

ICT이노베이션스퀘어 AI복합교육 고급 언어과정

# 자연어처리를 위한 Language Model

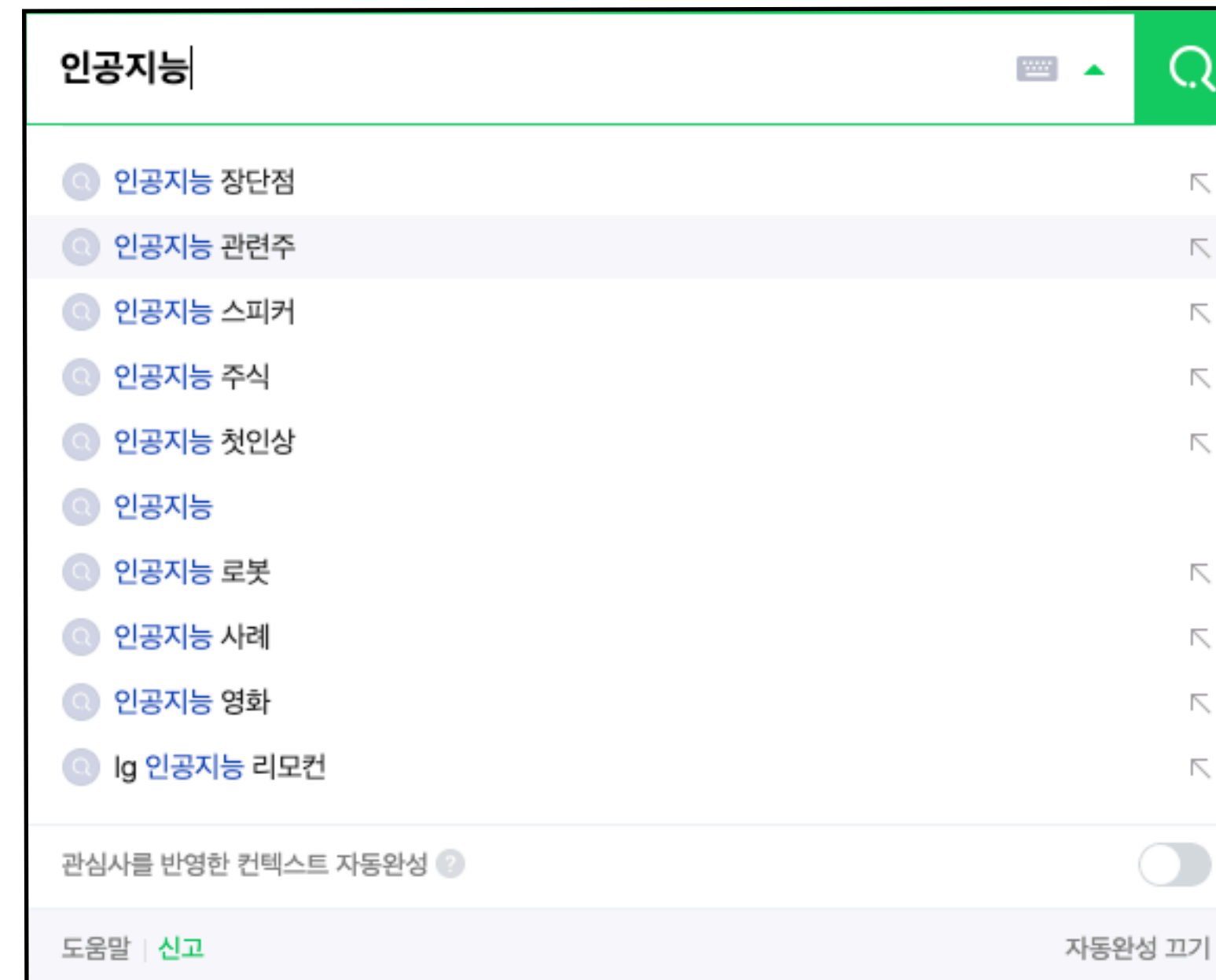
현청천

2021.04.19

# What is Language Model

**Language Model은 언어의 확률분포를 추정하는 것**

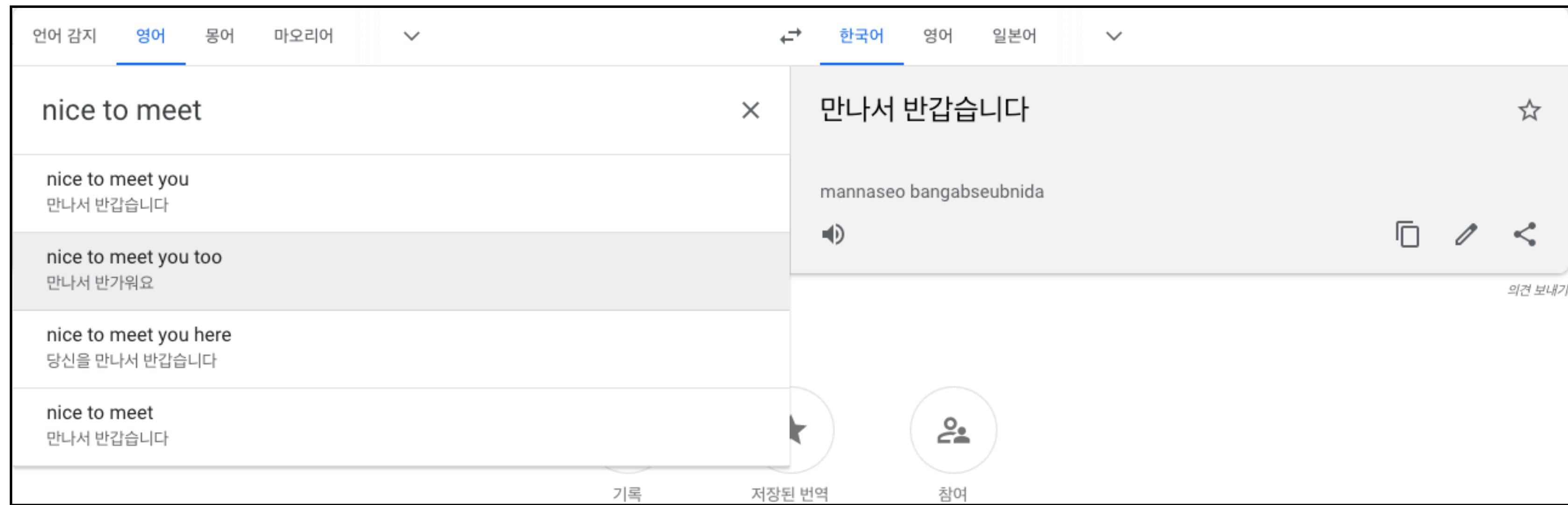
# What is Language Model



○ 자동 완성기능

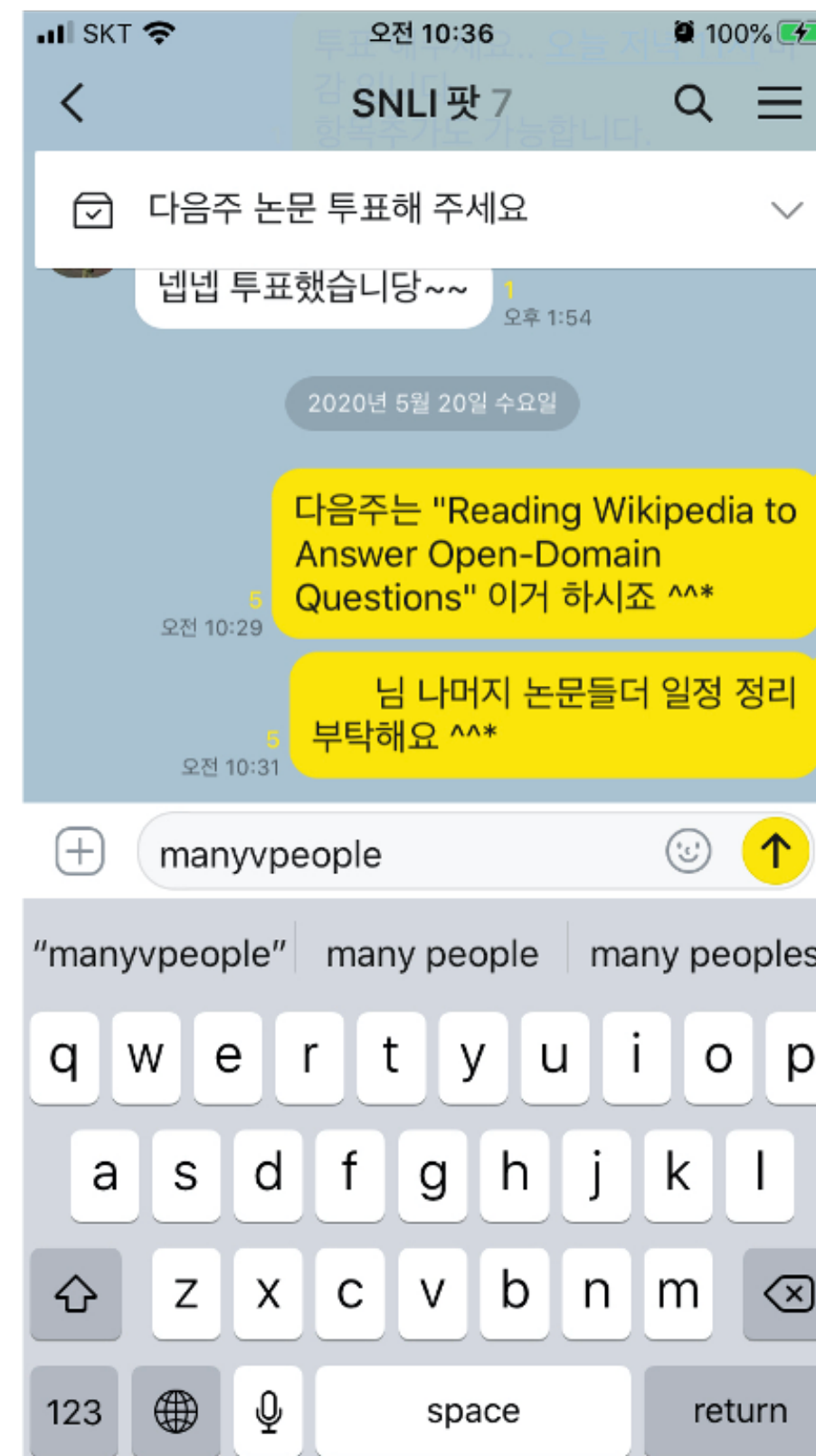
자동완성

# What is Language Model



자동완성

# What is Language Model



○ 후보군

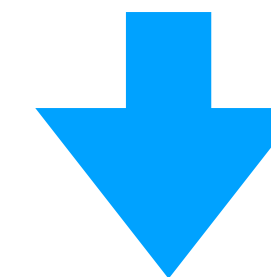
오타

# What is Language Model



아버지가 방에 들어가신다  
아버지 가방에 들어가신다  
아버지가방 들어가신다

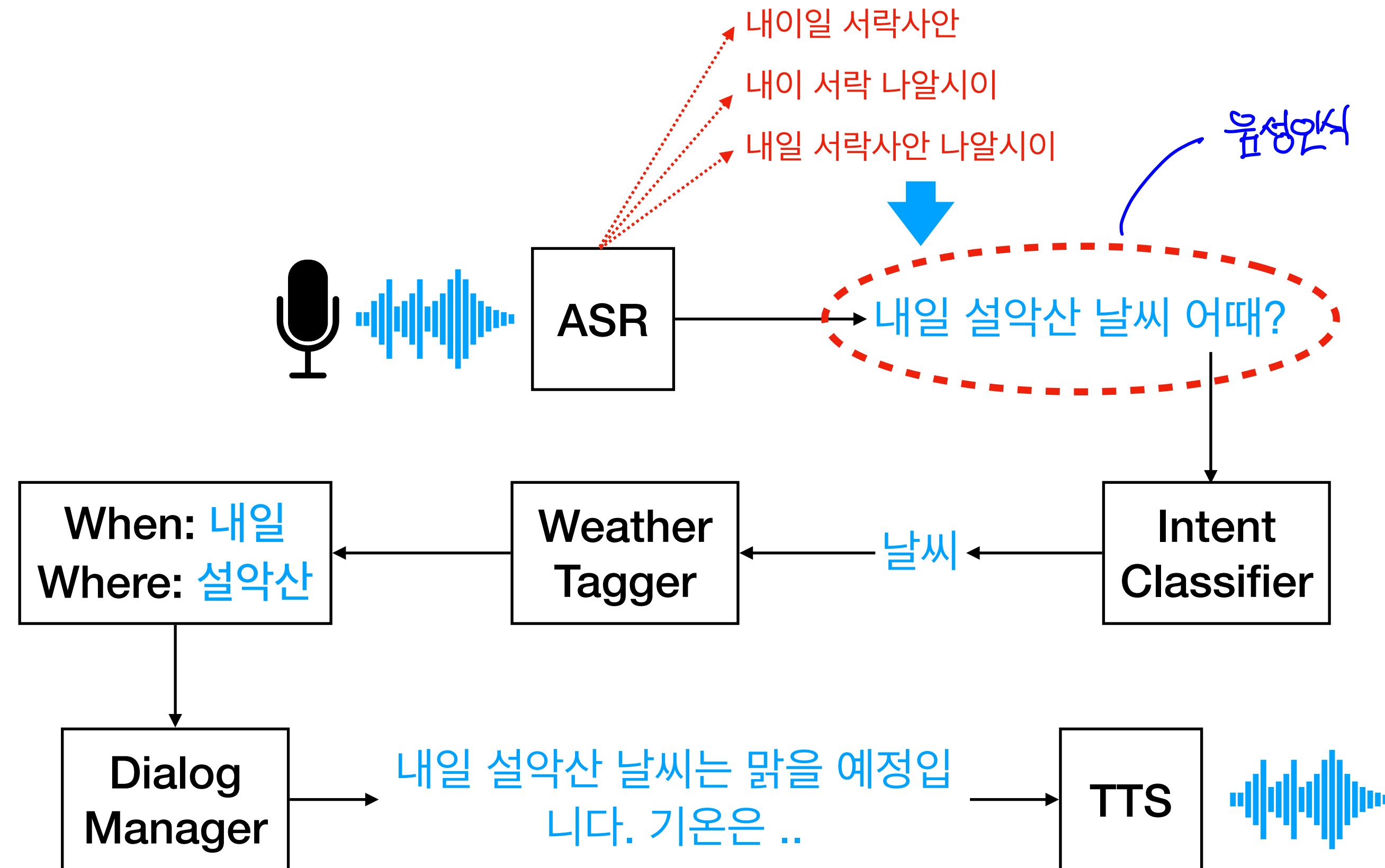
-----



아버지가 방에 들어가신다

사람간의 대화

# What is Language Model



음성인식에서 Language Model을 이용해 음성을 문자로 변환

# What is Language Model

자연어에서 발생할 확률

$p(\text{그는 사과를 보자 배고픔을 느꼈다}) > p(\text{그는 사과를 보자 외로움을 느꼈다})$

$p(\text{그녀는 운동을 열심히 한다}) > p(\text{그녀는 운동을 몇몇이 한다})$

확률된 확률표를 사용하는 거지.



