

Sequence to Sequence Models

Deep Learning

Woohwan Jung

Basic Sequence2Sequence Model



$x^{<1>}$ $x^{<2>}$ $x^{<3>}$ $x^{<4>}$ $x^{<5>}$
Jane visite l'Afrique en septembre

(문장 벡터 만들기)

Attention model / $\frac{0}{0}$

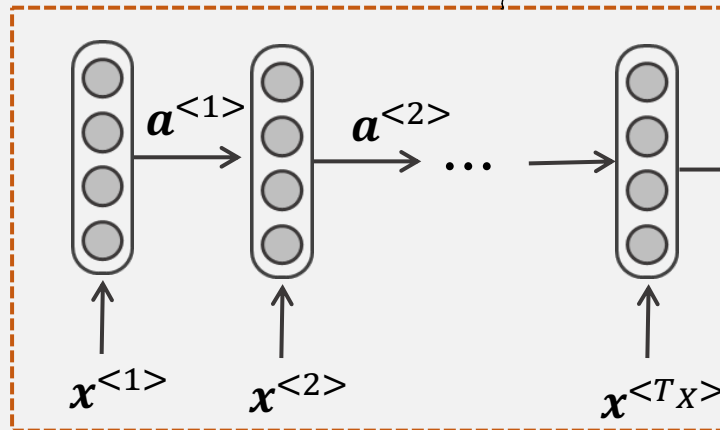
이(표)를 주

Jane is visiting Africa in September

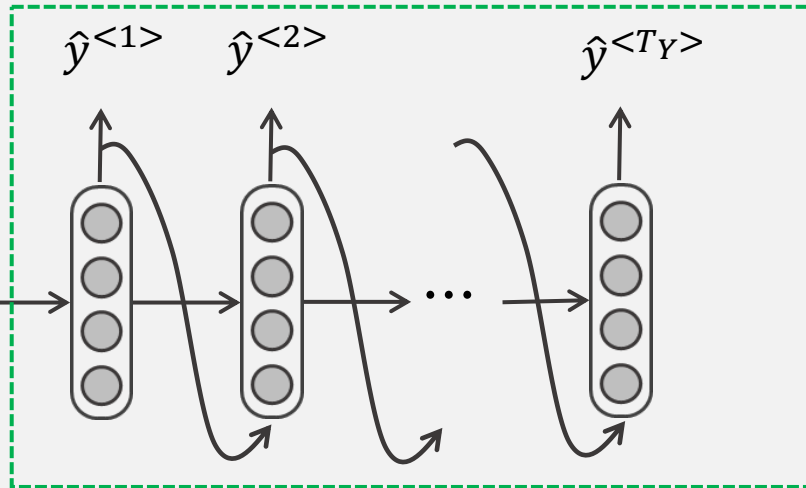
$y^{<1>}$ $y^{<2>}$ $y^{<3>}$ $y^{<4>}$ $y^{<5>}$ $y^{<6>}$

전 리정할 랑기

Encoder



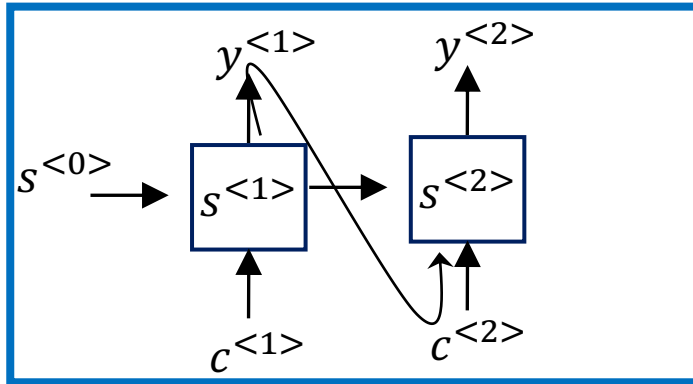
$\hat{y}^{<1>}$ $\hat{y}^{<2>}$ $\hat{y}^{<T_Y>}$



Decoder

Attention Model

Attention decoder



$$e^{<t,t'>} = v_a^T \tanh(W_a[s^{<t-1>}; a^{<t'>}])$$

$$\alpha^{<t,t'>} = \frac{\exp(e^{<t,t'>})}{\sum_{t'=1}^{T_x} \exp(e^{<t,t'>})}$$

