

# Introduction to Machine Learning

All for One

Choi Minjoo

# Choi Minjoo (Judy Choi)

- Careers
  - 2012 ~ 2019 **Malware Analyst**
  - 2018 **Working Holiday** in France
  - 2020 ~ **M.S** in Kangwon Univ
    - Intelligence Software Lab
      - **NLP (Machine Translation)**
  - 2021 ~ 2022 **Bering Lab** (NLP Researcher & Engineer)
- SNS
  - <https://www.facebook.com/minjoo.choi.562/>
  - <https://github.com/Judy-Choi>

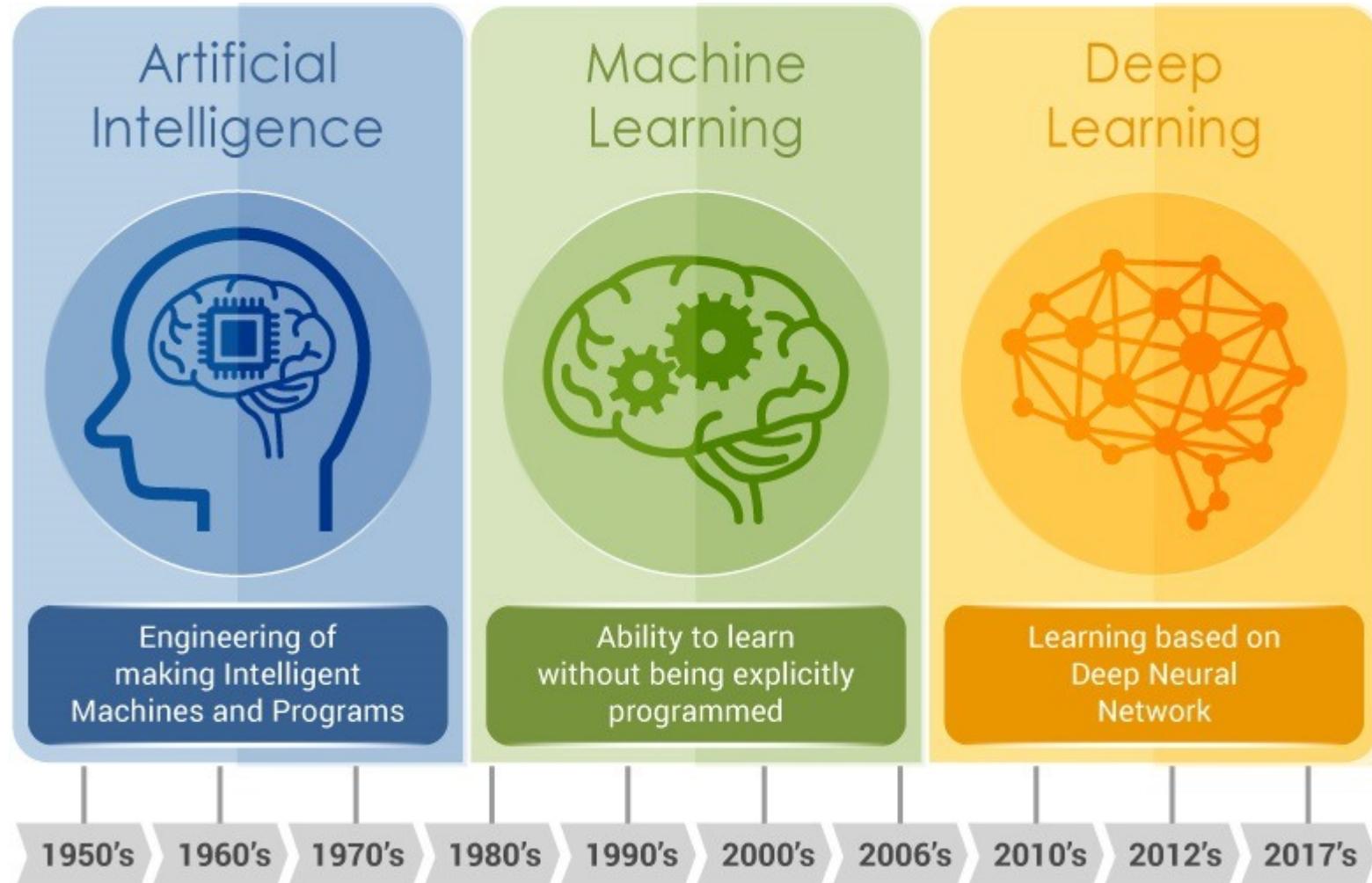


# Contents

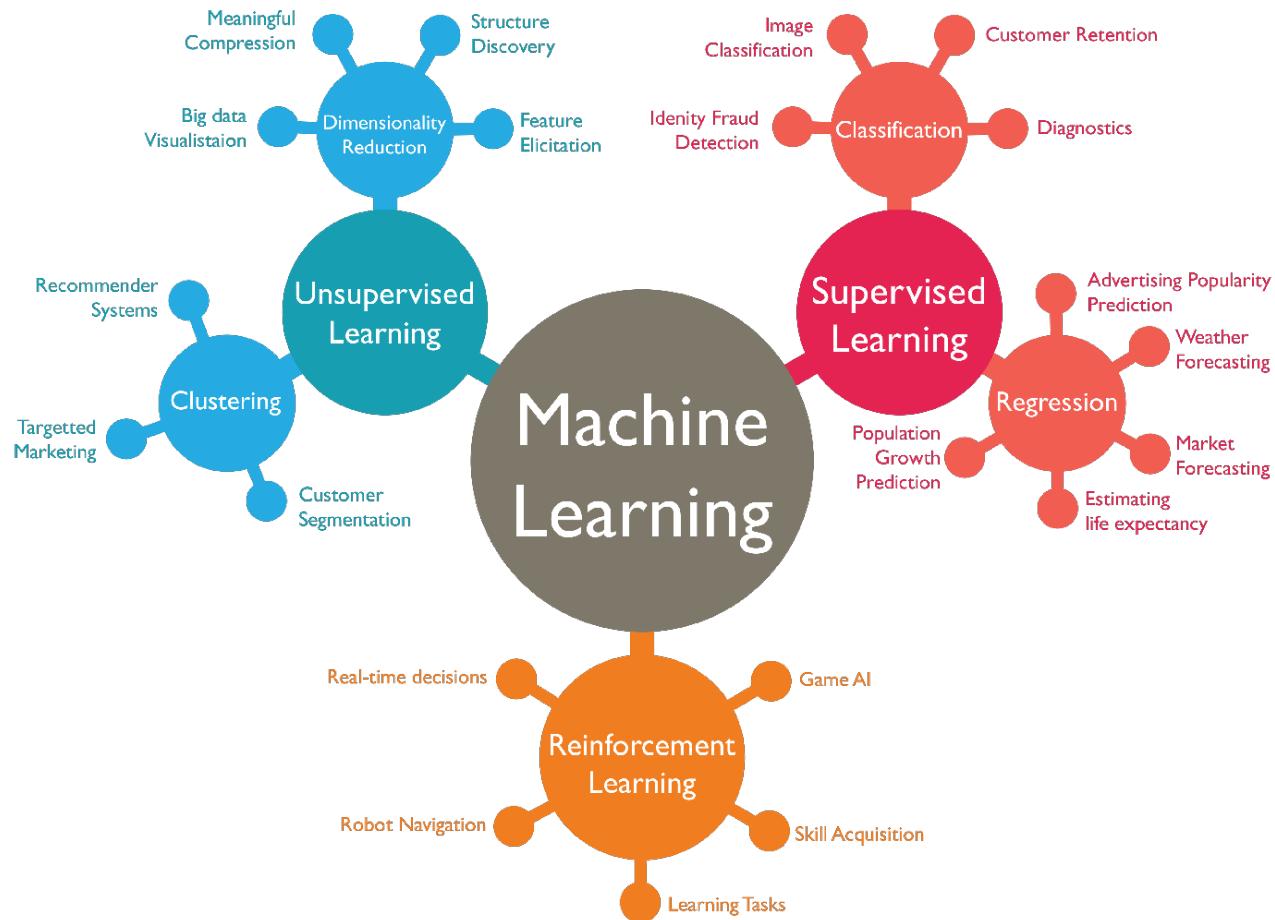
- Abstract of Deep Learning
- Natural Language Processing (NLP)
- Machine Translation

# Abstract of Deep Learning

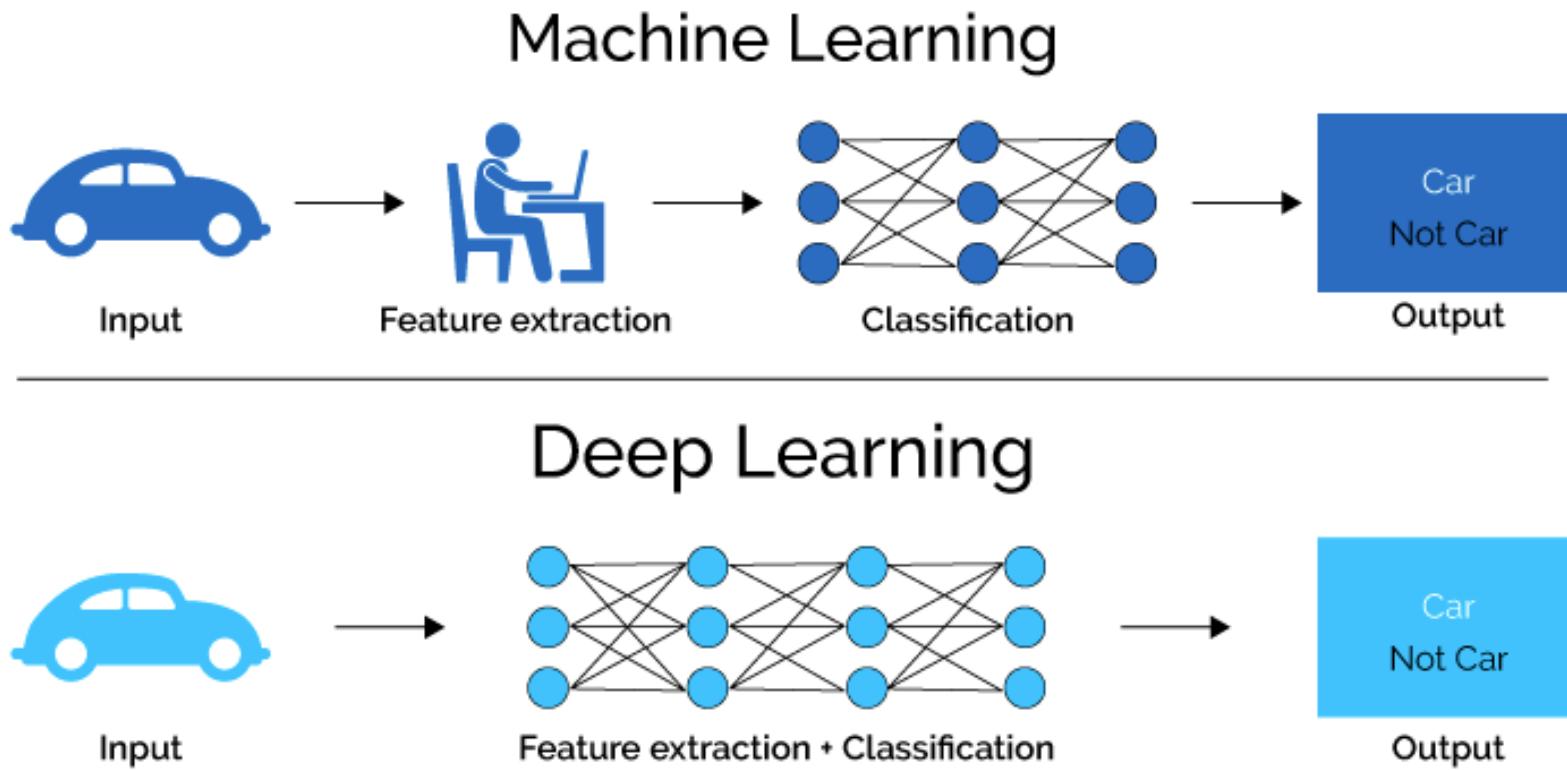
# from A.I to Deep Learning...



# Machine Learning



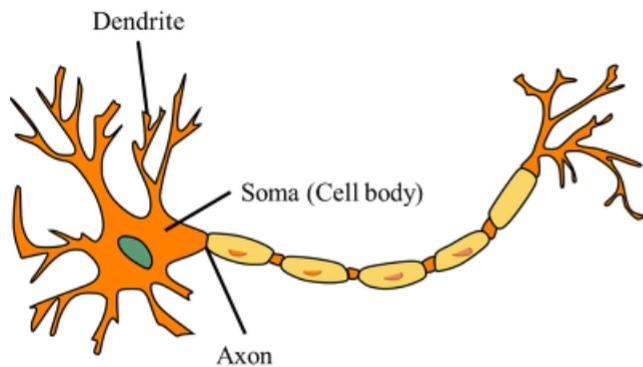
# Machine Learning vs Deep Learning



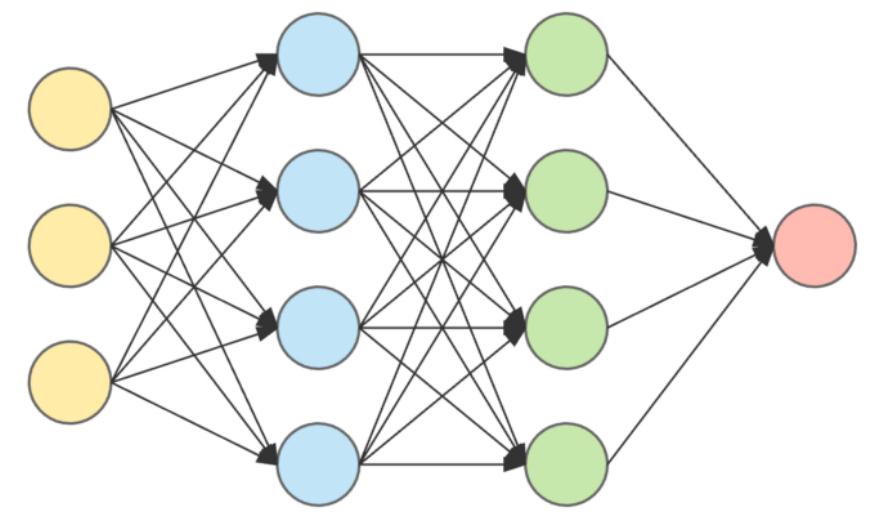
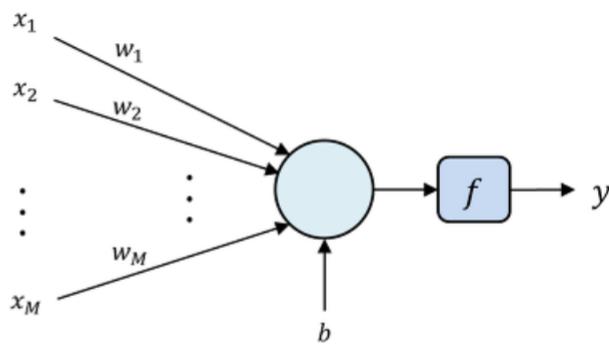
# Deep Learning

- Neural Network

- Artificial neural network mimicking neural networks in the human brain.
- Brain : Neuron = A.I : Perceptron



[그림 1] 생물체의 neuron (좌)과 artificial neuron (우)

