# [10주차]
# Recurrent Neural Networks
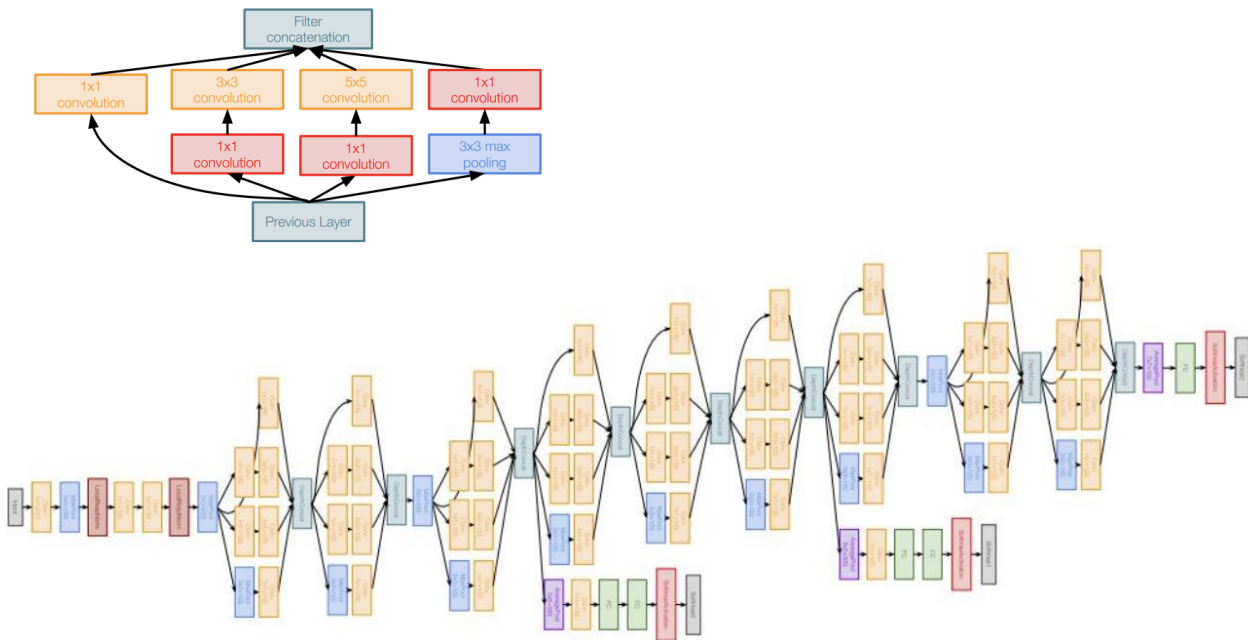
### 1기 강다연
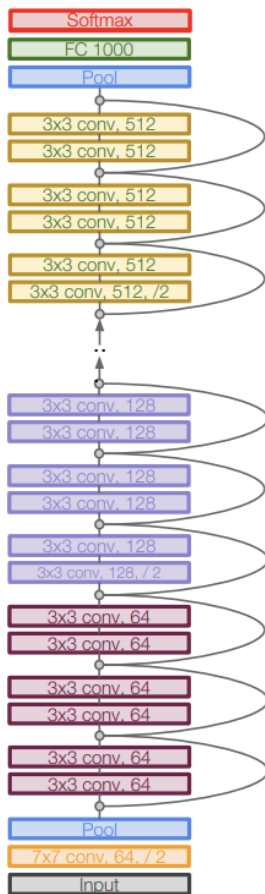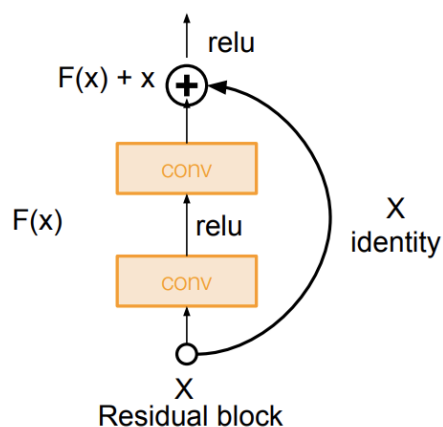### 1기 장예서

# 목차

# 9주차 Review

## VGGNet



VGG16      VGG19

## GoogLeNet
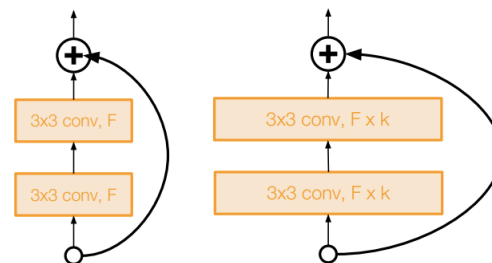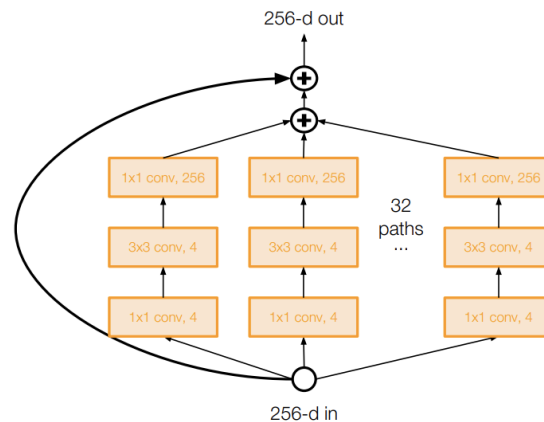
# 9주차 Review

## ResNet



## Wide Residual Network



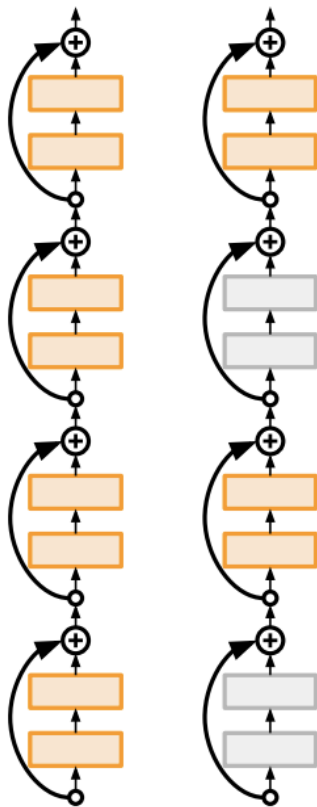Basic residual block    Wide residual block

## ResNeXt



EURON

# 9주차 Review

## Deep Networks with Stochastic Depth



## FractalNet



Figures copyright Larsson et al., 2017. Reproduced with permission.

# 9주차 Review

## Densely Connected Convolution Network



Dense Block

| |
|---|
| Softmax |
| FC |
| Pool |
| Dense Block 3 |
| Conv |
| Pool |
| Conv |
| Dense Block 2 |
| Conv |
| Pool |
| Conv |
| Dense Block 1 |
| Conv |
| Input |

## SqueezeNet



Figure copyright Iandola, Han, Moskewicz, Ashraf, Dally, Keutzer, 2017. Reproduced with permission.

# RNN Basics



$$h_t = f_W(h_{t-1}, x_t)$$

new state — some function with parameters W — old state   input vector at some time step

## Vanila RNN

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t)$$

$$y_t = W_{hy} h_t$$

# RNN Basics

# RNN Basics



one to one     one to many     many to one     many to many     many to many

# RNN Basics

## Many-to-many

# RNN Basics

## Many-to-one



## One-to-many

# RNN Basics

## Sequence-to-sequence

# Character-level language model



Training

Test time

target chars: "e" "l" "l" "o"

Sample "e" "l" "l" "o"

