

Visual Attention

2020. 02. 14

Data Mining & Quality Analytics Lab.

강현규

Contents

- **Introduction**
- **Attention**
 - Basics
 - Visual Attention
 - Conclusion
- **Self-Attention**
 - Basics
 - Visual Self-Attention
 - Conclusion

01 | Introduction

Seminar Topic

❖ Attention: Concept

정답률: 13%

26. So far as you are wholly concentrated on bringing about a certain result, clearly the quicker and easier it is brought about the better. Your resolve to secure a sufficiency of food for yourself and your family will induce you to spend weary days in tilling the ground and tending livestock; but if Nature provided food and meat in abundance ready for the table, you would thank Nature for sparing you much labor and consider yourself so much the better off. An executed purpose, in short, is a transaction in which the time and energy spent on the execution are balanced against the resulting assets, and the ideal case is one in which _____ . Purpose, then, justifies the efforts it exacts only conditionally, by their fruits. [3점]

- ① demand exceeds supply, resulting in greater returns
 - ② life becomes fruitful with our endless pursuit of dreams
 - ③ the time and energy are limitless and assets are abundant
 - ④ Nature does not reward those who do not exert efforts
- the former approximates to zero and the latter to infinity

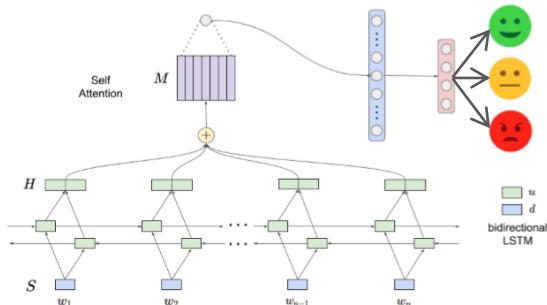
Attention Mechanism

The decoder decides parts of the source sentence to pay attention to. By letting the decoder have an **attention mechanism**, we relieve the encoder from the burden of having to encode all information in the source sentence into a **fixed length vector**. With this new approach the information can be spread throughout the sequence of annotations, which can be selectively retrieved by the decoder accordingly.

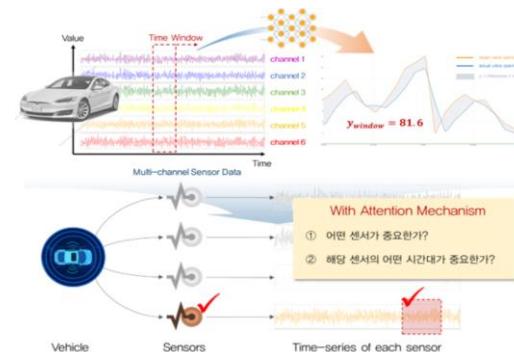
01 | Introduction Seminar Topic

❖ Attention: 19's DMQA

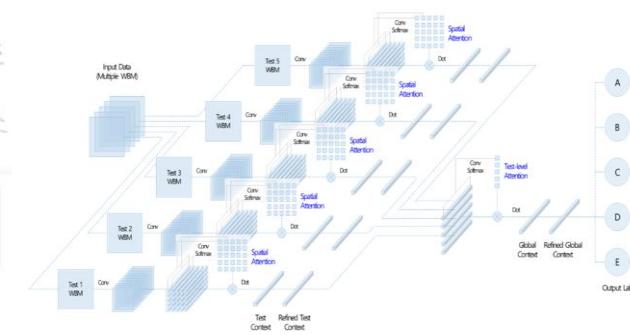
① Sentiment Analysis



② Anomaly Detection



③ Wafer Bin Map Classification



01 | Introduction Seminar Topic

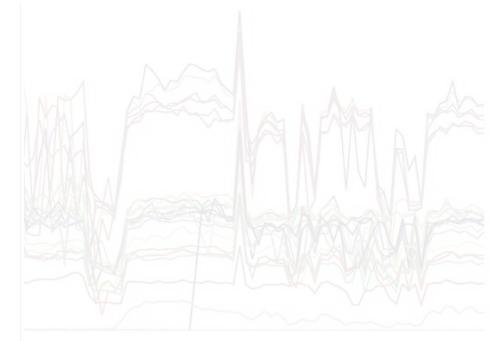
❖ Attention: Application

① Natural Language Processing



- Document Classification
- Machine Translation
- Question & Answering
- Document Summarization

② Multivariate Time Series



- Anomaly Detection
- Predictive Modeling

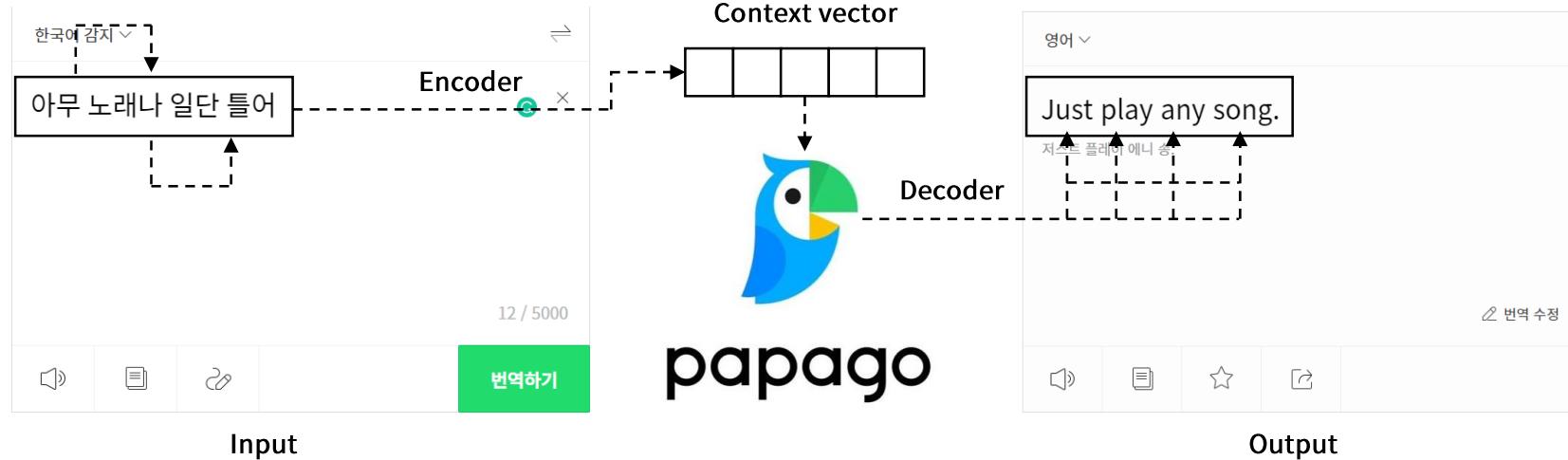
③ Computer Vision



- Image Classification
- Image Captioning
- Object Detection
- Visual QA

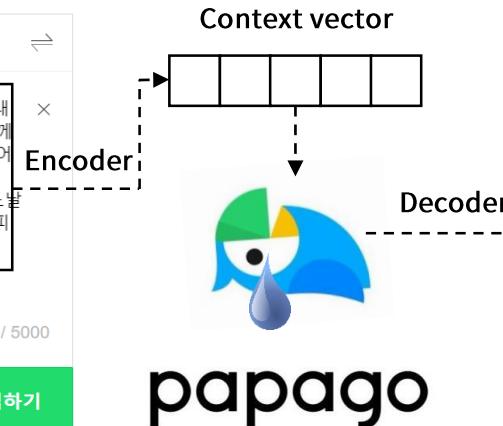
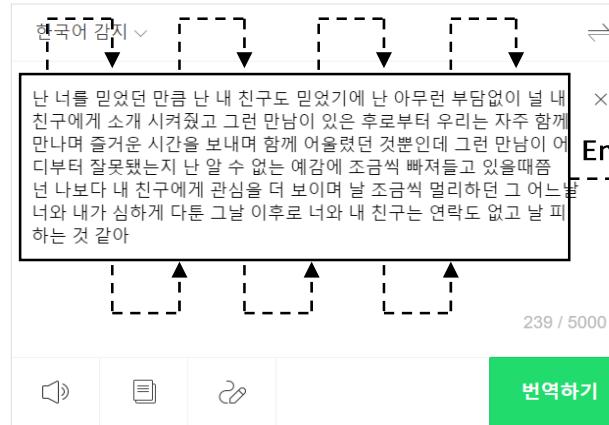
02 | Attention Basics

❖ Machine Translation: Seq2seq

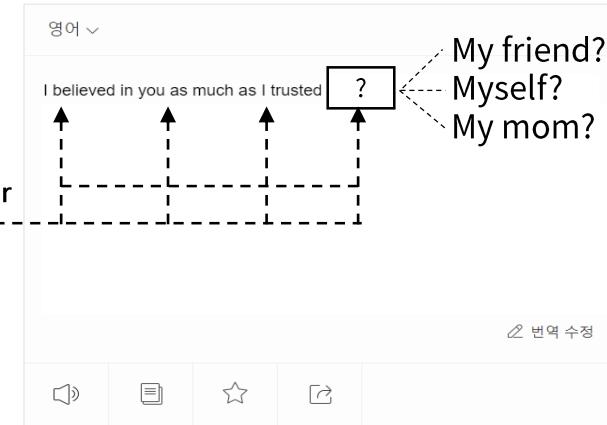


02 | Attention Basics

❖ Machine Translation: Seq2seq



Input



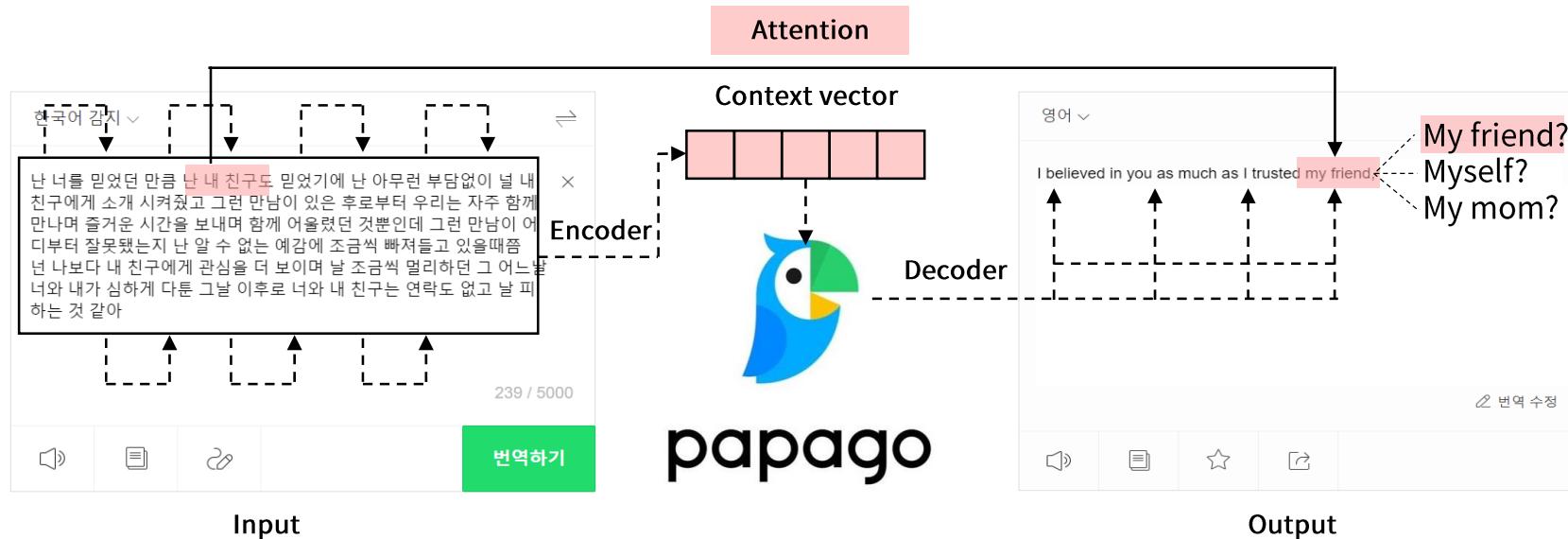
Output

▪ Long input & output sequence

- Long term dependency
- Vanishing gradient
- Limitation of fixed context vector

02 | Attention Basics

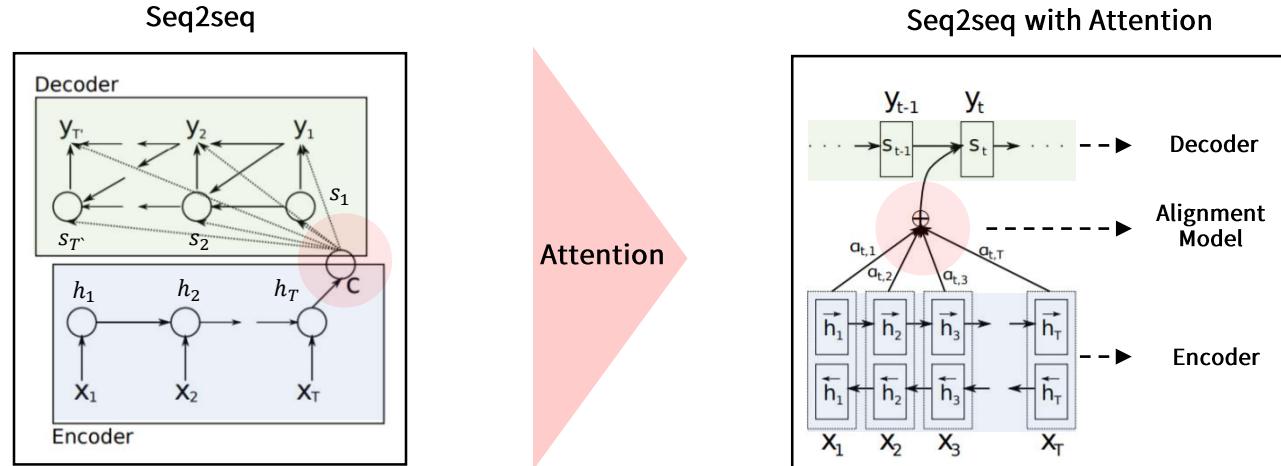
❖ Machine Translation: Seq2seq



▪ Attention

- Decoder가 특정 시점의 단어를 출력할 때 encoder의 정보 중 연관성이 있는 부분에 주목하여 adaptive context vector를 만든다.
= Relieve the encoder from the burden of having to encode all information into a fixed length vector.

