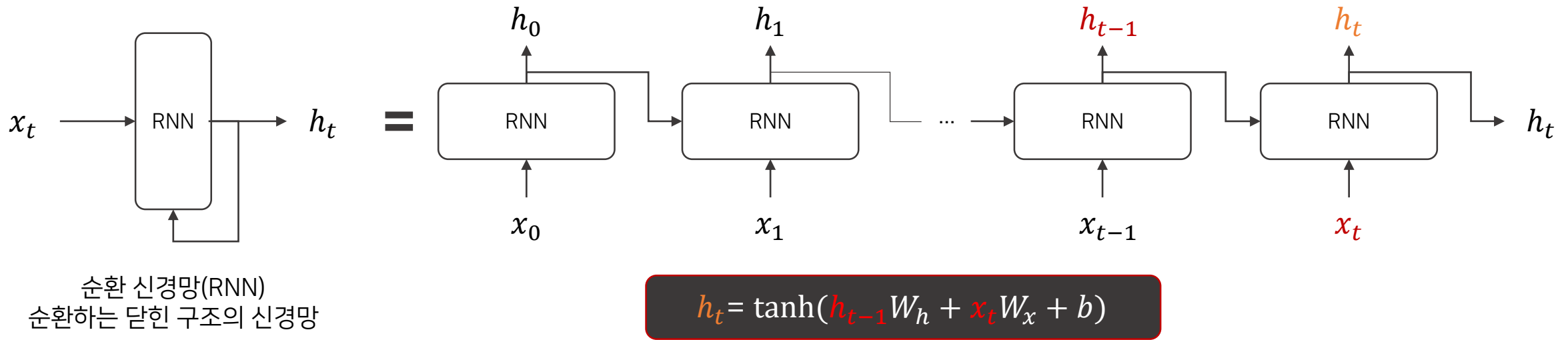




Attention Mechanism

Recurrent Neural Network (RNN)



$$\begin{aligned} h_0 &= \tanh(h_{init}W_h + x_0W_x + b) \longrightarrow \\ h_1 &= \tanh(h_0W_h + x_1W_x + b) \longrightarrow \\ &\dots \longrightarrow \\ h_{t-1} &= \tanh(h_{t-2}W_h + x_{t-1}W_x + b) \longrightarrow \\ h_t &= \tanh(h_{t-1}W_h + x_tW_x + b) \end{aligned}$$

Recurrent Neural Network (RNN)

그림 6-2 RNN 계층의 계산 그래프(MatMul 노드는 행렬 곱을 나타냄)

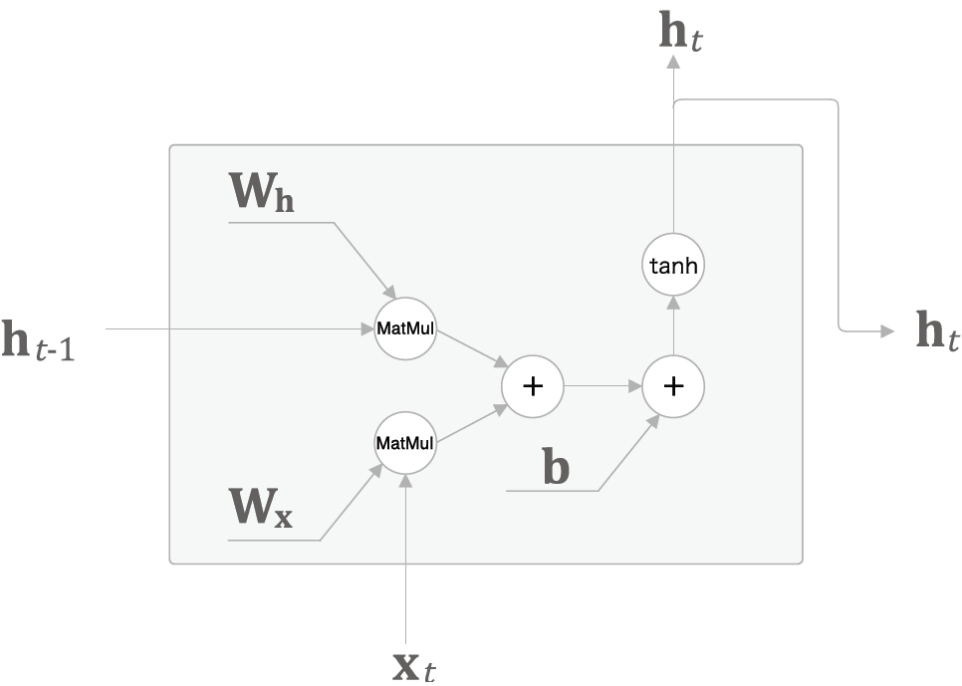
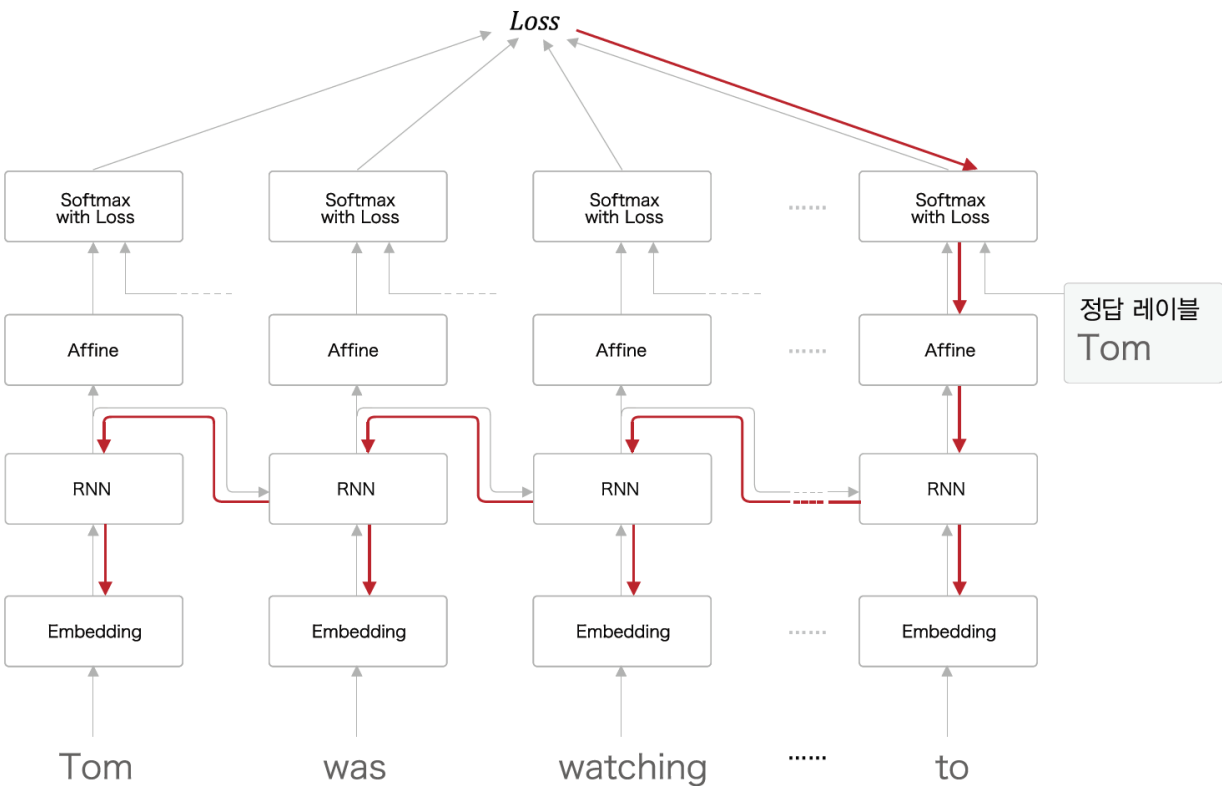