Softmax

EURON

# Character-level language model

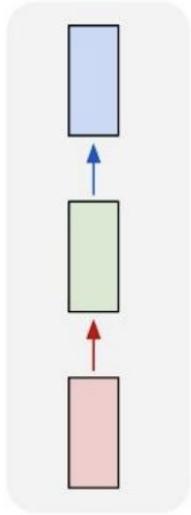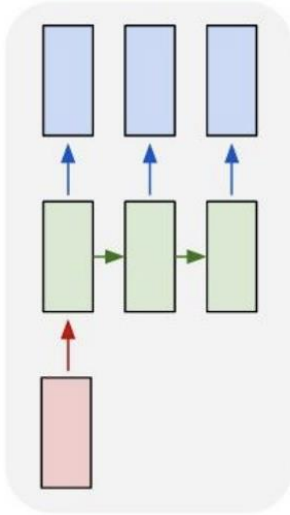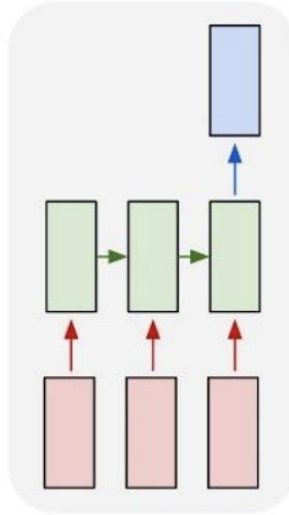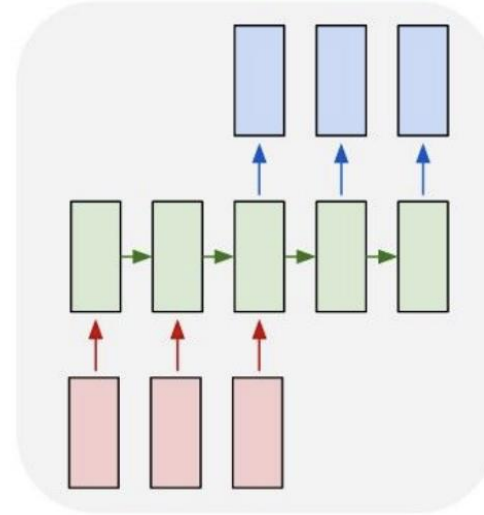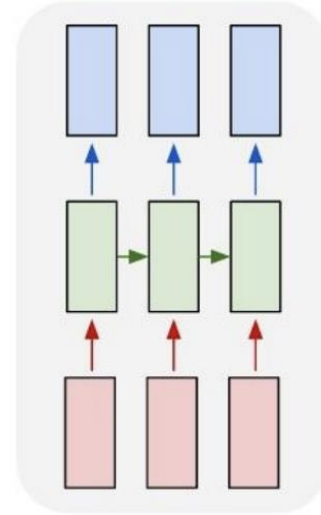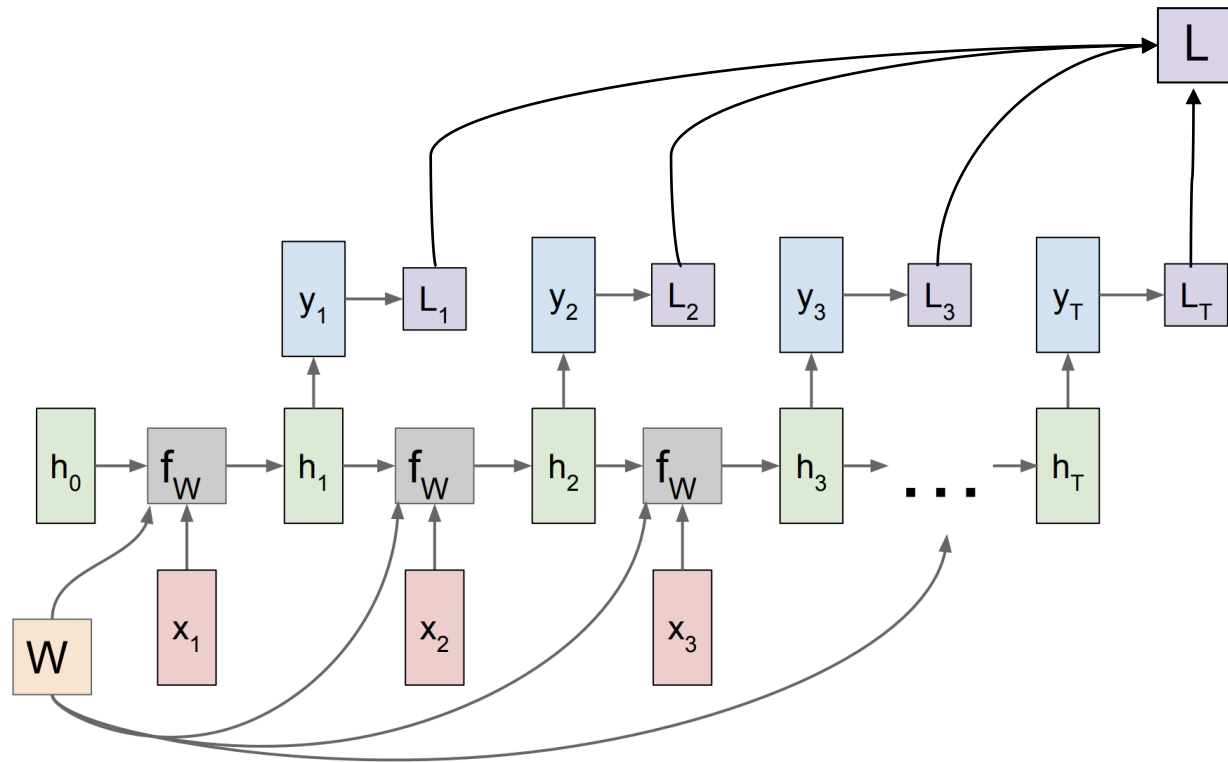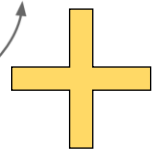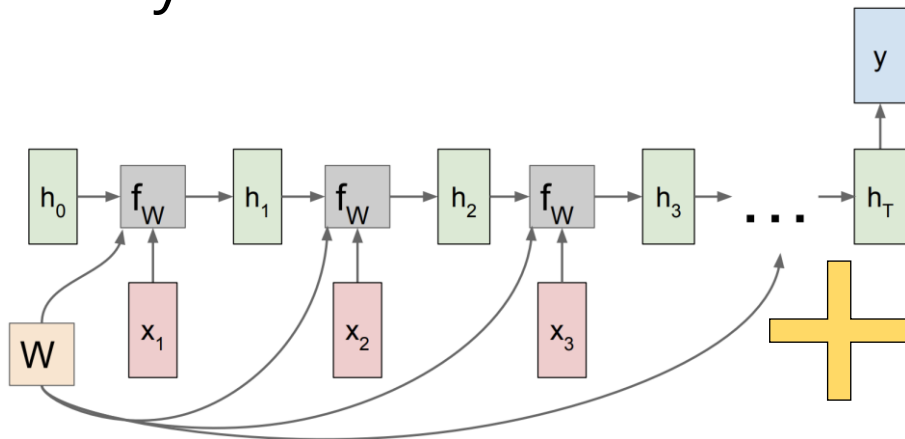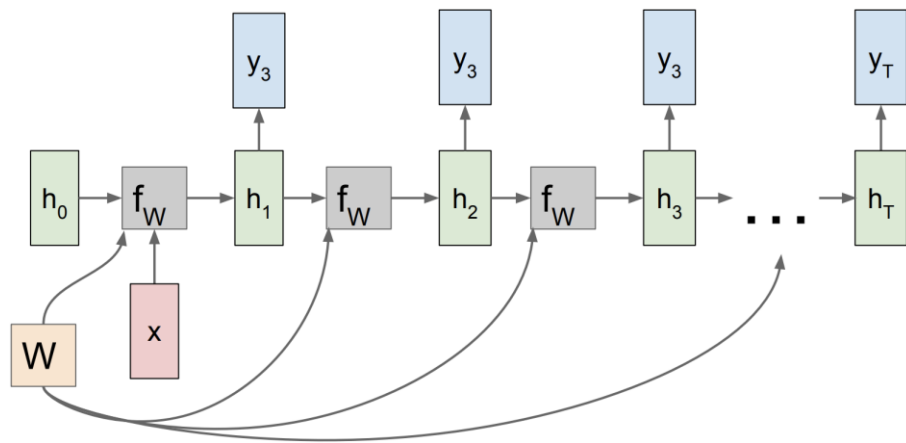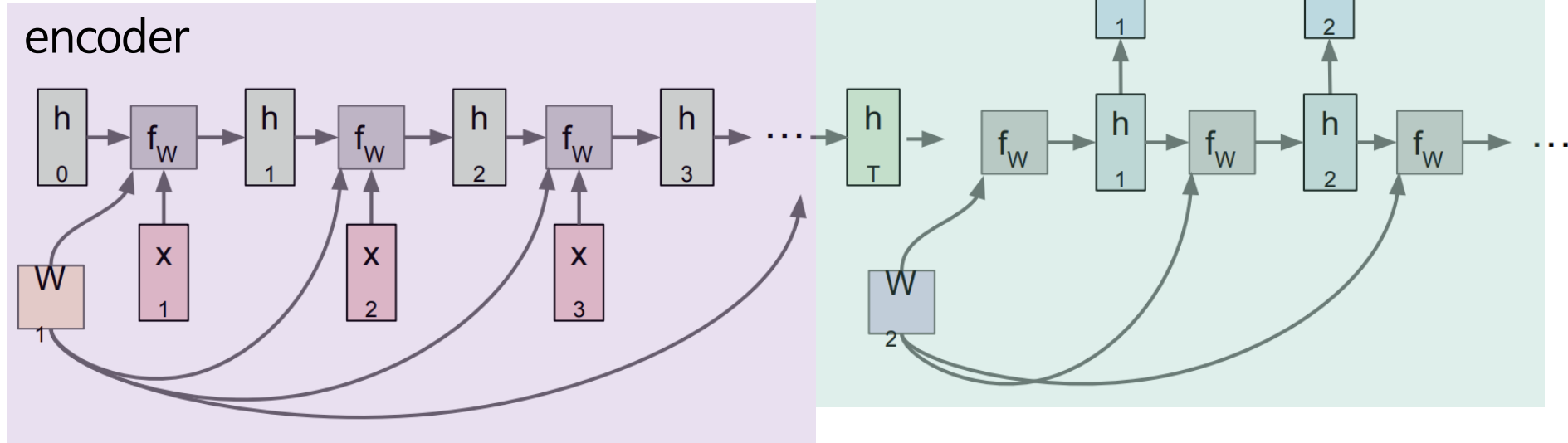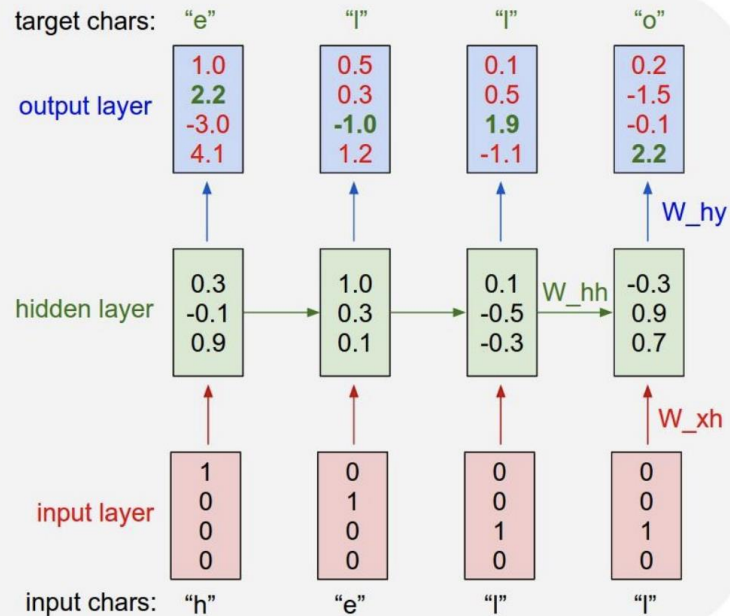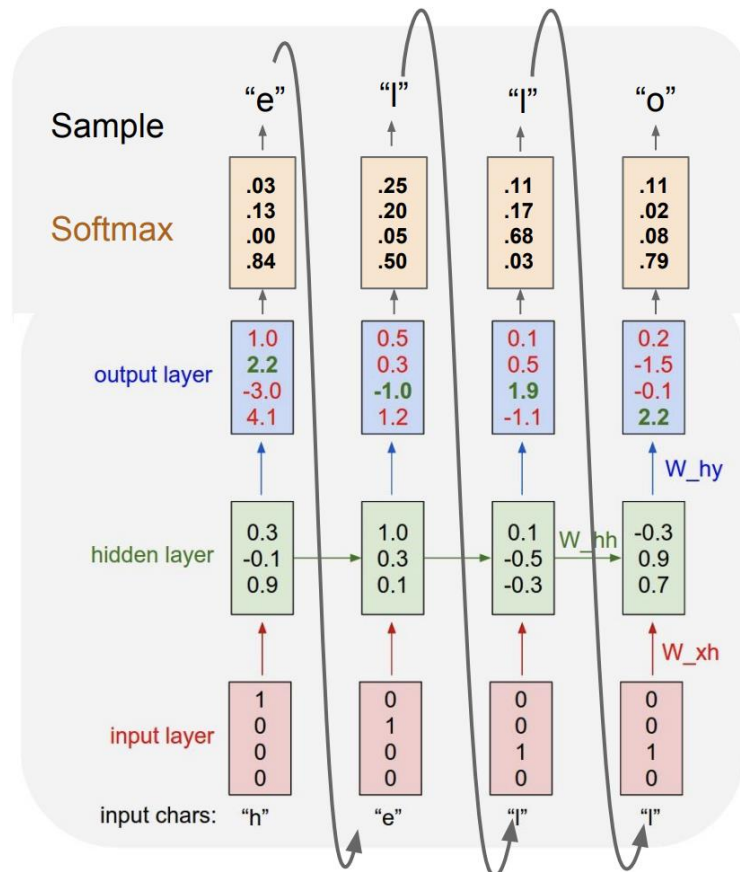"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

