

Sequence to sequence models

11-785 Introduction to Deep Learning

– lecture 16 –

TAVE Research DL001

Heeji Won

Contents

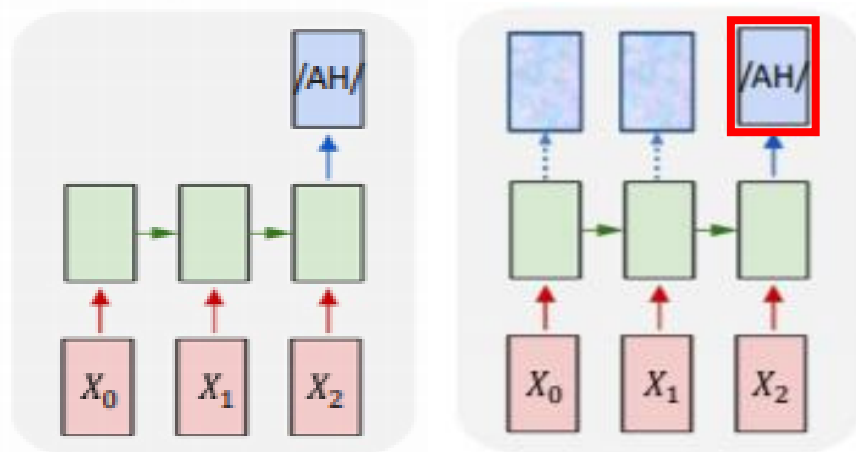
1. Many to one model
2. Sequence to sequence model

Contents

1. Many to one model
2. Sequence to sequence model

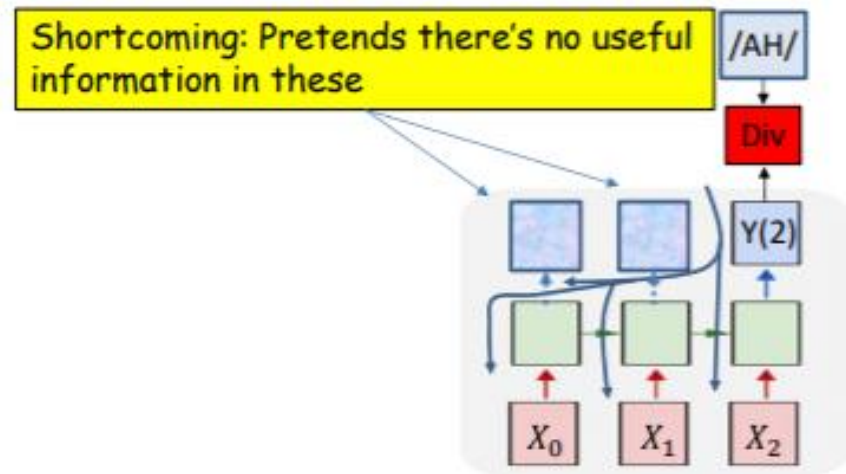
01. Many to one model

- Many to one
 - used for Q&A, Speech recognition

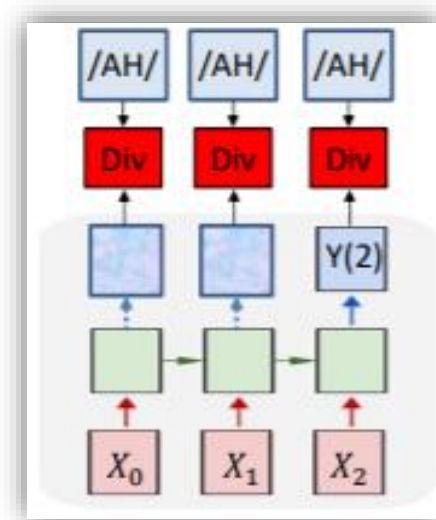


- Outputs are actually produced for every input
- But, we only read it at the end of the sequence

- Training



=> Exploit them! Assume the same output for the entire input



Define the divergence everywhere!