❖ Machine Translation: Seq2seq with Attention



$$p(y_t | y_{t-1}, \dots, y_1, c) = g(s_t, y_{t-1}, c)$$

$$p(y_t | y_{t-1}, \dots, y_1, c) = g(s_t, y_{t-1}, c_t), c_t = \sum_{i=1}^T a_i \vec{h}_i$$

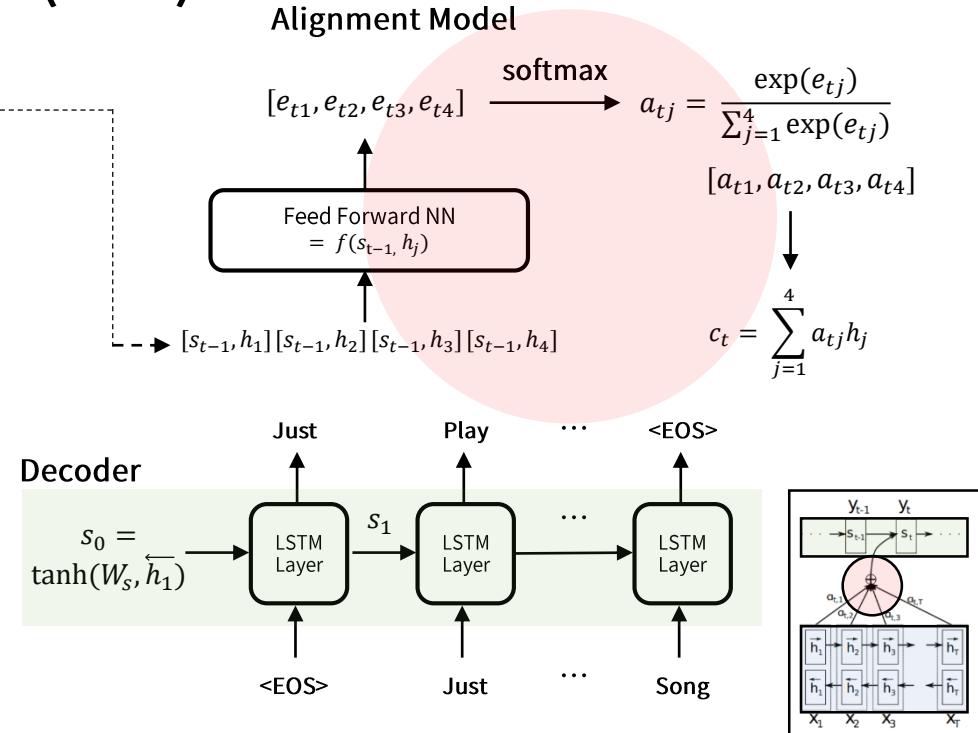
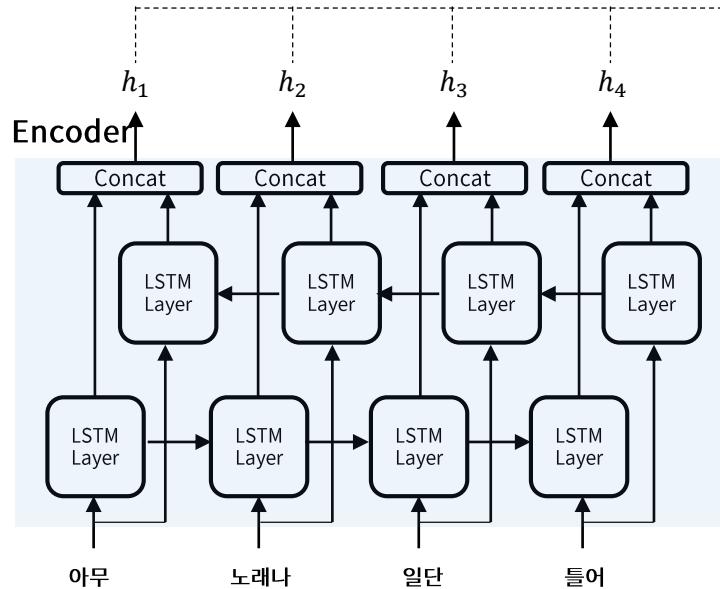
▪ Attention

- Decoder가 특정 시점의 단어를 출력할 때 encoder의 정보 중 연관성이 있는 부분에 주목하여 adaptive context vector를 만든다.
 = Relieve the encoder from the burden of having to encode all information into a fixed length vector.

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *arXiv preprint arXiv:1409.0473* (2014).

02 | Attention Basics

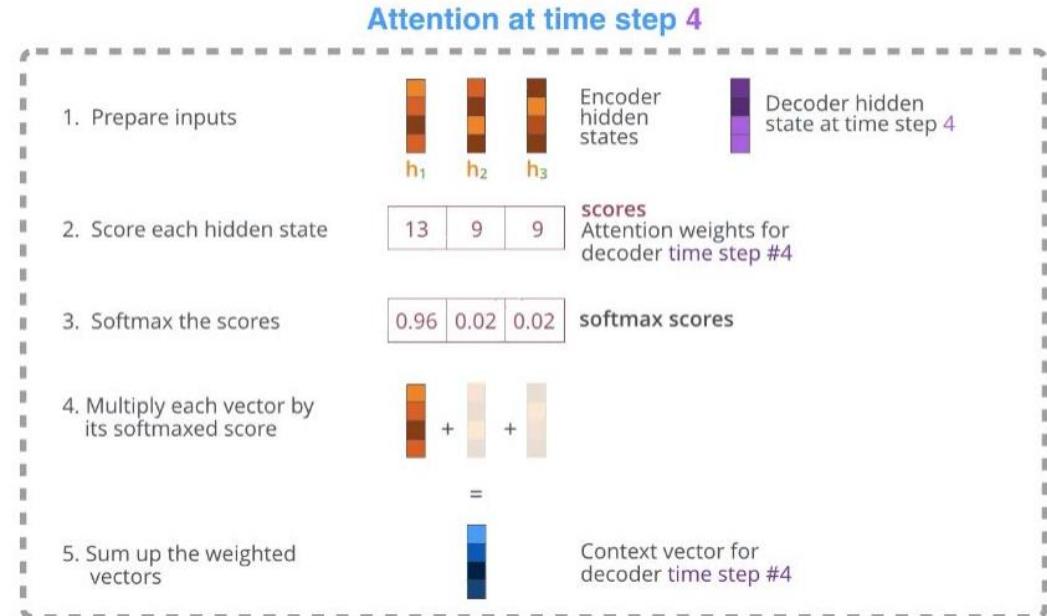
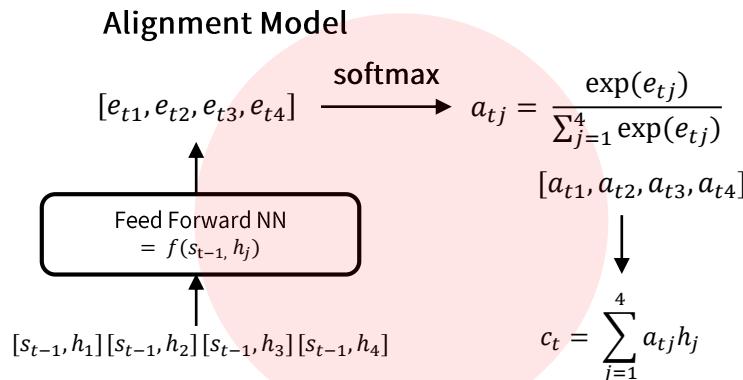
❖ Seq2seq with Attention – Bahdanau(2014)



Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *arXiv preprint arXiv:1409.0473* (2014).

02 | Attention Basics

❖ Seq2seq with Attention – Bahdanau(2014)



<https://jamalmar.github.io/visualizing-neural-machine-translation-mechanics-of-seq2seq-models-with-attention/>

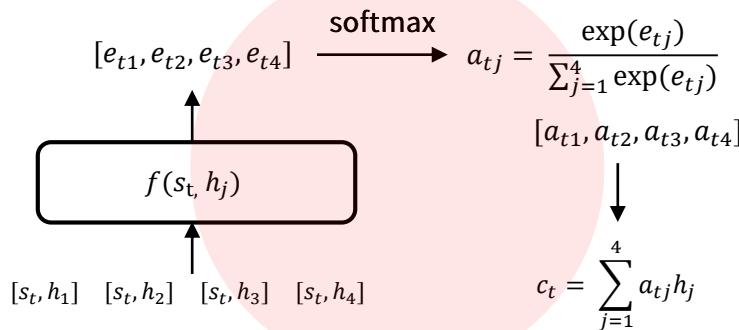
Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *arXiv preprint arXiv:1409.0473* (2014).

02 | Attention Basics

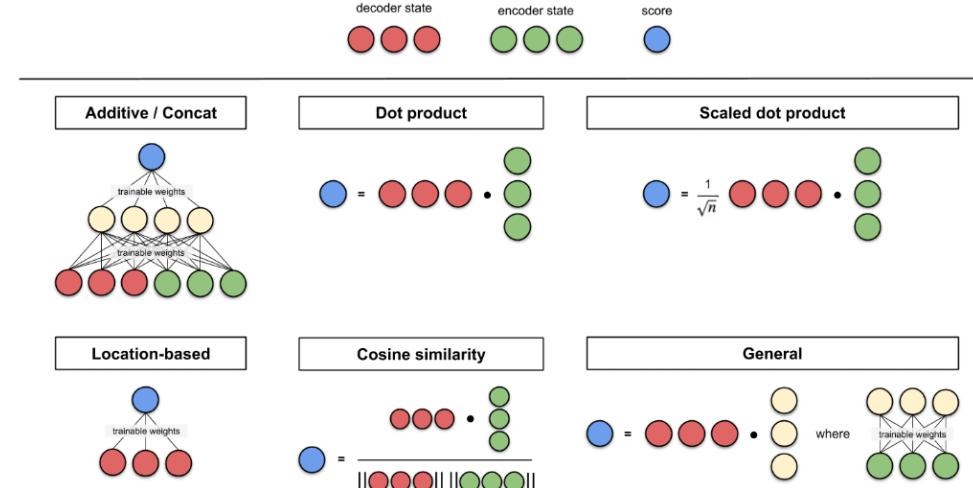
❖ Seq2seq with Attention – Luong(2015)

- Alignment model의 input으로 s_{t-1} 가 아닌 s_t 을 사용
- 다양한 similarity function $f(s_t, h_j)$ 제시

Alignment Model



$$f(s_t, h_j) = e_j$$



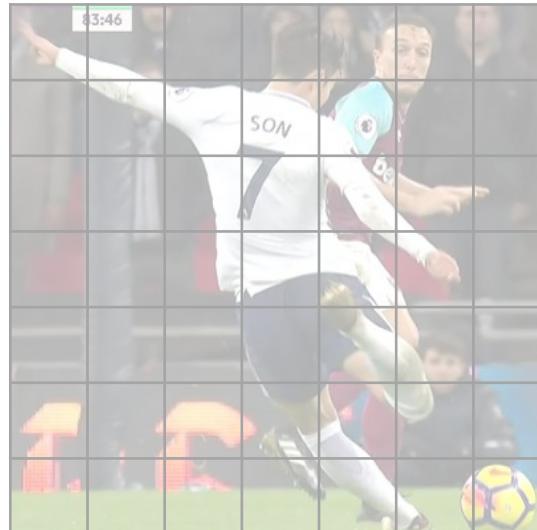
<https://towardsdatascience.com/attn-illustrated-attention-5ec4ad276ee3>

Luong, Minh-Thang, Hieu Pham, and Christopher D. Manning. "Effective approaches to attention-based neural machine translation." *arXiv preprint arXiv:1508.04025* (2015).

❖ Image Captioning

- 이미지를 입력 데이터로 사용하여 이미지를 설명하는 문장을 생성하는 문제

Input: Image



Output: Text

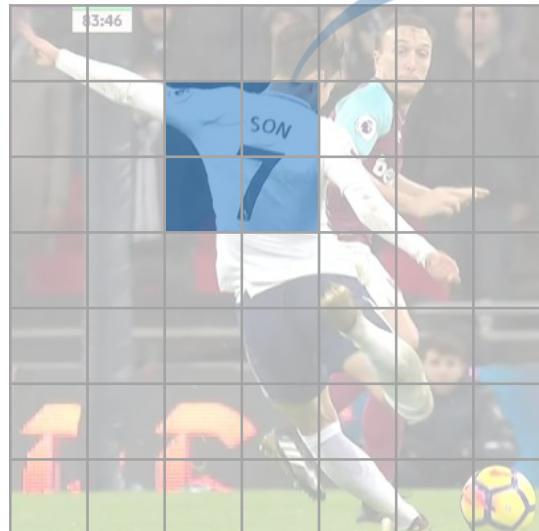
“Sonny is shooting in
front of the defender”

02 | Attention Visual Attention

❖ Image Captioning

- 이미지를 입력 데이터로 사용하여 이미지를 설명하는 문장을 생성하는 문제

Input: Image



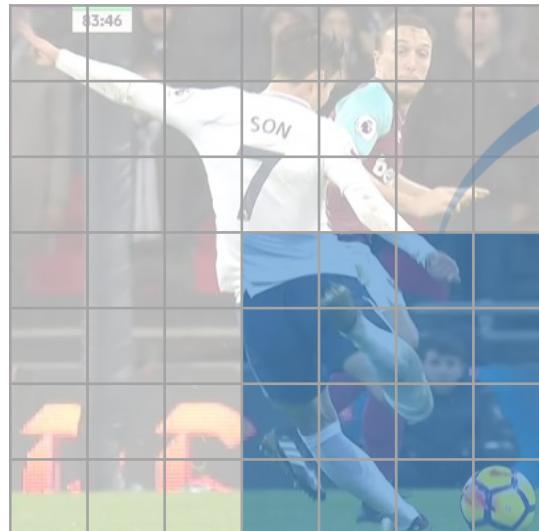
Output: Text

“Sonny is shooting in front of the defender”

❖ Image Captioning

- 이미지를 입력 데이터로 사용하여 이미지를 설명하는 문장을 생성하는 문제

Input: Image



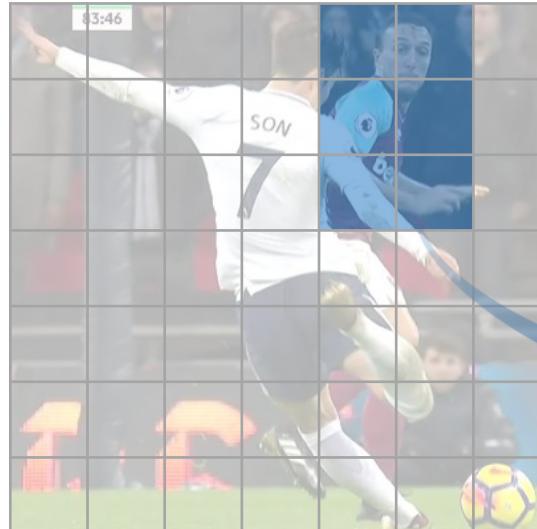
Output: Text

“Sonny is shooting in
front of the defender”

❖ Image Captioning

- 이미지를 입력 데이터로 사용하여 이미지를 설명하는 문장을 생성하는 문제

Input: Image



Output: Text

“Sonny is shooting in
front of the **defender**”



❖ Paper Review

- Show, Attend and Tell
- Conference on machine learning
- Published in 2015
- Xu, Kelvin, et al.
- Kyunghyun Cho, Yoshua Bengio…

[PDF] [Show, attend and tell: Neural image caption generation with visual attention](#)

[K Xu, J Ba, R Kiros, K Cho, A Courville... - ... conference on machine ... , 2015 - jmlr.org](#)

Inspired by recent work in machine translation and object detection, we introduce an attention based model that automatically learns to describe the content of images. We describe how we can train this model in a deterministic manner using standard ...

