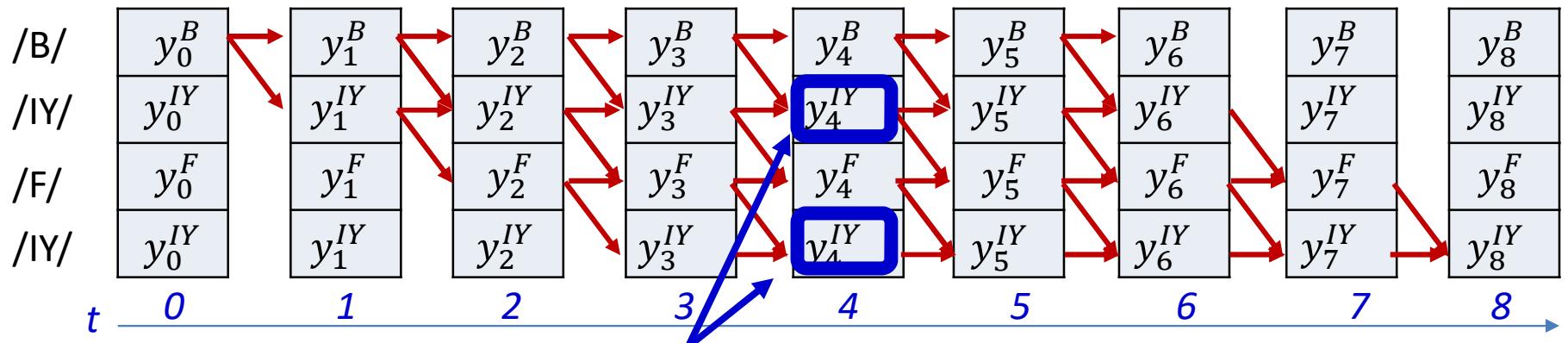
- The derivative of the divergence w.r.t the output  $Y_t$  of the net at any time:

$$\nabla_{Y_t} DIV = \left[ \frac{dDIV}{dy_t^1} \quad \frac{dDIV}{dy_t^2} \quad \dots \quad \frac{dDIV}{dy_t^S} \right]$$

- Components will be non-zero only for symbols that occur in the sequence.

Must compute these terms from here, keeping in mind that  $\gamma(t, r)$  is also a function of  $y_t^{S(r)}$

# The expected divergence



The derivatives at both these locations must be summed to get  $\frac{dDIV}{dy_4^5}$

$$DIV = - \sum_t \sum_r \gamma(t, r) \log y_t^{S(r)}$$

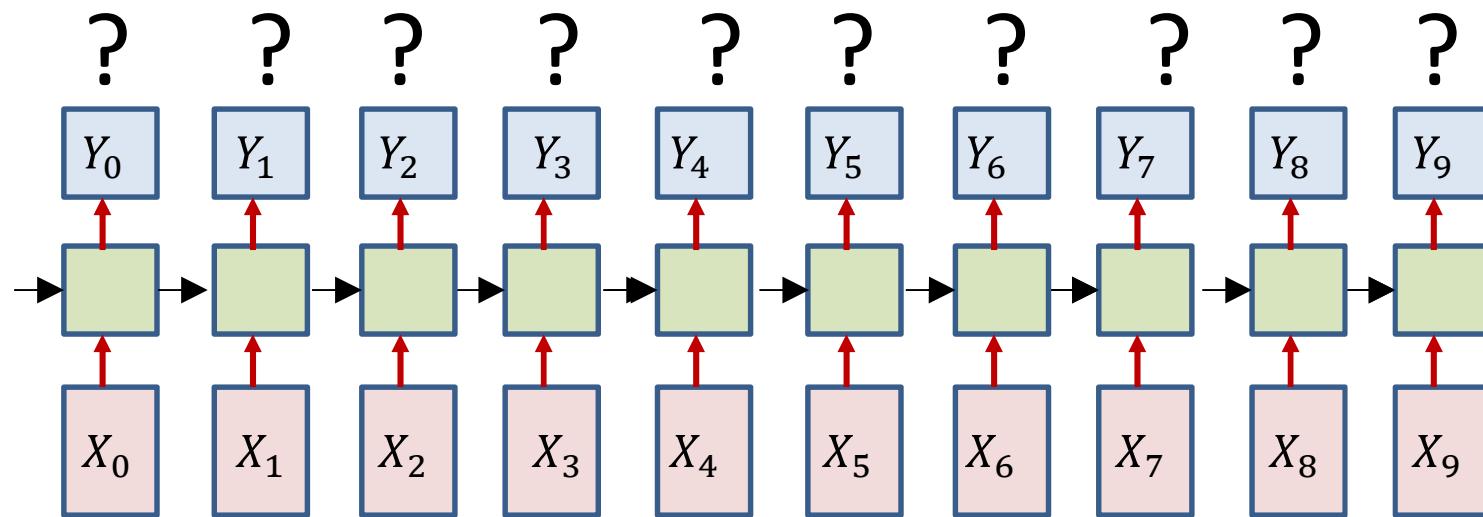
- The derivative of the divergence w.r.t any particular output of the network must sum over all instances of that symbol in the target sequence

$$\frac{dDIV}{dy_t^l} = - \sum_{r : S(r)=l} \frac{d}{dy_t^{S(r)}} (\gamma(t, r) \log y_t^{S(r)})$$

- E.g. the derivative w.r.t  $y_t^5$  will sum over both rows representing /IY/ in the above figure

# Overall training procedure for Seq2Seq case 1

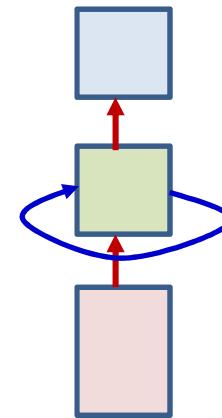
/B/ /IY/ /F/ /IY/



- Problem: Given input and output sequences without alignment, train models

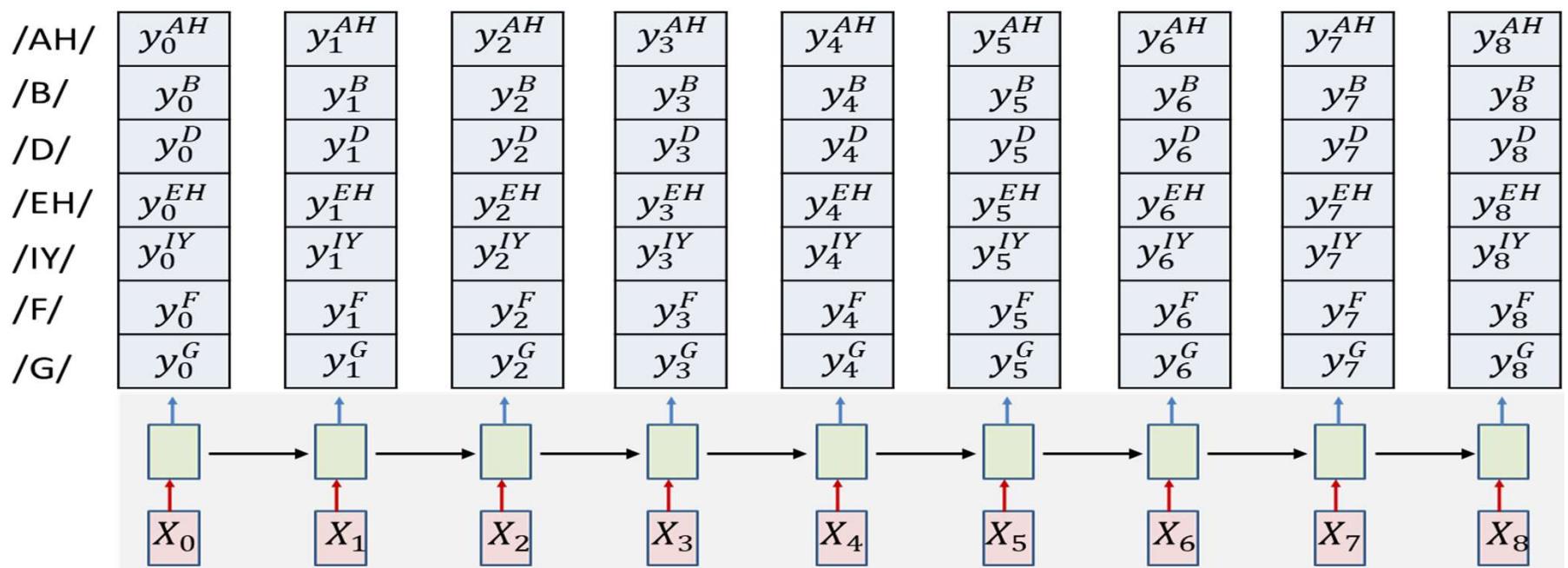
# Overall training procedure for Seq2Seq case 1

- **Step 1:** Setup the network
  - Typically many-layered LSTM
- **Step 2:** Initialize all parameters of the network

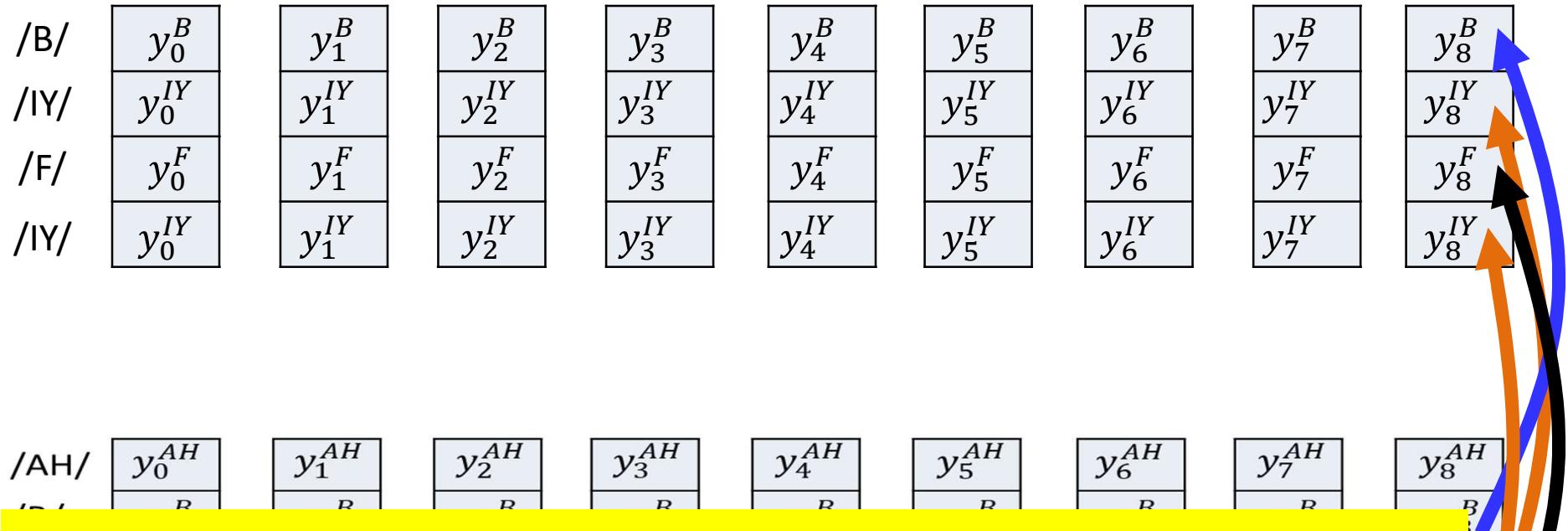


# Overall Training: Forward pass

- Foreach training instance
  - **Step 3:** Forward pass. Pass the training instance through the network and obtain all symbol probabilities at each time

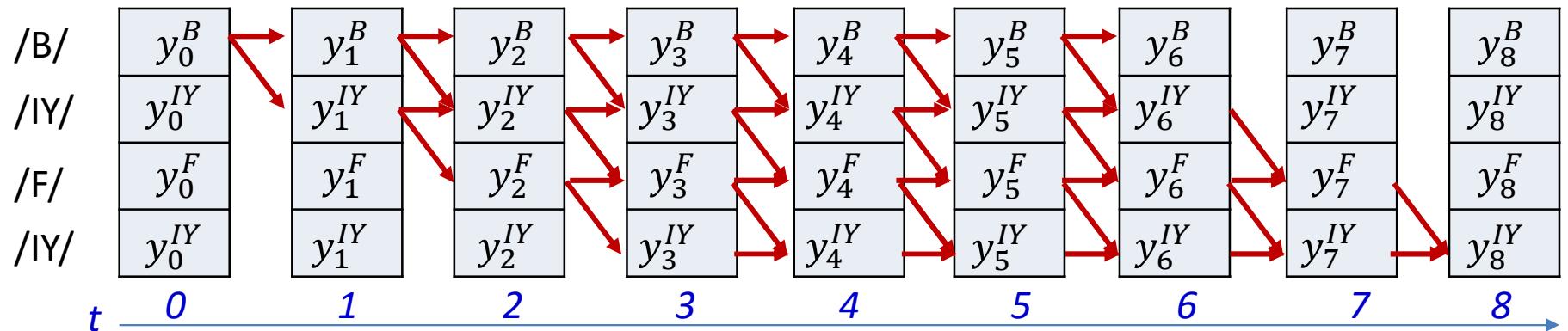


# Overall training: Backward pass



- Foreach training instance
  - **Step 3:** Forward pass. Pass the training instance through the network and obtain all symbol probabilities at each time
  - **Step 4:** Construct the graph representing the specific symbol sequence in the instance. This may require having multiple rows of nodes with the same symbol scores

# Overall training: Backward pass



- Fforeach training instance:
  - **Step 5:** Perform the forward backward algorithm to compute  $\alpha(t, r)$  and  $\beta(t, r)$  at each time, for each row of nodes in the graph
  - **Step 6:** Compute derivative of divergence  $\nabla_{Y_t} DIV$  for each  $Y_t$

# Overall training: Backward pass

- Foreach instance
  - **Step 6:** Compute derivative of divergence  $\nabla_{Y_t} DIV$  for each  $Y_t$

$$\nabla_{Y_t} DIV = \begin{bmatrix} \frac{dDIV}{dy_t^1} & \frac{dDIV}{dy_t^2} & \dots & \frac{dDIV}{dy_t^L} \end{bmatrix}$$

$$\frac{dDIV}{dy_t^l} = - \sum_{r : S(r)=l} \frac{d}{dy_t^{S(r)}} \left( \frac{\alpha(t, r) \beta(t, r)}{\sum_{r'} \alpha(t, r') \beta(t, r')} \log y_t^{S(r)} \right)$$

- **Step 7:** Aggregate derivatives over minibatch and update parameters

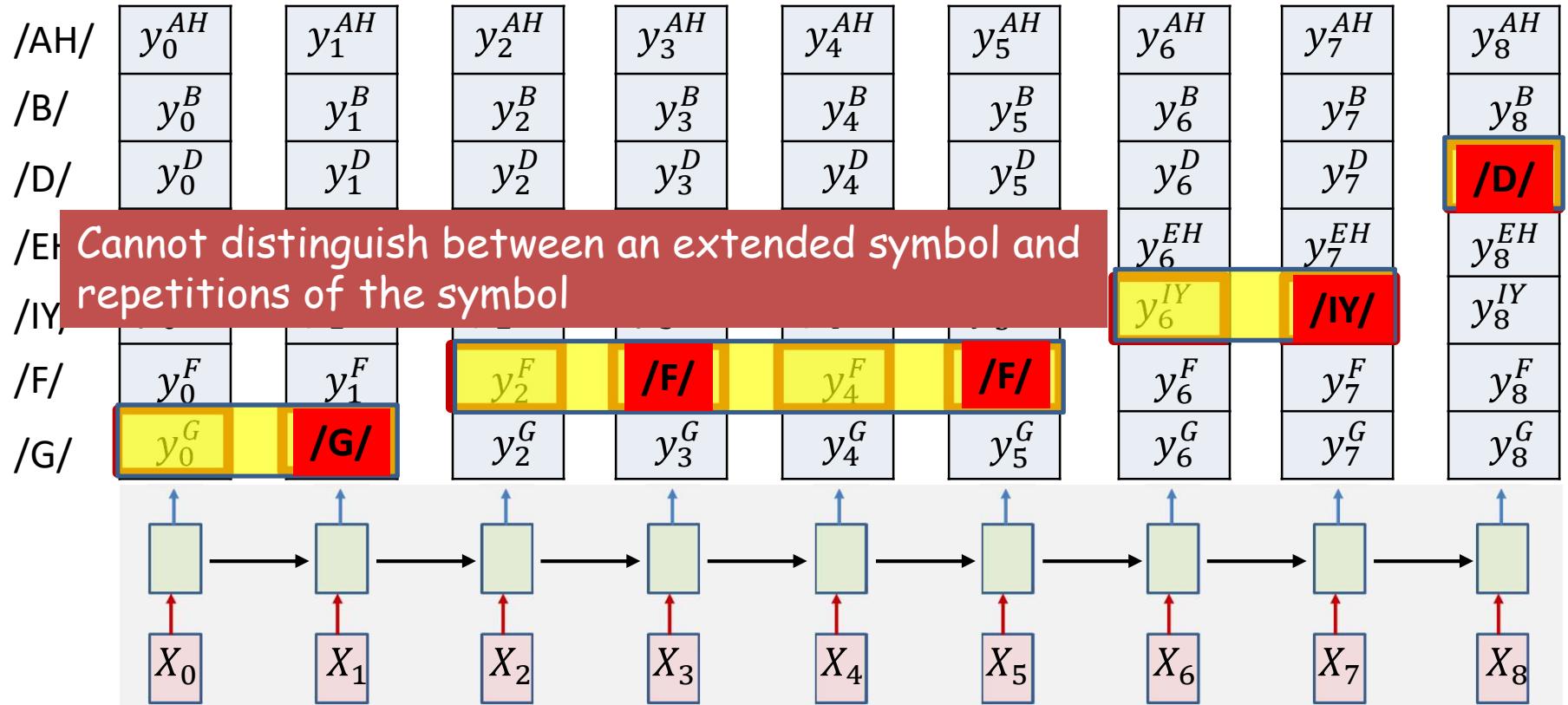
# Story so far: CTC models

- Sequence-to-sequence networks which irregularly output symbols can be “decoded” by Viterbi decoding
  - Which assumes that a symbol is output at each time and *merges* adjacent symbols
- They require alignment of the output to the symbol sequence for training
  - This alignment is generally not given
- Training can be performed by iteratively estimating the alignment by Viterbi-decoding and time-synchronous training
- Alternately, it can be performed by optimizing the expected error over *all* possible alignments
  - Posterior probabilities for the expectation can be computed using the forward backward algorithm

# A key *decoding* problem

- Consider a problem where the output symbols are characters
- We have a decode: R R R O O O O O D
- Is this the merged symbol sequence ROD or ROOD?

# We've seen this before



- /G/ /F/ /F/ /IY/ /D/ or /G/ /F/ /IY/ /D/ ?

