

Effective approaches to attention-based neural machine translation



논문리뷰
이도연

목차

1 Introduction

2 NMT (Neural Machine Translation)

3 Attention-based Models

- **Global Attention**

- **Local Attention**

- **Input-feeding Approach**

4 Experiments & Analysis

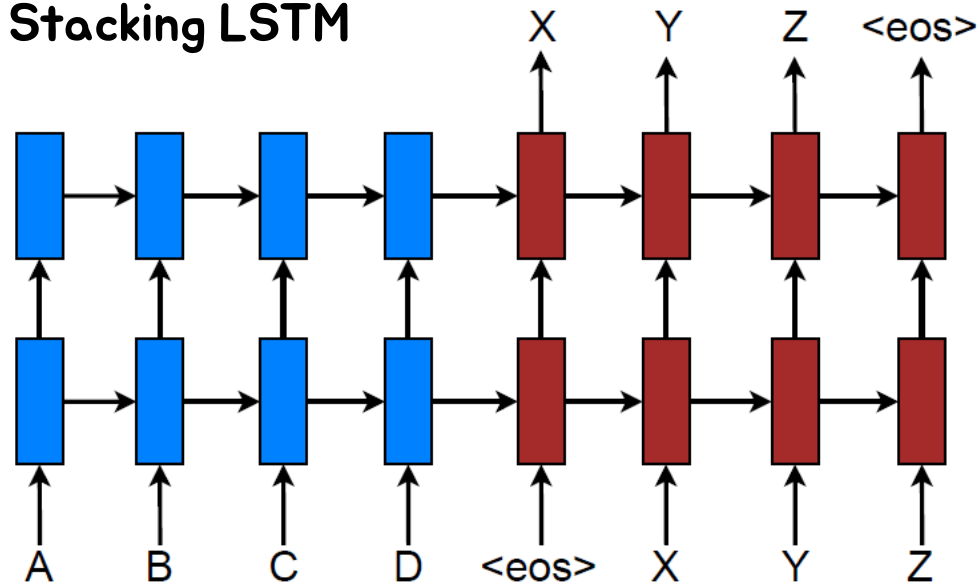
5 Conclusion

1. Introduction

- Attentional mechanism이 NMT의 성능 향상을 위해 사용됨
- NMT에 Attention을 더 효과적으로 사용하기 위한 방법 제시
 - > Global Attention / Local Attention
- SOTA for both WMT'14 and WMT'15

2. NMT (Neural Machine Translation)

Stacking LSTM



Source sentence x 에서 Target sentence y 가 생성될 확률

$$\log p(y|x) = \sum_{j=1}^m \log p(y_j | y_{<j}, s)$$

각 단어 y 를 decoding 할 확률

$$p(y_j | y_{<j}, s) = \text{softmax}(g(h_j))$$

이전 hidden state로부터 현재의 hidden state를 계산

$$h_j = f(h_{j-1}, s)$$

f 는 RNN, LSTM, GRU

Source sentence x_1, \dots, x_n Target sentence y_1, \dots, y_m ↓

Encoder - 각 source 문장을 대표하는 s 계산

Decoder - 한 번에 하나의 target word를 생성

$$J_t = \sum_{(x,y) \in \mathbb{D}} -\log p(y|x)$$

3. Attention-based Models

Common to these two types of models is the fact that at each time step t in the decoding phase, both approaches first take as input the hidden state h_t at the top layer of a stacking LSTM. The goal is then to derive a context vector c_t that captures relevant source-side information to help predict the current target word y_t . While these models differ in how the context vector c_t is derived, they share the same subsequent steps.

Specifically, given the target hidden state h_t and the source-side context vector c_t , we employ a simple concatenation layer to combine the information from both vectors to produce an attentional hidden state as follows:

$$\tilde{h}_t = \tanh(W_c[c_t; h_t]) \quad (5)$$

The attentional vector \tilde{h}_t is then fed through the softmax layer to produce the predictive distribution formulated as:

$$p(y_t|y_{<t}, x) = \text{softmax}(W_s \tilde{h}_t) \quad (6)$$

We now detail how each model type computes the source-side context vector c_t .

- Stacked LSTM layer의 맨 위층의 hidden state를 사용

- 현재 target hidden state h_t

- Target Word를 예측하기 위한 Source의 정보를 담고 있는 context vector c_t

- attentional hidden state vector $\tilde{h}_t = \tanh(W_c[c_t; h_t])$
 h_t, c_t concat