☆ 4286회 인용 관련 학술자료 전체 24개의 버전 ☆

Show, Attend and Tell: Neural Image Caption Generation with Visual Attention

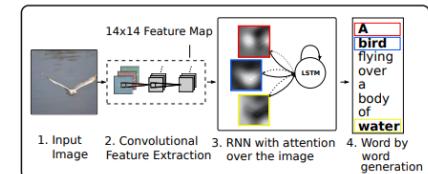
Kelvin Xu
Jimmy Lei Ba
Ryan Kiros
Kyunghyun Cho
Aaron Courville
Ruslan Salakhutdinov
Richard S. Zemel
Yoshua Bengio

KELVIN.XU@UMONTREAL.CA
JIMMY@PSI.UTORONTO.CA
RKIROS@CS.TORONTO.EDU
KYUNGHYUN.CHO@UMONTREAL.CA
AARON.COURVILLE@UMONTREAL.CA
RSALAKHUU@CS.UTORONTO.EDU
ZEMEL@CS.UTORONTO.EDU
FIND-ME@THE.WEB

Abstract

Inspired by recent work in machine translation and object detection, we introduce an attention based model that automatically learns to describe the content of images. We describe how we can train this model in a deterministic manner using standard backpropagation techniques and stochastically by maximizing a variational lower bound. We also show through visualization how the model is able to automatically learn to fix its gaze on salient objects while generating the corresponding words in the output sequence. We validate the use of attention with state-of-the-art performance on three benchmark datasets: Flickr8k, Flickr30k and MS COCO.

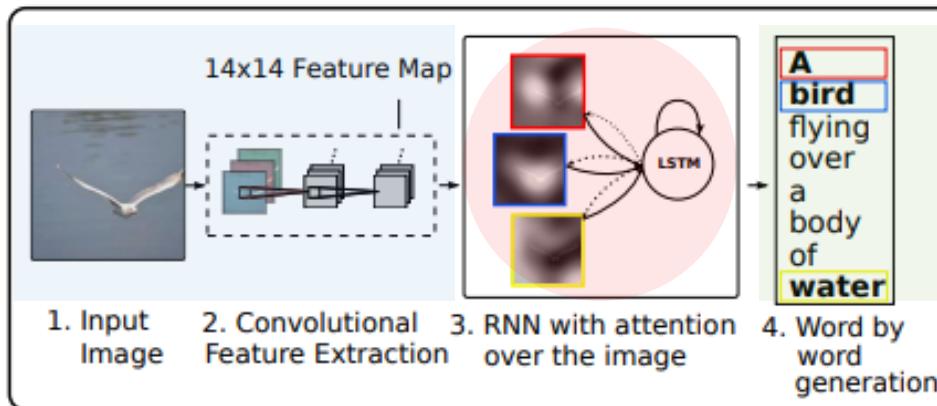
Figure 1. Our model learns a words/image alignment. The visualized attentional maps (3) are explained in section 3.1 & 5.4



has significantly improved the quality of caption generation using a combination of convolutional neural networks (convnets) to obtain vectorial representation of images and recurrent neural networks to decode those representations

02 | Attention Visual Attention

❖ Model Architecture



Encoder – Feature Extraction

Convolutional Neural Network 모델으로 이미지의 특징을 잘 요약하는 feature map 생성

Attention

Decoder에서 단어를 출력할 때 단어와 가장 유사한 feature map과의 attention score을 계산

Decoder – Word Generation

이전 시점에 출력된 단어의 정보, feature map, attention score를 decoder의 input으로 사용하여 순차적으로 단어를 출력

02 | Attention Visual Attention

❖ Model Architecture



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



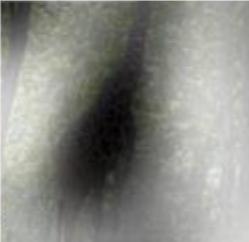
A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.

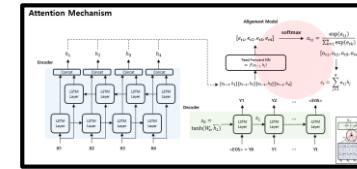
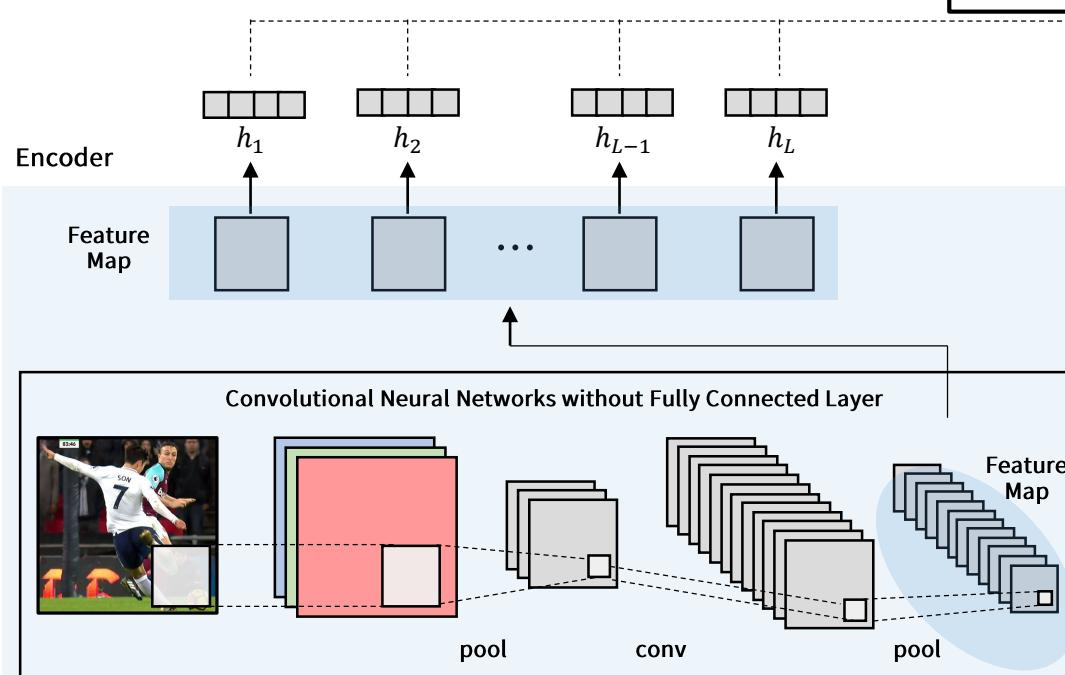


A giraffe standing in a forest with trees in the background.



02 | Attention Visual Attention

❖ Model Architecture



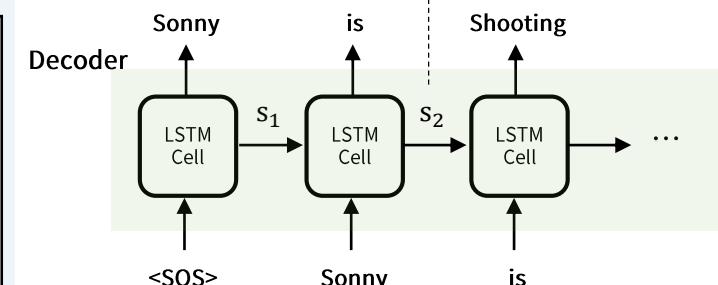
Attention

$$e_{ti} = f(h_i, s_t)$$

$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

$$\hat{z}_t = \phi(\{h_i\}, \{\alpha_i\})$$

$$p(y_t | fm, y_1, \dots, y_{t-1}) = g(s_t, y_{t-1}, \hat{z}_t)$$



Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." *International conference on machine learning*. 2015.

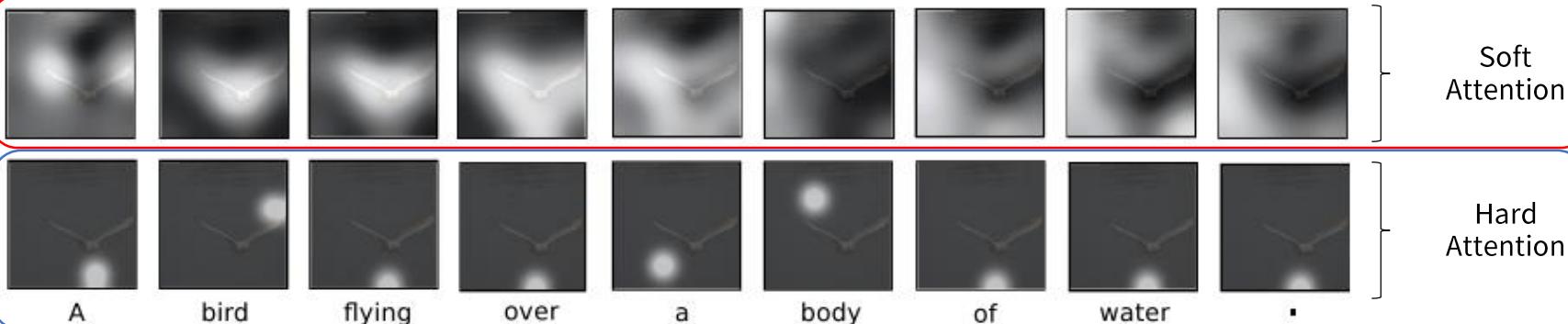
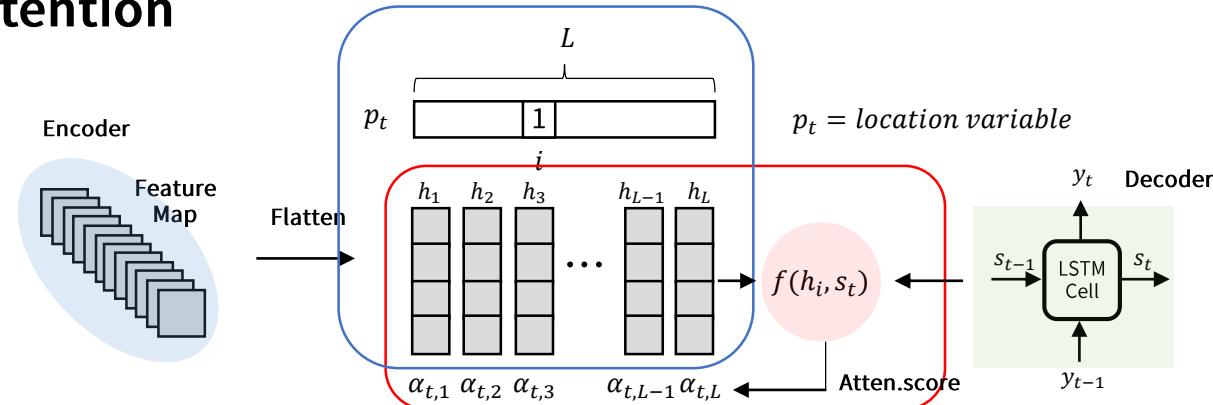
02 | Attention

Visual Attention

❖ Soft Attention vs Hard Attention

Soft Attention $\hat{z}_t = \phi(\{h_i\}, \{\alpha_i\}) = \sum_{k=1}^L \alpha_i h_i$

Hard Attention $\hat{z}_t = \phi(\{h_i\}, \{\alpha_i\}) = \sum_i p_{t,i} h_i$
 $p(p_{t,i} = 1 | p_{j < t}, h) = \alpha_{t,i}$



Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." *International conference on machine learning*. 2015.

❖ Soft Attention vs Hard Attention

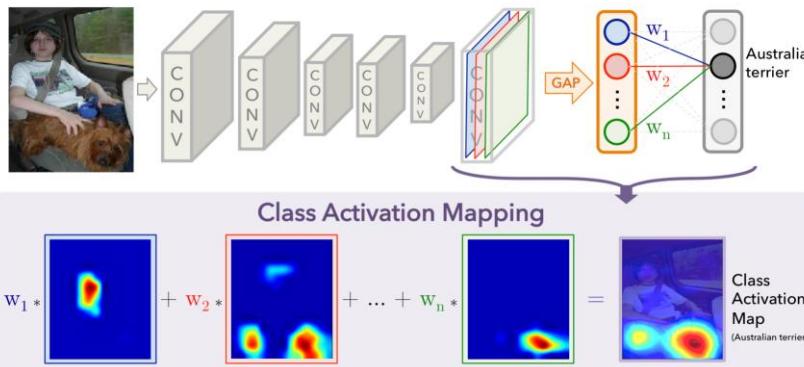
	Soft Attention	Hard Attention
Pros	<ul style="list-style-type: none">• Interpretability• End to end learning	<ul style="list-style-type: none">• Better performance than soft attention
Cons	<ul style="list-style-type: none">• Worse performance than hard attention	<ul style="list-style-type: none">• Hard to optimize<ul style="list-style-type: none">- Monte Carlo based sampling- REINFORCE

Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." *International conference on machine learning*. 2015.

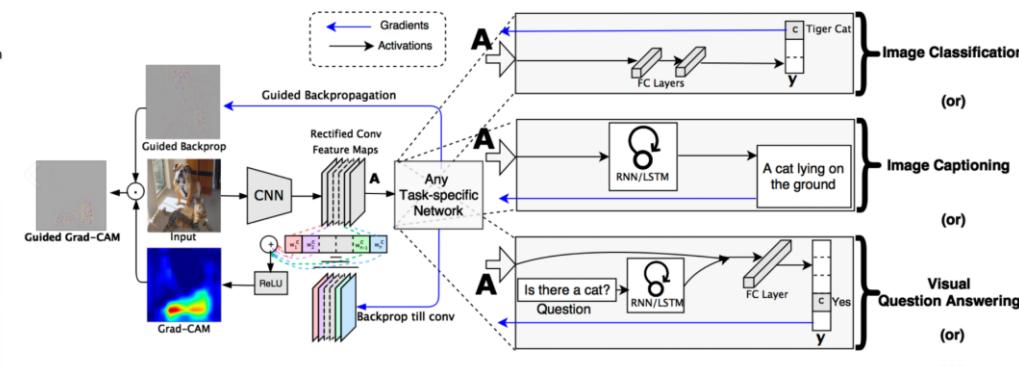
02 | Attention Conclusion

❖ Localization vs Attention

- Localization: Interpretability
- Attention: Interpretability + Model performance



Class Activation Map(CAM)

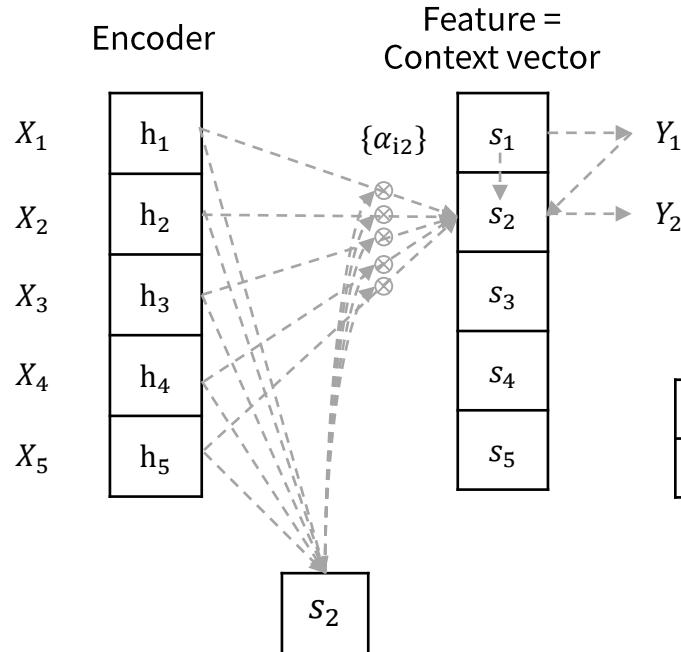


Grad_CAM

Zhou, Bolei, et al. "Learning deep features for discriminative localization." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.
Selvaraju, Ramprasaath R., et al. "Grad-cam: Visual explanations from deep networks via gradient-based localization." *Proceedings of the IEEE international conference on computer vision*. 2017.