# A key *decoding* problem

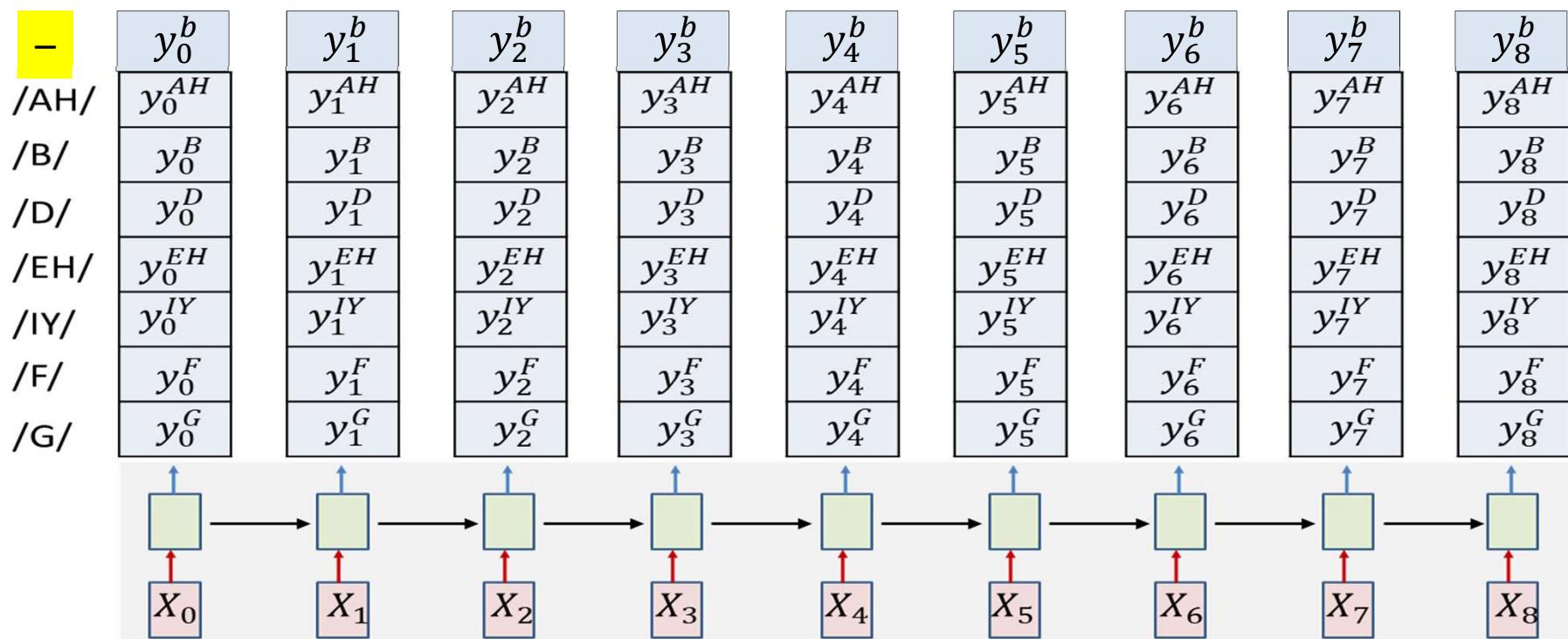
- Consider a problem where the output symbols are characters
- We have a decode: R R R O O O O O D
- Is this the symbol sequence ROD or ROOD?
- Note: This problem does not always occur, e.g. when symbols have sub symbols
  - E.g. If O is produced as O1 and O2 and must always begin with O1 and terminate with O2
    - A single O would be of the form O1 O1 .. O2 → O
    - Multiple Os would have the decode O1 .. O2.. O1..O2.. → oo

# A key *decoding* problem

- We have a decode: R R R O O O O O D
- Is this the symbol sequence ROD or ROOD?
- Solution: Introduce an explicit extra symbol which serves to separate discrete versions of a symbol
  - A “blank” (represented by “-”)
  - RRR---OO---DDD = ROD
  - RR-R---OO---D-DD = RRODD
  - R-R-R---O-O DD-D DDD-D = RRROODDD
    - The next symbol at the end of a sequence of blanks is always a new character
    - When a symbol repeats, there must be at least one blank between the repetitions
- The symbol set recognized by the network must now include the extra blank symbol
  - Which too must be trained

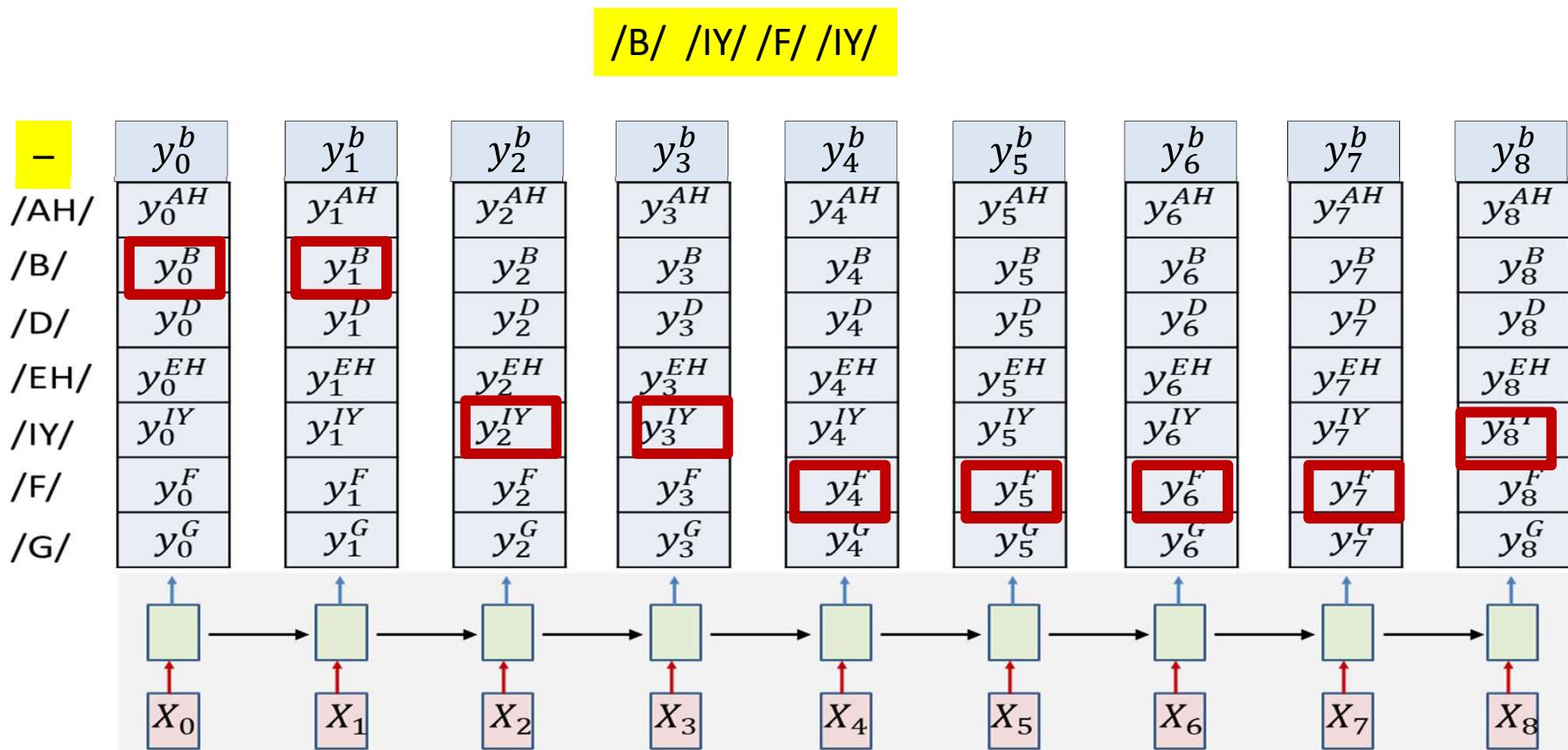
# The modified forward output

- Note the extra “blank” at the output



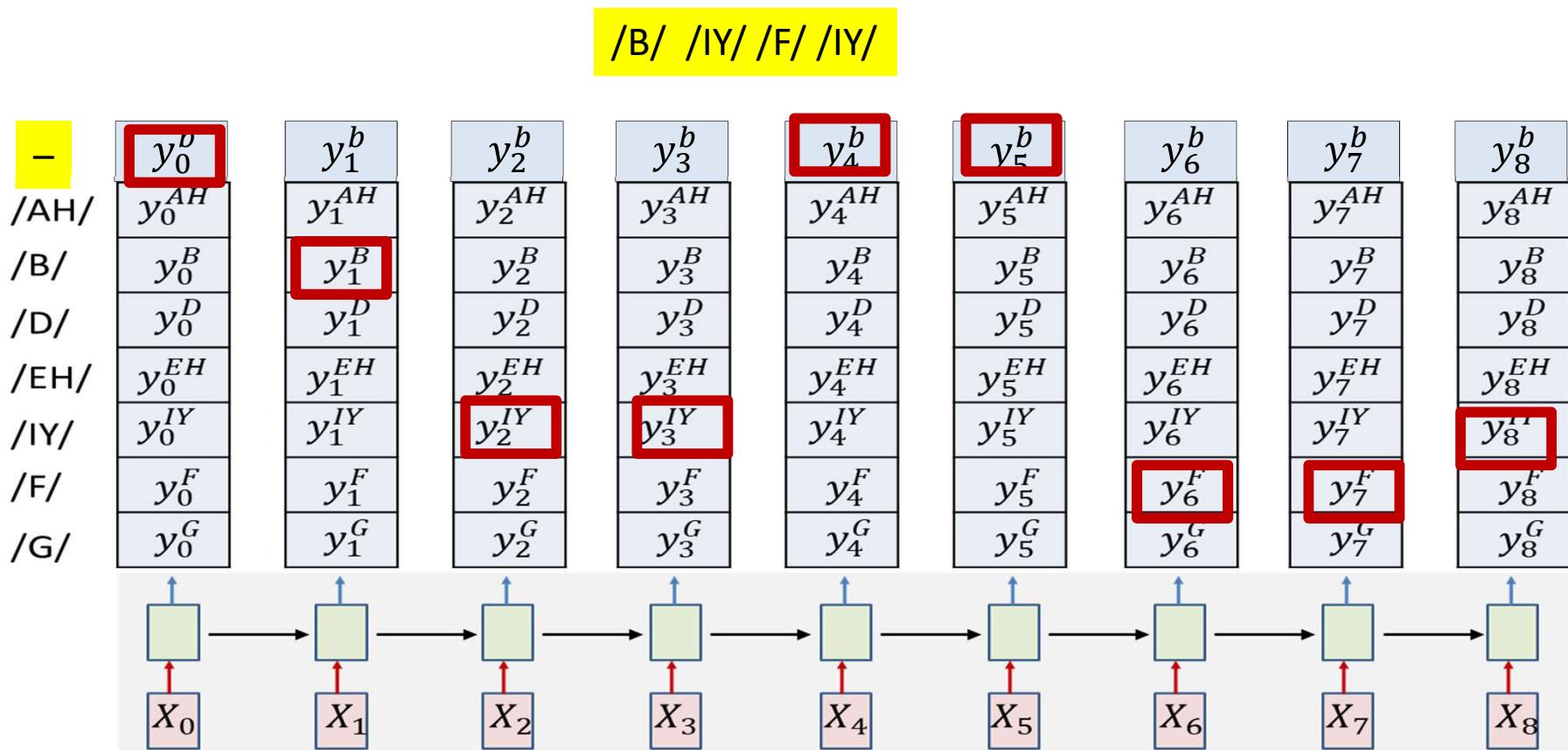
# The modified forward output

- Note the extra “blank” at the output



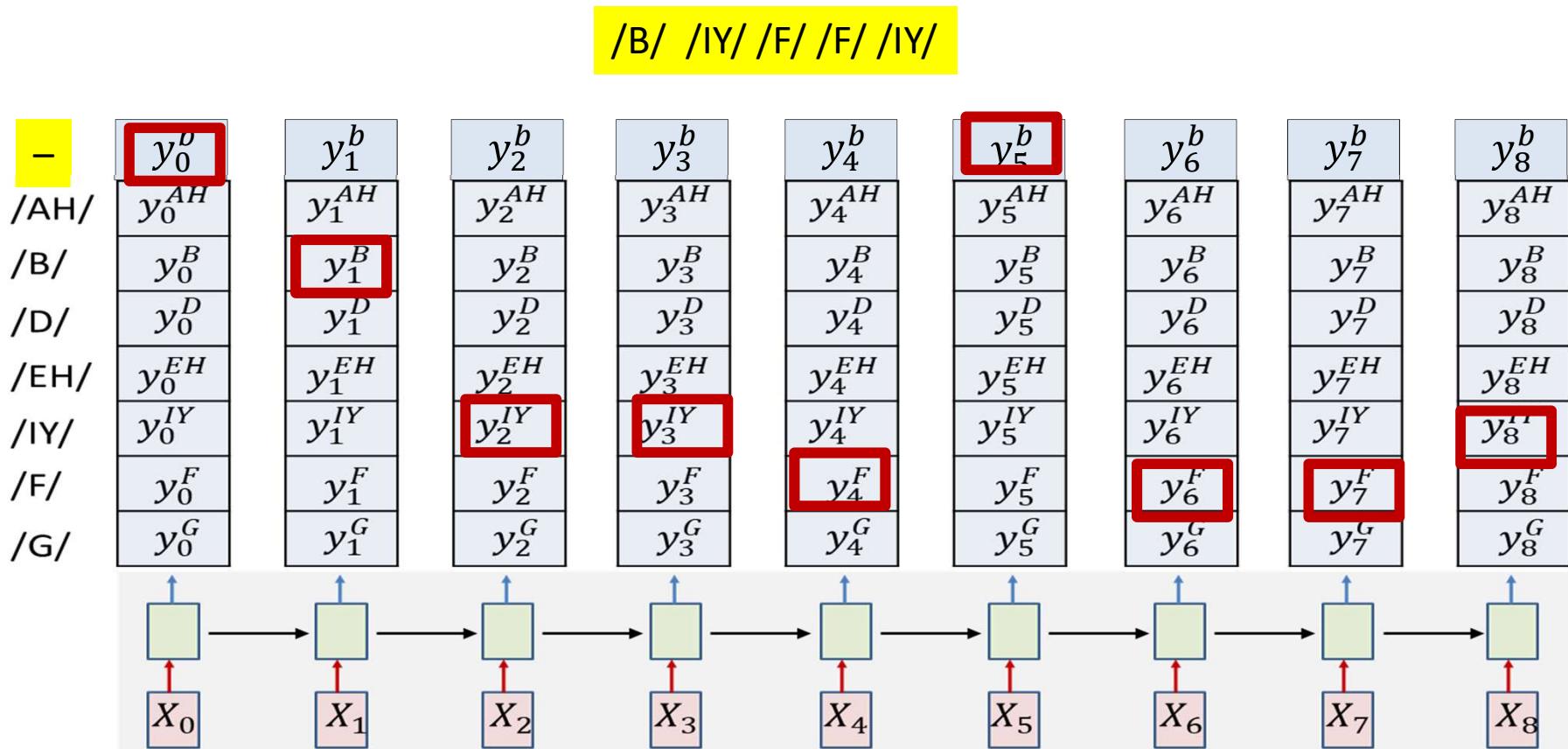
# The modified forward output

- Note the extra “blank” at the output

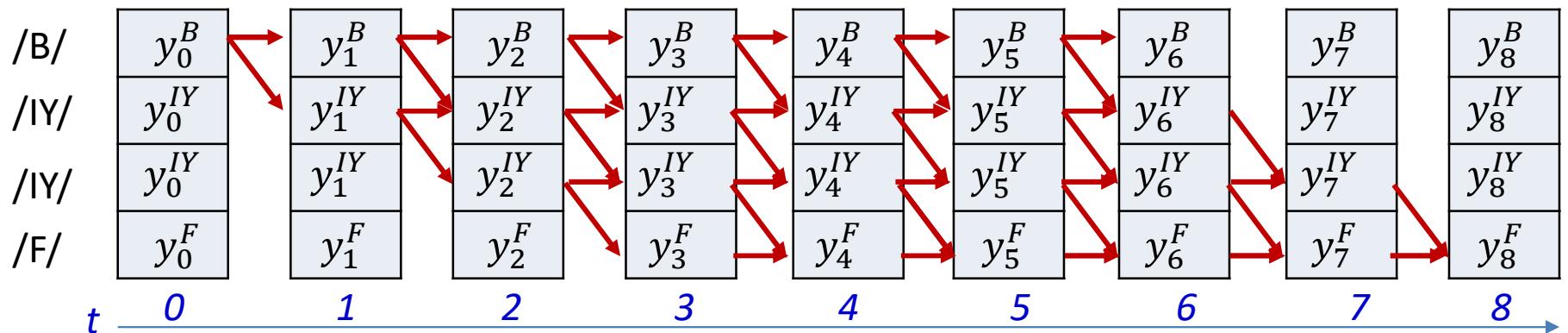


# The modified forward output

- Note the extra “blank” at the output



# Composing the graph for training



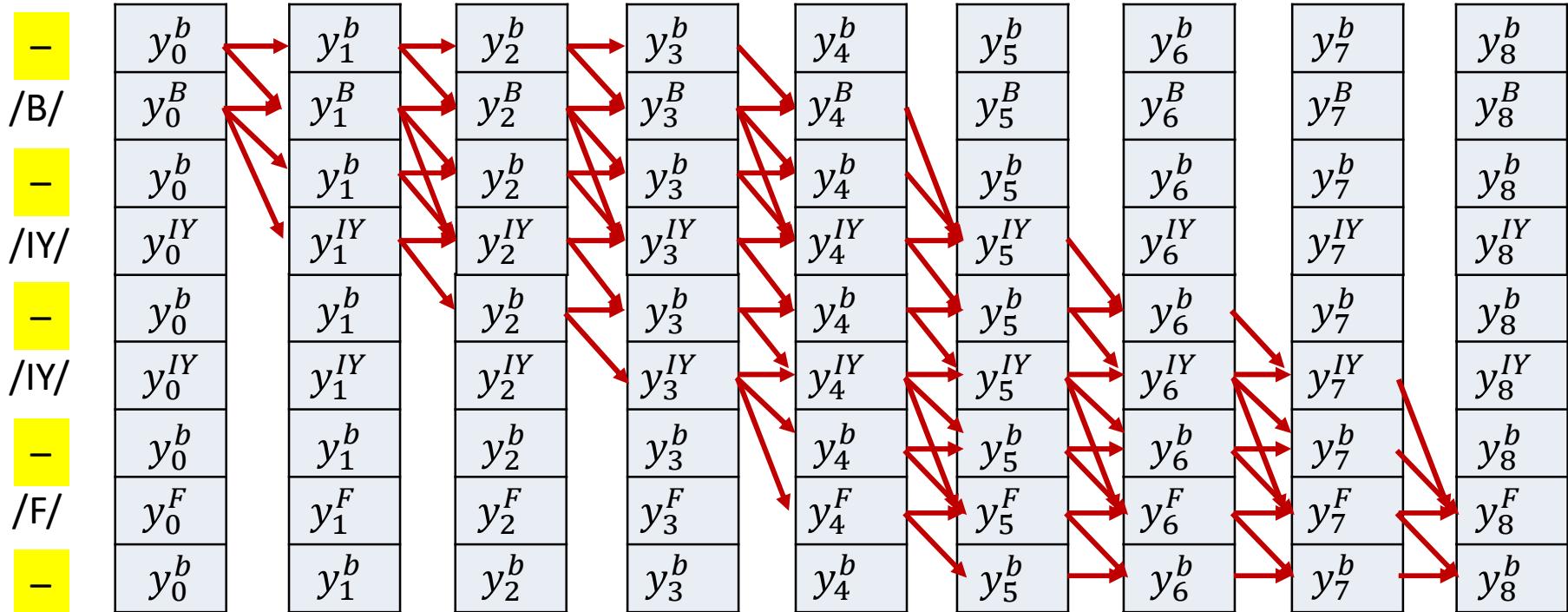
- The original method without blanks
- **Changing the example to /B/ /IY/ /IY/ /F/ from /B/ /IY/ /F/ /IY/ for illustration**

# Composing the graph for training

|      |            |            |            |            |            |            |            |            |            |
|------|------------|------------|------------|------------|------------|------------|------------|------------|------------|
| -    | $y_0^b$    | $y_1^b$    | $y_2^b$    | $y_3^b$    | $y_4^b$    | $y_5^b$    | $y_6^b$    | $y_7^b$    | $y_8^b$    |
| /B/  | $y_0^B$    | $y_1^B$    | $y_2^B$    | $y_3^B$    | $y_4^B$    | $y_5^B$    | $y_6^B$    | $y_7^B$    | $y_8^B$    |
| -    | $y_0^b$    | $y_1^b$    | $y_2^b$    | $y_3^b$    | $y_4^b$    | $y_5^b$    | $y_6^b$    | $y_7^b$    | $y_8^b$    |
| /IY/ | $y_0^{IY}$ | $y_1^{IY}$ | $y_2^{IY}$ | $y_3^{IY}$ | $y_4^{IY}$ | $y_5^{IY}$ | $y_6^{IY}$ | $y_7^{IY}$ | $y_8^{IY}$ |
| -    | $y_0^b$    | $y_1^b$    | $y_2^b$    | $y_3^b$    | $y_4^b$    | $y_5^b$    | $y_6^b$    | $y_7^b$    | $y_8^b$    |
| /IY/ | $y_0^{IY}$ | $y_1^{IY}$ | $y_2^{IY}$ | $y_3^{IY}$ | $y_4^{IY}$ | $y_5^{IY}$ | $y_6^{IY}$ | $y_7^{IY}$ | $y_8^{IY}$ |
| -    | $y_0^b$    | $y_1^b$    | $y_2^b$    | $y_3^b$    | $y_4^b$    | $y_5^b$    | $y_6^b$    | $y_7^b$    | $y_8^b$    |
| /F/  | $y_0^F$    | $y_1^F$    | $y_2^F$    | $y_3^F$    | $y_4^F$    | $y_5^F$    | $y_6^F$    | $y_7^F$    | $y_8^F$    |
| -    | $y_0^b$    | $y_1^b$    | $y_2^b$    | $y_3^b$    | $y_4^b$    | $y_5^b$    | $y_6^b$    | $y_7^b$    | $y_8^b$    |