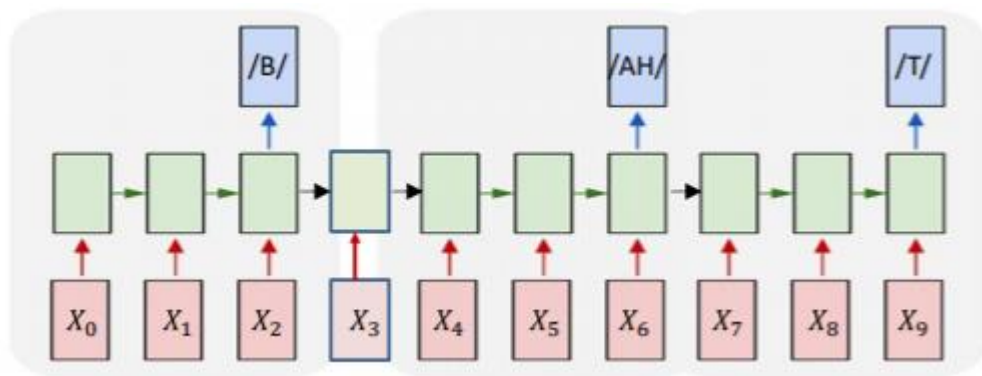
$$DIV(Y_{target}, Y) = \sum_t w_t X_{ent}(Y(t), Phoneme)$$

Contents

1. Many to one model
2. Sequence to sequence model

02. Sequence to sequence model

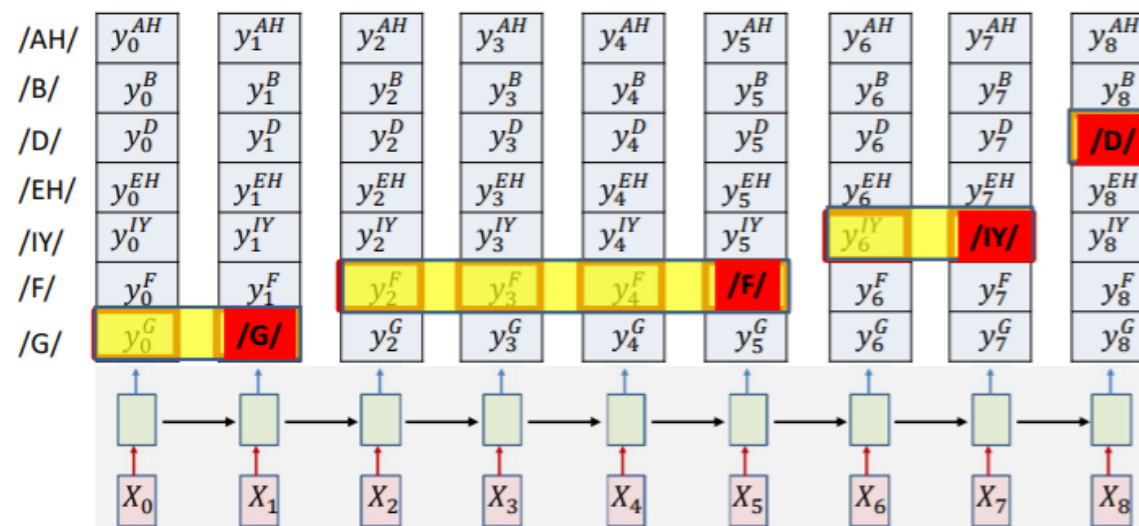
- Sequence to sequence
 - Order synchronous, but time asynchronous
 - E.g. phoneme recognition, speech recognition



- ✓ How do we know when to output symbols?
 - In fact, the network produces outputs at every time