02 | Attention Conclusion

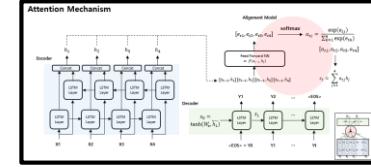
❖ Feature Representation



Feature =
Context vector

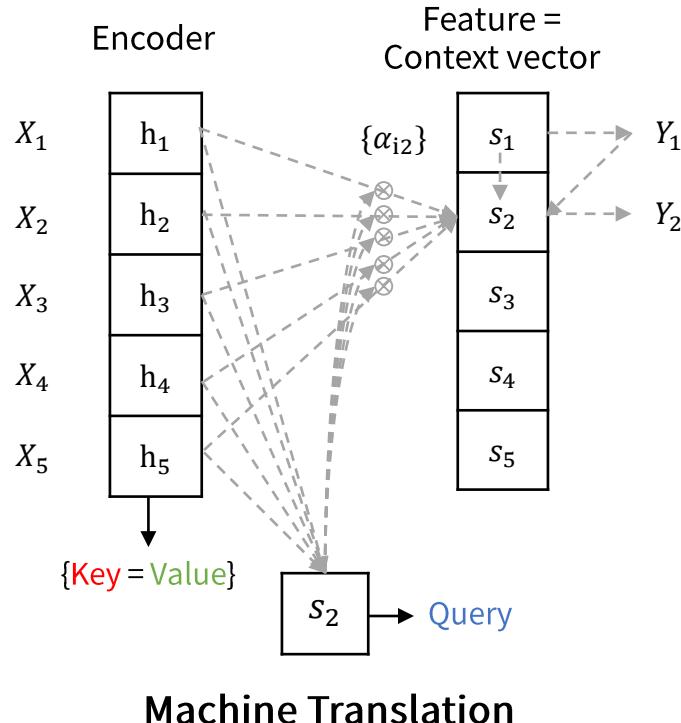
$$\{\alpha_{i2}\} = \text{softmax}(f(h_i, s_2))$$

$$\text{feature} = \sum_{i=1}^5 \alpha_{i2} * h_i$$



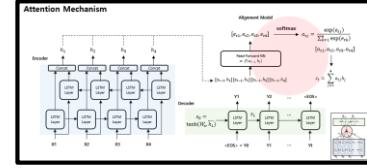
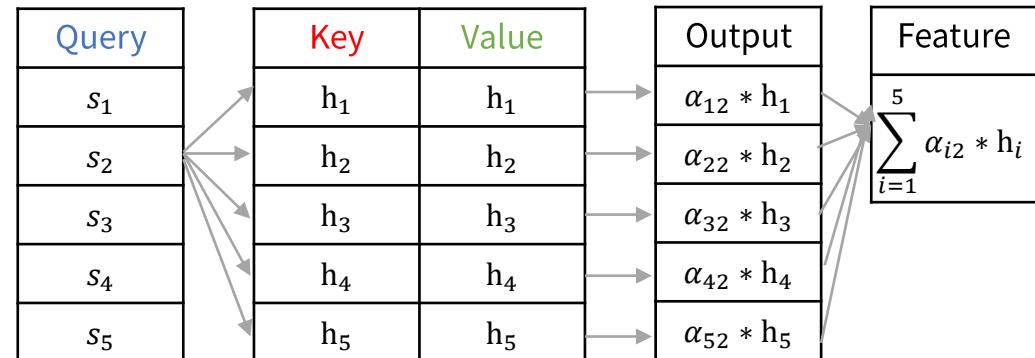
02 | Attention Conclusion

❖ Feature Representation



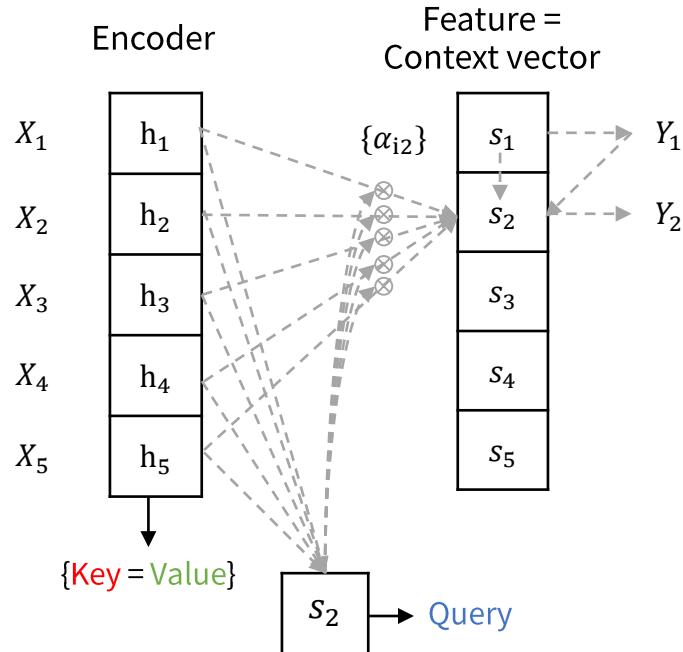
$$\{\alpha_{i2}\} = \text{softmax}(f(\mathbf{h}_i, \mathbf{s}_2))$$

$$\text{feature} = \sum_{i=1}^5 \alpha_{i2} * \mathbf{h}_i$$



02 | Attention Conclusion

❖ Feature Representation



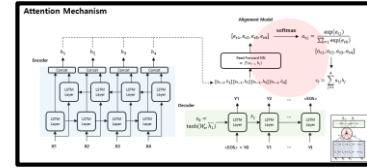
Machine Translation

$$\{\alpha_{i2}\} = \text{softmax}(f(\mathbf{h}_i, \mathbf{s}_2))$$

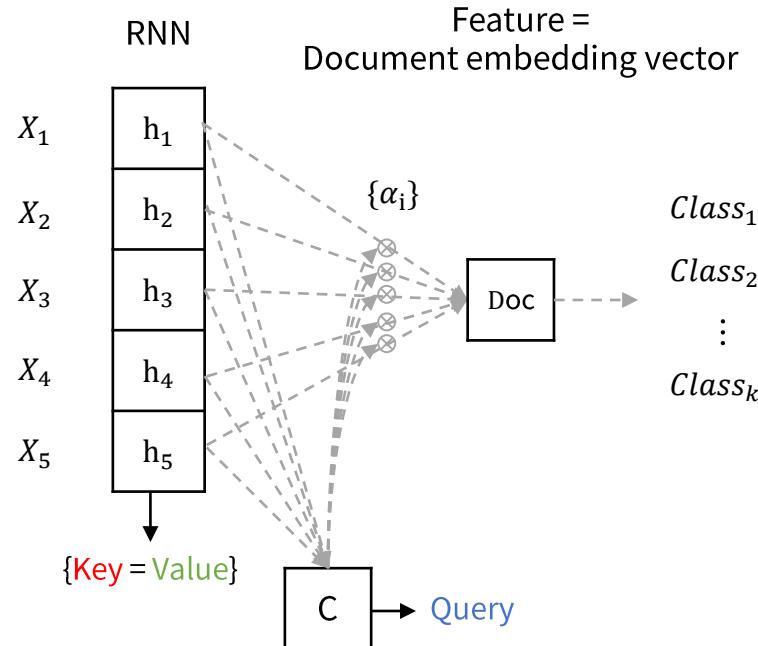
$$feature = \sum_{i=1}^5 \alpha_{i2} * \mathbf{h}_i$$



$$A(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_i \text{softmax}(f(\mathbf{K}, \mathbf{q})) \mathbf{V}$$



❖ Feature Representation



$$\{\alpha_i\} = \text{softmax}(f(\mathbf{h}_i, \mathbf{c}))$$

$$\text{feature} = \sum_{i=1}^5 \alpha_i * \mathbf{h}_i$$



$$A(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_i \text{softmax}(f(\mathbf{K}, \mathbf{q})) \mathbf{V}$$

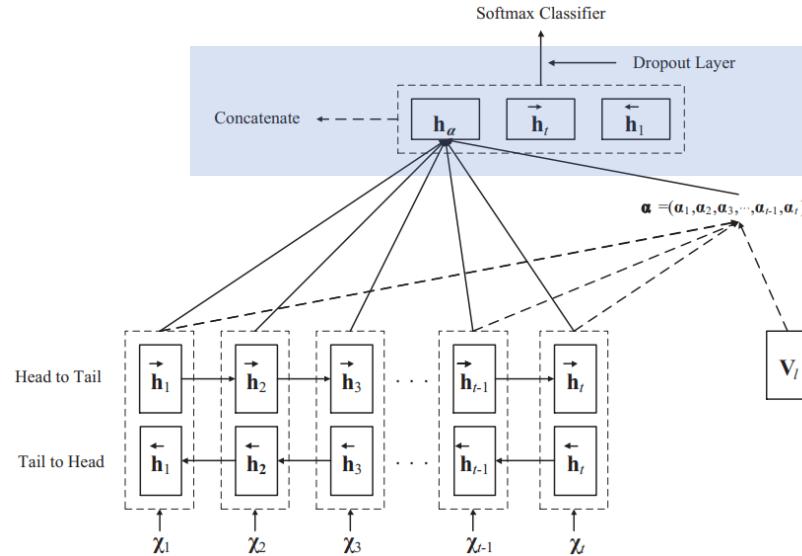
Document Classification

Generalized Attention Form

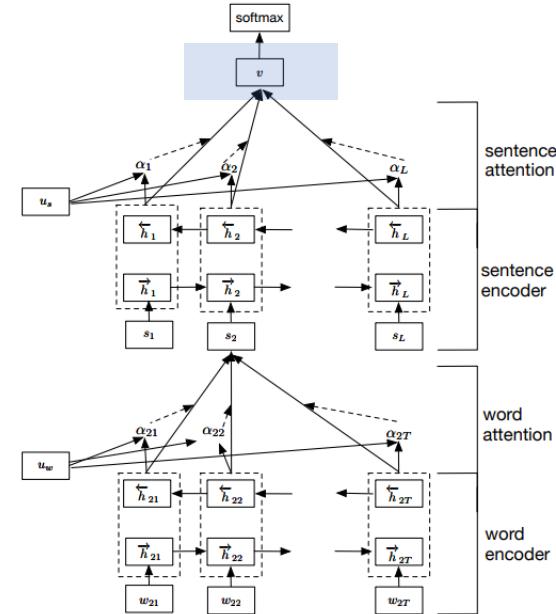
02 | Attention Conclusion

❖ Document Classification: RNN-based

- 문서 또는 문장을 보다 더 잘 표현하는 feature를 생성



Bidirectional RNN with Attention



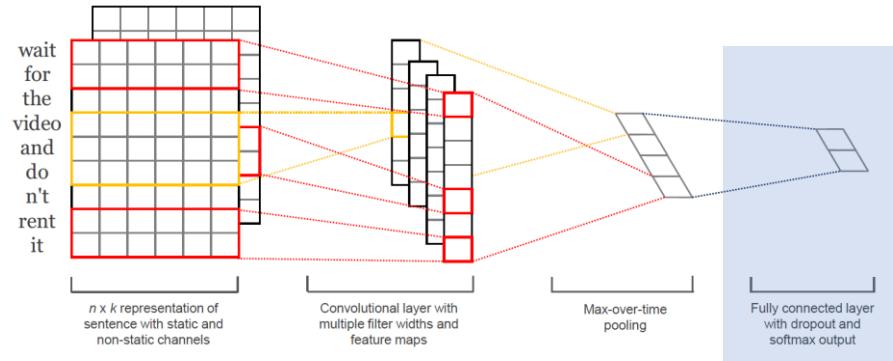
Hierarchical Attention Network

Liu, Tengfei, et al. "Recurrent networks with attention and convolutional networks for sentence representation and classification." *Applied Intelligence* 48.10 (2018): 3797-3806.
Yang, Zichao, et al. "Hierarchical attention networks for document classification." *Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics 2016*.

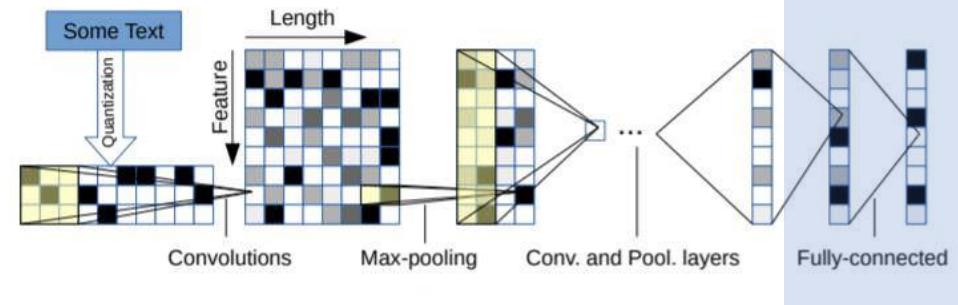
02 | Attention Conclusion

❖ Document Classification: CNN-based

- 문서 또는 문장을 보다 더 잘 표현하는 feature를 생성



TextCNN

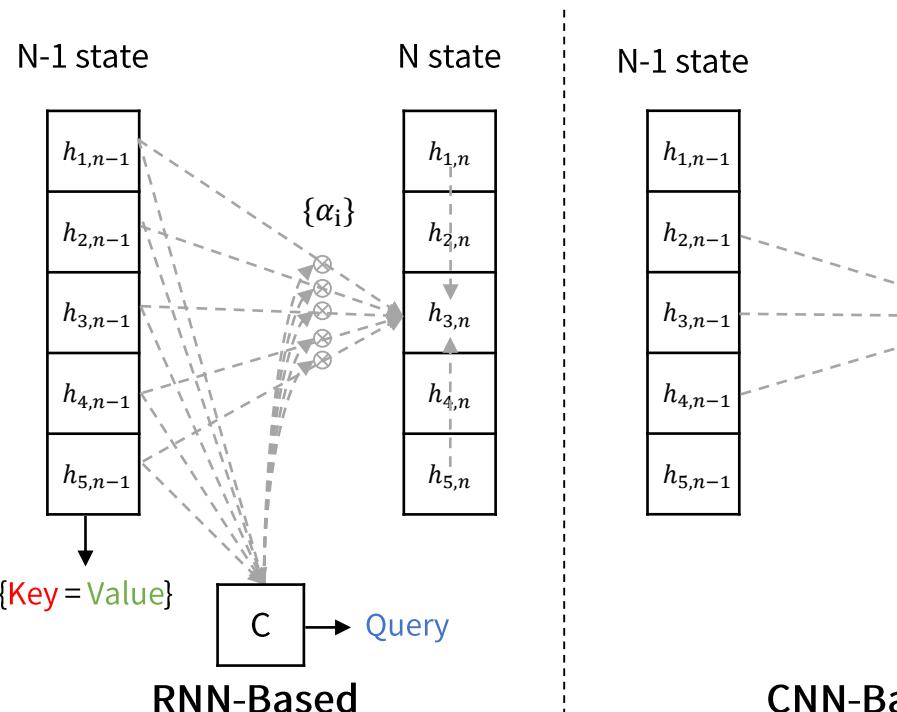


Character-level CNN

Kim, Yoon. "Convolutional neural networks for sentence classification." *arXiv preprint arXiv:1408.5882* (2014).
Zhang, Xiang, Junbo Zhao, and Yann LeCun. "Character-level convolutional networks for text classification." *Advances in neural information processing systems*. 2015.