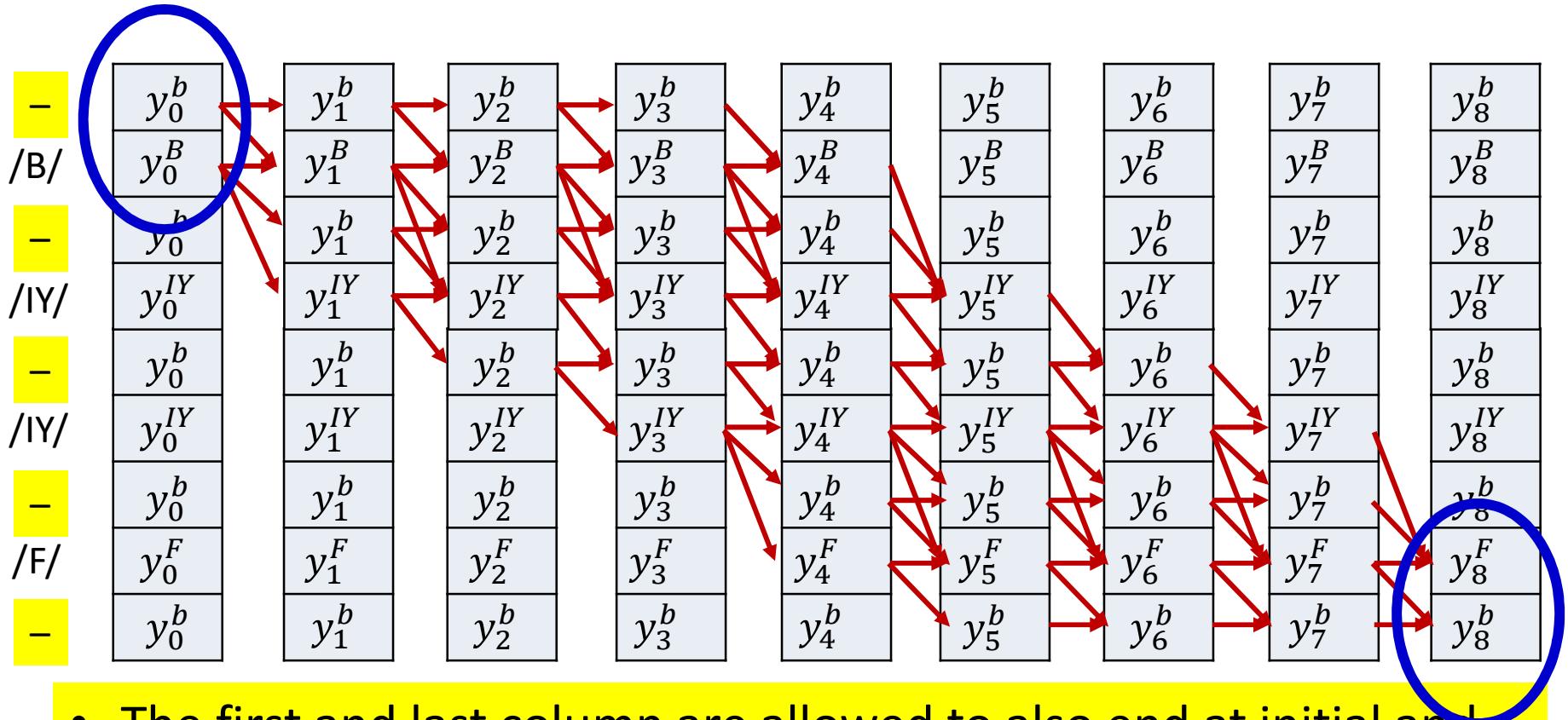
- With blanks
- Note: a row of blanks between any two symbols
- Also blanks at the very beginning and the very end

# Composing the graph for training



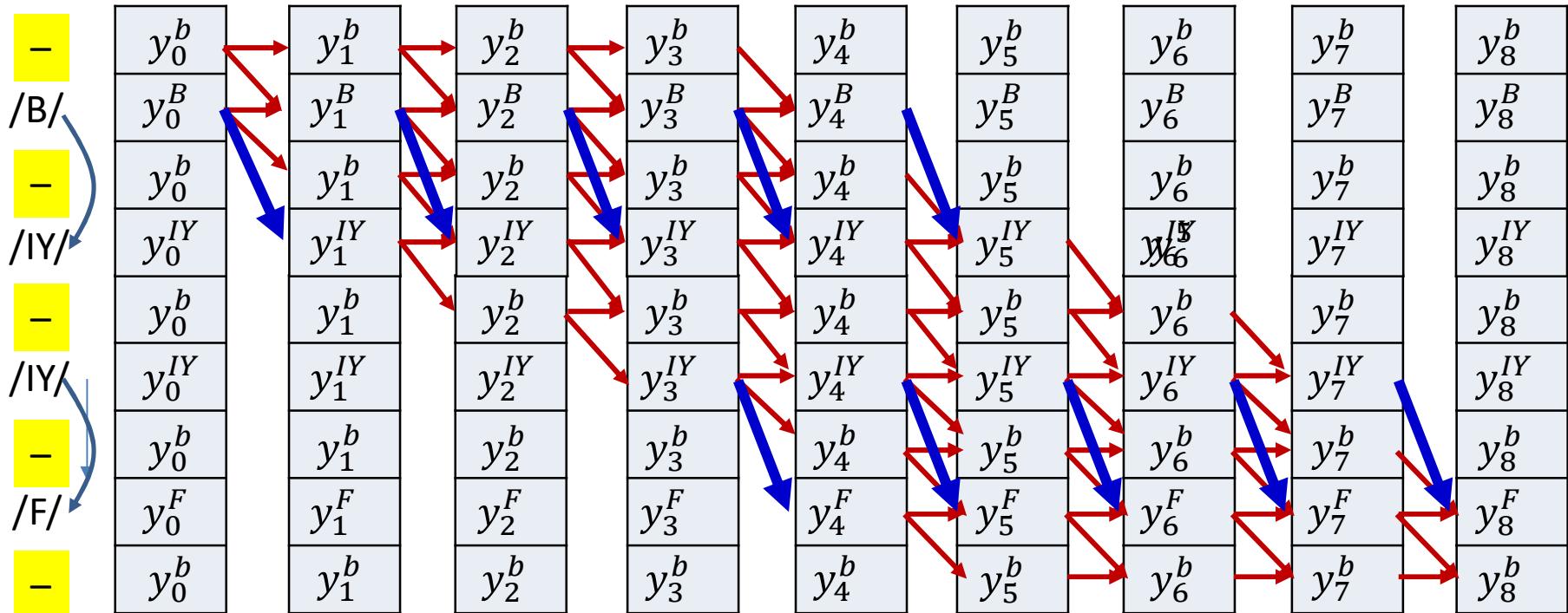
- Add edges such that all paths from initial node(s) to final node(s) unambiguously represent the target symbol sequence

# Composing the graph for training



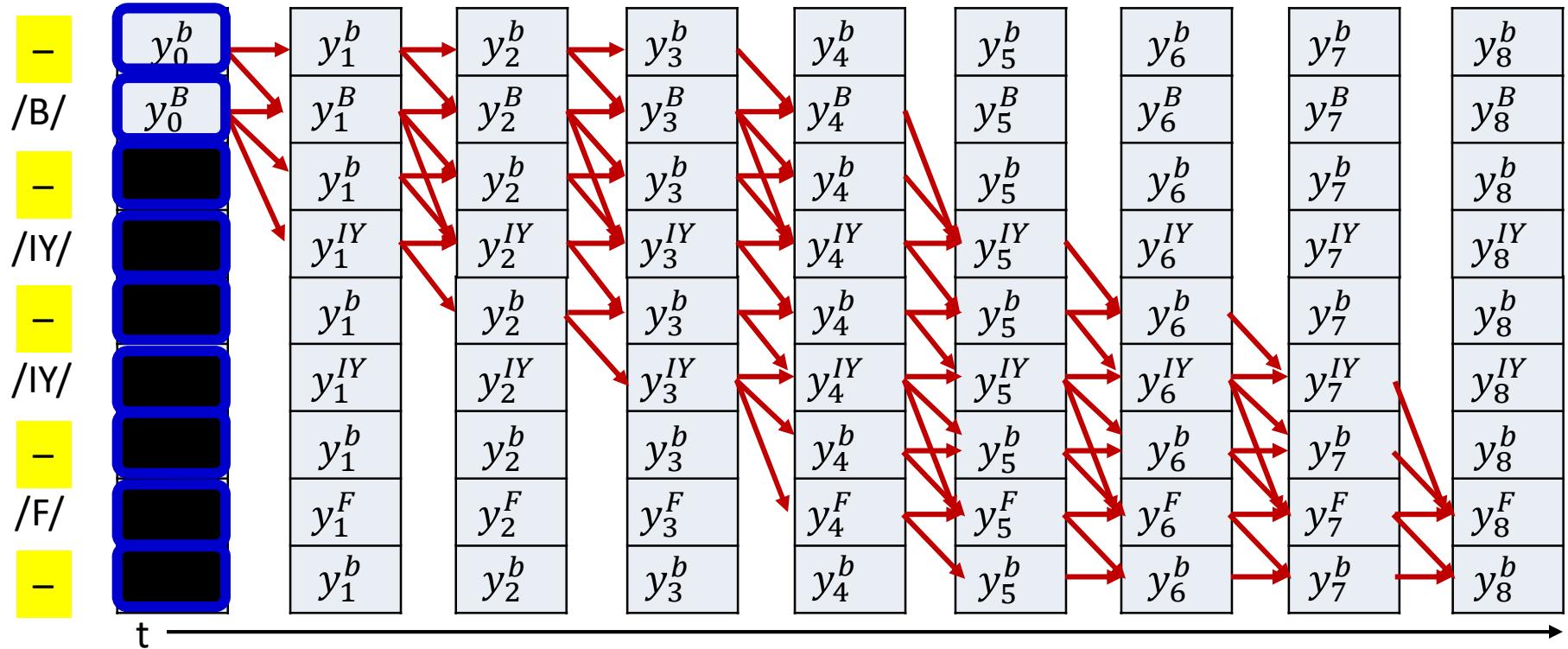
- The first and last column are allowed to also end at initial and final blanks

# Composing the graph for training



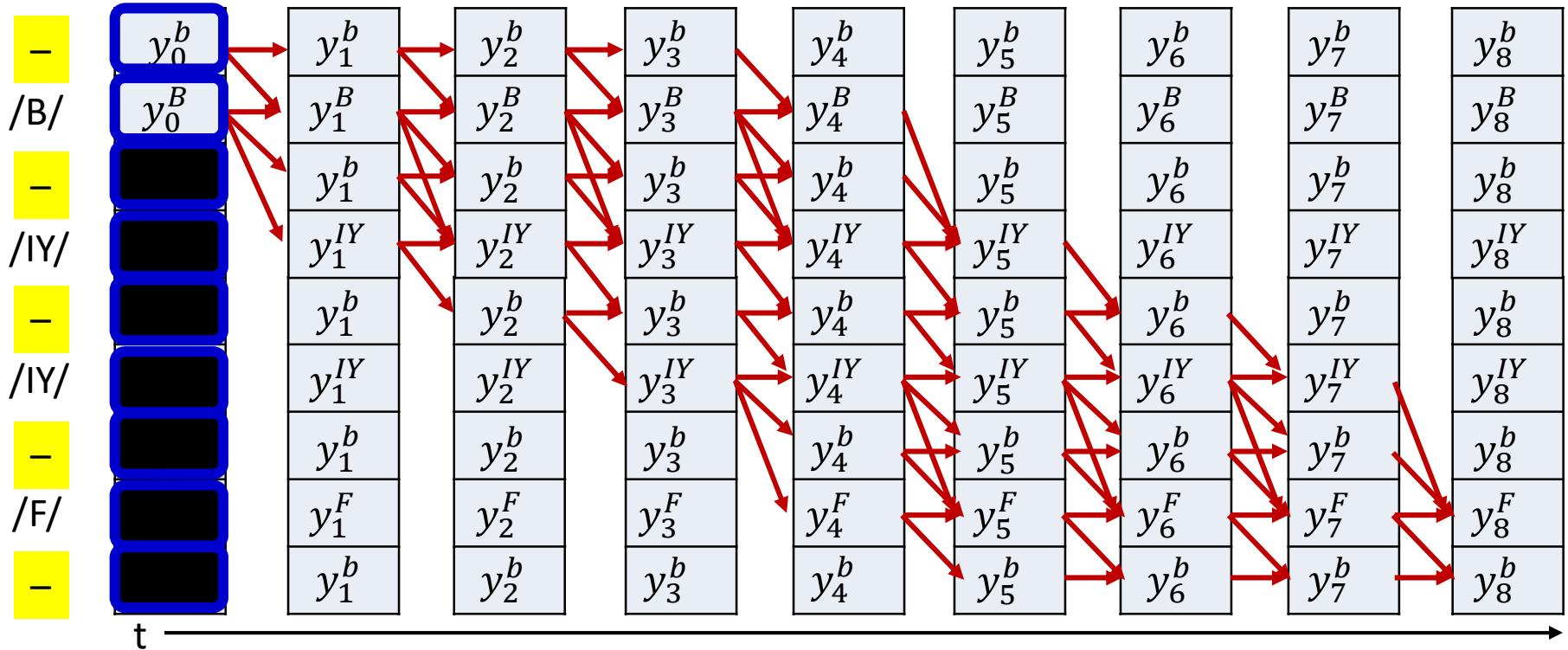
- The first and last column are allowed to also end at initial and final blanks
- Skips are permitted across a blank, but only if the symbols on either side are different
  - Because a blank is *mandatory between repetitions of a symbol* but *not required between distinct symbols*

# Modified Forward Algorithm



- Initialization:
  - $\alpha(0,0) = y_0^b, \alpha(0,1) = y_0^b, \alpha(0,r) = 0 \quad r > 1$

# Modified Forward Algorithm



- Iteration:

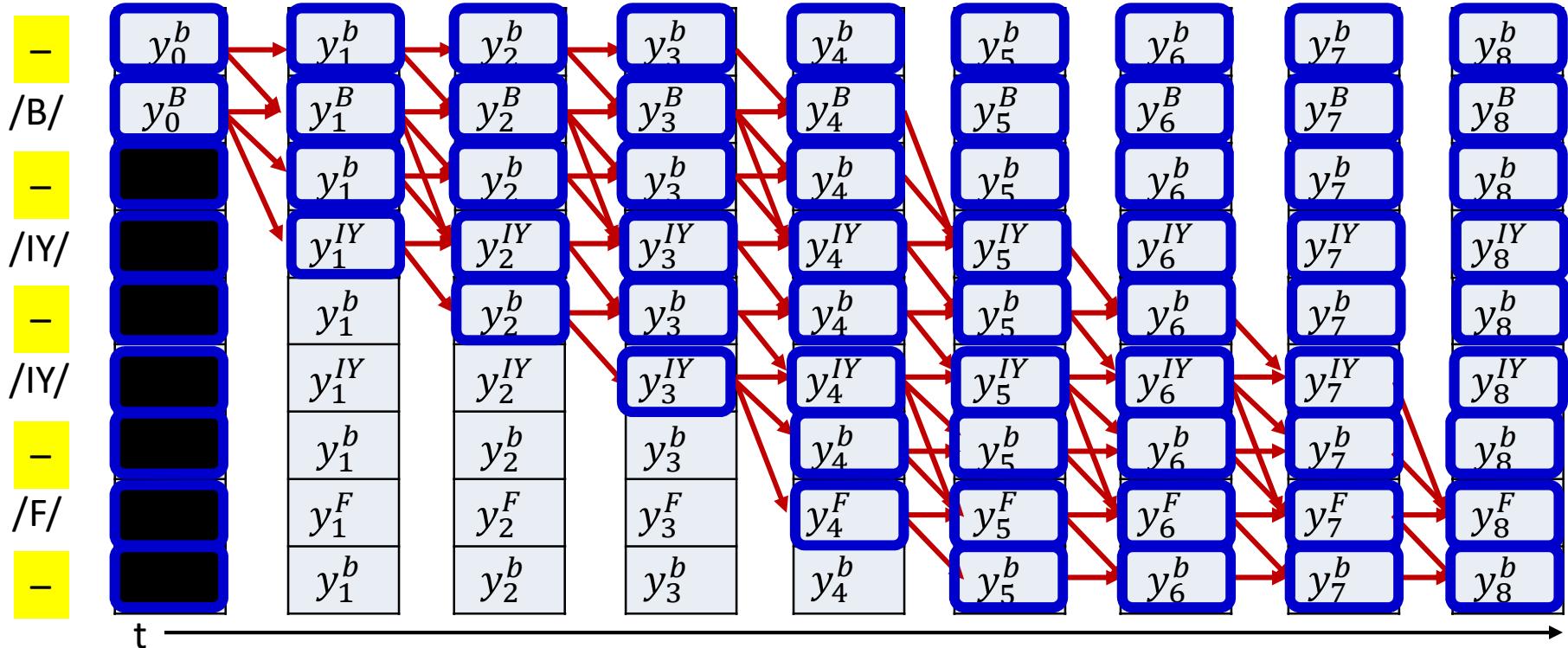
$$\alpha(t, r) = (\alpha(t-1, r) + \alpha(t-1, r-1))y_t^{S(r)}$$

- If  $S(r) = " - "$  or  $S(r) = S(r-2)$

$$\alpha(t, r) = (\alpha(t-1, r) + \alpha(t-1, r-1) + \alpha(t-1, r-2))y_t^{S(r)}$$

- Otherwise

# Modified Forward Algorithm



- Iteration:

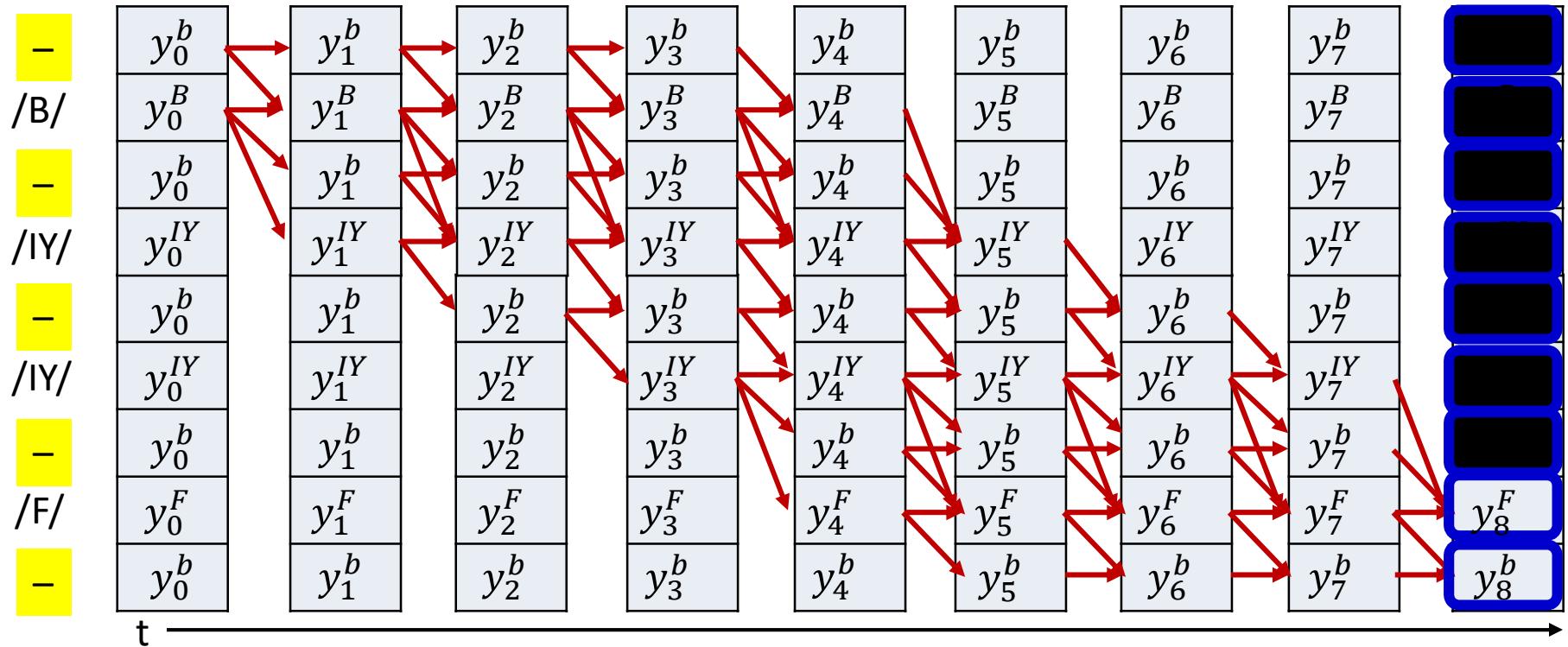
$$\alpha(t, r) = (\alpha(t-1, r) + \alpha(t-1, r-1))y_t^{S(r)}$$