# What is Language Model

실제 언어의 확률분포를 아는 것은 어려움

좋은 근사치를 제공하는 Language Model을 정의 할 수 있음

만들게 욕망이겠지.

# N-gram Language Model

48봉지 2620개의 M&M의 컬러 분포



372 + 544 + 369 + 483 + 481 + 371 = N

이 데이터로부터 확률 분포를 추론하는 방법은?

$$\underline{p(\text{color})} = \frac{\text{count}(\text{color})}{N}, \quad N = \sum_{\text{color}} \text{count}(\text{color})$$

# N-gram Language Model

Word sequence로부터 확률 분포를 추론하는 방법

$$s = (w^{(1)}, w^{(2)}, \dots, w^{(n)})$$

o 단어들의 순서의 배열 = 문장

많은 text corpus가 있다면

ex) 법률, 회계.. 각각에는  
bias가 있다..

이 corpus로부터 확률 분포를 추론할 수 있음

# N-gram Language Model

$$p(s = w^{(1)}, w^{(2)}, \dots, w^{(n)}) \quad p(s = \text{the cat slept quietly})$$

$$p(w^{(1)} = \text{the}, w^{(2)} = \text{cat}, w^{(3)} = \text{slept}, w^{(4)} = \text{quietly})$$

$$p(\text{quietly} | \text{the cat slept}) \cdot p(\text{slept} | \text{the cat}) \cdot p(\text{cat} | \text{the}) \cdot p(\text{the})$$

$$p(w^{(1)}, w^{(2)}, \dots, w^{(n)}) = \left( \prod_{i=1}^n p(w^{(i)} | w^{(1)}, \dots, w^{(i-1)}) \right)$$

$\text{the } i=1$   
 $\text{cat } i=2$   
 $\text{slept } i=3$

$w^1, w^2$

# N-gram Language Model

○ 너무 어려우니까 전체를 보지  
○ 몇개만 보자!

## Independent Assumption

단어의 분포는 고정된 몇 개의 이전 단어에 의존함

○ 너무

$$p(w^{(i)} | w^{(1)}, w^{(2)}, \dots, w^{(i-1)}) \dashrightarrow p(w^{(i)} | w^{(i-n+1)}, w^{(i-n+2)}, \dots, w^{(i-1)})$$

$$\text{Trigram: } p(w^{(i)} | w^{(1)}, w^{(2)}, \dots, w^{(i-1)}) \approx p(w^{(i)} | w^{(i-2)}, w^{(i-1)})$$

$n=3$

$i-3+1 = i-2$  까지

$$\text{bigram: } p(w^{(i)} | w^{(1)}, w^{(2)}, \dots, w^{(i-1)}) \approx p(w^{(i)} | w^{(i-1)})$$

$n=2$

$i-2+1 = i-1$  까지

$$\text{unigram: } p(w^{(i)} | w^{(1)}, w^{(2)}, \dots, w^{(i-1)}) \approx p(w^{(i)})$$

$n=1$

# N-gram Language Model

$$p(w^{(i)} | w^{(1)}, w^{(2)}, \dots, w^{(i-1)}) \dashrightarrow p(w^{(i)} | w^{(i-n+1)}, w^{(i-n+2)}, \dots, w^{(i-1)})$$
$$= \frac{p(w^{(i-n+1)}, w^{(i-n+2)}, \dots, w^{(i-1)}, w^{(i)})}{p(w^{(i-n+1)}, w^{(i-n+2)}, \dots, w^{(i-1)})}$$

○ 독립성인가?

N-gram과 (N-1)-gram의 확률 분포를 어떻게 구할 것인가?

➡ 큰 text corpus에서 개수를 세면 분포를 구할 수 있음

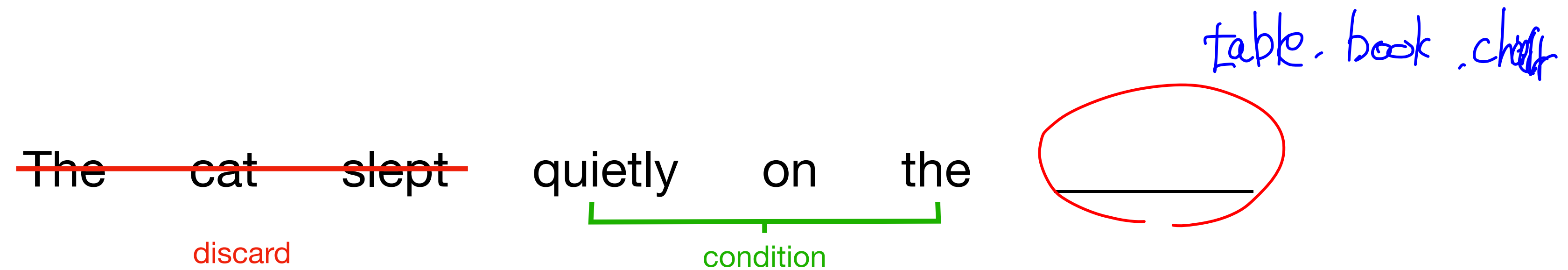
$$\frac{\text{count}(w^{(i-n+1)}, w^{(i-n+2)}, \dots, w^{(i-1)}, w^{(i)})}{\text{count}(w^{(i-n+1)}, w^{(i-n+2)}, \dots, w^{(i-1)})}$$

(Statistical approximation)

단어 배열 발생 횟수

여러 단어 배열이 발생한 횟수

# N-gram Language Model



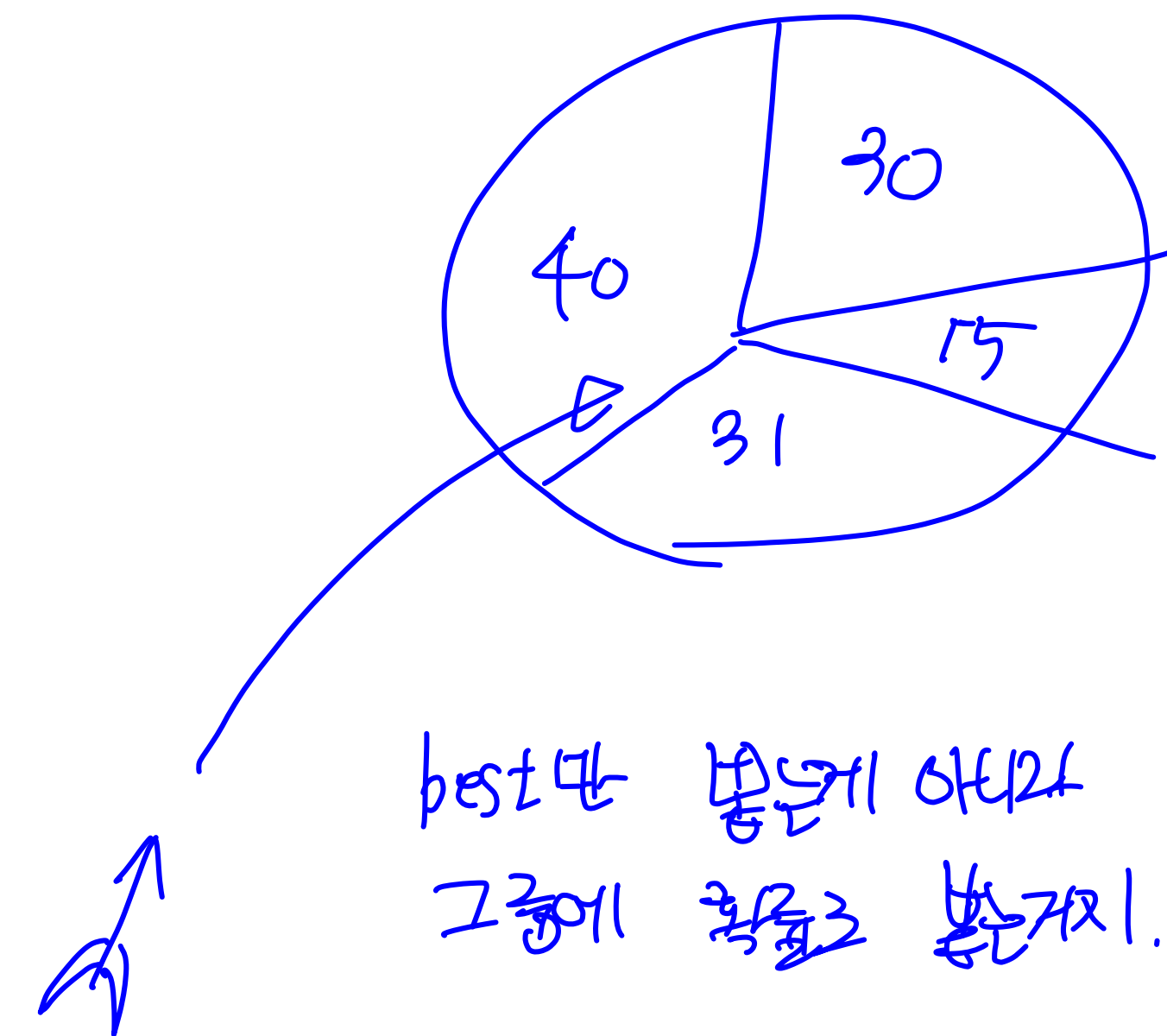
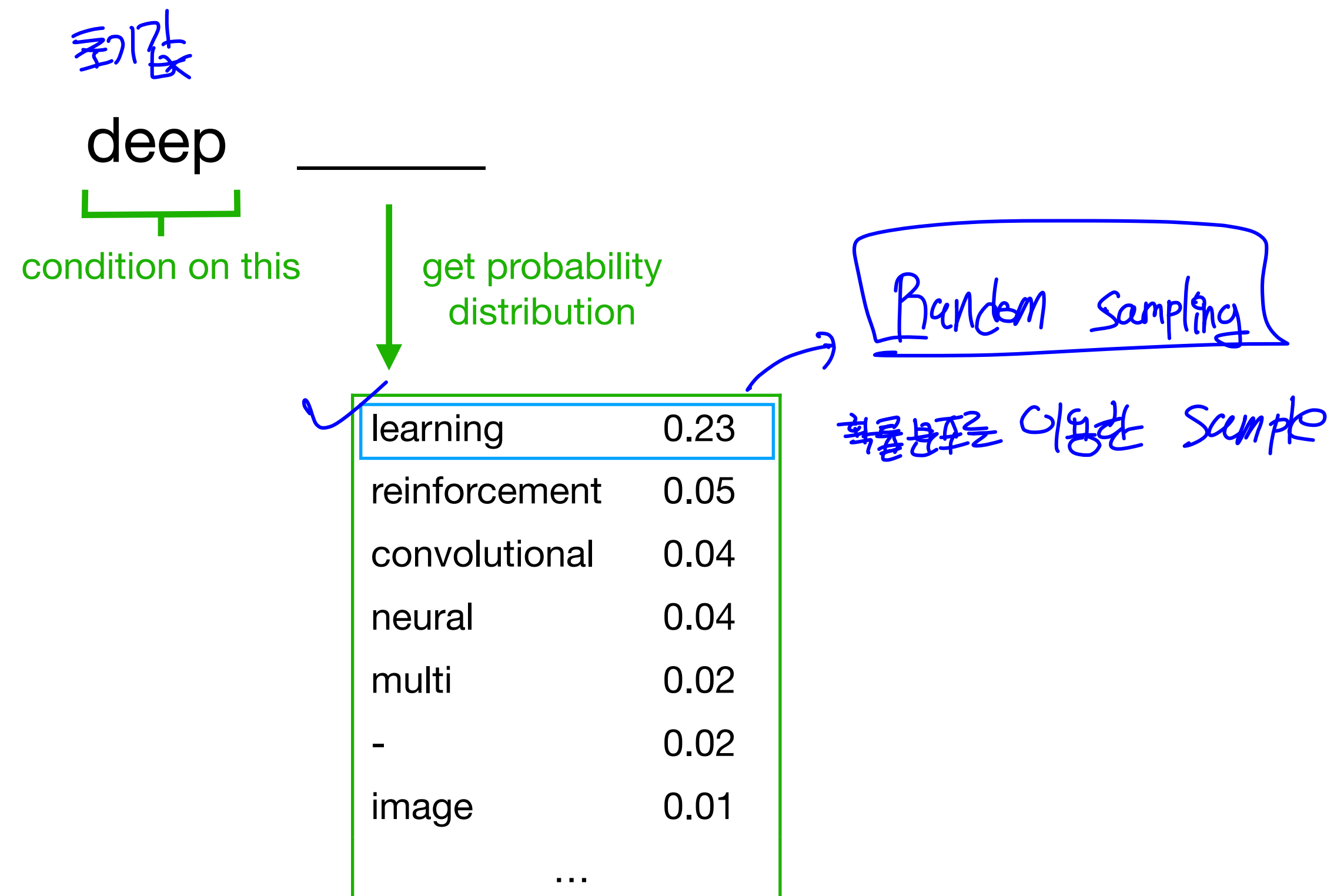
$$p(w^{(i)} | \text{The cat slept quietly on the}) \approx \frac{\text{count(quietly on the } w^{(i)})}{\text{count(quietly on the)}} \quad \begin{matrix} \text{table의 개수} \\ = 36 \\ 100 \end{matrix}$$