---

*Proof.* Omitted. □

**Lemma 0.1.** *Let $C$ be a set of the construction.*
*Let $C$ be a gerber covering. Let $F$ be a quasi-coherent sheaves of $O$-modules. We have to show that*

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

*Proof.* This is an algebraic space with the composition of sheaves $F$ on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where $\mathcal{G}$ defines an isomorphism $\mathcal{F} \to \mathcal{F}$ of $\mathcal{O}$-modules. □

**Lemma 0.2.** *This is an integer $Z$ is injective.*

*Proof.* See Spaces, Lemma ??. □

**Lemma 0.3.** *Let $S$ be a scheme. Let $X$ be a scheme and $X$ is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let $X$ be a scheme. Let $X$ be a scheme which is equal to the formal complex.*

*The following to the construction of the lemma follows.*

*Let $X$ be a scheme. Let $X$ be a scheme covering. Let*

$$b : X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$$

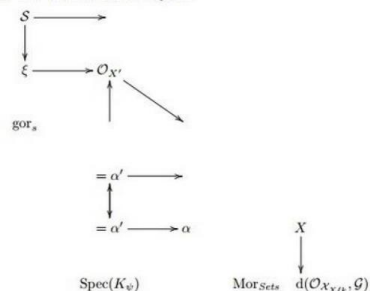*be a morphism of algebraic spaces over $S$ and $Y$.*

*Proof.* Let $X$ be a nonzero scheme of $X$. Let $X$ be an algebraic space. Let $F$ be a quasi-coherent sheaf of $\mathcal{O}_X$-modules. The following are equivalent

(1) $\mathcal{F}$ is an algebraic space over $S$.
(2) If $X$ is an affine open covering.

Consider a common structure on $X$ and $X$ the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

---

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram



is a limit. Then $\mathcal{G}$ is a finite type and assume $S$ is a flat and $\mathcal{F}$ and $\mathcal{G}$ is a finite type $f_*$. This is of finite type diagrams, and

• the composition of $\mathcal{G}$ is a regular sequence,
• $\mathcal{O}_{X'}$ is a sheaf of rings. □

*Proof.* We have see that $X = \operatorname{Spec}(R)$ and $\mathcal{F}$ is a finite type representable by algebraic space. The property $\mathcal{F}$ is a finite morphism of algebraic stacks. Then the cohomology of $X$ is an open neighbourhood of $U$. □

*Proof.* This is clear that $\mathcal{G}$ is a finite presentation, see Lemmas ??.
A *reduced above* we conclude that $U$ is an open covering of $C$. The functor $\mathcal{F}$ is a "field

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\overline{x}} \quad -1(\mathcal{O}_{X_{\text{étale}}}) \longrightarrow \mathcal{O}_{X_{\overline{x}}}^{-1}\mathcal{O}_{X_{\lambda}}(\mathcal{O}_{X_{\overline{x}}}^{\nabla})$$

is an isomorphism of covering of $\mathcal{O}_{X_i}$. If $\mathcal{F}$ is the unique element of $\mathcal{F}$ such that $X$ is an isomorphism.
The property $\mathcal{F}$ is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme $\mathcal{O}_X$-algebra with $\mathcal{F}$ are opens of finite type over $S$.
If $\mathcal{F}$ is a scheme theoretic image points. □

If $\mathcal{F}$ is a finite direct sum $\mathcal{O}_{X_{\lambda}}$ is a closed immersion, see Lemma ??. This is a sequence of $\mathcal{F}$ is a similar morphism.

# Another examples of RNN

## Examples of sequence data to apply RNN

| | | |
|---|---|---|
| Speech recognition | ⌇⌇⌇⌇⌇⌇⌇ | "The quick brown fox jumped over the lazy dog." |
| Music generation | ∅ (Nothing) | ♪♪♪♪ |
| Sentiment classification | "There is nothing to like in this movie." | ★☆☆☆☆ |
| DNA sequence analysis | AGCCCCTGTGAGGAACTAG | AGCCCCTGTGAGGAACTAG |
| Machine translation | Voulez-vous chanter avec moi? | Do you want to sing with me? |
| Video activity recognition | | Running |

# Image Captioning



"man in black shirt is playing guitar."

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

"...oung girls are playing with lego toy."

Input Image (224x224x3)

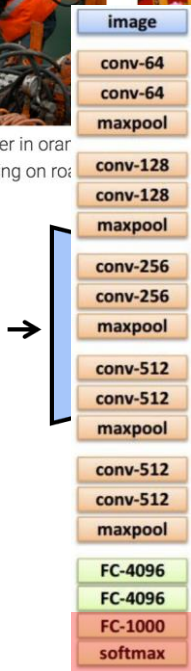image
conv-64
conv-64
maxpool
conv-128
conv-128
maxpool
conv-256
conv-256
maxpool
conv-512
conv-512
maxpool
conv-512
conv-512
maxpool
FC-4096
FC-4096
FC-1000
softmax

Feature vector at fc layer (1×1×2048)

Linear

Wih

$h = \tanh(Wxh * x + Whh * h + Wih * v)$

# Visual Question Answering



COCOQA
Q: What is the color of the desk?
A: white
Q: What are on the white desk?
A: computers

COCOQA
Q: What is the color of the dresses?
A: purple
Q: What are three women dressed up and on?
A: phones

VQA
Q: How many bikes are there?
A: 2
Q: What number is the bus?
A: 48

VQA
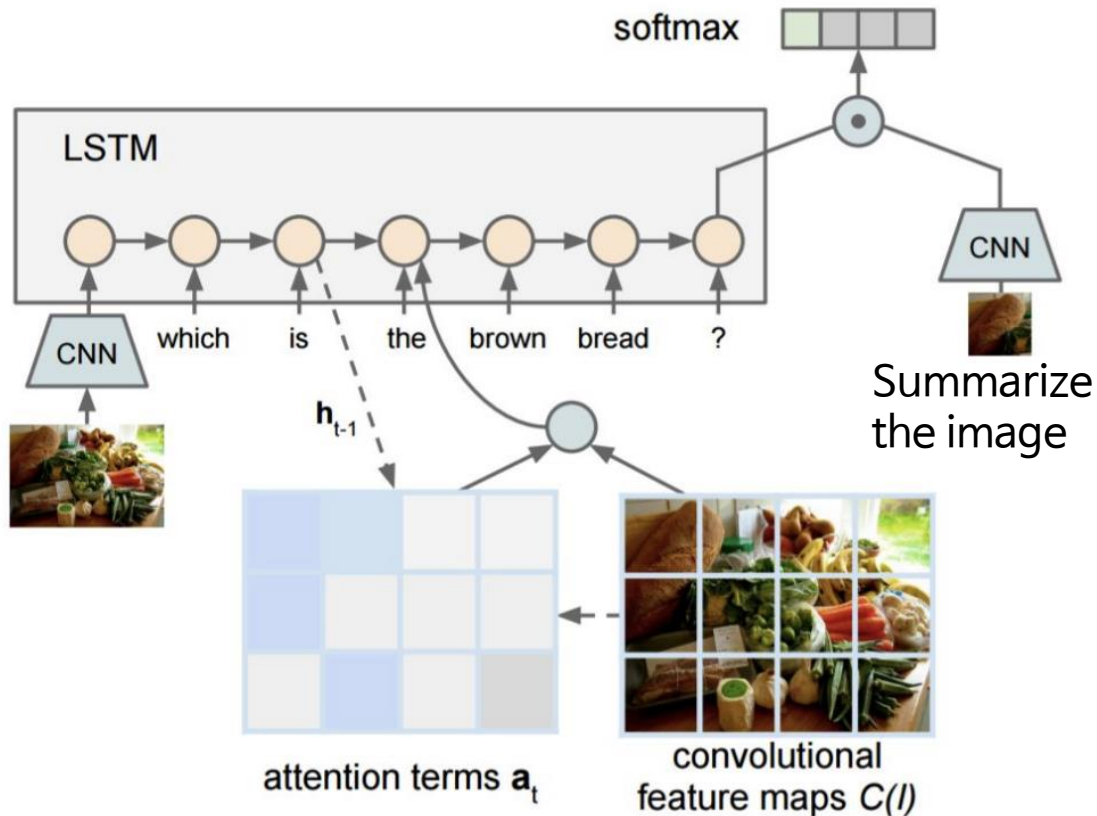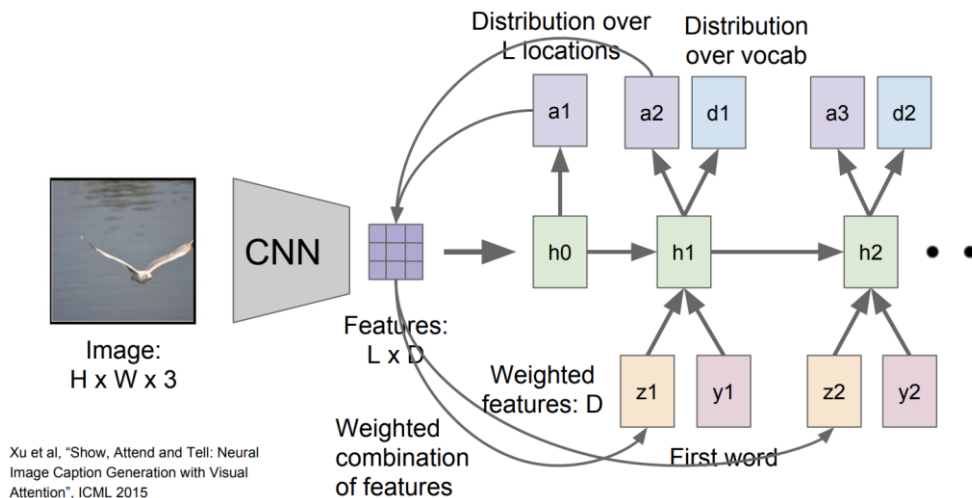Q: How many pickles are on the plate?
A: 1
Q: What is the shape of the plate?
A: round

LSTM

CNN    which    is    the    brown    bread    ?