- If  $S(r) = " - "$  or  $S(r) = S(r-2)$

$$\alpha(t, r) = (\alpha(t-1, r) + \alpha(t-1, r-1) + \alpha(t-1, r-2))y_t^{S(r)}$$

- Otherwise

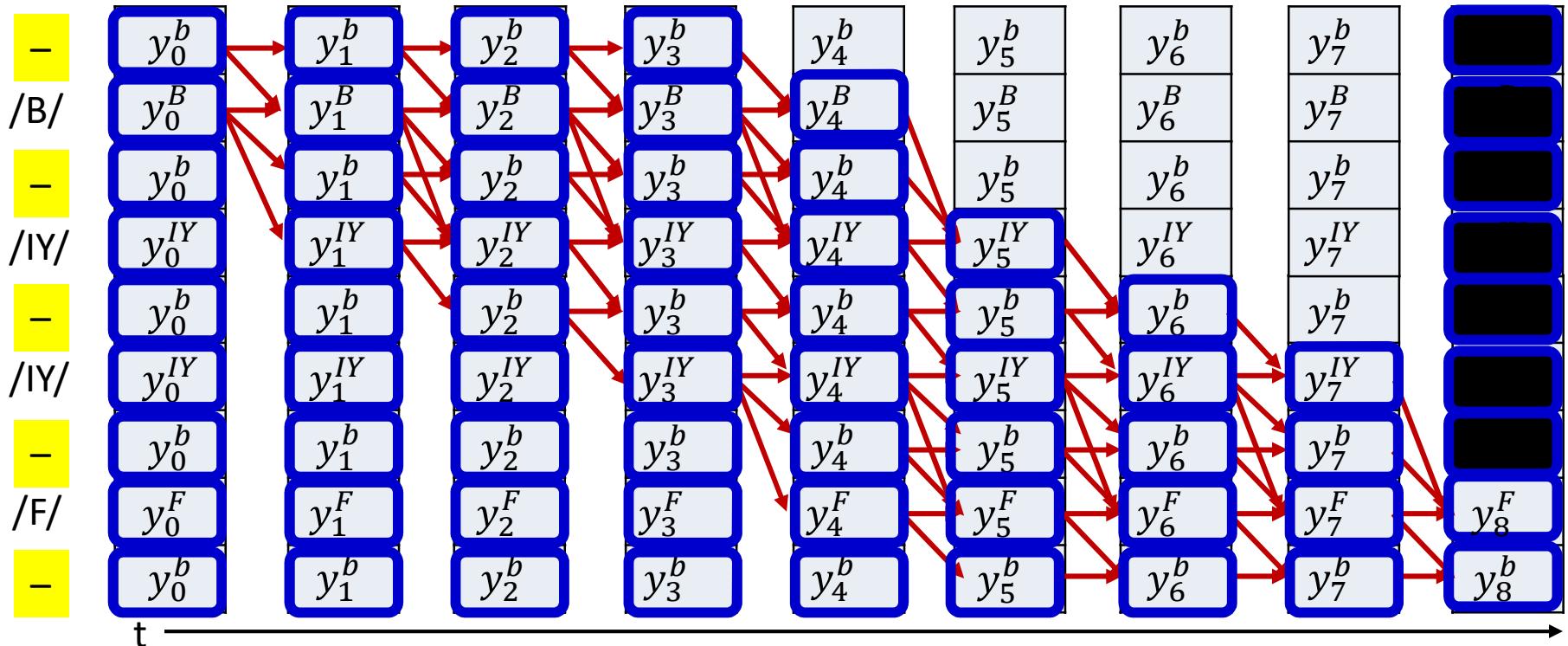
# Modified Backward Algorithm



- Initialization:

$$\begin{aligned}\beta(T-1, 2K) &= \beta(T-1, 2K-1) = 1 \\ \beta(T-1, r) &= 0 \quad r < 2K-1\end{aligned}$$

# Modified Backward Algorithm



- Iteration:

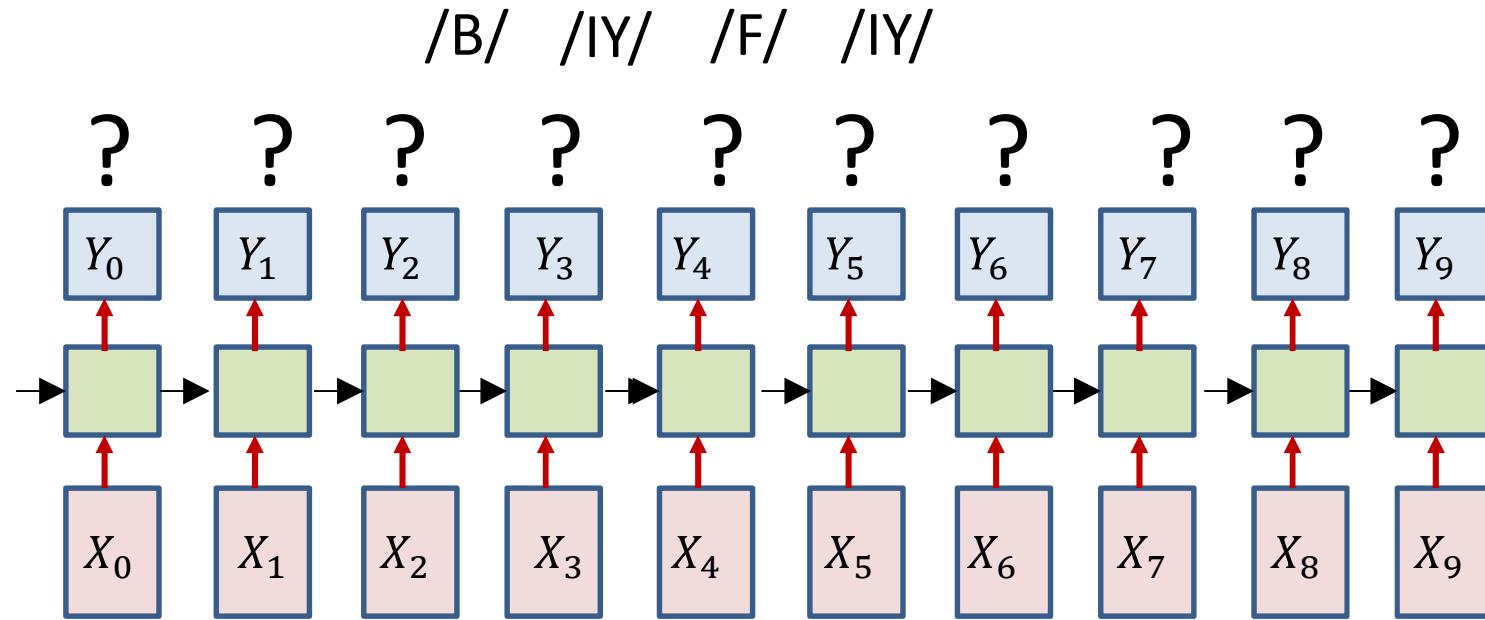
$$\beta(t, r) = \beta(t + 1, r)y_t^{S(r)} + \beta(t + 1, r + 1)y_t^{S(r+1)}$$

- If  $S(r) = " - "$  or  $S(r) = S(r + 2)$

$$\beta(t, r) = \beta(t + 1, r)y_t^{S(r)} + \beta(t + 1, r + 1)y_t^{S(r+1)} + \beta(t + 1, r + 2)y_t^{S(r+2)}$$

- Otherwise

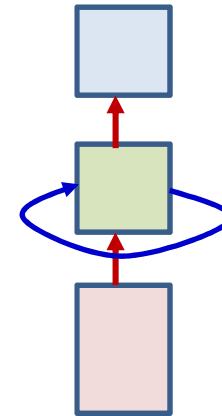
# Overall training procedure for Seq2Seq with blanks



- Problem: Given input and output sequences without alignment, train models

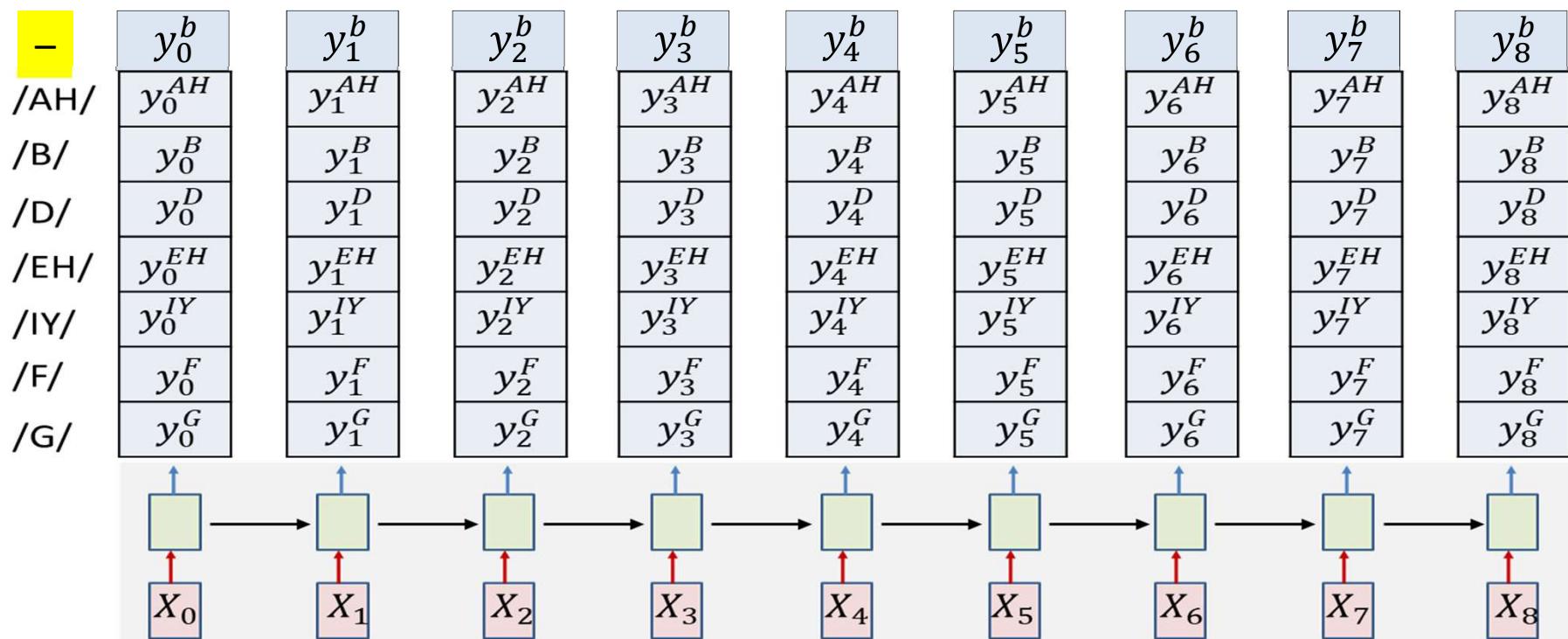
# Overall training procedure

- **Step 1:** Setup the network
  - Typically many-layered LSTM
- **Step 2:** Initialize all parameters of the network
  - Include a “blank” symbol in vocabulary

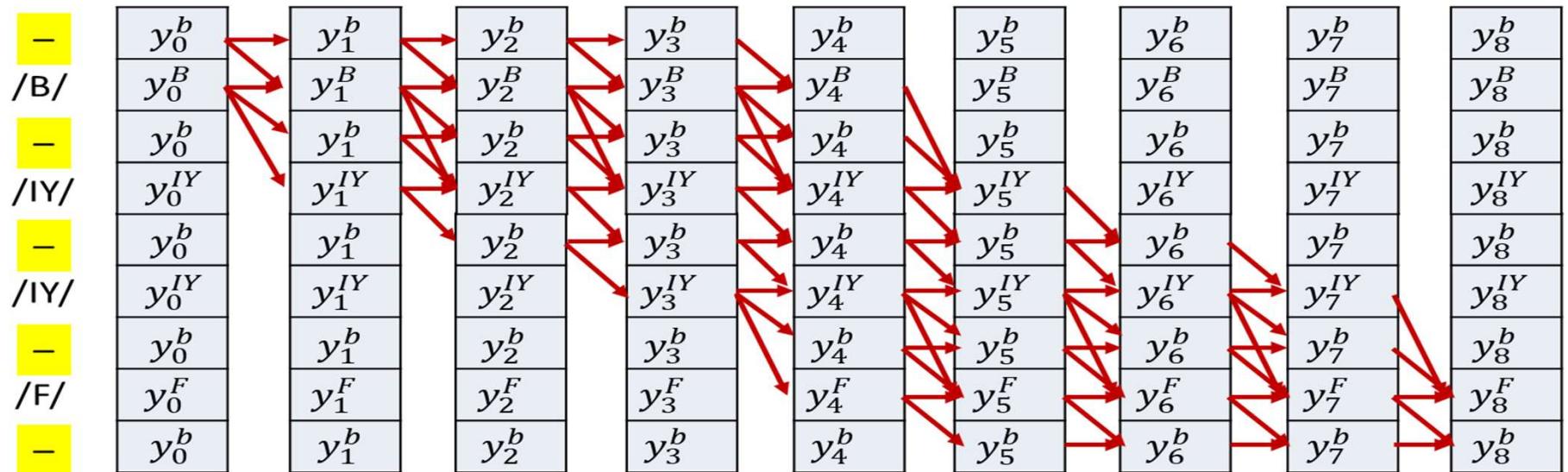


# Overall Training: Forward pass

- Foreach training instance
  - **Step 3:** Forward pass. Pass the training instance through the network and obtain all symbol probabilities at each time, including blanks

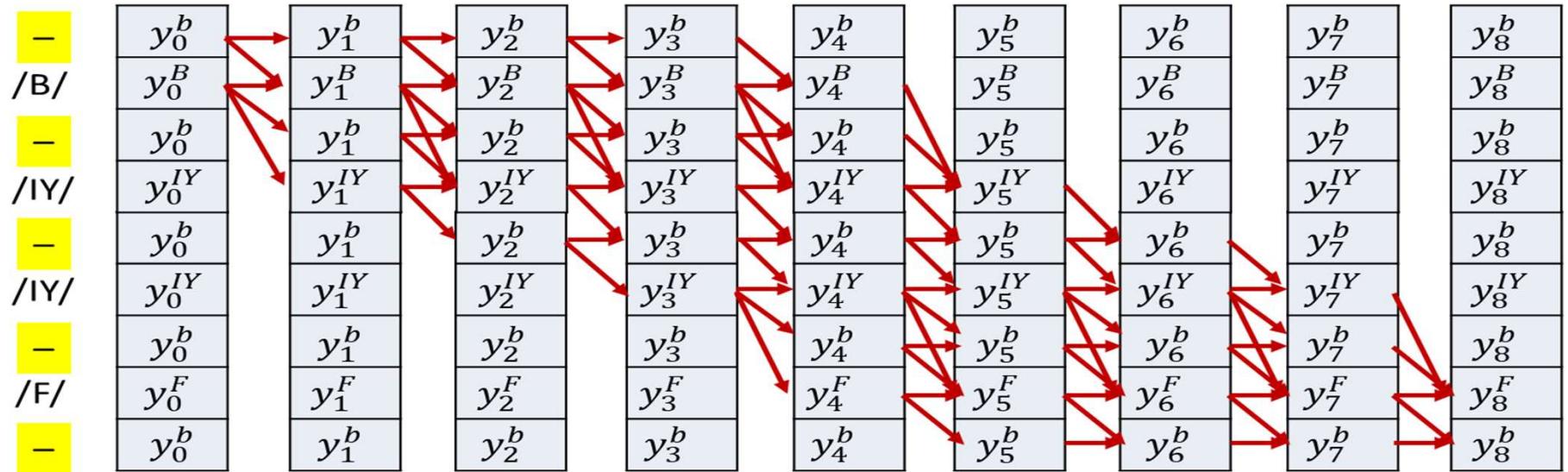


# Overall training: Backward pass



- Foreach training instance
  - **Step 3:** Forward pass. Pass the training instance through the network and obtain all symbol probabilities at each time
  - **Step 4:** Construct the graph representing the specific symbol sequence in the instance. Use appropriate connections if blanks are included

# Overall training: Backward pass



- Foreach training instance:
  - **Step 5:** Perform the forward backward algorithm to compute  $\alpha(t, r)$  and  $\beta(t, r)$  at each time, for each row of nodes in the graph using the modified forward-backward equations
  - **Step 6:** Compute derivative of divergence  $\nabla_{Y_t} DIV$  for each  $Y_t$

# Overall training: Backward pass

- Foreach instance
  - **Step 6:** Compute derivative of divergence  $\nabla_{Y_t} DIV$  for each  $Y_t$

$$\nabla_{Y_t} DIV = \begin{bmatrix} \frac{dDIV}{dy_t^1} & \frac{dDIV}{dy_t^2} & \dots & \frac{dDIV}{dy_t^L} \end{bmatrix}$$

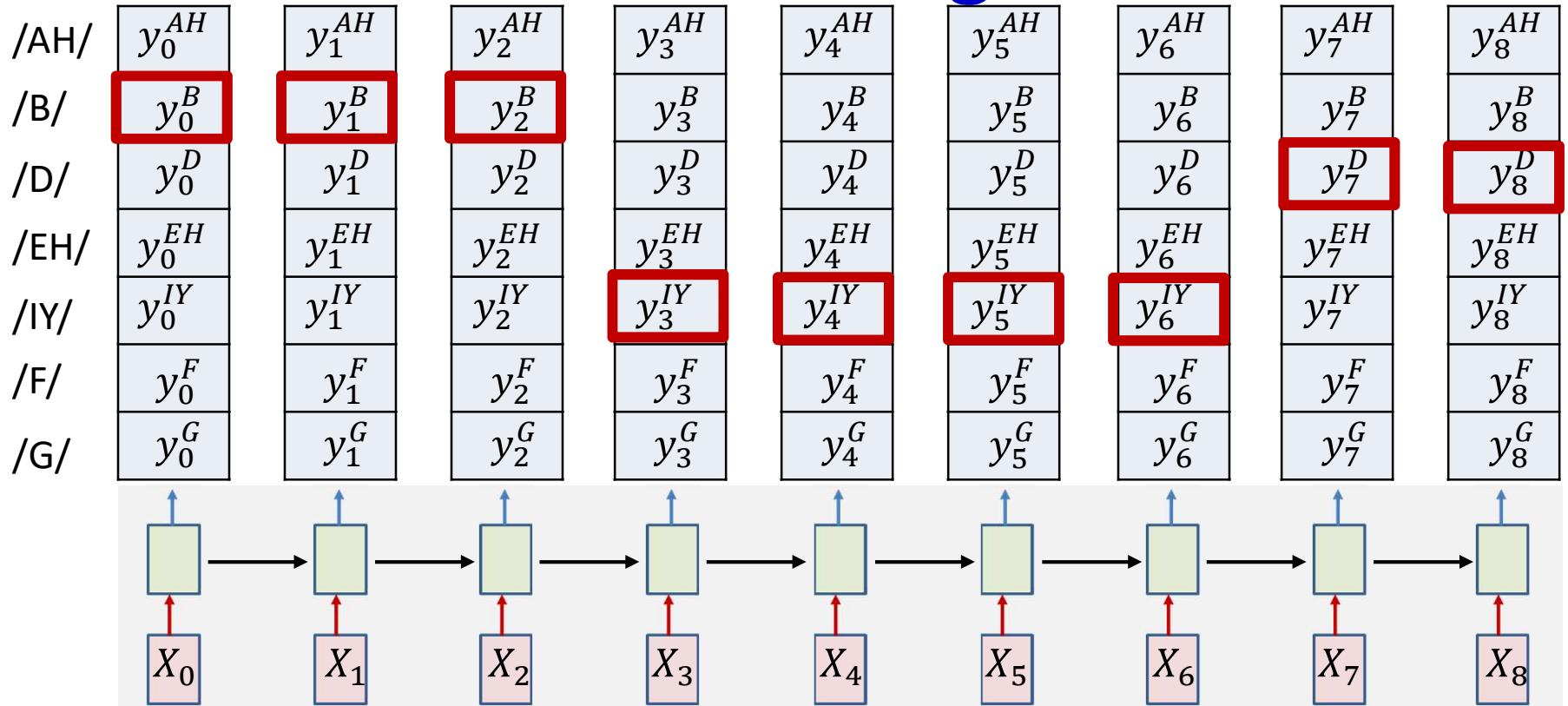
$$\frac{dDIV}{dy_t^l} = - \sum_{r : S(r)=l} \frac{d}{dy_t^{S(r)}} (\gamma(t, r) \log y_t^{S(r)})$$