$$p(y_t|y_{<t}, x) = \text{softmax}(W_s \tilde{h}_t)$$

c_t 를 구하는 방법에 따라 Global/Local Attention

3. Attention-based Models

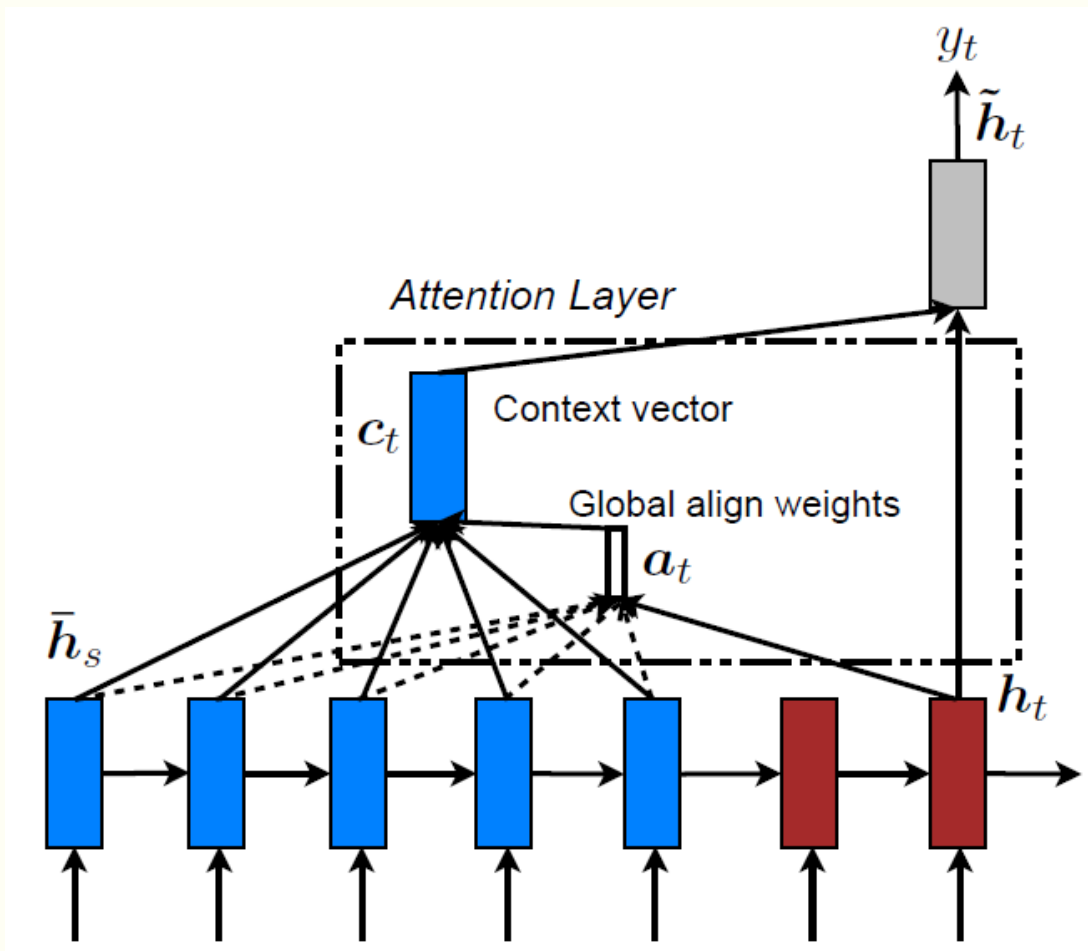
Global Attention

Encoder의 모든 hidden state를 고려

$$\begin{aligned} a_t(s) &= \text{align}(h_t, \bar{h}_s) && \text{Score Softmax} \\ &= \frac{\exp(\text{score}(h_t, \bar{h}_s))}{\sum_{s'} \exp(\text{score}(h_t, \bar{h}_{s'}))} \end{aligned}$$

$$\text{score}(h_t, \bar{h}_s) = \begin{cases} h_t^\top \bar{h}_s & \text{dot} \\ h_t^\top W_a \bar{h}_s & \text{general} \\ v_a^\top \tanh(W_a[h_t; \bar{h}_s]) & \text{concat} \end{cases}$$

$c_t \Leftarrow a_t$ 랑 \bar{h}_s weighted average



3. Attention-based Models

cf) 다양한 score function

content-based

$$\begin{aligned} a_t(s) &= \text{align}(h_t, \bar{h}_s) \\ &= \frac{\exp(\text{score}(h_t, \bar{h}_s))}{\sum_{s'} \exp(\text{score}(h_t, \bar{h}_{s'}))} \end{aligned}$$

$$\text{score}(h_t, \bar{h}_s) = \begin{cases} h_t^\top \bar{h}_s & \text{dot} \\ h_t^\top W_a \bar{h}_s & \text{general} \\ v_a^\top \tanh(W_a [h_t; \bar{h}_s]) & \text{concat} \end{cases}$$

location-based

$$a_t = \text{softmax}(W_a h_t) \quad \text{location}$$

Encoder의 hidden state가 사용되지 X

각각 이름과 식을 통해 생각해보면,
Encoder의 각 hidden state와 Decoder의 hidden state
각 content들을 비교한 content-based와
Decoder의 hidden state의 위치만을 고려한 location-based

3. Attention-based Models

cf) **Luong** Attention(Global Attention) vs **Bahdanau** Attention

Effective Approaches to Attention-based Neural Machine Translation

Minh-Thang **Luong** Hieu Pham Christopher D. Manning
Computer Science Department, Stanford University, Stanford, CA 94305
{lmthang, hyhieu, manning}@stanford.edu

NEURAL MACHINE TRANSLATION
BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

Dzmitry **Bahdanau**
Jacobs University Bremen, Germany

3. Attention-based Models

cf) **Luong** Attention(Global Attention) vs **Bahdanau** Attention

Hidden State

- Encoder, Decoder 모두 stacked LSTM layer의 맨 위층의 hidden state를 사용
- bi-directional Encoder의 hidden state와 non-stacked uni-directional Decoder의 hidden state를 사용

