



Transformer and its variants for NLP

Quan Thanh Tho
qttho@hcmut.edu.vn



Assoc. Prof. Quan Thanh Tho

Vice Dean

Faculty of Computer Science and Engineering

Ho Chi Minh City University of Technology (HCMUT)

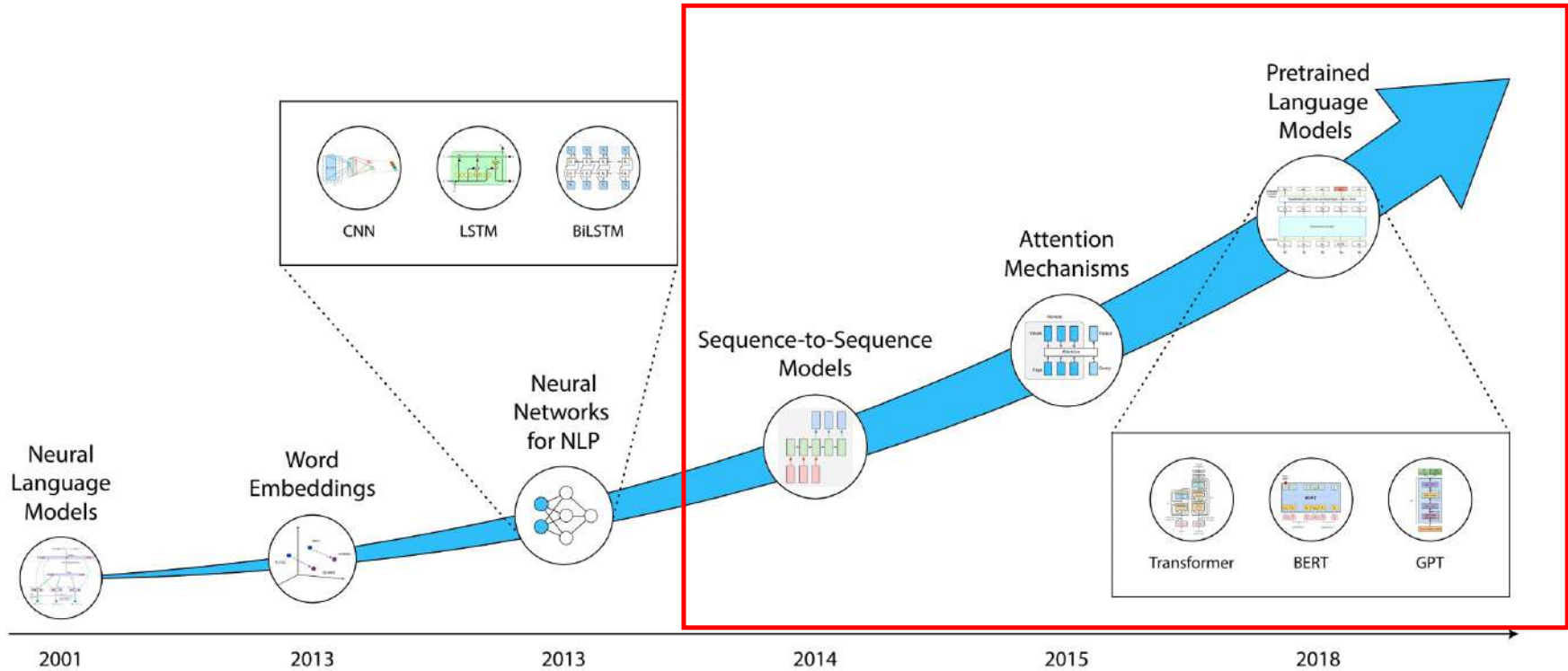
Vietnam National University - Ho Chi Minh City

qttho@hcmut.edu.vn

<http://www.cse.hcmut.edu.vn/qttho/>

- BEng, HCMUT, Vietnam, 1998
- PhD, NTU, Singapore, 2006
- Research Interests: Artificial Intelligence, Natural Language Processing, intelligent systems, formal methods

NLP Milestones



Quan Thanh Tho, "Modern Approaches in Natural Language Processing", *VNU Journal of Science: Computer Science and Communication Engineering*, 2022

Agenda

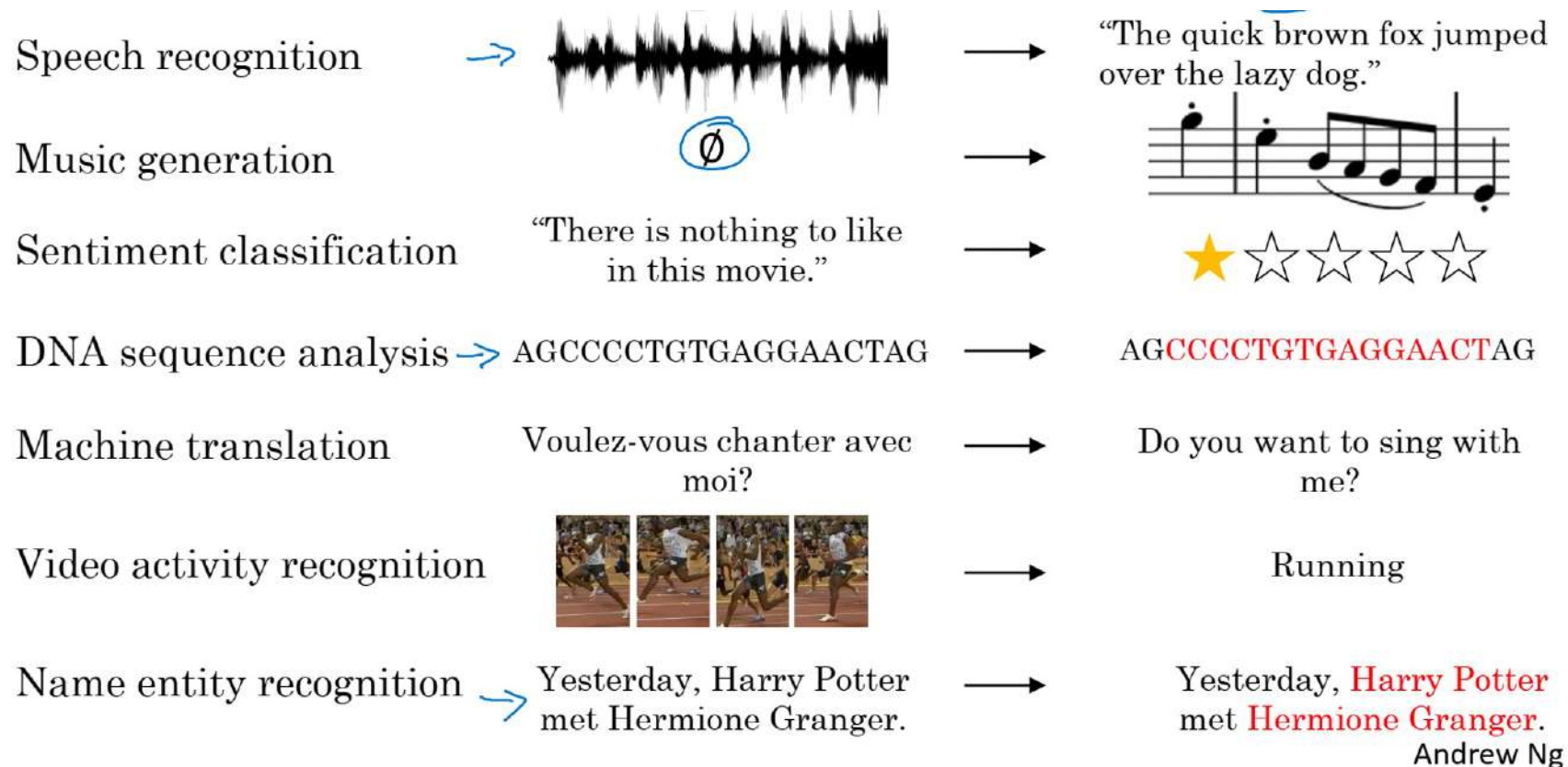
- Sequence data and sequence models
- Seq2Seq and attention
- Transformer model
- BERT and other variants
- Applications in NLP

Sequence data and sequence models

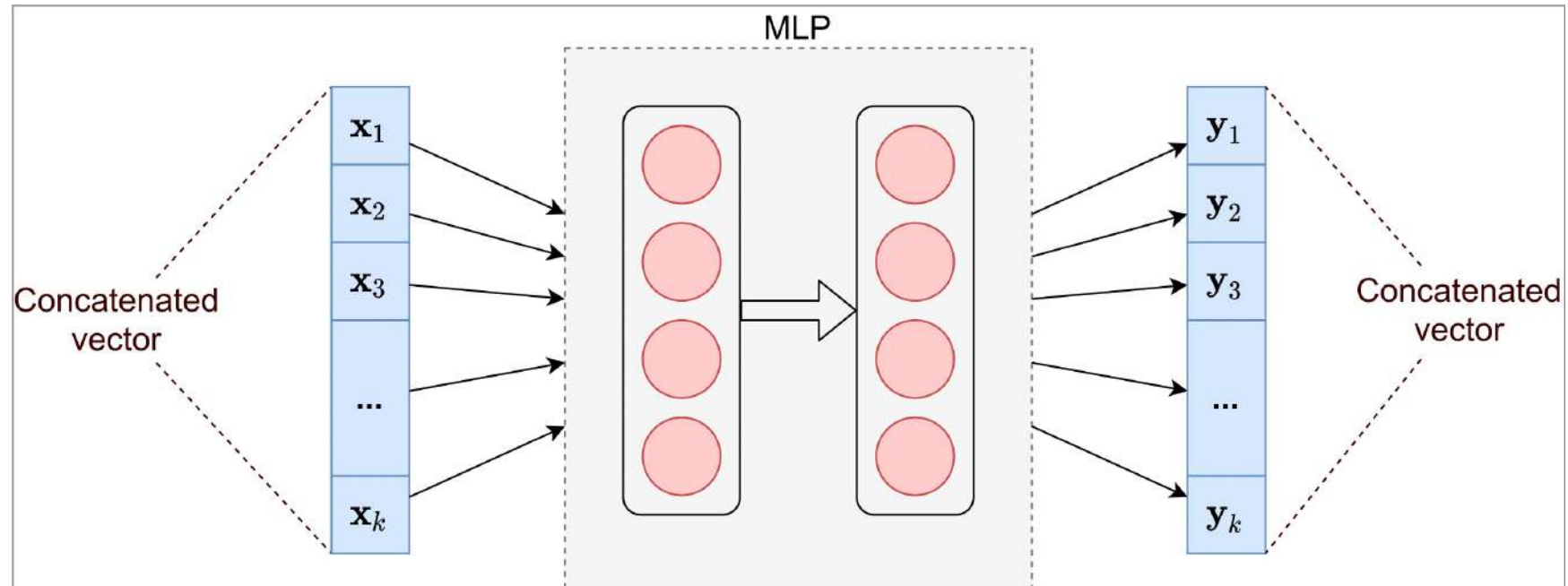
Sequence data

A series of data points whose points reliant on each other

- Length can be varied
- Positions matter

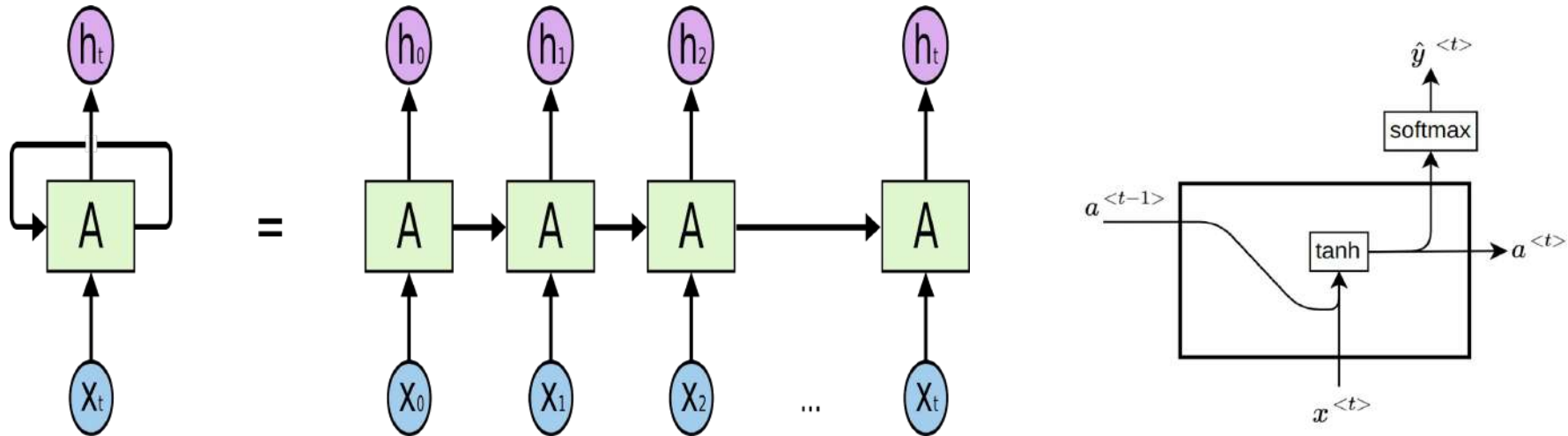


Problem of Standard Networks



- Inputs, outputs can be different lengths in different examples.
- Relations between positions are not well reflected

RNN comes as a rescue



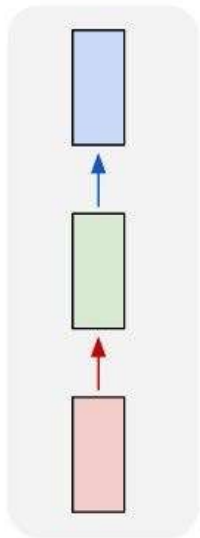
RNN: an architecture tailored for sequence data:

- 1) Doesn't depend on data length
- 2) Take advantage of past information

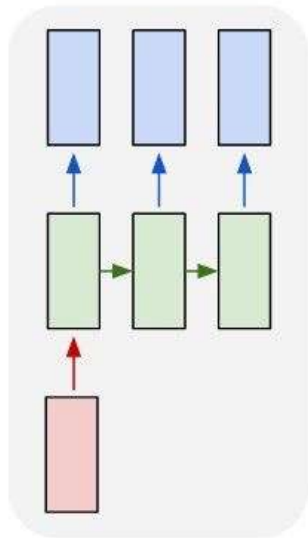
Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning internal representations by error propagation. In D. E. Rumelhart & J. L. McClelland (Eds.), *Parallel distributed processing: Explorations in the microstructure of cognition, Volume 1: Foundations* (pp. 318–362). MIT Press