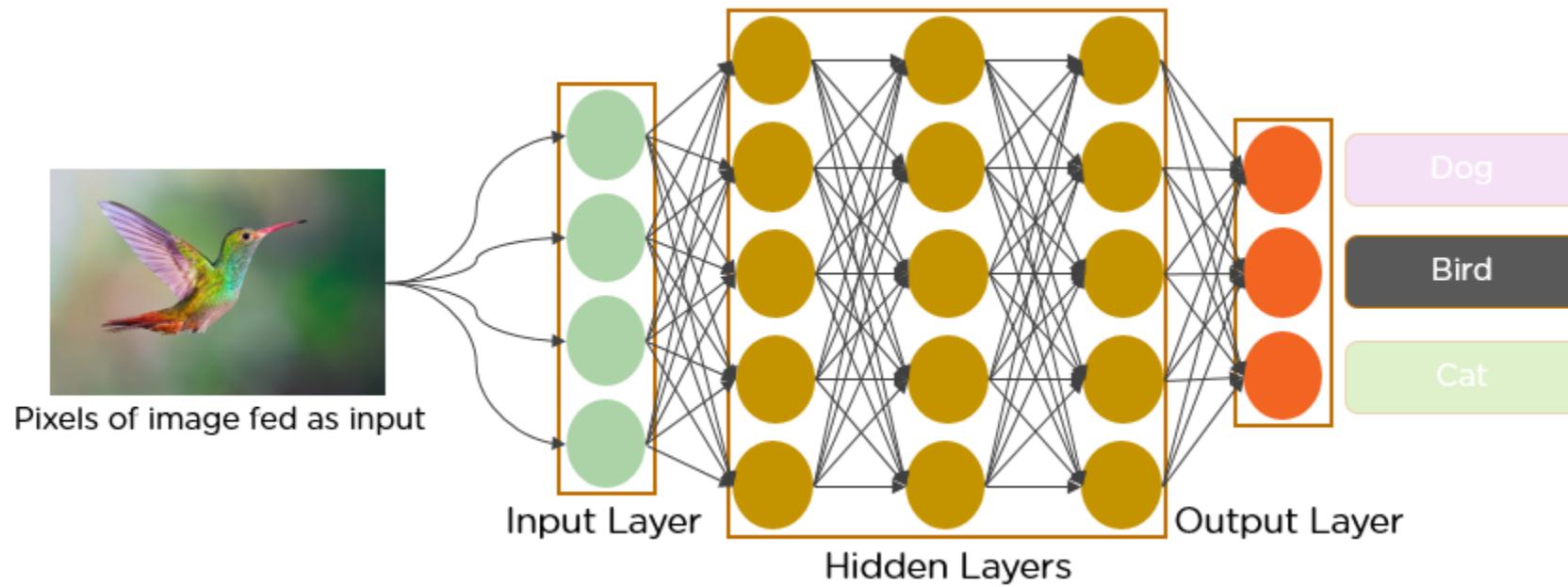
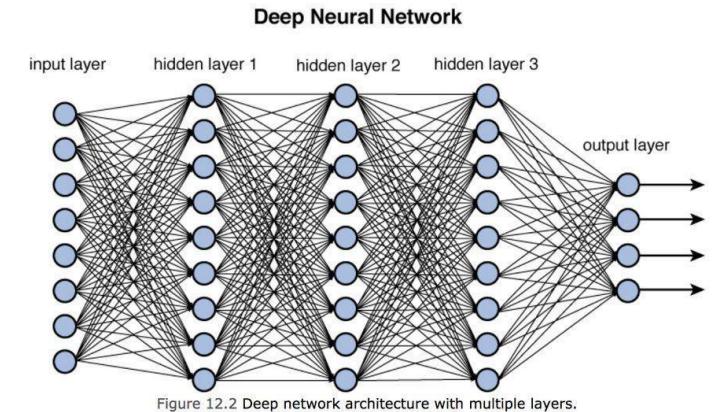


<https://smartstuartkim.wordpress.com/2019/01/27/history-of-neural-networks-1-perceptron/>

<https://gongster.medium.com/how-does-a-neural-network-work-intuitively-in-code-f51f7b2c1e3f>

# Deep Learning

- Deep Neural Network (DNN)
  - A network of multiple layers of perceptron layers

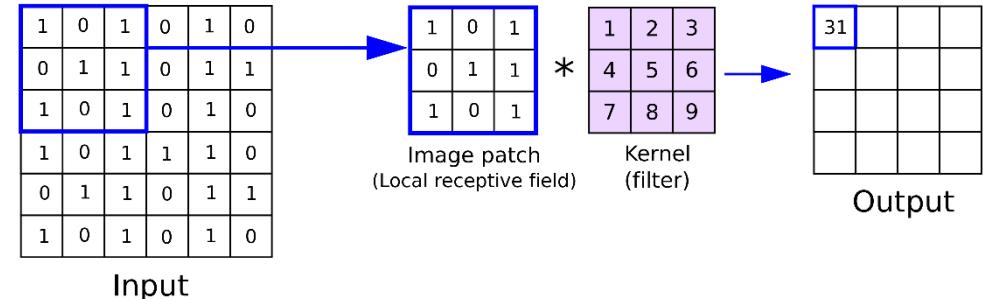
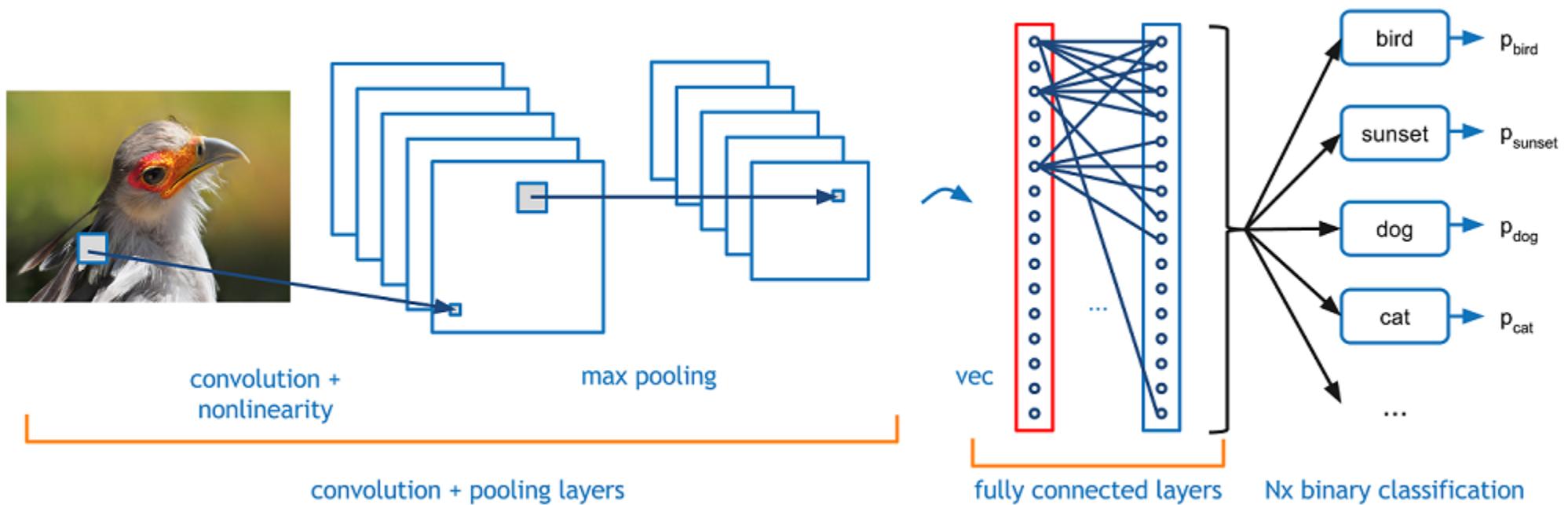


<https://towardsdatascience.com/training-deep-neural-networks-9fdb1964b964>

<https://www.analyticsvidhya.com/blog/2021/05/convolutional-neural-networks-cnn/>

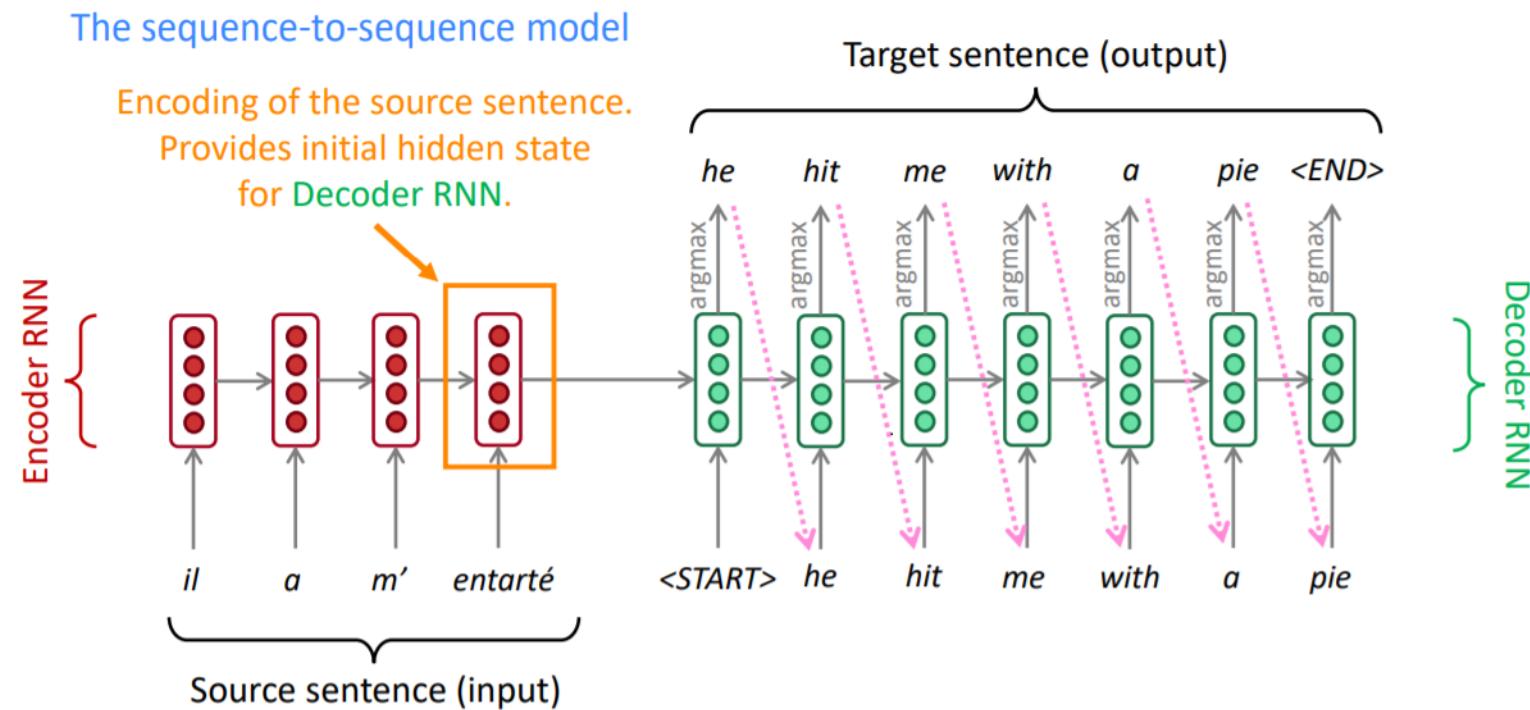
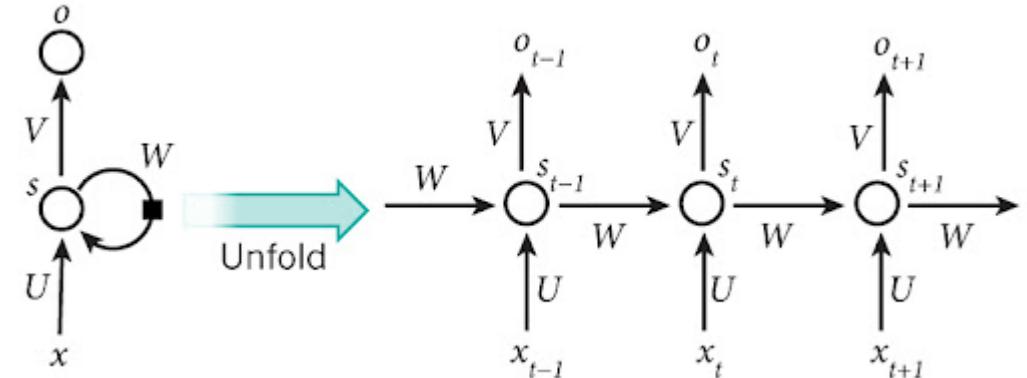
# Deep Learning

- CNN (Convolution Neural Network)



# Deep Learning

- RNN (Recursive Neural Network)



<https://aikorea.org/blog/rnn-tutorial-1/>

<https://medium.com/analytics-vidhya/neural-machine-translation-using-bahdanau-attention-mechanism-d496c9be30c3>

# Natural Language Processing (NLP)

# NLP?

- ‘Interactions between computers and human language’
- A subfield of linguistics, computer science, and artificial intelligence
- Programming computers to process and analyze large amounts of natural language data.
- The goal is a computer capable of "understanding" the contents of documents, including the contextual nuances of the language within them.

# Field of application

- Chatbot
  - Alexa, Siri, Bixby.
- Text Summarization
  - Google search
- Document Classification
  - Spam mail filtering
- Machine Translation
  - Google Translate



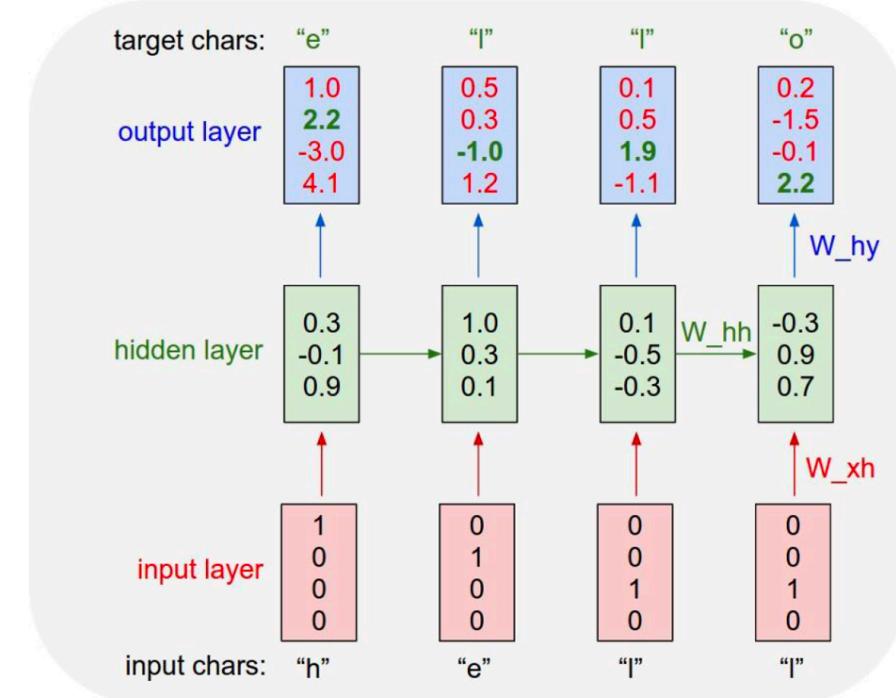
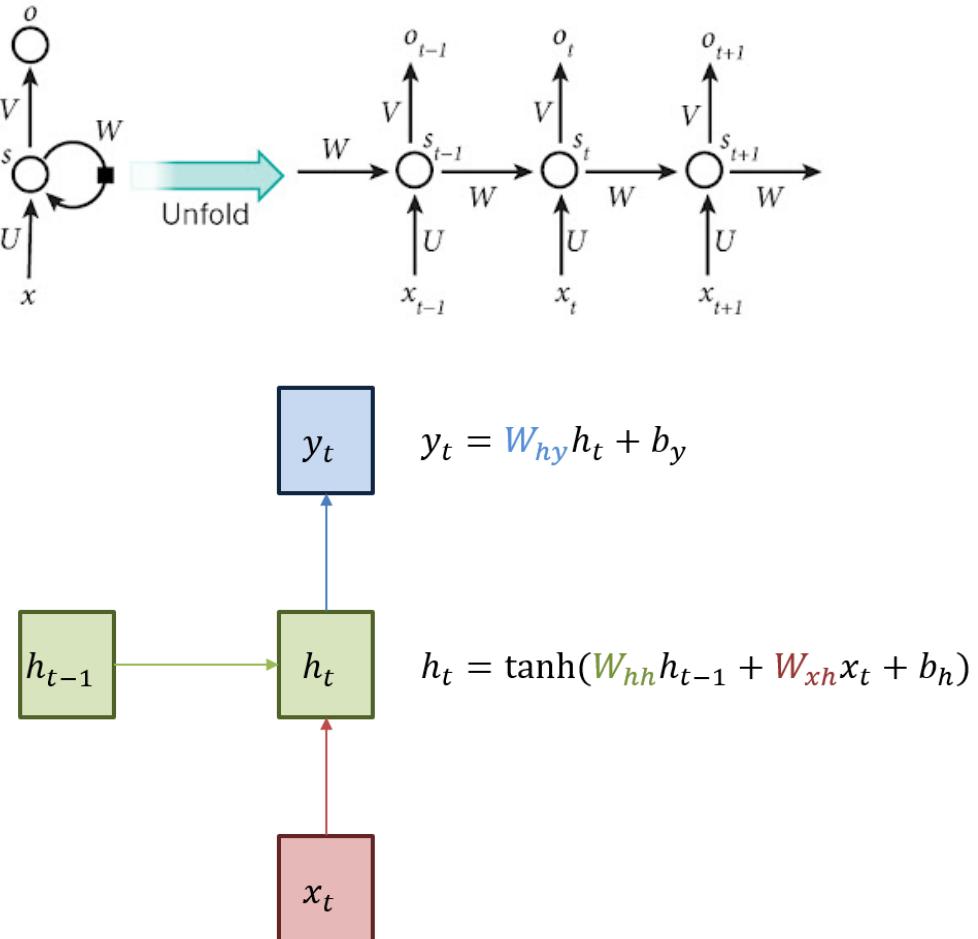
## Uzbekistan