➤ Where to output

- Option 1 : Simply select the most probable symbol at each time

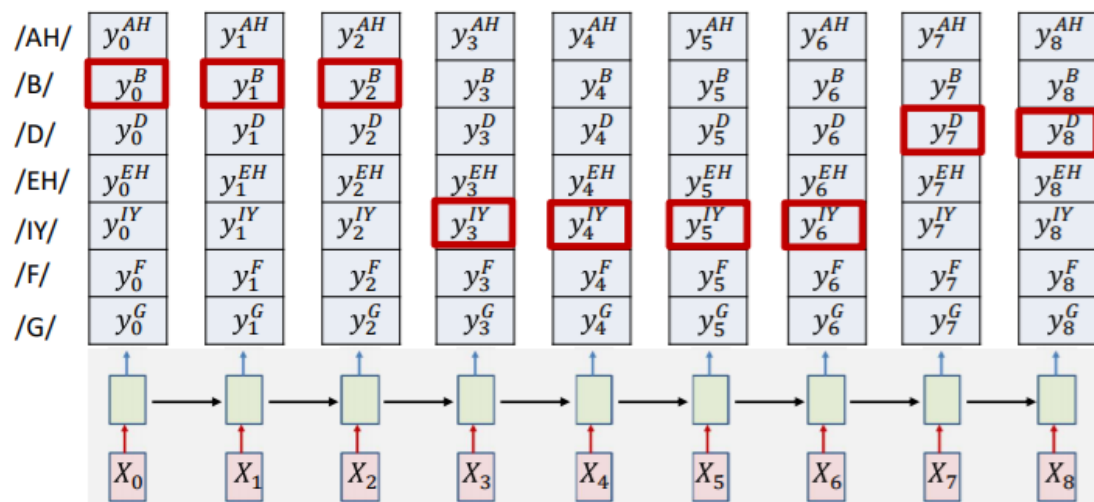


- Merge adjacent repeated symbols and place the actual emission of the symbol in the final instant
- But, resulting sequence may be meaningless

02. Sequence to sequence model

➤ Where to output

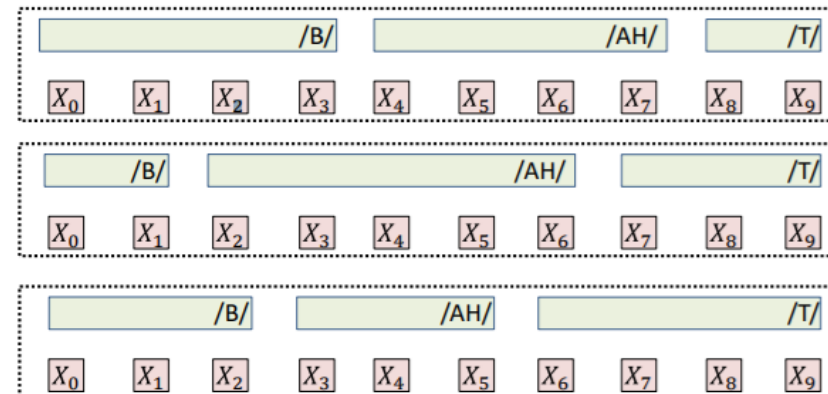
- Option 2 : Simply select the most probable symbol at each time



- E.g. only allow sequences corresponding to dictionary words
- This is a suboptimal decode that finds the most likely time-synchronous output sequence