RNN Revision

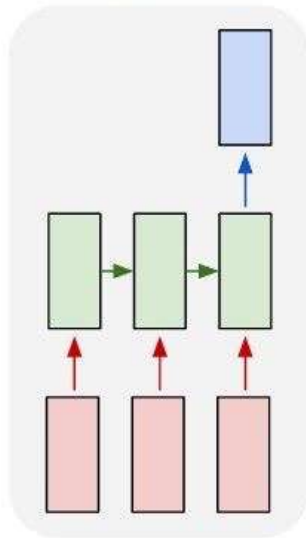
one to one



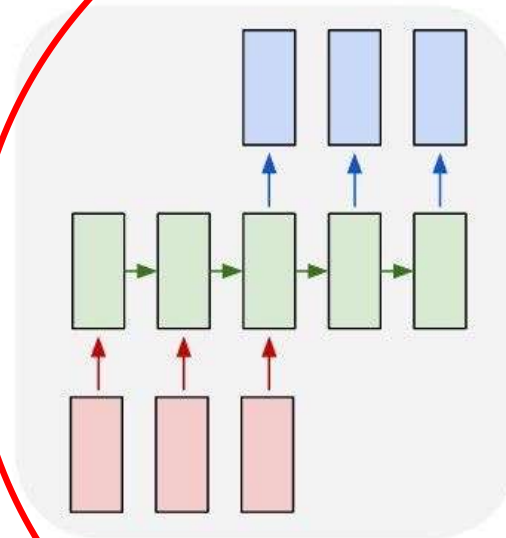
one to many



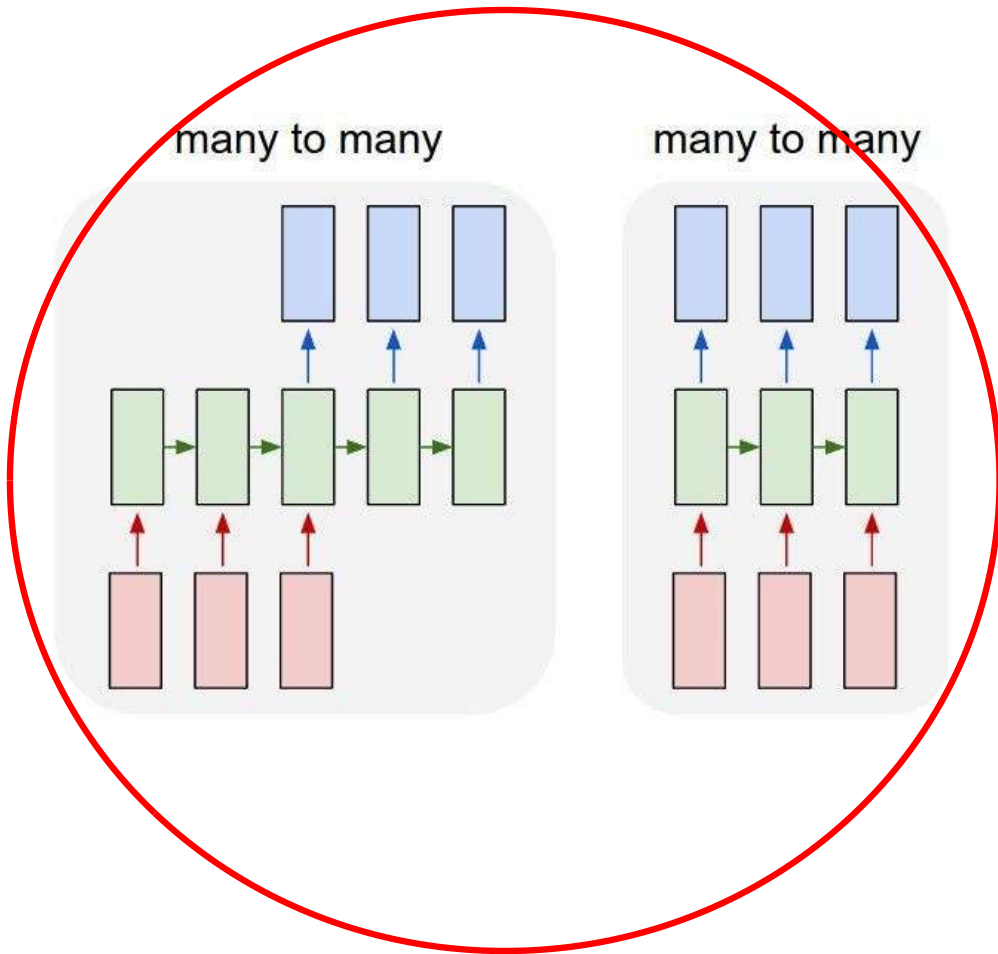
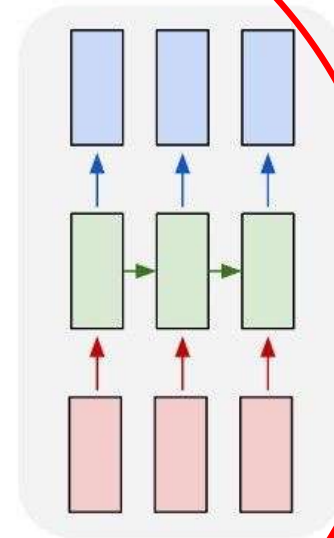
many to one



many to many



many to many



Seq2seq and Attention

Intuition

Take machine translation task as an example: human would first read some parts of the text and then start to do the translation

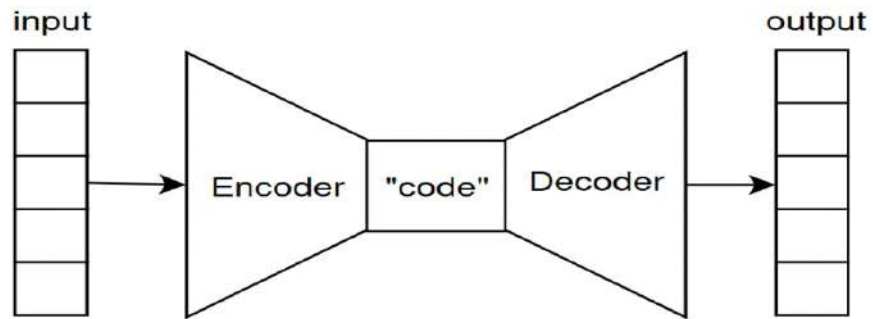
The cat likes to eat pizza



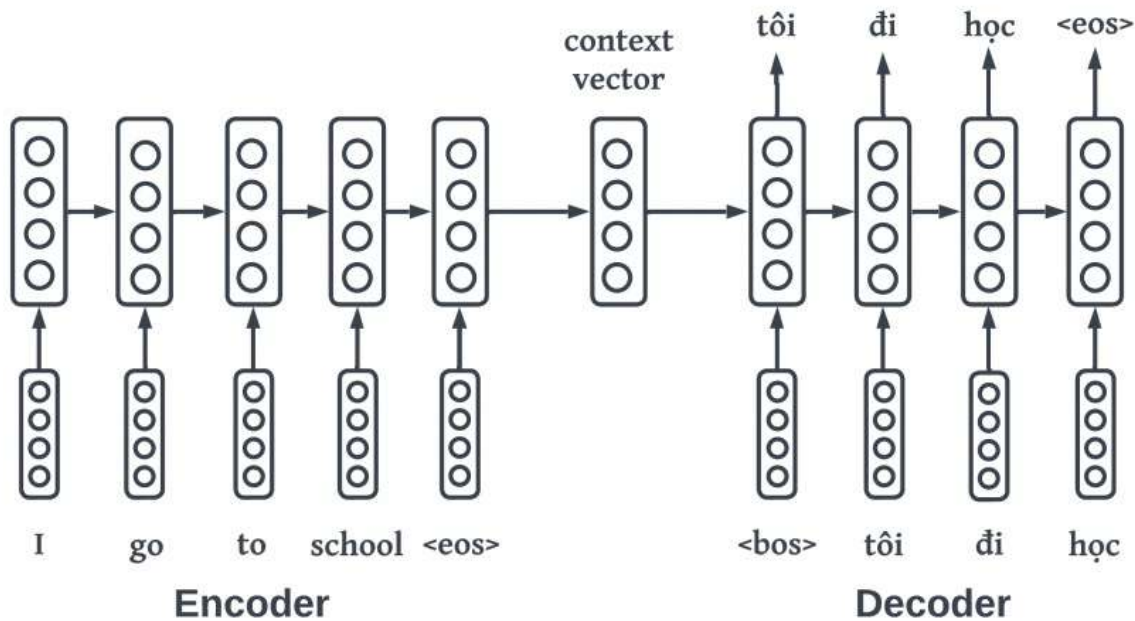
el gato le gusta comer pizza



Seq2Seq architecture



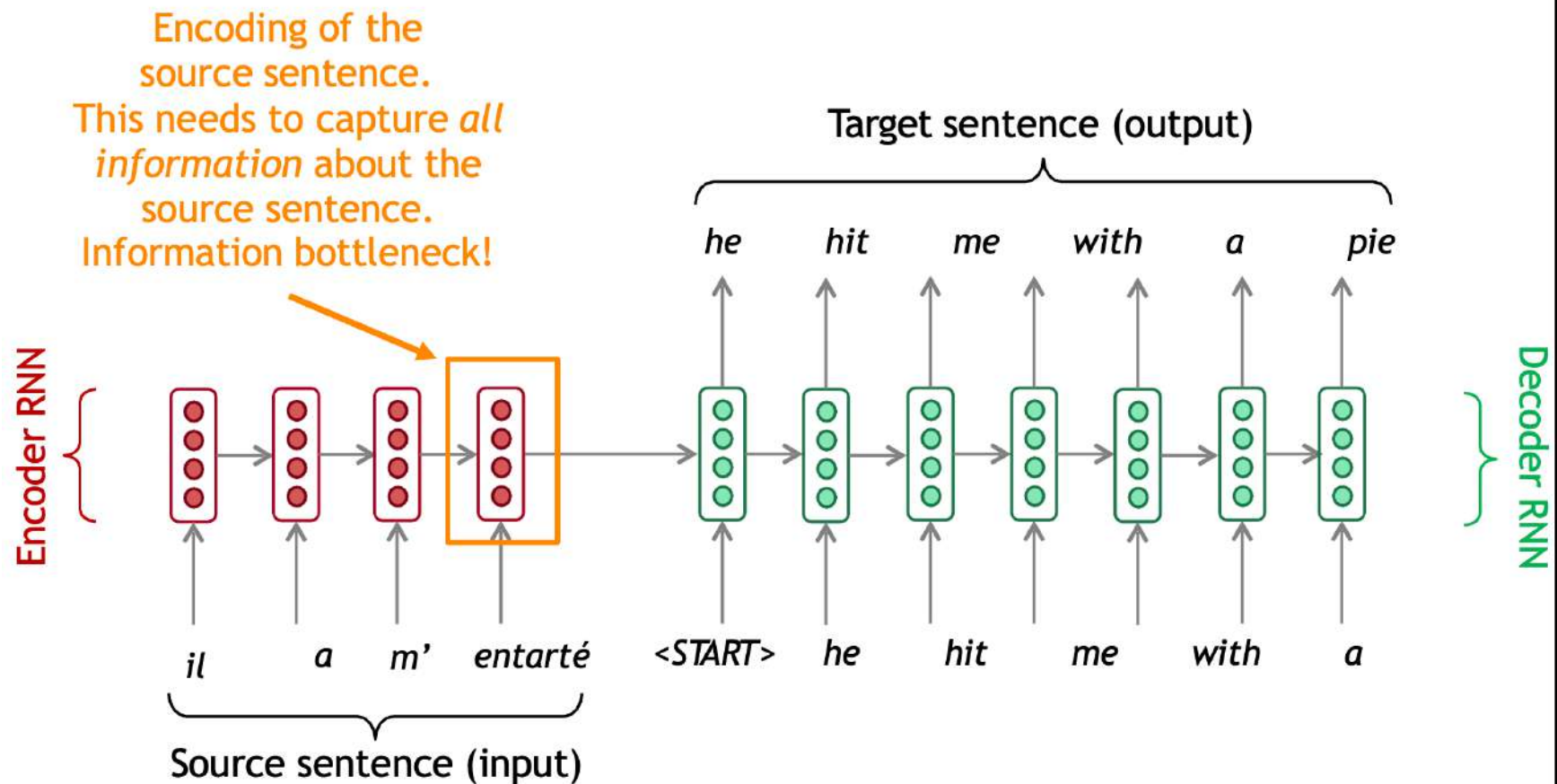
NLP researchers also employ that **idea** into designing a structure dubbed as Sequence-to-Sequence (Seq2Seq), which extends AutoEncoder architecture



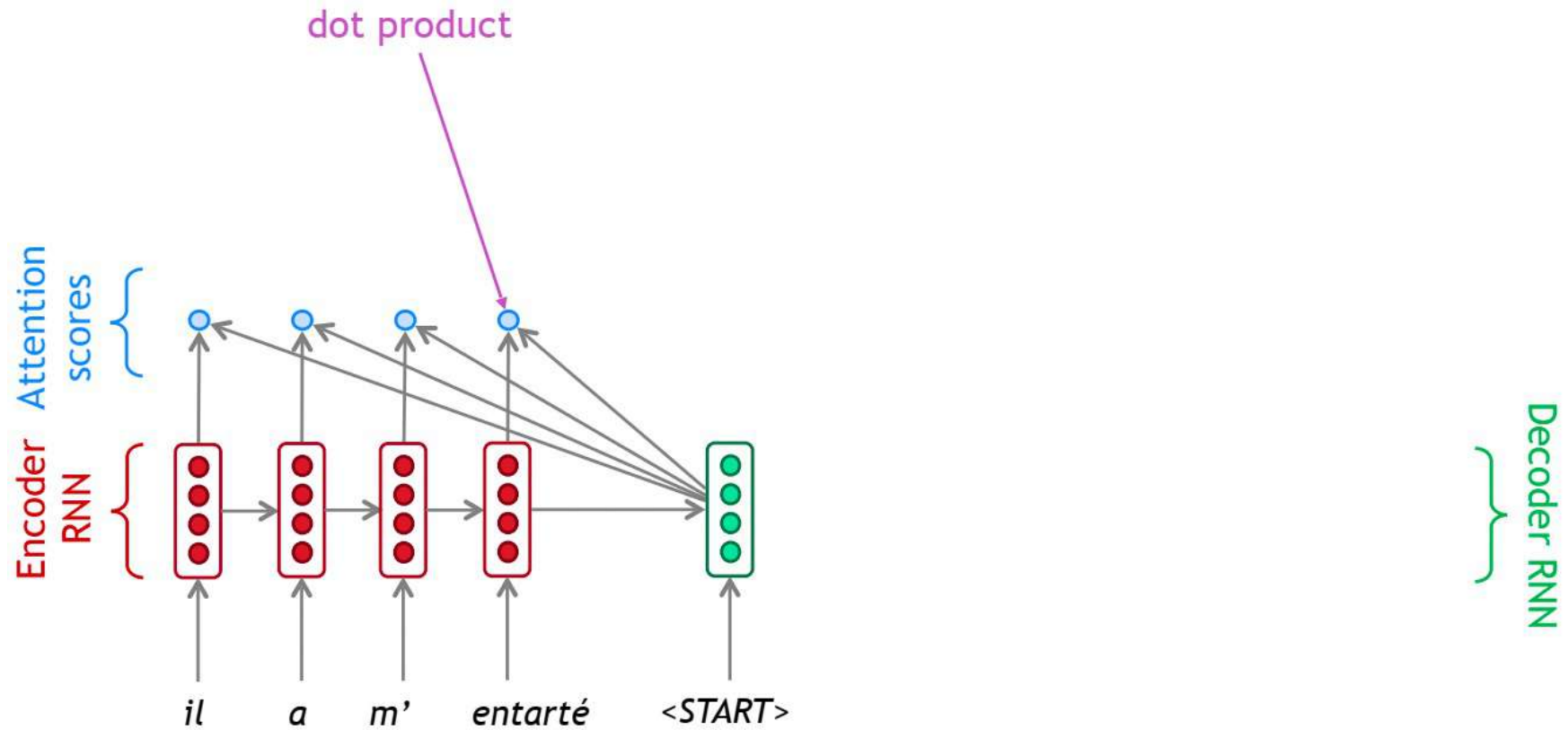
Kiros, R., Zhu, Y., Salakhutdinov, R. R., Zemel, R., Urtasun, R., Torralba, A., & Fidler, S. (2015). Skip-thought vectors. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information processing systems*. Curran Associates, Inc

Sutskever, I., Vinyals, O., & Le, Q. V. (2014). Sequence to sequence learning with neural networks. *Advances in neural information processing systems*, 27