문맥 벡터

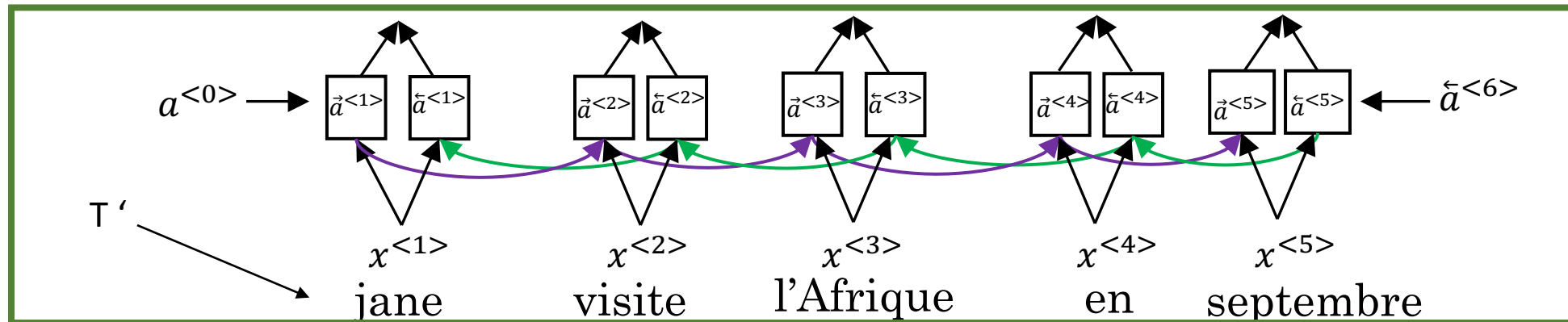
$$c^{<t>} = \sum_{t'} \alpha^{<t,t'>} a^{<t'>}$$