$h_{t-1}$

softmax

CNN

Summarize the image

attention terms $a_t$

convolutional feature maps $C(I)$

# RNN with Attention

## Image Captioning with Attention



Distribution over L locations

Distribution over vocab

Image: H x W x 3

CNN

Features: L x D

Weighted features: D

Weighted combination of features

First word

Xu et al, "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention", ICML 2015
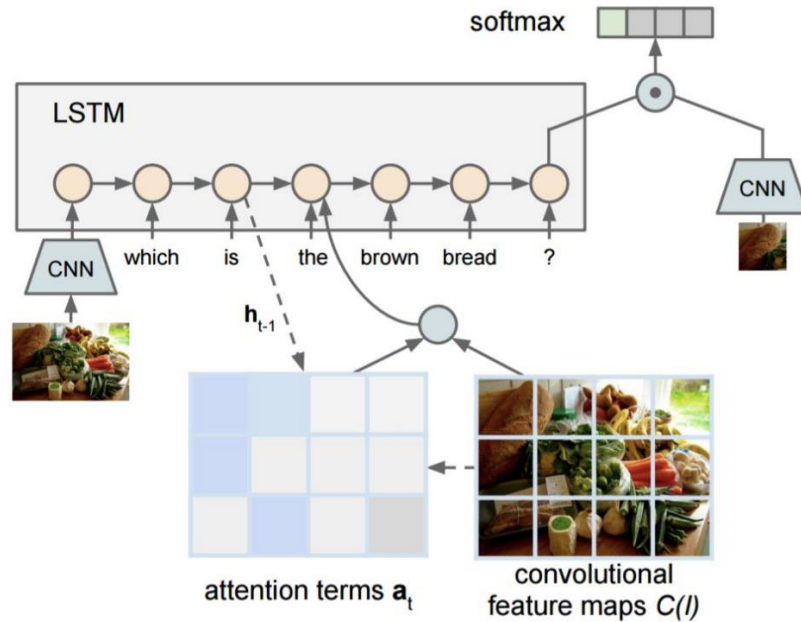
Soft attention

Hard attention

A    bird    flying    over    a    body    of    water    .

# RNN with Attention

## Visual Question Answering with Attention



softmax

LSTM

CNN  which  is  the  brown  bread  ?

$h_{t-1}$

CNN

attention terms $a_t$

convolutional
feature maps $C(I)$

Zhu et al, "Visual 7W: Grounded Question Answering in Images", CVPR 2016
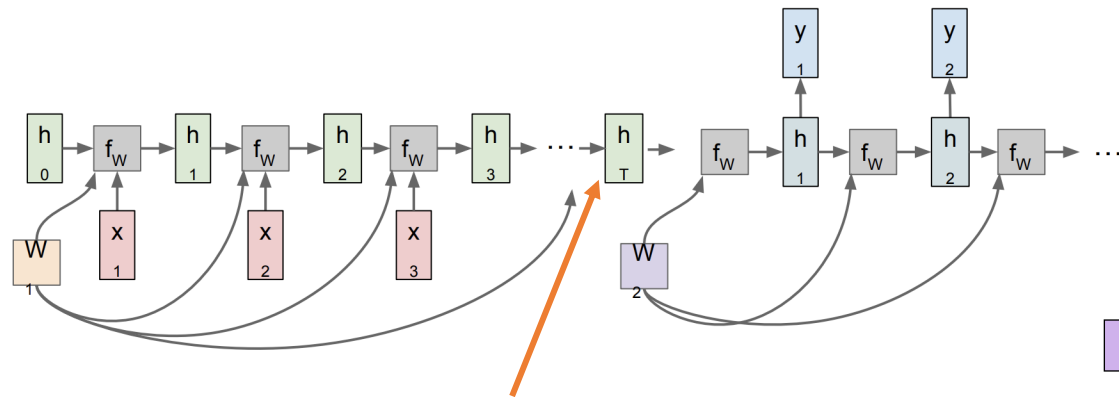Figures from Zhu et al, copyright IEEE 2016. Reproduced for educational purposes.

A
cat

What kind of animal is in the photo?
A **cat**.

B
cake

Why is the person holding a knife?
To cut the **cake** with.

# Attention

## RNN 기반 sequence-to-sequence model



하나의 고정된 크기의 벡터에 모든 정보 압축

**Attention**
Decoder에서 단어를 예측하는 매시점마다 예측할 단어와 연관이 있는 encoder의 입력 부분을 더 집중해서 다시 참고

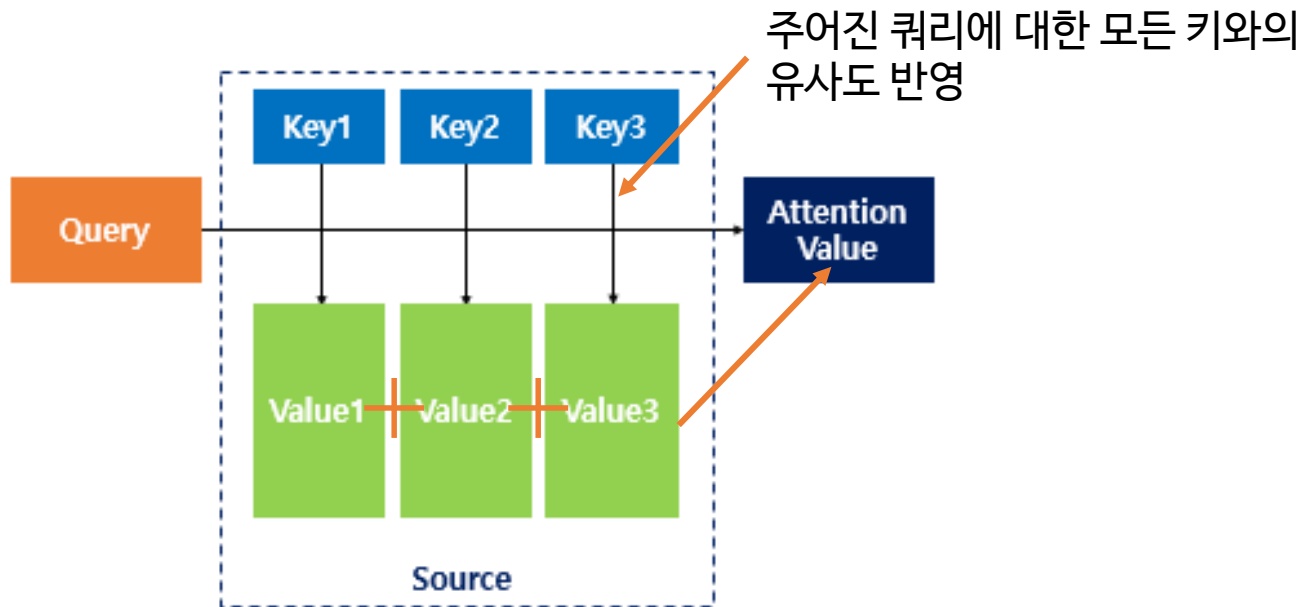**Problems:**
1.  정보 손실 발생
2.  RNN의 고질적인 문제 – Vanishing Gradient

# Attention

## Attention 함수
Attention(Q, K, V) = Attention Value

Q = Query: t 시점의 디코더 셀에서의 은닉 상태
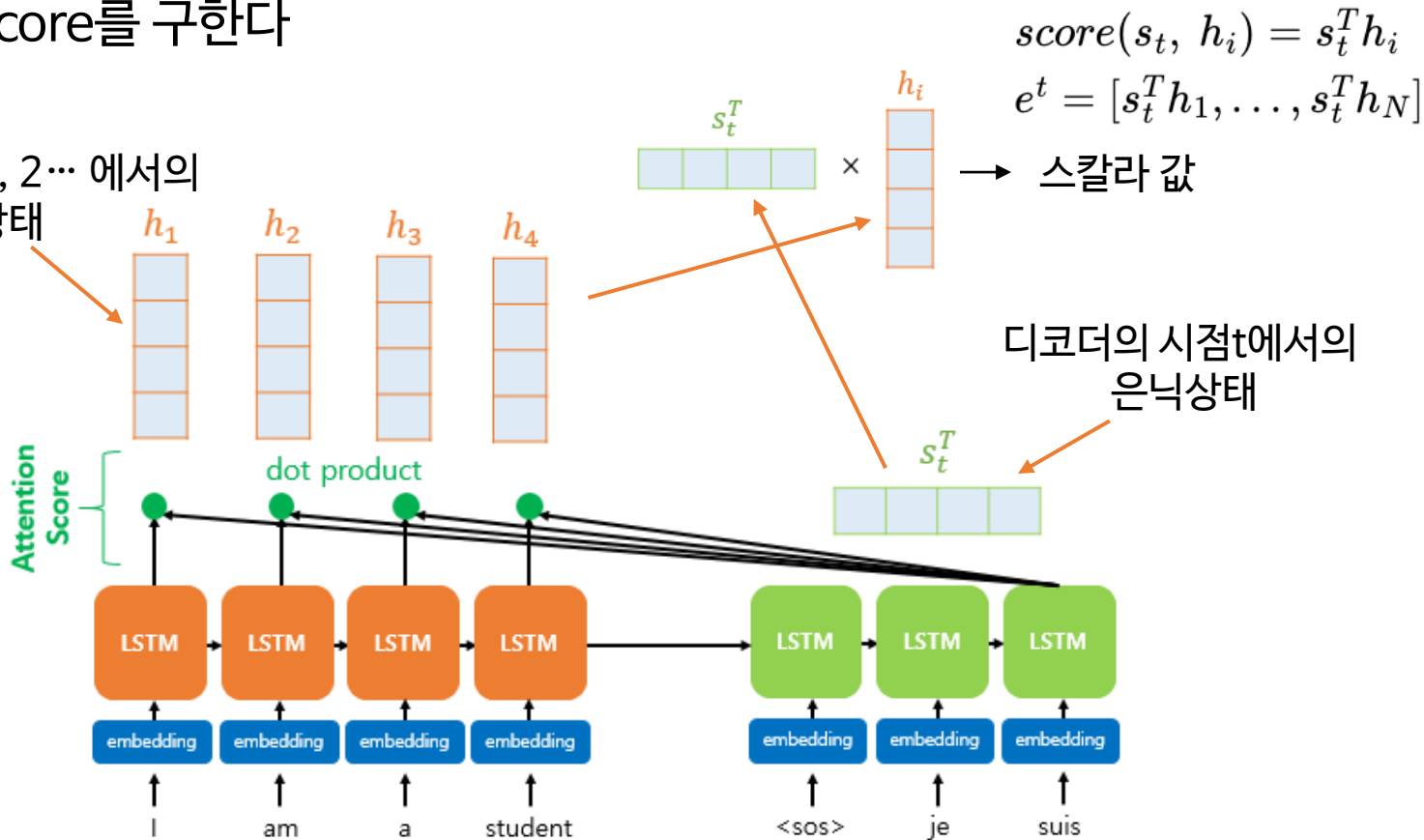K = Keys: 모든 시점의 인코더 셀의 은닉 상태들
V = Values: 모든 시점의 인코더 셀의 은닉 상태들