03 | Self-Attention Basics

❖ Attention to Self-Attention



$$A(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_i softmax(f(\mathbf{K}, \mathbf{q})) \mathbf{V}$$

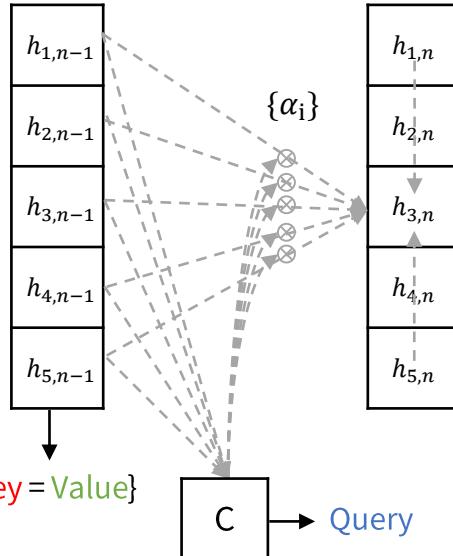
- **RNN-based**
 - 입력 시퀀스를 순차적으로 처리하여 병렬처리 어려움
 - 연산 시간, 계산 복잡도 多
- **CNN-based**
 - 병렬처리는 쉽지만 거리가 먼 단어와의 관계를 학습하기에는 여러 합성곱 연산 필요
- 더 나은 feature 계산 방법?
 - RNN, CNN 구조 사용 X
 - Key = Value = Query
 - $f = \text{dot product}$ 사용

03 | Self-Attention Basics

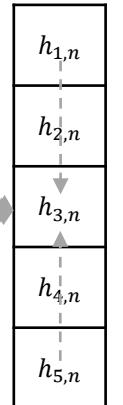
❖ Attention to Self-Attention

$$A(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_i softmax(f(\mathbf{K}, \mathbf{q})) \mathbf{V}$$

N-1 state



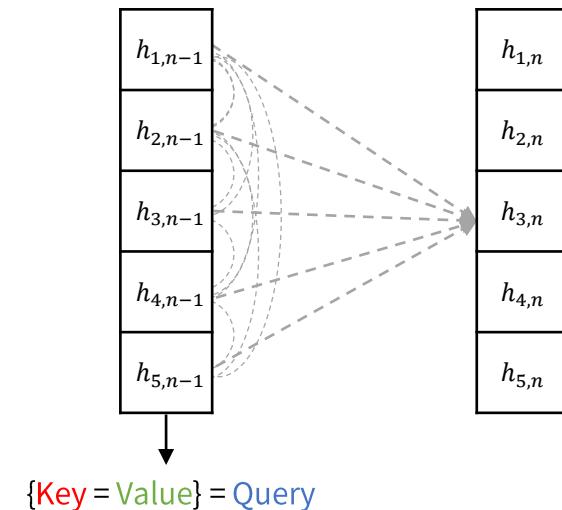
N state



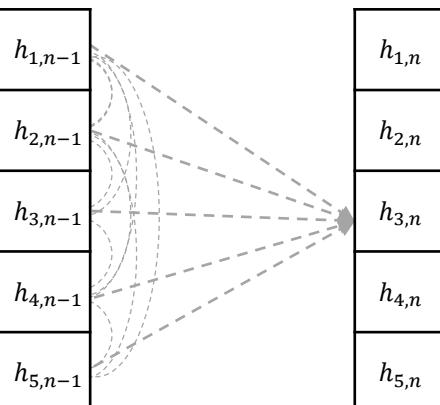
RNN-Based

CNN-Based

N-1 state



N state



Self-Attention

❖ Paper Review

- Attention is All You Need
- NIPS
- Published in 2017
- A Vaswani, et al, Google Brain, Research
- A Vaswani, N Parmar...

Attention is all you need

[A Vaswani](#), [N Shazeer](#), [N Parmar](#)... - *Advances in neural ...*, 2017 - [papers.nips.cc](#)

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks in an encoder and decoder configuration. The best performing such models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely.

☆ 99 5899회 인용 관련 학술자료 전체 20개의 버전 ☺

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

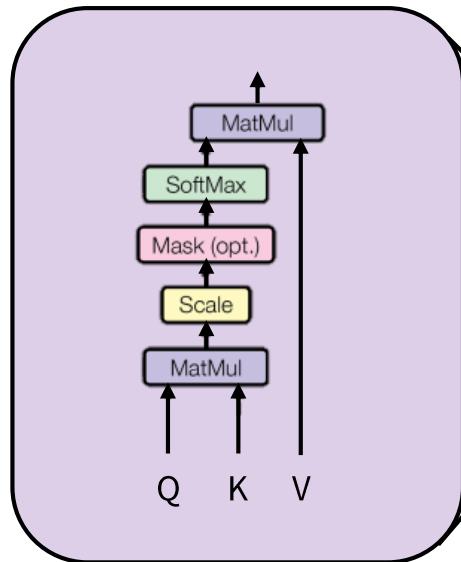
Illia Polosukhin* §
illia.polosukhin@gmail.com

Abstract

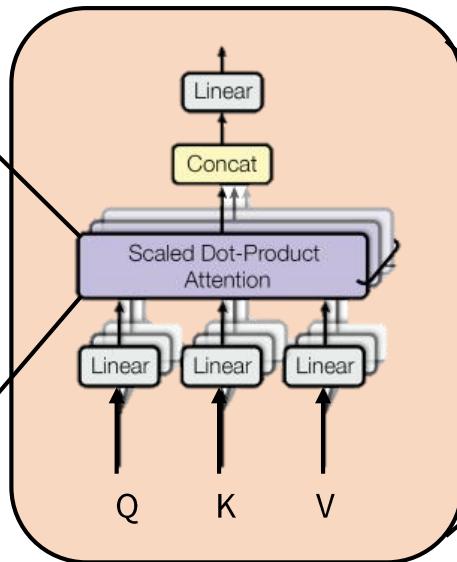
The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.0 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature.

03 | Self-Attention Basics

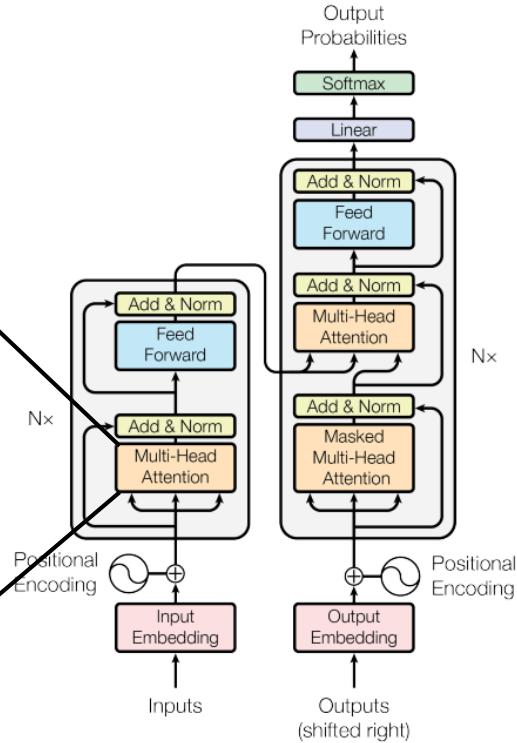
❖ Transformer



Scaled Dot-Product Attention



Multi-Head Attention



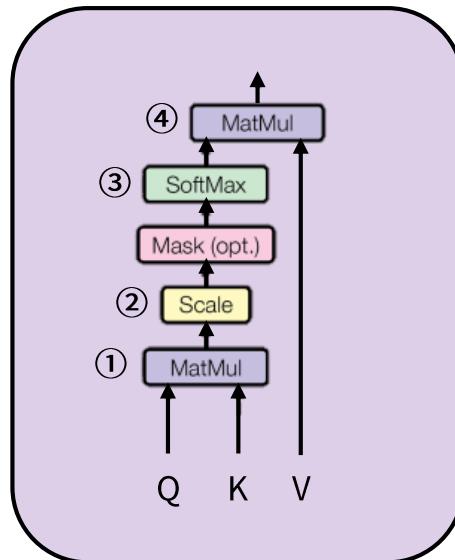
Transformer

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

<http://jalammar.github.io/>

03 | Self-Attention Basics

❖ Transformer



Scaled Dot-Product
Attention

Generalized
Attention Form

$$A(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_i \text{softmax}(f(\mathbf{K}, \mathbf{q})) \mathbf{V}$$

①

MatMul

$$f(\mathbf{K}, \mathbf{Q}) = \mathbf{Q}\mathbf{K}^T$$

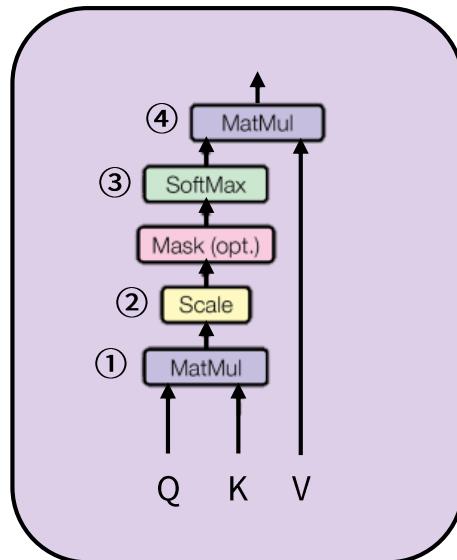
While the two are similar in theoretical complexity, dot-product attention is much faster and more space-efficient in practice, since it can be implemented using highly optimized matrix multiplication code. (…)

<http://jalammar.github.io/>

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

03 | Self-Attention Basics

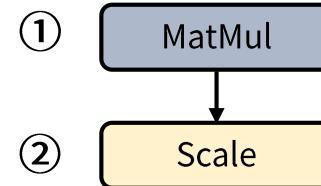
❖ Transformer



Scaled Dot-Product
Attention

Generalized
Attention Form

$$A(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_i softmax(f(\mathbf{K}, \mathbf{q})) \mathbf{V}$$



$$f(\mathbf{K}, \mathbf{Q}) = \mathbf{Q}\mathbf{K}^T$$

$$\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}$$

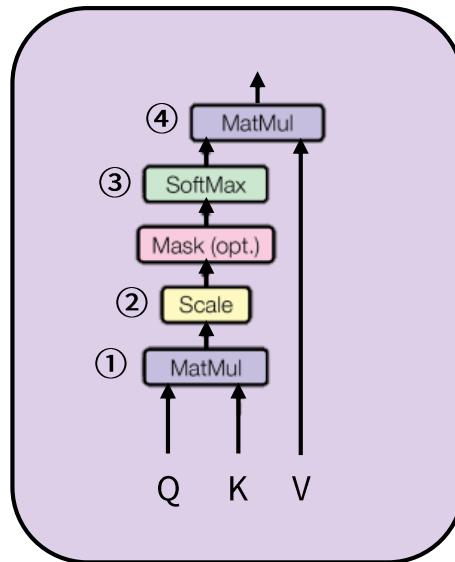
We suspect that for large values of d_k , the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients. To counteract this effect, we scale the dot products (...)

<http://jalammar.github.io/>

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

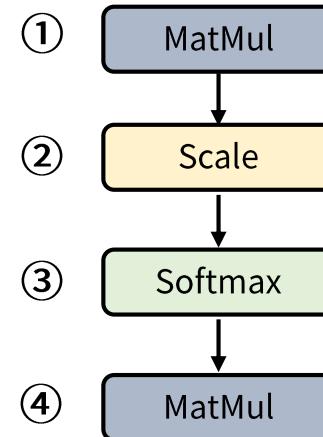
03 | Self-Attention Basics

❖ Transformer



Scaled Dot-Product
Attention

Generalized
Attention Form



Self-
Attention

$$A(\mathbf{q}, \mathbf{K}, \mathbf{V}) = \sum_i softmax(f(\mathbf{K}, \mathbf{q})) \mathbf{V}$$

$$f(\mathbf{K}, \mathbf{Q}) = \mathbf{Q}\mathbf{K}^T$$

$$\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}$$

$$softmax\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)$$

$$softmax\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

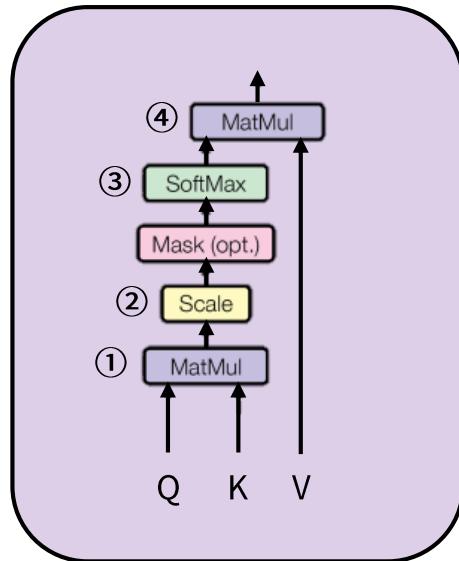
$$A(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = softmax\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}$$