Computation path

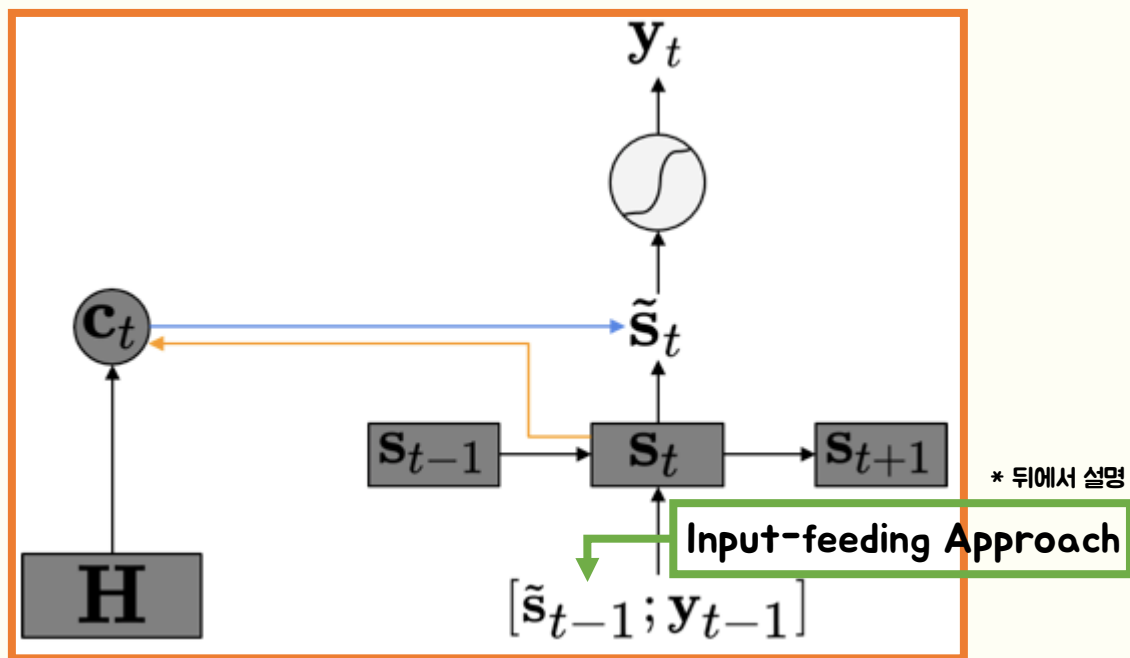
$$h_t \rightarrow a_t \rightarrow c_t \rightarrow \tilde{h}_t$$

$$h_{t-1} \rightarrow a_t \rightarrow c_t \rightarrow h_t$$

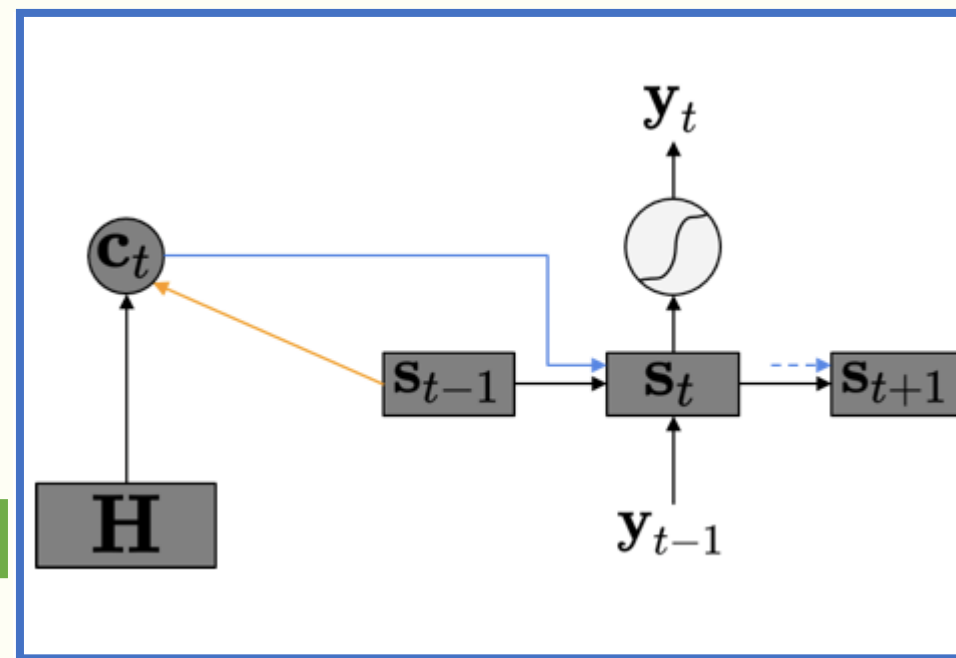
3. Attention-based Models

cf) **Luong** Attention(Global Attention) vs **Bahdanau** Attention

$$h_t \rightarrow a_t \rightarrow c_t \rightarrow \tilde{h}_t$$



$$h_{t-1} \rightarrow a_t \rightarrow c_t \rightarrow h_t$$



Bahdanau에서는 RNN의 재귀 연산이 수행되는 도중에 c_t 가 구해질 때까지 기다려야 하지만

Luong의 경우에는 y_t 를 구하는 부분과 RNN의 재귀 연산이 수행되는 부분을 분리할 수 있어 Computation path가 간소화!