그림 6-4 정답 레이블이 “Tom”임을 학습할 때의 기울기 흐름



Recurrent Neural Network (RNN)

Tom was watching TV in his room. Mary came into the room. Mary said hi to ?

Recurrent Neural Network (RNN)

Tom was watching TV in his room. Mary came into the room. Mary said hi to

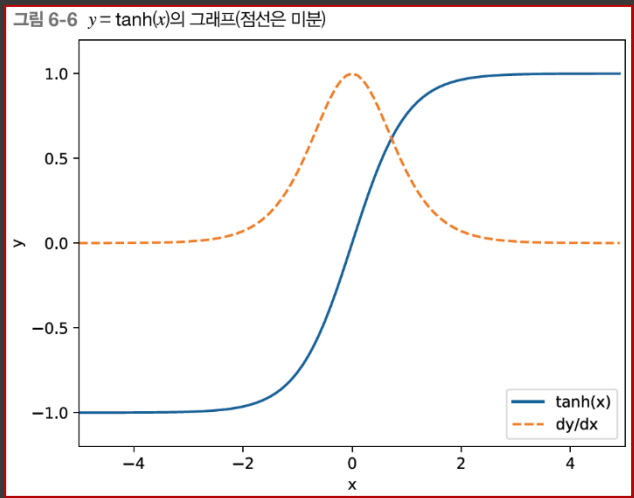
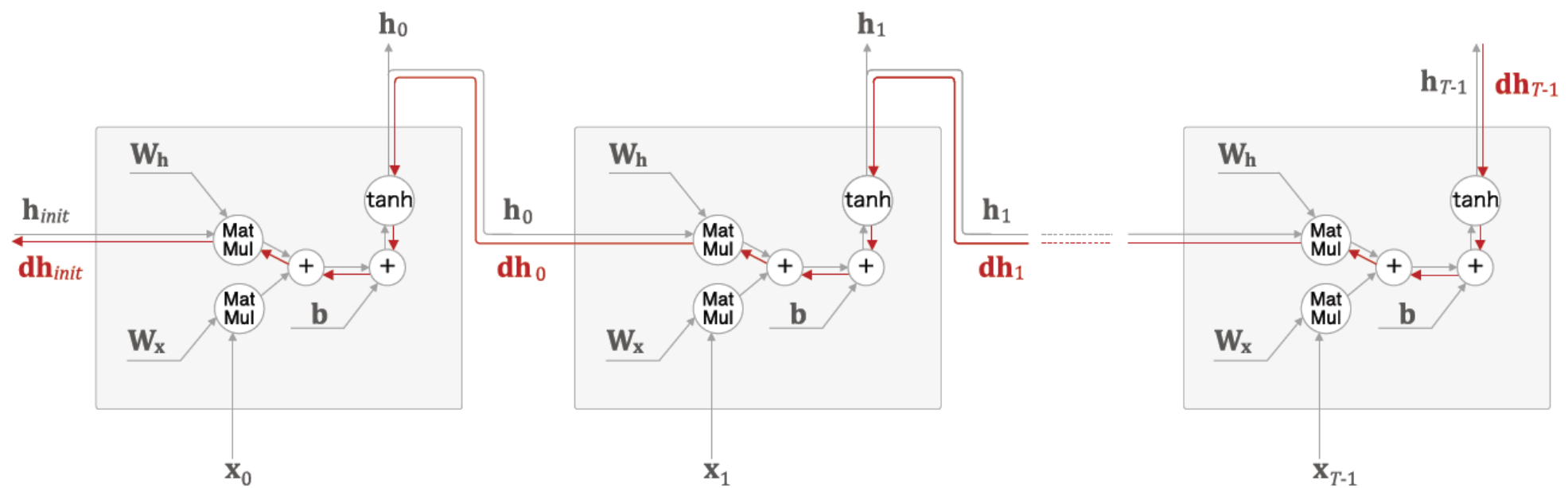
?