Country in Central Asia

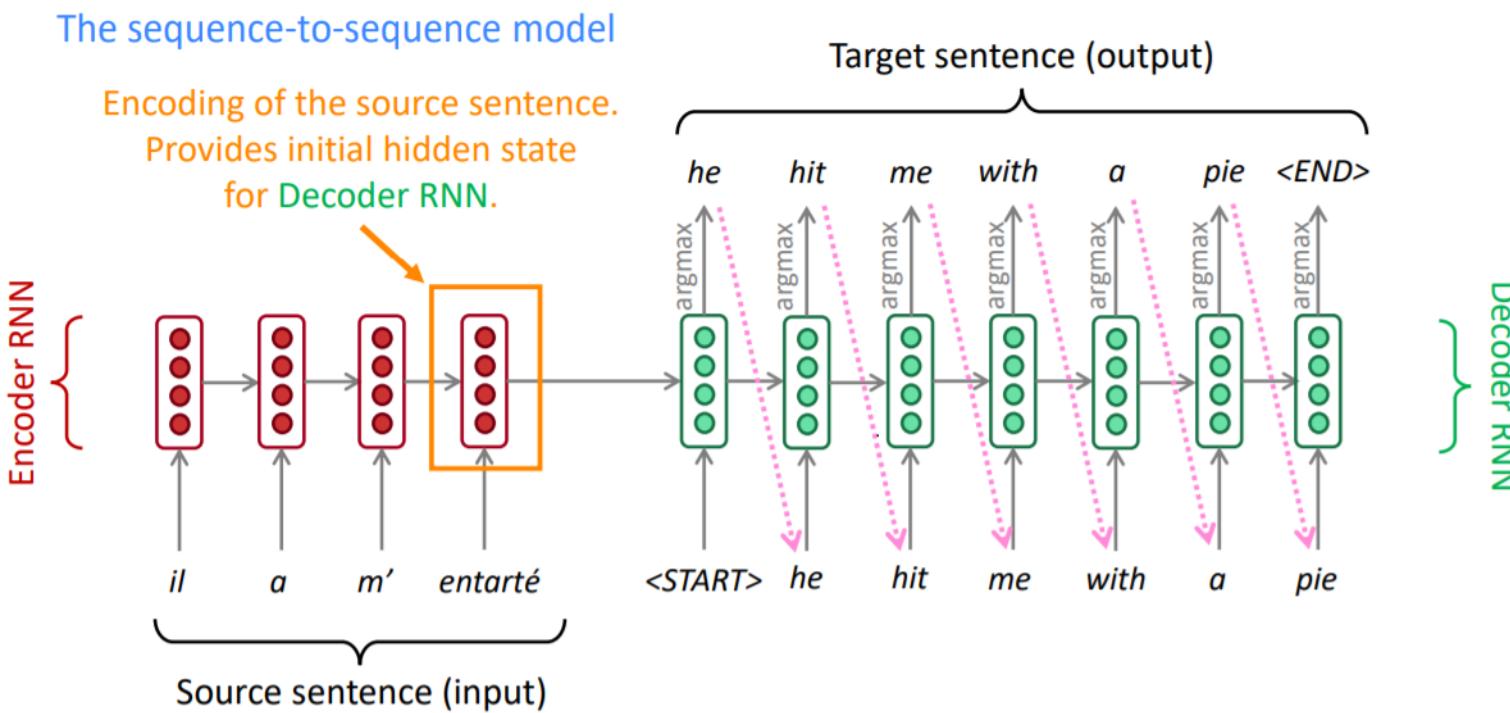
Uzbekistan is a Central Asian nation and former Soviet republic. It's known for its mosques, mausoleums and other sites linked to the Silk Road, the ancient trade route between China and the Mediterranean. Samarkand, a major city on the route, contains a landmark of Islamic architecture: the Registan, a plaza bordered by 3 ornate, mosaic-covered religious schools dating to the 15th and 17th centuries. — Google



# RNN

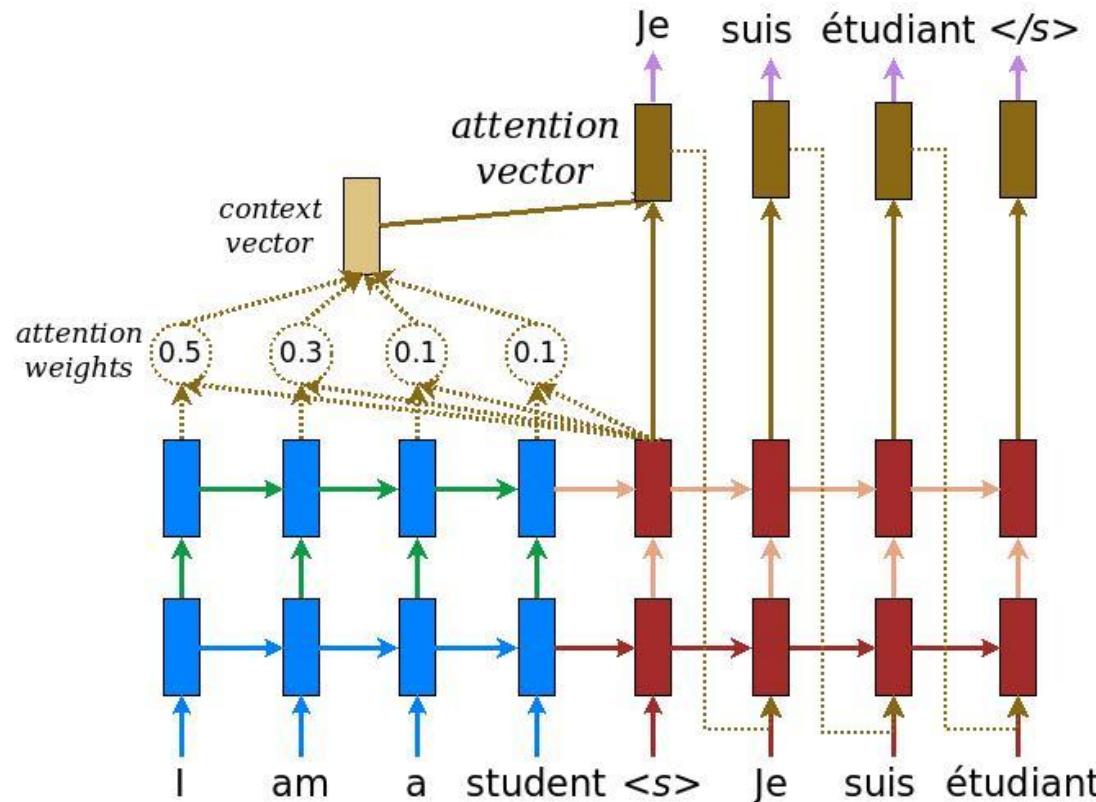


# Seq2Seq



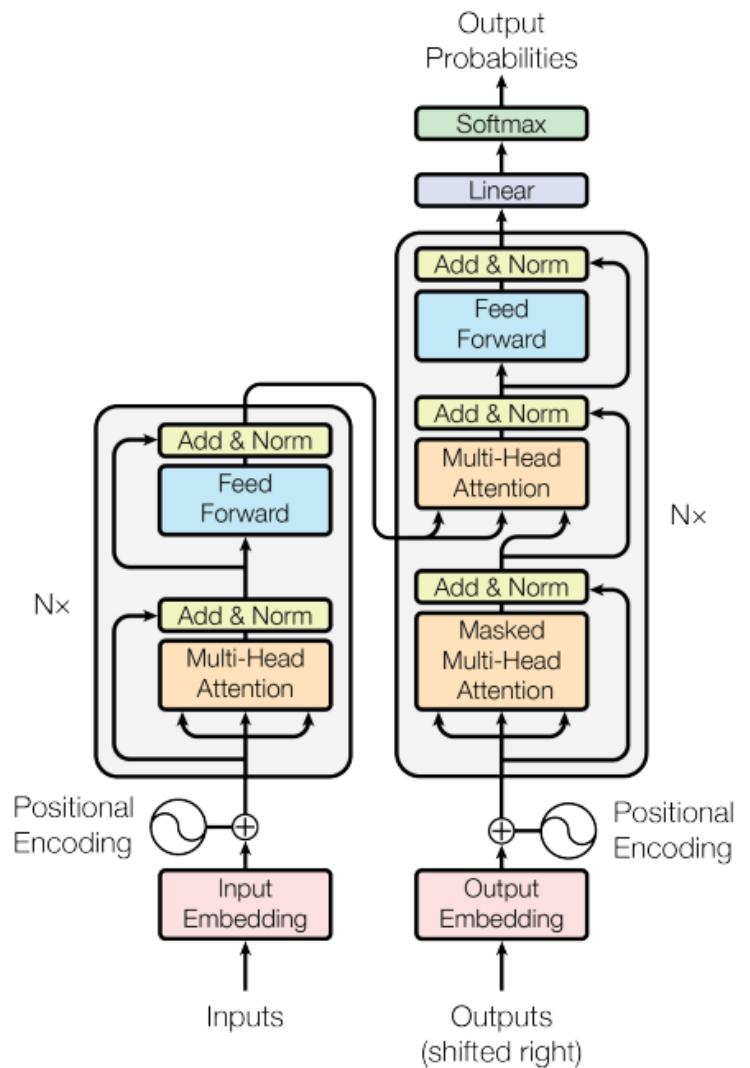
# Seq2Seq + ‘Attention’

- Decoder learn to focus over a specific range of the input sequence



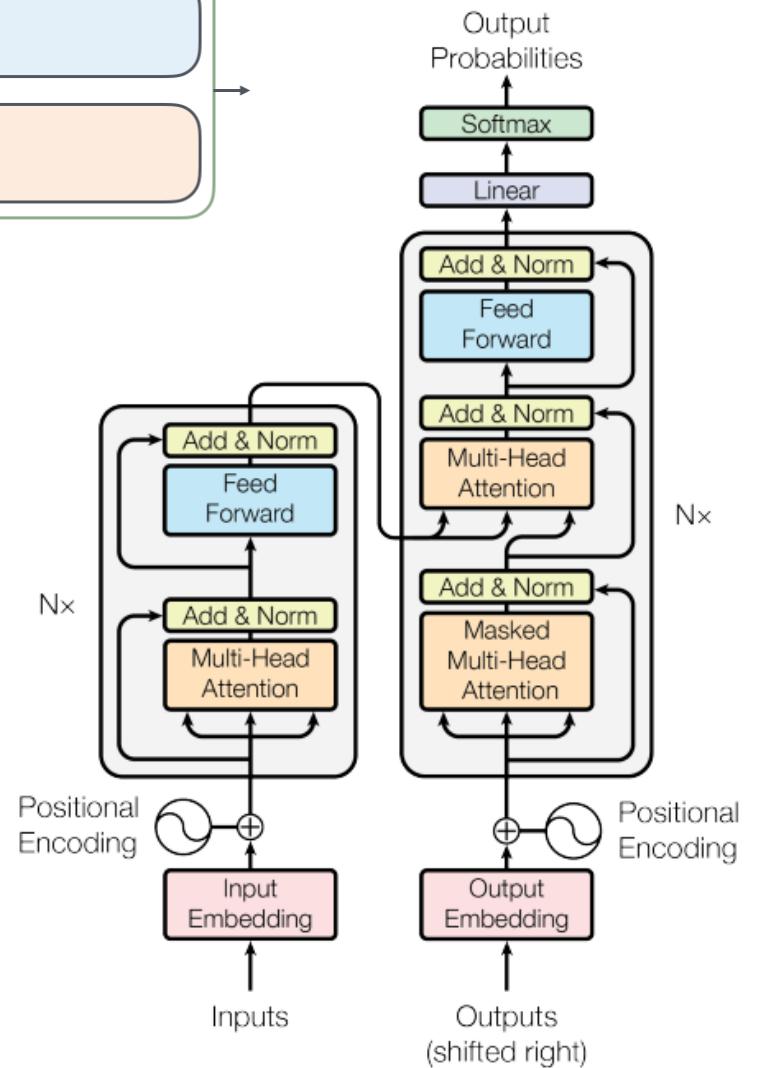
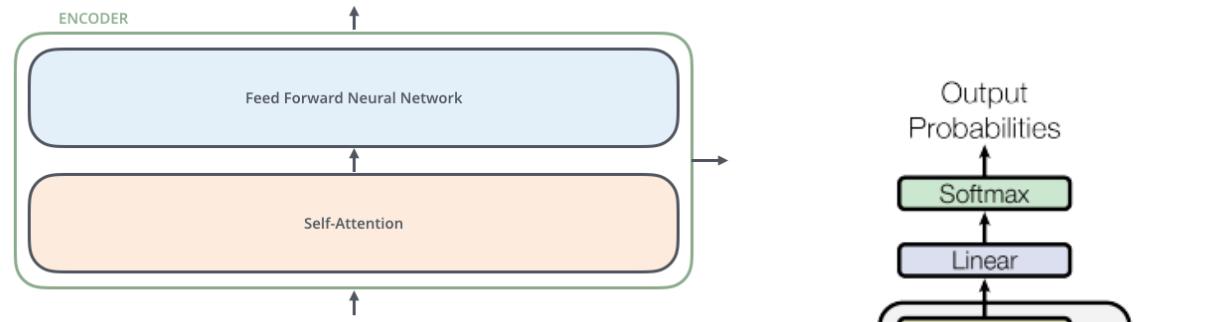
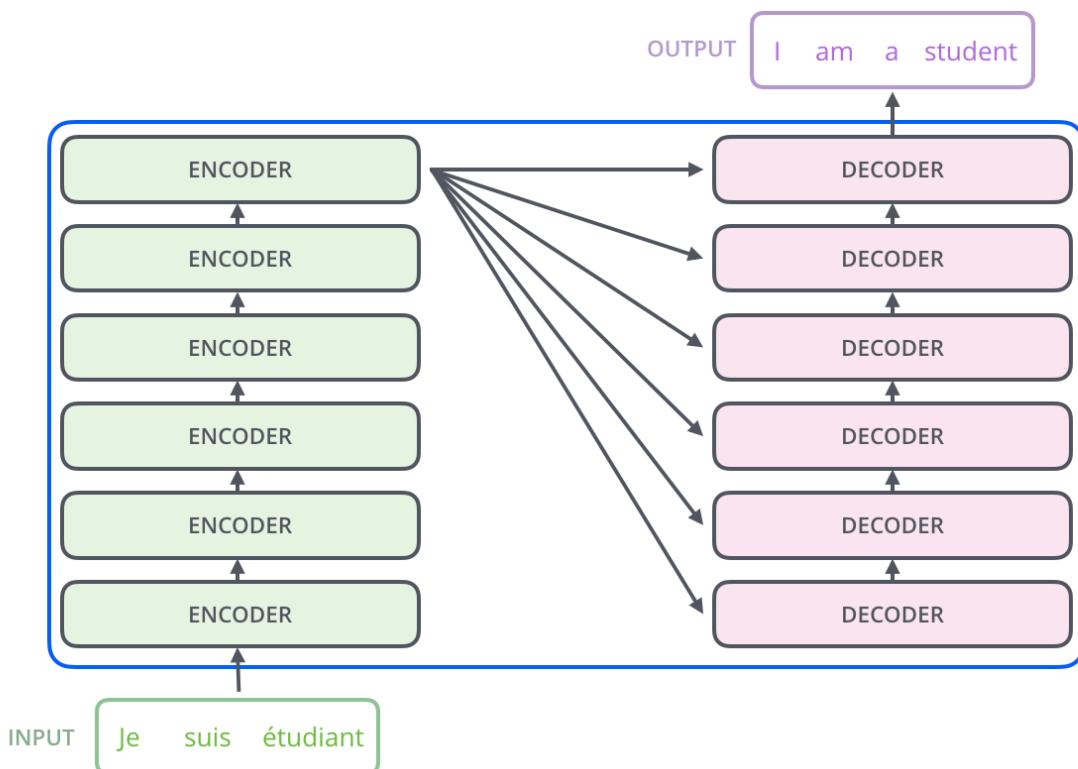
# Transformer

- ‘Attention is All you Need’