3. Attention-based Models

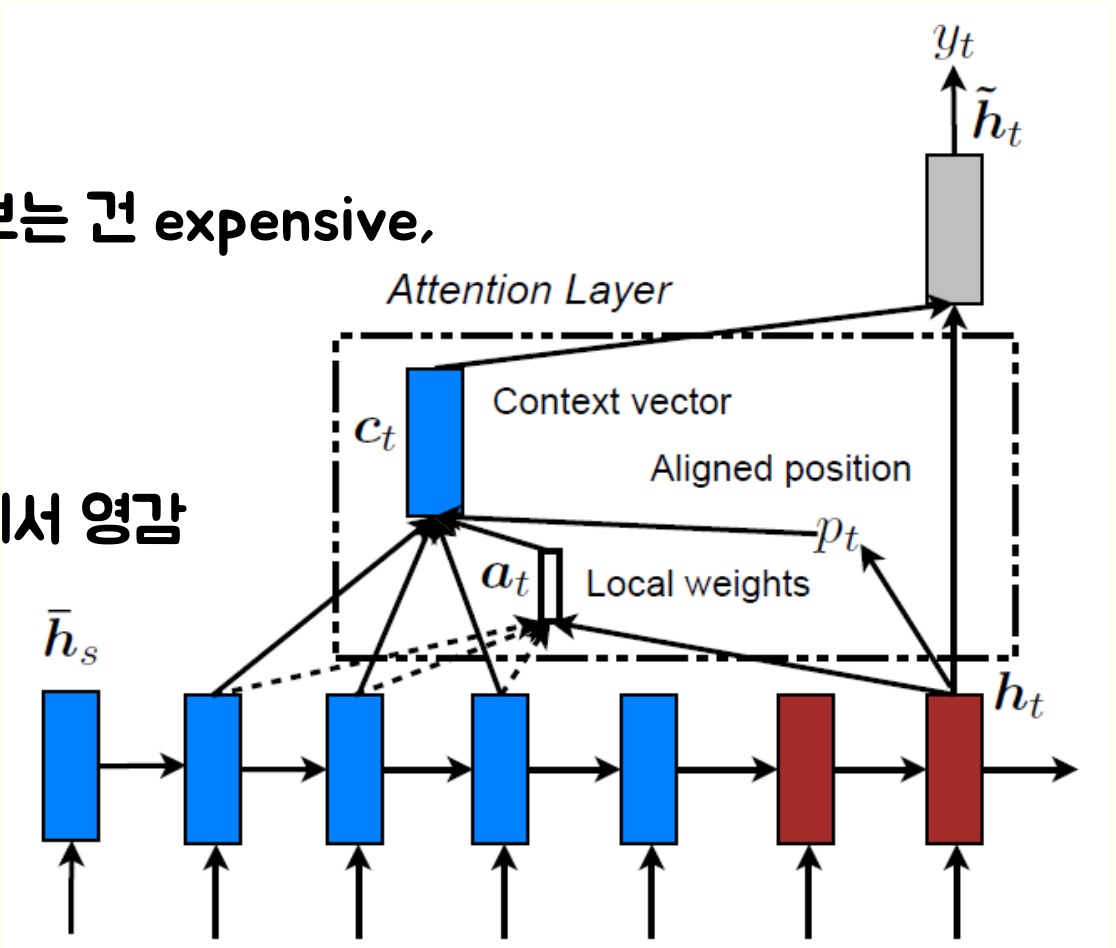
Local Attention

Global Attention에서 source의 모든 단어를 보는 건 expensive,
긴 sequence에서 실용적이지 못하다

-> source position의 small subset에 집중

Soft and Hard Attention 사이의 tradeoff에서 영감

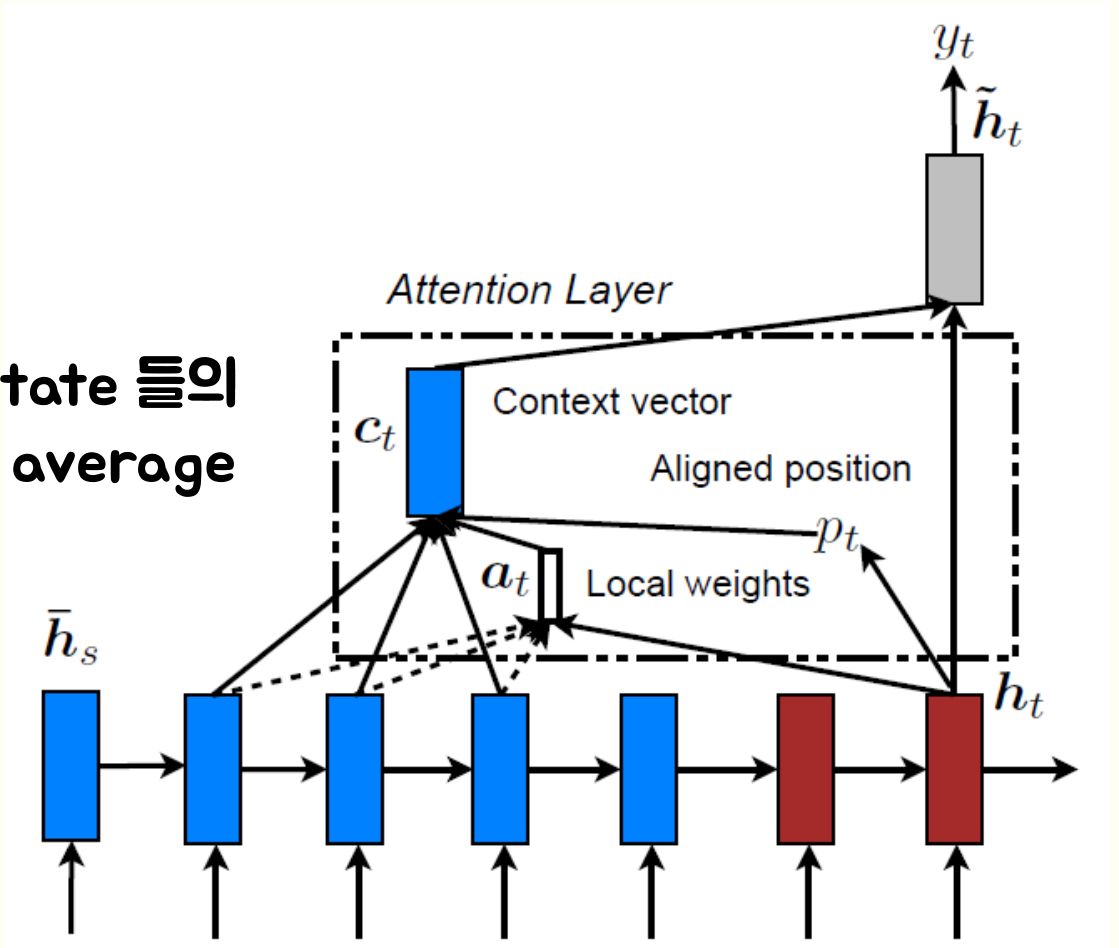
<https://arxiv.org/abs/1502.03044>



3. Attention-based Models

Local Attention

- Aligned Position p_t
- Context vector