주어진 쿼리에 대한 모든 키와의
유사도 반영

# Dot-Product Attention



각 단어가 출력 단어를 예측할 때
도움이 되는 정도

Decoder로 하나의 정보로 전송

Attention 사용

étudiant

softmax

Dense

softmax

LSTM LSTM LSTM LSTM    LSTM LSTM LSTM

embedding embedding embedding embedding    embedding embedding embedding

I    am    a    student    <sos>    je    suis

# Dot-Product Attention

## 1) attention score를 구한다

$$score(s_t, h_i) = s_t^T h_i$$
$$e^t = [s_t^T h_1, \ldots, s_t^T h_N]$$

인코더의 시점 1, 2··· 에서의
은닉상태

$s_t^T$     $h_i$

$h_1$    $h_2$    $h_3$    $h_4$

× → 스칼라 값

디코더의 시점t에서의
은닉상태

Attention Score

dot product

$s_t^T$

LSTM LSTM LSTM LSTM    LSTM LSTM LSTM

embedding embedding embedding embedding    embedding embedding embedding

I    am    a    student     \<sos\>    je    suis

# Dot-Product Attention

## 2) 소프트맥스 함수를 통해 어텐션 분포를 구한다



Attention Weight
모든 attention weight의 합은 1

어텐션 분포
$e^t$에 소프트맥스 함수
적용한 확률분포

$$\alpha^t = softmax(e^t)$$

# Dot-Product Attention

3) 각 인코더의 어텐션 weight와 은닉 상태를 가중합하여 어텐션 value를 구한다

Context vector

**Attention Value $a_t$**

$$a_t = \sum_{i=1}^{N} \alpha_i^t h_i$$

$h_1$    $h_2$    $h_3$    $h_4$

softmax

LSTM   LSTM   LSTM   LSTM    LSTM   LSTM   LSTM

embedding   embedding   embedding   embedding    embedding   embedding   embedding

I    am    a    student    &lt;sos&gt;    je    suis

# Dot-Product Attention

4) Attention value와 디코더의 t 시점의 은닉 상태를 연결한다 (concatenate)



예측 연산의 입력으로 사용
(인코더로부터 얻은 정보 활용)

# Dot-Product Attention

5) 출력층 연산의 입력이 되는 $s_t$를 계산한다



가중치 matrix

$W_c$

$v_t$

$\tilde{s_t}$

$$\tilde{s}_t = \tanh\left(\mathbf{W_c}[a_t; s_t] + b_c\right)$$

6) $\tilde{s}_t$를 출력층의 입력으로 사용한다

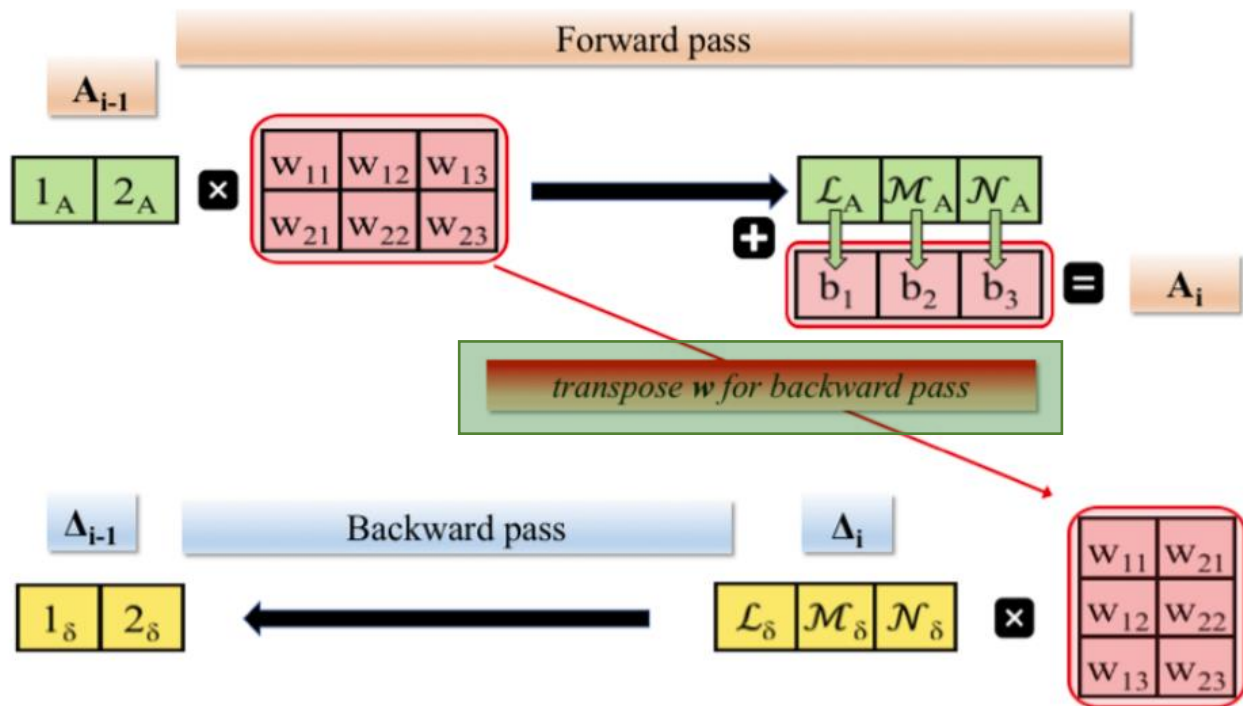$$\hat{y}_t = \text{Softmax}\left(W_y \tilde{s}_t + b_y\right)$$
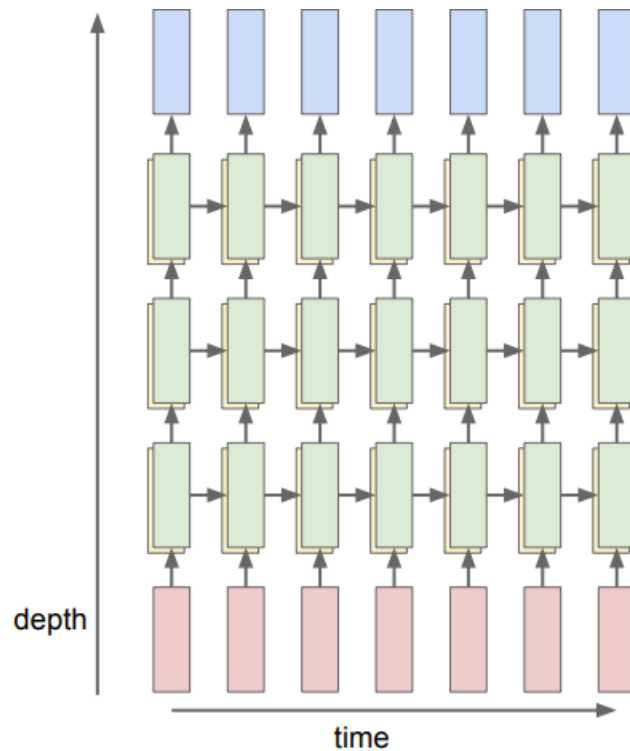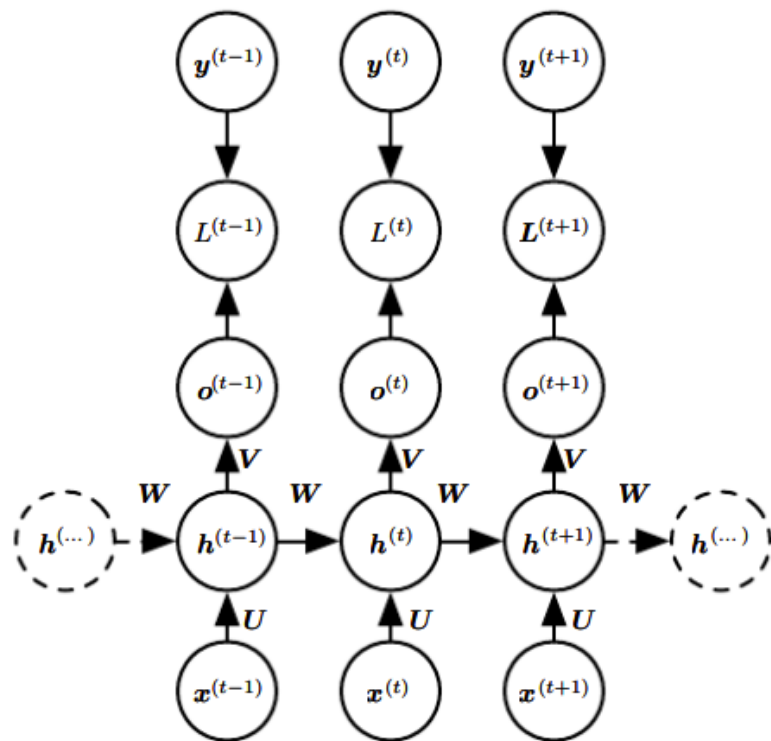
예측 벡터

# Backpropagation of RNN



Backpropagation from $h_t$ to $h_{t-1}$ multiplies by W (actually $W_{hh}^T$)