장기 의존성 문제(기울기 소실 / 기울기 폭발 문제)

- ✓ 학습 과정에서 기울기를 역전파할 때, Step을 여러 번 거치면서 기울기가 사라지거나 폭발하는 현상
- ✓ 기울기를 통한 가중치의 갱신이 어려움

Recurrent Neural Network (RNN)

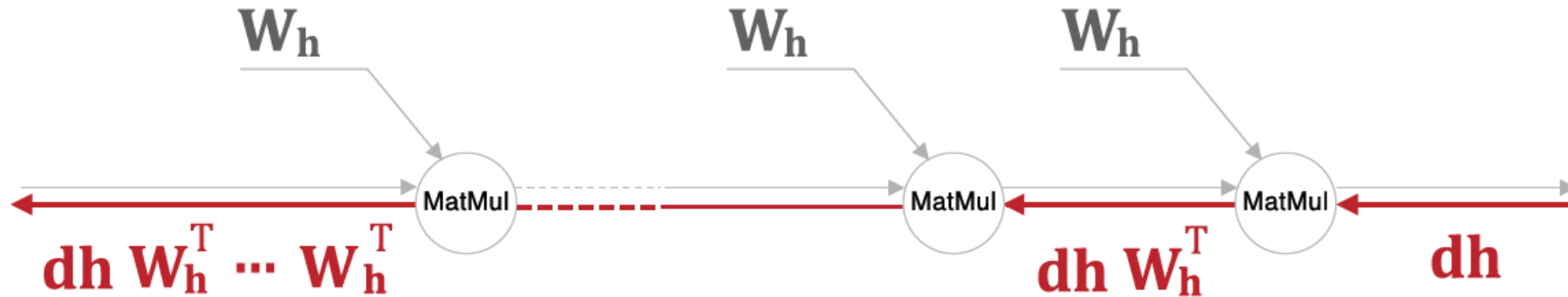


$\tanh(x)$ 에서 기울기 소실 발생

- $\tanh(x)$ 의 미분값은 0이상 1이하의 값
- 역전파에서 문장의 단어 길이(T)만큼 \tanh 를 통과하기 때문에, T 만큼 기울기 값이 계속 감소함
- 기울기가 너무 작아져 소실되면 가중치 매개변수가 더 이상 갱신되지 않음

Recurrent Neural Network (RNN)

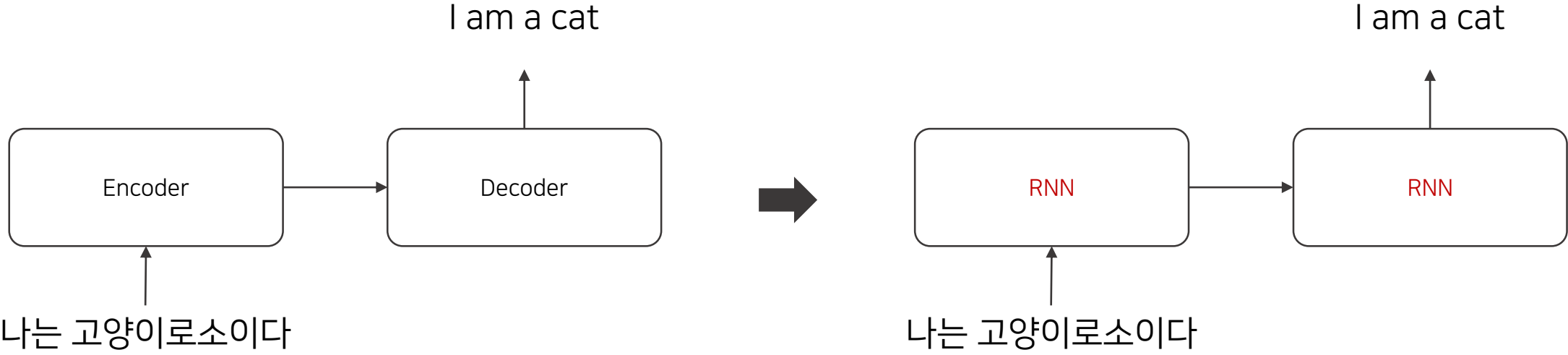
그림 6-7 RNN 계층의 행렬 곱에만 주목했을 때의 역전파의 기울기



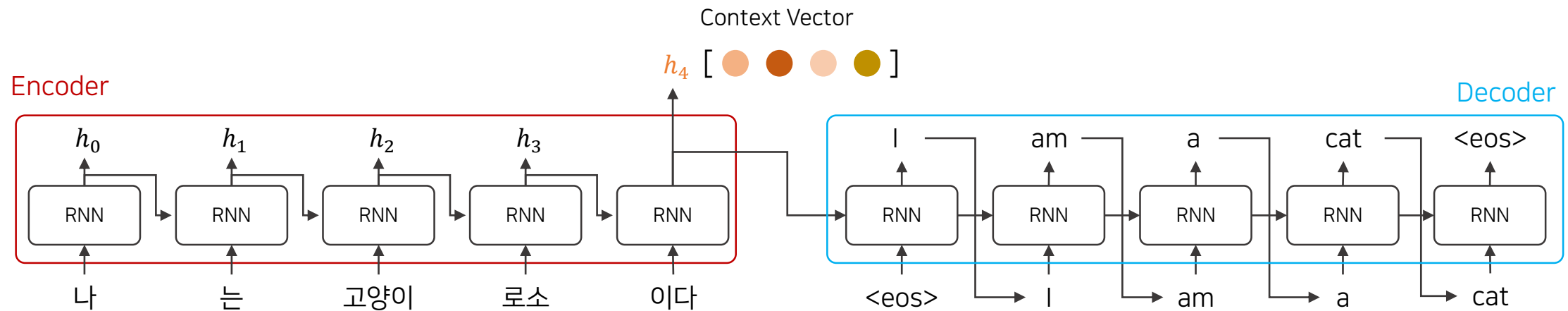
MatMul에서 기울기 소실 혹은 폭발 발생

- MatMul의 역전파는 **전치 행렬의 곱**이며, 단어 개수(T)만큼 똑같은 행렬을 곱하기 때문에 값이 한쪽 방향으로 커지거나 작아짐
- 기울기 폭발: 값이 너무 커지면 오버플로우를 일으켜 NaN같은 오류 발생
- 기울기 소실: 값이 너무 작아지면 매개변수가 더 이상 갱신되지 않아 장기 의존 관계를 학습할 수 없음