 $[p_t - D, p_t + D]$ 사이의 source hidden state 들의 weighted average
- Local alignment vector a_t 는
fixed-dimensional 2D+1
- +) global 에서 a_t 는 source의
timestamp 개수와 동일한 크기



3. Attention-based Models

Local Attention

$[p_t - D, p_t + D]$ Alignment

1. Monotonic Alignment (local-m)

$p_t = t$ +시점 같은 위치에 있는 단어끼리 연관이 클 것이라는 생각

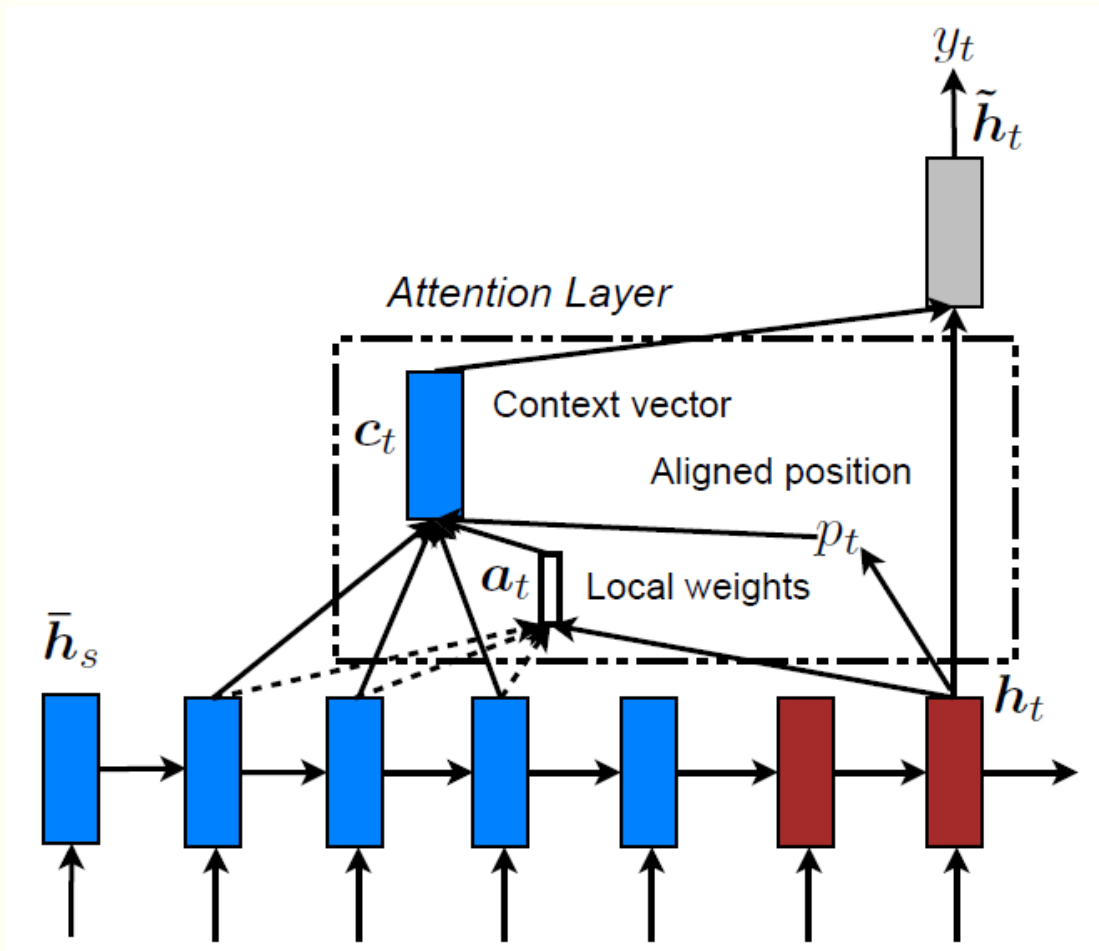
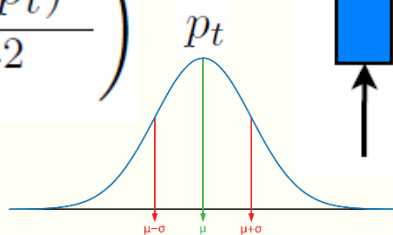
2. Predictive Alignment (local-p)