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$= \tanh\left( \begin{pmatrix} W_{hh} & W_{hx} \end{pmatrix} \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

$$= \tanh\left( W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

# Backpropagation of RNN



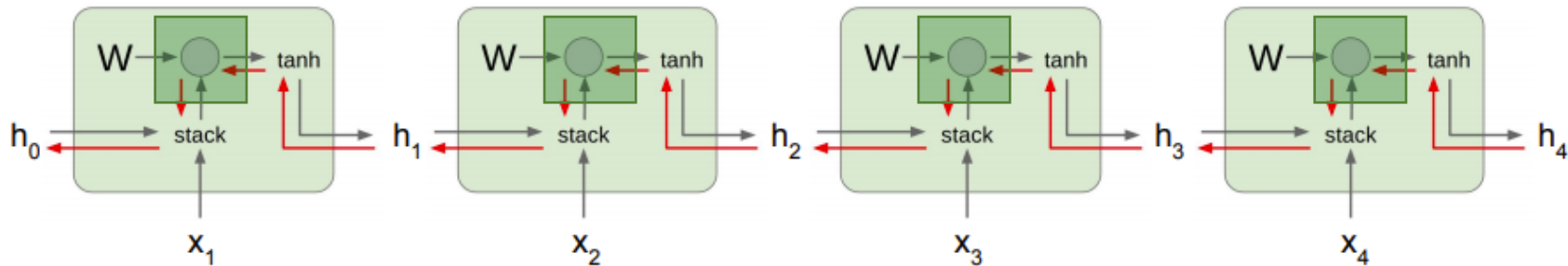Deep learning for pedestrians: backpropagation in CNNs Laurent Bou´e

# Backpropagation of RNN: Multi-Layer RNN

# Backpropagation of RNN: Truncated Backpropagation

# Backpropagation of RNN: Multi-Layer RNN
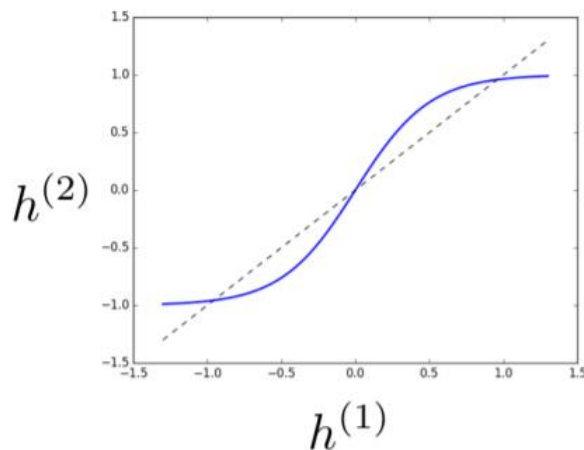


Cell 하나를 **통과할 때마다 mul gate 존재**
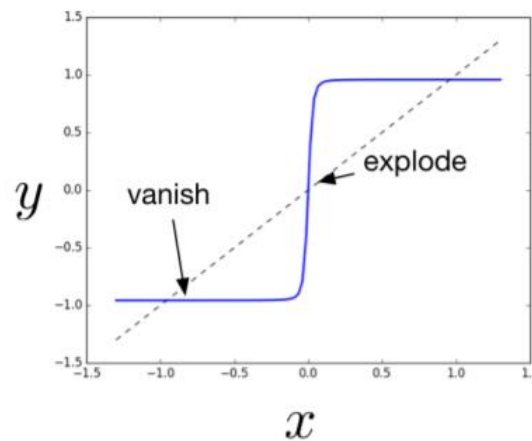
**각 Cell**의 W transpose factor 관여

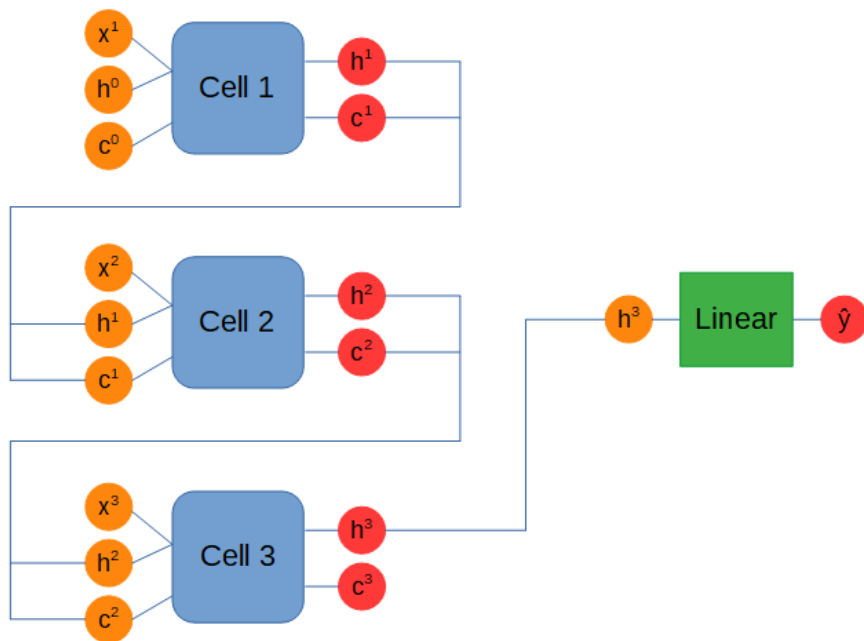**매우 많은 W가 관여**하게 됨

# Backpropagation of RNN



Largest Singular Value $>$ 1:
Exploding Gradients

Largest Singular Value $<$ 1:
Vanishing Gradients

# ▶ Change Architecture!
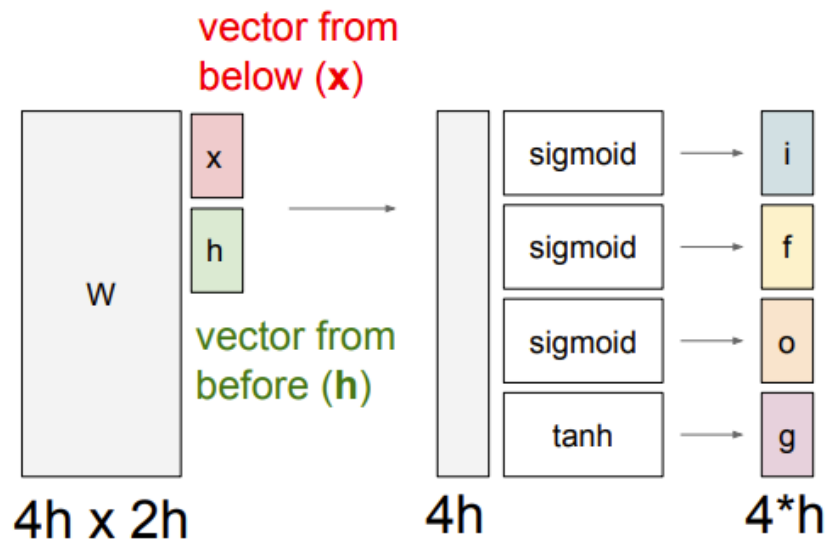
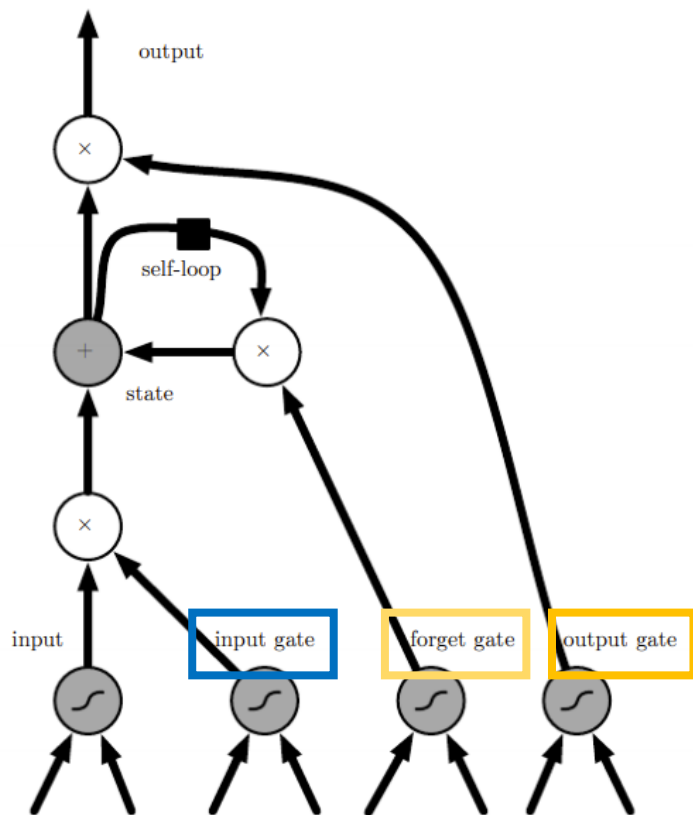EURON

# LSTM-Long Short Term Memory



cell당 하나의 hidden states ▶ cell당 두 개의 hidden states

EURON

# LSTM-Long Short Term Memory



output

self-loop

state

input

input gate

forget gate

output gate

vector from below (**x**)

x

h

W

vector from before (**h**)