예)  $\frac{\text{아름다운} - \text{코스모스}}{100\text{번}} = \frac{20}{100} \%$

$\frac{\text{바라} - \text{10번}}{100\text{번}} = \frac{10}{100} \%$

지금까지 count기반

# N-gram Language Model (Text Generation 3-gram)





# N-gram Language Model (Text Generation 3-gram)

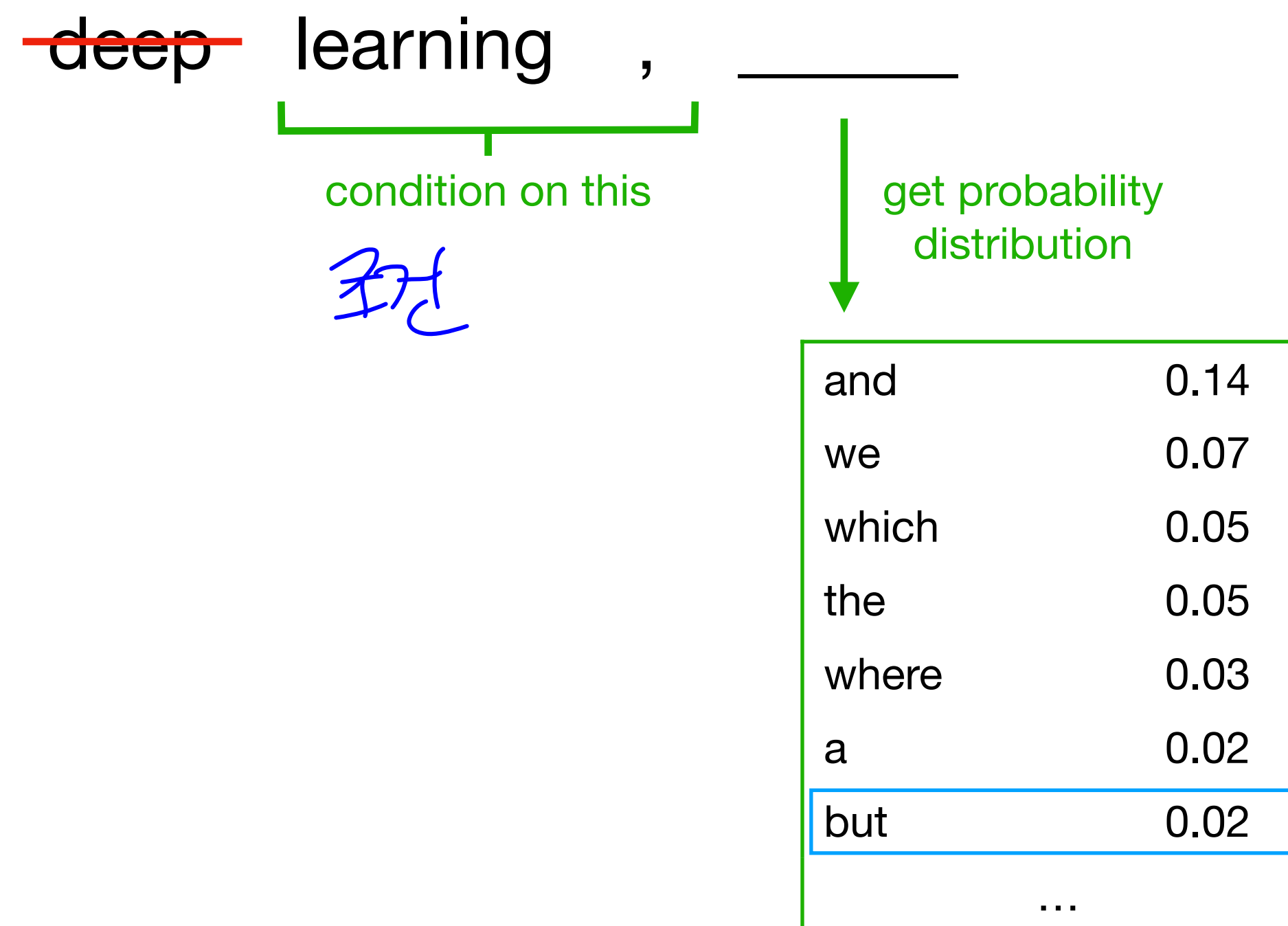
deep learning

condition on this

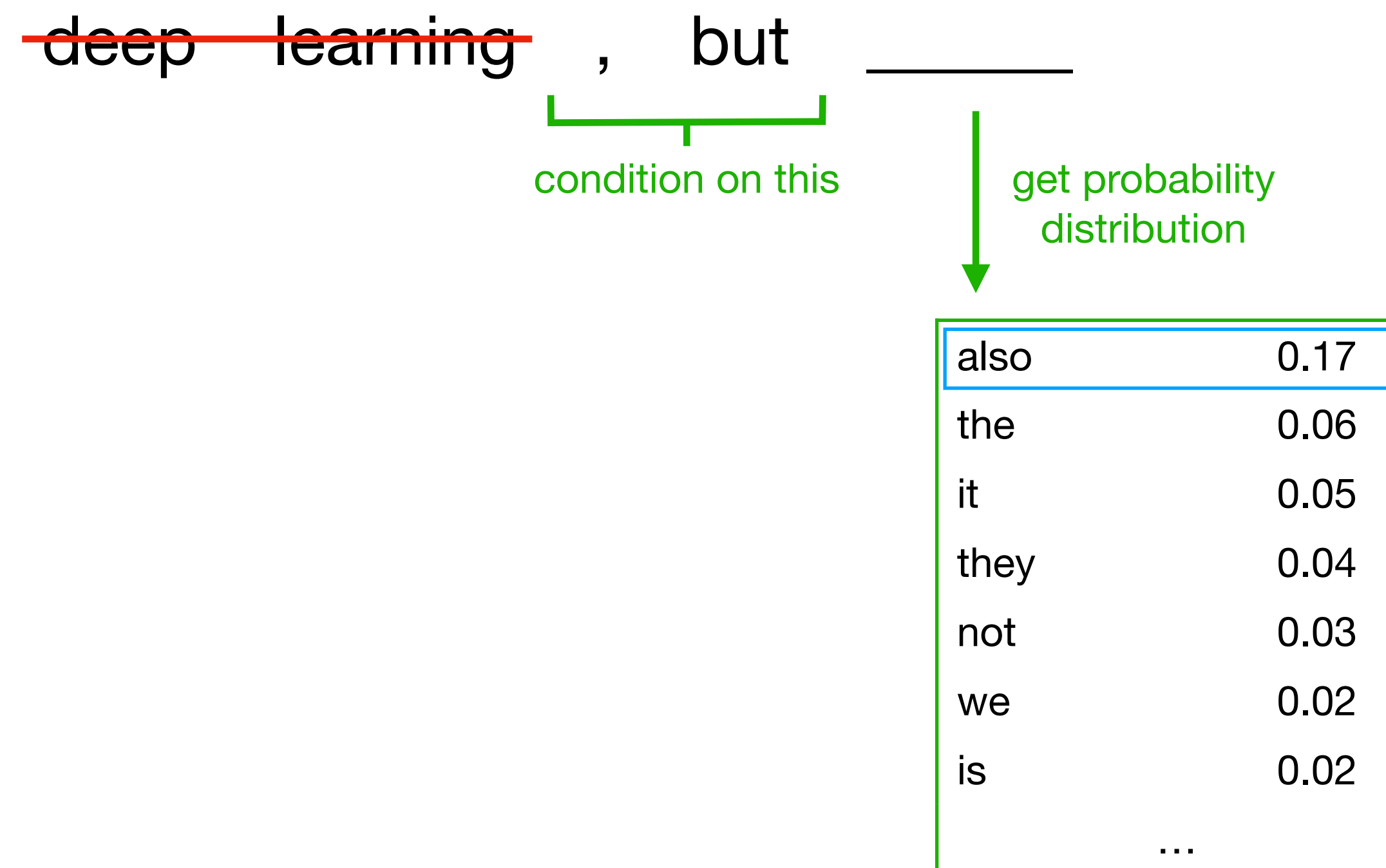
get probability distribution

models	0.06
.	0.05
;	0.05
based	0.05
,	0.05
for	0.04
methods	0.04
...	

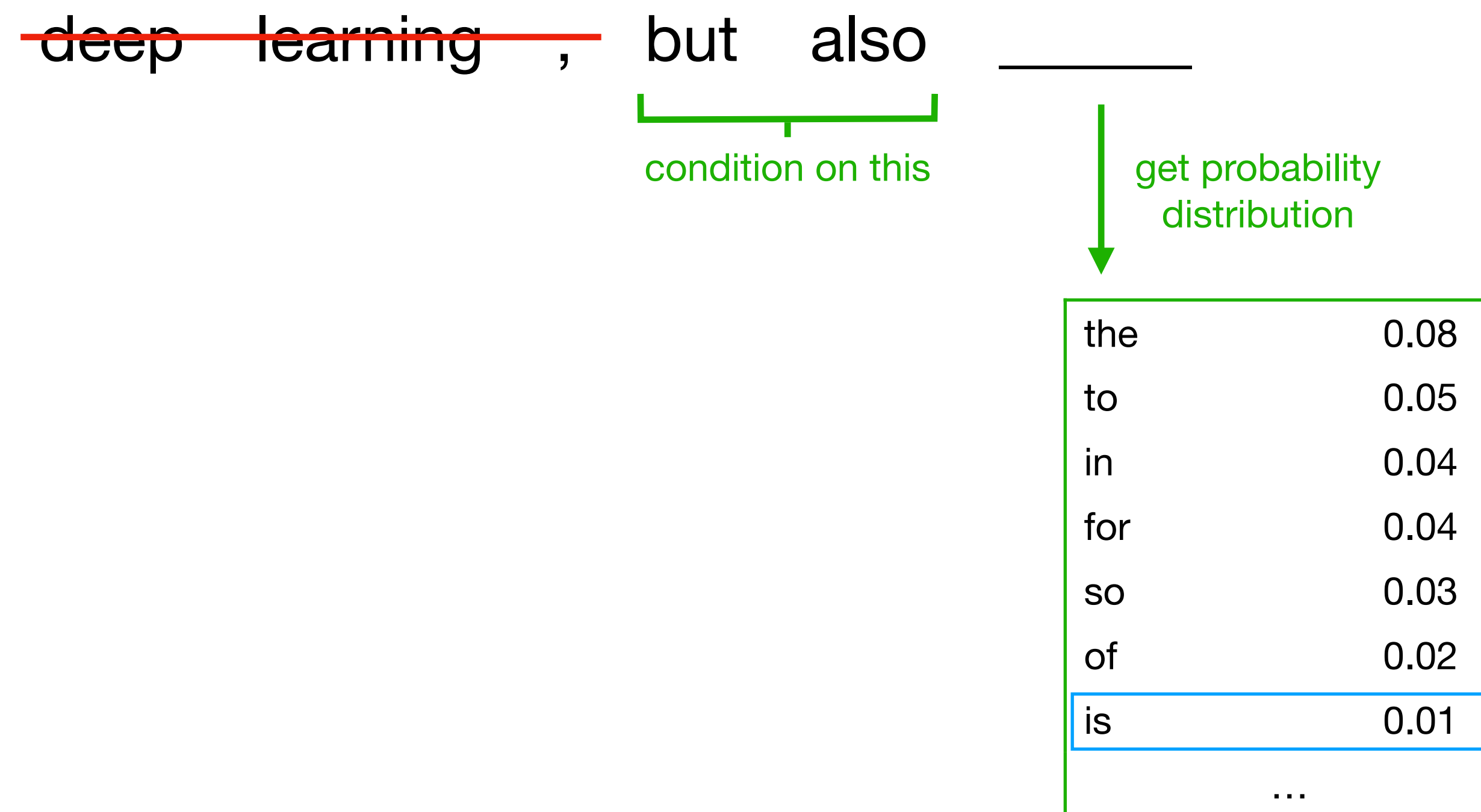
# N-gram Language Model (Text Generation 3-gram)



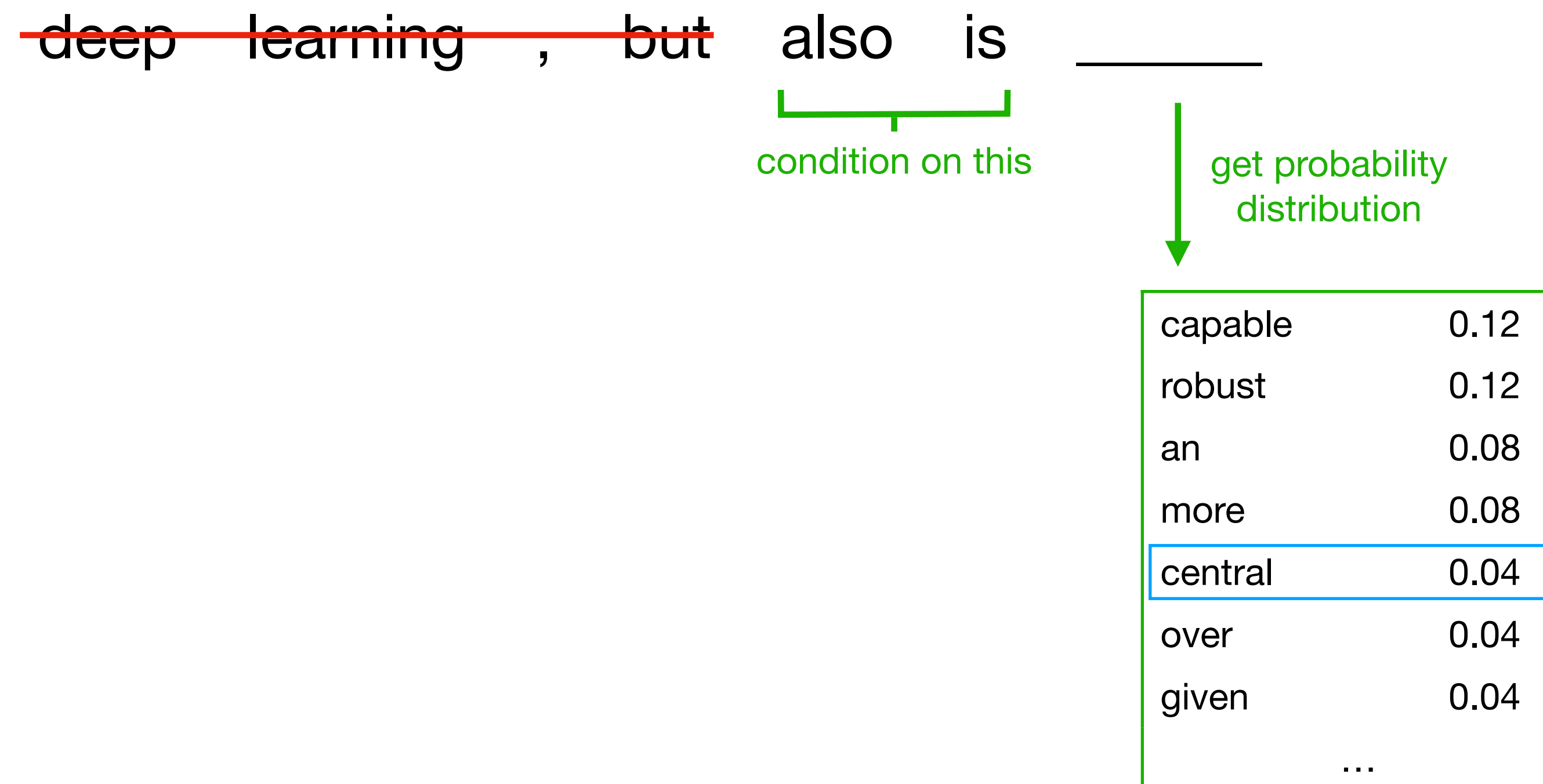
# N-gram Language Model (Text Generation 3-gram)



# N-gram Language Model (Text Generation 3-gram)



# N-gram Language Model (Text Generation 3-gram)



# N-gram Language Model (Text Generation 3-gram)

deep learning , but also is central to human. performance . however , using structural similarity index measure than other partitioned sampling schemes , while making the approach with empirical data has the effect of phonetics has received little attention within the context of information on ...