$p_t \in [0, S]$
 $p_t = \underline{S} \cdot \text{sigmoid}(\mathbf{v}_p^\top \tanh(\mathbf{W}_p \mathbf{h}_t))$
문장의 길이

$$\sigma = \frac{D}{2}$$

$$a_t(s) = \text{align}(\mathbf{h}_t, \bar{\mathbf{h}}_s) \exp\left(-\frac{(s - p_t)^2}{2\sigma^2}\right)$$

p_t 를 기준으로 Gaussian 분포에 따라 주변 단어들이 의미를 가질 것이라는 생각



3. Attention-based Models

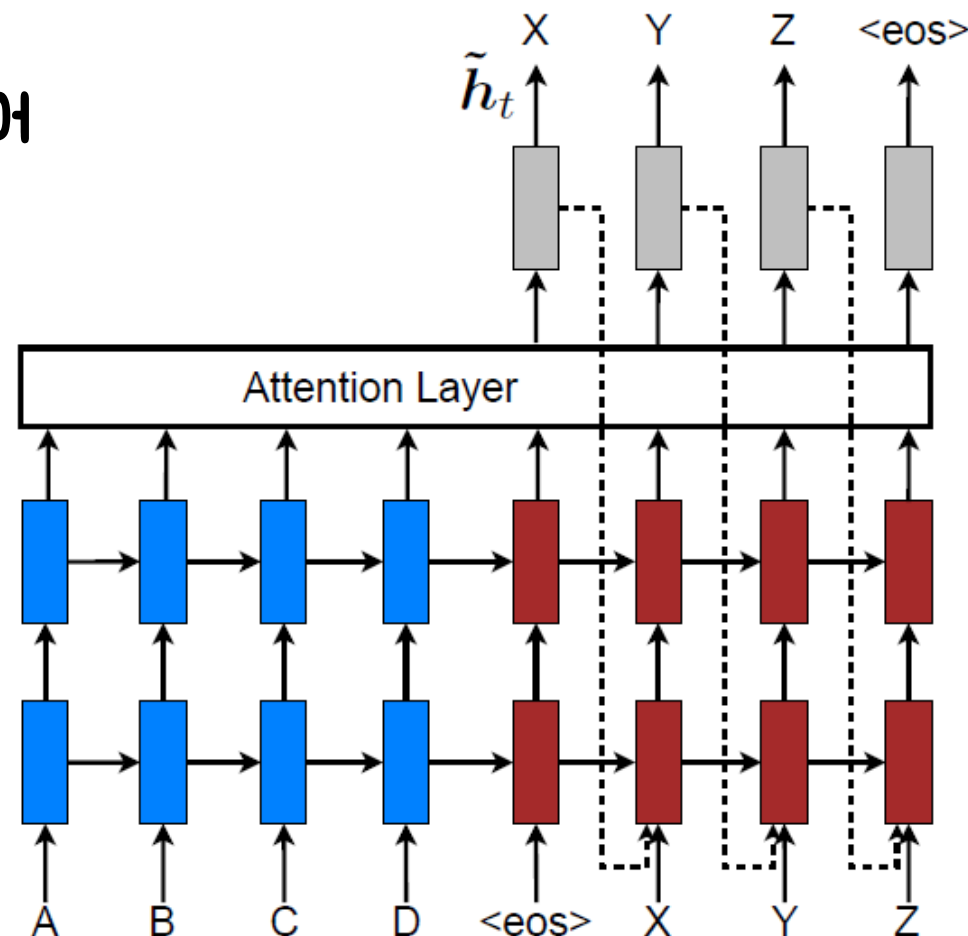
Input-feeding Approach

Standard MT에서는 coverage set이 계속 유지되어 어떤 source word가 번역되었는지 계속 따라간다

Attentional NMT에서도 이전 alignment 정보가 지금 alignment decision에 고려되도록 한다!

Attentional vector \tilde{h}_t 를 다음 시점의 input과 concat

-> 이전 Alignment 정보를 잘 알 수 있도록 하고, 수평/수직적으로 매우 깊은 네트워크를 만든다



4. Experiments & Analysis

System	Ppl	BLEU
Winning WMT'14 system – <i>phrase-based</i> + <i>large LM</i> (Buck et al., 2014)		20.7
<i>Existing NMT systems</i>		
RNNsearch (Jean et al., 2015)		16.5
RNNsearch + unk replace (Jean et al., 2015)		19.0
RNNsearch + unk replace + large vocab + <i>ensemble</i> 8 models (Jean et al., 2015)		21.6
<i>Our NMT systems</i>		
Base	10.6	11.3
Base + reverse	9.9	12.6 (+1.3)
Base + reverse + dropout	8.1	14.0 (+1.4)
Base + reverse + dropout + global attention (<i>location</i>)	7.3	16.8 (+2.8)
Base + reverse + dropout + global attention (<i>location</i>) + feed input	6.4	18.1 (+1.3)
Base + reverse + dropout + local-p attention (<i>general</i>) + feed input	5.9	19.0 (+0.9)
Base + reverse + dropout + local-p attention (<i>general</i>) + feed input + unk replace		20.9 (+1.9)
<i>Ensemble</i> 8 models + unk replace		23.0 (+2.1)

Table 1: **WMT'14 English-German results** – shown are the perplexities (ppl) and the *tokenized* BLEU scores of various systems on newstest2014. We highlight the **best** system in bold and give *progressive* improvements in italic between consecutive systems. *local-p* refers to the local attention with predictive alignments. We indicate for each attention model the alignment score function used in parentheses.

SOTA for both WMT'14 and WMT'15

System	BLEU
Top – <i>NMT</i> + <i>5-gram rerank</i> (Montreal)	24.9
Our ensemble 8 models + unk replace	25.9

Table 2: **WMT'15 English-German results**
NIST BLEU scores of the winning entry in WMT'15 and our best one on newstest2015.