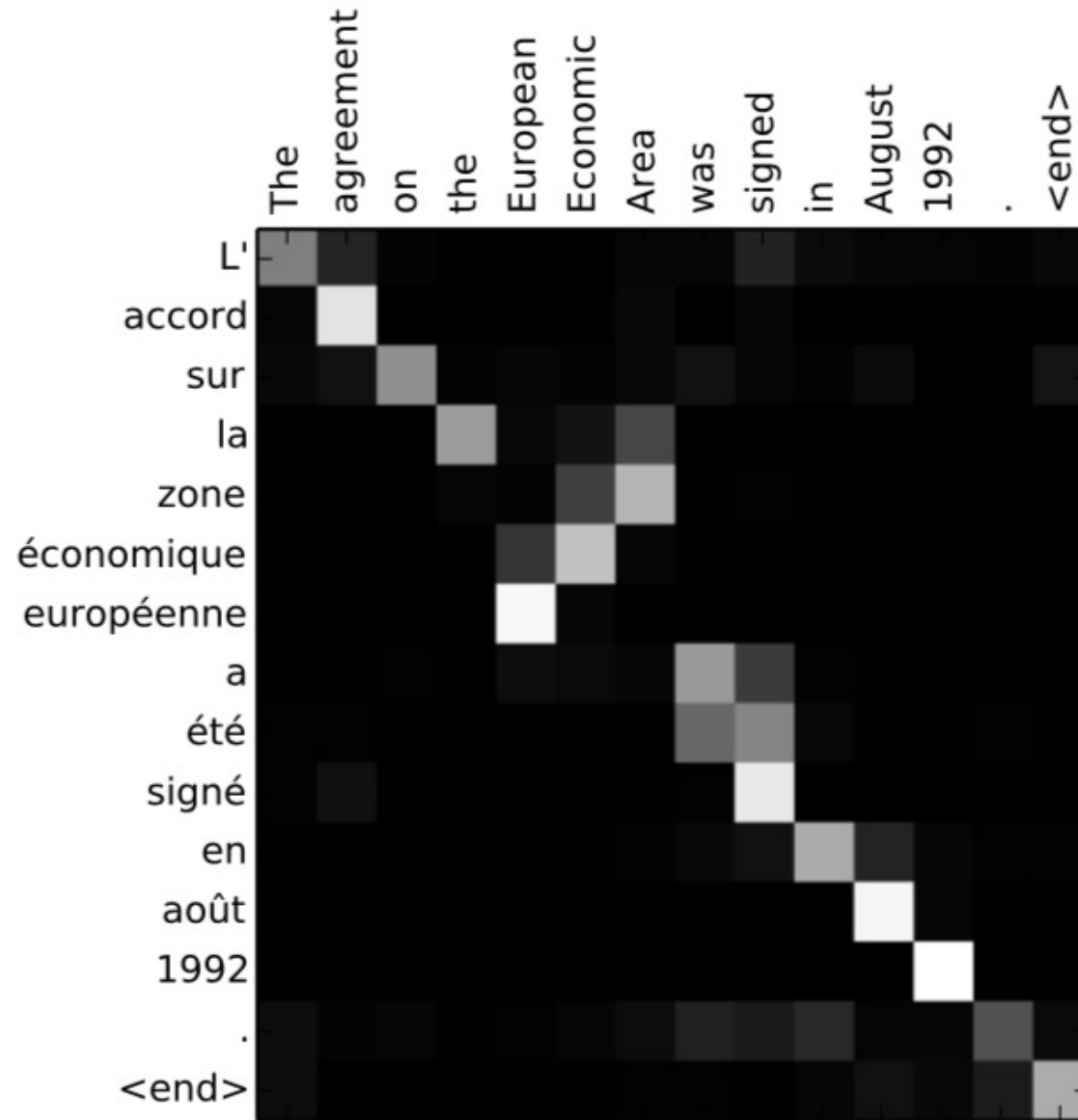
Attention

$$s_t = \text{RNNCell}(s_{t-1}, y_{t-1}, c_t)$$

+ 어텐션 디코더 Basic Encoder De

(기본 S2S)
(Bidirectional) encoder $a^{<t'>} = (\vec{a}^{<t'>}, \tilde{a}^{<t'>})$



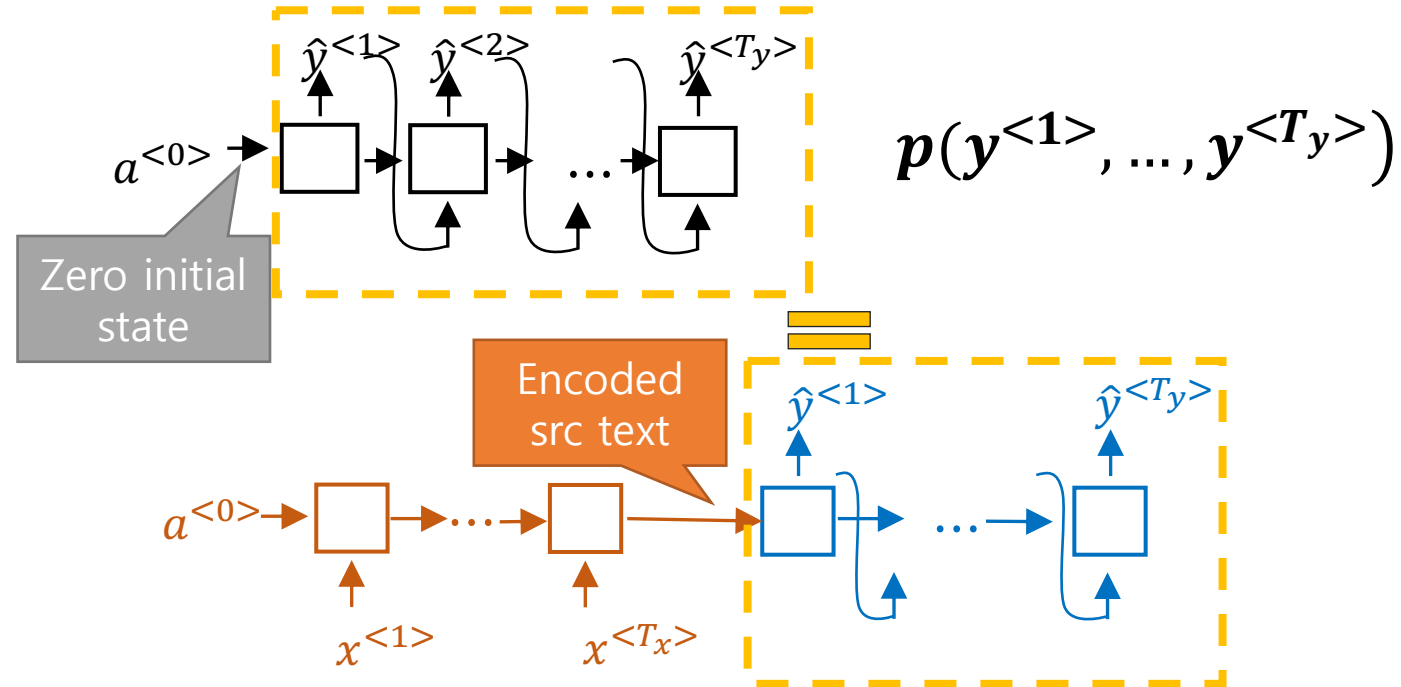


**Generating the most likely
sentence**

Machine Translation as Building a Conditional Language Model

2/0/21