내용의 일관성이 전혀 없음

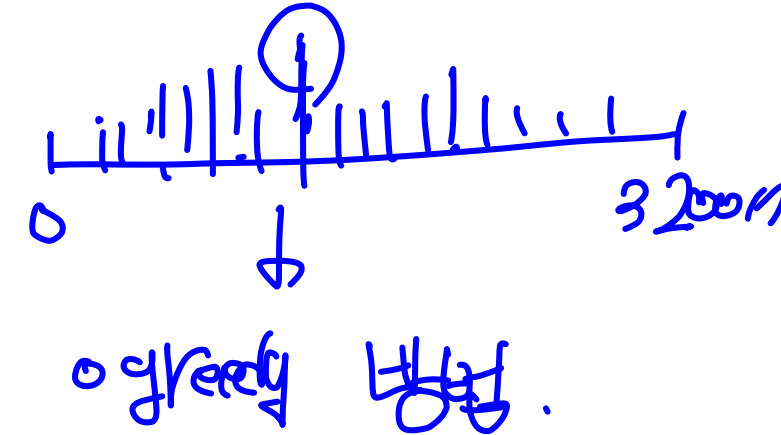
# Neural Language Model (Fixed Window)

단어의 분포

Output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

$\hat{y} \Rightarrow 32.00\%$ 가 일어나는 확률



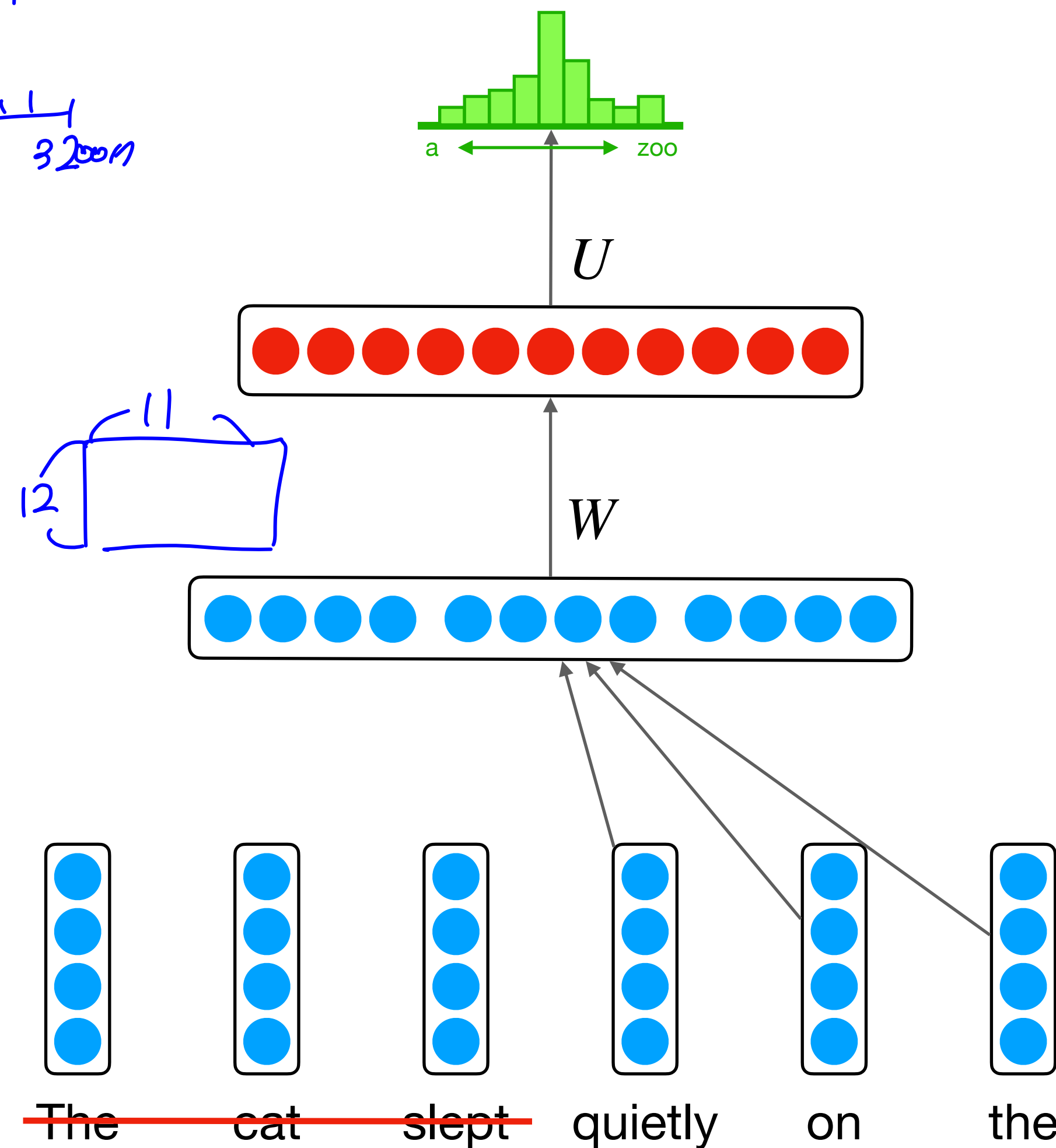
Hidden layer

$$h = f(Wx + b_1)$$

Concatenate word Embedding

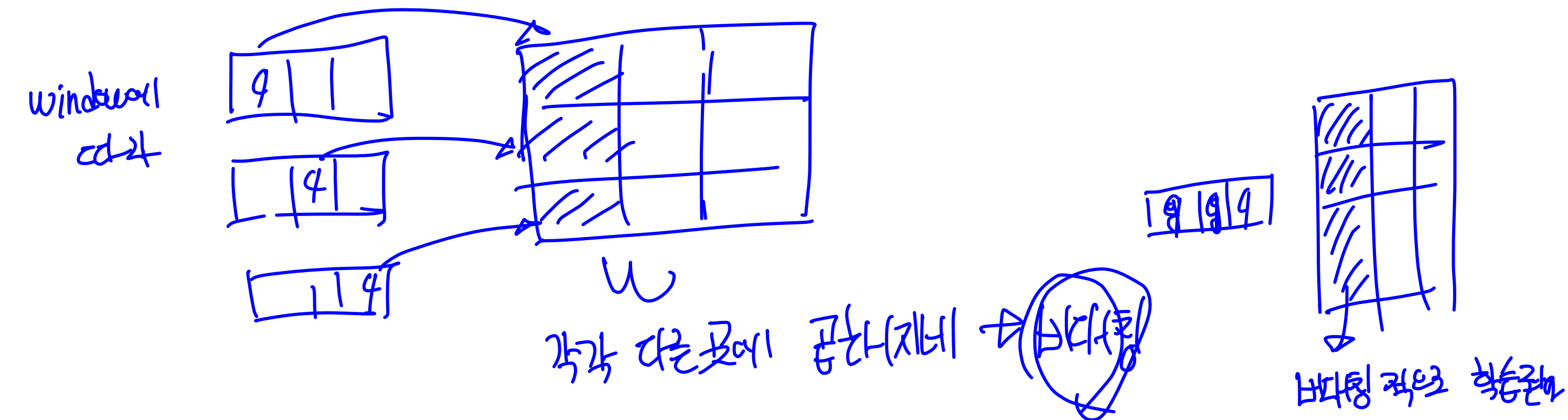
$$x = (x^{(i-3)}; x^{(i-2)}; x^{(i-1)})$$

Word Embedding

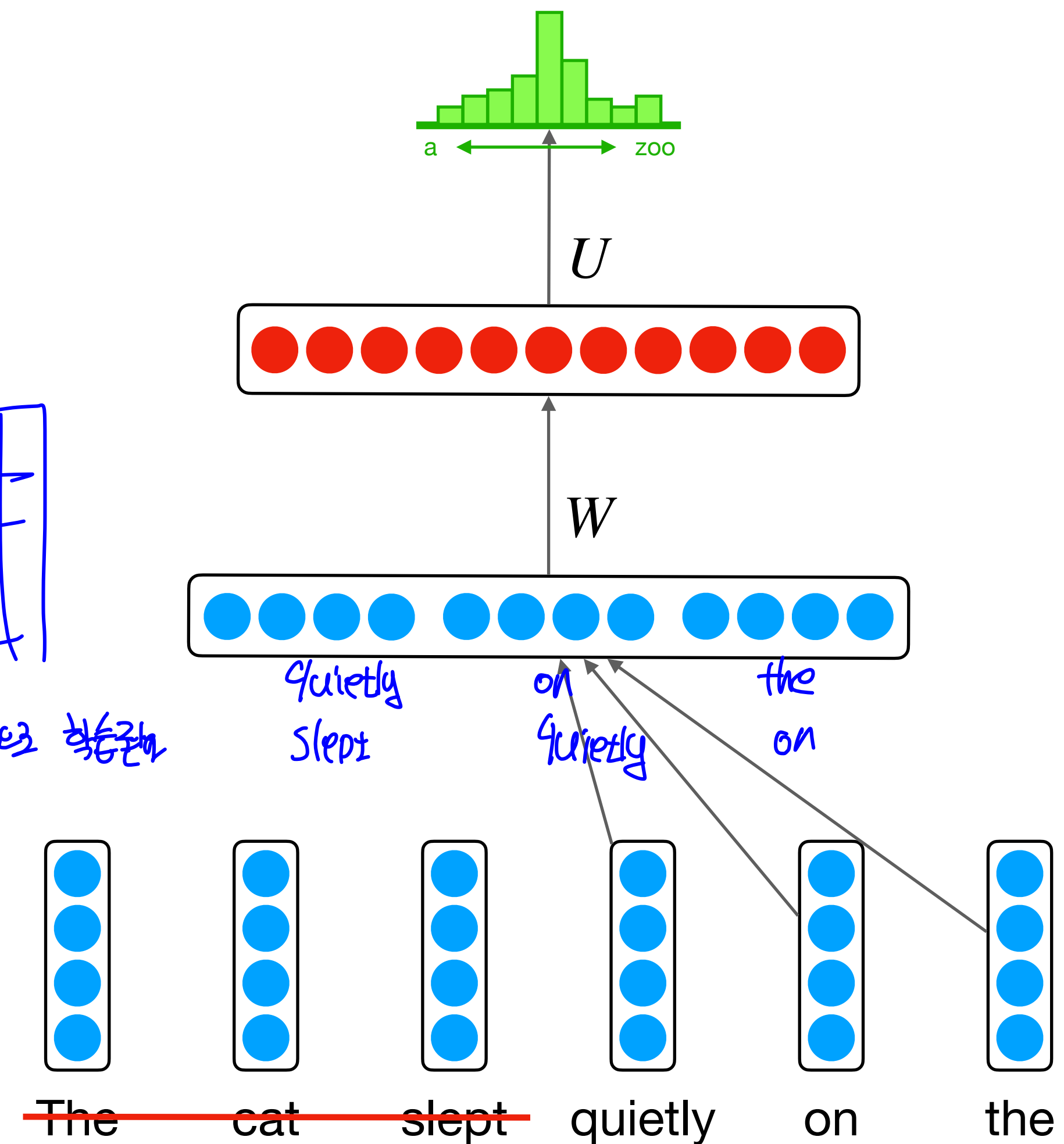


# Neural Language Model (Fixed Window)

- 고정된 Window는 자연어를 처리하는데 크기가 부족함
- $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ 은 window 위치에 따라 다른 weight를 사용 함 (비 대칭)