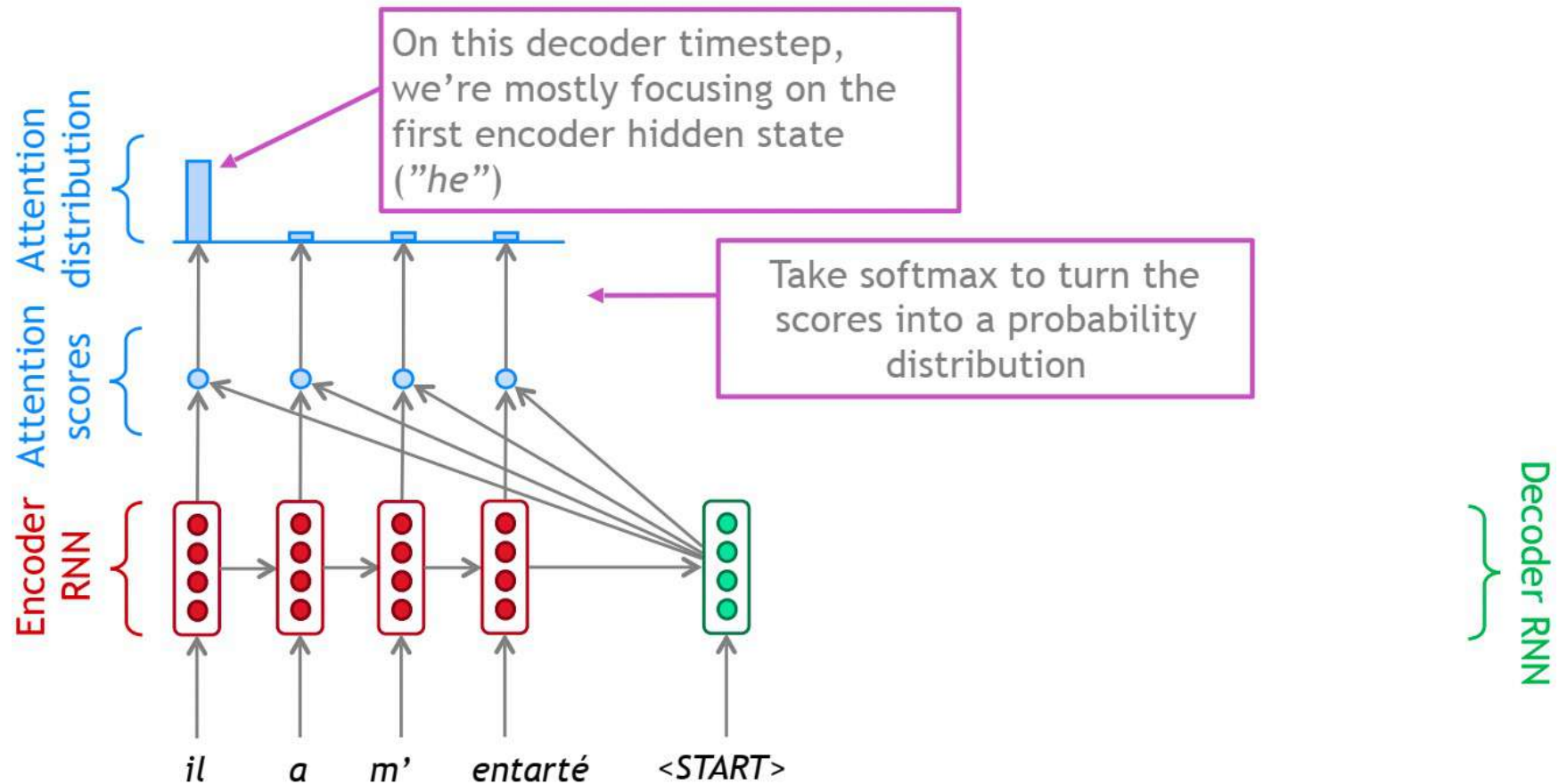
Seq2Seq: The bottle neck problem



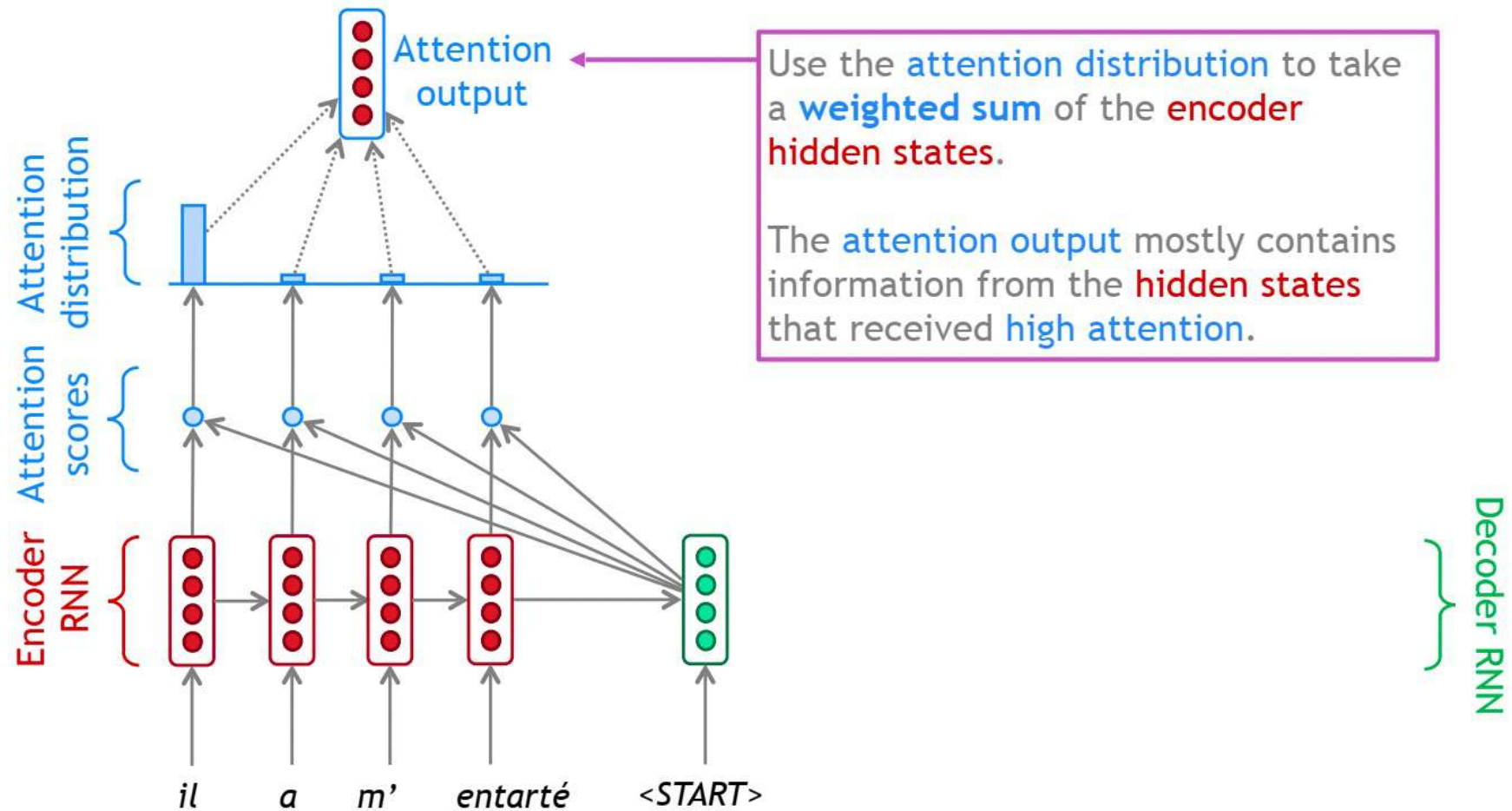
Seq2Seq with attention



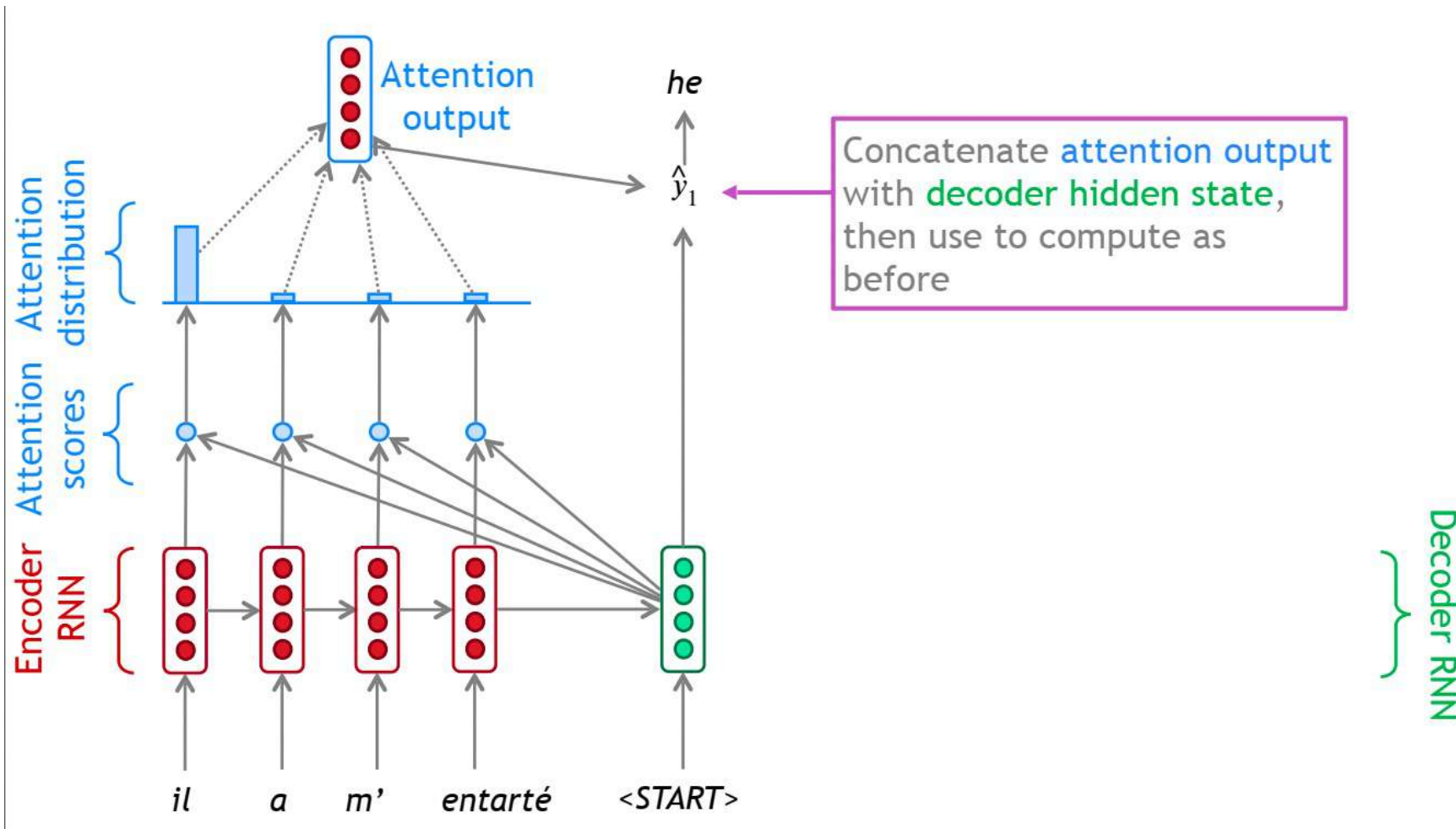
Seq2Seq with attention



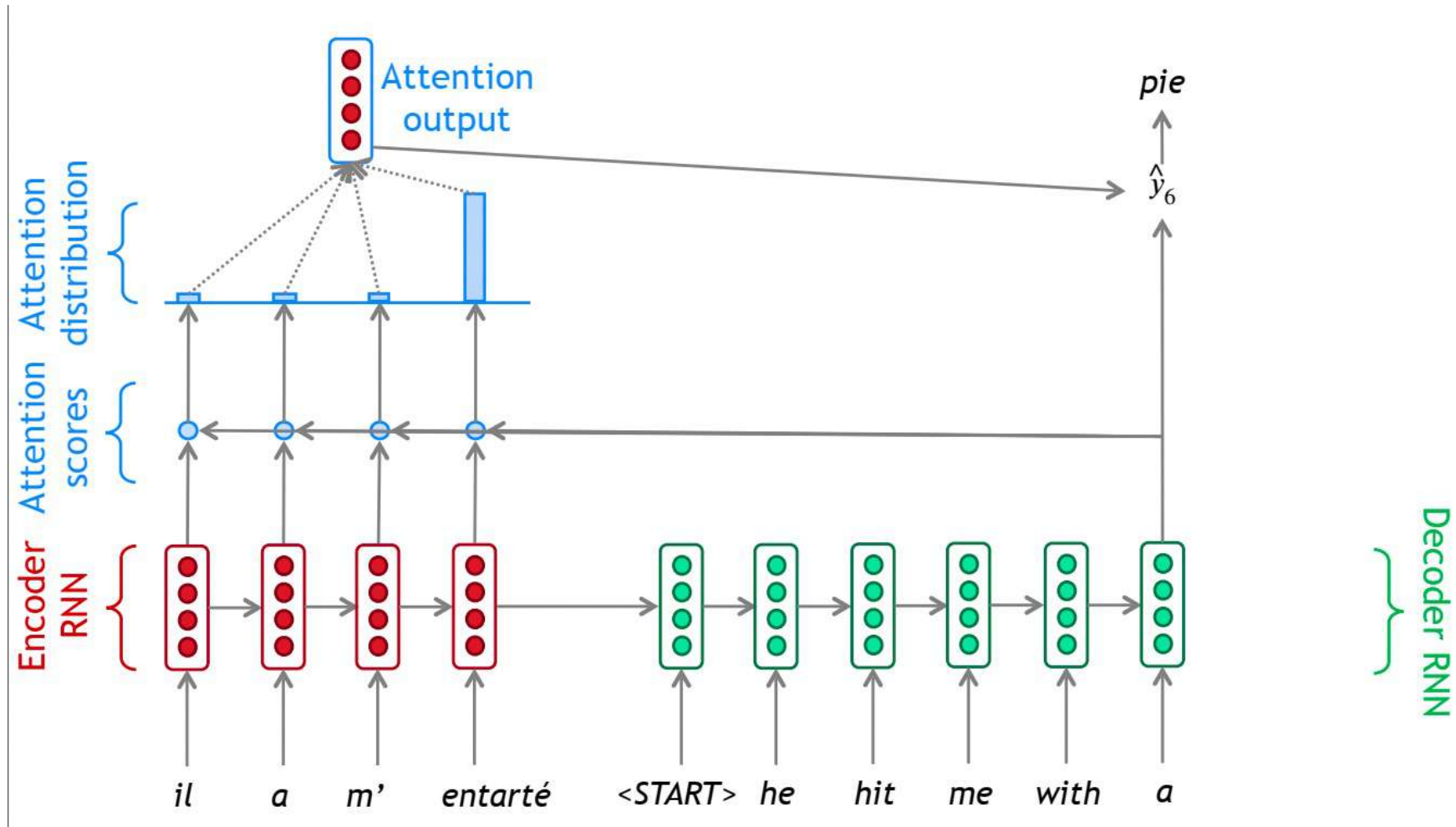
Seq2Seq with attention



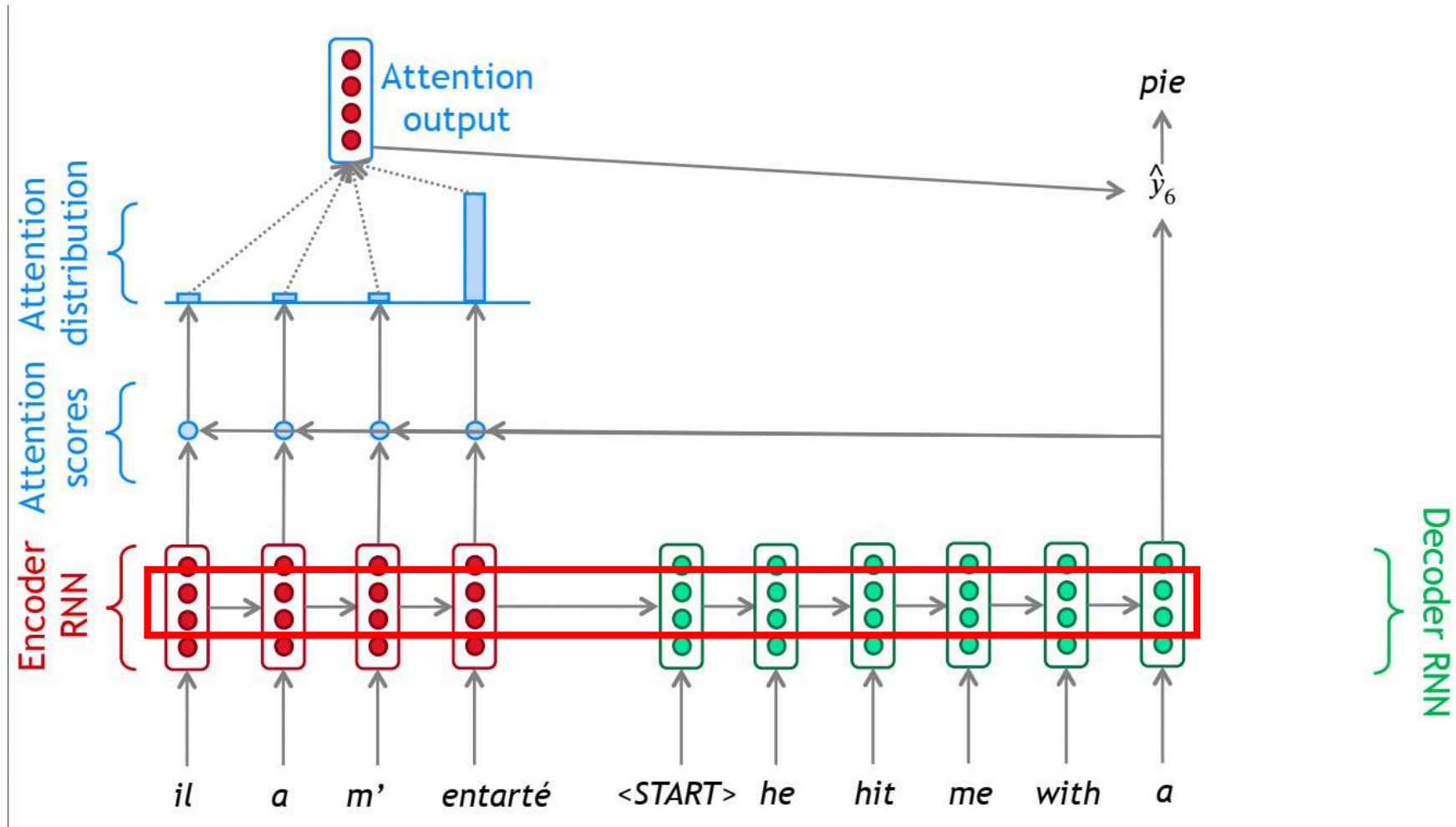
Seq2Seq with attention

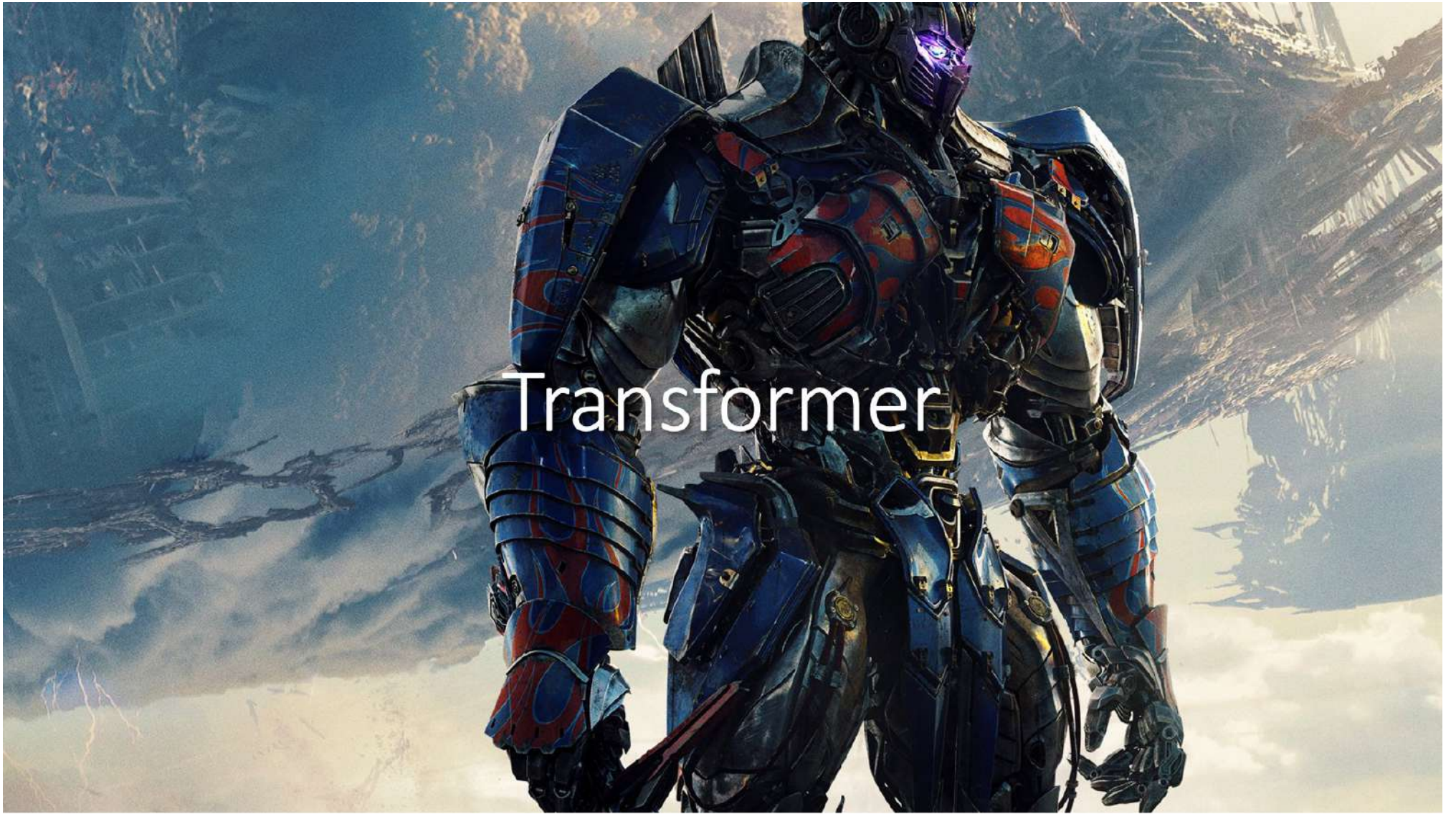


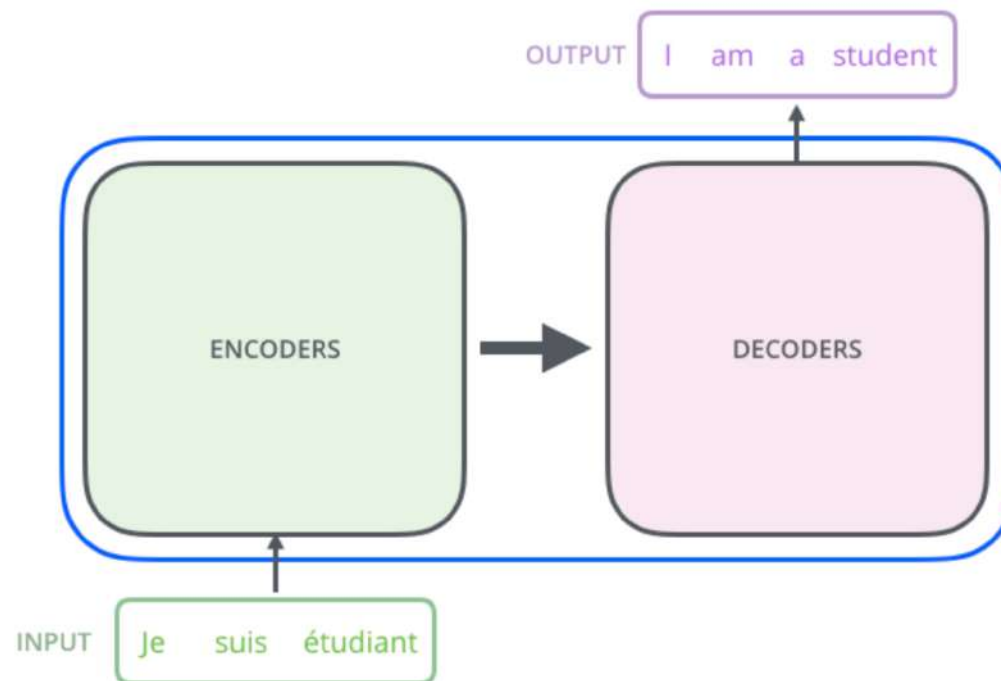