<http://jalammar.github.io/>

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

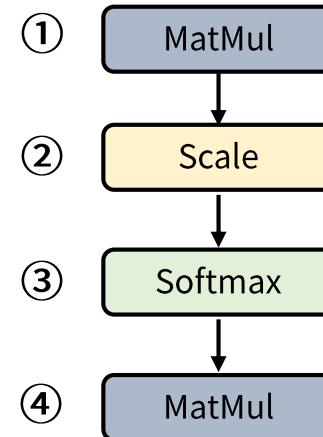
03 | Self-Attention Basics

❖ Transformer

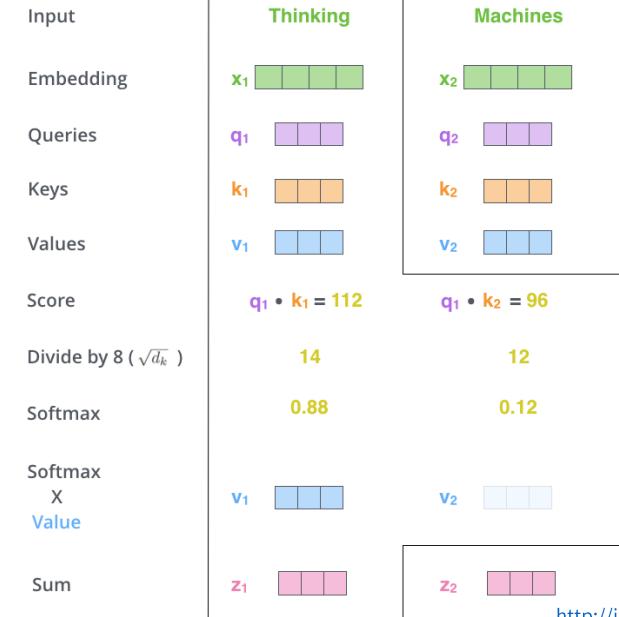


Scaled Dot-Product
Attention

Self-Attention



$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

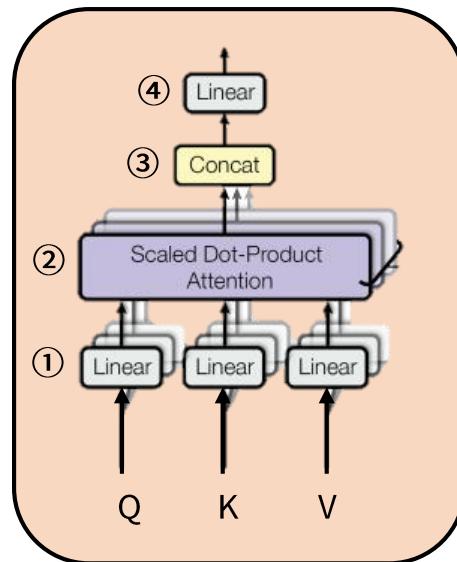


Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

<http://jalammar.github.io/>

03 | Self-Attention Basics

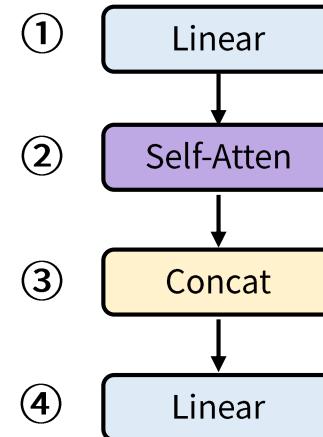
❖ Transformer



Multi-Head
Attention

Self-Attention

$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Multi-Head
Attention

$$\begin{aligned} Q' &= QW_i^Q & K' &= KW_i^K & V' &= VW_i^V \\ \text{head}_i &= A(Q', K', V') \end{aligned}$$

$$[\text{head}_1, \text{head}_2, \dots, \text{head}_h]$$

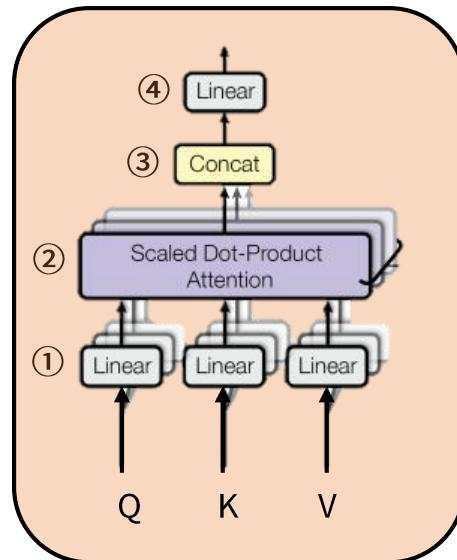
$$[\text{head}_1, \text{head}_2, \dots, \text{head}_h]W^O$$

$$\begin{aligned} \text{MultiHead}(Q, K, V) \\ = [\text{head}_1, \text{head}_2, \dots, \text{head}_h]W^O \end{aligned}$$

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017.. <http://jalamar.github.io/>

03 | Self-Attention Basics

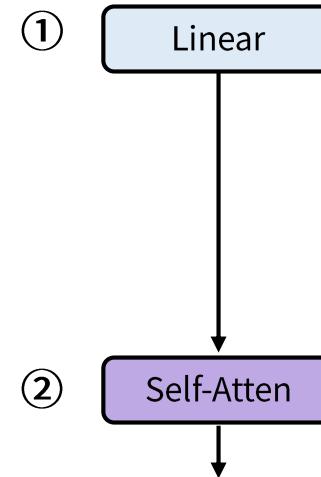
❖ Transformer



Multi-Head
Attention

Self-Attention

$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



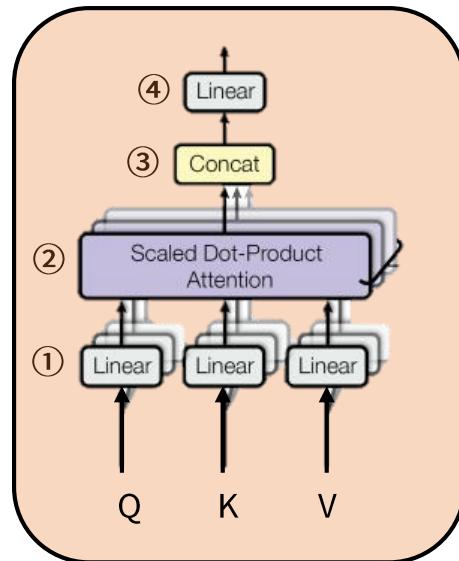
$$\begin{aligned} x &\times W^Q = Q \\ x &\times W^K = K \\ x &\times W^V = V \\ Q &\times K^T \\ \text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) &= Z \end{aligned}$$

Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

<http://jalammar.github.io/>

03 | Self-Attention Basics

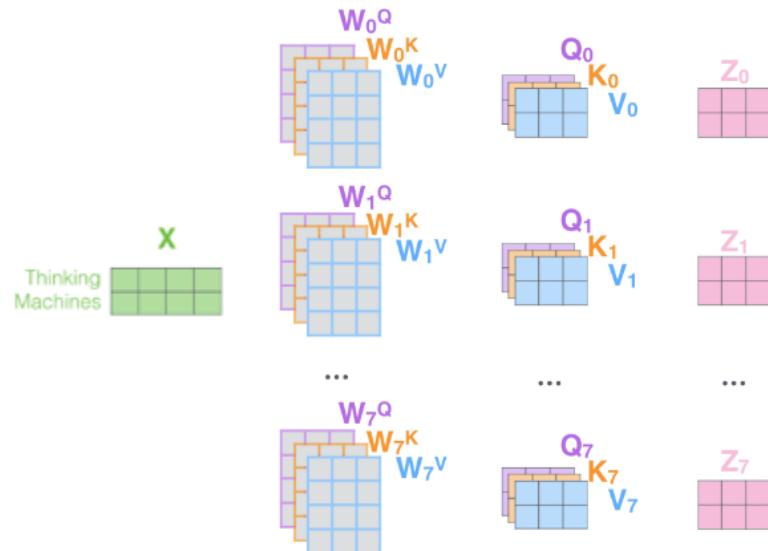
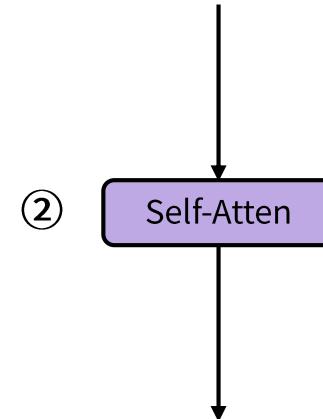
❖ Transformer



Multi-Head
Attention

Self-Attention

$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

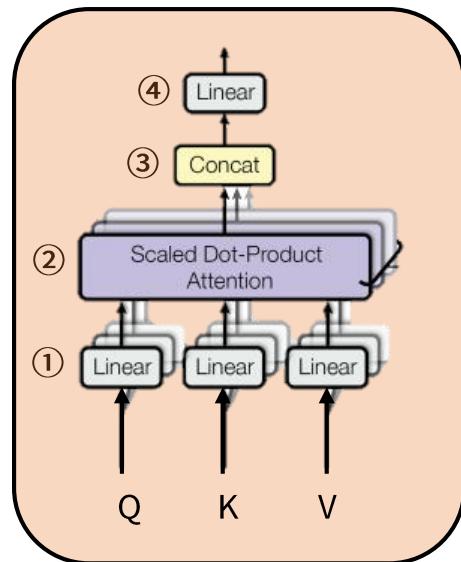


<http://jalammar.github.io/>

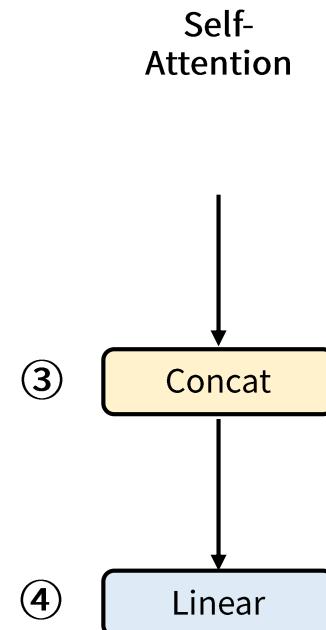
Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

03 | Self-Attention Basics

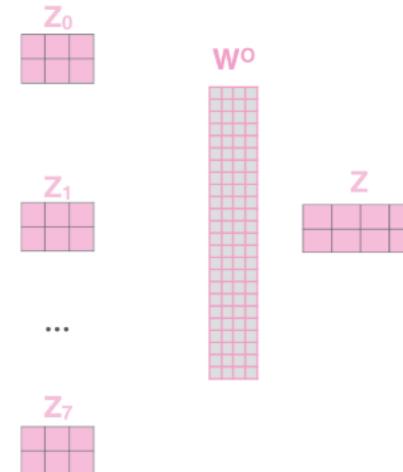
❖ Transformer



Multi-Head
Attention



$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

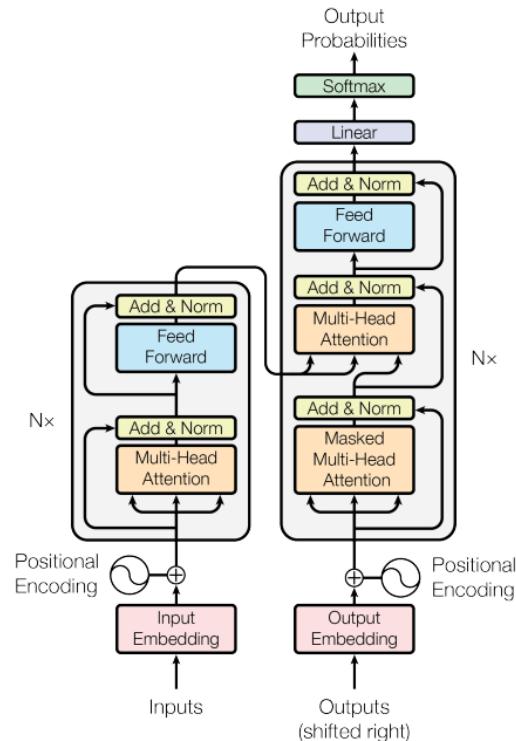


Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017.. <http://jalammar.github.io/>

03 | Self-Attention Basics

❖ Contributions

- One is the **total computational complexity** per layer. (⋯)
- Another is the amount of computation that can **be parallelized**, (⋯)
- The third is the path length between long-range dependencies in the network. **Learning long-range dependencies** is a key challenge in many sequence transduction tasks. (⋯)
- As side benefit, self-attention could yield **more interpretable models**.



Transformer <http://jalammar.github.io/>

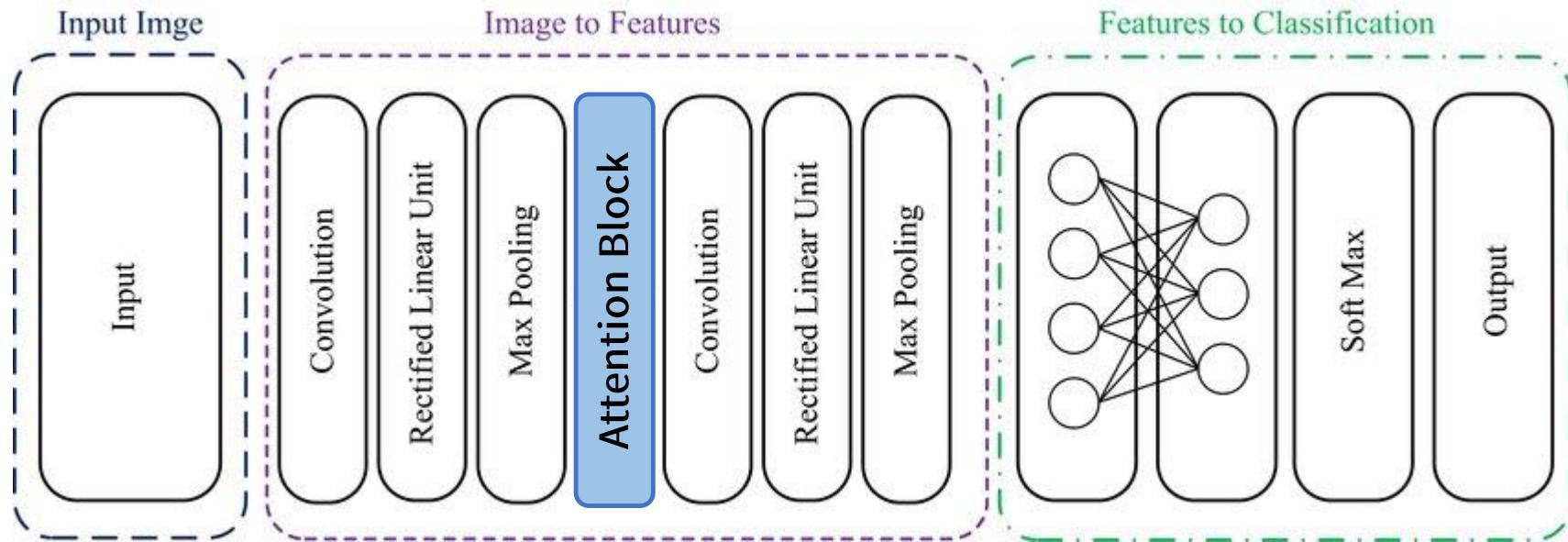
Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017..