Language model:



Machine translation:

Conditional language model

$$p(y^{<1>}, \dots, y^{<T_y>} \mid x^{<1>}, \dots, x^{<T_x>})$$

Finding the Most likely Translation

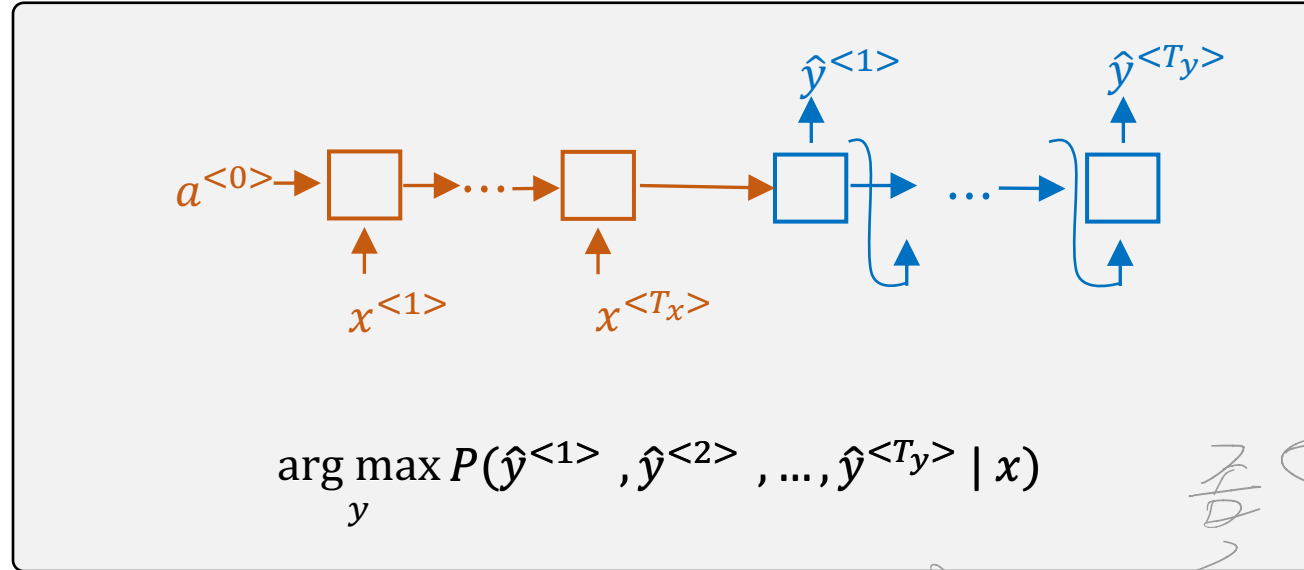
Jane visite l'Afrique en septembre.

$$\frac{P(y^{<1>}, \dots, y^{<T_y>} | x)}{\text{English}} \quad \begin{array}{c} \nearrow \\ \text{French} \end{array}$$

- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.
- In September, Jane will visit Africa.
- Her African friend welcomed Jane in September.

$$\arg \max_{y^{<1>}, \dots, y^{<T_y>}} P(\hat{y}^{<1>}, \hat{y}^{<2>}, \dots, y^{<T_y>} | x)$$

Why not a Greedy Search? 그리디는 안됨



- Jane is visiting Africa in September.
- Jane is going to be visiting Africa in September.

$$P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$$

그리드는 이거 못라게 할

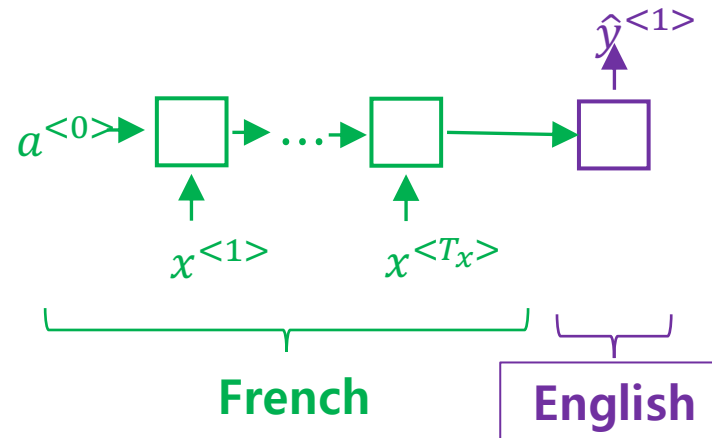
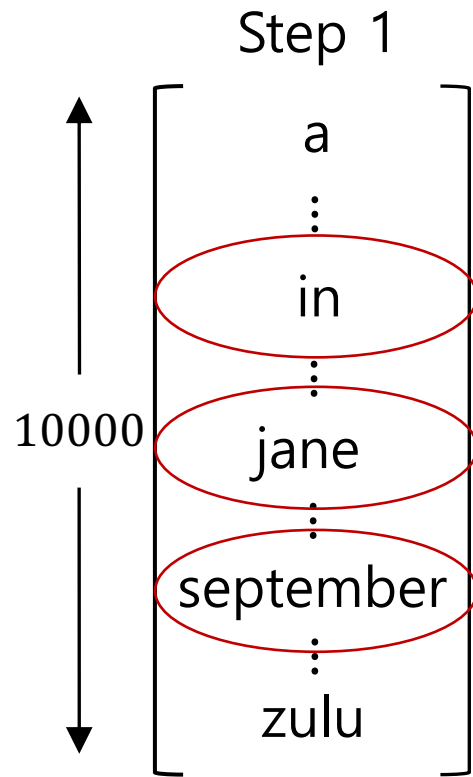
그리드는
이거 못라게 할