Sequence to Sequence (Seq2Seq)



Sequence to Sequence (Seq2Seq)



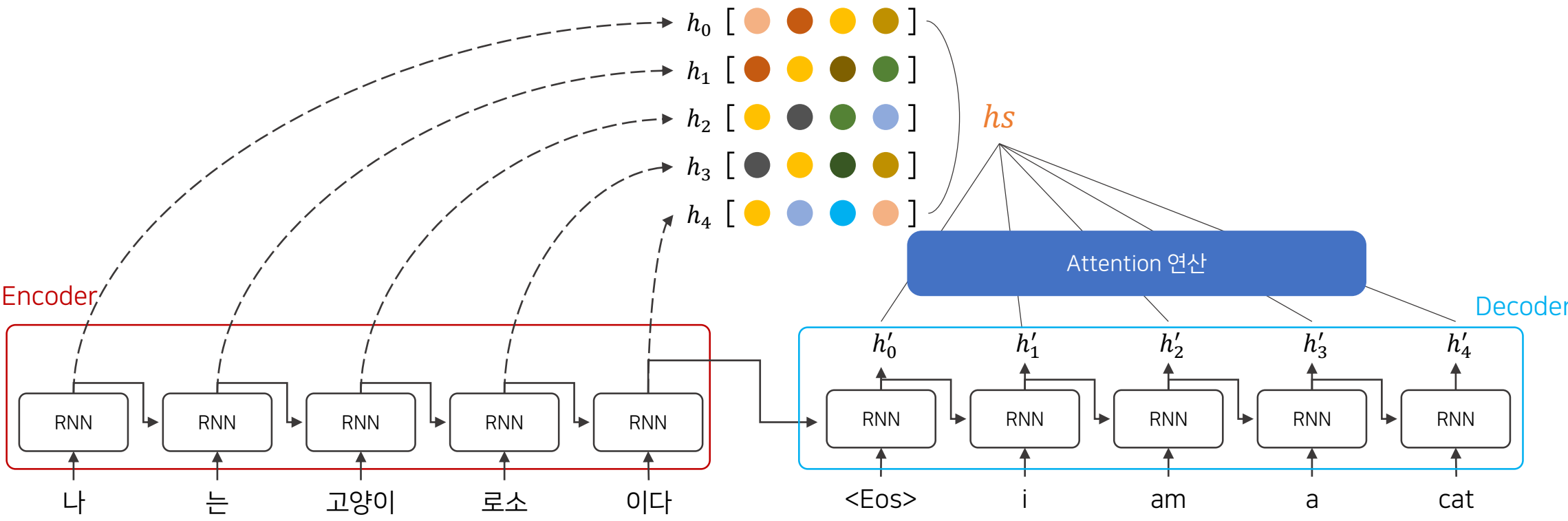
- Seq2Seq의 내부 구조는 RNN이기 때문에 RNN의 **장기 의존성** 문제를 그대로 가지고 있음
- 여러 단어로 이루어진 긴 문장을 하나의 벡터로 압축하기 때문에 **정보의 손실**이 일어남
- h_4 를 **Context Vector**로 활용함

$$h_{<eos>} = \tanh(h_4 W'_h + x_{<eos>} W'_x + b)$$

⇒ 하나의 레이어를 더 거쳐 "I" 라는 단어 생성

지금부터 Attention Mechanism 설명합니다.

Seq2Seq with Attention



Seq2Seq with Attention

일반적으로 내적 사용

$$\begin{matrix} h_0 [& \text{orange} & \text{dark orange} & \text{yellow} & \text{dark yellow} &] \\ h_1 [& \text{dark orange} & \text{yellow} & \text{dark yellow} & \text{green} &] \\ h_2 [& \text{yellow} & \text{dark grey} & \text{green} & \text{blue} &] \\ h_3 [& \text{dark grey} & \text{yellow} & \text{dark green} & \text{dark yellow} &] \\ h_4 [& \text{yellow} & \text{blue} & \text{cyan} & \text{orange} &] \end{matrix}$$

\otimes

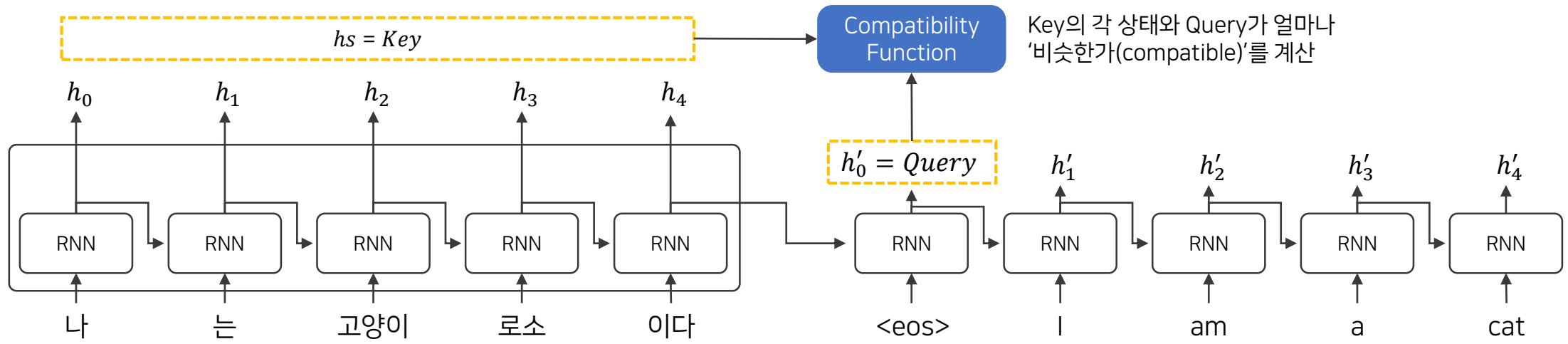
$$\begin{matrix} h'_0 [& \text{orange} &] \\ & \text{dark orange} &] \\ & \text{yellow} &] \\ & \text{dark grey} &] \end{matrix}$$