03 | Self-Attention

Visual Self-Attention

❖ Attention: CNN Architecture

- 이미지의 feature을 보다 더 잘 생성하기 위해 convolution layer 사이에 attention block 삽입



https://www.researchgate.net/publication/333112125_An_islanding_detection_method_based_on_image_classification_with_convolution_neural_network

03 | Self-Attention

Visual Self-Attention



❖ Various Attention Algorithms: CNN

- 이미지의 feature을 보다 더 잘 생성하기 위해 convolution layer 사이에 attention block 삽입

Squeeze-and-excitation networks

J Hu, L Shen, G Sun... - ... of the IEEE conference on computer ..., 2018 - openaccess.thecvf.com

Convolutional neural networks are built upon the convolution operation, which extracts informative features by fusing spatial and channel-wise information together within local receptive fields. In order to boost the representational power of a network, several recent ...

☆ 99 2267회 인용 관련 학술자료 전체 16개의 버전 ☰

Cbam: Convolutional block attention module

S Woo, J Park, JY Lee... - Proceedings of the ..., 2018 - openaccess.thecvf.com

Abstract We propose Convolutional Block Attention Module (CBAM), a simple and effective attention module that can be integrated with any feed-forward convolutional neural networks. Given an intermediate feature map, our module sequentially infers attention maps ...

☆ 99 378회 인용 관련 학술자료 전체 9개의 버전 ☰

Non-local neural networks

X Wang, R Girshick, A Gupta... - Proceedings of the IEEE ..., 2018 - openaccess.thecvf.com

Both convolutional and recurrent operations are building blocks that process one local neighborhood at a time. In this paper, we present non-local operations as a generic family of building blocks for capturing long-range dependencies. Inspired by the classical non-local ...

☆ 99 871회 인용 관련 학술자료 전체 10개의 버전 ☰

Stand-alone self-attention in vision models

P Ramachandran, N Parmar, A Vaswani, I Bello... - arXiv preprint arXiv ..., 2019 - arxiv.org

Convolutions are a fundamental building block of modern computer vision systems. Recent approaches have argued for going beyond convolutions in order to capture long-range dependencies. These efforts focus on augmenting convolutional models with content-based ...

☆ 99 10회 인용 관련 학술자료 전체 5개의 버전 ☰

- Jie Hu et al., University of Oxford
- Computer Vision and Pattern Recognition(CVPR) in 2018
- Channel attention
- Winner of ILSVRC 2017
- S Woo, J Park et al., KAIST, LUNIT
- European Conference on Computer Vision(ECCV) in 2018
- Channel attention + Spatial attention
- X Wang, Abhinav Gupta, Carnegie Mellon University
- Ross Girshick, Kaiming He, Facebook AI Research
- Computer Vision and Pattern Recognition(CVPR) in 2018
- Non-local self-attention
- X Wang, Abhinav Gupta, Carnegie Mellon University
- P. Ramachandran, A. Vaswani et al., Google Research, Brain
- arXiv in 2019
- Local self-attention

03 | Self-Attention Visual Self-Attention

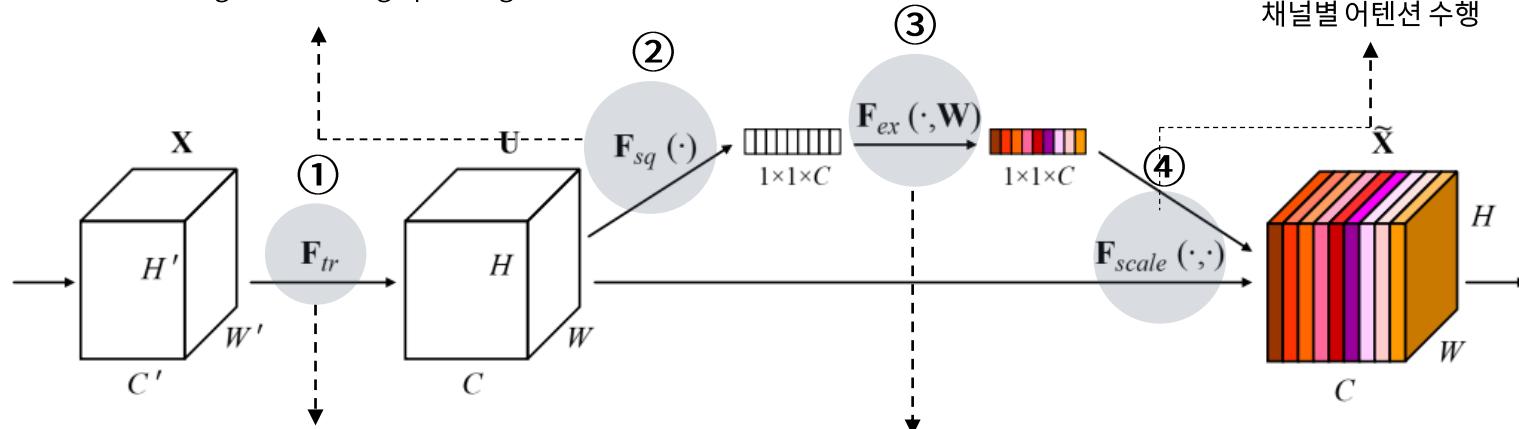
❖ Squeeze-and-Excitation Networks(SENNet) - 2018

$$F_{sq}(u_c) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W u_c(i, j) = z_c$$

② Channel-wise global average pooling을 통해 채널별 특성 추출

$$F_{scale}(u_c, s_c) = s_c u_c = \tilde{x}_c$$

④ 이전 layer에 채널 특성 벡터의 dot product를 통한
채널별 어텐션 수행



$$F_{tr}(X) = v_c * X = u_c$$

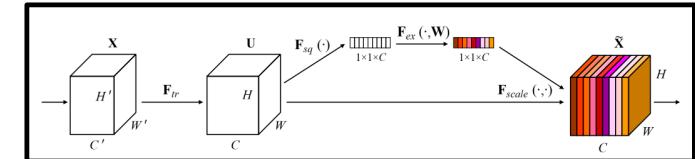
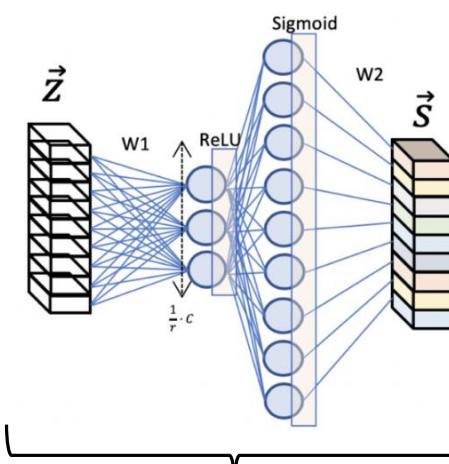
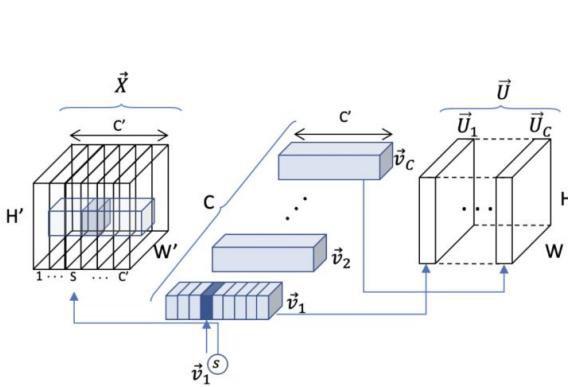
① Conv 연산을 통해 차원 변환

$$F_{ex}(z, W) = \text{sigmoid}(W_2 \text{Relu}(W_1 z)) = s$$

③ 채널별 특성의 관계를 학습하는 과정을 거쳐 채널 특성 벡터를 추출

03 | Self-Attention Visual Self-Attention

❖ Squeeze-and-Excitation Networks(SENNet) - 2018



① Transform

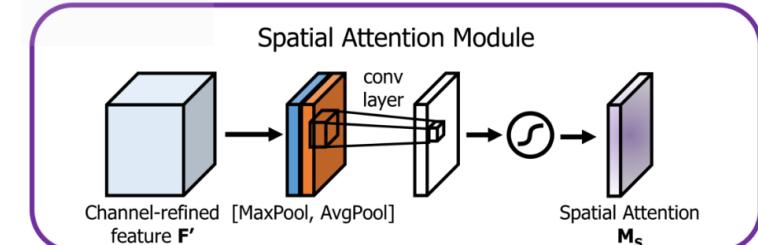
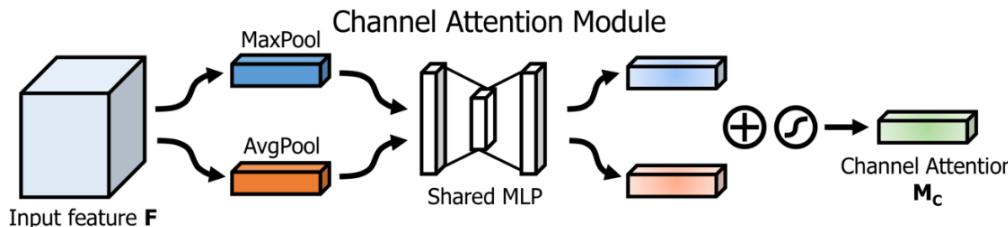
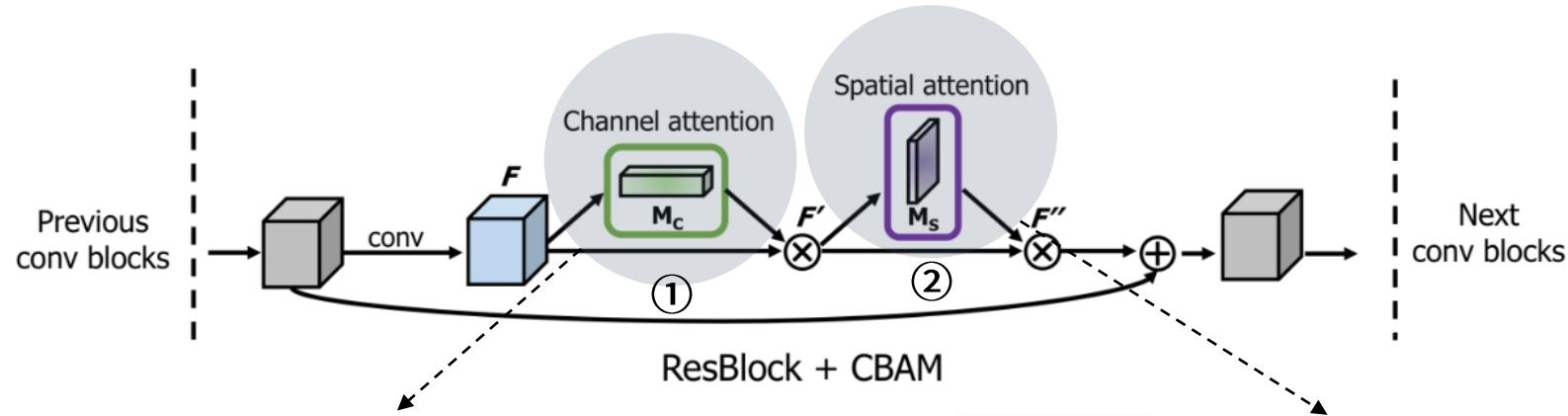
② Squeeze ③ Excitation

④ Scaling

<https://towardsdatascience.com/understanding-and-visualizing-se-nets-154aff0fc68>

03 | Self-Attention Visual Self-Attention

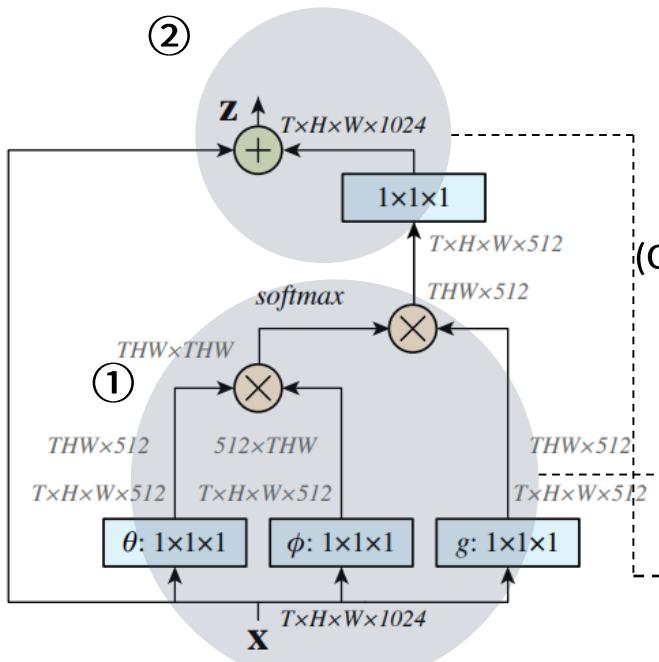
❖ Convolutional Block Attention Module(CBAM) - 2018



Woo, Sanghyun, et al. "Cbam: Convolutional block attention module." *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018.

03 | Self-Attention Visual Self-Attention

❖ Non-Local Neural Networks - 2018



Local Layer
(Convolution)

$$y_{ij} = \sum_{a,b \in N} W_{i-a,j-b}, x_{ab}$$

Non-local
layer

$$\textcircled{1} \quad y_i = \frac{1}{C(x)} \sum_{\forall j} softmax(f(\theta(x_i), \phi(x_j))) g(x_j)$$

$$\textcircled{2} \quad z_i = W_z y_i + x_i$$

- The third is the path length between long-range dependencies in the network. [Learning long-range dependencies](#) is a key challenge in many sequence transduction tasks. (...)

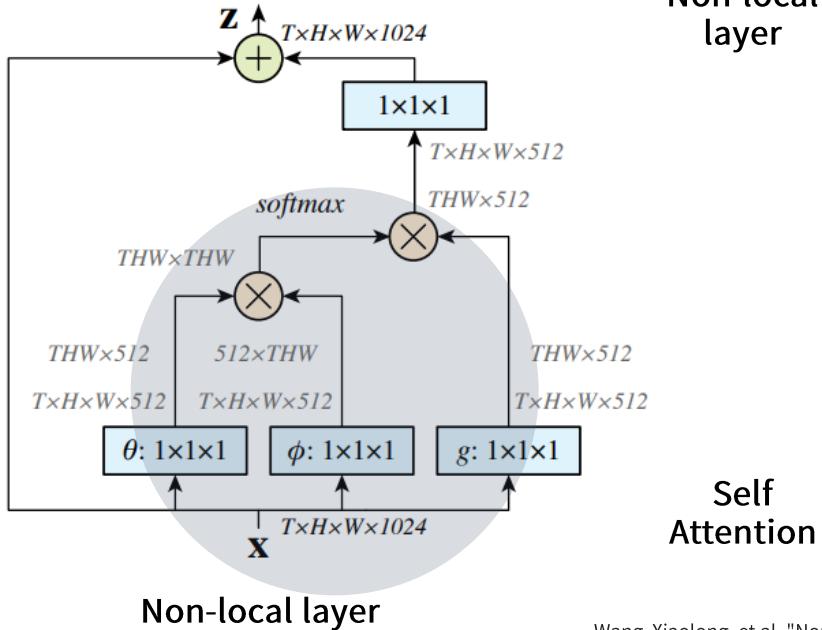
Wang, Xiaolong, et al. "Non-local neural networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.