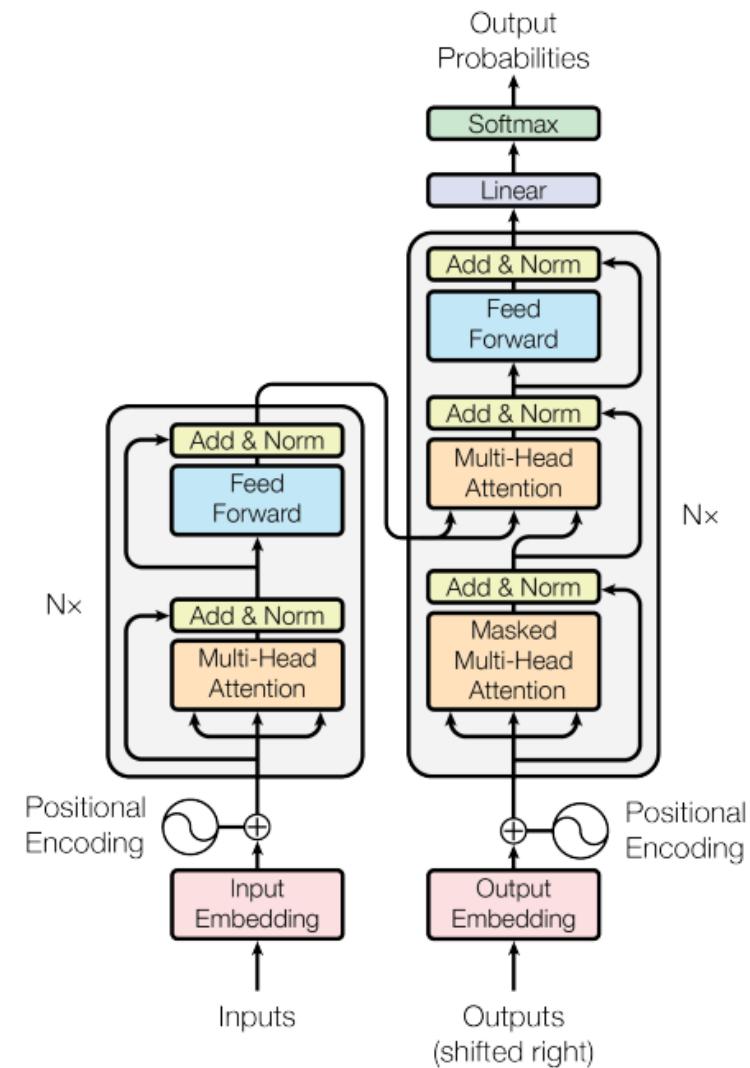
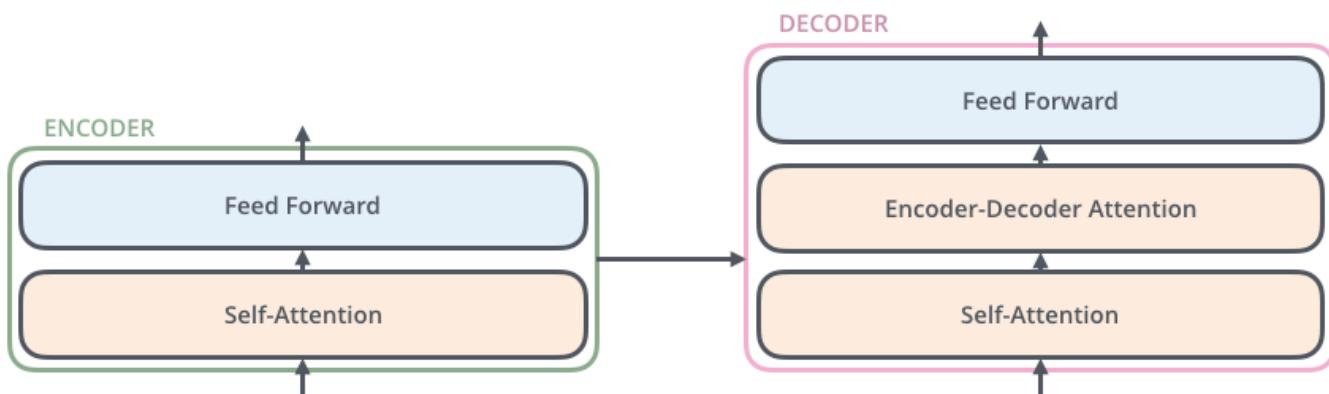
# Transformer

- Encoder – Decoder (6-6)



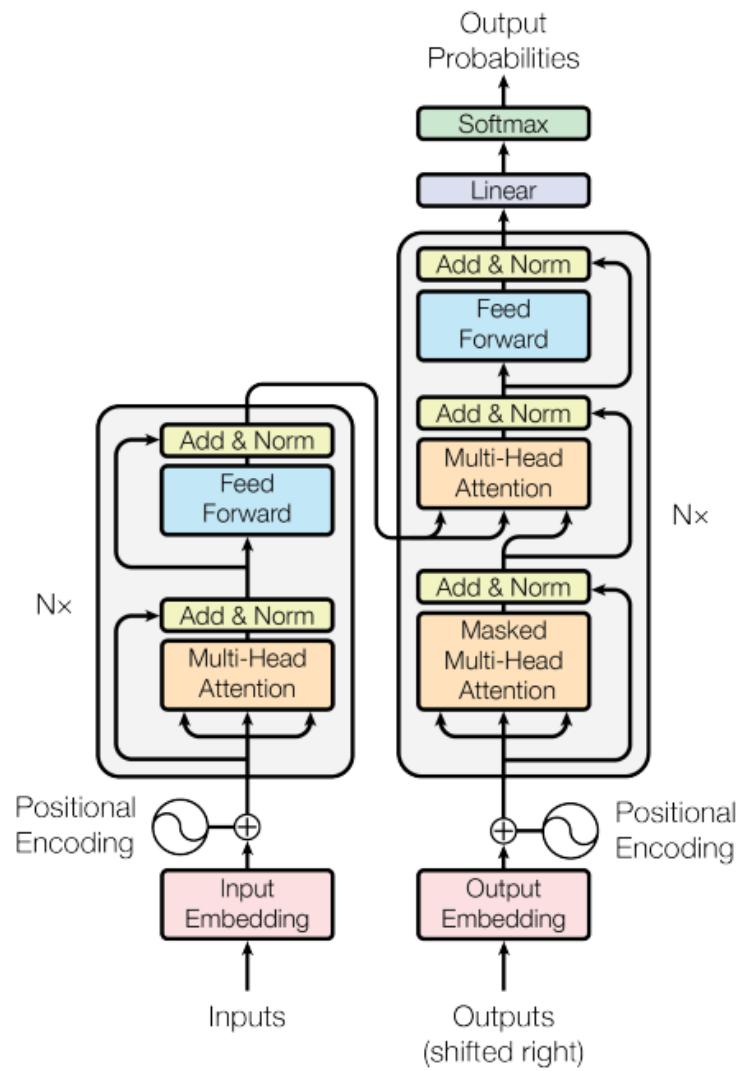
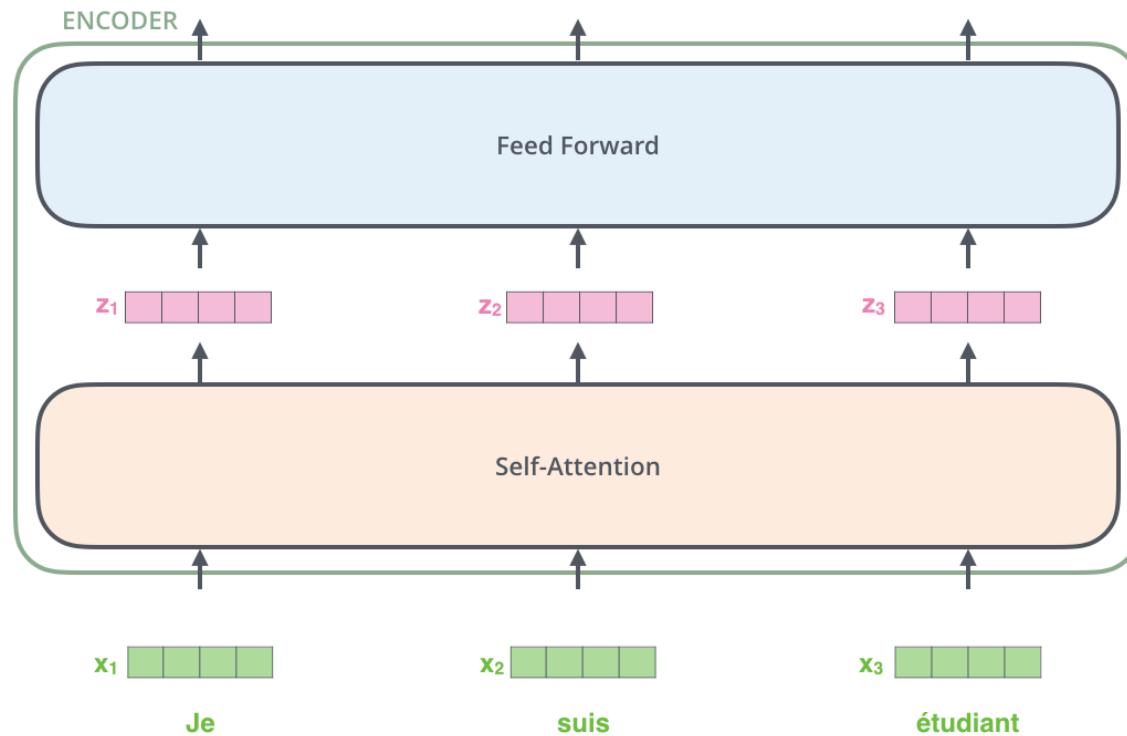
# Transformer

- FFNN & Self-Attn in Encoder



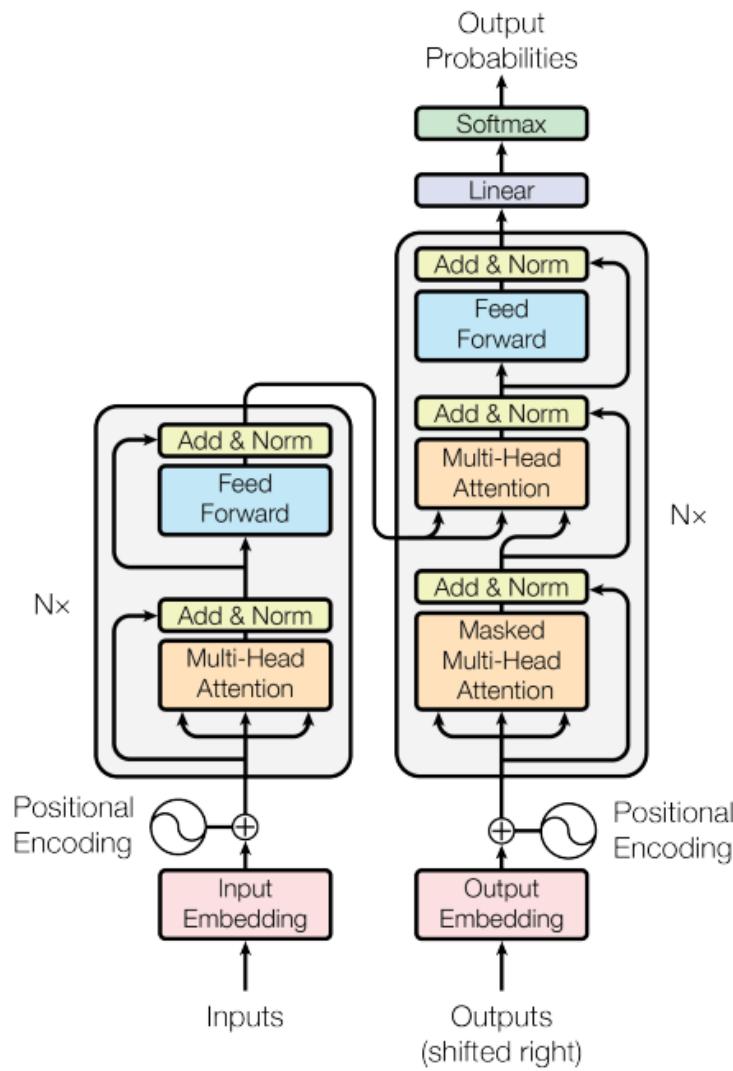
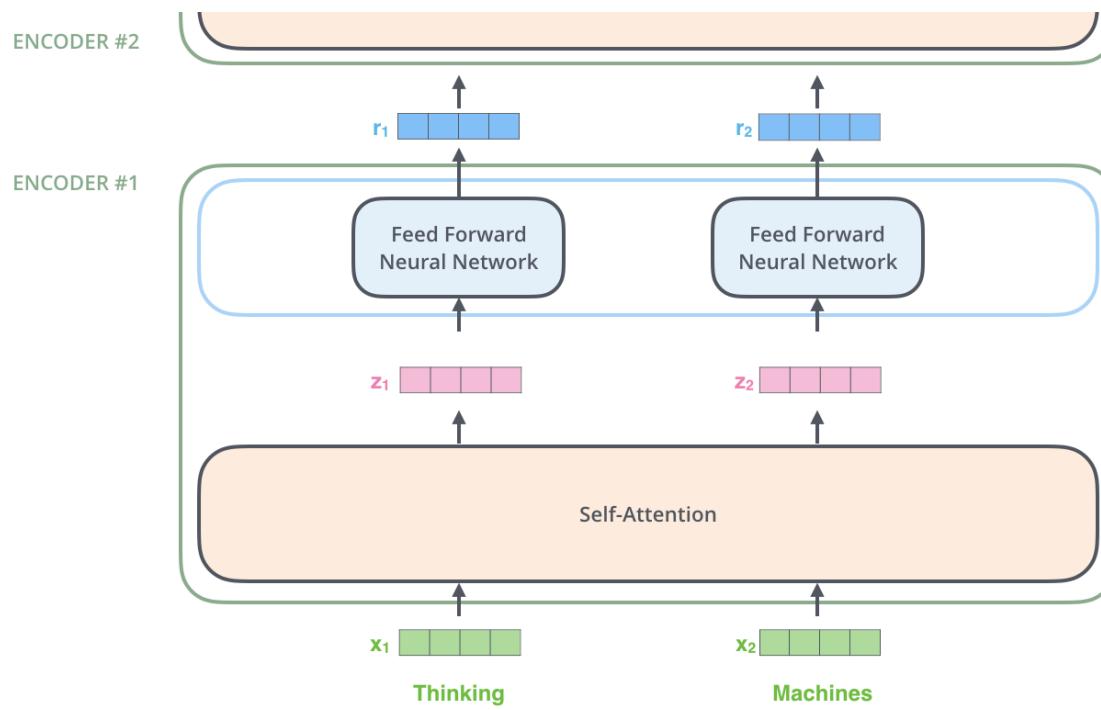
# Transformer

- Input embedding



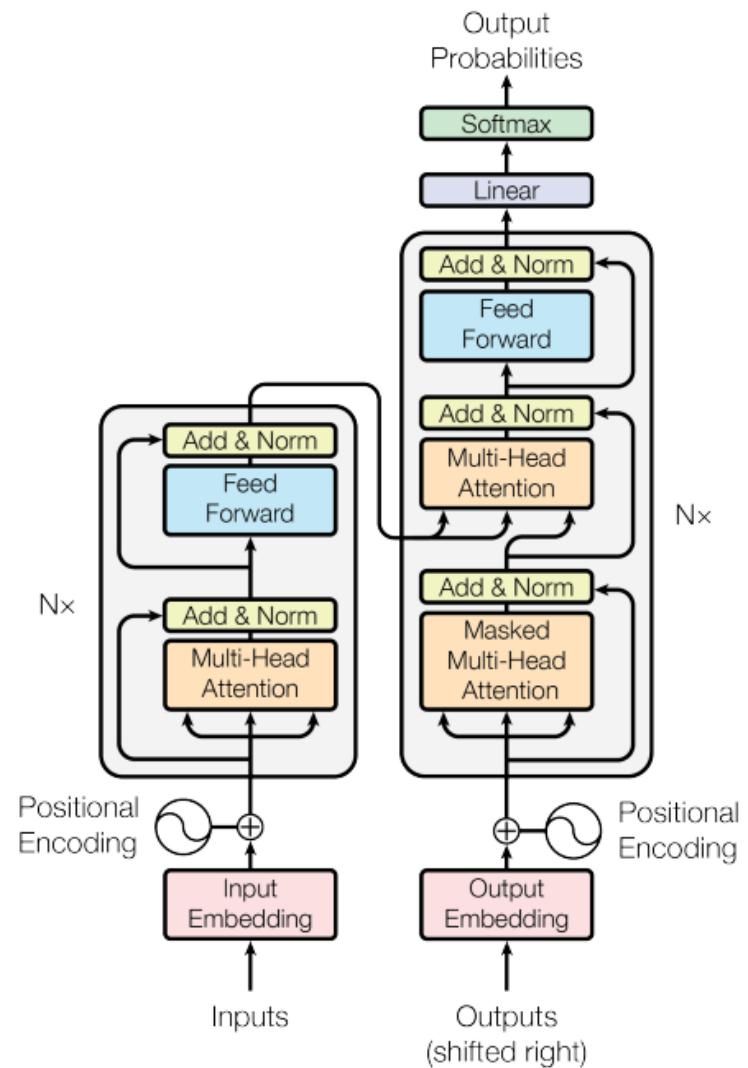
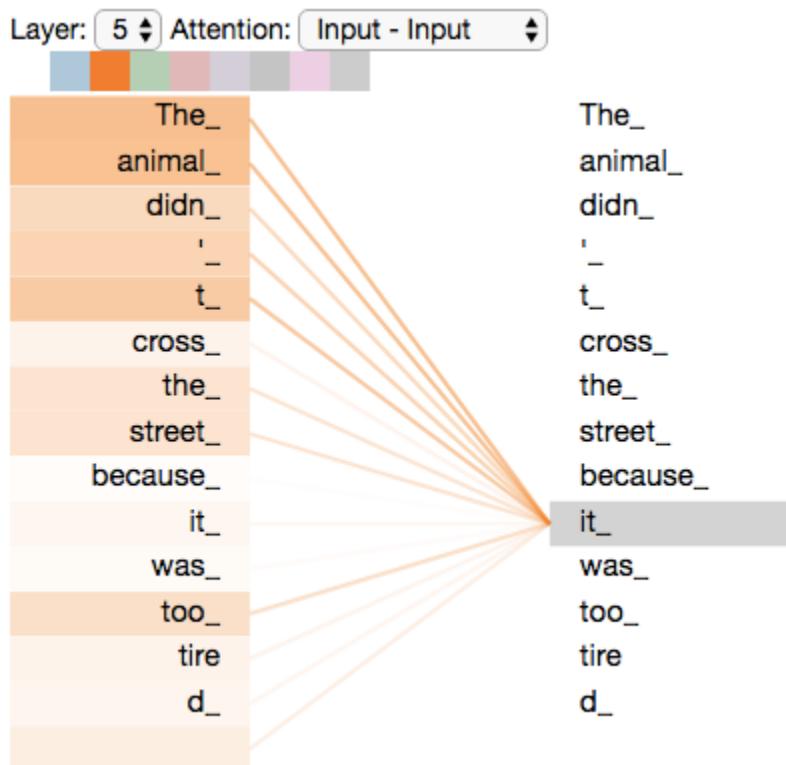
# Transformer