\equiv

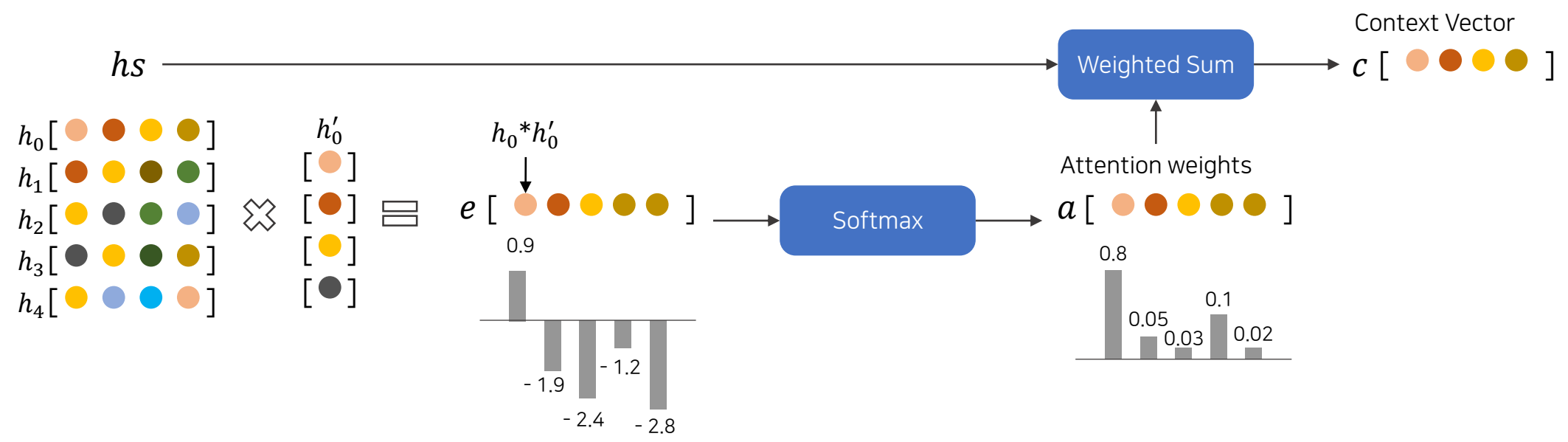
Energy scores

$$e [\text{orange} \text{ dark orange} \text{ yellow} \text{ dark yellow} \text{ dark green}]$$