➤ Training

- There can be various alignment of labels

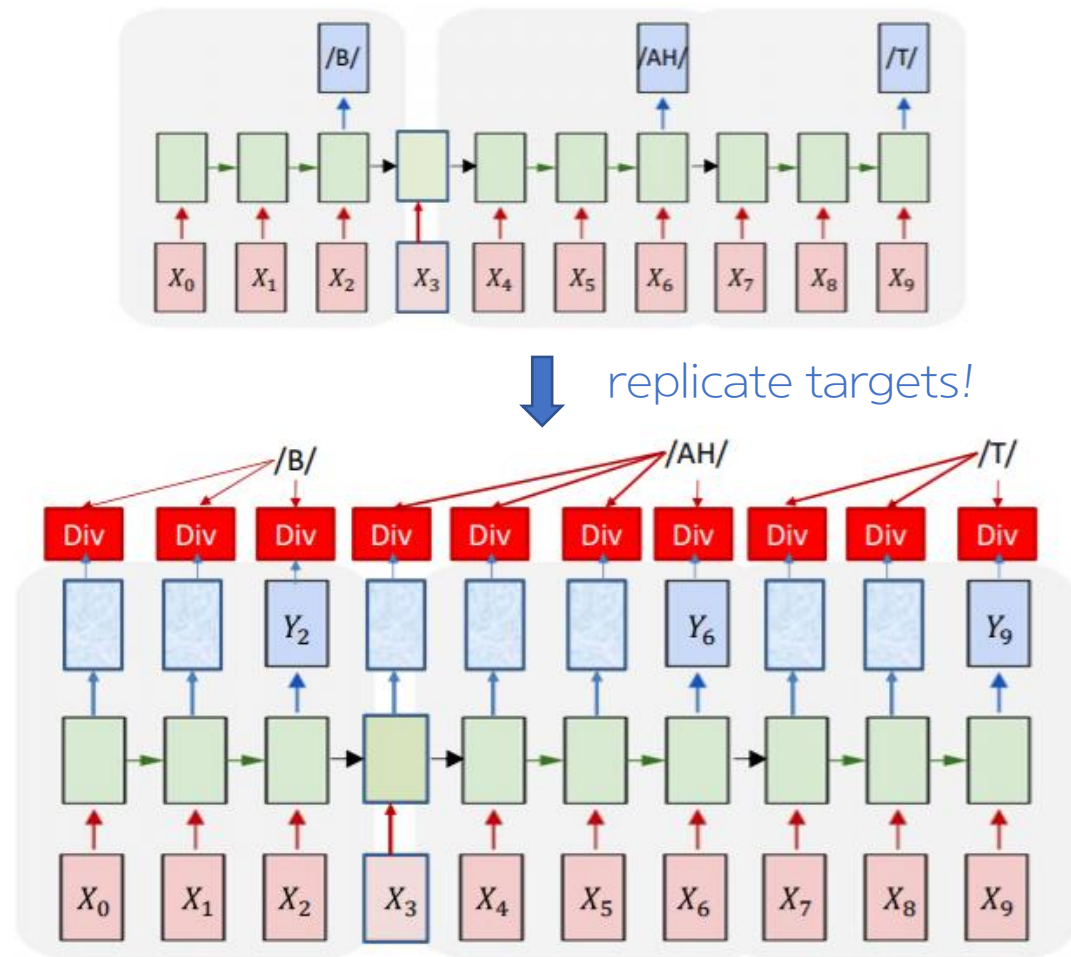


Situation 1. When we know the alignment

Situation 2. When we don't know the alignment

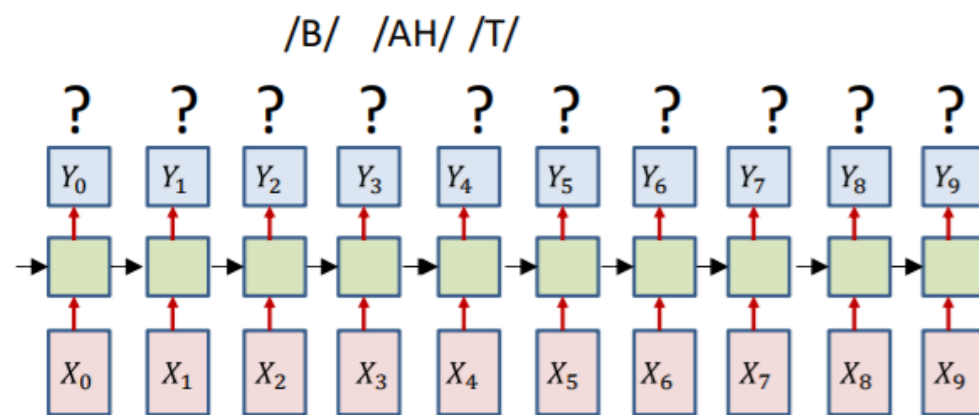
02. Sequence to sequence model

- Training with alignment



Time synchronous expansion of order synchronous seq.

- Training without alignment

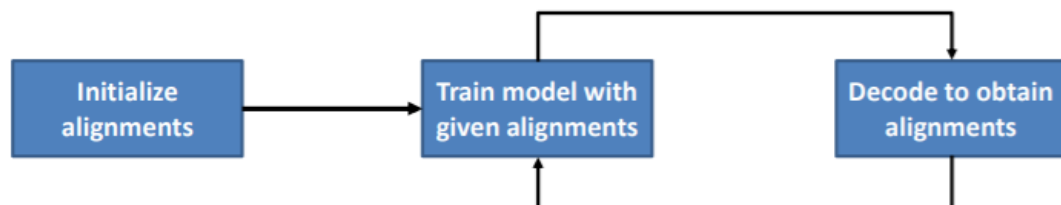


Situation 1. Guess the alignment

Situation 2. Consider all possible alignment

02. Sequence to sequence model

- Solution 1 : Guess the alignment
 - Guess an initial alignment and iteratively refine as the model improves
 - Initialize either randomly, based on some heuristic, or any other rationale