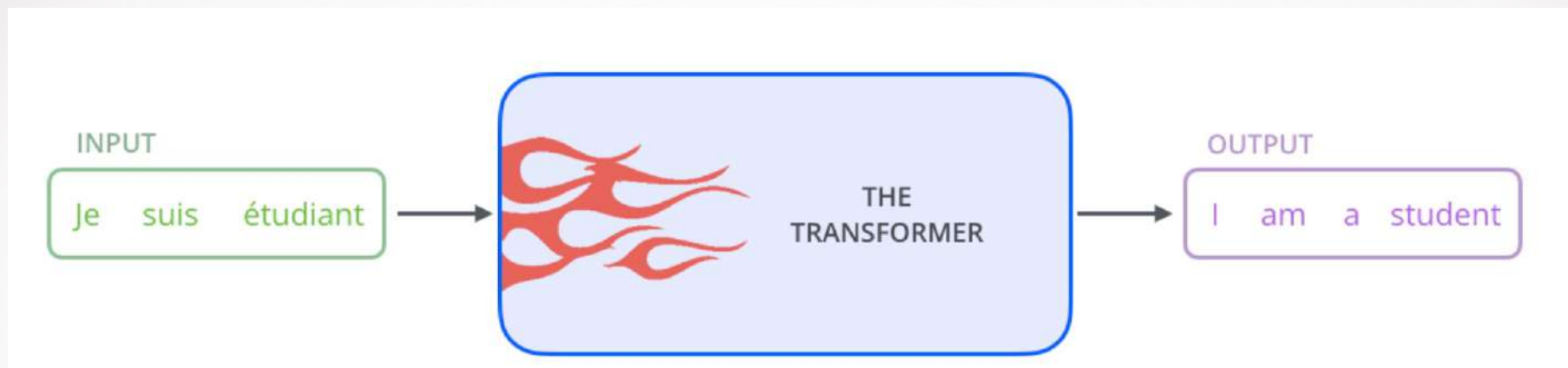
Seq2Seq with attention



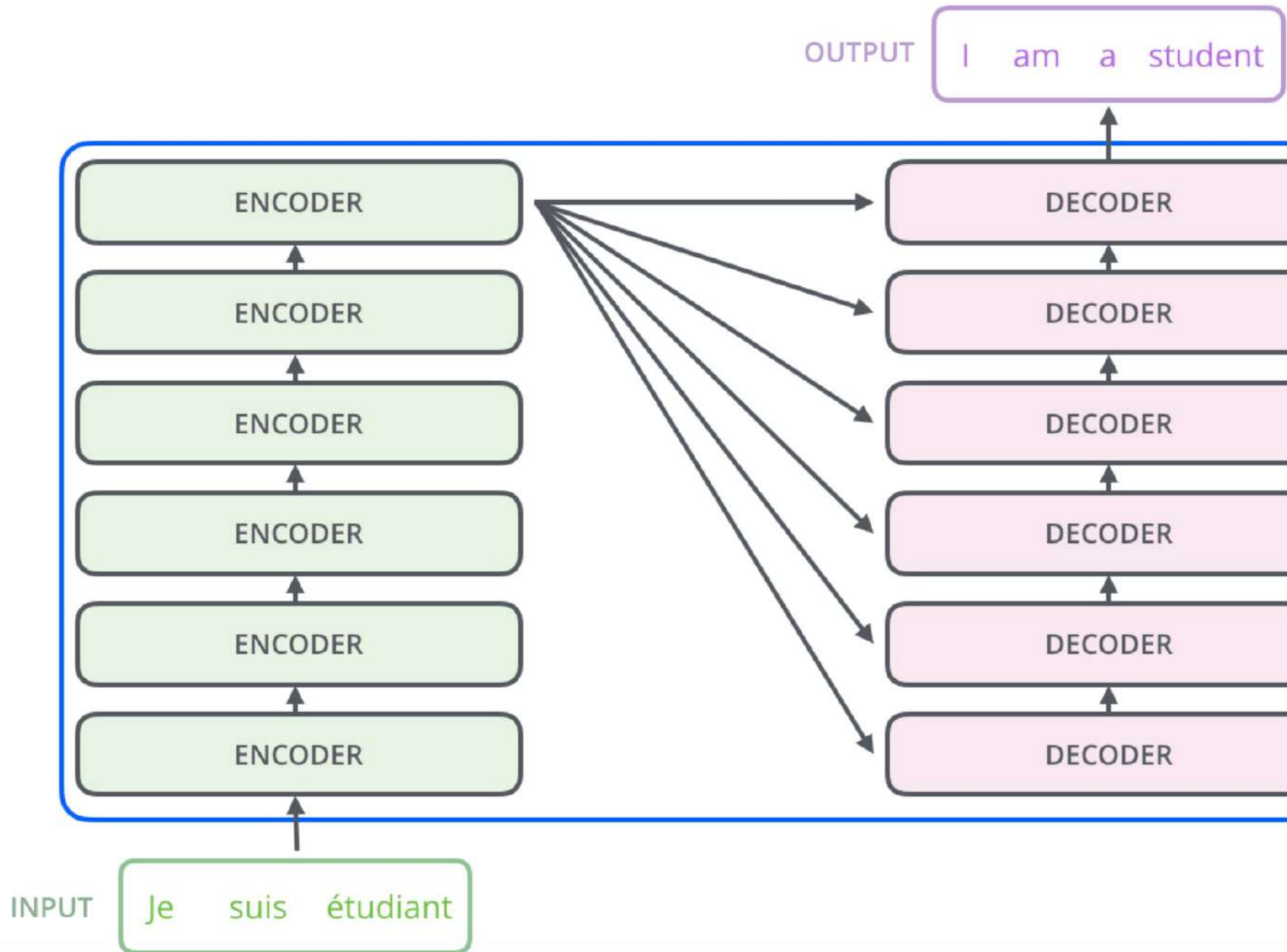
Seq2Seq with another bottleneck



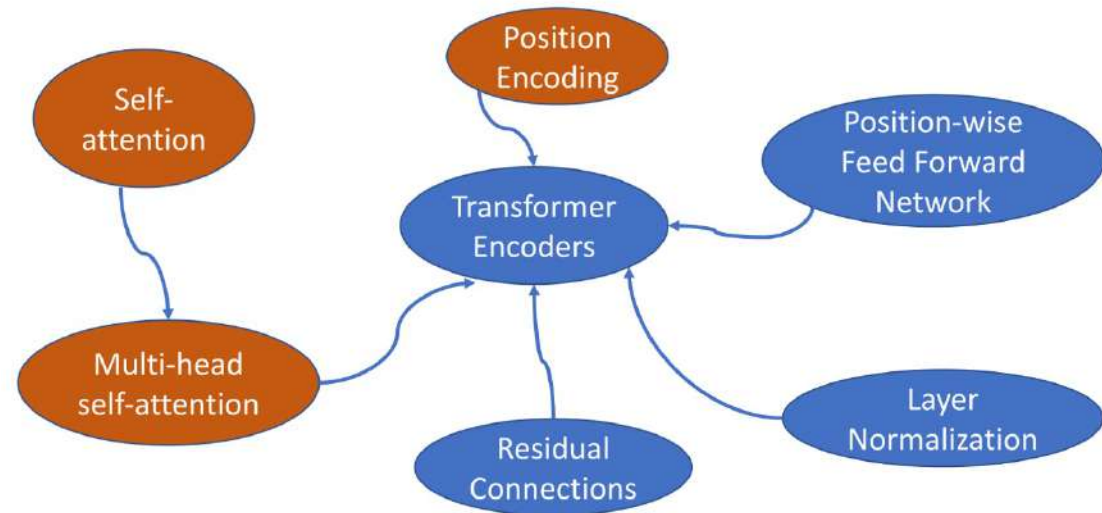
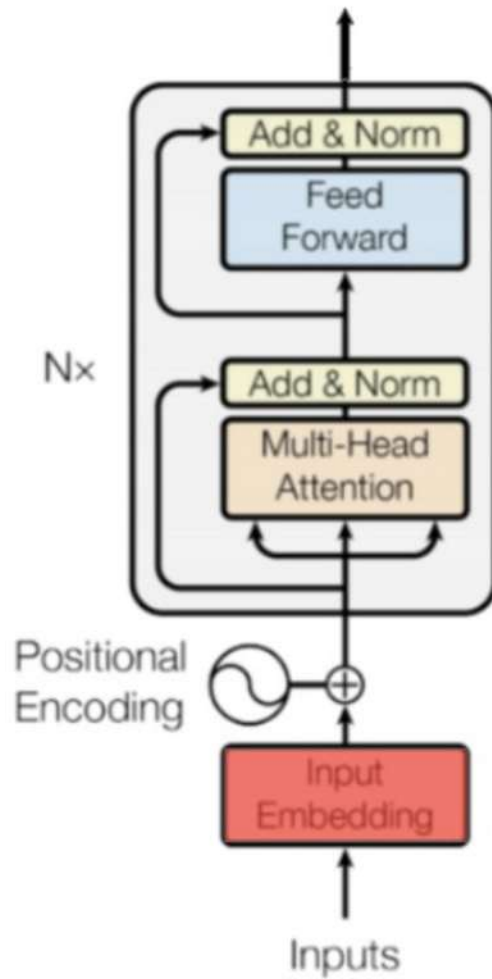




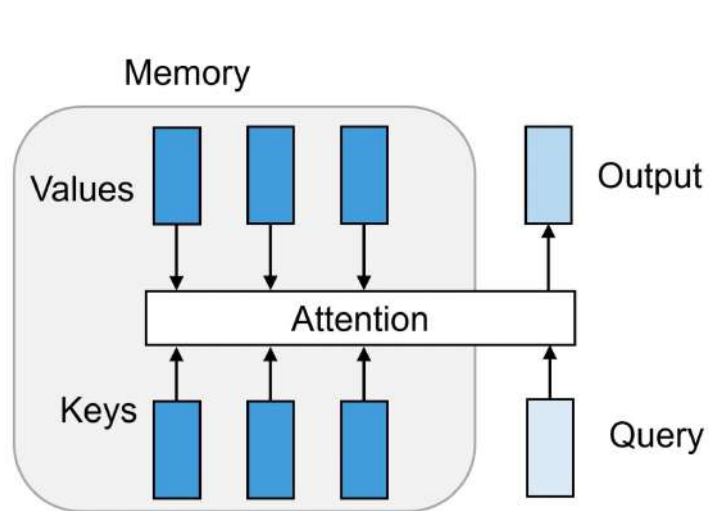
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, & R. Garnett (Eds.), *Advances in neural information processing systems*. Curran Associates, Inc