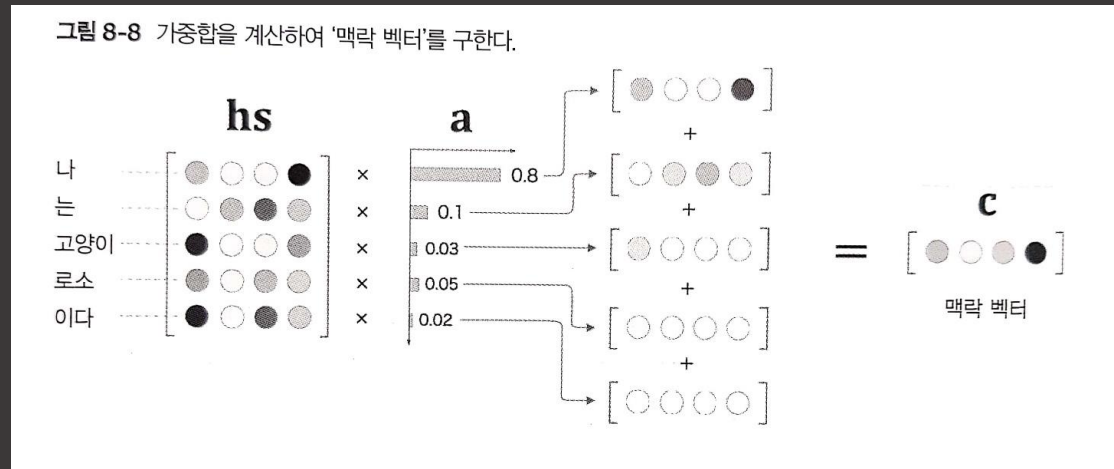
Encoder



Seq2Seq with Attention



Weighted Sum



Seq2Seq with Attention

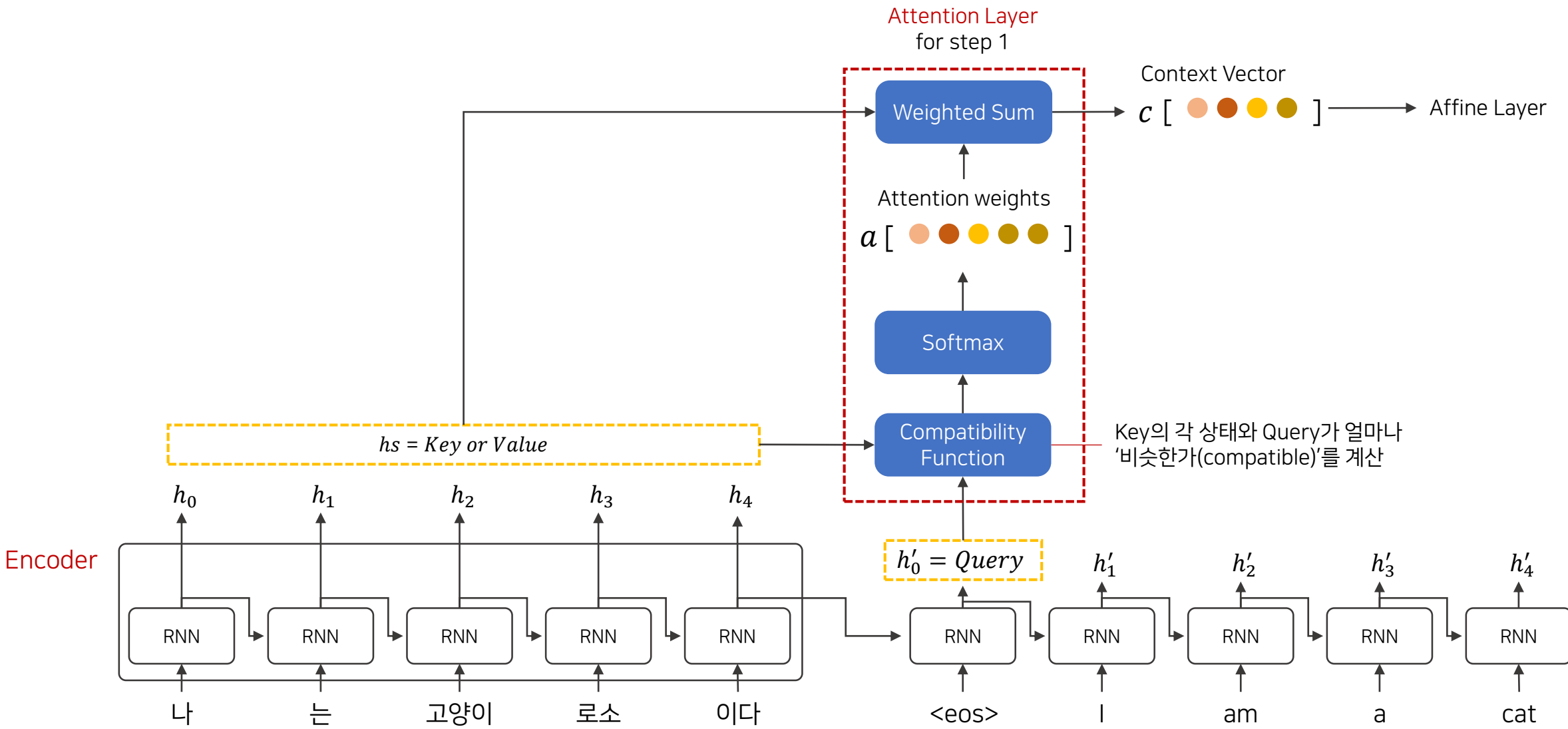
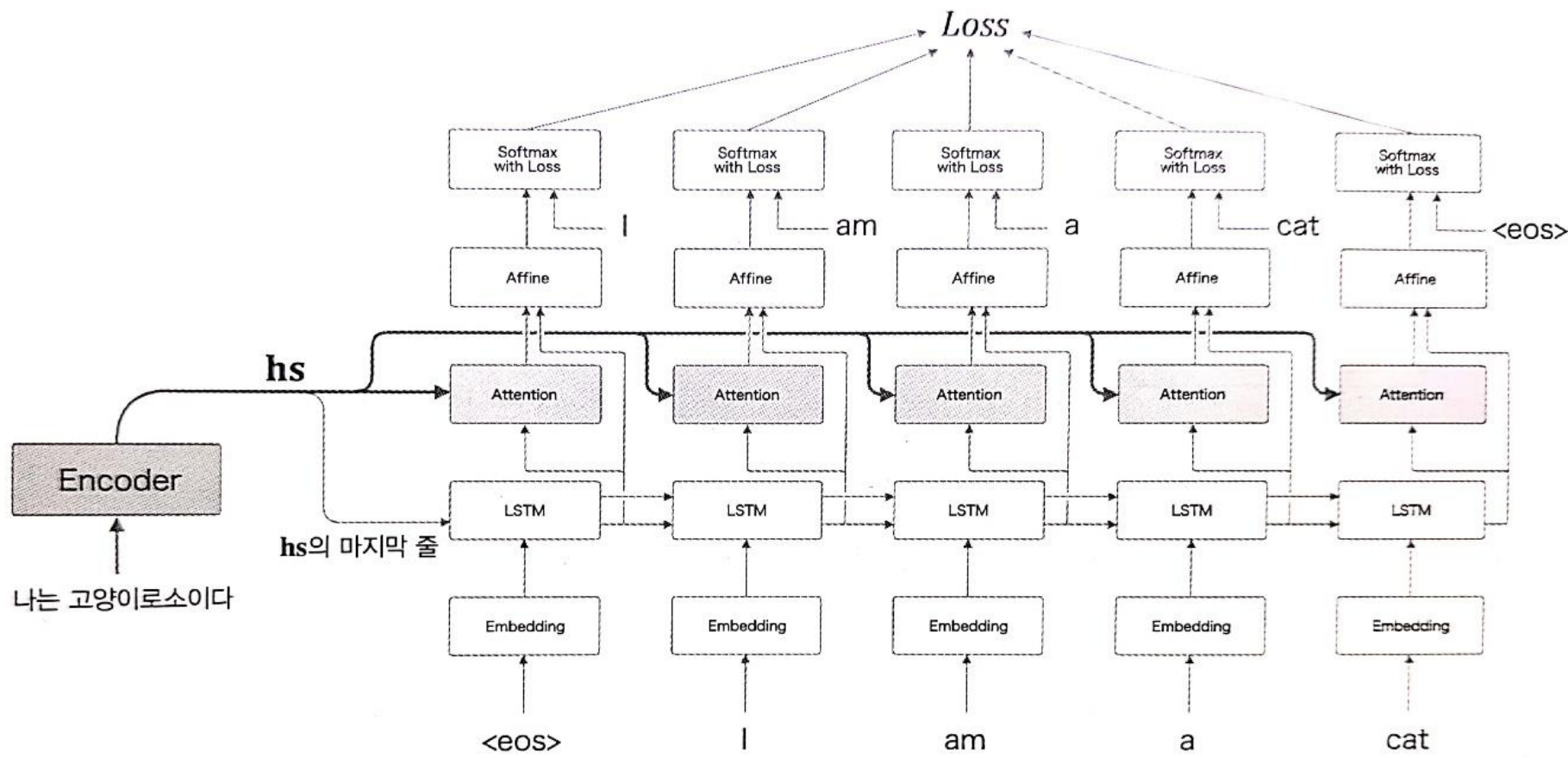