- Encoding



# Transformer

- Self-Attention



# Transformer

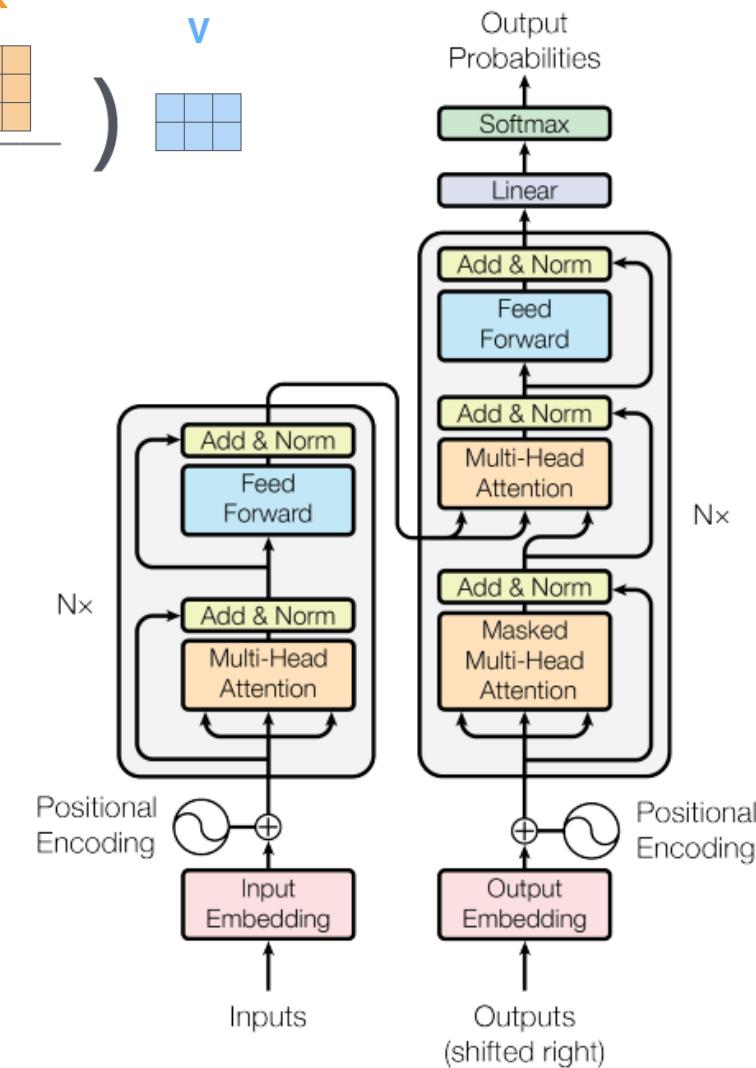
- Self-Attention

Input	Thinking		Machines	
Embedding	$x_1$		$x_2$	
Queries	$q_1$		$q_2$	
Keys	$k_1$		$k_2$	
Values	$v_1$		$v_2$	
Score	$q_1 \cdot k_1 = 112$		$q_1 \cdot k_2 = 96$	
Divide by 8 ( $\sqrt{d_k}$ )	14		12	
Softmax	0.88		0.12	
Softmax X Value	$v_1$		$v_2$	
Sum	$z_1$		$z_2$	

$$\text{softmax} \left( \frac{Q \times K^T}{\sqrt{d_k}} \right) V = Z$$

Matrix Calculation of Self-Attention

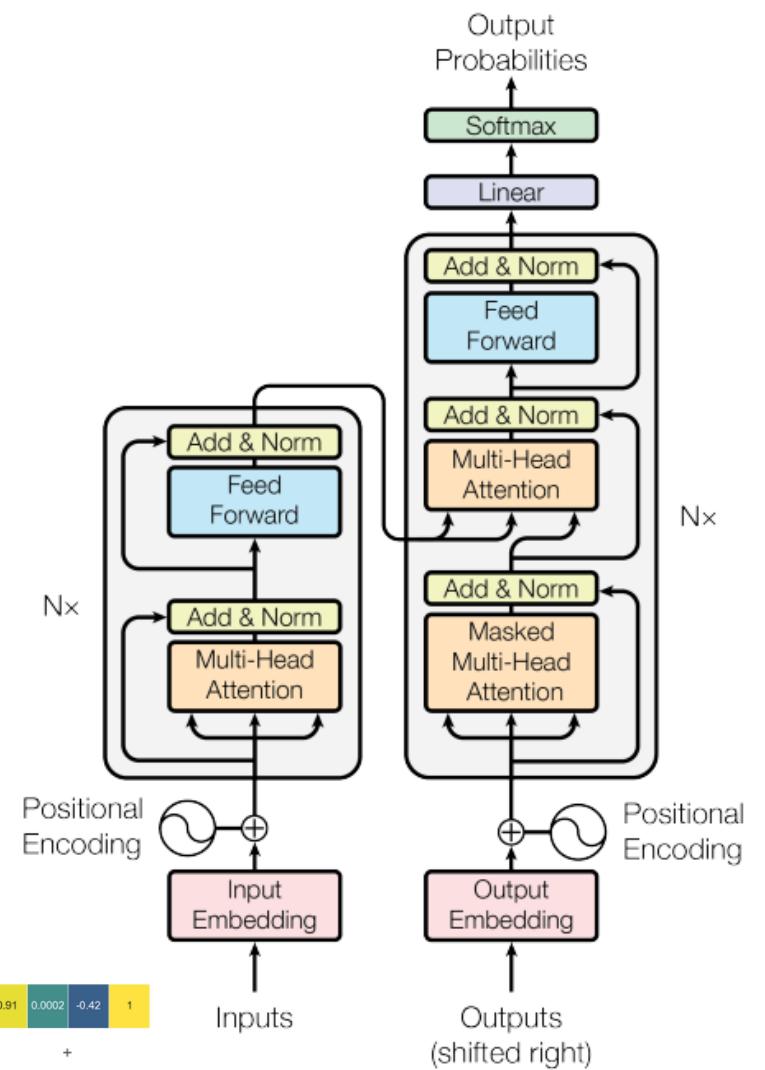
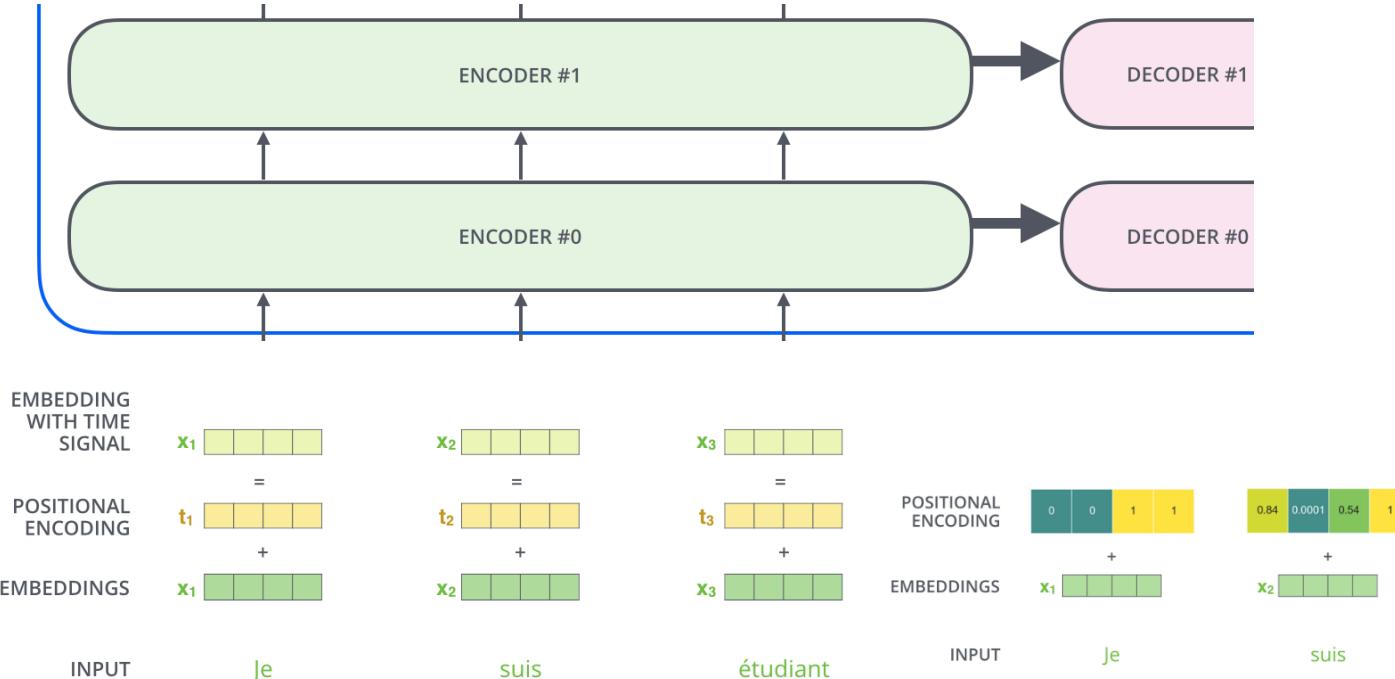
$$\begin{aligned} X &\times W^Q = Q \\ X &\times W^K = K \\ X &\times W^V = V \end{aligned}$$



<https://jalammar.github.io/illustrated-transformer/>

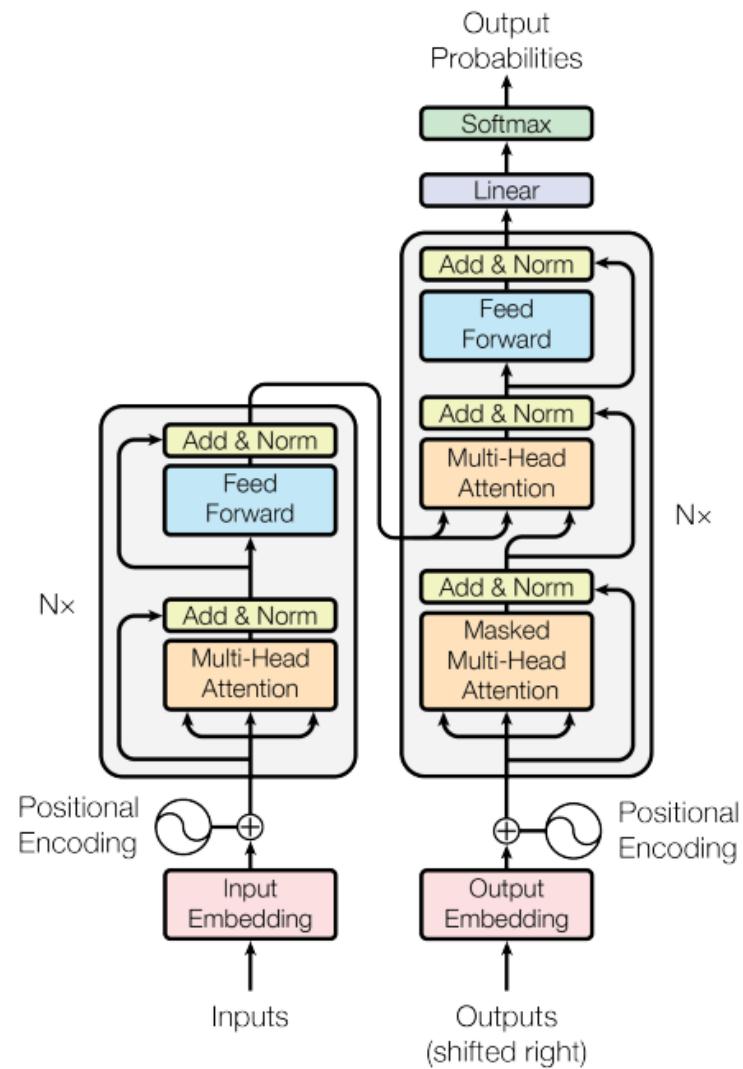
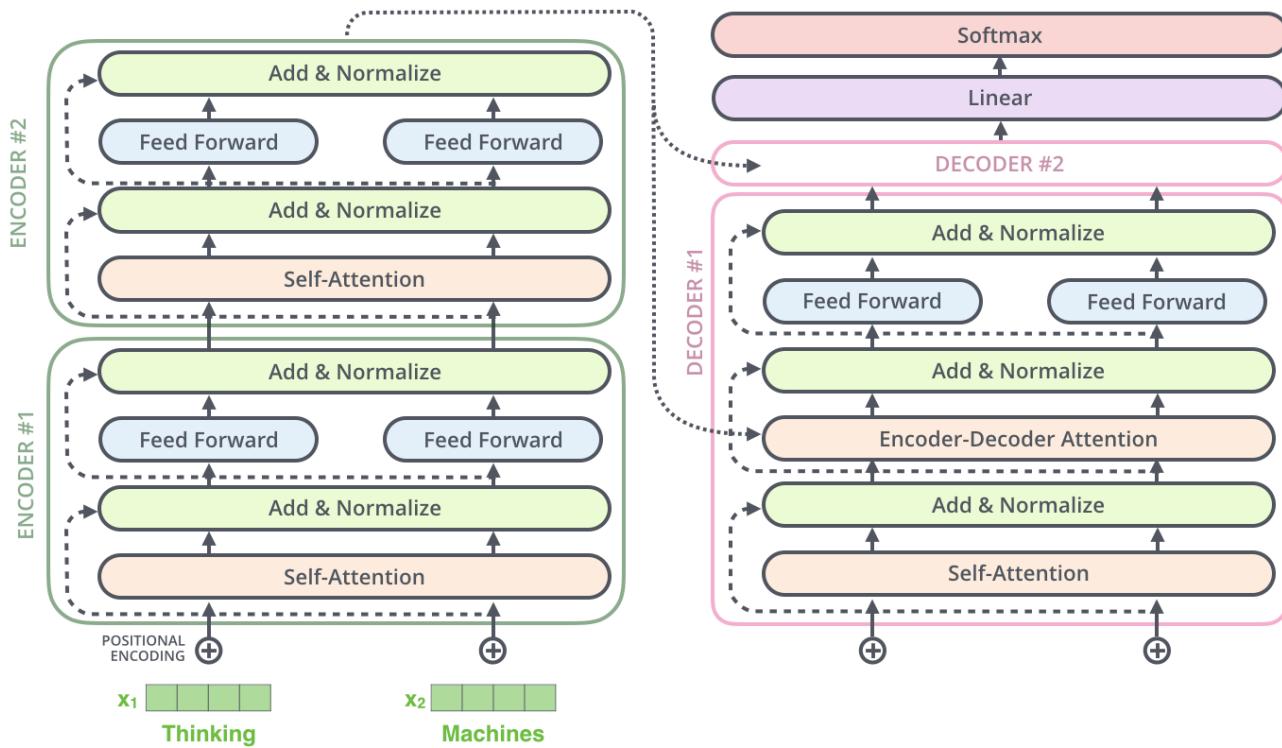
# Transformer