길이에 상관없이 처리 가능한  
Neural Network가 필요 함





# Neural Language Model (RNN)

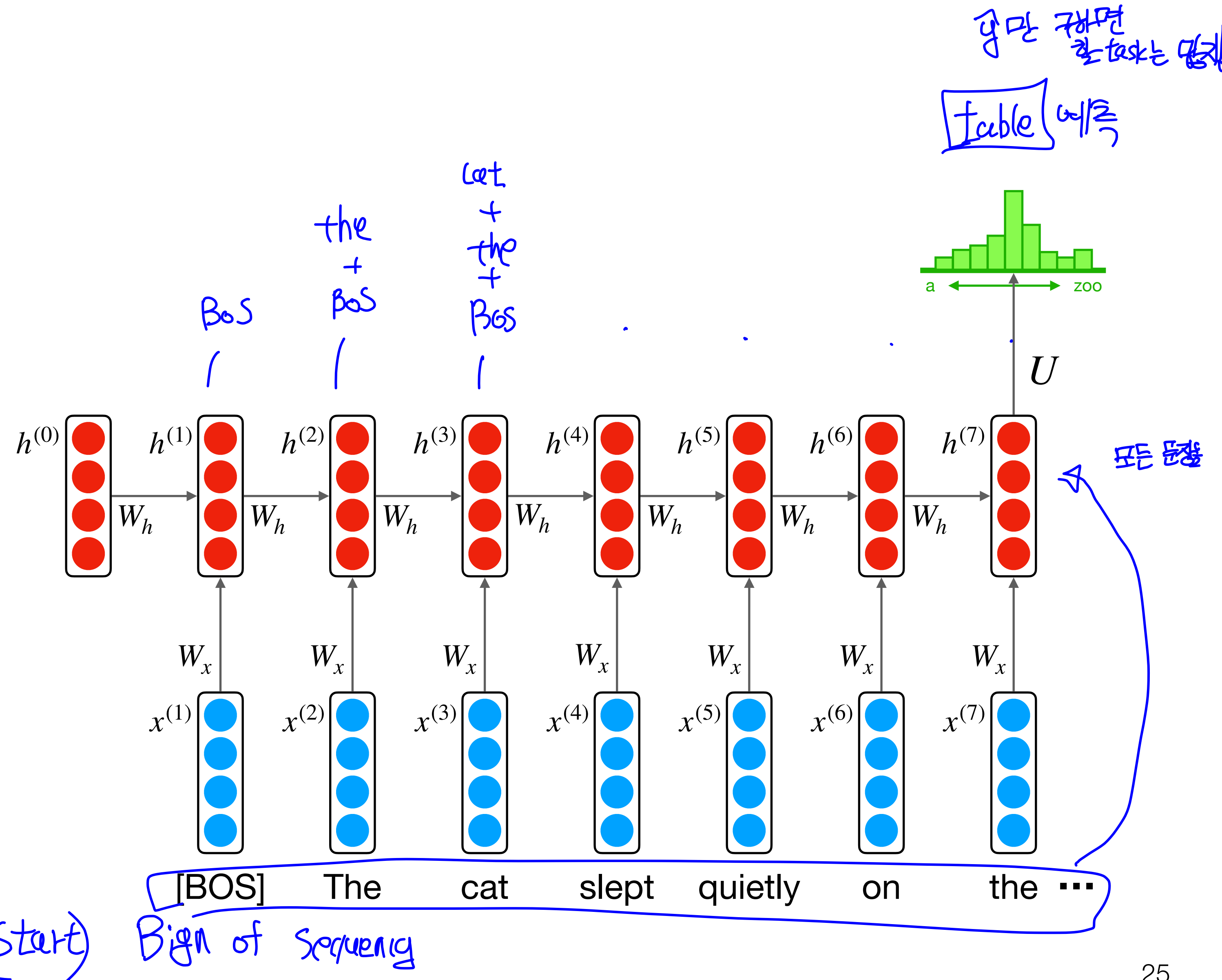
Output distribution

$$\hat{y} = \text{softmax}(Uh + b_2) \in \mathbb{R}^{|V|}$$

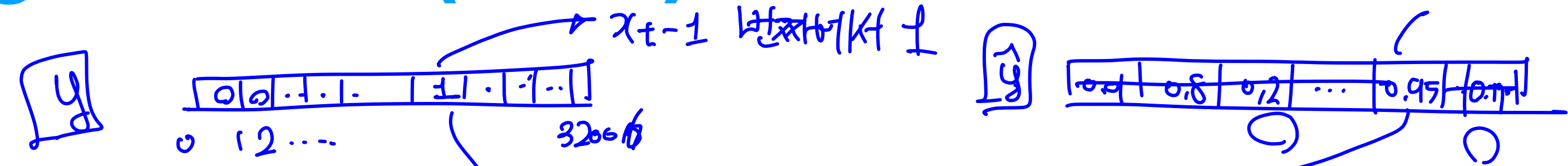
Hidden state

$$h^{(t)} = \tanh(W_h h^{(t-1)} + W_x x^{(t)} + b_1)$$

Word Embedding



# Neural Language Model (RNN)



$$J^{(t)}(\theta) = \underbrace{CE}_{\text{one-hot}}(\underbrace{y^{(t)}}_{\text{true label}}, \underbrace{\hat{y}^{(t)}}_{\text{predicted label}}) = - \sum_{w \in V} \underbrace{y_w^{(t)}}_{\text{one-hot}} \log(\underbrace{\hat{y}_w^{(t)}}_{\text{predicted probability}}) = - \log \hat{y}_{x_{t+1}}^{(t)}$$

$\hat{y}_{x_{t+1}}^{(t)}$  다음 단어를 예측하는 확률을 계산해 줍니다.

$\propto \log \hat{y}_{x_{t+1}}^{(t)}$

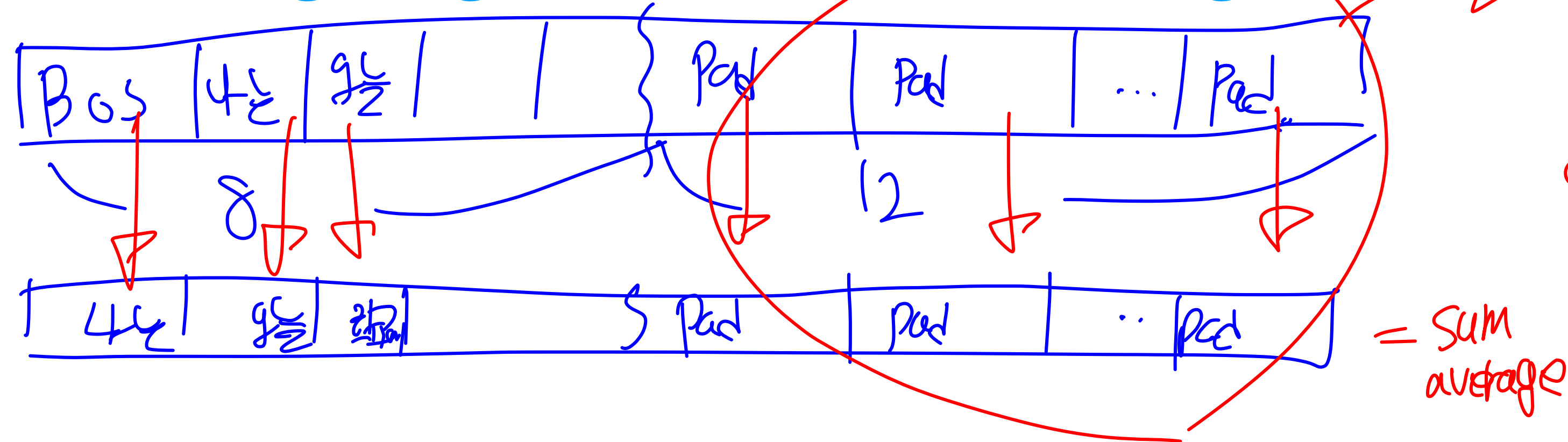
$\frac{1}{\alpha} - \log \hat{y}_{x_{t+1}}^{(t)}$

$$\textcircled{J(\theta)} = \frac{1}{T} \sum_{t=1}^T \underbrace{-\log \hat{y}_{x_{t+1}}^{(t)}}_{\text{negative log likelihood를 줄여줍니다.}}$$

= MLE

= Negative log likelihood

# Neural Language Model (Training)



End of sequence

Labels →	The	cat	slept	quietly	on	the		.	[EOS]
Inputs →	[BOS]	The	cat	slept	quietly	on	the		.

다음 단어를 예측하는 데이터를 만들면 되지.

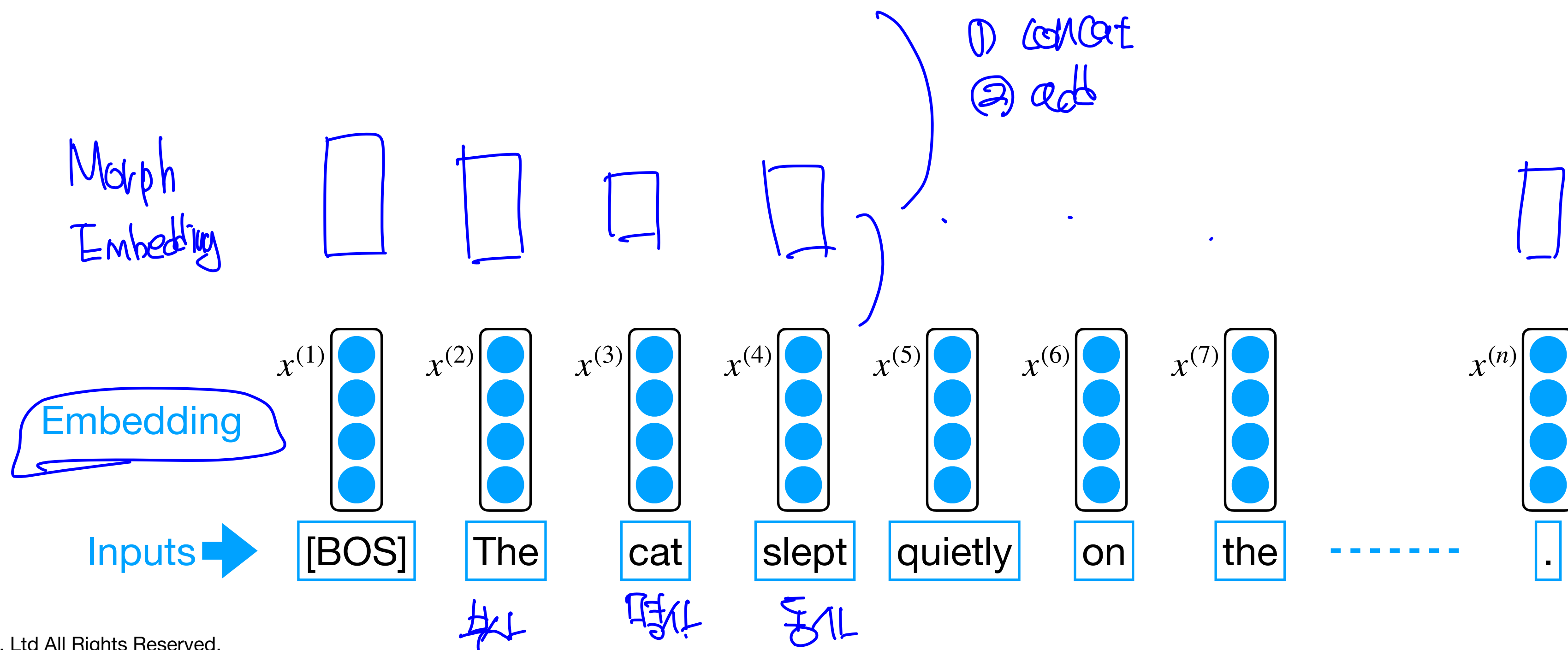
# Neural Language Model (Training)

Labels → The cat slept quietly on the ..... . [EOS]

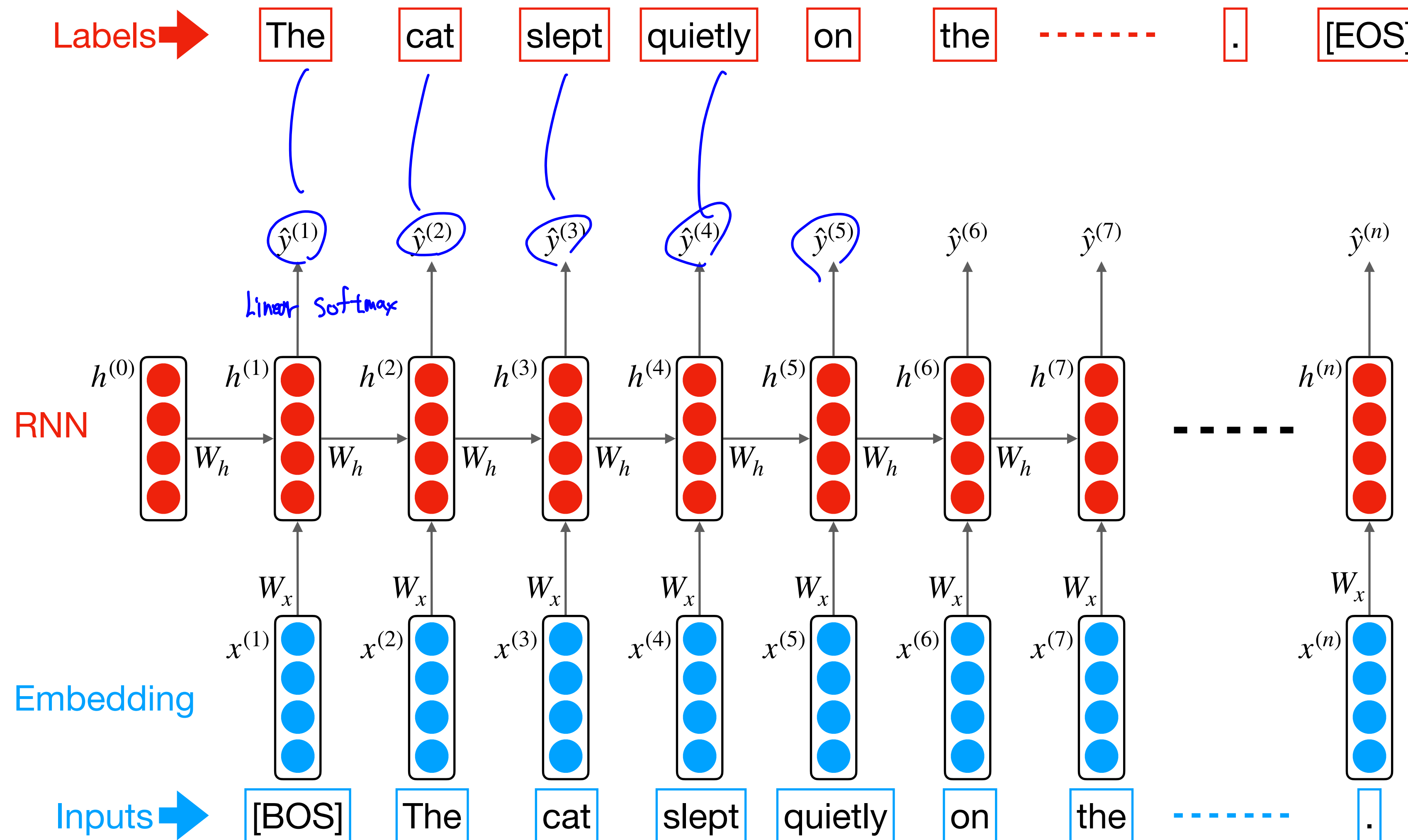
Inputs → [BOS] The cat slept quietly on the ..... .