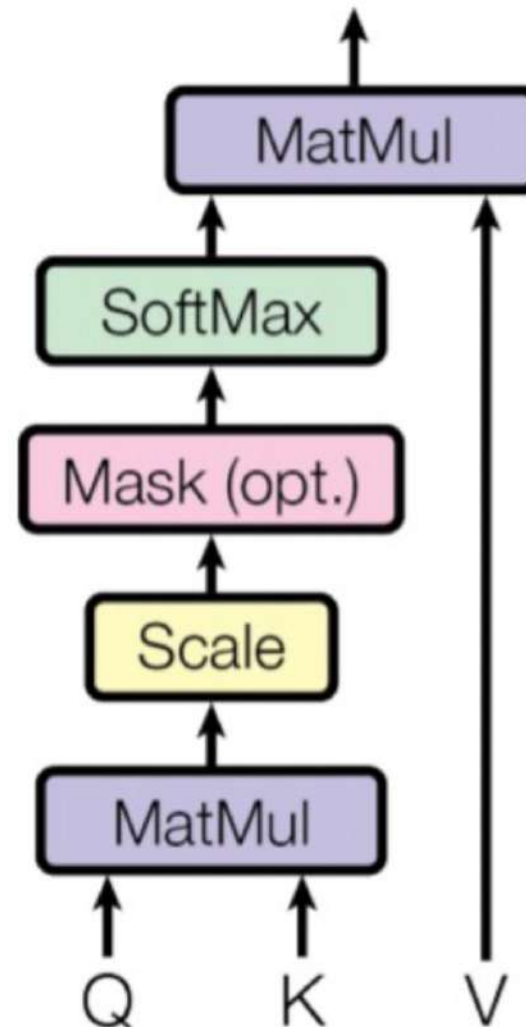
Inside an Encoder Block



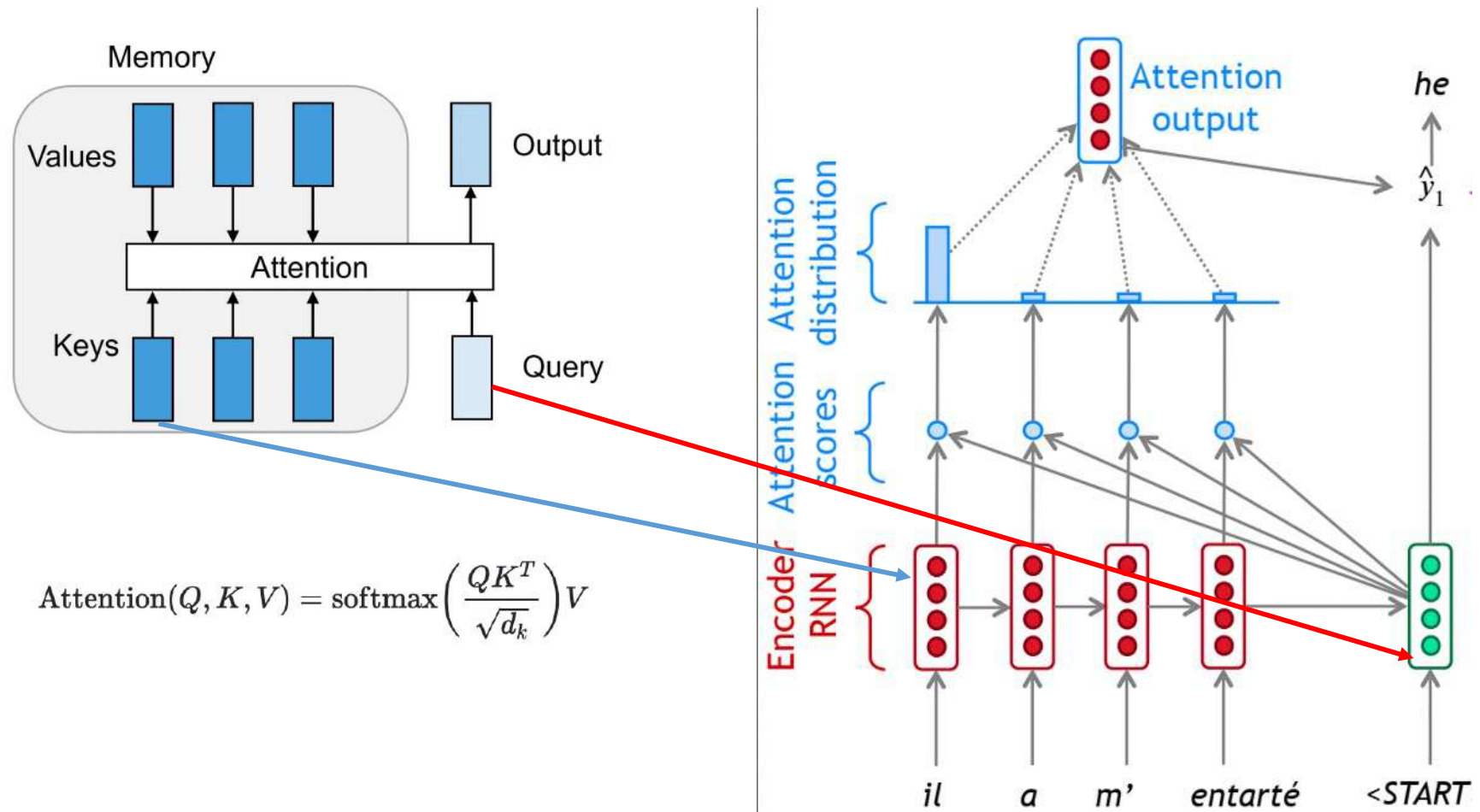
Scaled Dot Product Attention



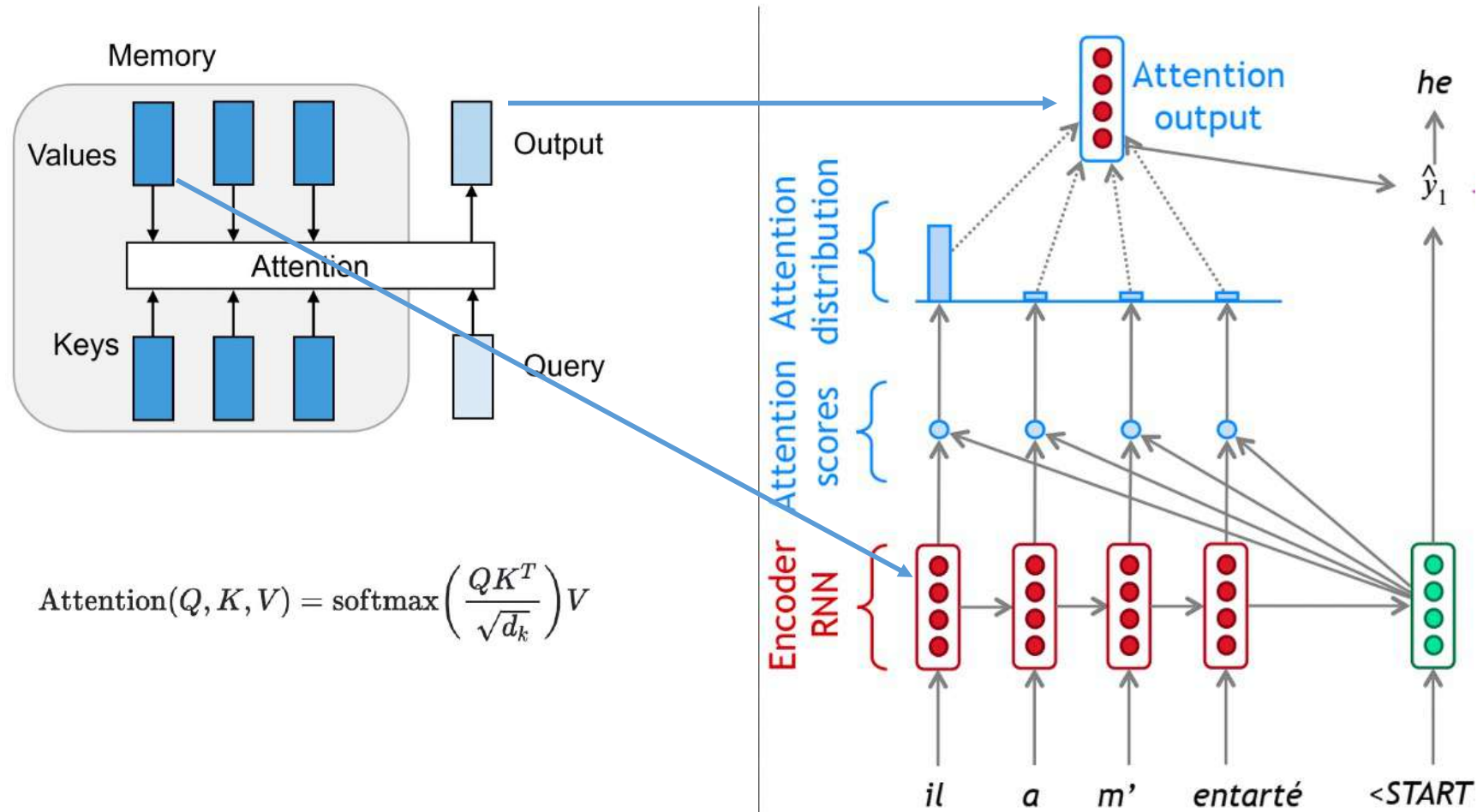
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



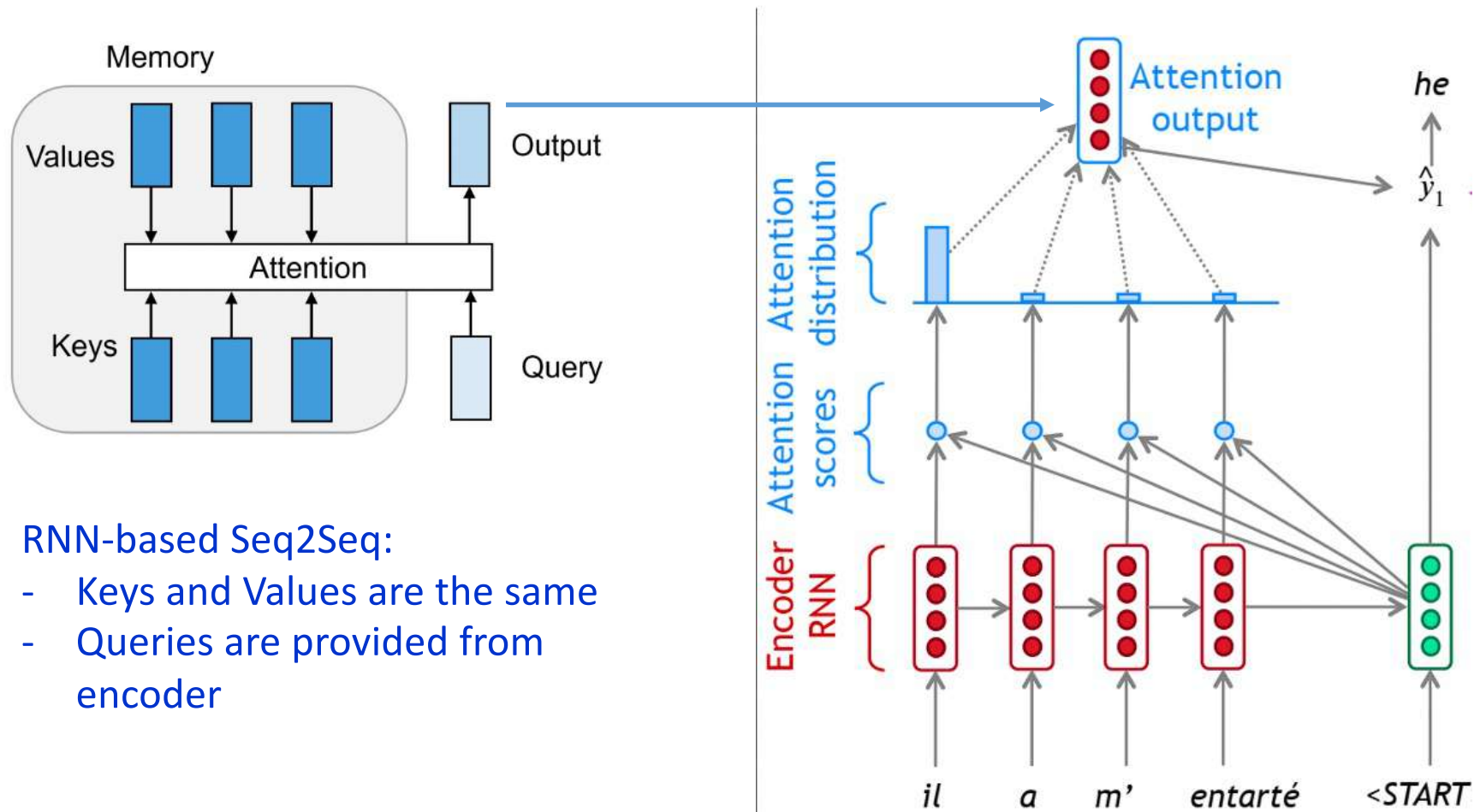
Scaled Dot Product Attention



Scaled Dot Product Attention



Scaled Dot Product Attention

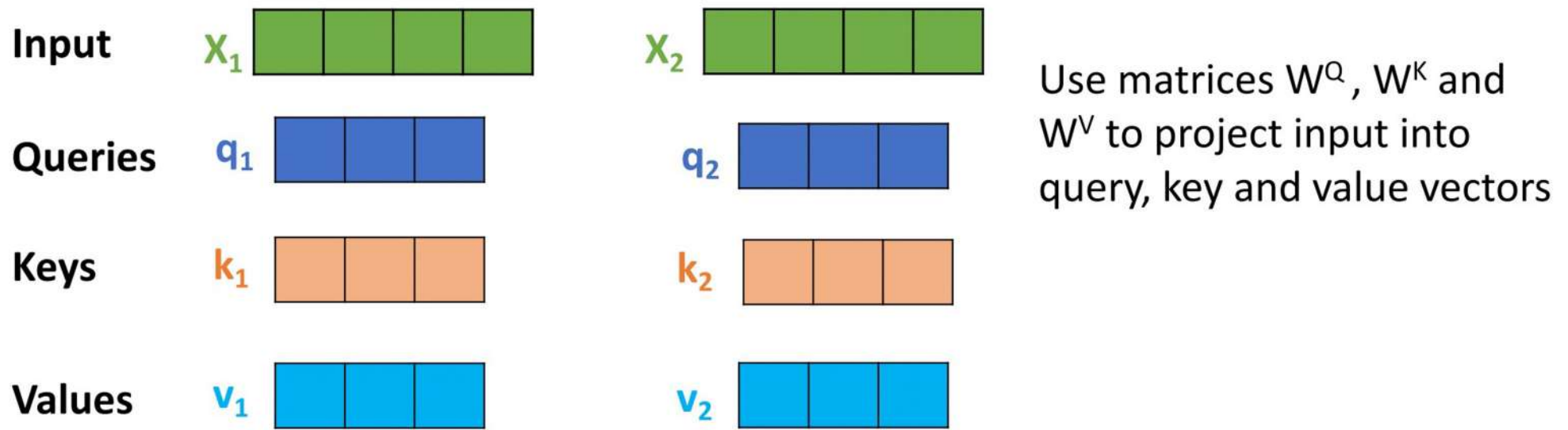


RNN-based Seq2Seq:

- Keys and Values are the same
- Queries are provided from encoder

Self-Attention in Transformer

- Attention maps a query and a set of key-value pairs to an output
 - query, keys, and output are all vectors



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

d_k is the dimension of key vectors

Self-Attention

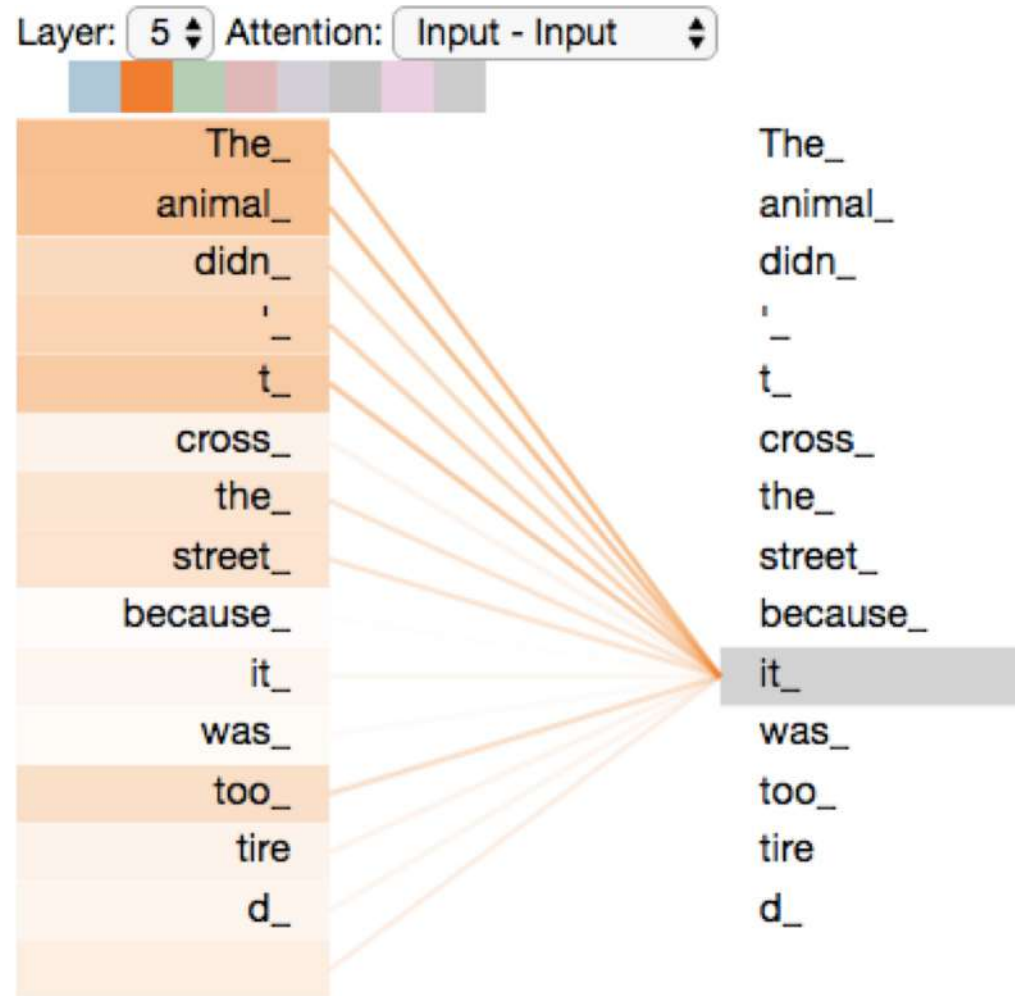


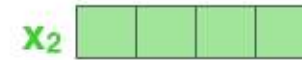
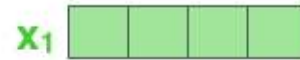
Image source:
<https://jalammr.github.io/illustrated-transformer/>