Beam Search



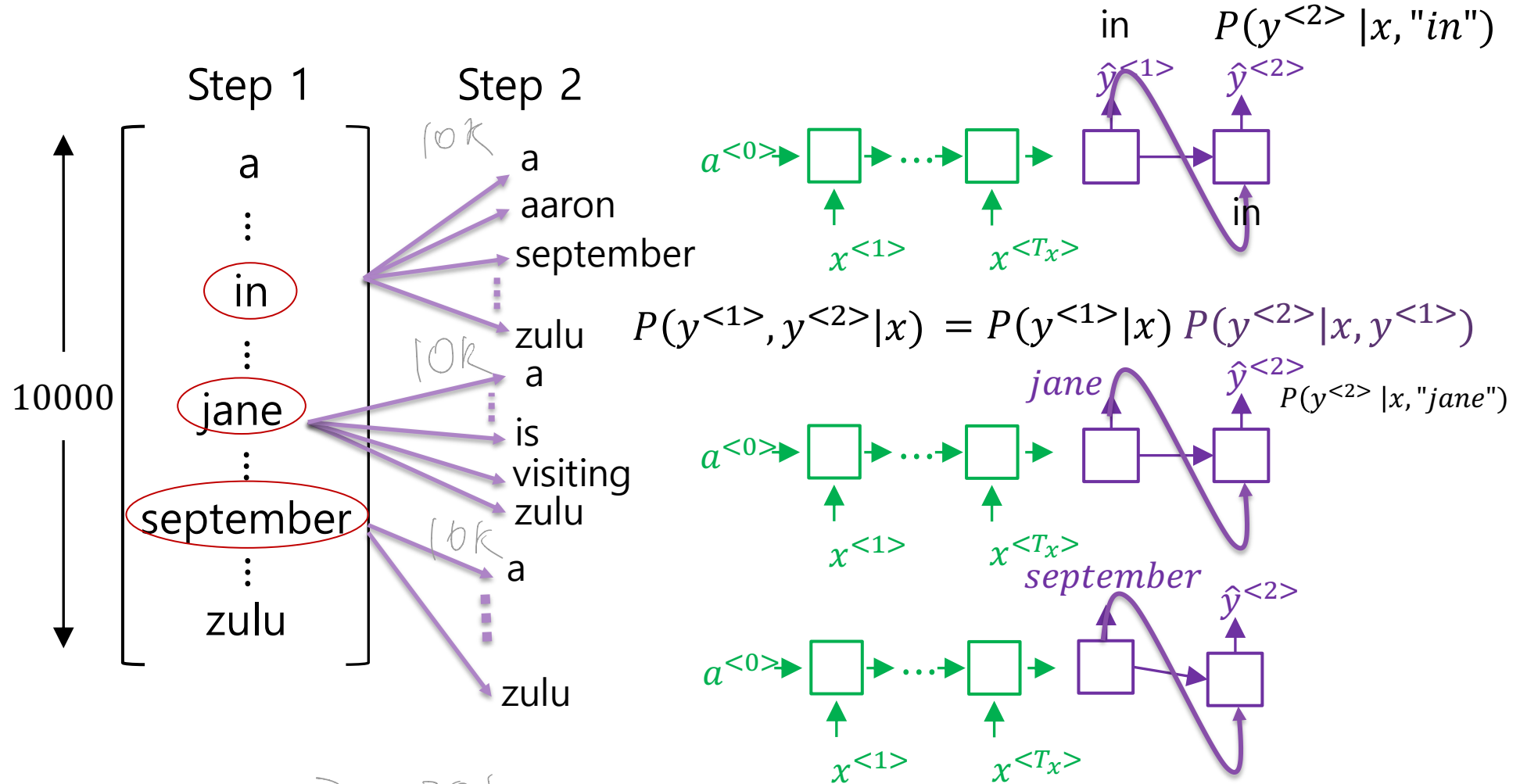
Beam Search

B = 3 (Beam width)



