

## 트랜스포머 기반 영문 인공지능 PDF 논문 번역을 위한 기계번역 모델

함지율 이현도  
강남대학교 데이터사이언스 전공

E-mail: hramsm@naver.com, hanse199911@naver.com

## Introduction 서론

시중에 네이버 파파고, 구글 번역기 등 여러 번역기가 나와 있습니다. 하지만 특정 도메인 지식을 반영하여 번역하지 못한다는 단점을 파악하였습니다. 예를 들어, 이상 탐지 알고리즘 논문에 **contamination**을 ‘오염’이라고 번역합니다. 하지만 도메인 지식에 따르면 ‘데이터의 불순도’로 의역할 수 있습니다. 이러한 점을 반영하여 (1) 영문 인공지능 논문을 한국어로 번역하여 학생들이 더욱 편하게 읽을 수 있도록 하고자 주제를 선정하였습니다. 그리고 논문의 문장을 하나하나 스크롤 하여 번역기에 입력하는게 불편하였습니다. 저희는 PDF의 텍스트 또한 편리하게 (2) PDF를 입력하면 이미지로 변환하고 변환된 이미지의 텍스트를 인식하여 한영 번역을 해주는 기계번역 모델(NMT)를 구축하는 것을 목표로 프로젝트를 진행했습니다.

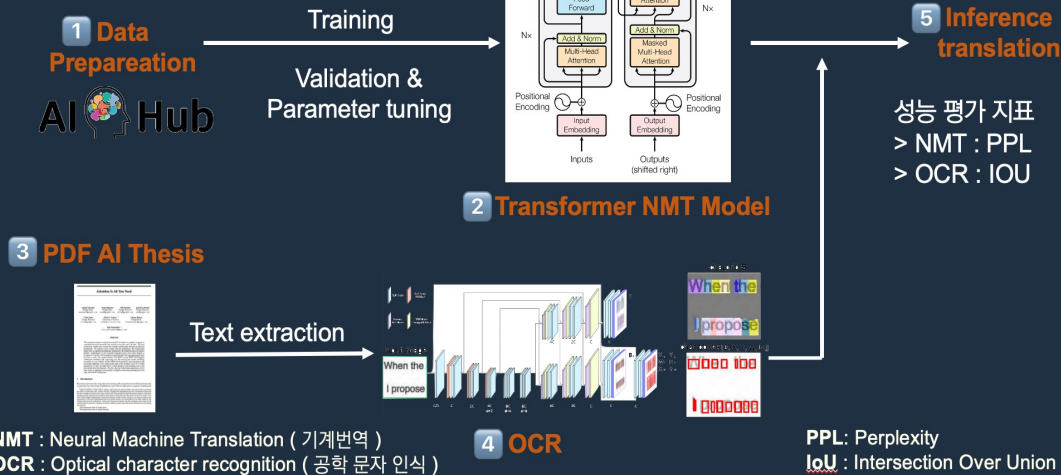
NMT : Neural Machine Translation (기계번역)

OCR : Optical Character Recognition (광학 문자 인식)

## Method 방법론

- 1) 수집한 데이터를 학습 데이터로 사용할 수 있도록 준비
- 2) Transformer를 통해 기술과학 논문 영한 번역 데이터와 일상 구어체 데이터를 통해 학습
- 3,4) OCR(Optical character recognition) 기술로 PDF의 텍스트를 추출
- 5) 추출한 텍스트를 transformer 모델을 통해 추론(번역)

## Project Overview



## Result 결과

- 1) Transformer를 소개한 ‘Attention Is All You Need’ 논문 PDF에서 문서 추출
  - 2) 문서 추출된 텍스트를 번역 모델에 입력값으로 넣어 번역
- 번역 결과가 특정 단어 수준을 잡아내지만 문맥적인 요소와 자연스러운 문장 번역에는 한계가 있음

The image shows the title page of the paper "Attention Is All You Need" by Ashish Vaswani et al. and a corresponding code snippet for the transformer implementation.

**Attention Is All You Need**

**Abstract**

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that restrict the number of parallel computation paths. 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 requiring less parameters and computing resources than those of the state-of-the-art models. On WMT 2014 English-to-French translation task, the Transformer achieves 28.4 BLEU on the WMT 2014 English-to-French translation task, surpassing the best performing sequence-to-sequence model 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 41.8 after waiting for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**1 Introduction**

Recent neural networks, long short-term memory [13] and gated recurrent [7] neural networks in particular, have been firmly established as state of the art approaches in sequence modeling and "sequence-to-sequence" tasks. In this paper, 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 requiring less parameters and computing resources than those of the state-of-the-art models. On WMT 2014 English-to-French translation task, the Transformer achieves 28.4 BLEU on the WMT 2014 English-to-French translation task, surpassing the best performing sequence-to-sequence model 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 41.8 after waiting for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.

**1.1 Hardware and Software**

We trained our models on a single NVIDIA V100 GPU. Each training step took about 0.4 seconds. We trained the base model for a total of 100,000 steps or 12 hours. For our big model, described on subsection line of table 1, step time was 1.8 seconds. The big models were trained for 100,000 steps or 12 days.

**1.2 Replication**

We used the Allen2 optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.98$  and  $\epsilon = 1e-9$ .

We varied the learning rate over the course of training, according to the formula  $lr = lr_{init} \cdot \sqrt{\text{step\_num} / \text{step\_max}}$ .

We also corresponded to increasing the learning rate linearly for the first  $\text{warm\_up\_steps}$  training steps, and decreasing it thereafter proportionally to the inverse square root of the  $\text{step\_num} - \text{warm\_up\_steps} + 1$ .

**1.3 Regularization**

We apply dropout to the output of each sub-layer, before it is added to the sub-layer input.

## Discussion 토의

- 성능 개선을 위한 추가적인 방안은 어떤 것이 있을까?
- 모델 학습시 하이퍼 파라미터를 어떻게 지정하는 것이 가장 효과적일까?
- 기계 번역에 State-of-the-art 성능을 보이는 다른 모델을 통해 학습하면 어떨까?

## References 참고

1. Attention Is All You Need <https://arxiv.org/abs/1706.03762>
2. Transformer 코드 구현 깃허브 <https://github.com/navnoes/pytorch-transformer>
3. OCR 기술 블로그 <https://wandukong.tistory.com/9>
4. Bert Tokenizer <https://huggingface.co/bert-base-multilingual-cased>