그림 8-18 Attention 계층을 갖춘 Decoder의 계층 구성



Attention, please! A Critical Review of Neural Attention Models in Natural Language Processing

- Compatibility functions

q와 K 비교(Comparing)
유사도 비교 등...

q와 K 결합(Combining)

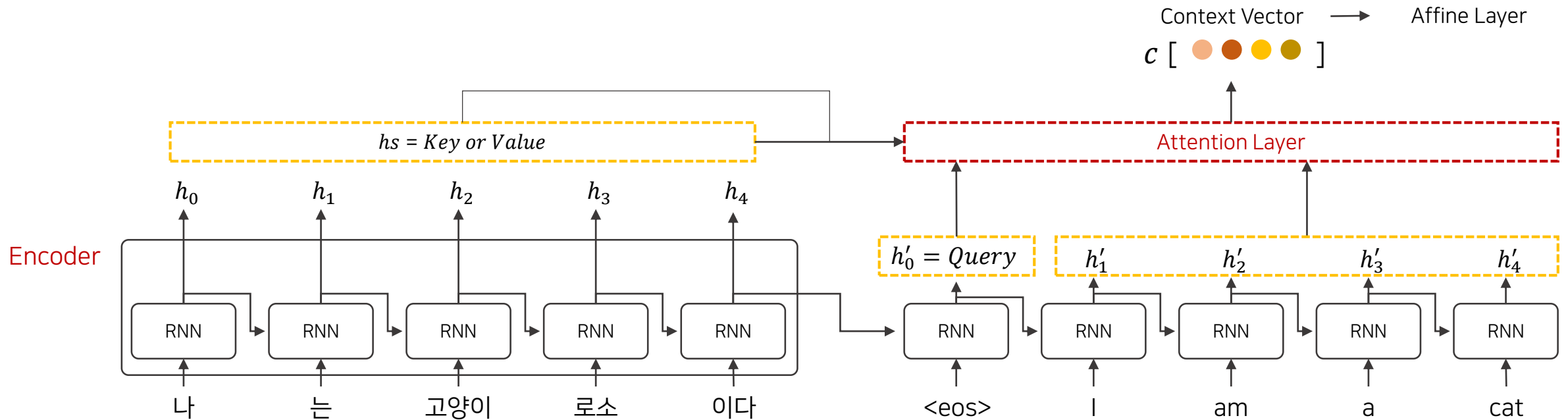
Self-Attention

Name	Equation	Reference
<i>similarity</i>	$f(q, K) = \text{sim}(q, K)$	Graves et al., 2014
<i>multiplicative or dot</i>	$f(q, K) = q^\top K$	Luong et al., 2015
<i>scaled multiplicative</i>	$f(q, K) = \frac{q^\top K}{\sqrt{d_k}}$	Vaswani et al., 2017
<i>general or bilinear</i>	$f(q, K) = q^\top W K$	Luong et al., 2015
<i>biased general</i>	$f(q, K) = K^\top (W q + b)$	Sordoni et al., 2016
<i>activated general</i>	$f(q, K) = \text{act}(q^\top W K + b)$	Ma et al., 2017
<i>concat</i>	$f(q, K) = w_{\text{imp}}^\top \text{act}(W[K; q] + b)$	Luong et al., 2015
<i>additive</i>	$f(q, K) = w_{\text{imp}}^\top \text{act}(W_1 K + W_2 q + b)$	Bahdanau et al., 2015
<i>deep</i>	$f(q, K) = w_{\text{imp}}^\top E^{(L-1)} + b^L$ $E^{(l)} = \text{act}(W_l E^{(l-1)} + b^l)$ $E^{(1)} = \text{act}(W_1 K + W_0 q + b^1)$	Pavlopoulos et al., 2017
<i>location-based</i>	$f(q, K) = f(q)$	Luong et al., 2015

Table 3: Summary of compatibility functions found in literature. W , W_0 , W_1 , ..., and b are learnable parameters.