4h x 2h
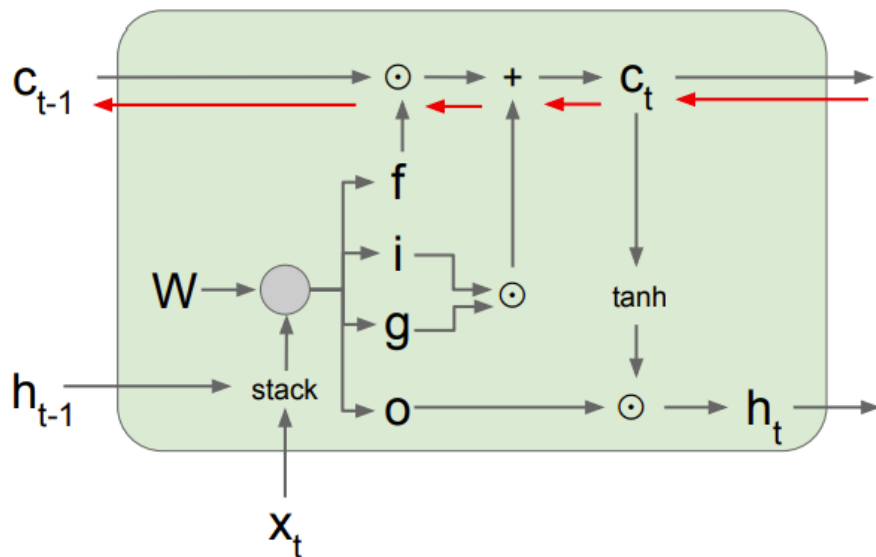
sigmoid → i

sigmoid → f
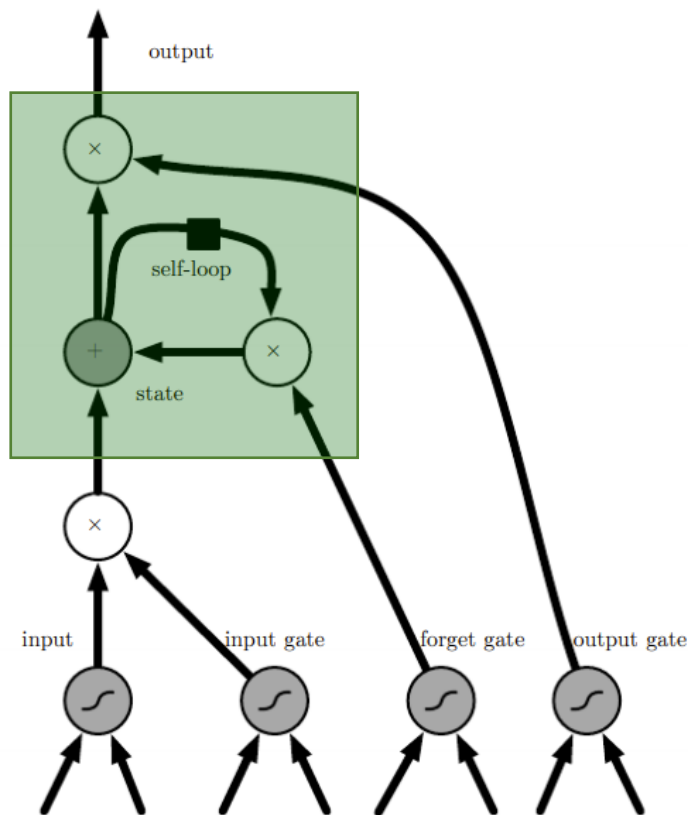
sigmoid → o

tanh → g

4h

4*h

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

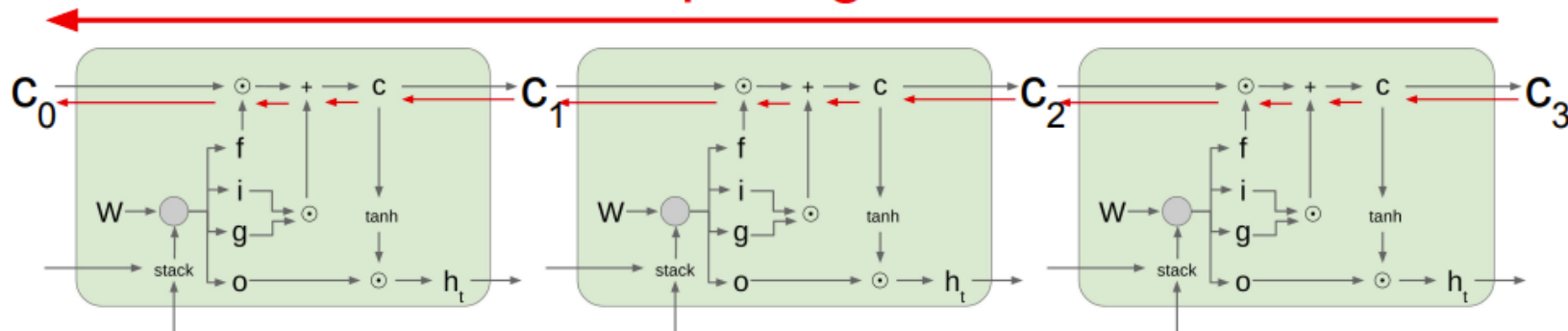$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

# LSTM-Long Short Term Memory



Addition Operation
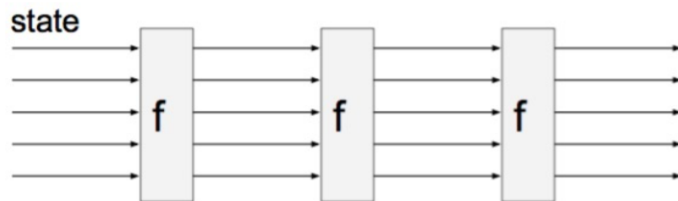▶ C의 backpropagation 과정에서
**W에 대한 multiplication 존재 X**

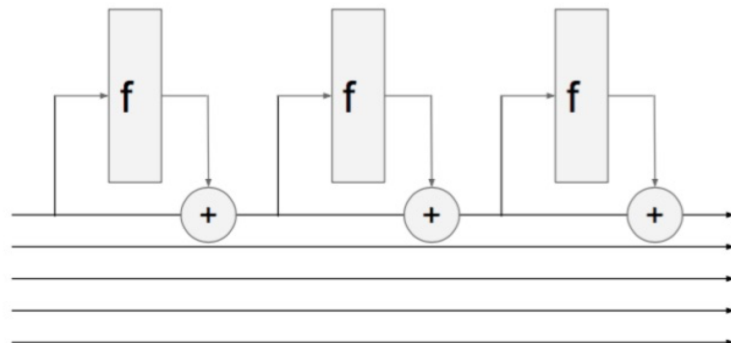# LSTM-Long Short Term Memory
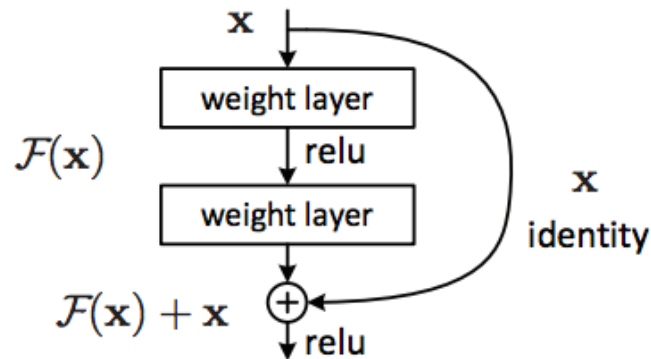


Uninterrupted gradient flow!

# LSTM: and ResNet



ResNet

Element Wise Multiplication
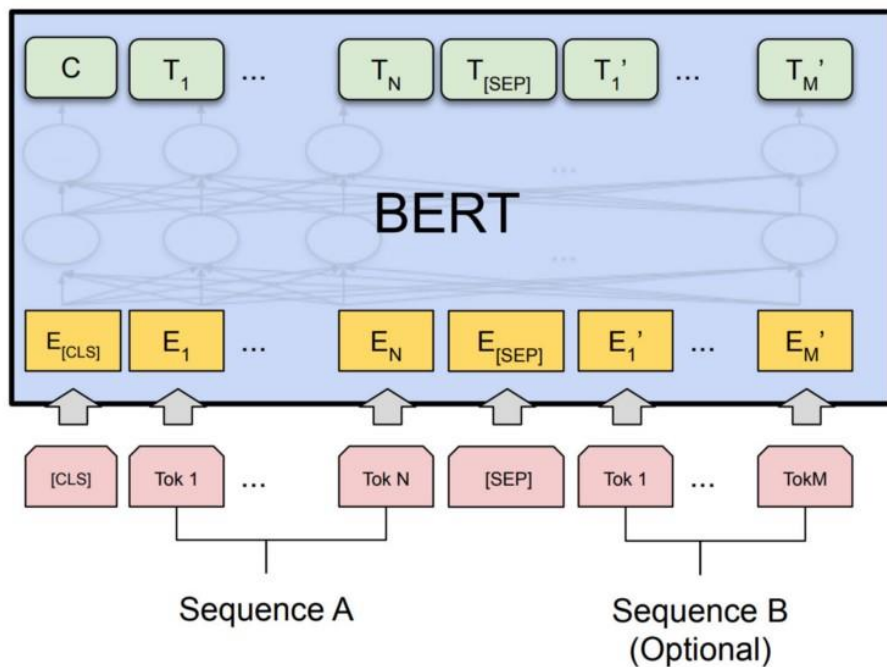▶ 매 step마다 다른 forget gate와
곱해질 수 있음

Activation을 직접적으로
transformation하는 weight 학습 X ▶
input과 ouput의 차이인 residual을 학습

# BERT: Bidirectional Encoder Representations from Transformers

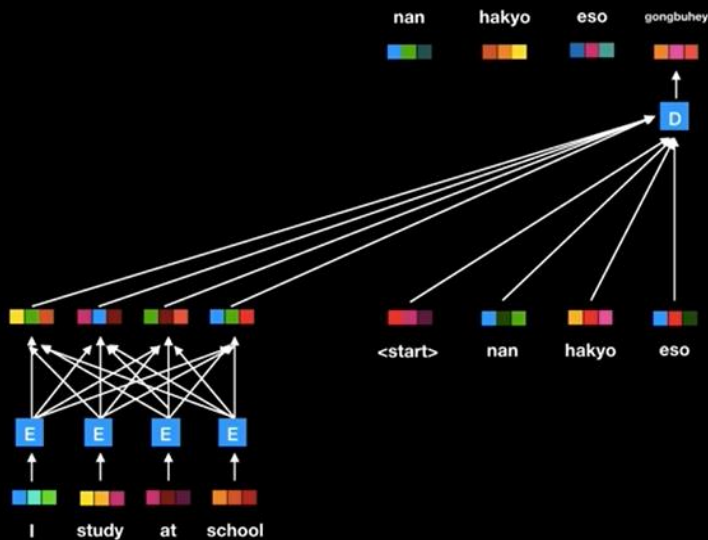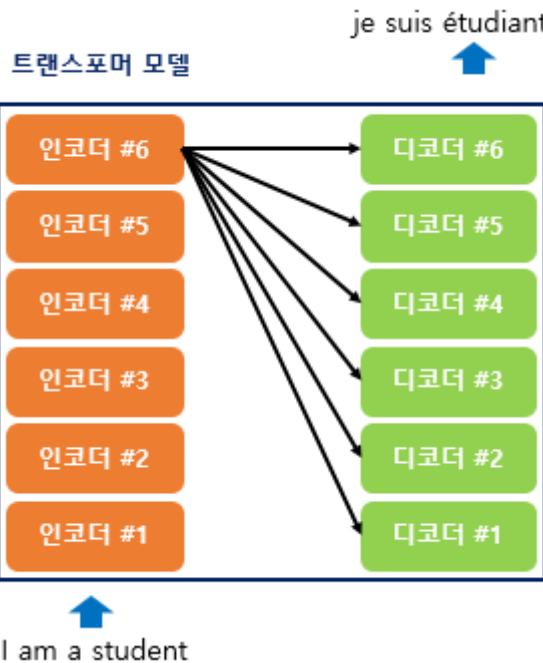## Pre-trained Model