- **Step 7:** Aggregate derivatives over minibatch and update parameters

# CTC: Connectionist Temporal Classification

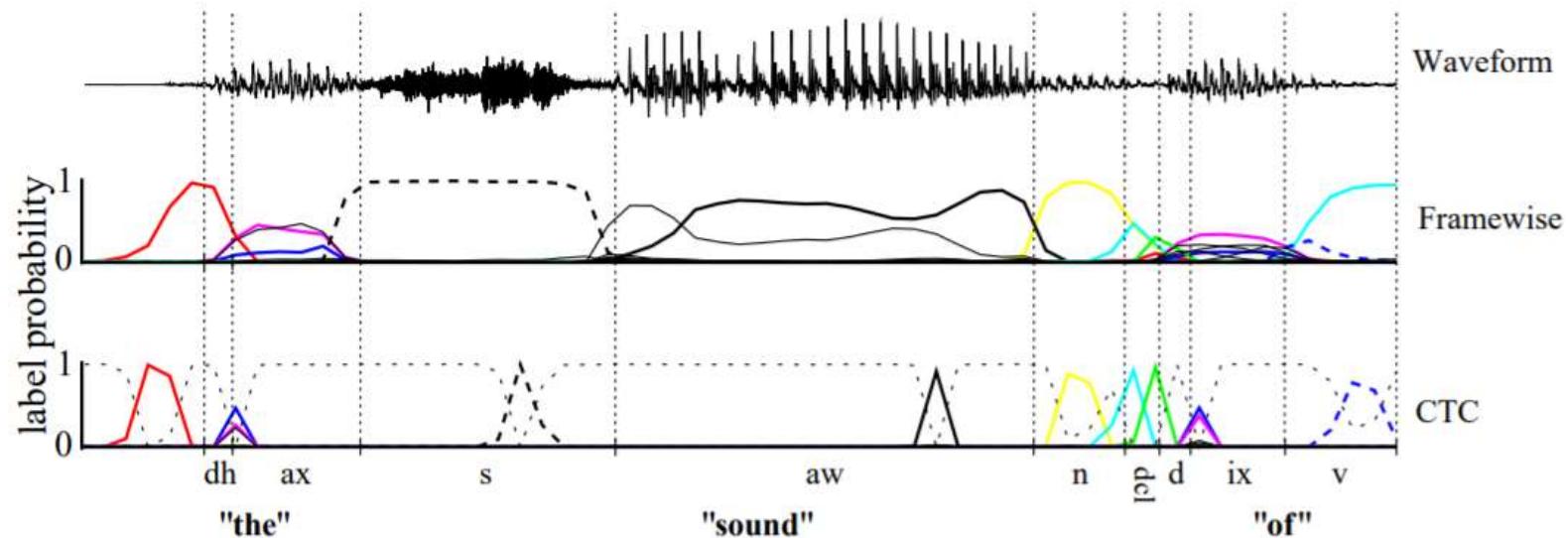
- The overall framework we saw is referred to as CTC
  - Applies when “duplicating” labels at the output is considered acceptable, and when output sequence length < input sequence length

# Returning to an old problem: Decoding



- Using a trained network, we can “decode” the symbol sequence for an input by tracing the most likely symbol at each frame and merging
- This is in fact a *suboptimal* decode that actually finds the most likely *synchronous* output sequence
  - Which is not necessarily the most likely *asynchronous* sequence
- Nevertheless it is effective with CTC models that have blanks**

# What a CTC system outputs



**Figure 1. Framewise and CTC networks classifying a speech signal.** The shaded lines are the output activations, corresponding to the probabilities of observing phonemes at particular times. The CTC network predicts only the sequence of phonemes (typically as a series of spikes, separated by ‘blanks’, or null predictions), while the framewise network attempts to align them with the manual segmentation (vertical lines). The framewise network receives an error for misaligning the segment boundaries, even if it predicts the correct phoneme (e.g. ‘dh’). When one phoneme always occurs beside another (e.g. the closure ‘dcl’ with the stop ‘d’), CTC tends to predict them together in a double spike. The choice of labelling can be read directly from the CTC outputs (follow the spikes), whereas the predictions of the framewise network must be post-processed before use.

- Ref: Graves
- Symbol outputs peak at the ends of the sounds
  - But are smeared..
- Better ability to find most likely symbol sequence, and can handle repetitions
- But this is still suboptimal..

# Actual objective of decoding

- Want to find most likely asynchronous symbol sequence
  - /R/ /EH/ /D/
- What Viterbi finds: most likely synchronous symbol sequence
  - /R/ /R/ /R//EH//EH//EH//D/
  - Which must be compressed
- The likelihood of an asynchronous symbol sequence  $\mathbf{S} = S_1, \dots, S_K$ , given an input  $\mathbf{X} = X_1, \dots, X_T$ , is given by the forward algorithm

$$P(\mathbf{S}|\mathbf{X}) = P(\mathbf{S}, s_T = S_K | \mathbf{X}) = \alpha_{\mathbf{S}}(S_K, T)$$

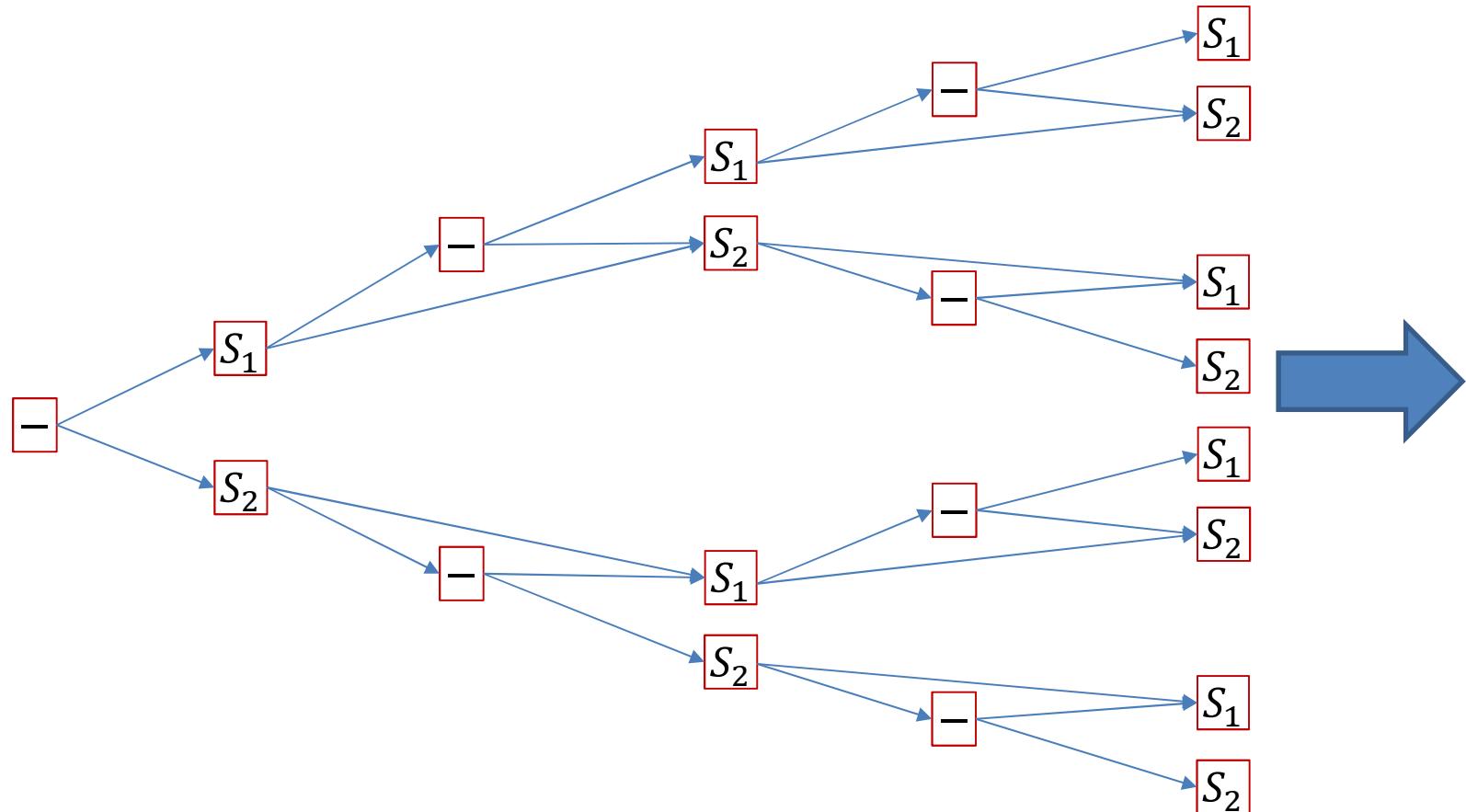
# Actual decoding objective

- Find the most likely (asynchronous) symbol sequence

$$\hat{S} = \operatorname*{argmax}_S \alpha_S(S_K, T)$$

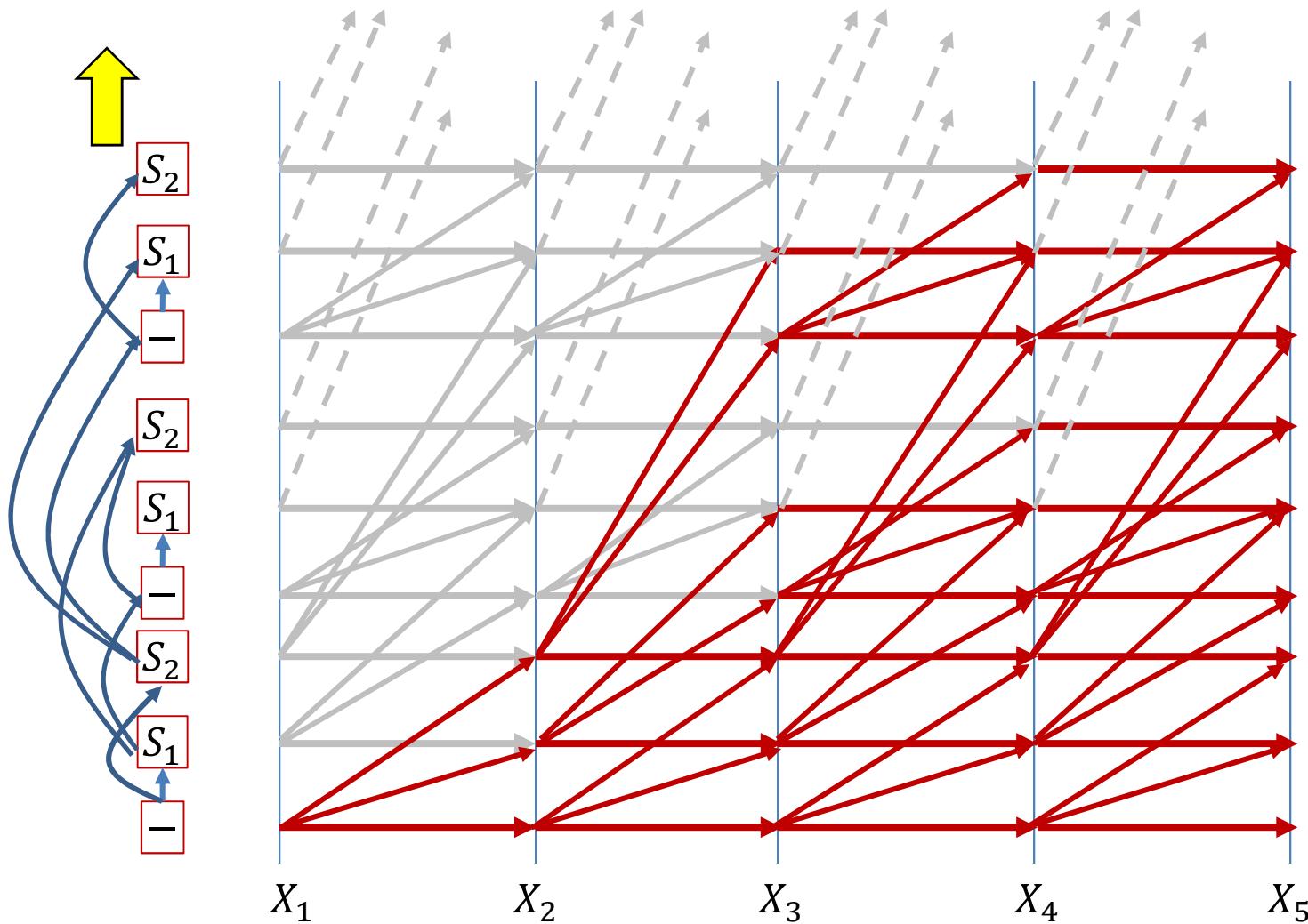
- Unfortunately, explicit computation of this will require evaluate of an exponential number of symbol sequences
- Solution: Organize all possible symbol sequences as a (semi)tree

# Hypothesis semi-tree



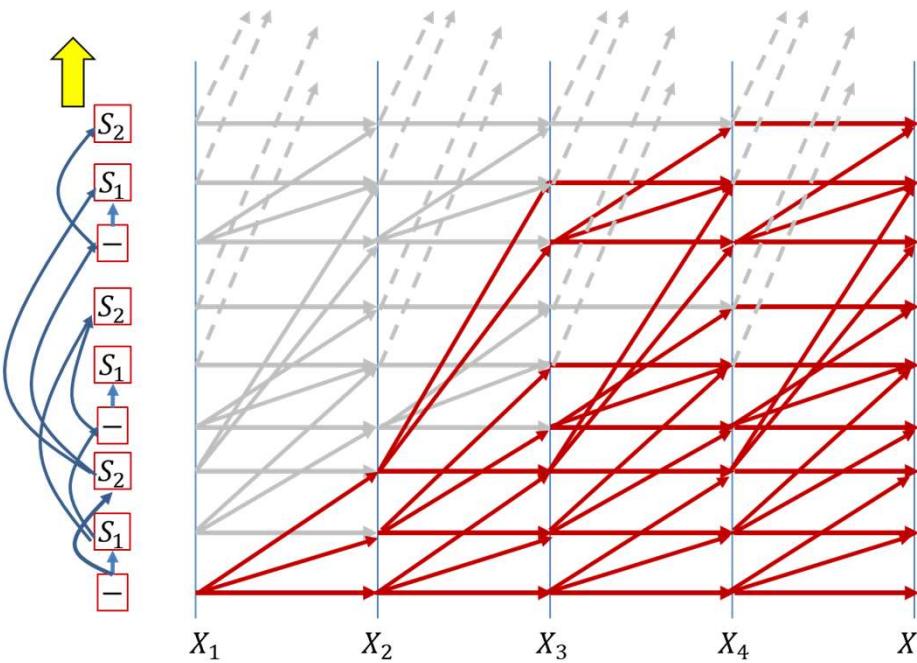
- The semi tree of hypotheses (assuming only 3 symbols in the vocabulary)
- Every symbol connects to every symbol other than itself
  - It also connects to a blank, which connects to every symbol including itself
- The simple structure repeats recursively
- Each node represents a unique (partial) symbol sequence!

# The decoding graph for the tree



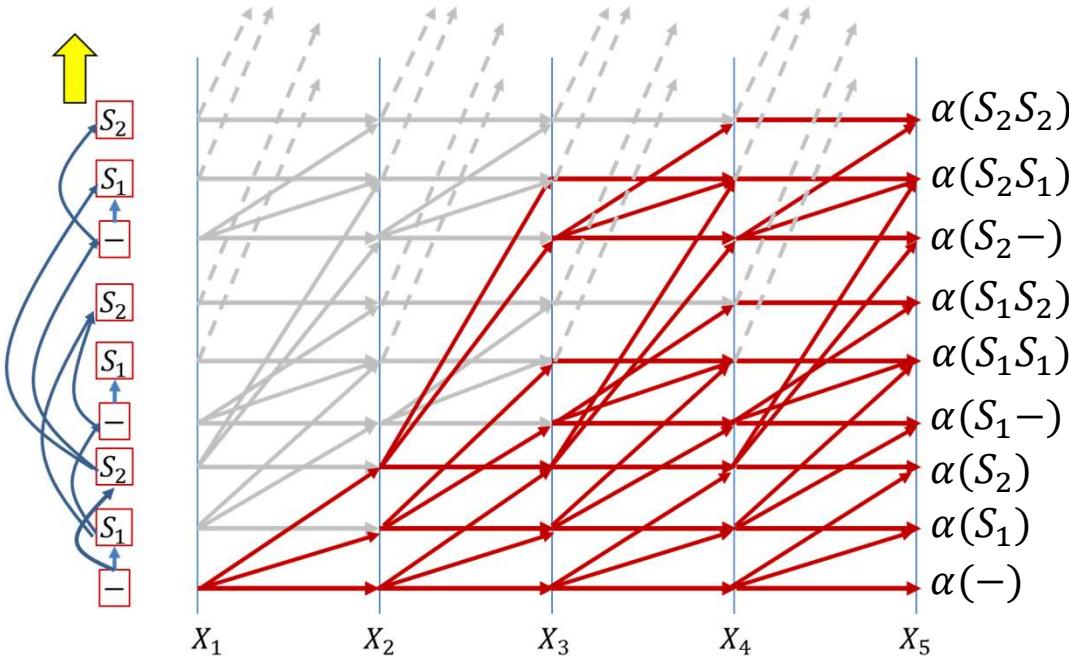
- Graph with more than 2 symbols will be similar but much more cluttered and complicated

# The decoding graph for the tree



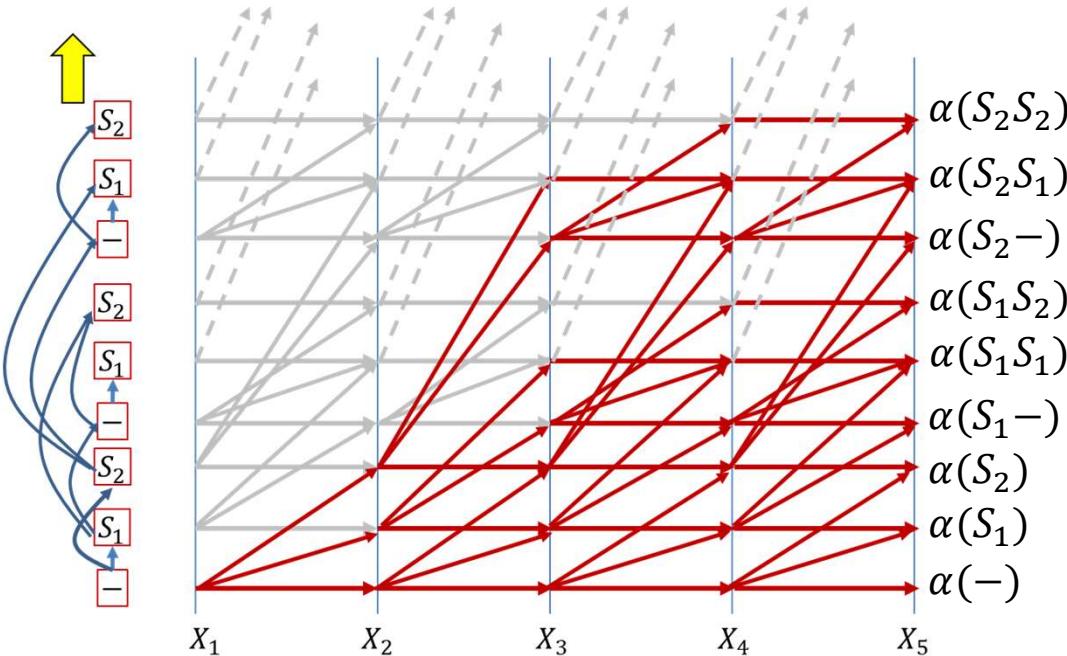
- The figure to the left is the tree, drawn in a vertical line
- The graph is just the tree unrolled over time
  - For a vocabulary of  $V$  symbols, every node connects out to  $V$  other nodes at the next time
- Every node in the graph represents a unique symbol sequence

# The decoding graph for the tree



- The forward score  $\alpha(r, T)$  at the final time represents the full forward score for a unique symbol sequence (including sequences terminating in blanks)
- Select the symbol sequence with the largest alpha
  - Some sequences may have two alphas, one for the sequence itself, one for the sequence followed by a blank
  - Add the alphas before selecting the most likely

