

Transformer: Attention is All You Need

Ki Hyun Kim

nlp.with.deep.learning@gmail.com

Since 2016,

- In 2016, Google published their neural machine translation system(GNMT), which outperforms previous traditional MT system.

Google's Neural Machine Translation System: Bridging the Gap
between Human and Machine Translation

Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi
yonghui,schuster,zhifengc,qvl,mnorouzi@google.com

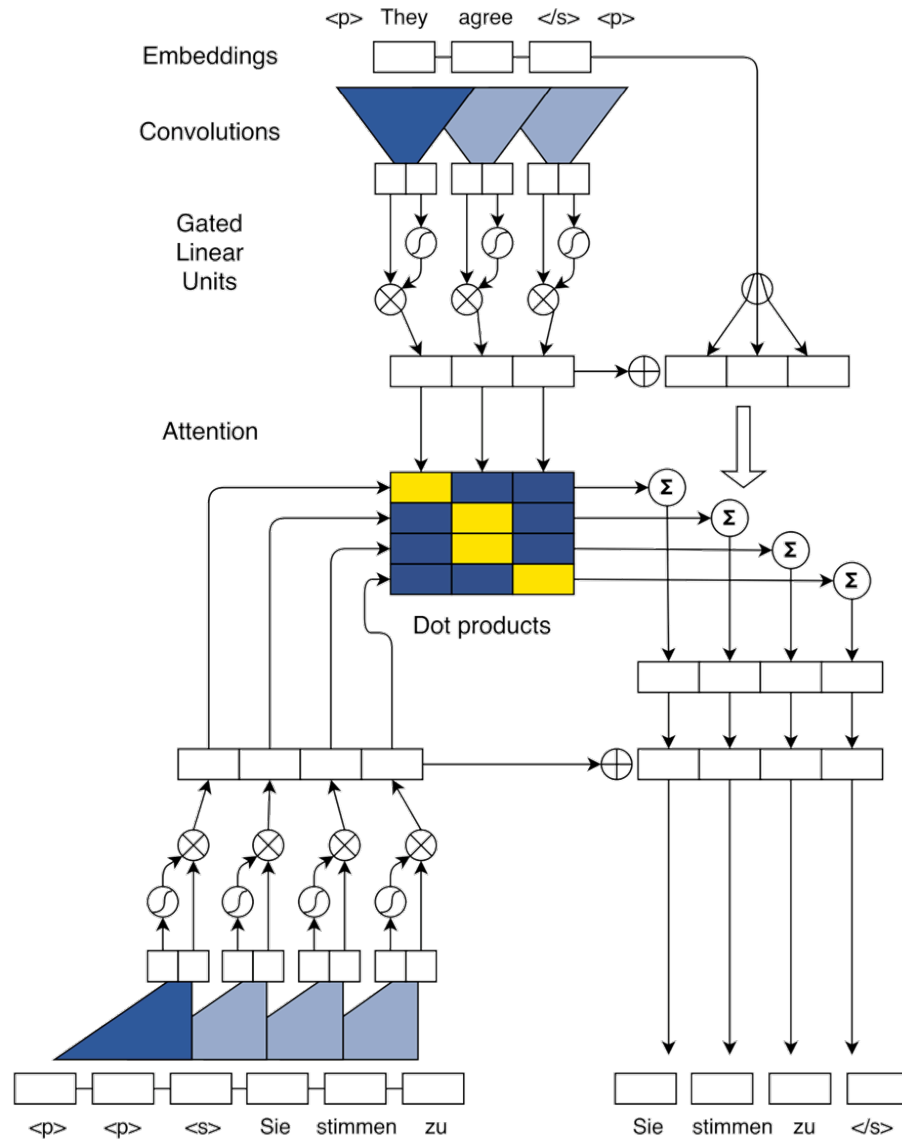
Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey,
Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser,
Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens,
George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa,
Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, Jeffrey Dean

- However, RNN based Sequence-to-sequence reveals its limitations.