Input

Thinking

Machines

Embedding



Queries



Keys



Values



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



Sum



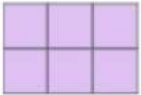
X

Thinking
Machines



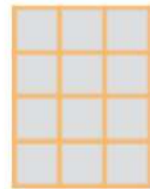
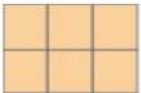
ATTENTION HEAD #0

Q_0



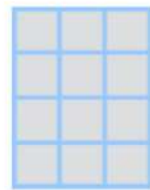
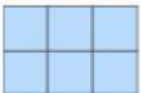
W_0^Q

K_0



W_0^K

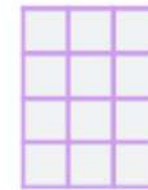
V_0



W_0^V

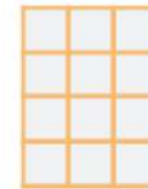
ATTENTION HEAD #1

Q_1



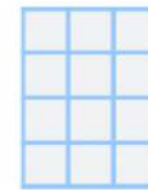
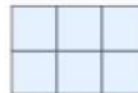
W_1^Q

K_1



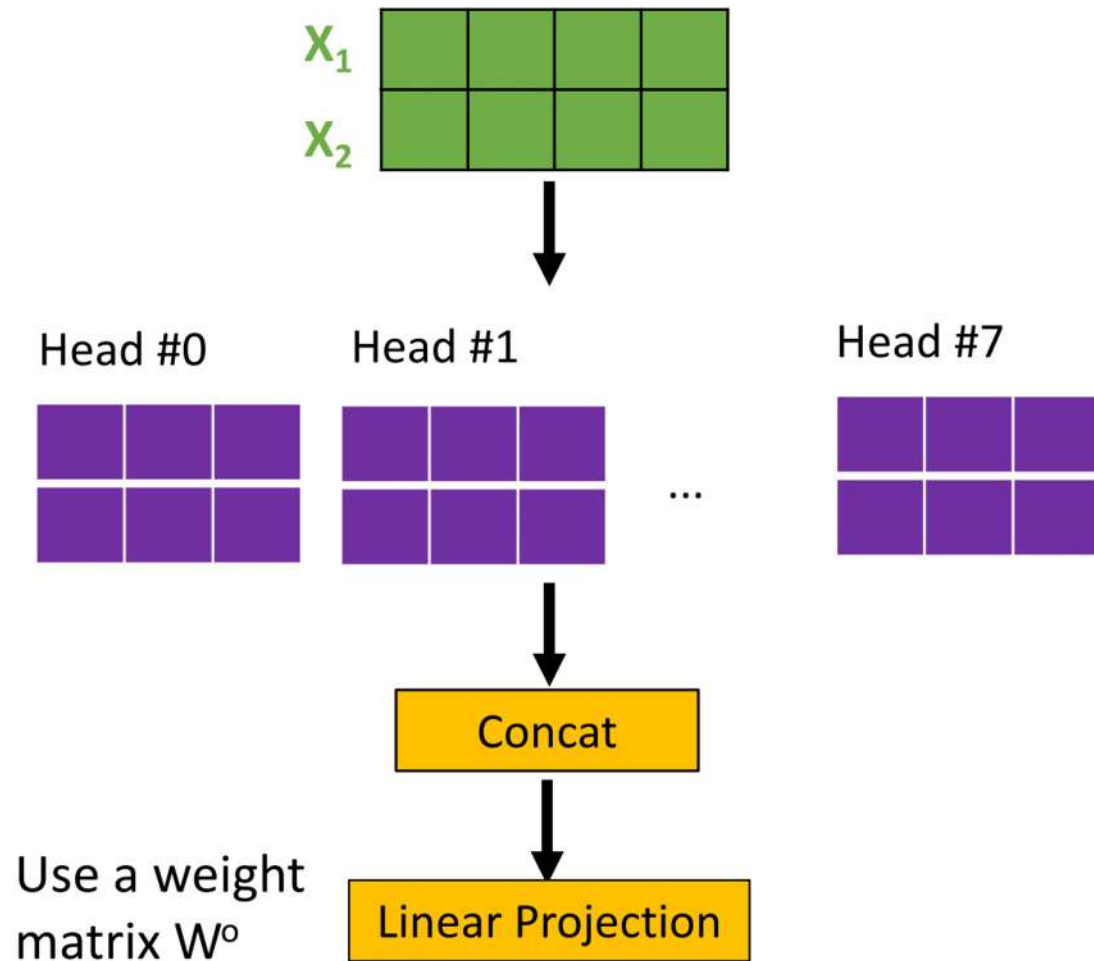
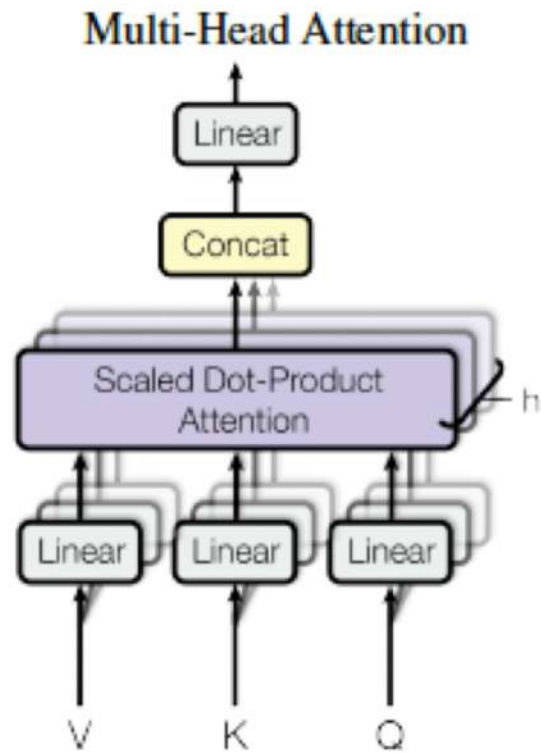
W_1^K

V_1



W_1^V

Multi-Head Attention



1) This is our input sentence*

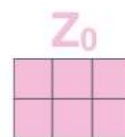
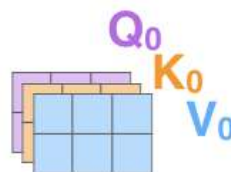
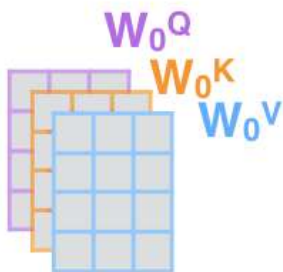
2) We embed each word*

3) Split into 8 heads. We multiply X or R with weight matrices

4) Calculate attention using the resulting $Q/K/V$ matrices

5) Concatenate the resulting Z matrices, then multiply with weight matrix W^O to produce the output of the layer

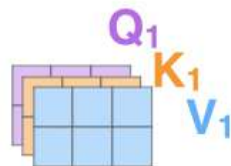
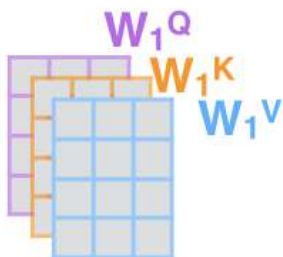
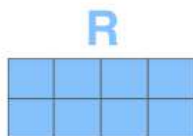
Thinking
Machines



W^O



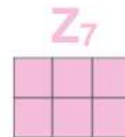
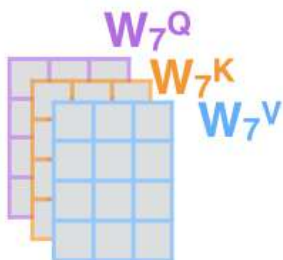
* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



...

...

...



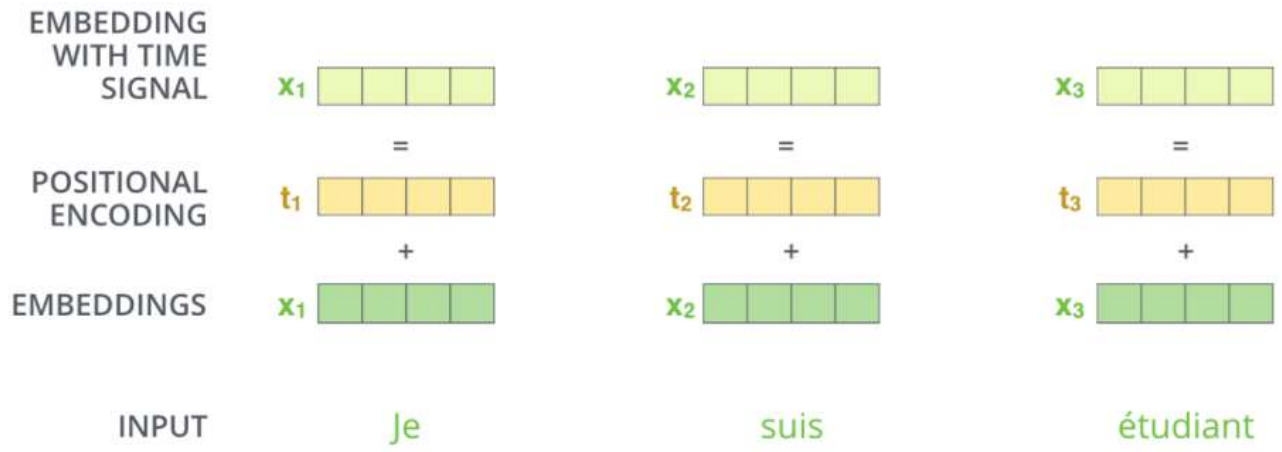
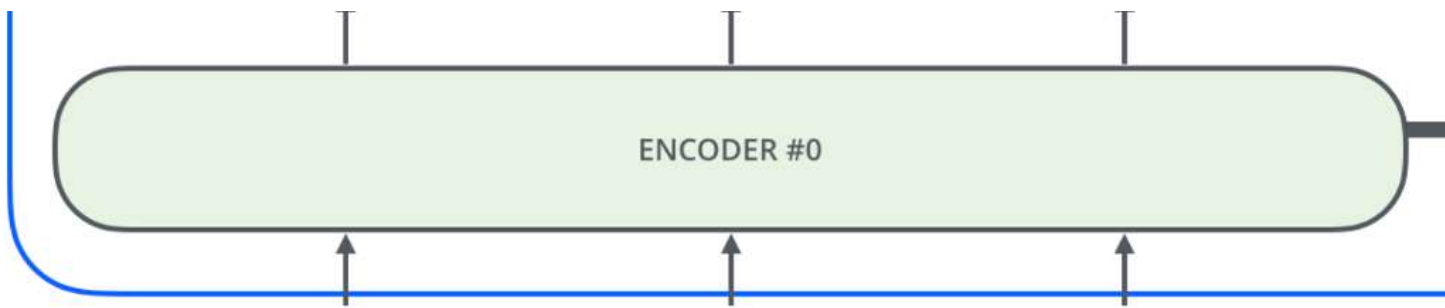
Position Encoding

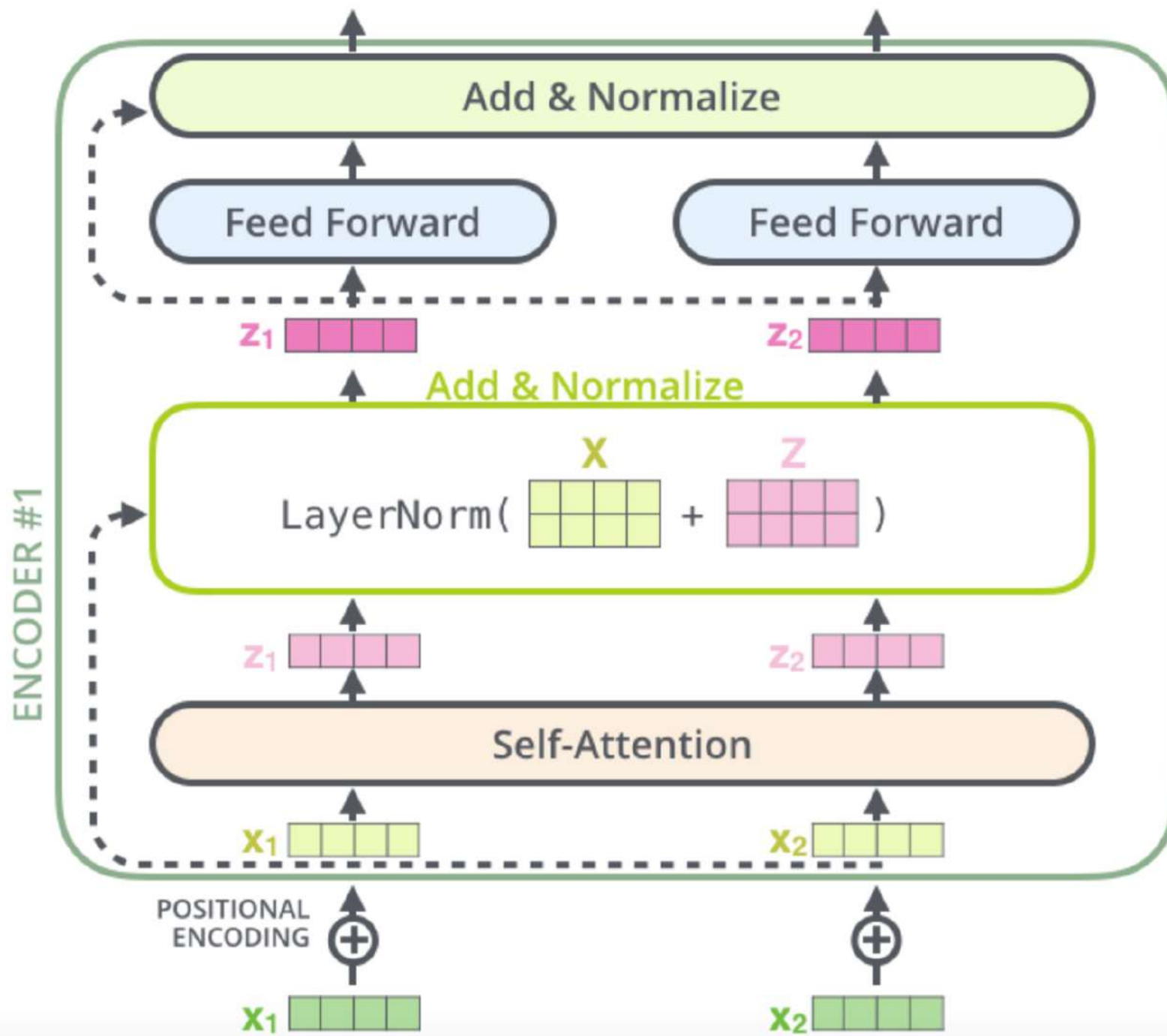
- Position Encoding is used to make use of the order of the sequence
 - Since the model contains no recurrence and no convolution
- In Vawasni et al., 2017, authors used sine and cosine functions of different frequencies

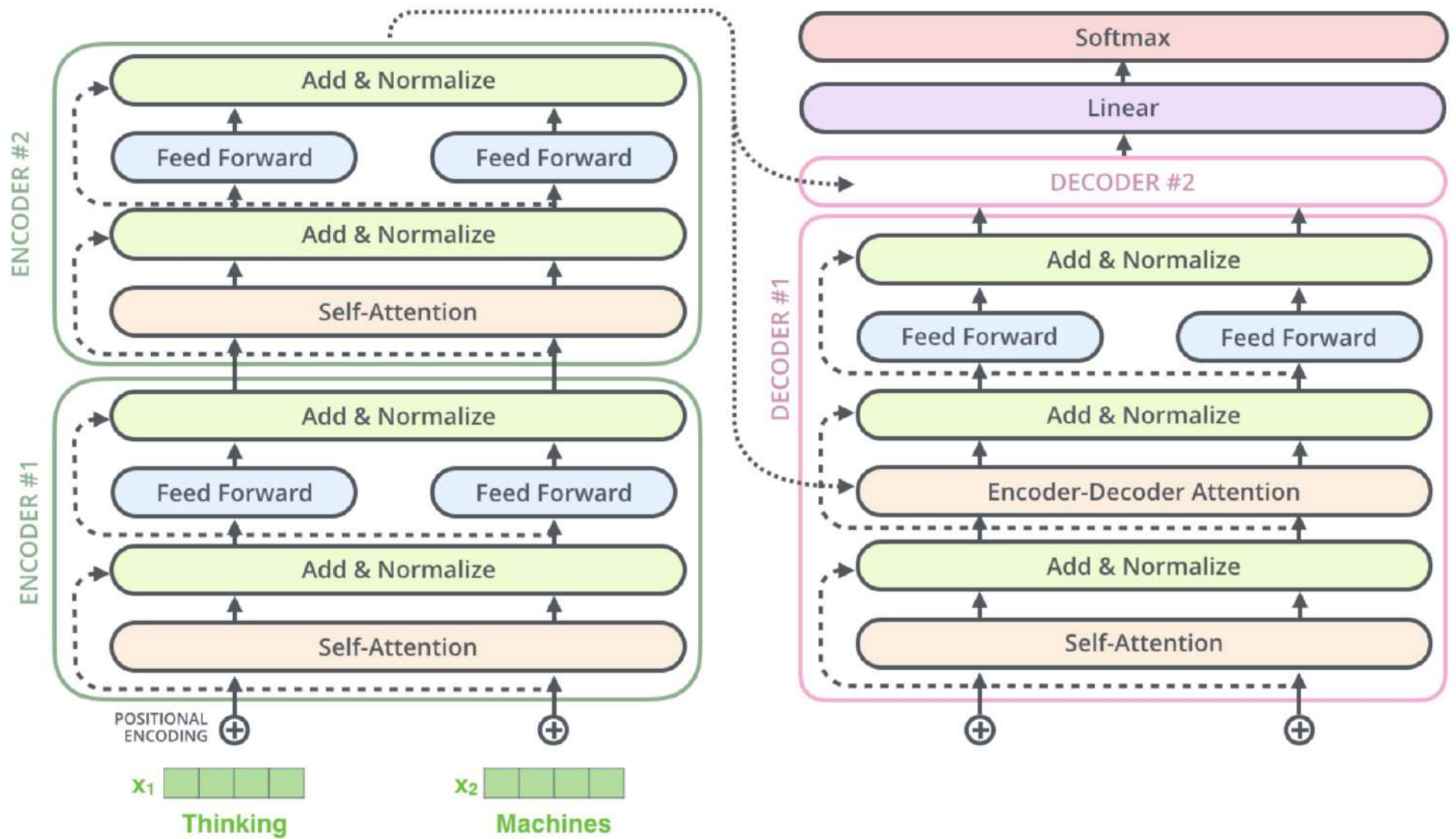
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