# CTC decoding



- This is the “theoretically correct” CTC decoder
- In practice, the graph gets exponentially large very quickly
- To prevent this pruning strategies are employed to keep the graph (and computation) manageable
  - This may cause suboptimal decodes, however
  - The fact that CTC scores peak at symbol terminations minimizes the damage due to pruning

# Story so far: CTC models

- Sequence-to-sequence networks which irregularly produce output symbols can be trained by
  - Iteratively aligning the target output to the input and time-synchronous training
  - Optimizing the expected error over *all* possible alignments: CTC training
- Distinct repetition of symbols can be disambiguated from repetitions representing the extended output of a single symbol by the introduction of blanks
- Decoding the models can be performed by
  - Best-path decoding, i.e. Viterbi decoding
  - Optimal CTC decoding based on the application of the forward algorithm to a tree-structured representation of all possible output strings

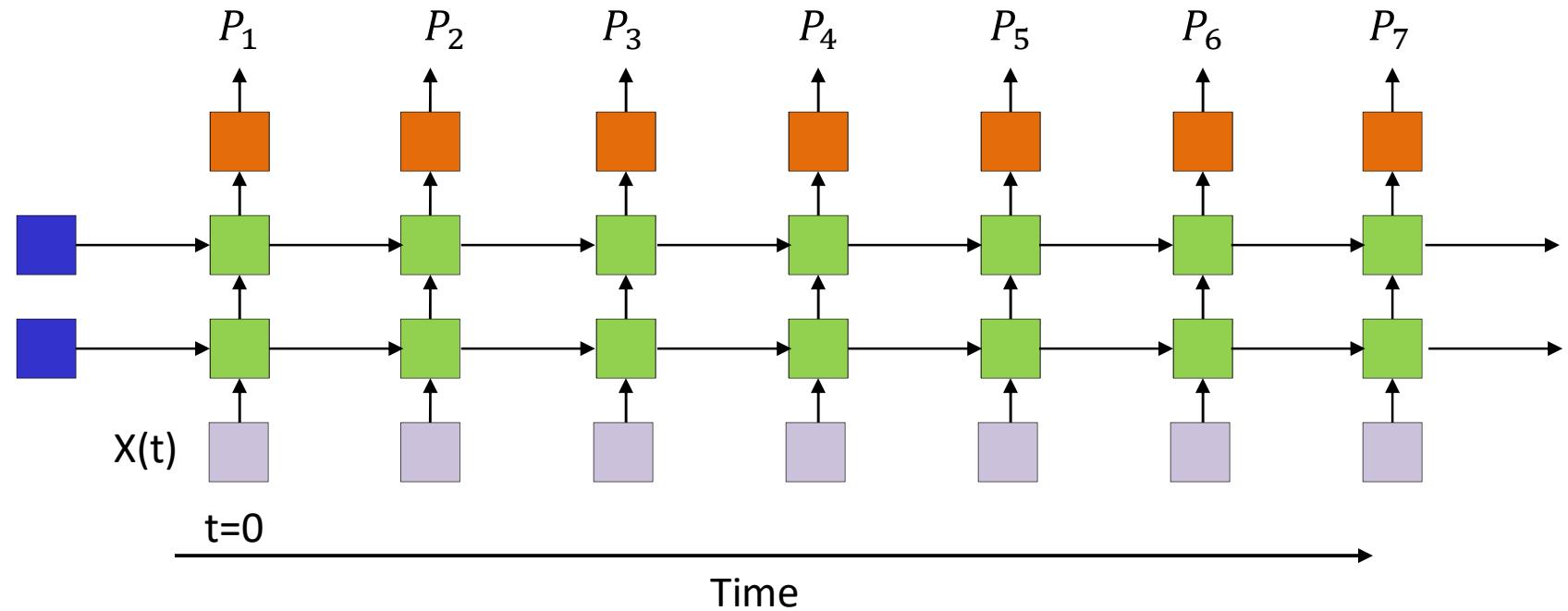
# CTC caveats

- The “blank” structure (with concurrent modifications to the forward-backward equations) is only one way to deal with the problem of repeating symbols
- Possible variants:
  - Symbols partitioned into two or more sequential subunits
    - No blanks are required, since subunits must be visited in order
  - Symbol-specific blanks
    - Doubles the “vocabulary”
  - CTC can use *bidirectional* recurrent nets
    - And frequently does
  - Other variants possible..

# Most common CTC applications

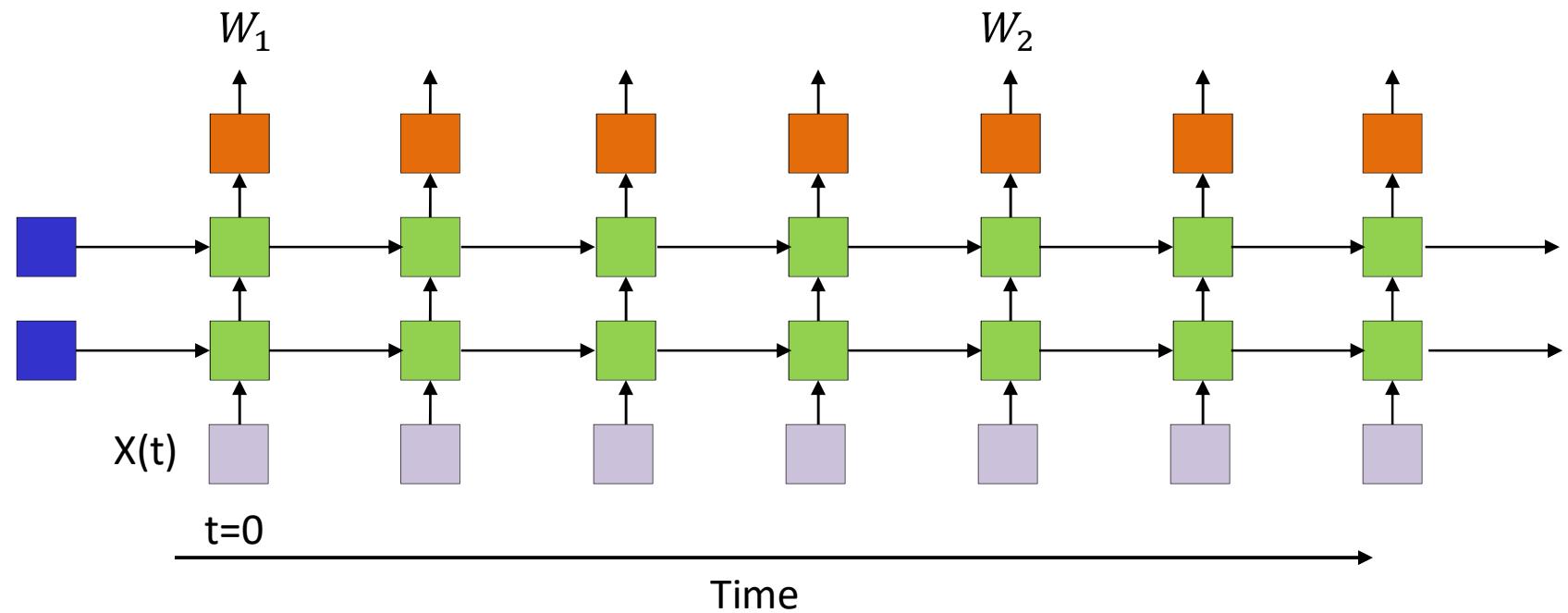
- Speech recognition
  - Speech in, phoneme sequence out
  - Speech in, character sequence (spelling out)
- Handwriting recognition

# Speech recognition using Recurrent Nets



- Recurrent neural networks (with LSTMs) can be used to perform speech recognition
  - Input: Sequences of audio feature vectors
  - Output: Phonetic label of each vector

# Speech recognition using Recurrent Nets



- Alternative: Directly output phoneme, character or word sequence

# Brief detour: Language models

- Modelling language using time-synchronous nets
- More generally language models and embeddings..