- Positional Encoding



# Transformer

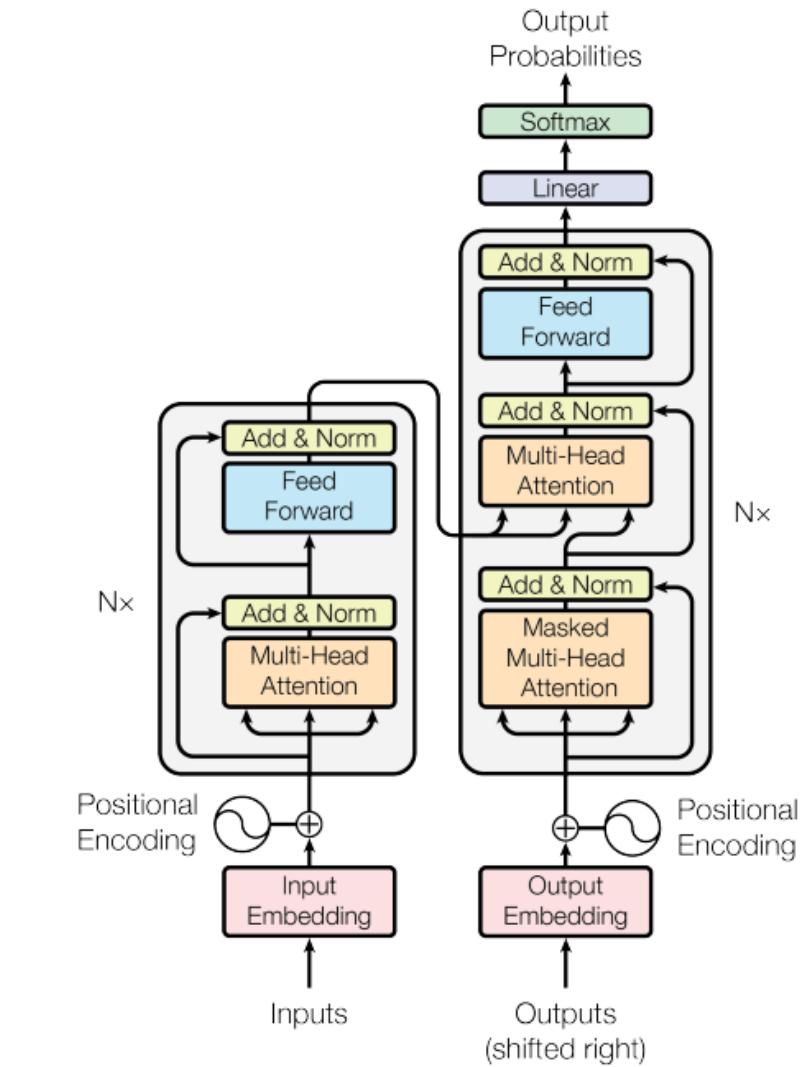
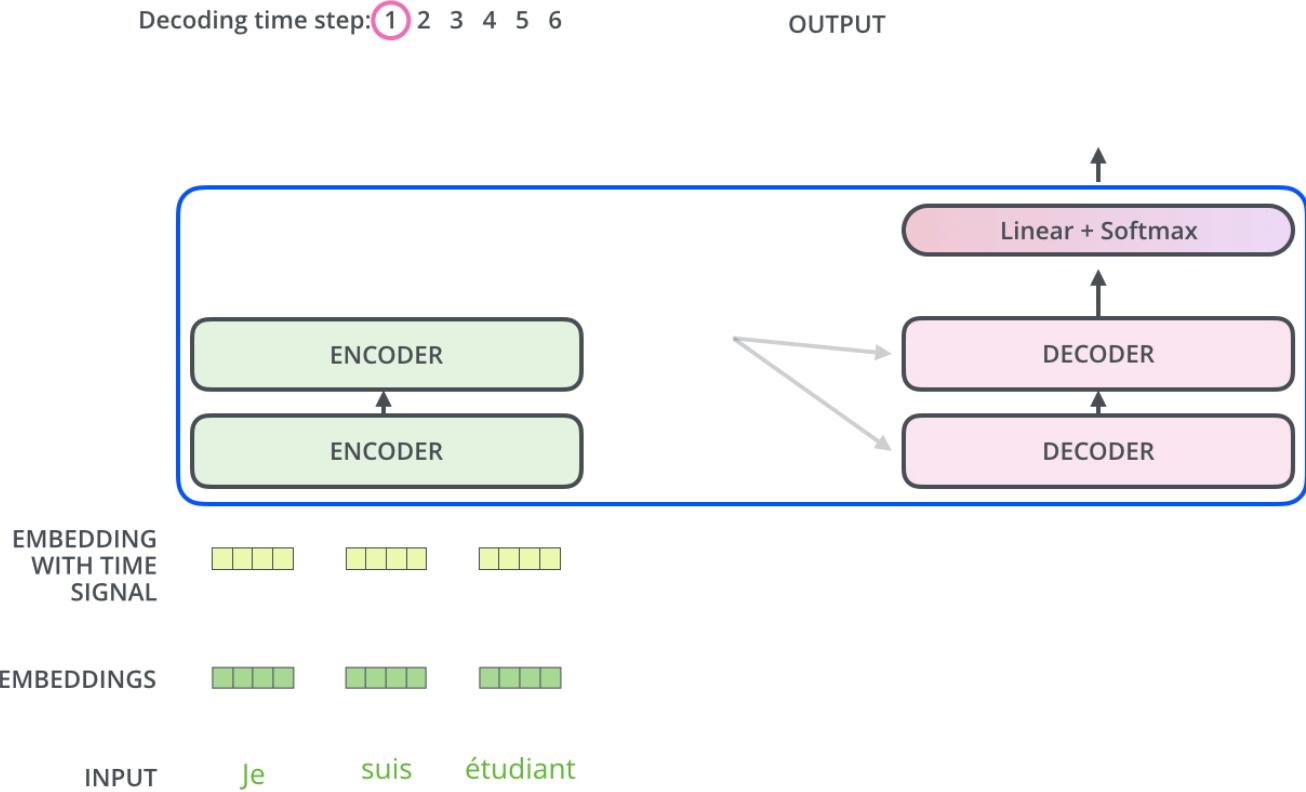
- Detail process in encoder & decoder



# Transformer

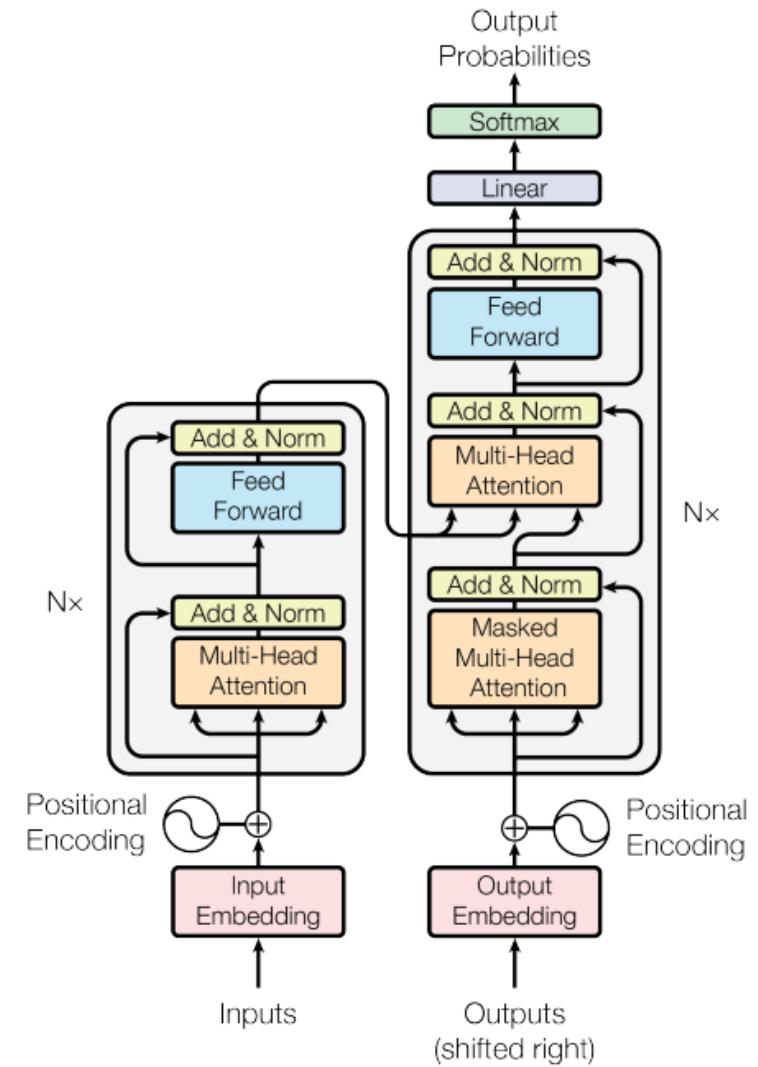
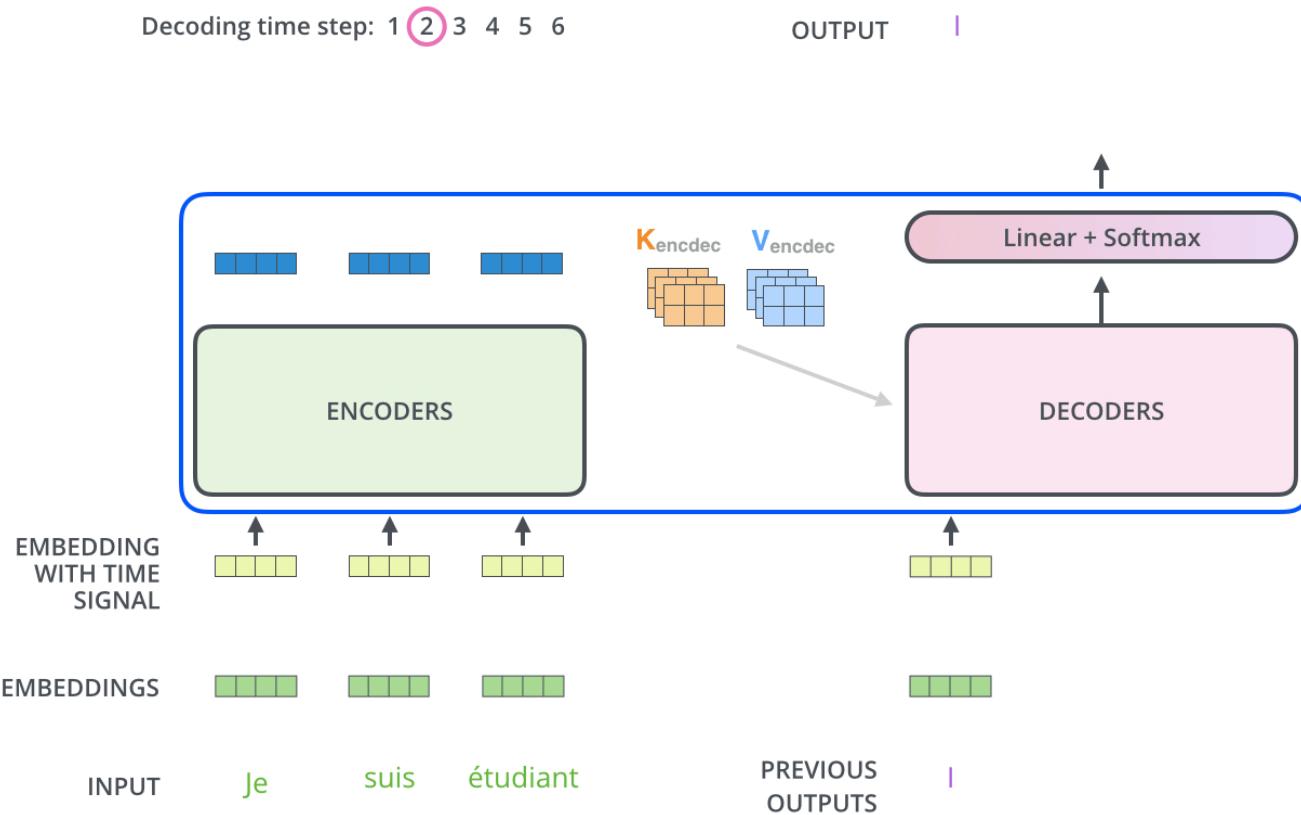
- Decoder

Decoding time step: 1 2 3 4 5 6



# Transformer

- Decoder



# BERT



- Quiz
  - Minjoo \_\_\_\_\_ kimchi
    - ① run
    - ② read
    - ③ eat



# BERT

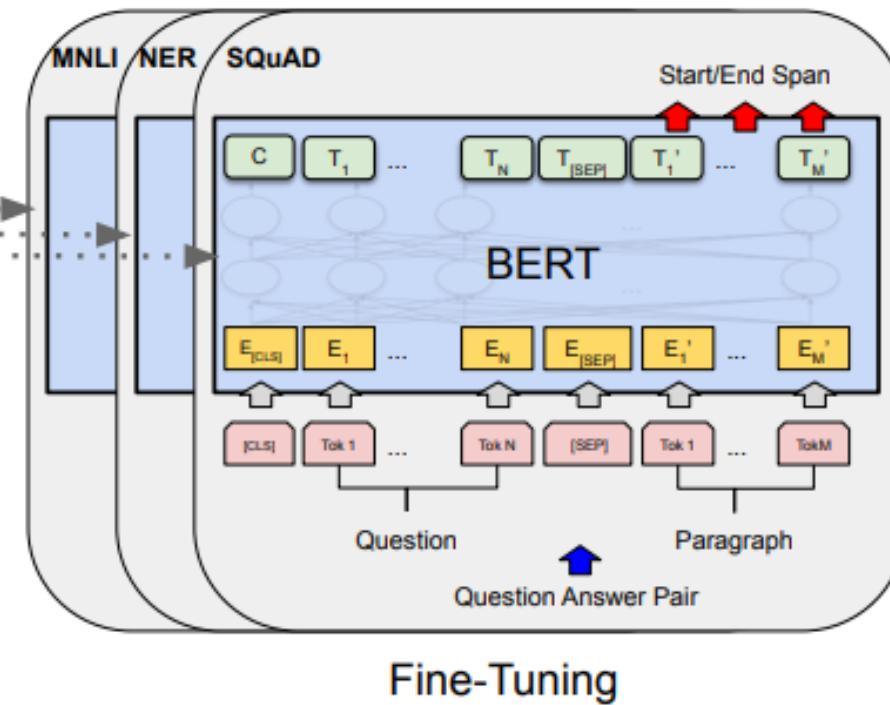
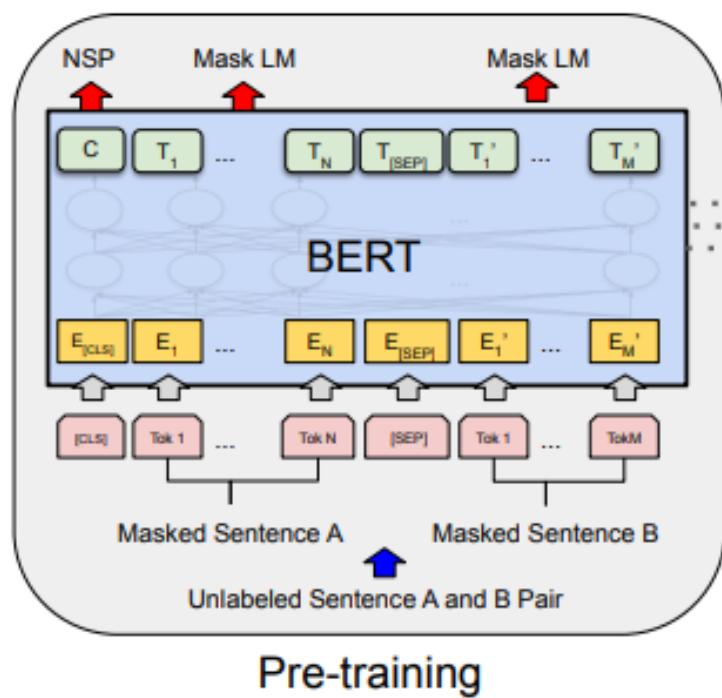


- Quiz
  - Minjoo \_\_\_\_\_ kimchi
  - ① run
  - ② read
  - ③ **eat**



# BERT

- Bidirectional Encoder Representations from Transformers



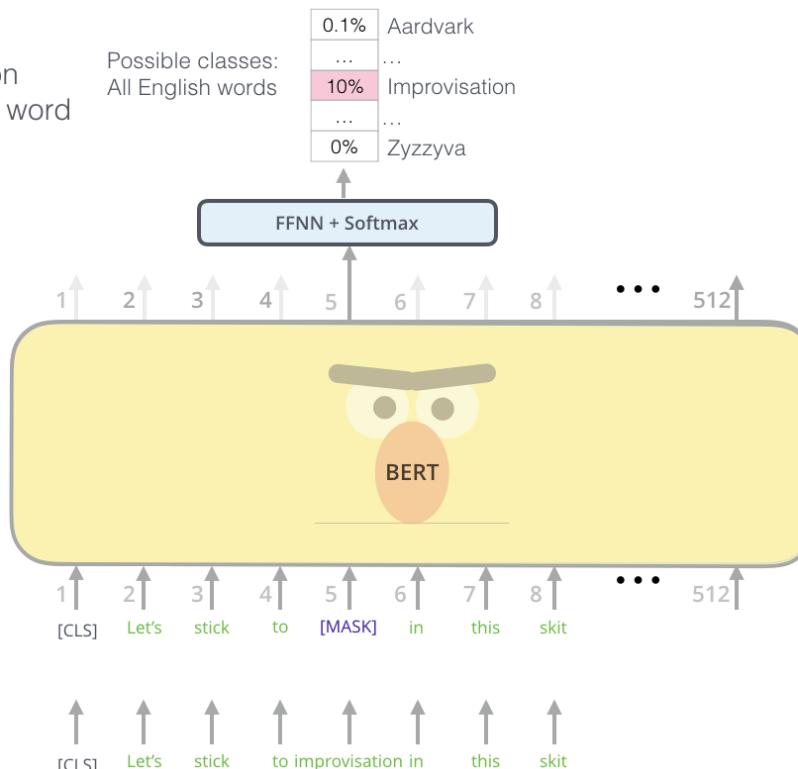
# BERT

- Masked Language Model

Use the output of the masked word's position to predict the masked word

Randomly mask 15% of tokens

Input



# BERT

- Pre-trained model

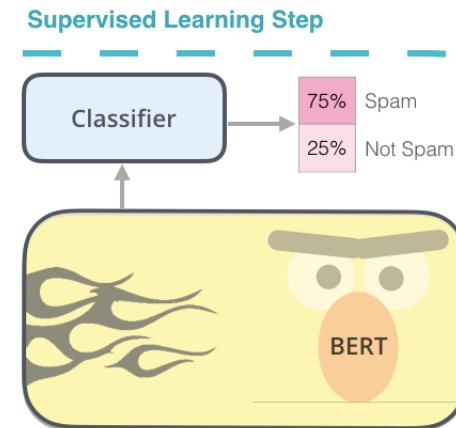
1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