System	Ppl.	BLEU
<i>WMT'15 systems</i>		
SOTA – <i>phrase-based</i> (Edinburgh)		29.2
NMT + 5-gram rerank (MILA)		27.6
<i>Our NMT systems</i>		
Base (reverse)	14.3	16.9
+ global (<i>location</i>)	12.7	19.1 (+2.2)
+ global (<i>location</i>) + feed	10.9	20.1 (+1.0)
+ global (<i>dot</i>) + drop + feed		22.8 (+2.7)
+ global (<i>dot</i>) + drop + feed + unk	9.7	24.9 (+2.1)

Table 3: **WMT'15 German-English results** – performances of various systems (similar to Table 1). The *base* system already includes source reversing on which we add *global* attention, *dropout*, input *feeding*, and *unk* replacement.

4. Experiments & Analysis

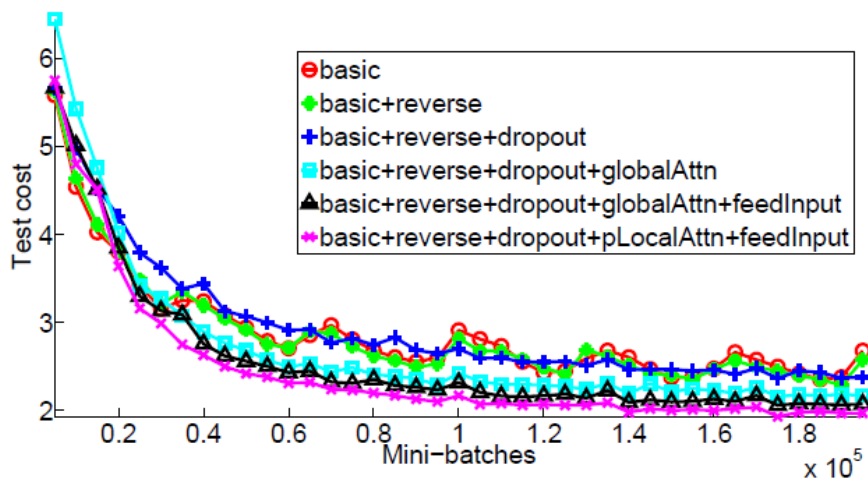


Figure 5: **Learning curves** – test cost (ln perplexity) on newstest2014 for English-German NMTs as training progresses.

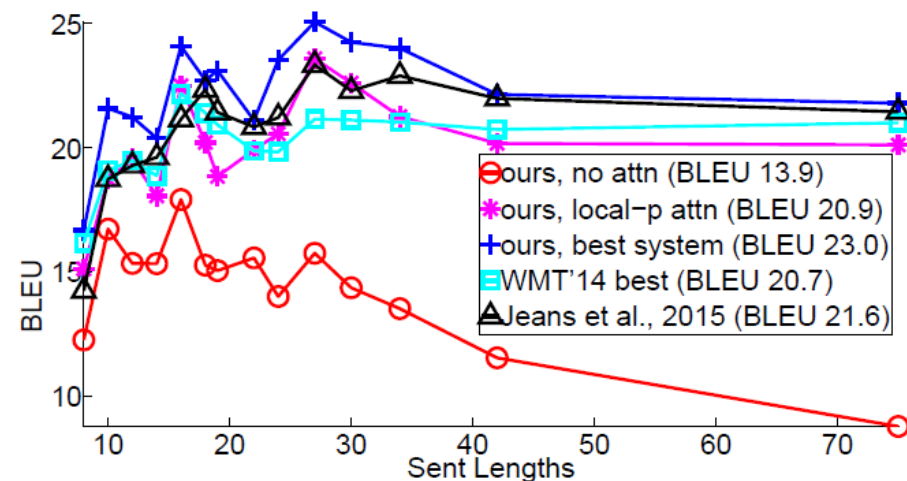


Figure 6: **Length Analysis** – translation qualities of different systems as sentences become longer.

Attention 적용했을 때 학습의 수렴 속도가 빠르다
문장의 길이가 길어질수록 번역 퀄리티 차이 보인다

4. Experiments & Analysis

System	Ppl	BLEU	
		Before	After unk
global (location)	6.4	18.1	19.3 (+1.2)
global (dot)	6.1	18.6	20.5 (+1.9)
global (general)	6.1	17.3	19.1 (+1.8)
local-m (dot)	>7.0	x	x
local-m (general)	6.2	18.6	20.4 (+1.8)
local-p (dot)	6.6	18.0	19.6 (+1.9)
local-p (general)	5.9	19	20.9 (+1.9)

Table 4: **Attentional Architectures** – performances of different attentional models. We trained two local-m (dot) models; both have ppl > 7.0.

Method	AER
global (location)	0.39
local-m (general)	0.34
local-p (general)	0.36
ensemble	0.34
Berkeley Aligner	0.32

Table 6: **AER scores** – results of various models on the RWTH English-German alignment data.