# Which open source project?

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECON
    return segtable;
}
```

# Language modelling using RNNs

Four score and seven years ???

ABRAHAM LINCOL??

- Problem: Given a sequence of words (or characters) predict the next one

# Language modelling: Representing words

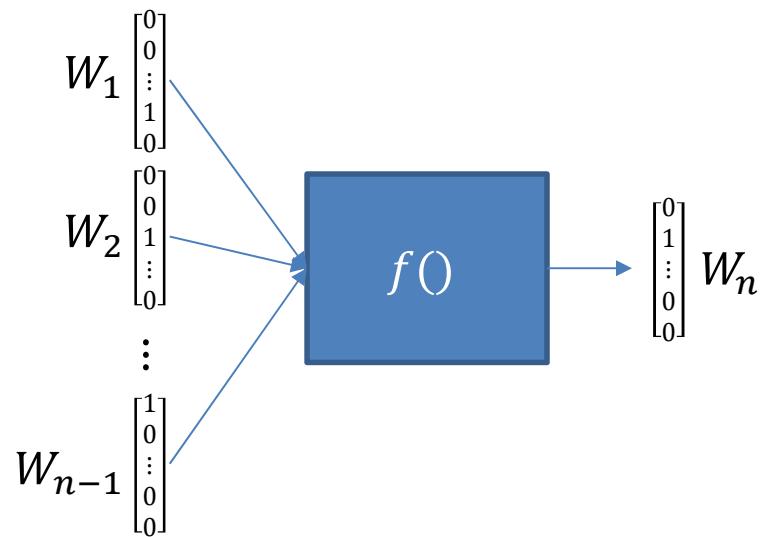
- Represent words as one-hot vectors
  - Pre-specify a vocabulary of N words in fixed (e.g. lexical) order
    - *E.g. [A AARDVARK AARON ABACK ABACUS... ZZYP]*
  - Represent each word by an N-dimensional vector with N-1 zeros and a single 1 (in the position of the word in the ordered list of words)
    - E.g. “AARDVARK” → [0 1 0 0 0 ...]
    - E.g. “AARON” → [0 0 1 0 0 0 ...]
- Characters can be similarly represented
  - English will require about 100 characters, to include both cases, special characters such as commas, hyphens, apostrophes, etc., and the space character

# Predicting words

Four score and seven years ???

$$W_n = f(W_{-1}, W_1, \dots, W_{n-1})$$

Nx1 one-hot vectors



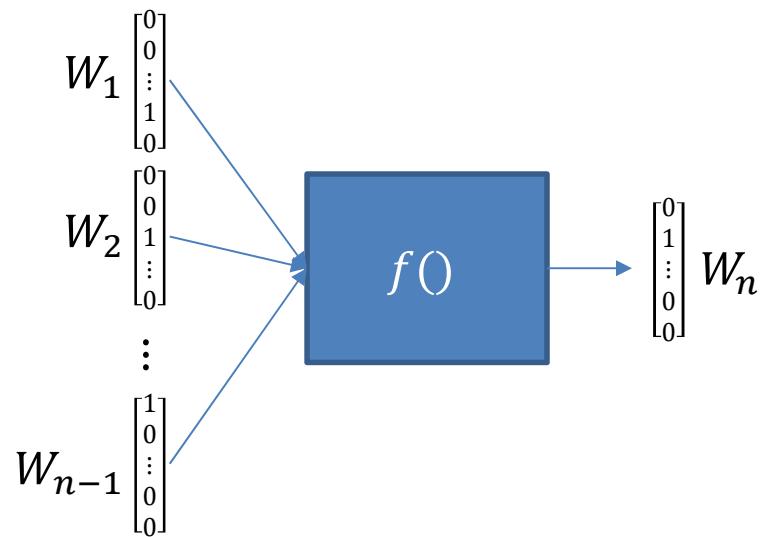
- Given one-hot representations of  $W_1 \dots W_{n-1}$ , predict  $W_n$

# Predicting words

Four score and seven years ???

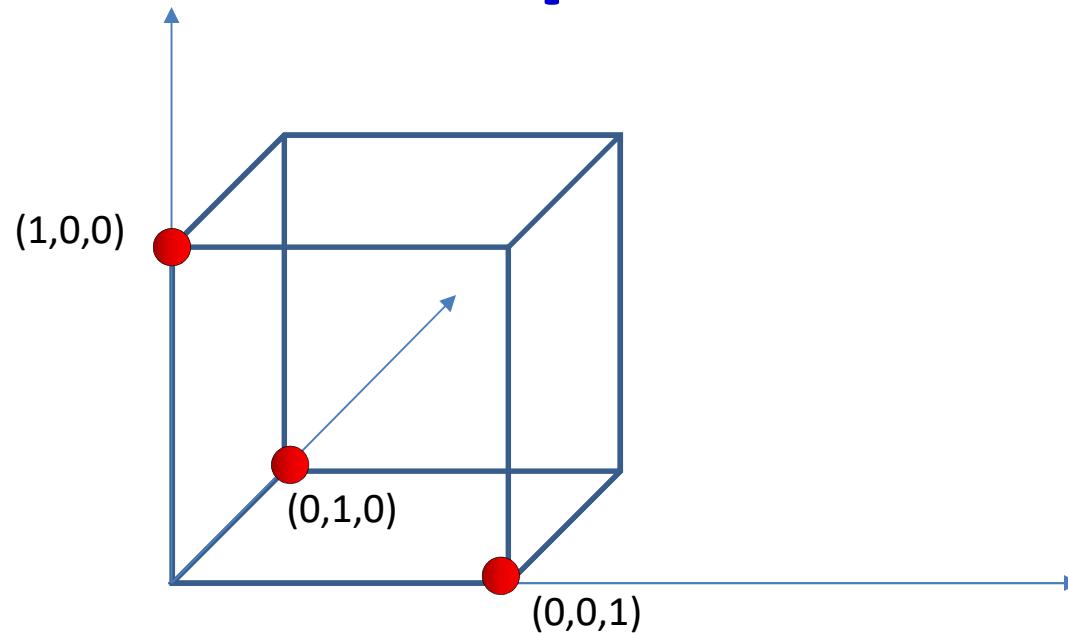
$$W_n = f(W_{-1}, W_1, \dots, W_{n-1})$$

Nx1 one-hot vectors



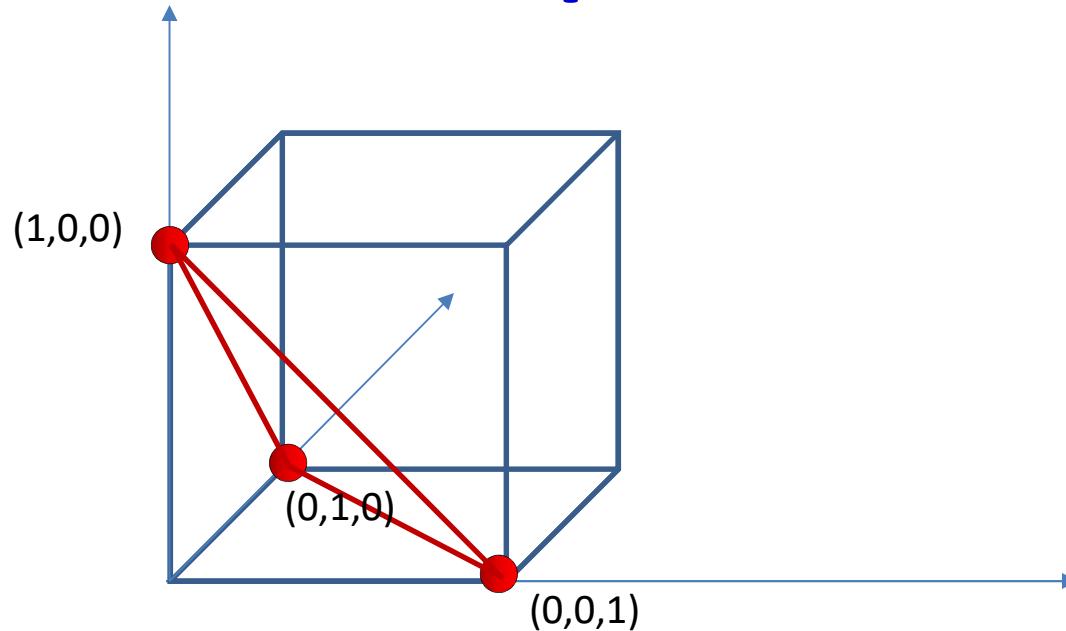
- Given one-hot representations of  $W_1 \dots W_{n-1}$ , predict  $W_n$
- Dimensionality problem:** All inputs  $W_1 \dots W_{n-1}$  are both very high-dimensional and very sparse

# The one-hot representation



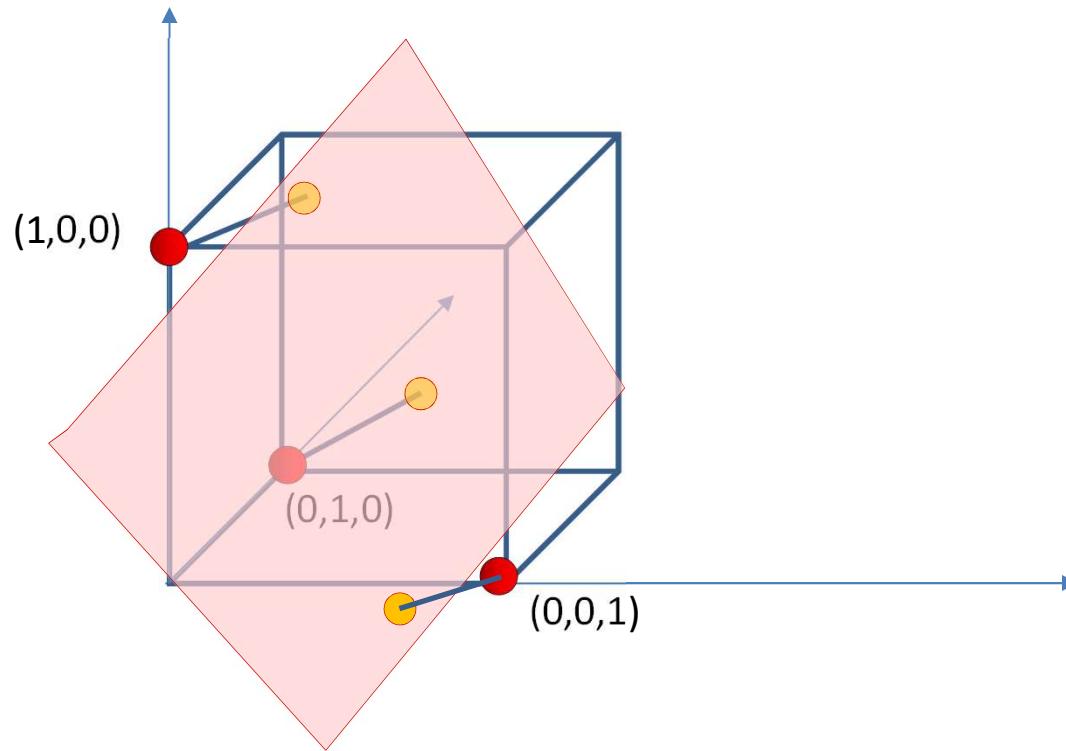
- The one hot representation uses only  $N$  corners of the  $2^N$  corners of a unit cube
  - Actual volume of space used = 0
    - $(1, \varepsilon, \delta)$  has no meaning except for  $\varepsilon = \delta = 0$
  - Density of points:  $\mathcal{O}\left(\frac{N}{r^N}\right)$
- This is a tremendously inefficient use of dimensions

# Why one-hot representation



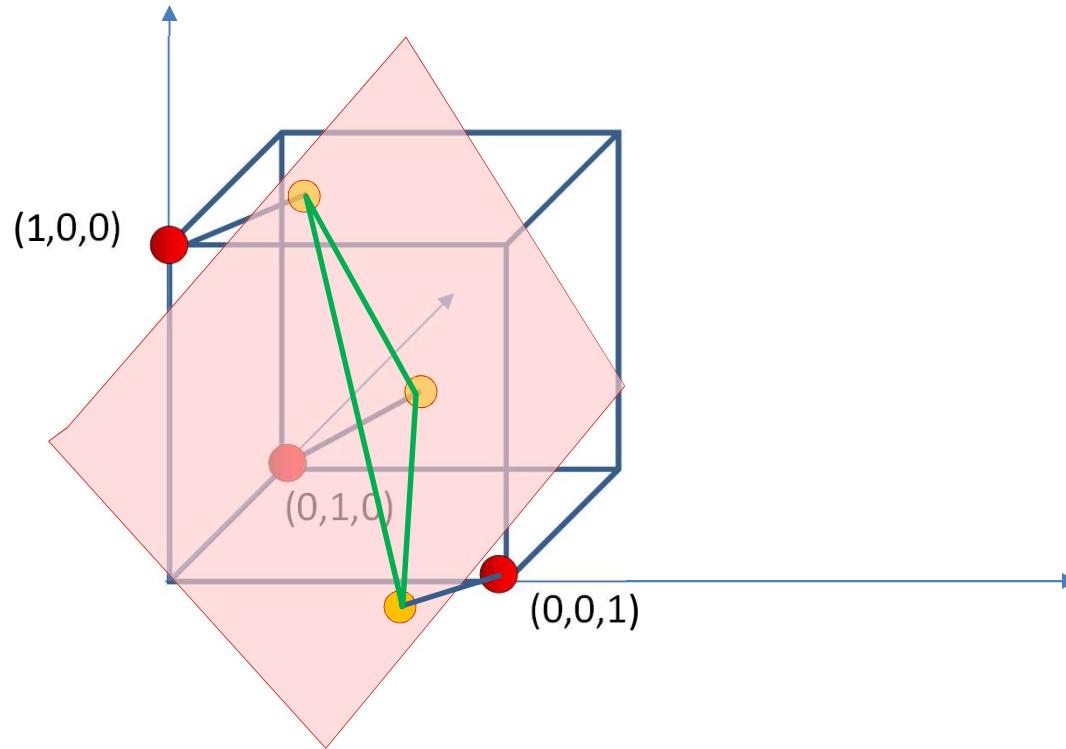
- The one-hot representation makes no assumptions about the relative importance of words
  - All word vectors are the same length
- It makes no assumptions about the relationships between words
  - The distance between every pair of words is the same

# Solution to dimensionality problem



- Project the points onto a lower-dimensional subspace
  - The volume used is still 0, but density can go up by many orders of magnitude
    - Density of points:  $\mathcal{O}\left(\frac{N}{r^M}\right)$

# Solution to dimensionality problem

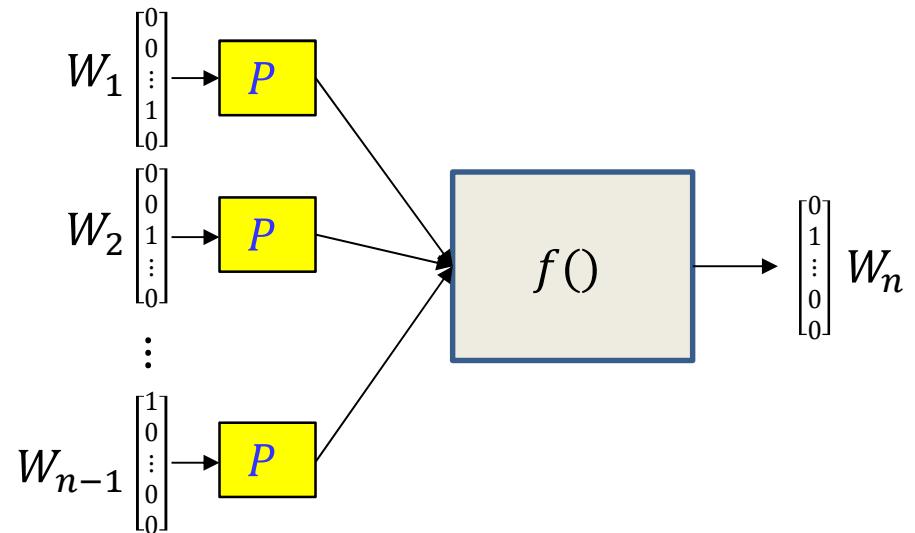
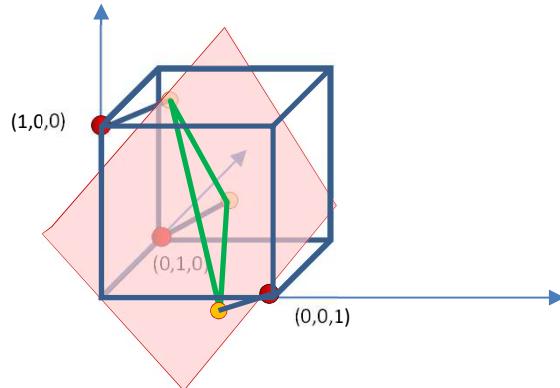


- Project the points onto a lower-dimensional subspace
  - The volume used is still 0, but density can go up by many orders of magnitude
    - Density of points:  $\mathcal{O}\left(\frac{N}{r^M}\right)$
  - If properly learned, the distances between projected points will capture semantic relations between the words
    - This will also require linear transformation (stretching/shrinking/rotation) of the subspace

# The *Projected* word vectors

Four score and seven years ???

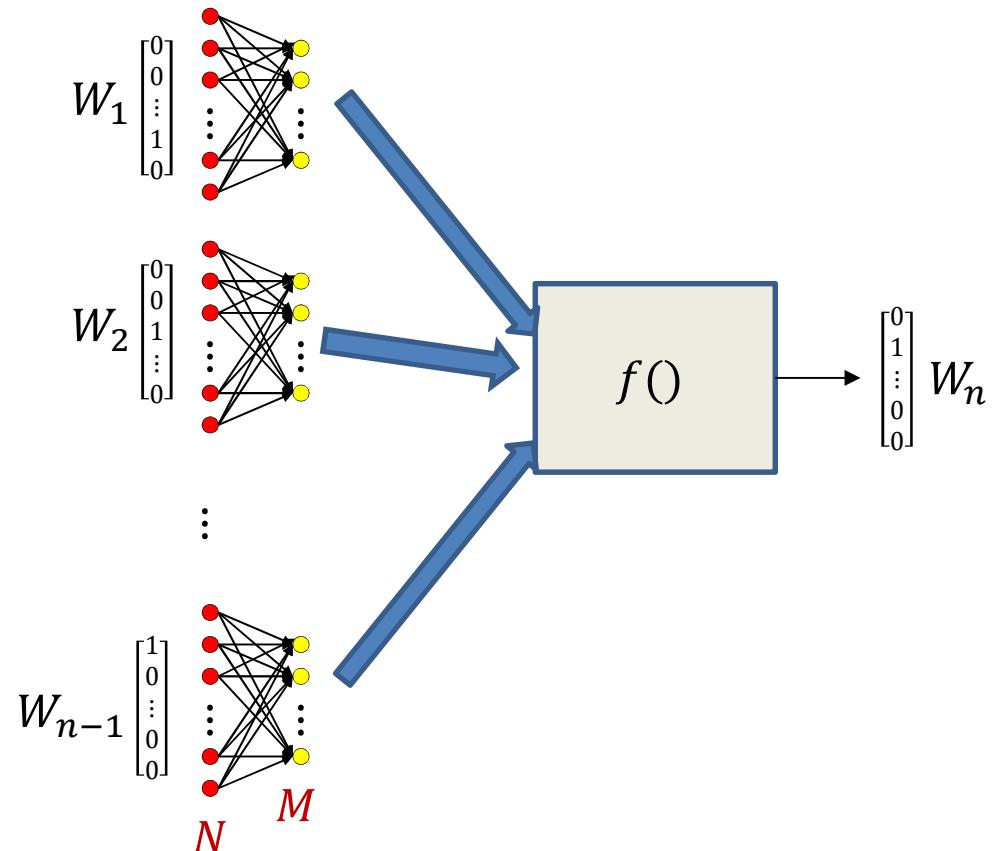
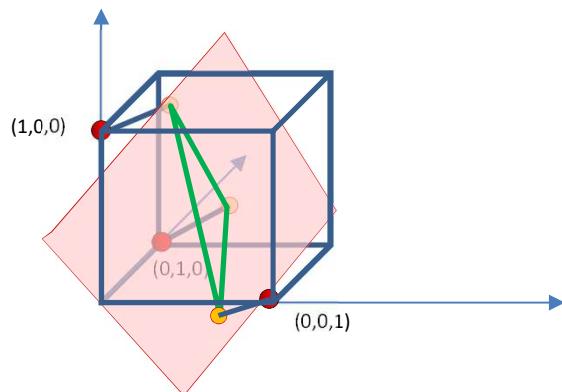
$$W_n = f(PW_1, PW_2, \dots, PW_{n-1})$$



- Project the N-dimensional one-hot word vectors into a lower-dimensional space
  - Replace every one-hot vector  $W_i$  by  $PW_i$
  - $P$  is an  $M \times N$  matrix
  - $PW_i$  is now an  $M$ -dimensional vector
  - Learn  $P$  using an appropriate objective
    - Distances in the projected space will reflect relationships imposed by the objective

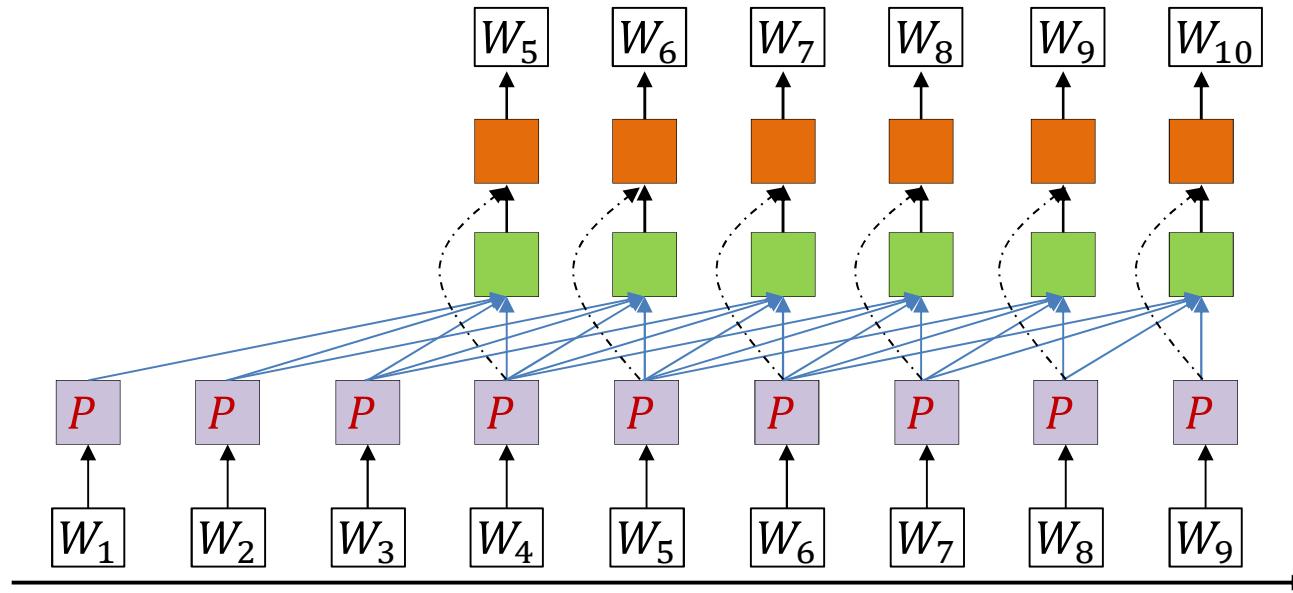
# “Projection”

$$W_n = f(PW_1, PW_2, \dots, PW_{n-1})$$



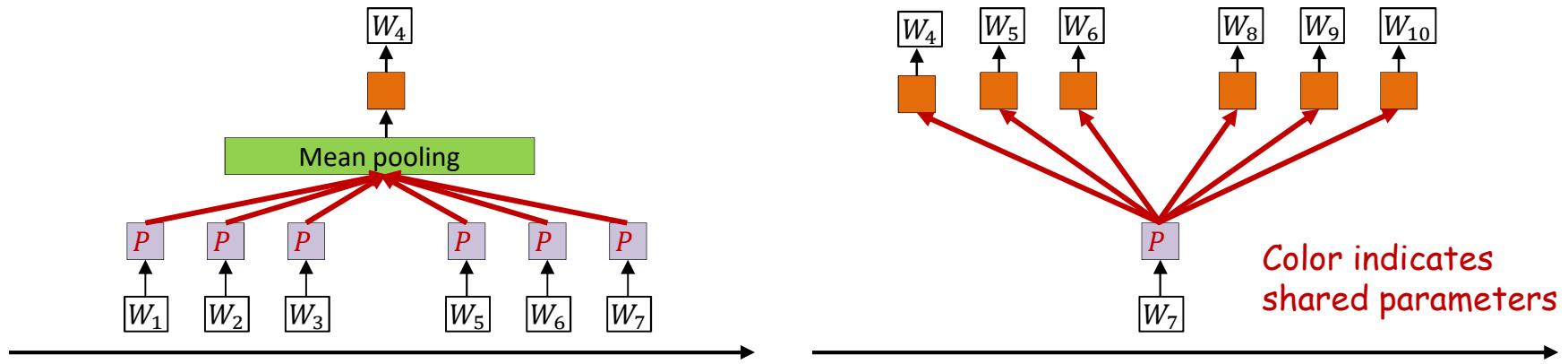
- $P$  is a simple linear transform
- A single transform can be implemented as a layer of  $M$  neurons with linear activation
- The transforms that apply to the individual inputs are all  $M$ -neuron linear-activation subnets with tied weights

# Predicting words: The TDNN model



- Predict each word based on the past N words
  - “A neural probabilistic language model”, Bengio et al. 2003
  - Hidden layer has  $\text{Tanh}()$  activation, output is softmax
- One of the outcomes of learning this model is that we also learn low-dimensional representations  $PW$  of words

# Alternative models to learn projections



- Soft bag of words: Predict word based on words in immediate context
  - Without considering specific position
- Skip-grams: Predict adjacent words based on current word

# Embeddings: Examples

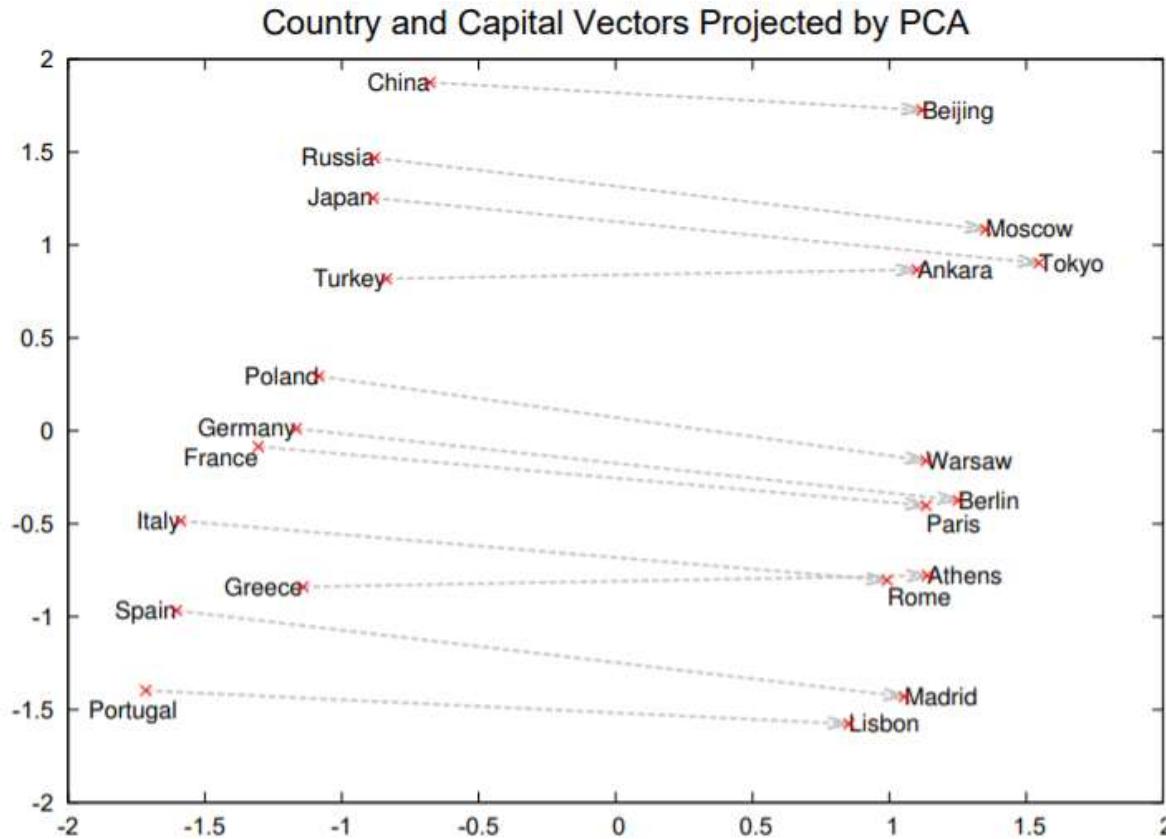
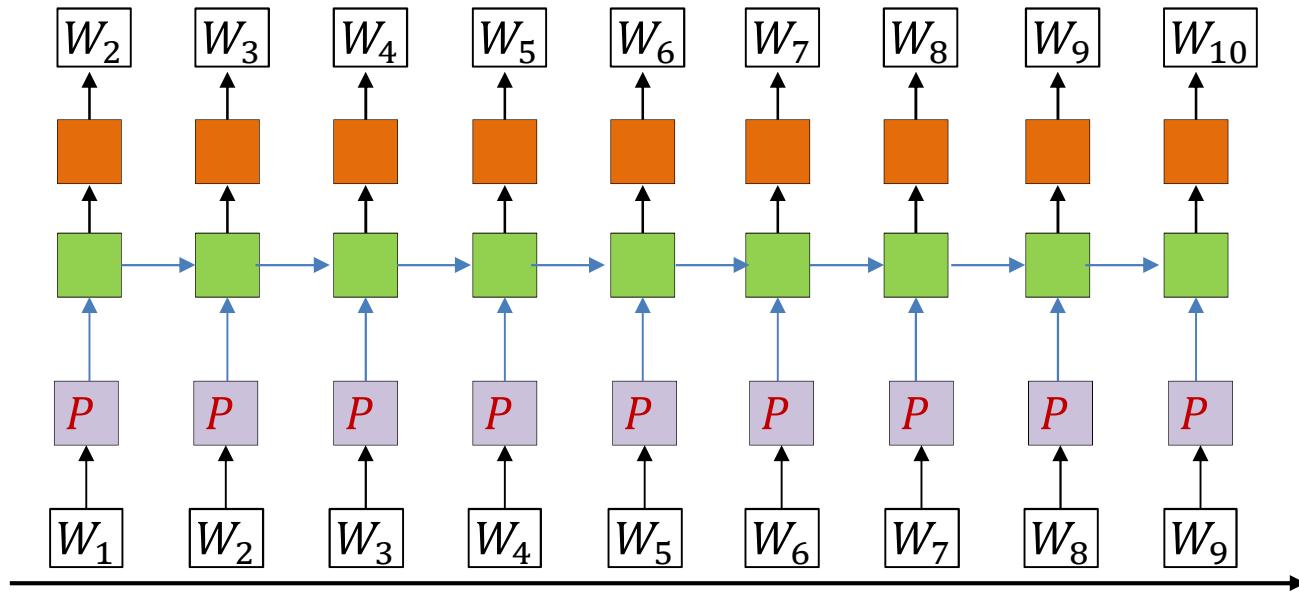


Figure 2: Two-dimensional PCA projection of the 1000-dimensional Skip-gram vectors of countries and their capital cities. The figure illustrates ability of the model to automatically organize concepts and learn implicitly the relationships between them, as during the training we did not provide any supervised information about what a capital city means.

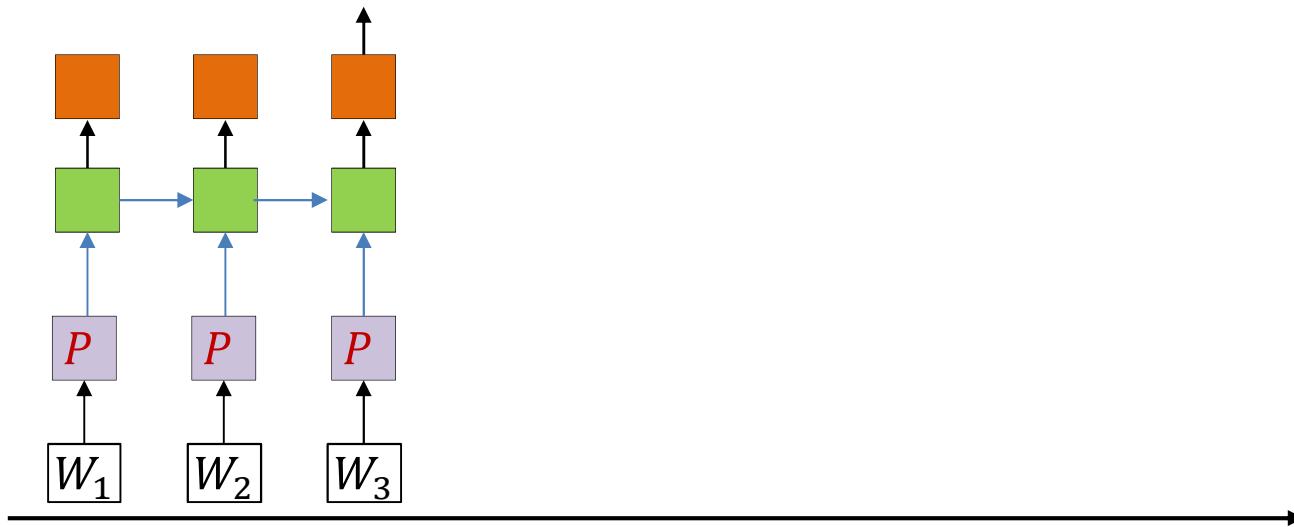
- From Mikolov et al., 2013, “Distributed Representations of Words and Phrases and their Compositionality”

# Generating Language: The model



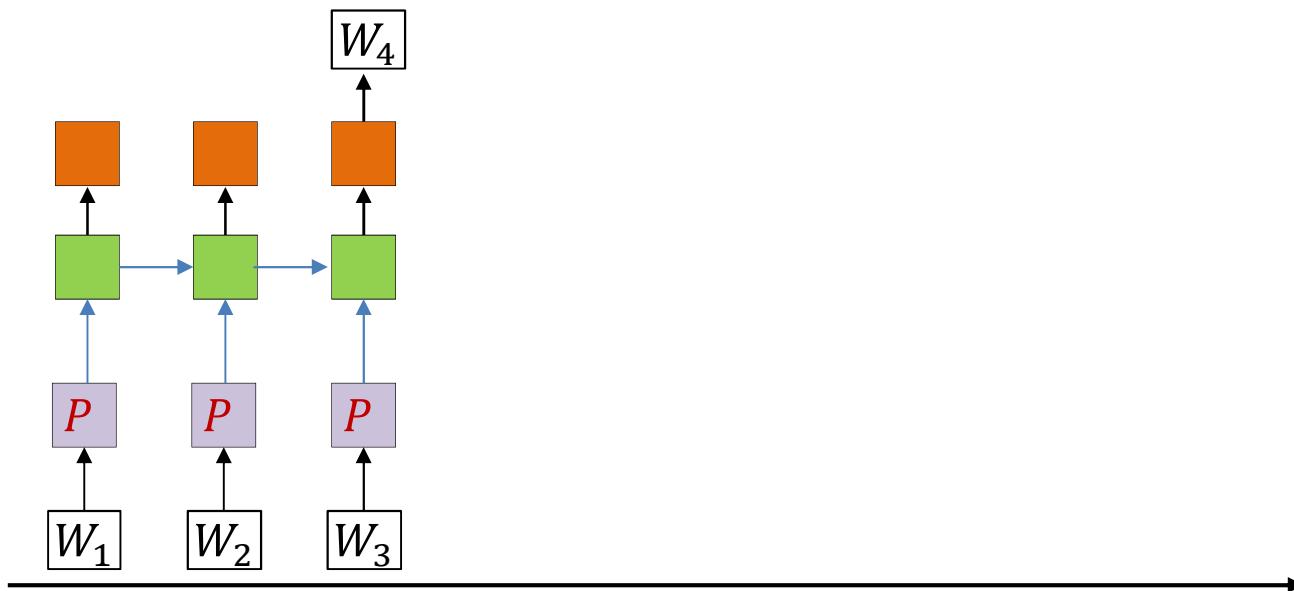
- The hidden units are (one or more layers of) LSTM units
- Trained via backpropagation from a lot of text

# Generating Language: Synthesis



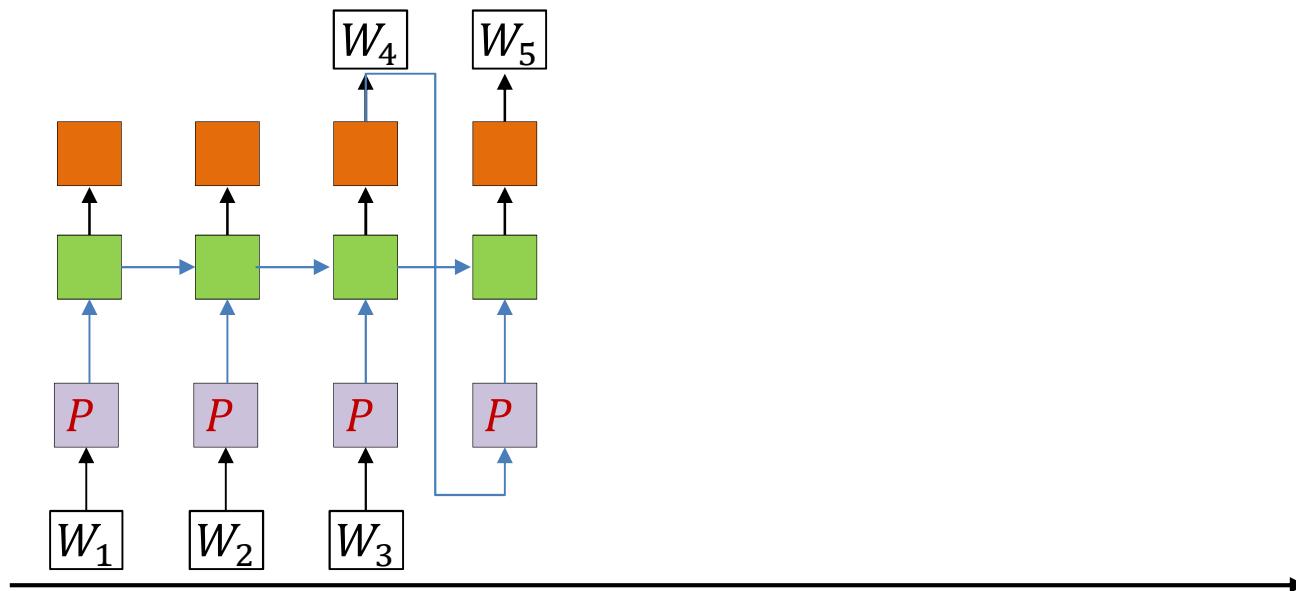
- On trained model : Provide the first few words
  - One-hot vectors
- After the last input word, the network generates a probability distribution over words
  - Outputs an N-valued probability distribution rather than a one-hot vector

# Generating Language: Synthesis



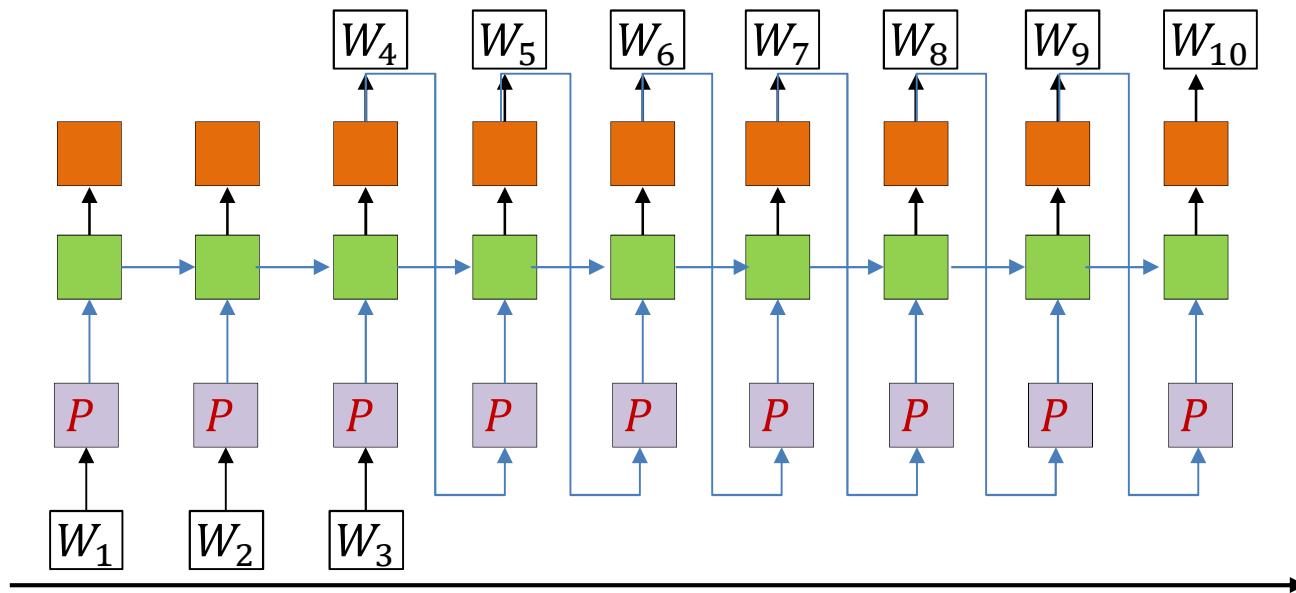
- On trained model : Provide the first few words
  - One-hot vectors
- After the last input word, the network generates a probability distribution over words
  - Outputs an N-valued probability distribution rather than a one-hot vector
- Draw a word from the distribution
  - And set it as the next word in the series

# Generating Language: Synthesis



- Feed the drawn word as the next word in the series
  - And draw the next word from the output probability distribution

# Generating Language: Synthesis



- Feed the drawn word as the next word in the series
  - And draw the next word from the output probability distribution
- Continue this process until we terminate generation
  - In some cases, e.g. generating programs, there may be a natural termination

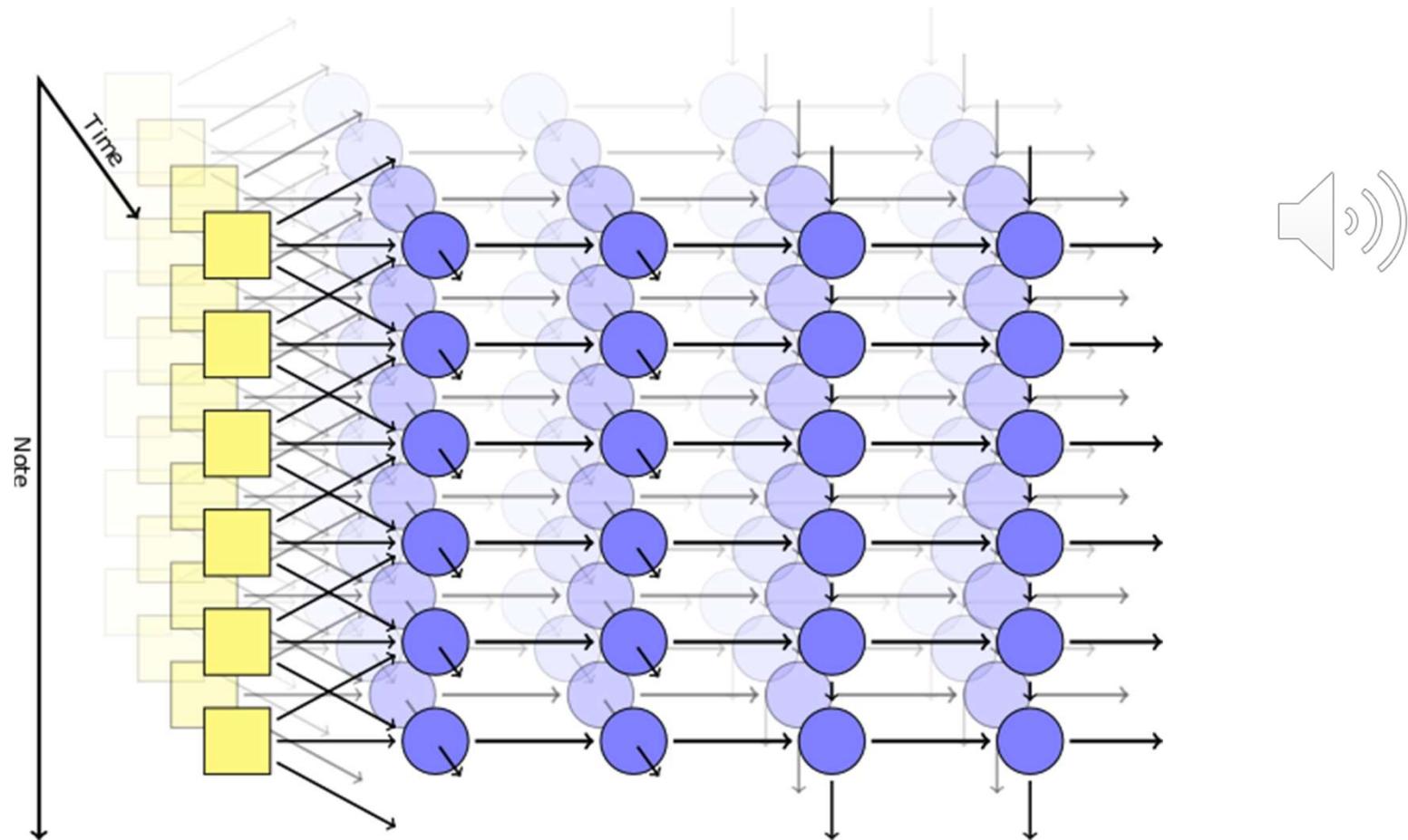
# Which open source project?