# Neural Language Model (Training)

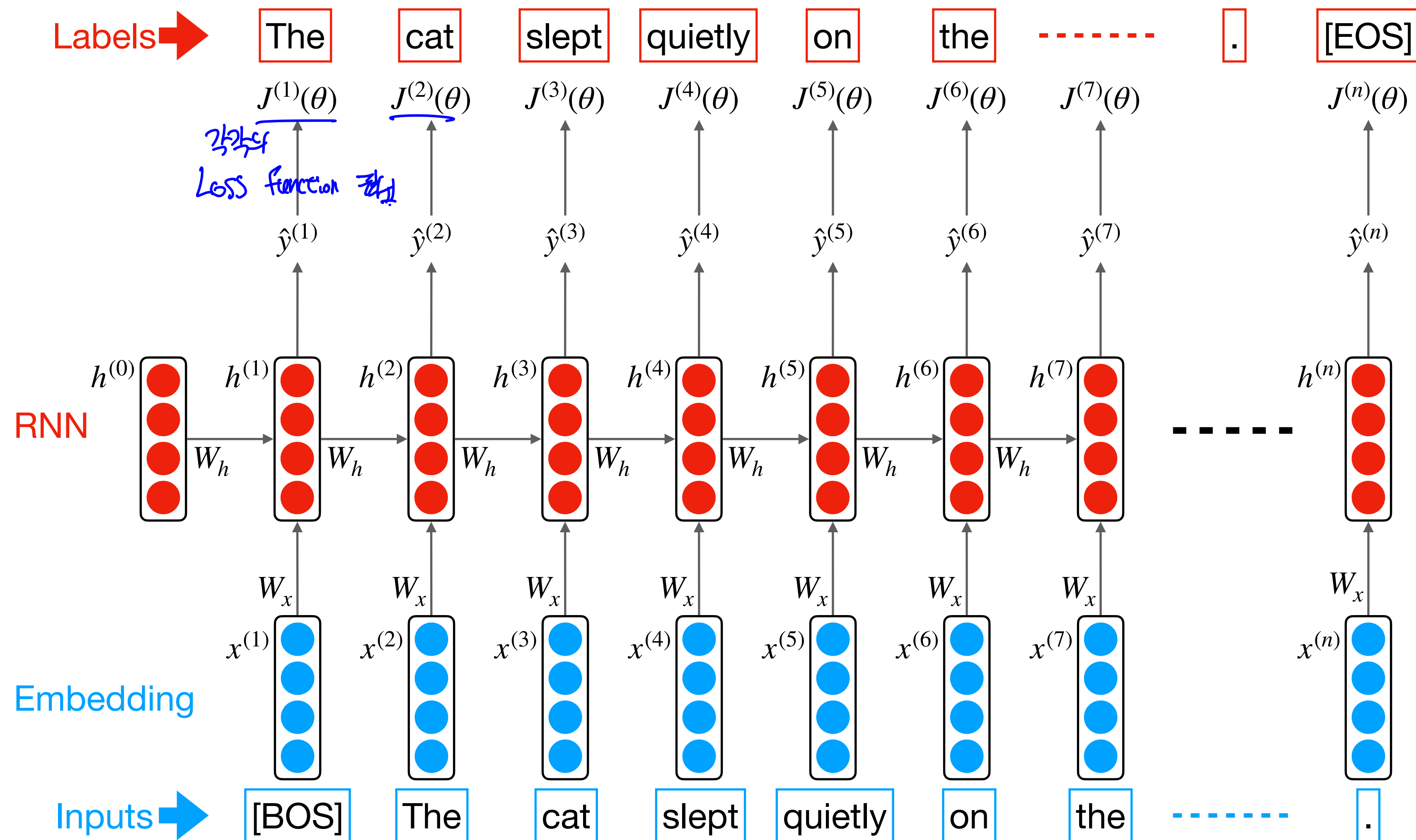
Labels → The cat slept quietly on the ..... . [EOS]



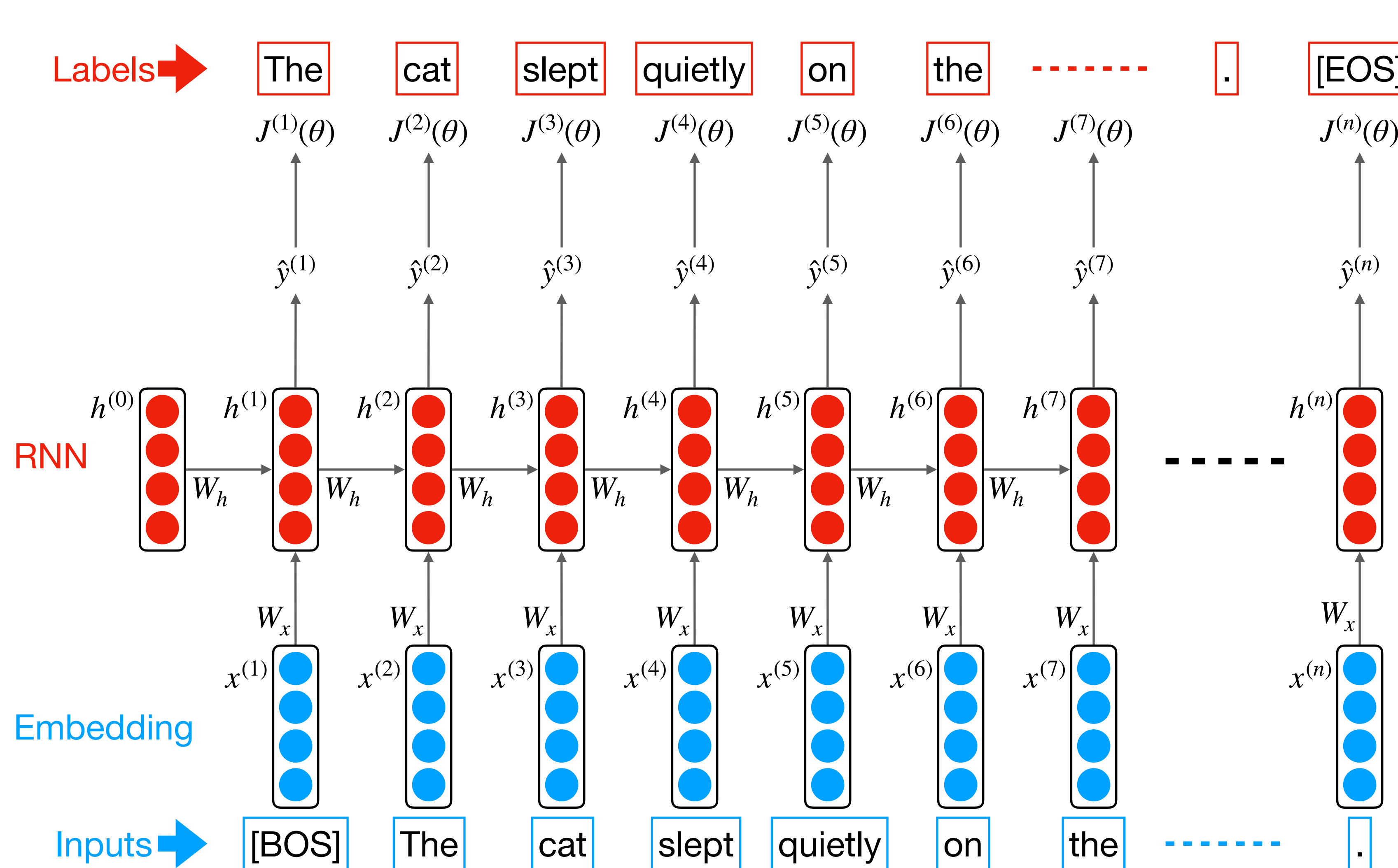
# Neural Language Model (Training)



# Neural Language Model (Training)



# Neural Language Model (Training)



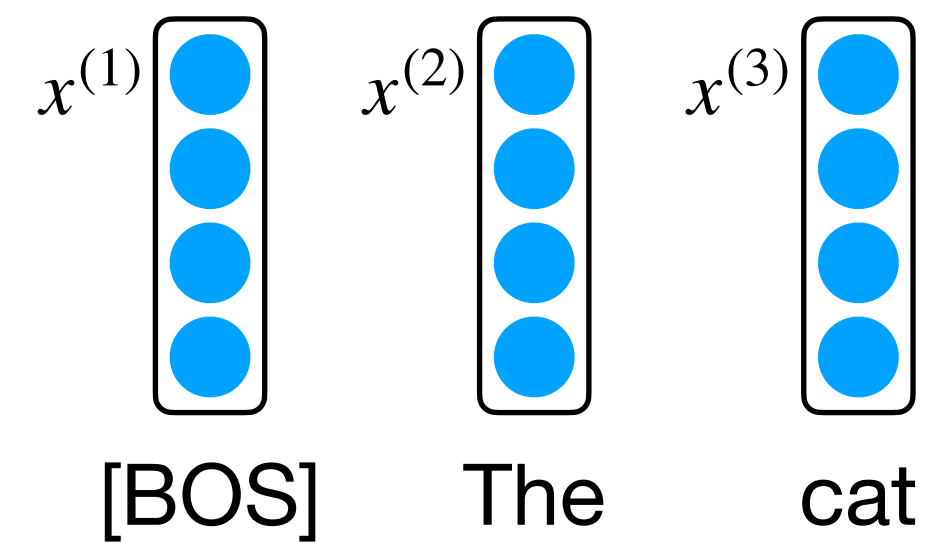
$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^t(\theta)$$

Mean of  
negative log prob

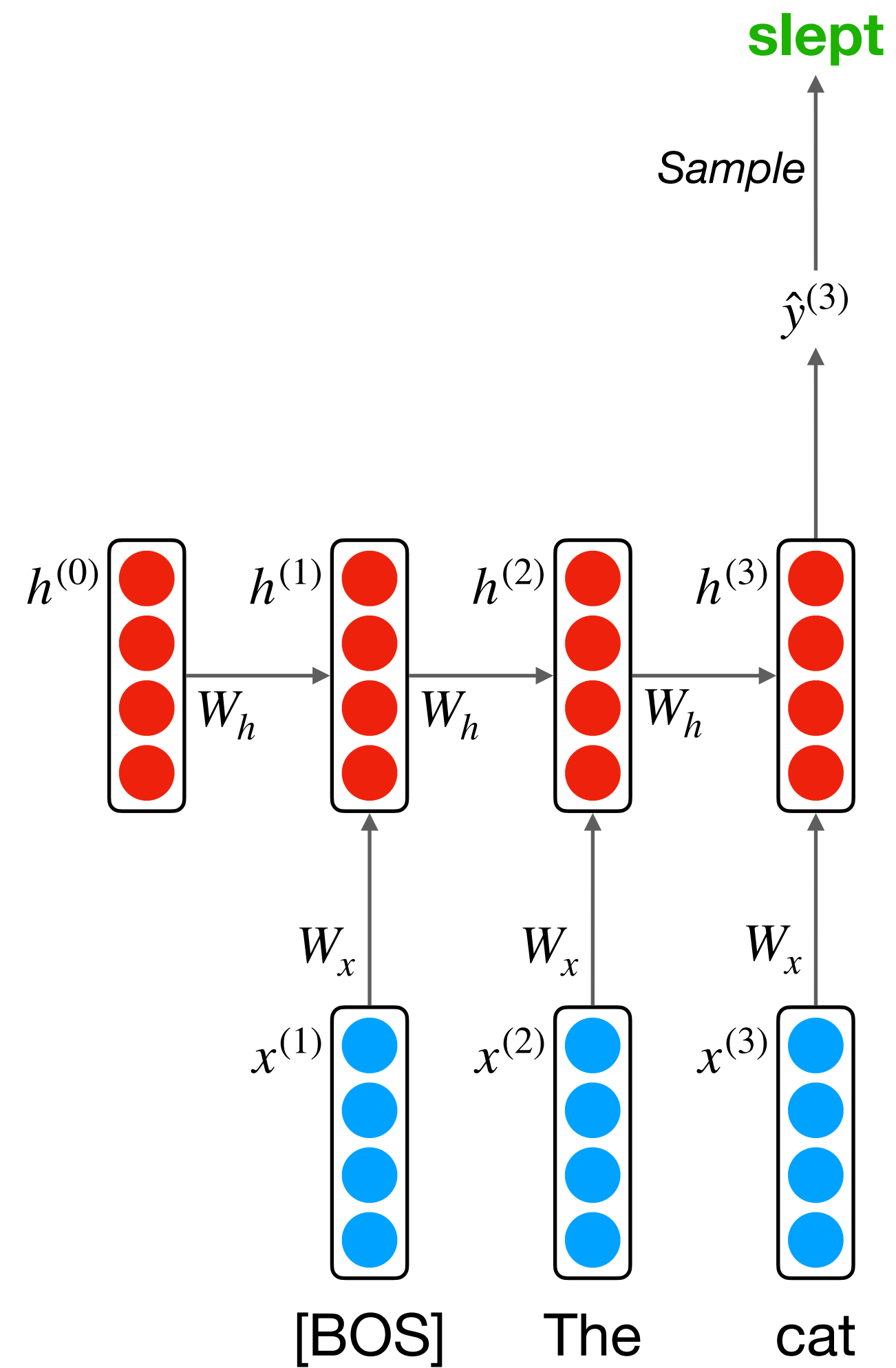
$J = \frac{1}{T} \sum_{t=1}^T J^t(\theta)$