Self- Attention Mechanism

Seq2Seq with Attention



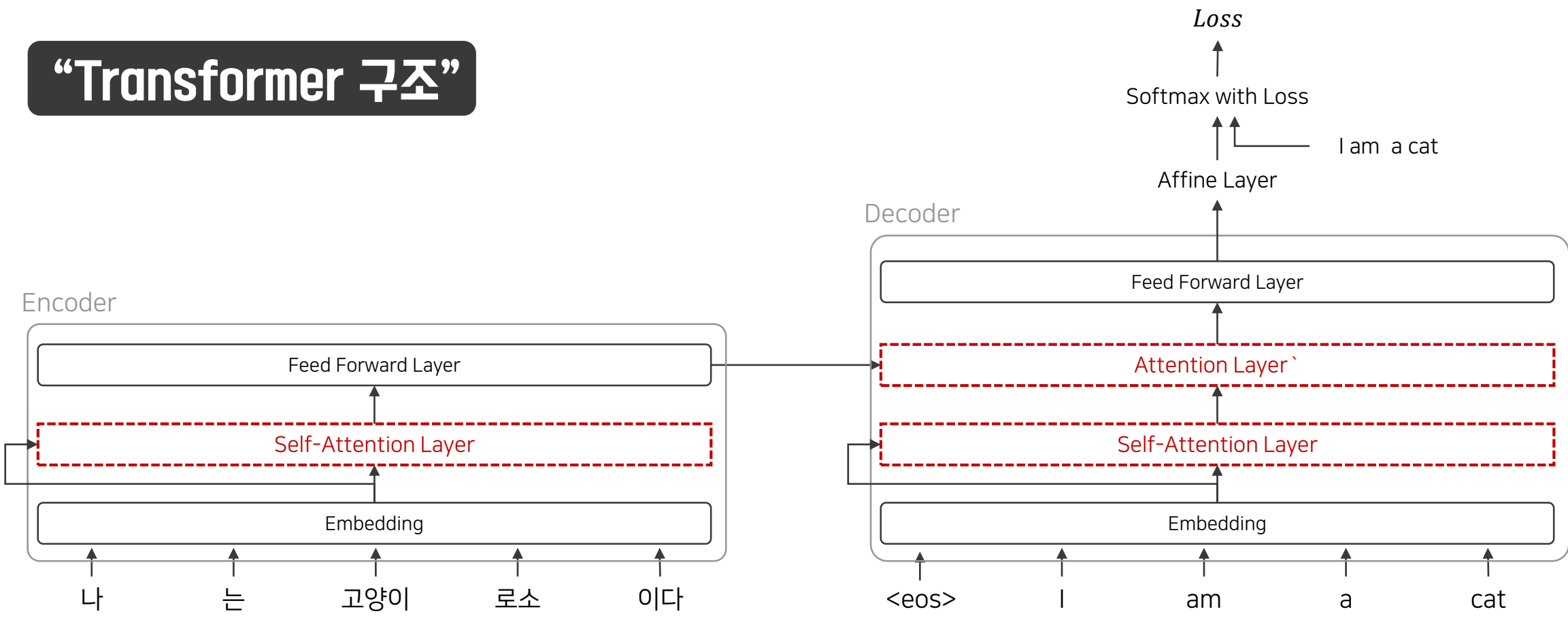
- RNN 구조의 근본적인 한계는 그대로 가지고 있음
 - ✓ 장기 의존성 문제
 - ✓ 병렬 연산이 불가능함 → 학습 속도에서 병목

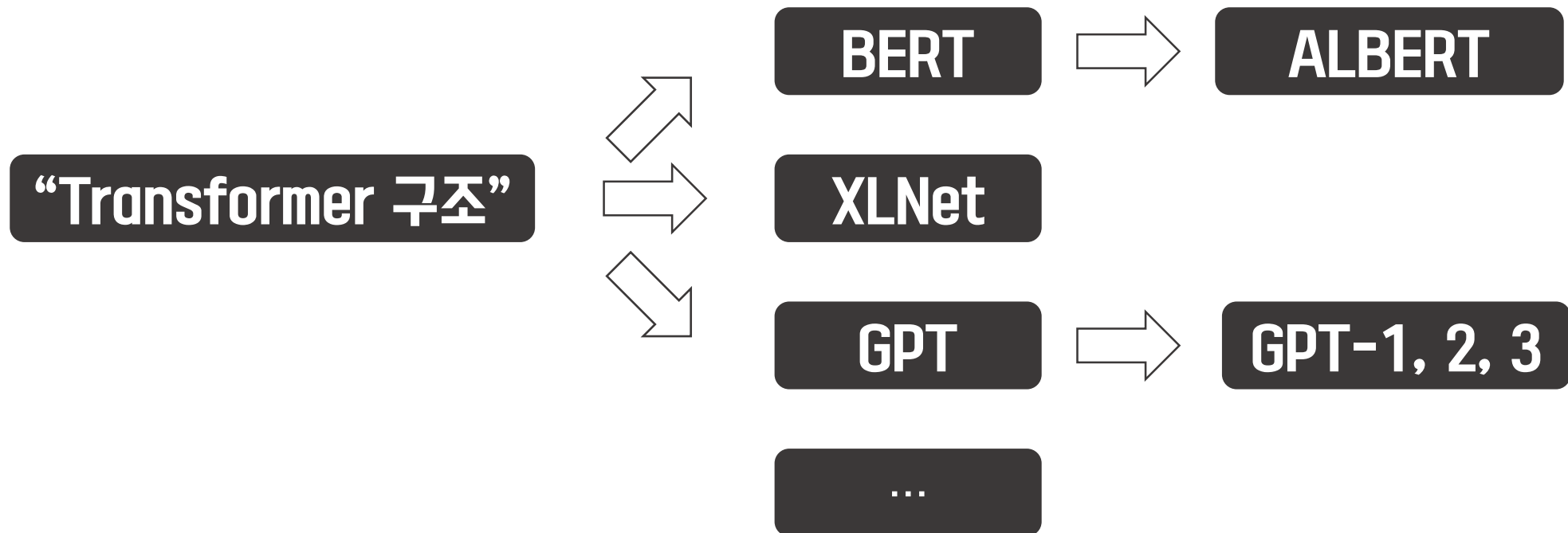
Attention is all you need !!!

Attention is all you need

Aurthors: Ashish Vaswani(Google Brain) 외 7명,
Conference: Advances in Neural Information Processing Systems 30 (NIPS 2017)

“Transformer 구조”





감사합니다.