Alignment Error Rate

Attention 구조에 따른 번역 성능과 Alignment Quality

4. Experiments & Analysis

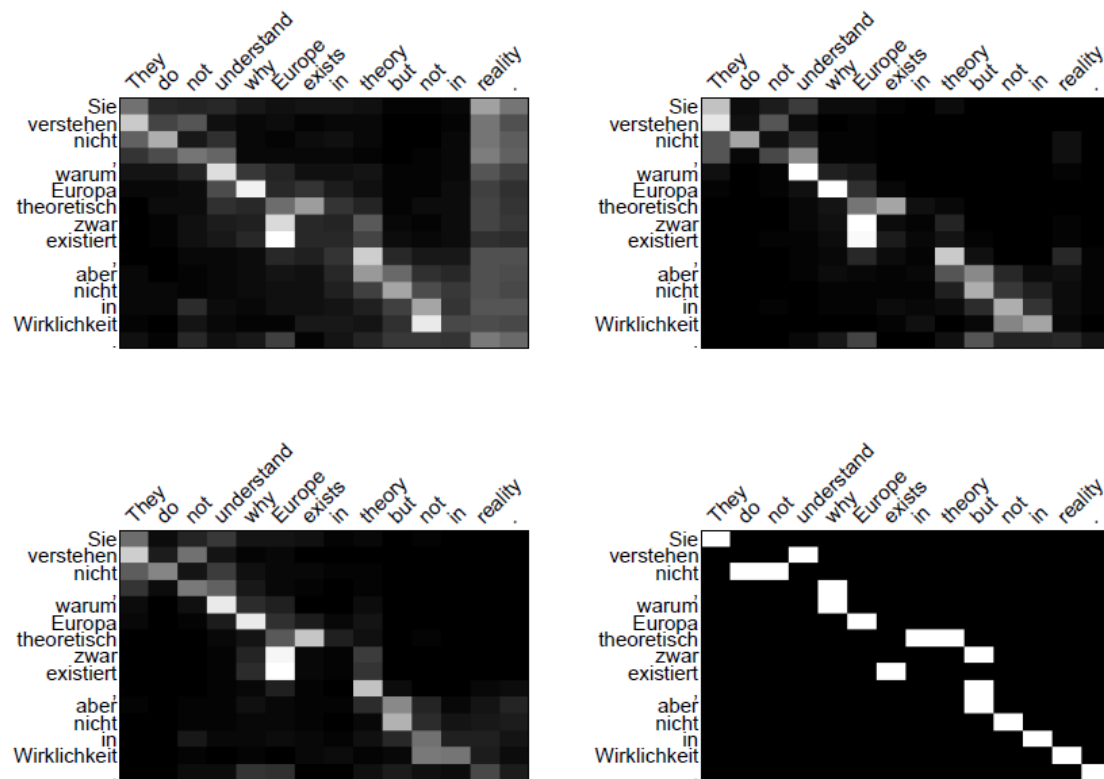


Figure 7: **Alignment visualizations** – shown are images of the attention weights learned by various models: (top left) global, (top right) local-m, and (bottom left) local-p. The *gold* alignments are displayed at the bottom right corner.

4. Experiments & Analysis

English-German translations

src	Orlando Bloom and Miranda Kerr still love each other
ref	Orlando Bloom und <i>Miranda Kerr</i> lieben sich noch immer
best	Orlando Bloom und <i>Miranda Kerr</i> lieben einander noch immer .
base	Orlando Bloom und Lucas Miranda lieben einander noch immer .
src	" We ' re pleased the FAA recognizes that an enjoyable passenger experience is not incompatible with safety and security , " said Roger Dow , CEO of the U.S. Travel Association .
ref	" Wir freuen uns , dass die FAA erkennt , dass ein angenehmes Passagiererlebnis nicht im Widerspruch zur Sicherheit steht " , sagte <i>Roger Dow</i> , CEO der U.S. Travel Association .
best	" Wir freuen uns , dass die FAA anerkennt , dass ein angenehmes ist nicht mit Sicherheit und Sicherheit <i>unvereinbar</i> ist " , sagte <i>Roger Dow</i> , CEO der US - die .
base	" Wir freuen uns über die <unk> , dass ein <unk> <unk> mit Sicherheit nicht vereinbar ist mit Sicherheit und Sicherheit " , sagte <i>Roger Cameron</i> , CEO der US - <unk> .

German-English translations

src	In einem Interview sagte Bloom jedoch , dass er und Kerr sich noch immer lieben .
ref	However , in an interview , Bloom has said that he and <i>Kerr</i> still love each other .
best	In an interview , however , Bloom said that he and <i>Kerr</i> still love .
base	However , in an interview , Bloom said that he and Tina were still <unk> .
src	Wegen der von Berlin und der Europäischen Zentralbank verhängten strengen Sparpolitik in Verbindung mit der Zwangsjacke , in die die jeweilige nationale Wirtschaft durch das Festhalten an der gemeinsamen Währung genötigt wird , sind viele Menschen der Ansicht , das Projekt Europa sei zu weit gegangen
ref	The <i>austerity imposed by Berlin and the European Central Bank</i> , coupled with the straitjacket imposed on national economies through adherence to the common currency , has led many people to think Project Europe has gone too far .
best	Because of the strict <i>austerity measures imposed by Berlin and the European Central Bank in connection with the straitjacket</i> in which the respective national economy is forced to adhere to the common currency , many people believe that the European project has gone too far .
base	Because of the pressure imposed by the European Central Bank and the Federal Central Bank with the strict austerity imposed on the national economy in the face of the single currency , many people believe that the European project has gone too far .

Table 5: **Sample translations** – for each example, we show the source (*src*), the human translation (*ref*), the translation from our best model (*best*), and the translation of a non-attentional model (*base*). We italicize some *correct* translation segments and highlight a few **wrong** ones in bold.

5. Conclusion

- NMT에 Attention을 더 효과적으로 사용하기 위한 방법 제시
 - > Global Attention / Local Attention
- 다양한 Alignment Function 비교
- Attention-based NMT 짱!
 - ex) Translating names and Handling long sentences
- SOTA for both WMT'14 and WMT'15

Q&A

참고

<https://hcnoh.github.io/2018-12-11-bahdanau-attention> Bahdanau Attention 개념 정리

<https://hcnoh.github.io/2019-01-01-luong-attention> Luong Attention 개념 정리