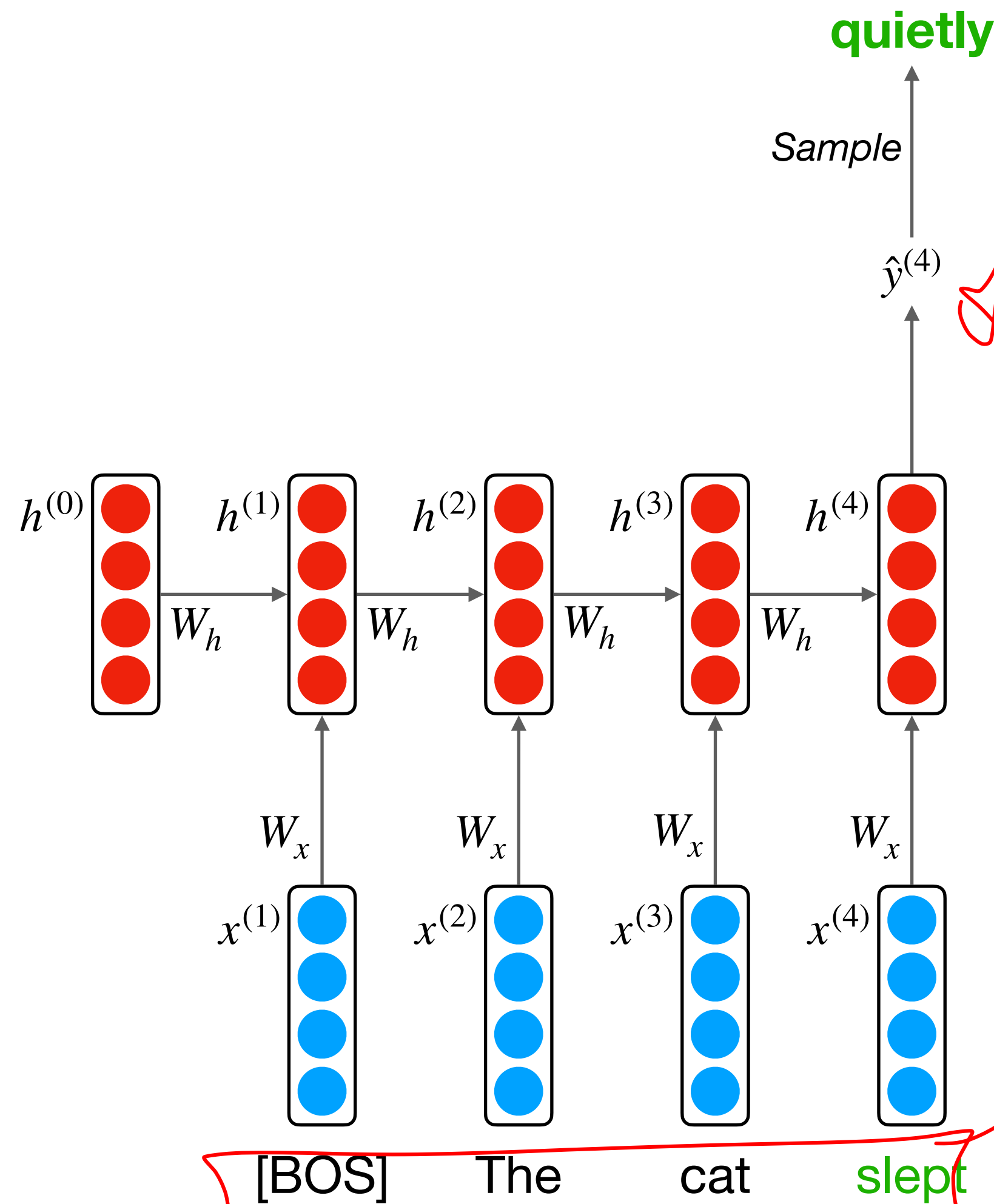
# Neural Language Model (Generate Text)



# Neural Language Model (Generate Text)



# Neural Language Model (Generate Text)



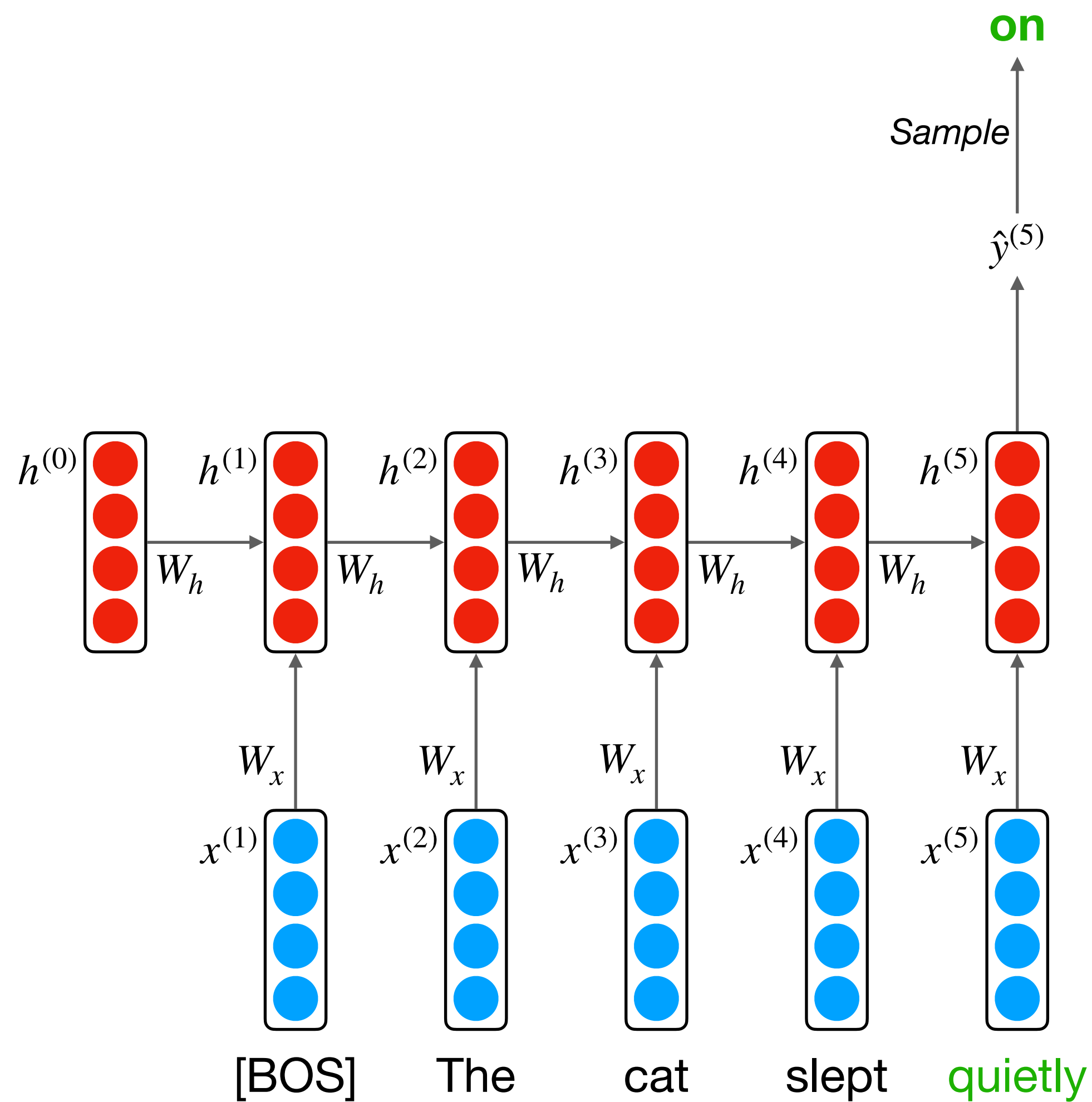
다음 단어 예측하도록 함

→ 이진수에 Sample 출력 넣어서

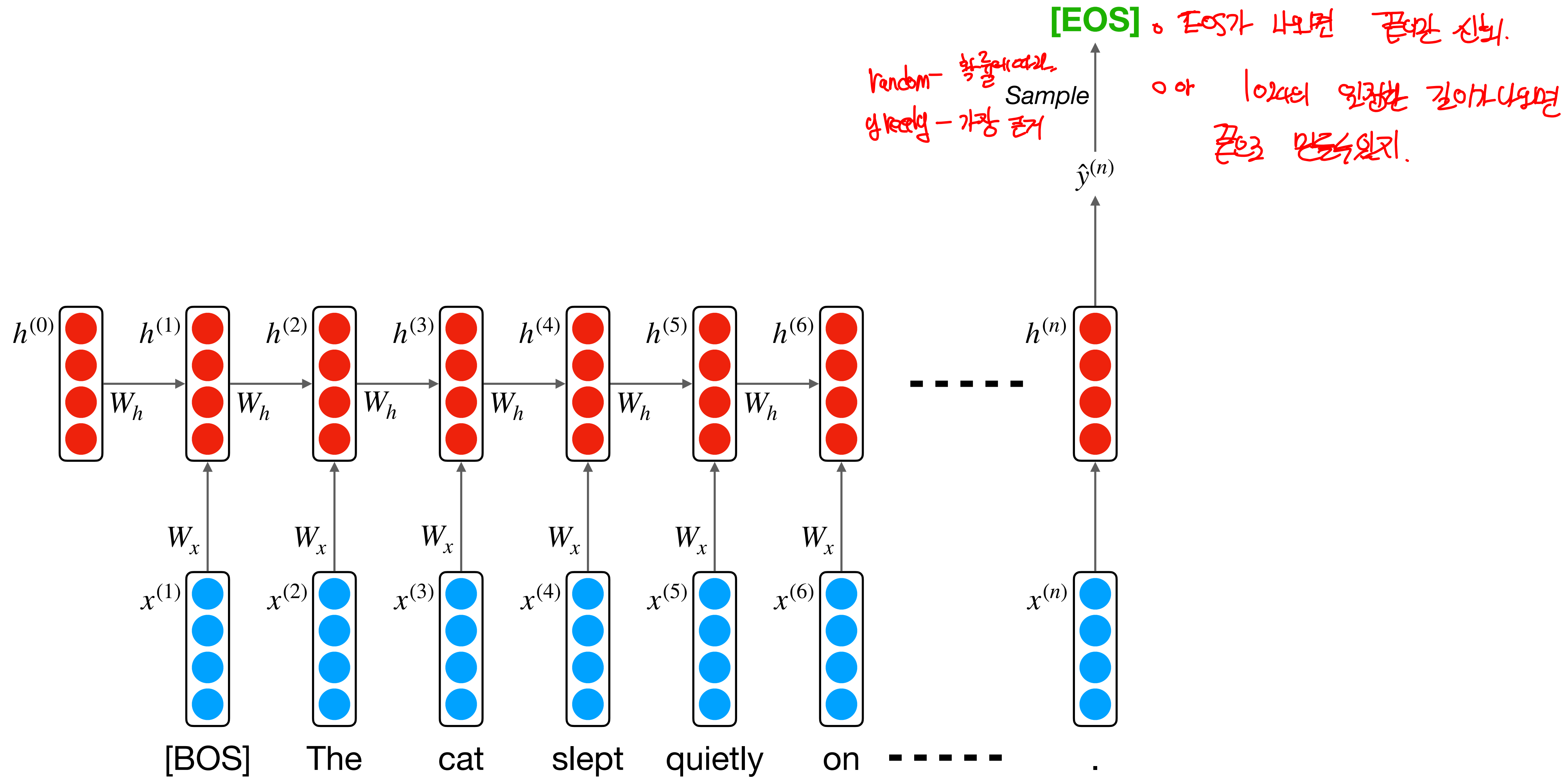
⇒ 이때 random sample 인거고

입력 4개.

# Neural Language Model (Generate Text)



# Neural Language Model (Generate Text)



# Neural Language Model (Subcomponent)

Pretrained Language model  
Bert,

Task,

- Predicting Typing ○ 자동완성
- Speech recognition ○ 음성인식
- Handwriting recognition ○ 글자 확인
- Spelling/grammar correction ○ 문법 수정
- Authorship identification ○ 작가 확인
- Machine Translation
- Summarization
- Dialog
- etc.

이들



# Neural Language Model (Metric)

모든 문장에서

$$\text{perplexity} = \prod_{t=1}^T \left( \frac{1}{P_{LM}(x^{(t+1)} | x^{(1)}, x^{(2)}, \dots, x^{(t)})} \right)^{\frac{1}{T}}$$

↑

Normalize by number of words  
○ T번 곱해나가  
[ ]  $\frac{1}{T}$  해서 normalize

Inverse probability of corpus  
역수.