### "Self-Attention"
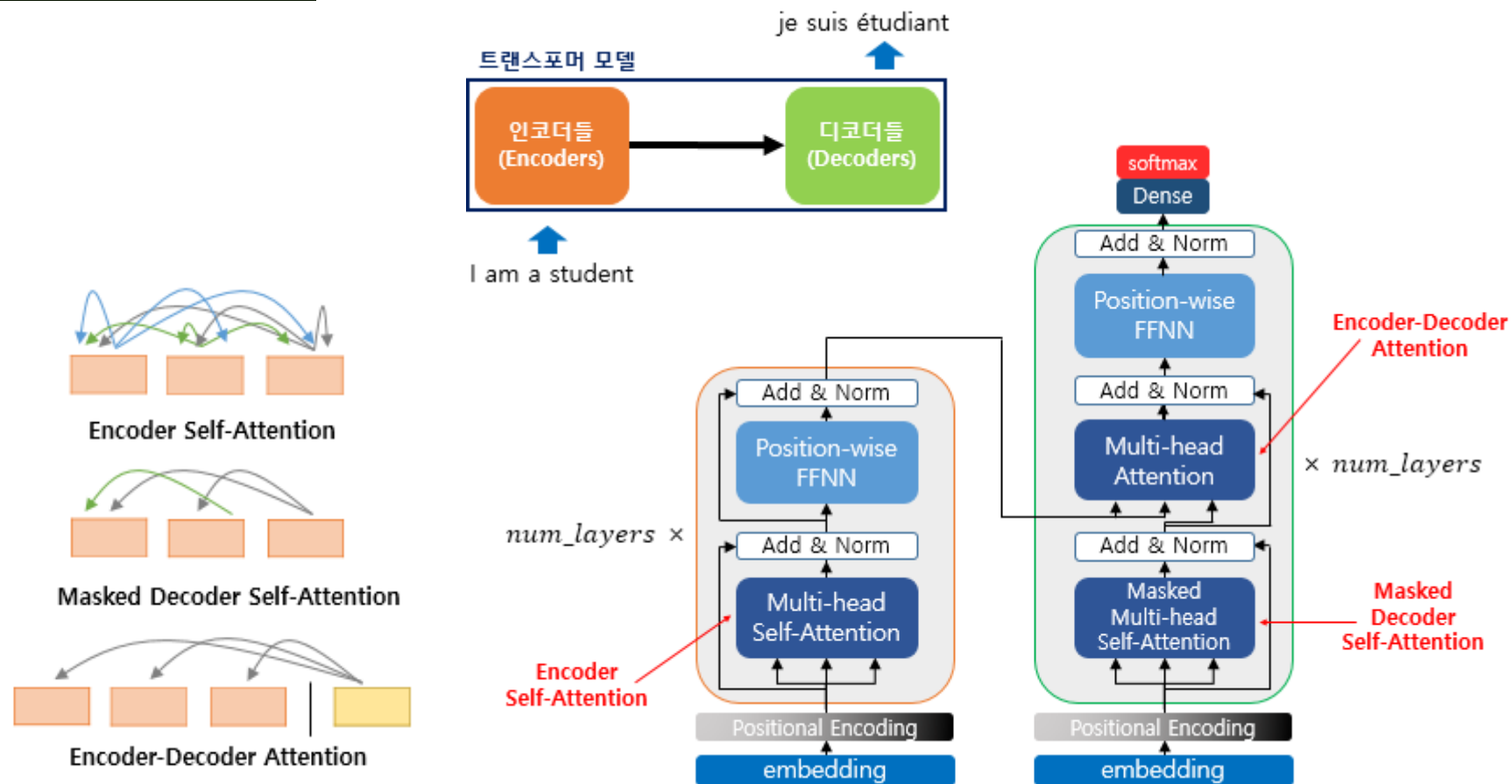
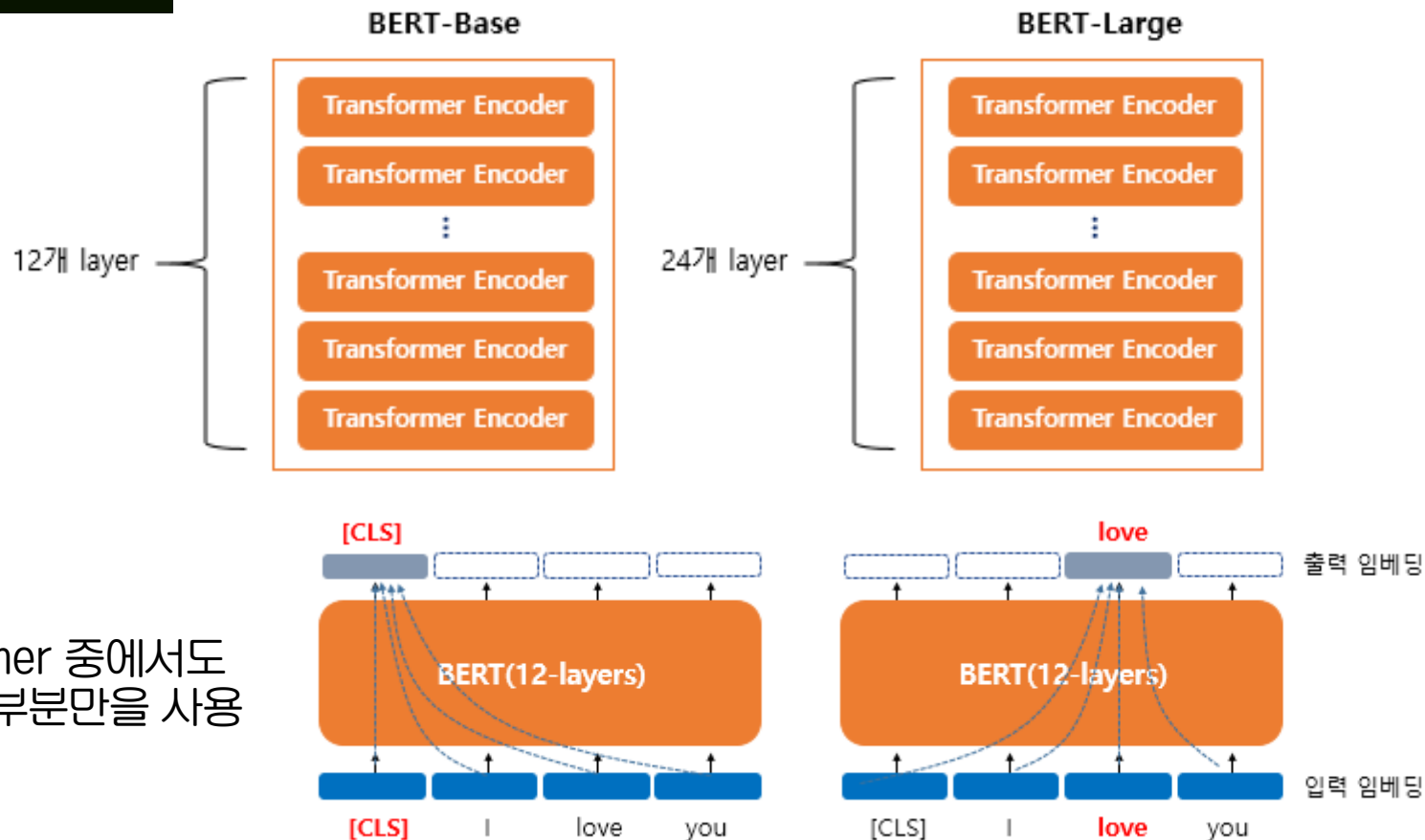# BERT-Transformer: All you need is Attention



" RNN 구조를 사용하지 않고,
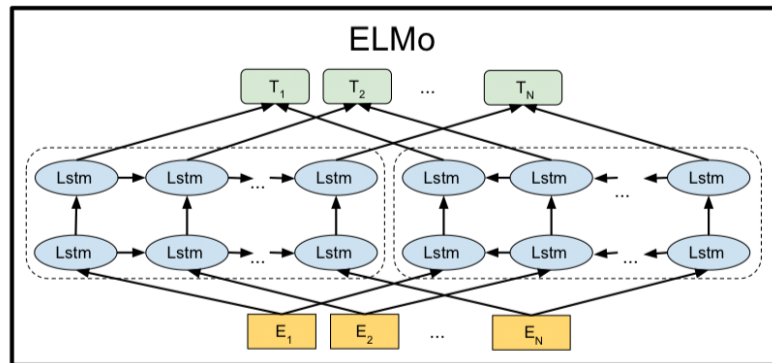Attention만을 사용해도
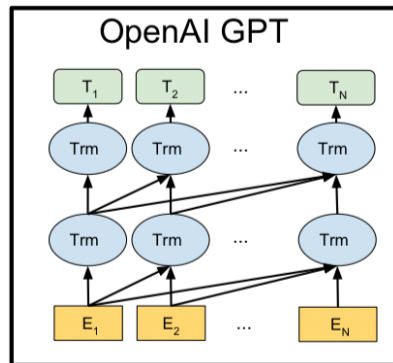원하는 결과를 얻어낼 수 있지 않을까?"

# BERT-Transformer

# BERT



transformer 중에서도
encoder 부분만을 사용

# BERT



classification token $\in \mathbb{R}^H$    $\in \mathbb{R}^H$    sentence separator    word piece embeddings (30,000 WordPiece vocabulary)

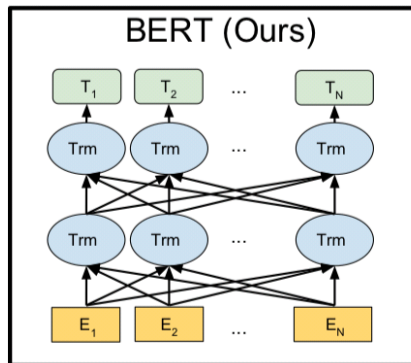| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| | + | + | + | + | + | + | + | + | + | + | + |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

# BERT

새로운 unsupervised prediction task로 pre-training을 수행

감사합니다

# Q&A