- How to estimate an alignment

Find

$$\operatorname{argmax} P(s_0, s_1, \dots, s_{N-1} | S_0, S_1, \dots, S_K, X_0, X_1, \dots, X_{N-1})$$

$$\operatorname{compress}(s_0, s_1, \dots, s_{N-1}) \equiv S_0, S_1, \dots, S_K$$

➤ Decoding

- Unconstrained decoding

/AH/	y_0^{AH}	y_1^{AH}	y_2^{AH}	y_3^{AH}	y_4^{AH}	y_5^{AH}	y_6^{AH}	y_7^{AH}	y_8^{AH}
/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
/D/	y_0^D	y_1^D	y_2^D	y_3^D	y_4^D	y_5^D	y_6^D	y_7^D	y_8^D
/EH/	y_0^{EH}	y_1^{EH}	y_2^{EH}	y_3^{EH}	y_4^{EH}	y_5^{EH}	y_6^{EH}	y_7^{EH}	y_8^{EH}
/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}
/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
/G/	y_0^G	y_1^G	y_2^G	y_3^G	y_4^G	y_5^G	y_6^G	y_7^G	y_8^G

- target : /B/ /IY/ /F/ /IY/
- output : /AH/ /AH/ ... /IY/

- ✓ Output may not correspond to an expansion of the desired symbol seq.

02. Sequence to sequence model

- Block out
 - Block out all rows that do not include symbols from the target sequence

/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}
/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



/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}
/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

- Only decode on reduced grid
- ✓ Still not assure that the decode sequence expands the target symbol seq.

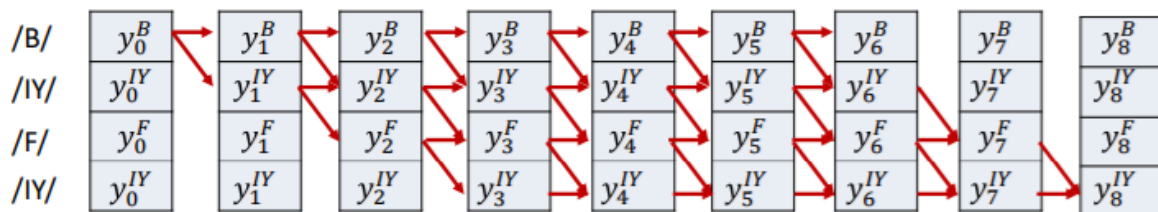
- Explicitly constrain alignment
 - Arrange the constructed table which has exact target sequence
 - the first symbol must be the top left block and the last symbol must be the bottom right
 - The rest of symbols must monotonically travel down from top left to bottom right

/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}
/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
/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}

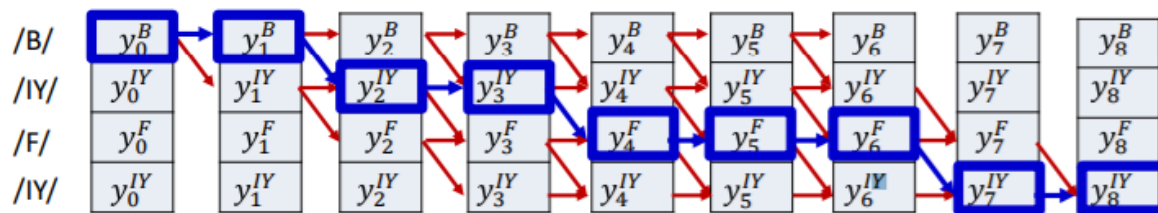
✓ The seq. is an expansion of the target seq.