(t)  
x<sub>t+1</sub>

# Neural Language Model (Metric)

$$\begin{aligned}
 \text{perplexity} &= \prod_{t=1}^T \left( \frac{1}{P_{LM}(x^{(t+1)} | x^{(1)}, x^{(2)}, \dots, x^{(t)})} \right)^{\frac{1}{T}} \\
 &= \prod_{t=1}^T \left( \frac{1}{\hat{y}_{x_{t+1}}^{(t)}} \right)^{\frac{1}{T}} = \exp \left( \log \left( \prod_{t=1}^T \left( \frac{1}{\hat{y}_{x_{t+1}}^{(t)}} \right)^{\frac{1}{T}} \right) \right) = \exp \left( \frac{1}{T} \sum_{t=1}^T -\log \hat{y}_{x_{t+1}}^{(t)} \right) \\
 &= \exp(J(\theta))
 \end{aligned}$$

Handwritten notes in blue ink:
 

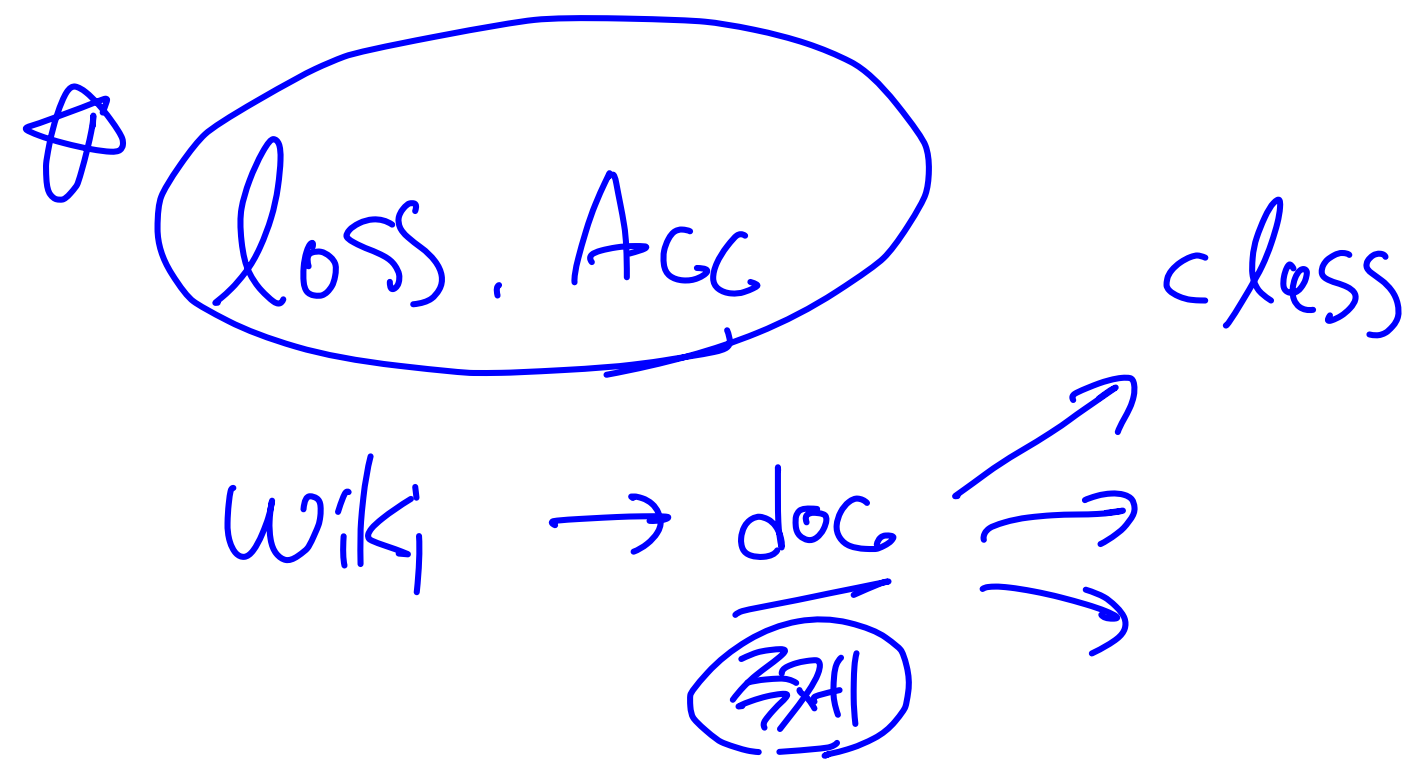
- A blue arrow points from the denominator  $P_{LM}(x^{(t+1)} | \dots)$  in the first equation to  $\hat{y}_{x_{t+1}}^{(t)}$  in the second equation.
- Below the  $\exp$  and  $\log$  in the third equation, there is a note:  $\frac{1}{24(32)}$  and  $\times 1 \text{ or } 2$ .
- Below the  $J(\theta)$  in the final equation, there is a note:  $\text{목적함수}$  (objective function).

**Low perplexity is better !!!**

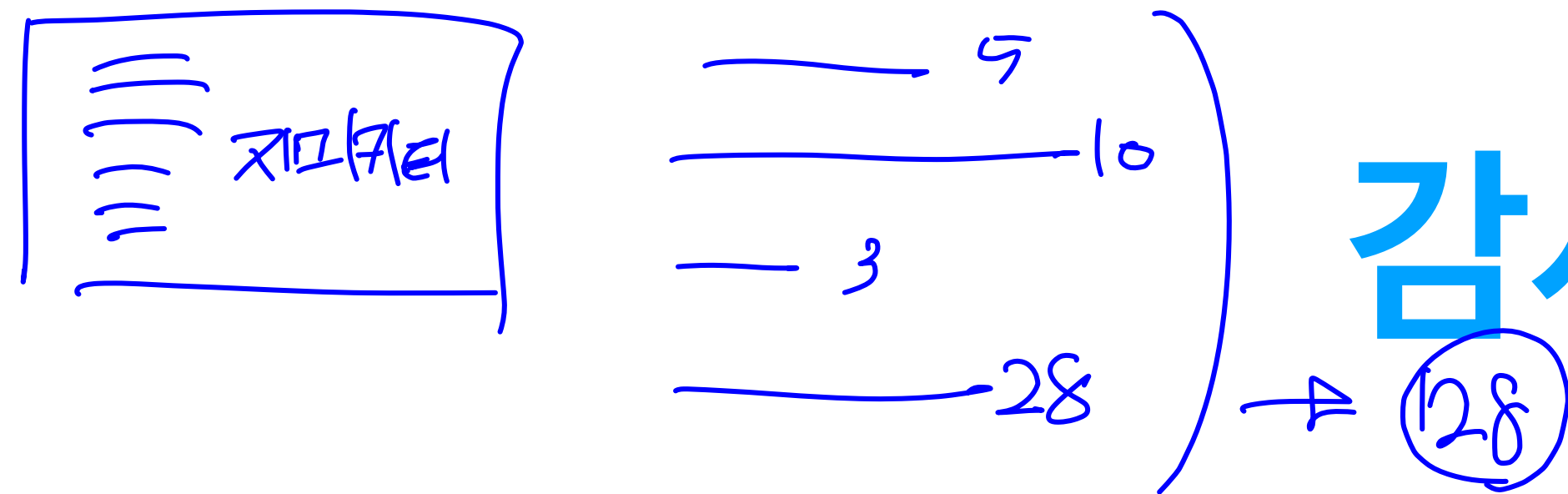
이것은 값이 목적함수가 좋다.  
= loss를 줄이는게 좋다.



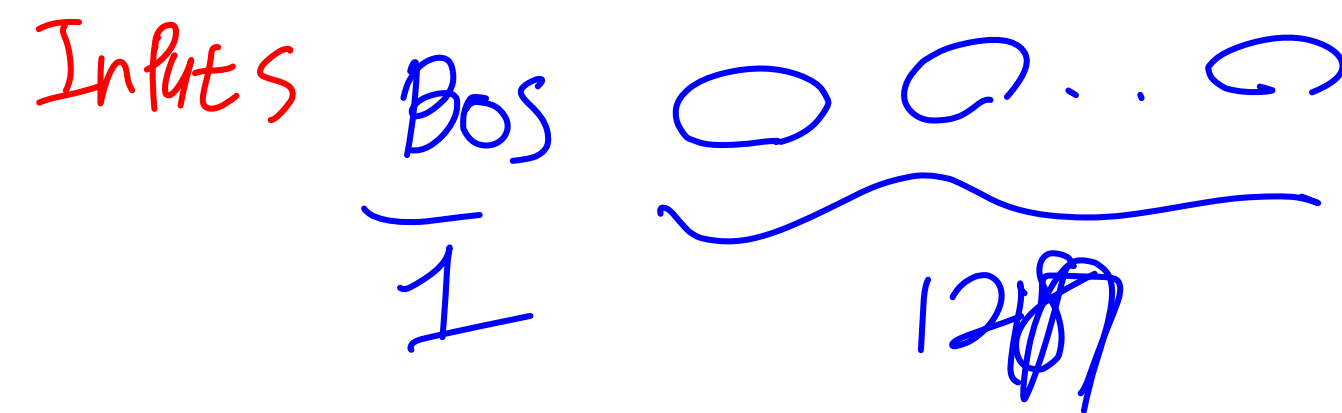
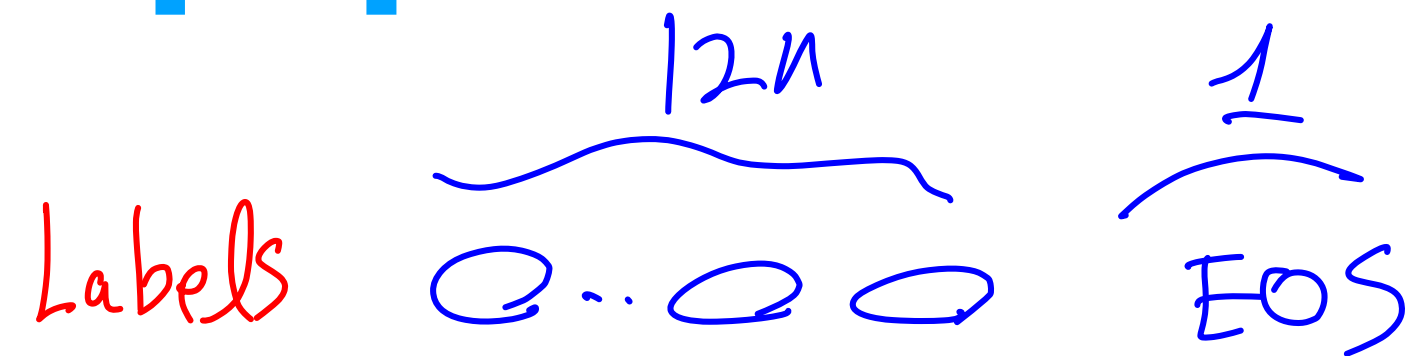
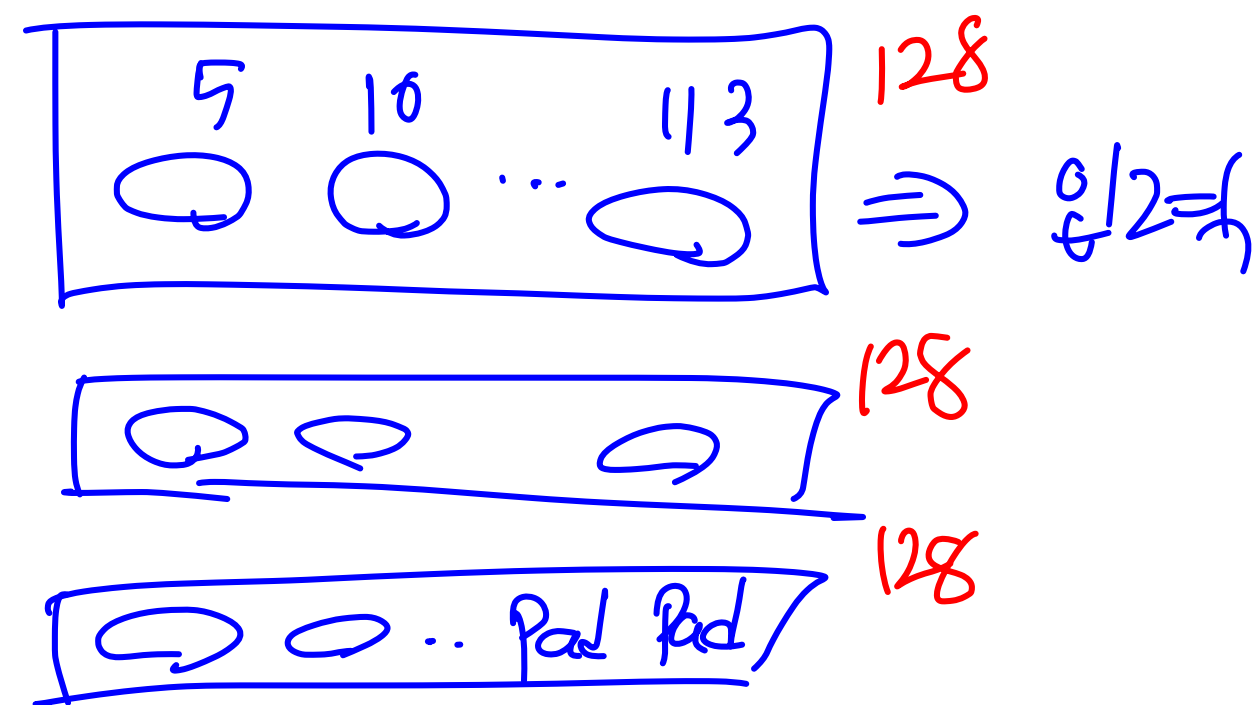
**감사합니다.**



⇒ 학습용 데이터



감사합니다.



하루 4시간

다들 스터디 팀이

하루 4시간

다들 스터디 팀이

3중 평행선과 2중 평행선

$$\text{match-size} = 1024$$

16PU 128

4 Capacity limit

big mad/cy #2 water full

→ 한글자씩 자르고 tokenizer

- subword tokenize → 정량화 도표

BP E

→ 중 ~~중~~ 지표수 분류

— Morpheme

9. 유망한 기업

○ 일대다 누정방지 서풍을 밝힌다(고려).

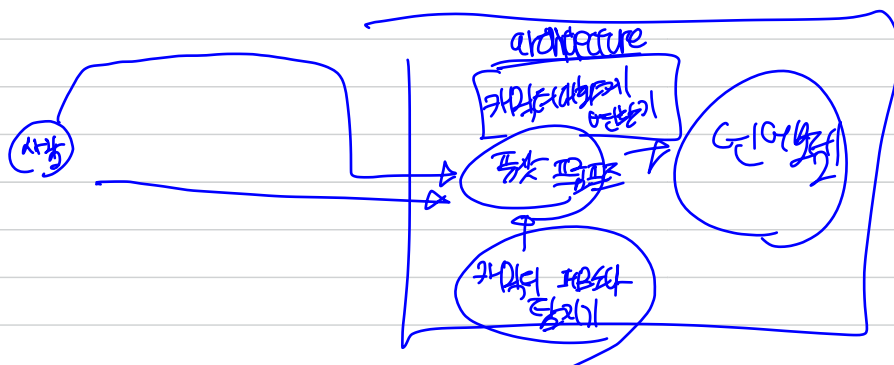
Bule score 22 확률 미정함.

9. Vocabulary 측정치는 perplexity로 측정된 이 모델의

9 real vs false .

→ 알파인 카운터 테이퍼 + 테이퍼링

Q 카탈루냐 분리



P-turning

## Prompt-tuning

Prompt - Control - Unit