```
/*
 * Increment the size file of the new incorrect UI_FILTER group information
 * of the size generatively.
 */
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clear1(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECON
    return segtable;
}
```

Trained on linux source code

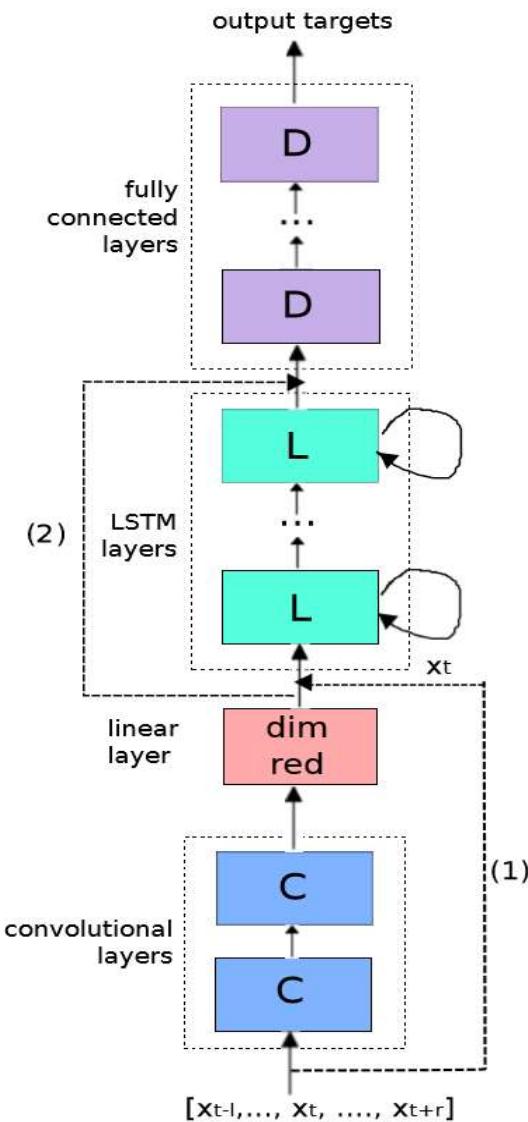
Actually uses a *character-level* model (predicts character sequences)

# Composing music with RNN



# Next up: Attention models

# CNN-LSTM-DNN for speech recognition

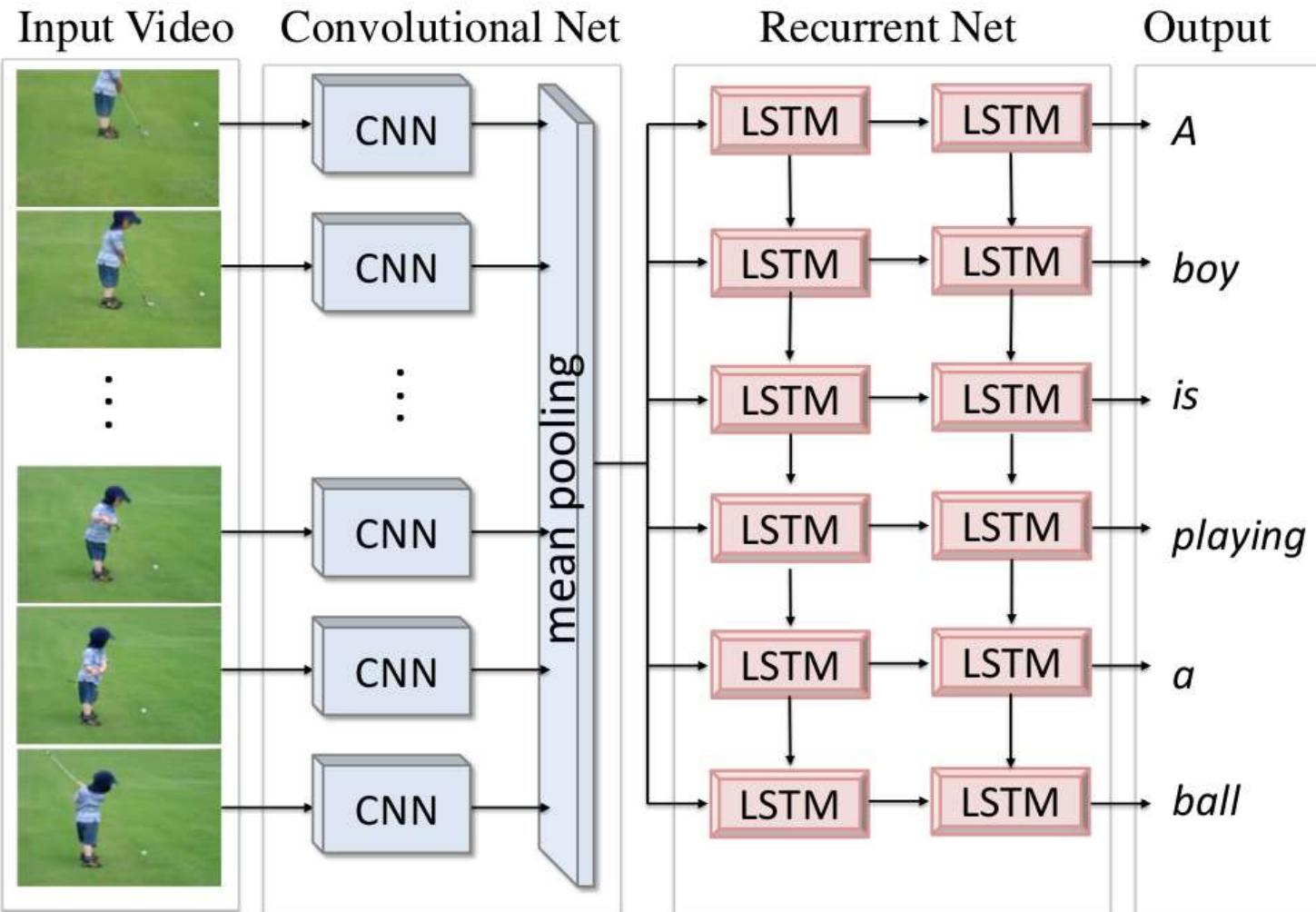


**Ensembles of RNN/LSTM, DNN, & Conv Nets (CNN) :**

T. Sainath, O. Vinyals, A. Senior, H. Sak.  
“Convolutional, Long Short-Term Memory,  
Fully Connected Deep Neural Networks,”  
ICASSP 2015.

Fig. 1. CLDNN Architecture

# Translating Videos to Natural Language Using Deep Recurrent Neural Networks

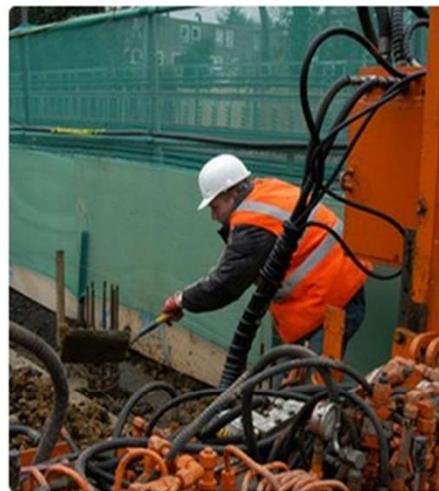


Translating Videos to Natural Language Using Deep Recurrent Neural Networks

Subhashini Venugopalan, Huijun Xu, Jeff Donahue, Marcus Rohrbach, Raymond Mooney, Kate Saenko <sup>139</sup>  
North American Chapter of the Association for Computational Linguistics, Denver, Colorado, June 2015.



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"a young boy is holding a baseball bat."



"a cat is sitting on a couch with a remote control."



"a woman holding a teddy bear in front of a mirror."



"a horse is standing in the middle of a road."

# Not explained

- Can be combined with CNNs
  - Lower-layer CNNs to extract features for RNN
- Can be used in tracking
  - Incremental prediction