Table 10: Mean of side-by-side scores on production data

	PBMT	GNMT	Human	Relative Improvement
English → Spanish	4.885	5.428	5.504	87%
English → French	4.932	5.295	5.496	64%
English → Chinese	4.035	4.594	4.987	58%
Spanish → English	4.872	5.187	5.372	63%
French → English	5.046	5.343	5.404	83%
Chinese → English	3.694	4.263	4.636	60%

Fully Convolutional Seq2Seq [Gehring et al., 2017]



WMT'16 English-Romanian		BLEU
Sennrich et al. (2016b) GRU (BPE 90K)		28.1
ConvS2S (Word 80K)		29.45
ConvS2S (BPE 40K)		30.02
WMT'14 English-German		BLEU
Luong et al. (2015) LSTM (Word 50K)		20.9
Kalchbrenner et al. (2016) ByteNet (Char)		23.75
Wu et al. (2016) GNMT (Word 80K)		23.12
Wu et al. (2016) GNMT (Word pieces)		24.61
ConvS2S (BPE 40K)		25.16
WMT'14 English-French		BLEU
Wu et al. (2016) GNMT (Word 80K)		37.90
Wu et al. (2016) GNMT (Word pieces)		38.95
Wu et al. (2016) GNMT (Word pieces) + RL		39.92
ConvS2S (BPE 40K)		40.51

Table 1. Accuracy on WMT tasks compared to previous work. ConvS2S and GNMT results are averaged over several runs.

Attention is All You Need

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 ...

facebook

🔖 10181회 인용 관련 학술자료 전체 21개의 버전 🔖



Transformer

- Encoder + Decoder

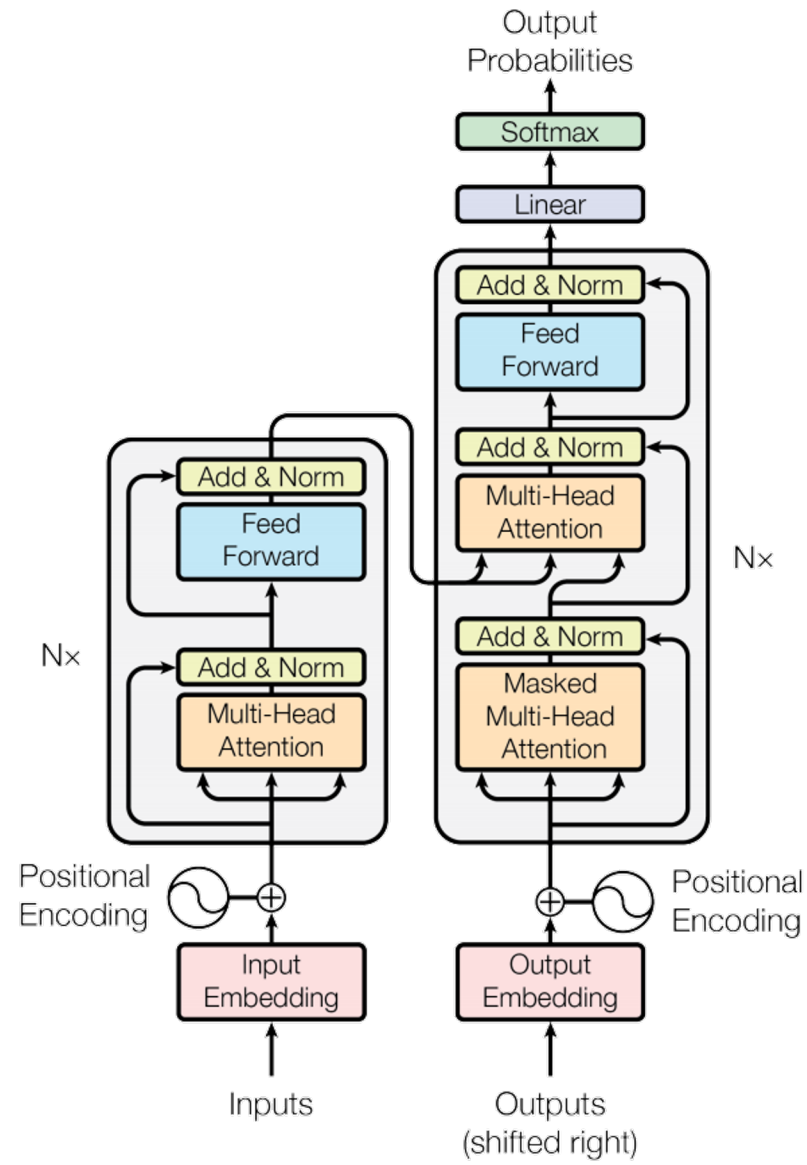


Figure 1: The Transformer - model architecture.

Transformer

- 일타쌍피: 성능과 속도 모두 기존 모델을 압도

Table 2: The Transformer achieves better BLEU scores than previous state-of-the-art models on the English-to-German and English-to-French newstest2014 tests at a fraction of the training cost.

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	41.29	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	$3.3 \cdot 10^{18}$	
Transformer (big)	28.4	41.8	$2.3 \cdot 10^{19}$	

Transformer and NLP

- Pretraining and finetuning (Transfer learning) with Big-LM.