Predict the masked word  
(language modeling)

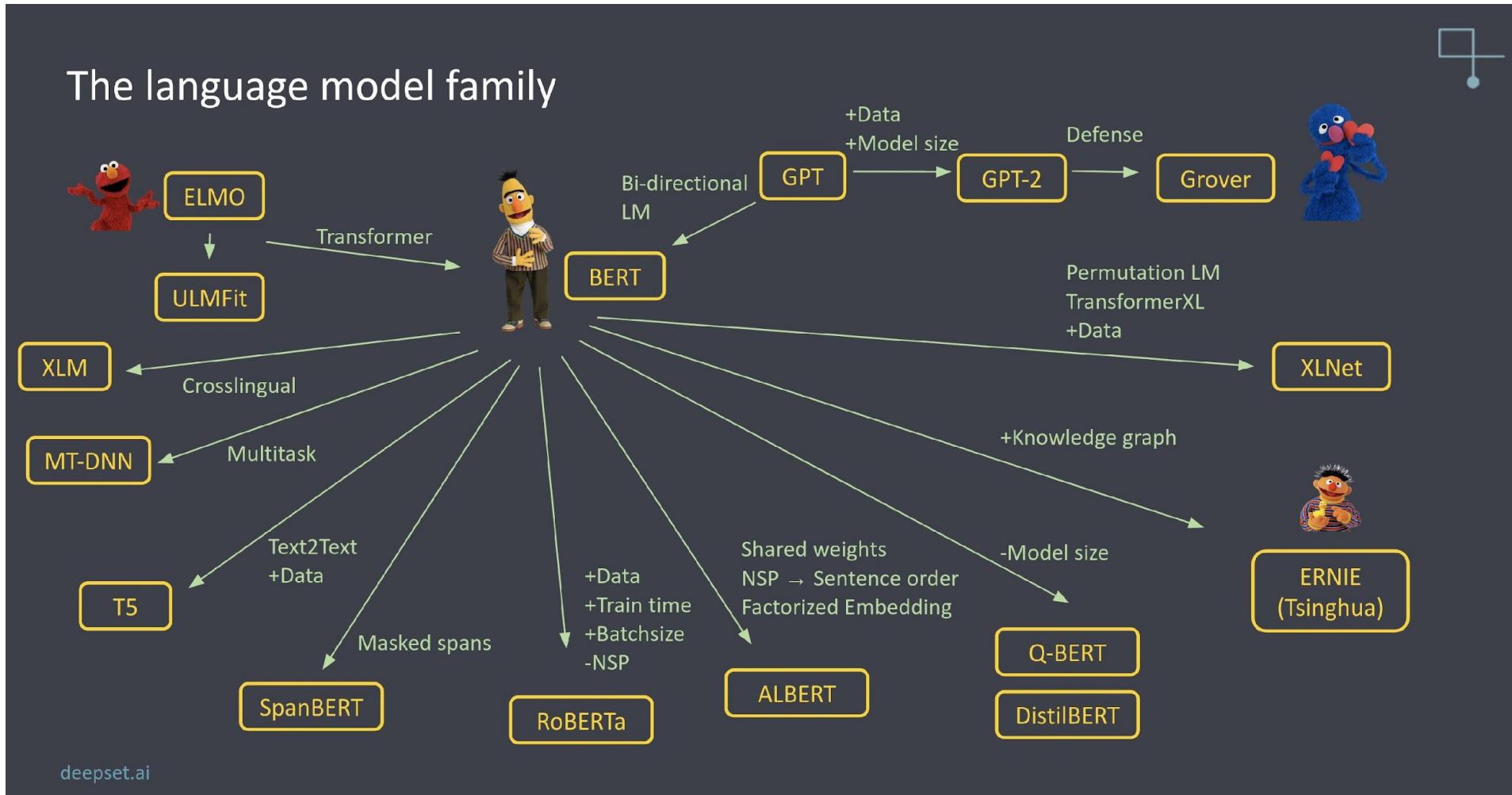
2 - **Supervised** training on a specific task with a labeled dataset.



Dataset:

Email message	Class
Buy these pills	Spam
Win cash prizes	Spam
Dear Mr. Atreides, please find attached...	Not Spam

# Family of BERT



# Machine Translation

# Machine Translation?

- Process when a computer software translates text from one language to another without human involvement



# History



RBMT

- Rule Based Machine Translation
- = PBMT (Phrase Based Machine Translation)



SMT

- Statistical Machine Translation



NMT

- Neural Machine Translation (by AI)
- Trained with sequence

# RBMT : Rule Based Machine Translation



## 규칙 기반 자동번역의 동작 예 [영어=>한국어]

I'm looking for the cosmetics section.

i/PRP be/VBP look/VBG for/IN  
the/DT cosmetics/NNS  
section/NN

look\_for/V[mood=decl,tense=pres,aspect=prog]  
subj I/R  
dobj cosmetics\_section/N  
punc ./S

- (R1:subj) look\_for/V2 (N3:dobj)[sem=활동영사] => (R1:subj) (N3:dobj) 기대하/V2
- (R1:subj) look\_for/V2 (N3:dobj) => (R1:subj) (N3:dobj) 찾/V2
- i/R[sem=사람] => 니/R
- cosmetics\_section/N[sem=장소] => 화장품\_코너/N

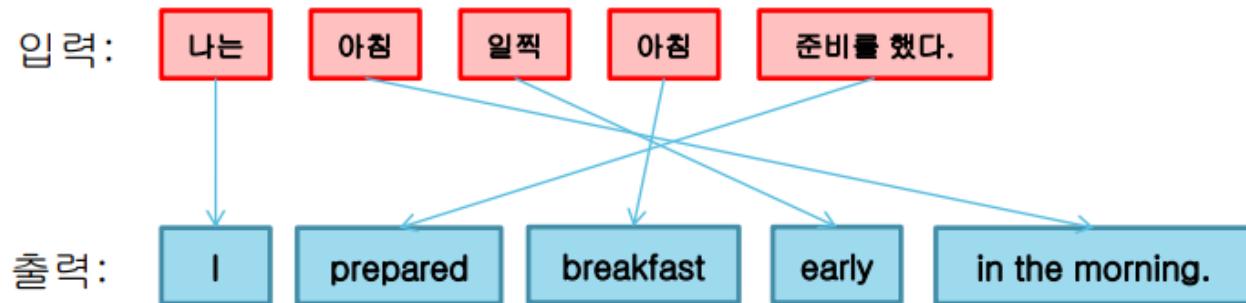
찾/V[mood=decl,tense=pres,aspect=prog]  
subj 니/R  
dobj 화장품\_코너/N  
punc ./S

니/NP는/FX 화장품\_코너/NN  
를/FO 찾/VB 고있/VX  
습니다/EE ./SF

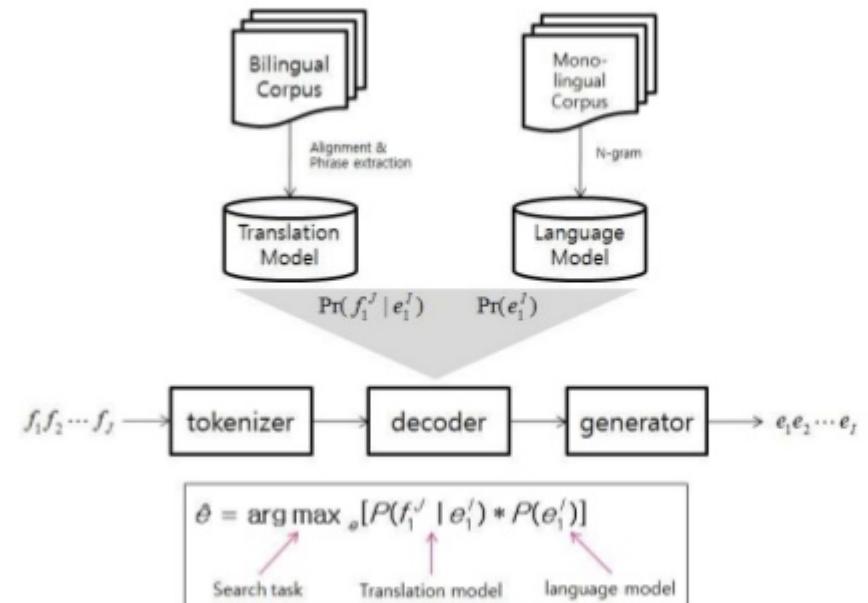
나는 화장품 코너를 찾고 있습니다.

# SMT : Statistical Machine Translation

- 작동 방식

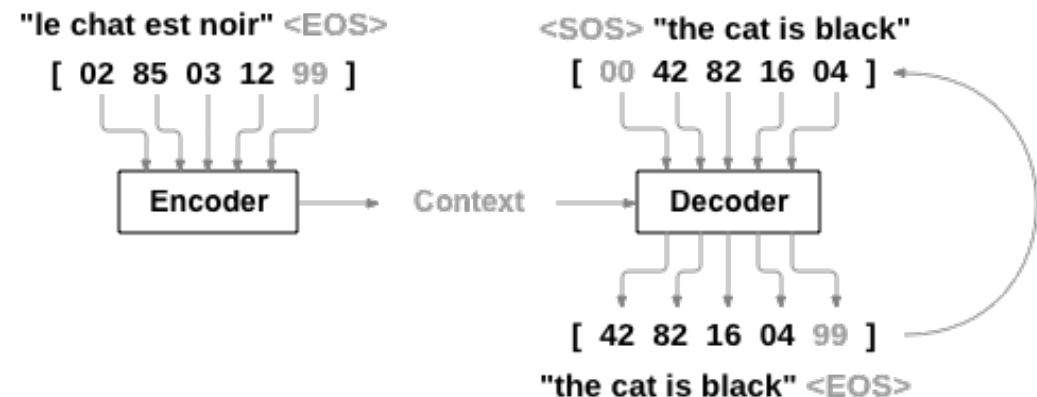


- 입력 문장의 일부분(Word 혹은 Phrase)를 보고 번역 수행  
Translate using word or phrase of input sentence
- 번역 사전을 사용하여, 확률에 기반한 번역을 수행  
Statistical translation using a dictionary
- 어순이 같은 언어 (한국어/일본어, 영어/프랑스어)
  - 괜찮은 성능.  
Similar order language (Korean/Japanise, English/French) : **Not Bad**
- 어순이 다른 언어 (한국어/영어)
  - 성능 부족.  
Different order language (Korean/English) : **BAD**



# Neural Machine Translation (NMT)

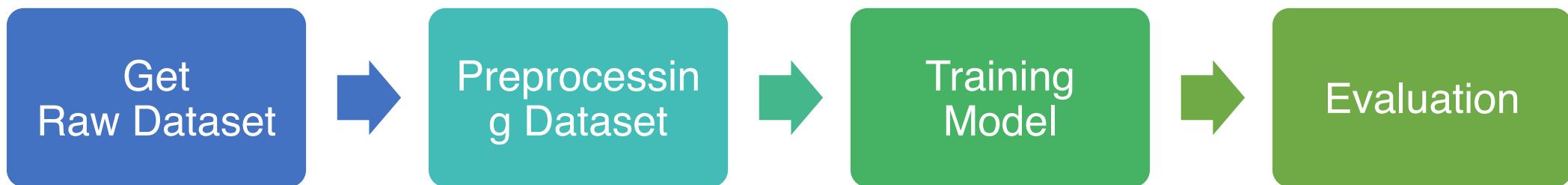
- Train translation model with language sentence pair
  - Dataset feature : [Source sentence] - [Target sentence]
    - Don't need to build phrase analyzer / dictionary
    - Translate sentence by sentence -> Translate in context is possible



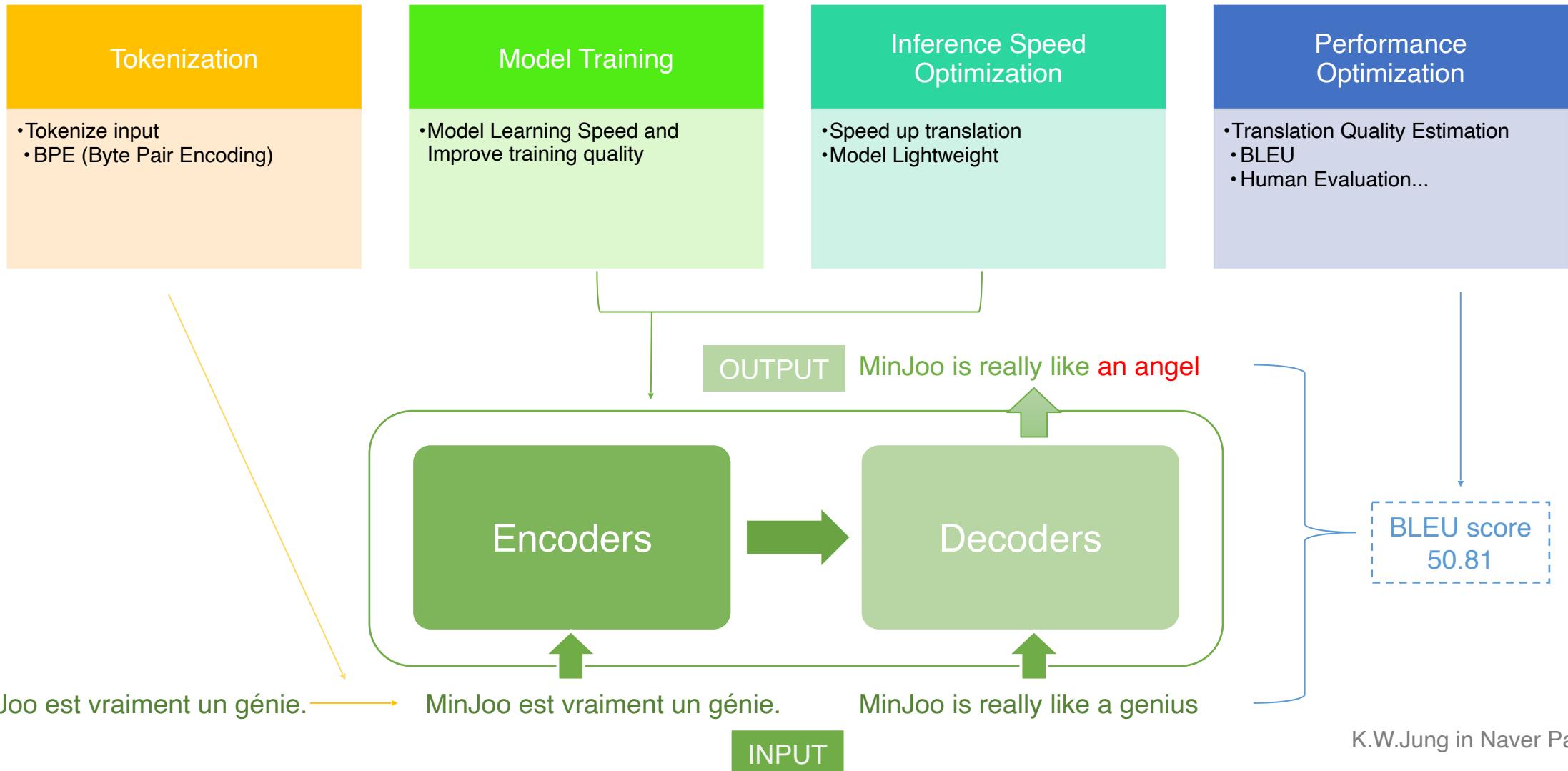
[https://tutorials.pytorch.kr/\\_images/seq2seq.png](https://tutorials.pytorch.kr/_images/seq2seq.png)

# Process to Develop NMT

- NMT (Neural Machine Translation)



# 4 Elements of Machine Translation



# (1/4) Tokenization

- Tokenize input sentence
  - 내용이 정말 거지같다. → 내용/NNG 이/JKS 정말/MAG 거지/NNG 같/VA 다/EF ./SF
  - Tokenizer example : (KOR) Mecab, Okt / (EN) Moses / (ZH) Jieba / ...
- BPE (Byte Pair Encoding)
  - Subword segmentation algorithm
  - OOV (Out Of Vocabulary) problem solved.
    - lowest -> low + est
  - <https://github.com/rsennrich/subword-nmt>