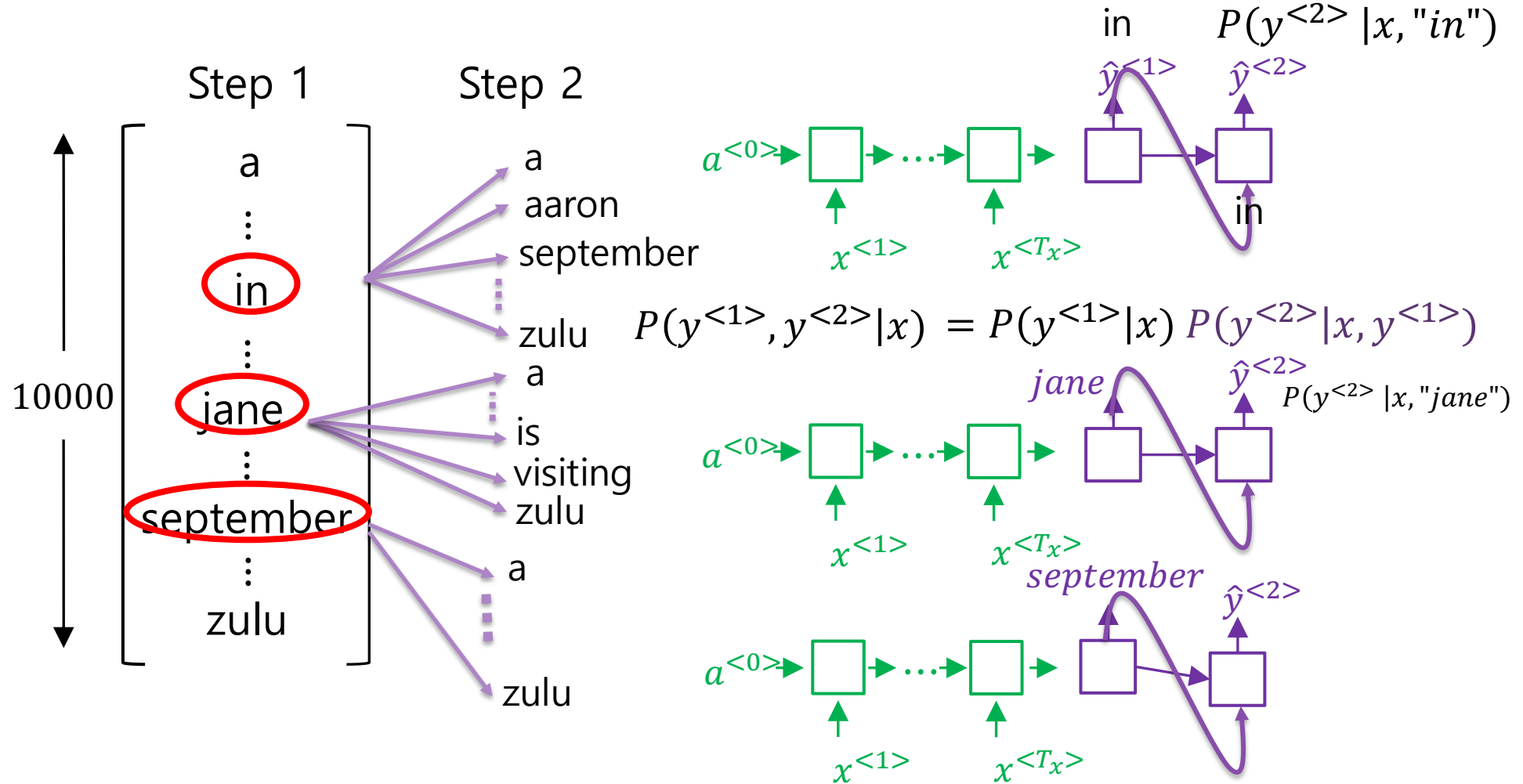
$$P(y^{<1>} | x)$$

Beam Search (B=3)