03 | Self-Attention

Visual Self-Attention

❖ Non-Local Neural Networks - 2018

- The third is the path length between long-range dependencies in the network. [Learning long-range dependencies](#) is a key challenge in many sequence transduction tasks. (…)



Non-local layer

$$y_i = \frac{1}{C(x)} \sum_{\forall j} \text{softmax}(f(\theta(x_i), \phi(x_j)))g(x_j)$$

$C(x) = 1$ $\theta, \phi, g = \text{Linear function}$ $f = \text{dot product}$

$y_i, x_i, x_j \rightarrow Y, X, X$

$$Y = \text{softmax}(X^T W_\theta^T W_\phi X) W_g X$$

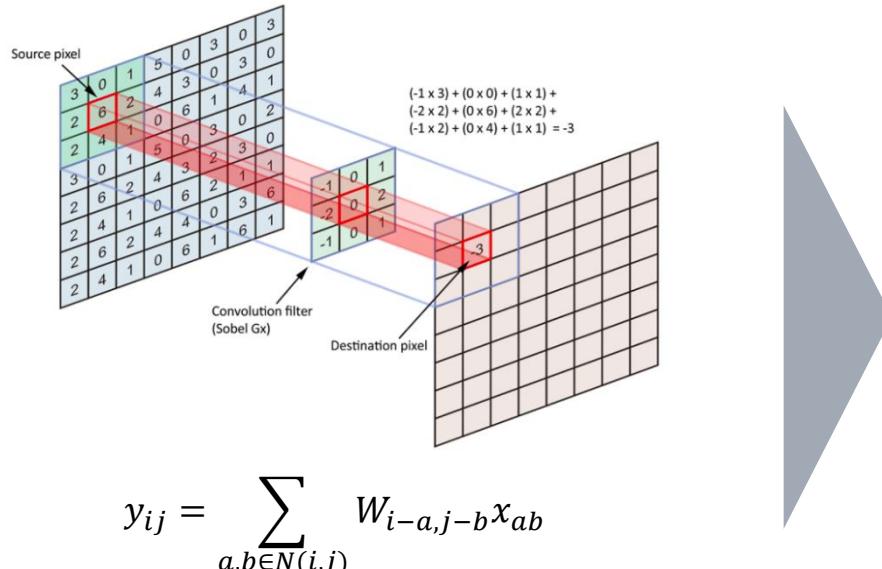
$$A(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Wang, Xiaolong, et al. "Non-local neural networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.

03 | Self-Attention

Visual Self-Attention

❖ Stand-Alone Self-Attention - 2019



Convolution

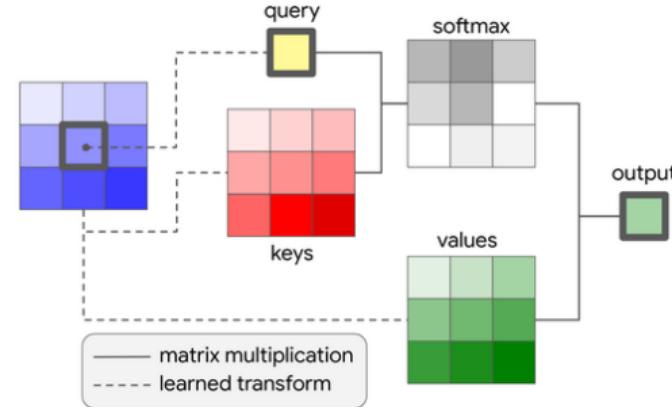
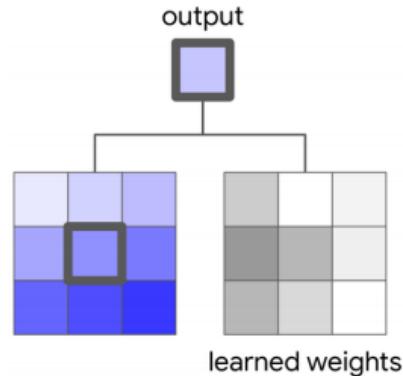
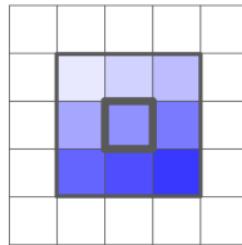
Self-Attention

Ramachandran, Prajit, et al. "Stand-alone self-attention in vision models." *arXiv preprint arXiv:1906.05909* (2019).

03 | Self-Attention

Visual Self-Attention

❖ Stand-Alone Self-Attention - 2019



$$y_{ij} = \sum_{a,b \in N(i,j)} W_{i-a,j-b} x_{ab}$$

Convolution

$$y_{ij} = \sum_{a,b \in N(i,j)} \text{softmax}_{ab}(q_{ij}^T k_{ab}) v_{ab}$$
$$q_{ij} = W_Q x_{ij} \quad k_{ab} = W_K x_{ab} \quad v_{ab} = W_V x_{ab}$$

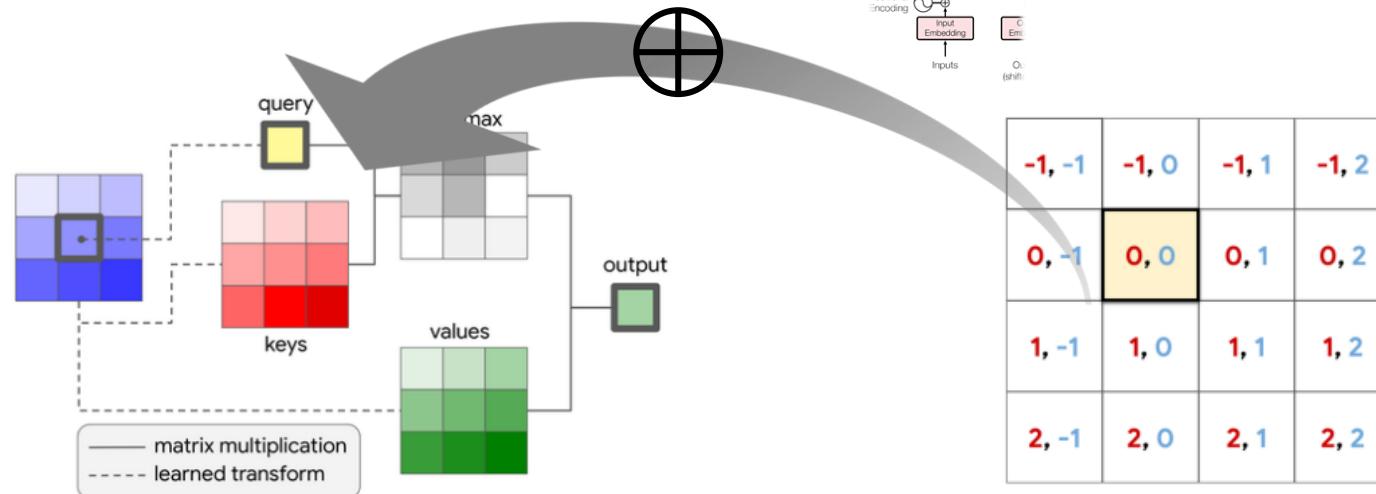
Self-Attention

Ramachandran, Prajit, et al. "Stand-alone self-attention in vision models." *arXiv preprint arXiv:1906.05909* (2019).

03 | Self-Attention

Visual Self-Attention

❖ Stand-Alone Self-Attention - 2019



$$y_{ij} = \sum_{a,b \in N(i,j)} softmax_{ab}(q_{ij}^T k_{ab}) v_{ab}$$

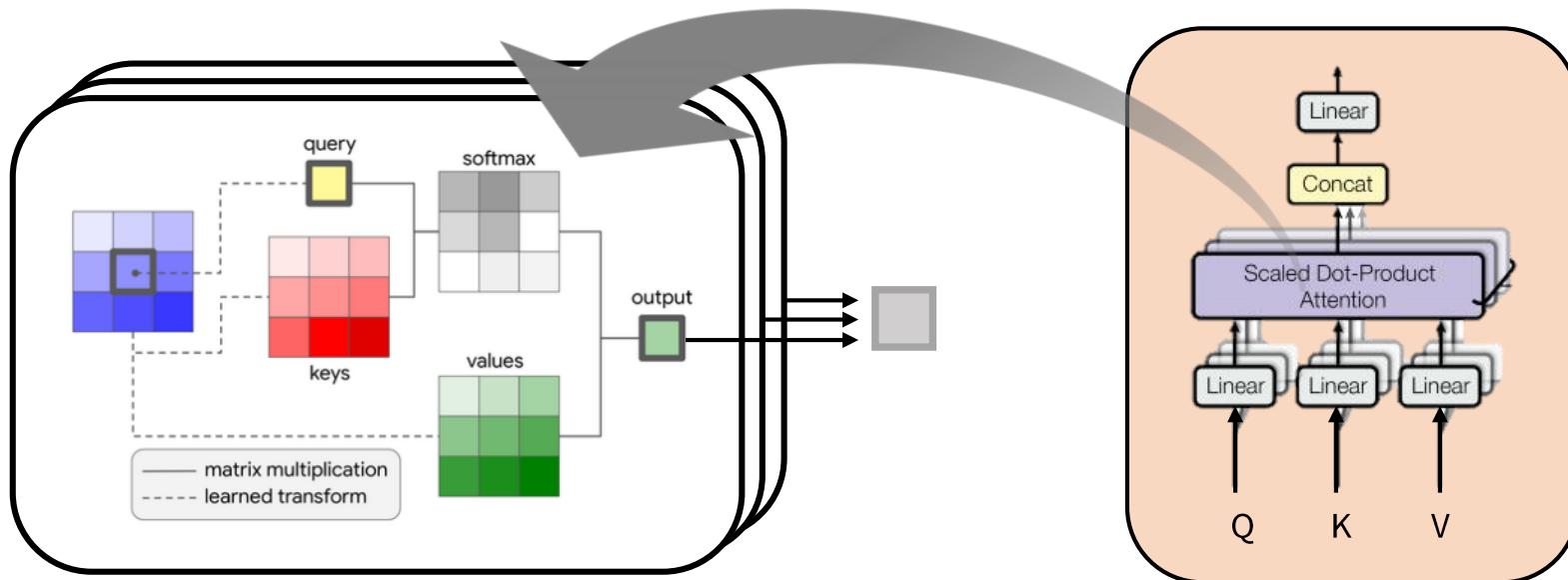
$$y_{ij} = \sum_{a,b \in N(i,j)} softmax_{ab}(q_{ij}^T k_{ab} + q_{ij}^T r_{a-i,b-j}) v_{ab}$$

Ramachandran, Prajit, et al. "Stand-alone self-attention in vision models." *arXiv preprint arXiv:1906.05909* (2019).

03 | Self-Attention

Visual Self-Attention

❖ Stand-Alone Self-Attention - 2019



Multi-Head
Attention

Ramachandran, Prajit, et al. "Stand-alone self-attention in vision models." *arXiv preprint arXiv:1906.05909* (2019).

03 | Self-Attention

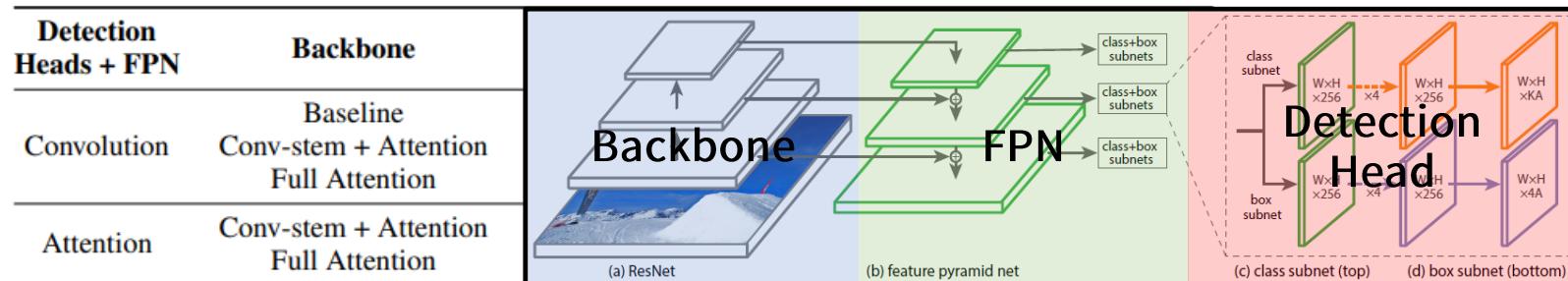
Visual Self-Attention

❖ Stand-Alone Self-Attention - 2019

- Image classification – ImageNet dataset

	ResNet-26			ResNet-38			ResNet-50		
	FLOPS (B)	Params (M)	Acc. (%)	FLOPS (B)	Params (M)	Acc. (%)	FLOPS (B)	Params (M)	Acc. (%)
Baseline	4.7	13.7	74.5	6.5	19.6	76.2	8.2	25.6	76.9
Conv-stem + Attention	4.5	10.3	75.8	5.7	14.1	77.1	7.0	18.0	77.4
Full Attention	4.7	10.3	74.8	6.0	14.1	76.9	7.2	18.0	77.6

- Object detection – COCO dataset



Ramachandran, Prajit, et al. "Stand-alone self-attention in vision models." *arXiv preprint arXiv:1906.05909* (2019).

❖ Conclusion

- Attention은 ‘모델이 더 나은 성능을 위해 집중적으로 학습해야 할 부분까지도 학습하도록 하자.’라는 관점에서 딥러닝 모델링에 필수적인 방법론
- 해석력 + 모델 성능 향상의 두 가지 이점
- Visual attention은 computer vision 분야 중 가장 활발하게 연구되고 있는 분야 중 하나
- 목적은 동일하지만 다양한 형태의 attention 방법론이 존재함
- 최근 CNN 모델의 block 중 하나로 활용되는 것이 아니라 convolution 연산까지 대체할 수 있다는 가능성 확인
- 텍스트, 이미지 뿐만 아니라 시그널 등 다양한 데이터의 모델링에도 효과적일 것으로 생각되어 연구실에서 진행하고 있는 다양한 영역의 연구 또는 프로젝트에서 활용하였으면 좋겠음.

❖ Reference

▪ Attention

- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio. "Neural machine translation by jointly learning to align and translate." *arXiv preprint arXiv:1409.0473* (2014).
- Luong, Minh-Thang, Hieu Pham, and Christopher D. Manning. "Effective approaches to attention-based neural machine translation." *arXiv preprint arXiv:1508.04025* (2015).
- Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention." *International conference on machine learning*. 2015.

▪ Self-Attention

- Vaswani, Ashish, et al. "Attention is all you need." *Advances in neural information processing systems*. 2017.
- Hu, Jie, Li Shen, and Gang Sun. "Squeeze-and-excitation networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
- Woo, Sanghyun, et al. "Cbam: Convolutional block attention module." *Proceedings of the European Conference on Computer Vision (ECCV)*. 2018.
- Wang, Xiaolong, et al. "Non-local neural networks." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.
- Ramachandran, Prajit, et al. "Stand-alone self-attention in vision models." *arXiv preprint arXiv:1906.05909* (2019).

감사합니다