# (2/4) Model Training

- Model
  - Fast model training strategy / High-quality strategy

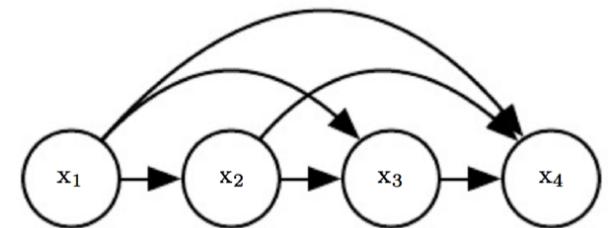
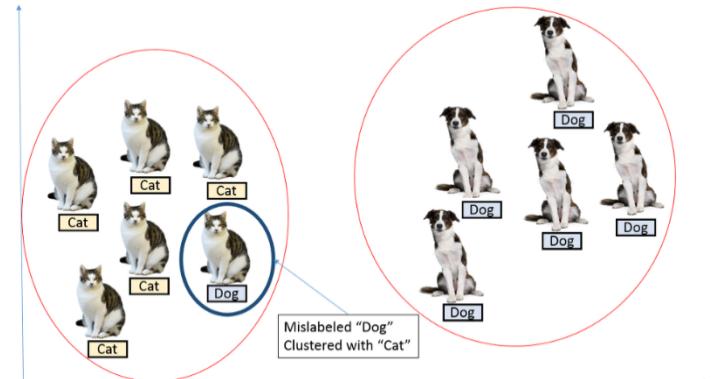
- Dataset
  - Data Augmentation
    - Back-Translation



<http://dsba.korea.ac.kr/seminar/?mod=document&uid=1328>

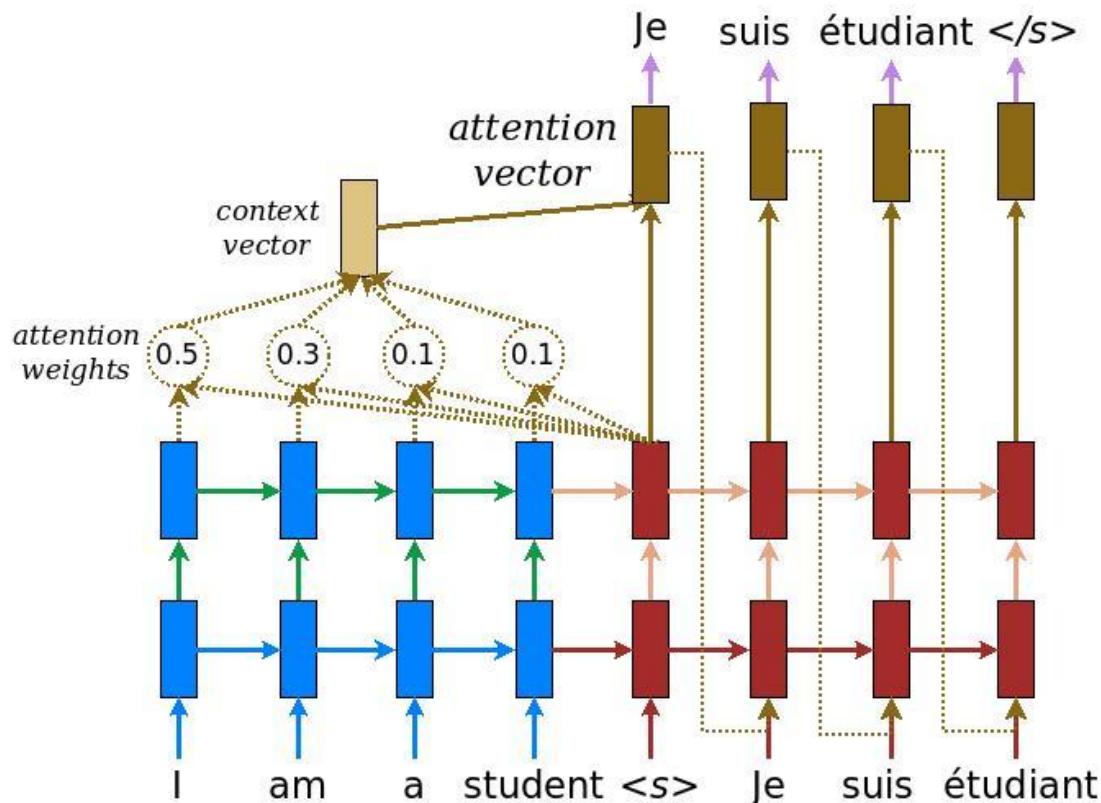
# (2/4) Model Training

- Optimization
  - Label Smoothing
  - Curriculum learning
- Modeling
  - Non-Autoregressive
  - LSTM (Seq2Seq + Attention) → Transformer (Self-Attention)



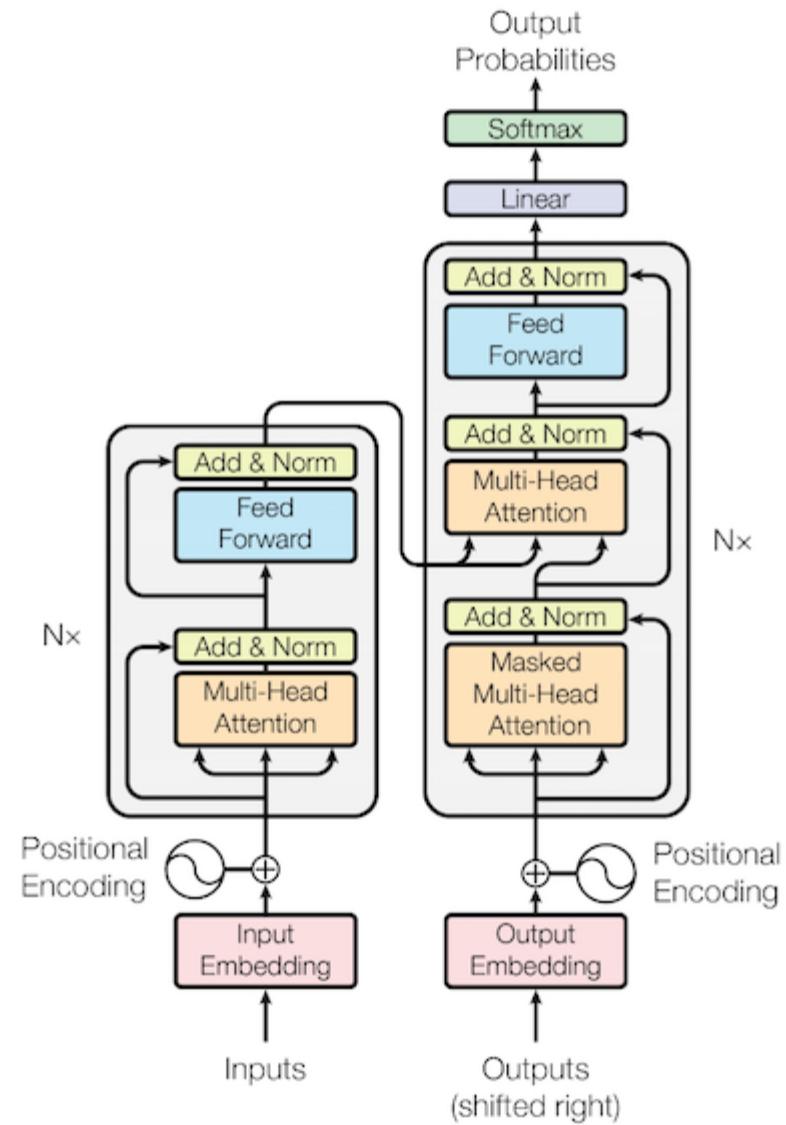
<https://3months.tistory.com/465>  
<https://ratsgo.github.io/generative%20model/2018/01/31/AR/>

# (2/4) Model Training



Seq2Seq+ Attention

[https://www.tensorflow.org/text/tutorials/nmt\\_with\\_attention](https://www.tensorflow.org/text/tutorials/nmt_with_attention)



Transformer (Self-Attention)

<https://paul-hyun.github.io/transformer-03/>

## (3/4) Performance Optimization

- Measure performance of translation models
- Evaluate translation quality
  - Compare translation results with reference
- **BLEU**, ROUGE, RIBES...
  - Research about evaluation method better than BLEU
- More fast, reliable Human Evaluation

# (3/4) Performance Optimization

- BLEU score (Bilingual Evaluation Understudy Score)
  - Evaluation method by N-gram

$$p_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n\text{-gram} \in C} Count_{clip}(n\text{-gram})}{\sum_{C' \in \{Candidates\}} \sum_{n\text{-gram}' \in C'} Count(n\text{-gram}')}$$

```
1 # 1 word different
2 from nltk.translate.bleu_score import sentence_bleu
3 reference = [['the', 'quick', 'brown', 'fox', 'jumped', 'over', 'the', 'lazy', 'dog']]
4 candidate = ['the', 'fast', 'brown', 'fox', 'jumped', 'over', 'the', 'lazy', 'dog']
5 score = sentence_bleu(reference, candidate)
6 print("{:.2f}".format(score))
```

0.75

- Moses multi-bleu perl
  - <https://github.com/moses-smt/mosesdecoder/blob/master/scripts/generic/multi-bleu.perl>

```
1 # 2 word different
2 from nltk.translate.bleu_score import sentence_bleu
3 reference = [['the', 'quick', 'brown', 'fox', 'jumped', 'over', 'the', 'lazy', 'dog']]
4 candidate = ['the', 'fast', 'brown', 'fox', 'jumped', 'over', 'the', 'sleepy', 'dog']
5 score = sentence_bleu(reference, candidate)
6 print("{:.2f}".format(score))
```

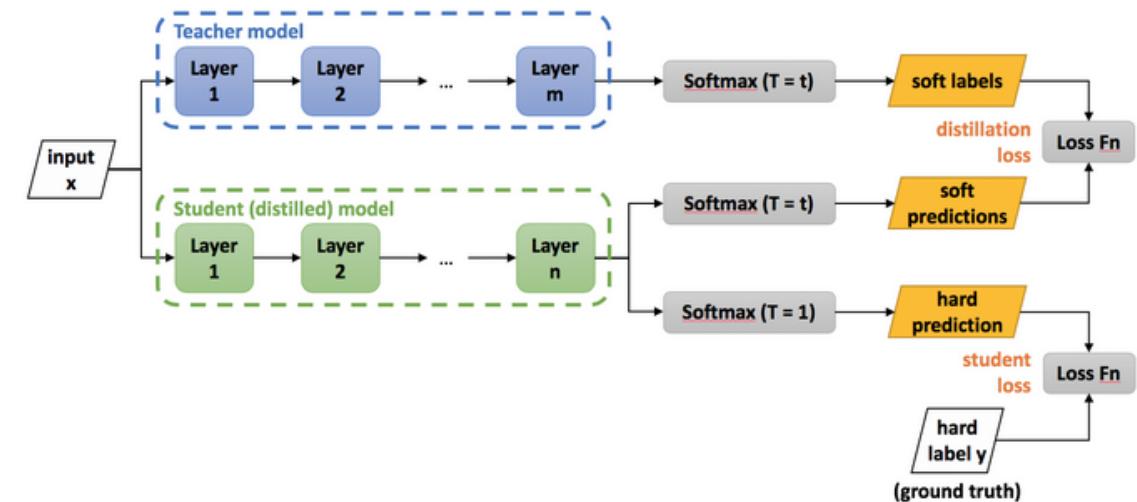
0.49

# (4/4) Inference Speed Optimization

- Improve translation speed
  - GPU inference, CPU inference
- Lighten Model, On-Device NMT
  - Decrease Encoder or Decoder Layer of NMT
    - Decrease model size, increase speed
  - Decrease model capacity (Algorithm) or change model structure
  - Decrease or share parameter
  - \* Often a slight decrease in translation performance (trade-off)

# (4/4) Inference Speed Optimization

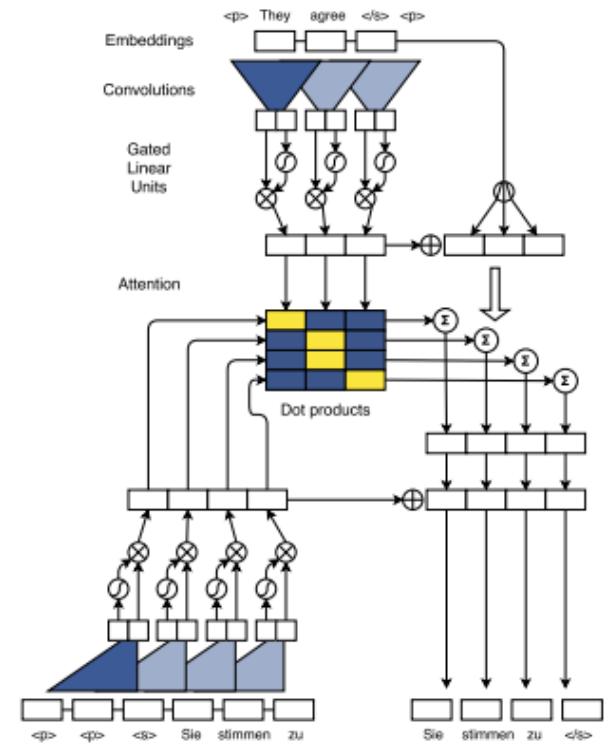
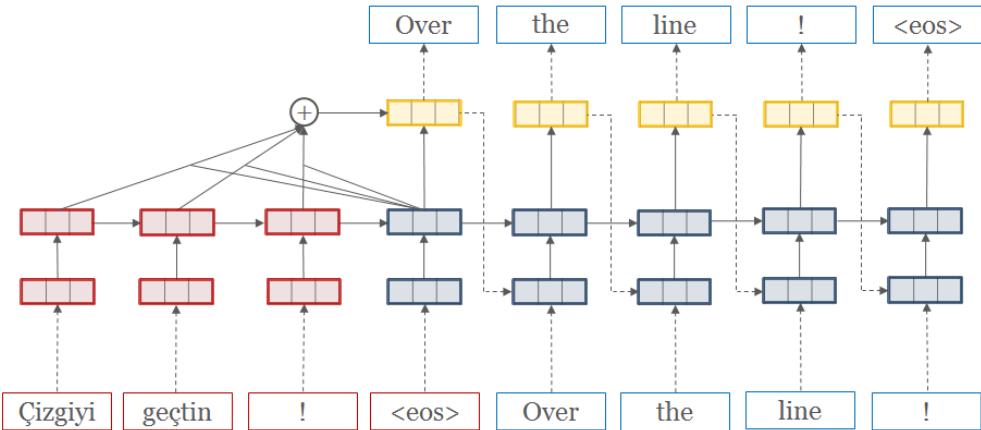
- (Additional) Knowledge Distillation
  - Transfer pre-trained knowledge (softer softmax) from teacher (big model) to student (small capacity model)



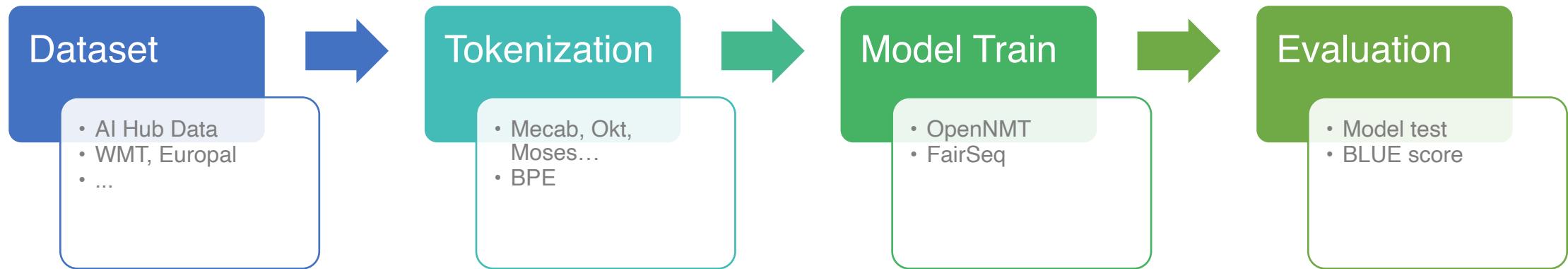
- <https://light-tree.tistory.com/196>
- [https://intellabs.github.io/distiller/knowledge\\_distillation.html](https://intellabs.github.io/distiller/knowledge_distillation.html)
- <http://dmqm.korea.ac.kr/activity/seminar/304>
- <https://blog.lunit.io/2018/03/22/distilling-the-knowledge-in-a-neural-network-nips-2014-workshop/>

# Open Source

- OpenNMT
  - by Havard NLP group & Systran
  - tensorflow, pytorch
- FairSeq (Facebook AI Research Sequence-to-Sequence Toolkit)
  - by Facebook AI Research
  - pytorch



# Example



Original (KO)	Reference (EN)	OpenNMT Translation Result
우리는 일주일에 5 일을 학교에 갑니다.	We go to school 5 days a week.	We go to school <b>for</b> 5 days a week.
다음 주말에는 함께 술을 마시기로 해요.	Let's drink together next weekend.	I <b>will</b> drink together next weekend.
난 당신이 행복해지기를 너무나 원해요.	I so want you to be happy.	I <b>really</b> want you to be happy.

# Q & A

<https://www.facebook.com/minjoo.choi.562/>

<https://github.com/Judy-Choi>