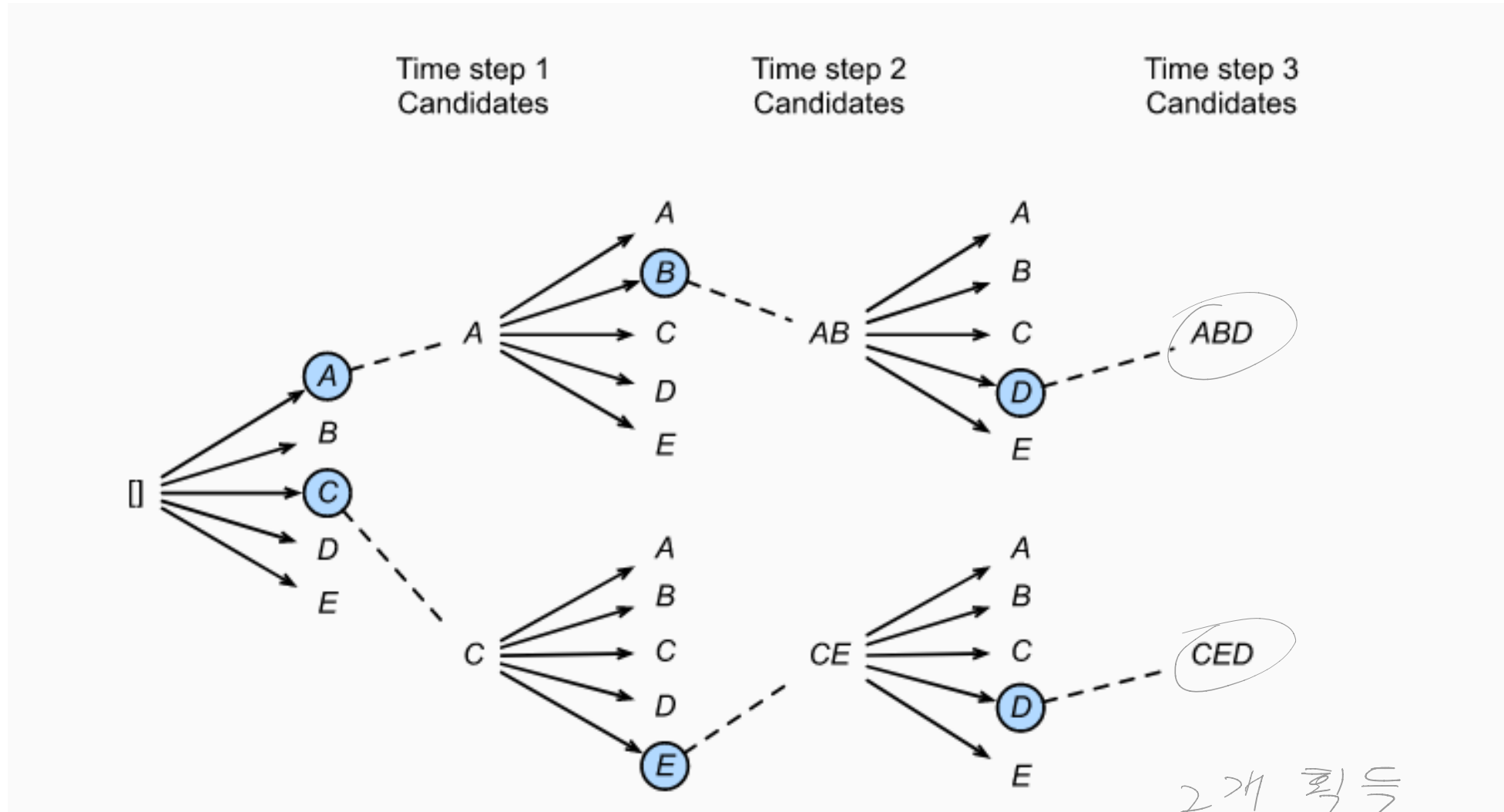


30k
top 371 12 (5)

Beam Search (B=3)



Beam Search (B=2)



[0개 중

2개 선택

2개 확 등

Refinements to Beam Search : Length Normalization

길이 정규화

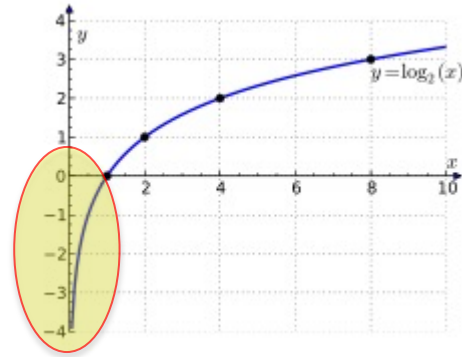
$$P(y^{<1>}, \dots, y^{<T_y>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>}) \dots P(y^{<T_y>} | x, y^{<1>}, \dots, y^{<T_y-1>})$$

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

log P

$$\arg \max_y \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

<1



log 값 0보다 작음 → 길이가 짧은게 골라질 확률 높음

$$\frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{<t>} | x, y^{<1>}, \dots, y^{<t-1>})$$

정규화

$\alpha = 0.7$ 적절.

$\alpha = 0?$ $\alpha = 1?$

short 문장 long 문장

Beam Width B

- The greater the beam width, the fewer states are pruned
 - Large B: Better result, slower
 - Small B: Worse result, faster

$B = \infty \rightarrow \text{품질} \uparrow$

$B = 1 \rightarrow \text{그리디, 시간 효율적}$

Preview - Transformers

Writing a deep learning paper

- You create a new deep learning model and write a paper about it. What content should it include?

- 모델 구조

- 경험치

- 문제 정의 ($x \mapsto y$)

Transformers: The Architecture Behind ChatGPT



ChatGPT

- Task
 - Chatbot
- The architecture
 - Transformer (we'll study this next week)
- Training
 - Unsupervised Pretraining. (our last topic)
 - Reinforcement Learning