02. Sequence to sequence model

- The graph representing all path

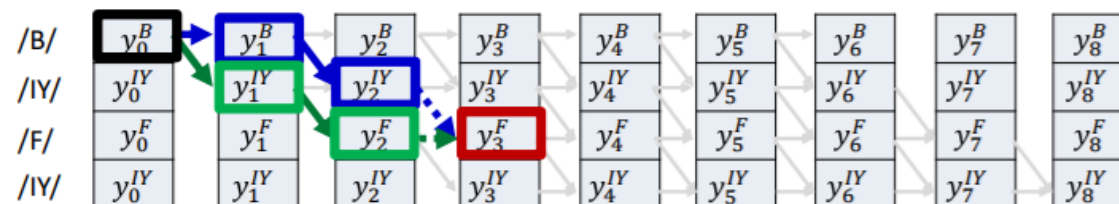


- Find the most probable path using any dynamic programming algorithm like the Viterbi algorithm



➤ Viterbi algorithm

- The best path to any node must be an extension of the best path to one of its parent nodes
- Dynamically track the best path



- Initialization

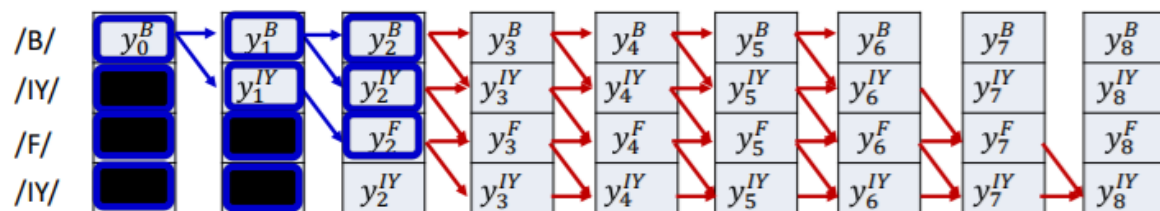
BP := Best Parent
Bscr := Bestpath Score to node

$$BP(0, i) = \text{null}, i = 0 \dots K - 1$$

$$Bscr(0, 0) = y_0^{S(0)}, Bscr(0, i) = -\infty, i = 1 \dots K - 1$$

02. Sequence to sequence model

➤ Verterbi algorithm



for $t = 1 \dots T - 1$

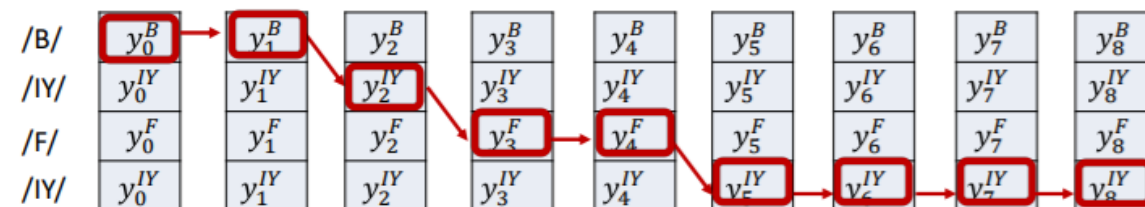
$BP(t, 0) = 0; Bscr(t, 0) = Bscr(t - 1, 0) \times y_t^{S(0)}$

for $l = 1 \dots K - 1$

- $BP(t, l) = (\text{if } (Bscr(t - 1, l - 1) > Bscr(t - 1, l)) \text{ } l - 1; \text{ else } l)$

- $Bscr(t, l) = Bscr(BP(t, l)) \times y_t^{S(l)}$

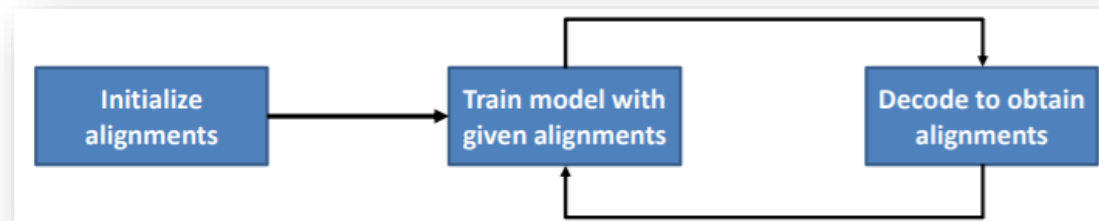
- Gradients from the alignment



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

$$\nabla_{Y_t} DIV = \begin{bmatrix} 0 & 0 & \dots & \frac{-1}{Y(t, symbol_t^{bestpath})} & 0 & \dots & 0 \end{bmatrix}$$

- The gradient is 0 except the component corresponding to the target (estimated alignment)



Thank you