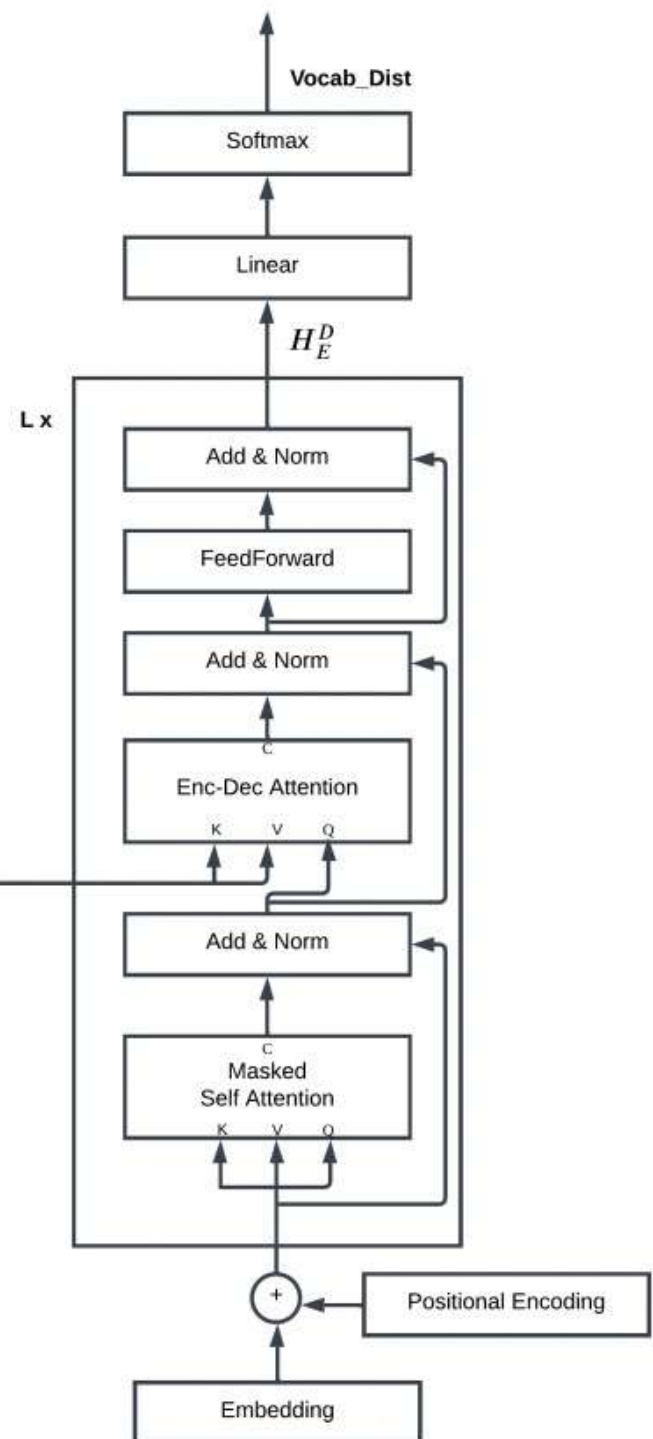
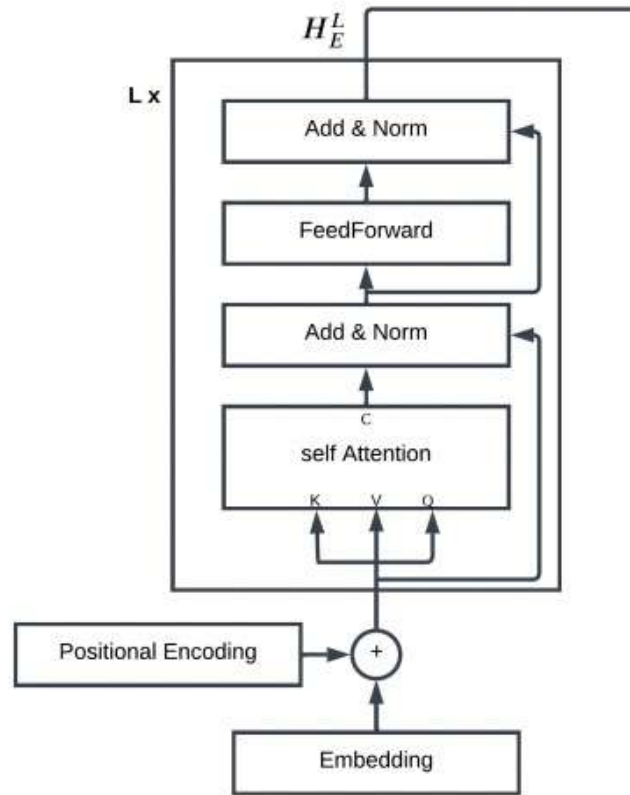
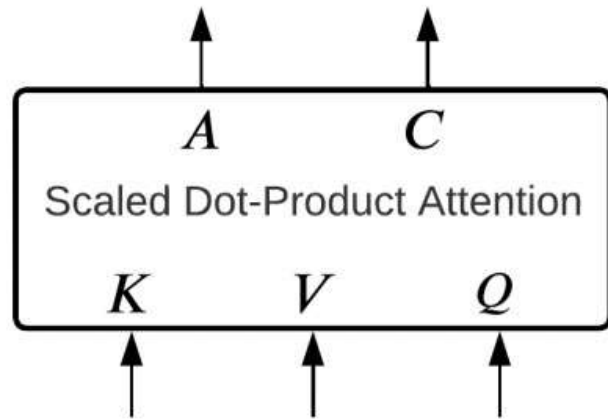
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

- pos is the position and i is the dimension



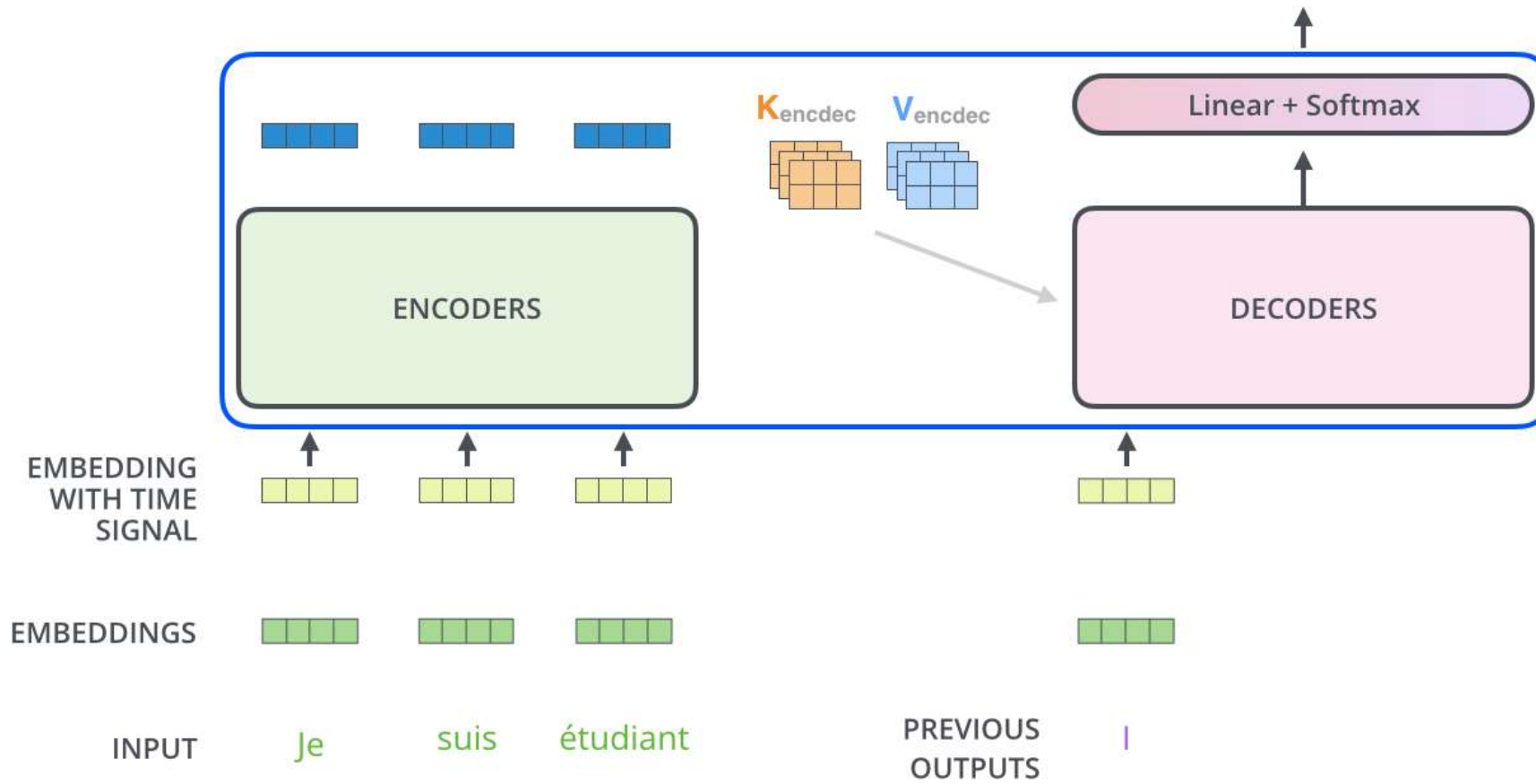






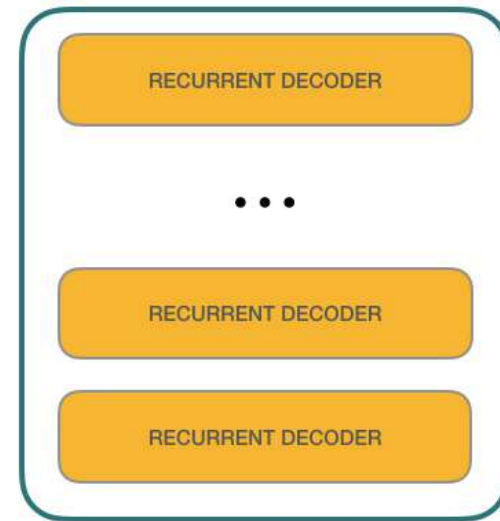
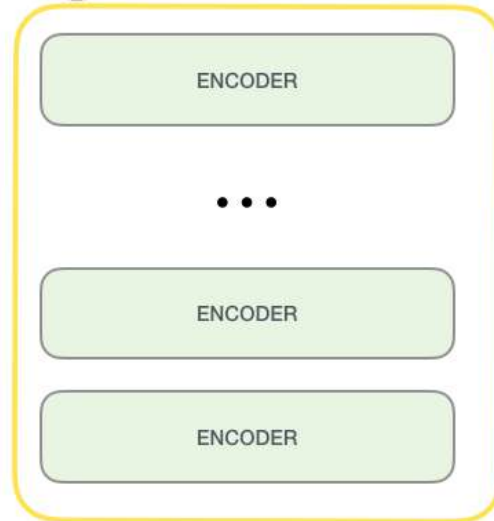
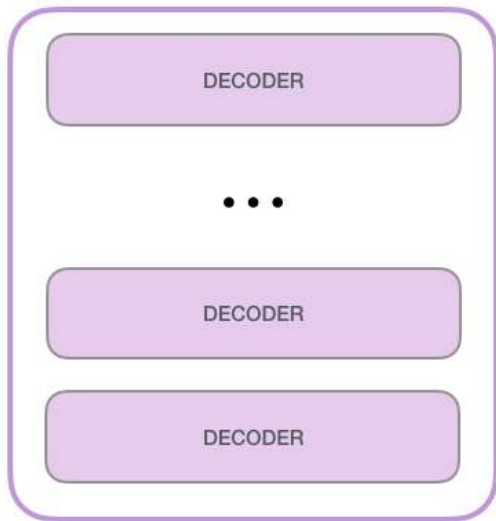
Decoding time step: 1 2 3 4 5 6

OUTPUT |



BERT and other variants

Transformer-based Language Models



BERT

- **Bidirectional Encoder Representations from Transformers.**
- Use the Transformer Encoder architecture.
- Introduced in 2018 by Google AI.

Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2019). *BERT: pretraining of deep bidirectional transformers for language understanding* (J. Burstein, C. Doran, & T. Solorio, Eds.).

unsupervised

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.

Semi-supervised Learning Step

Model:



Dataset:



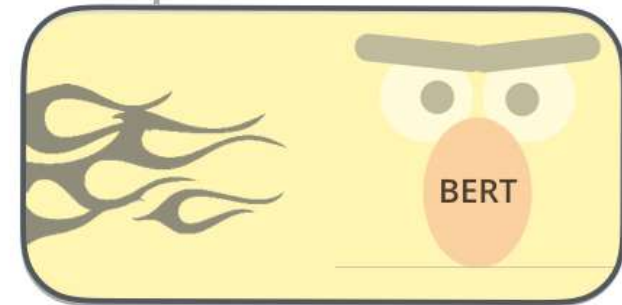
Objective:

Predict the masked word (language modeling)

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

Supervised Learning Step

Model:
(pre-trained in step #1)

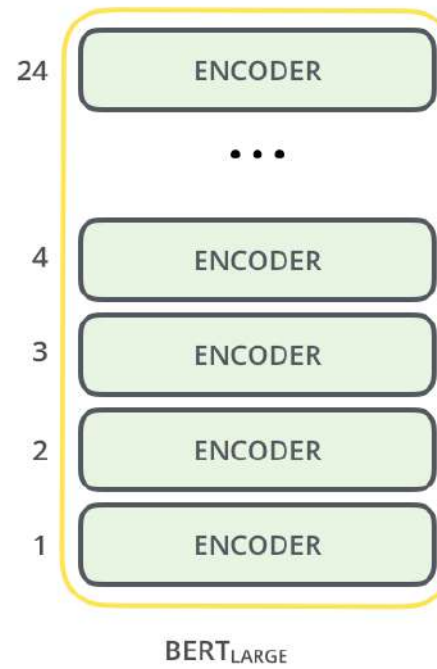
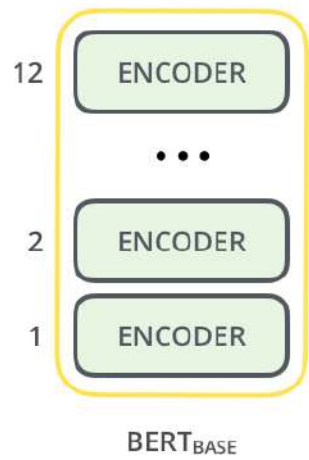


75% Spam
25% Not Spam

Dataset:

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

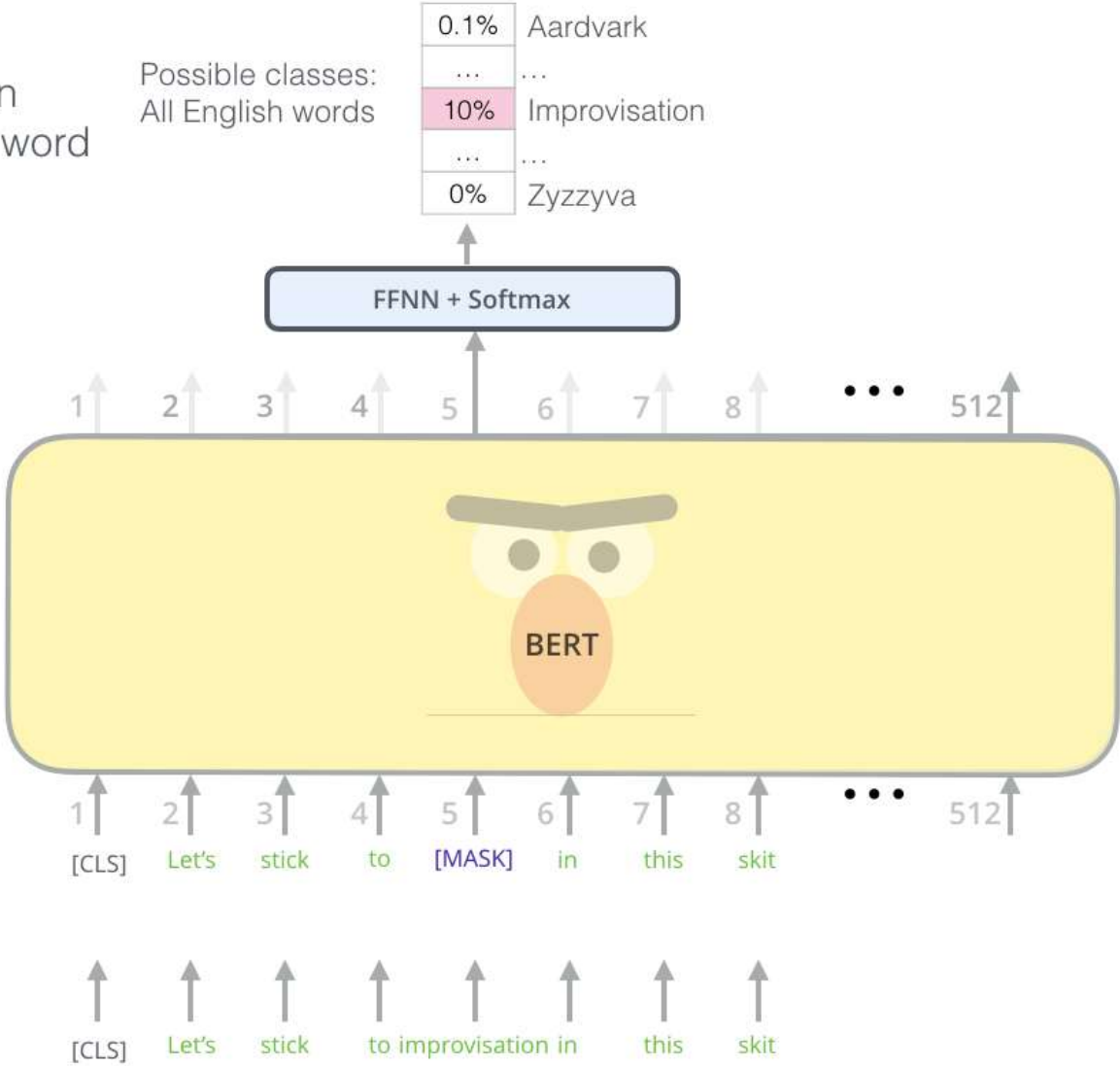
Architecture



Pretraining

- Two unsupervised tasks:
 1. Masked Language Model
 2. Next Sentence Prediction

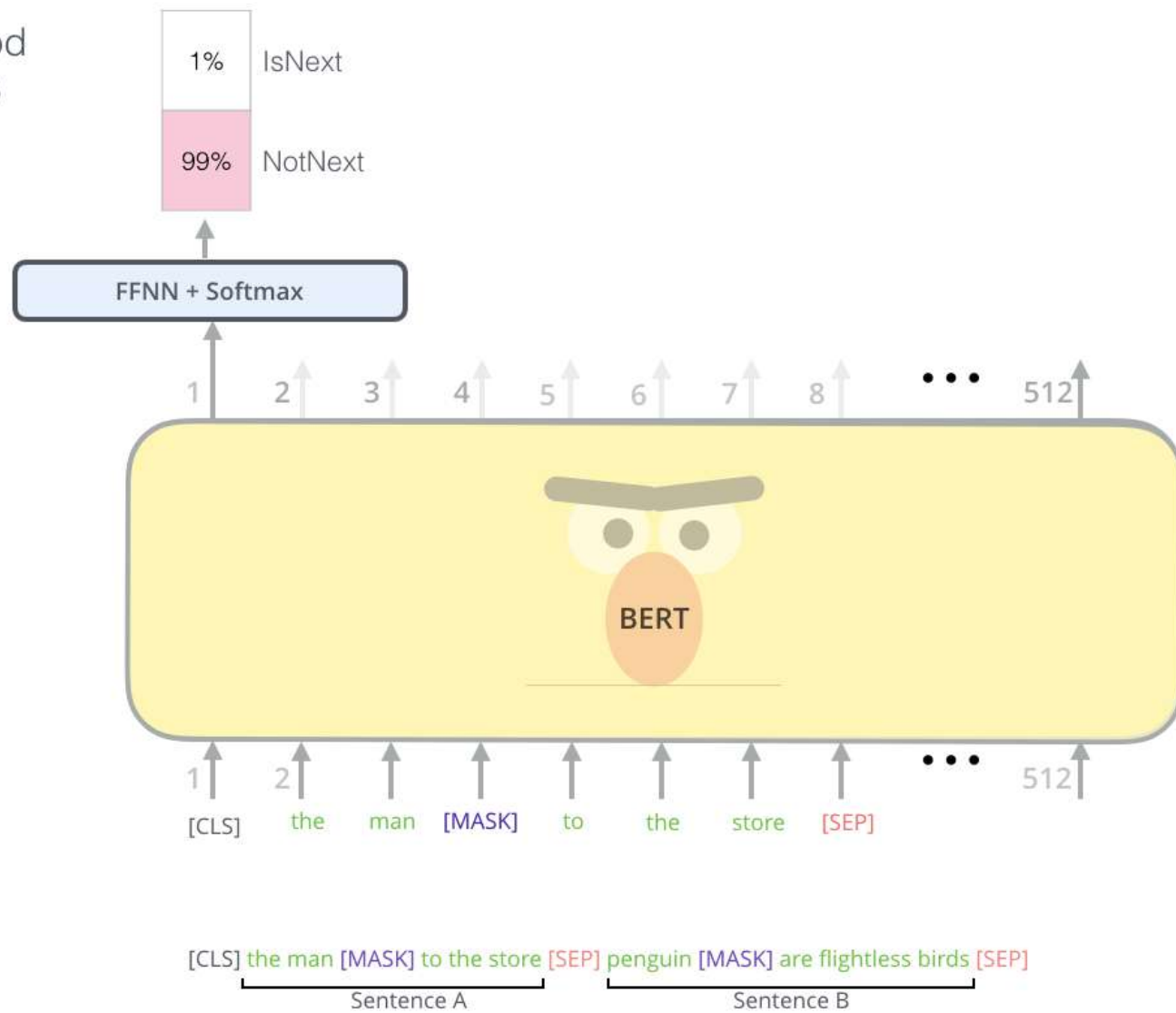
Use the output of the masked word's position to predict the masked word



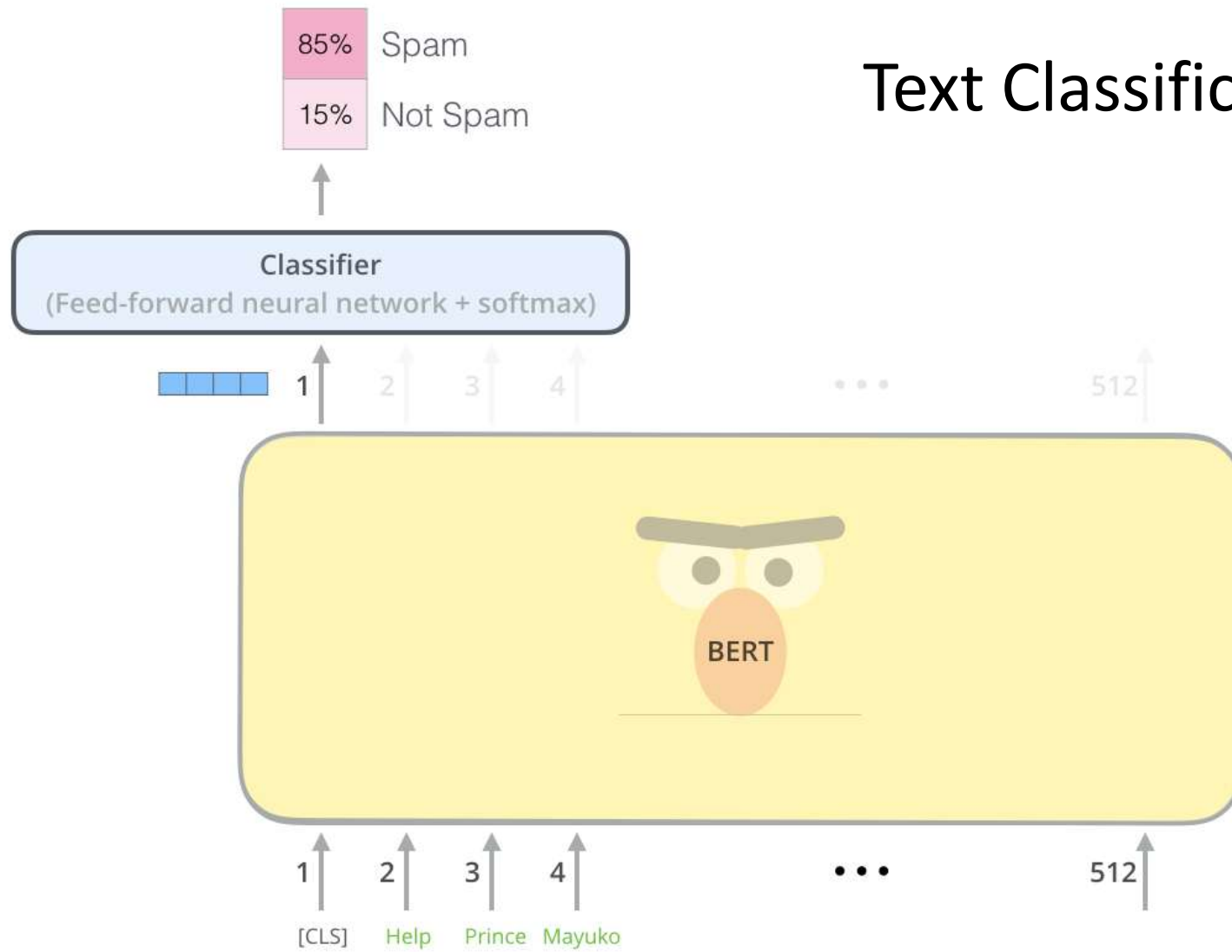
Randomly mask 15% of tokens

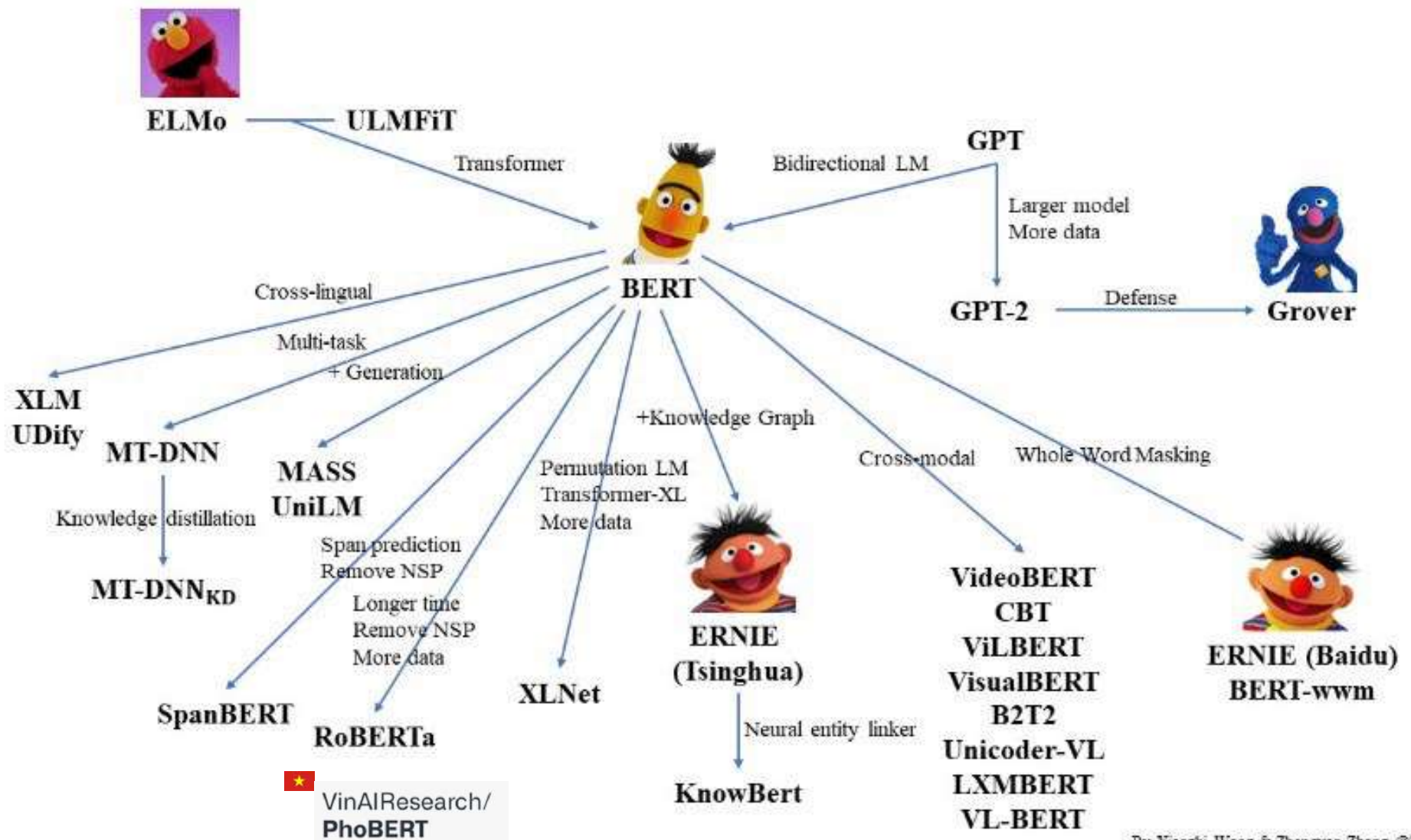
Input

Predict likelihood that sentence B belongs after sentence A



Text Classification





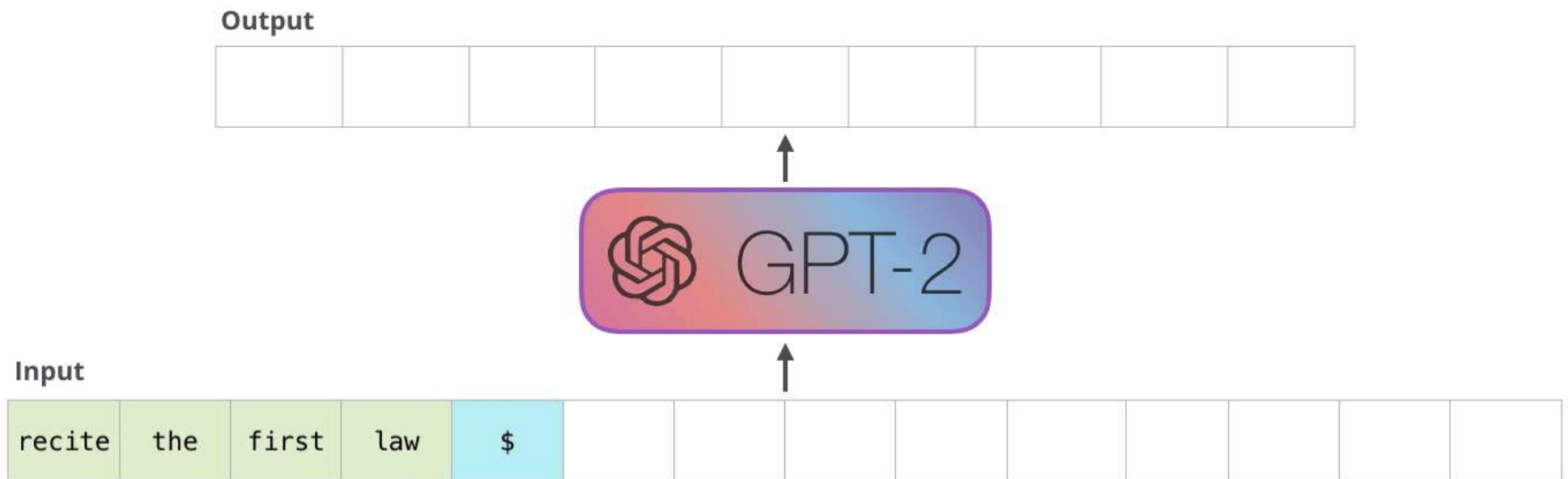
GPT

- Generative Pre-trained Transformer
- Use the Transformer Decoder architecture.
- Introduced in 2018 by OpenAI.

Model	Number of parameters	Training data size	Year
GPT	110M	4GB	2018
GPT-2	1.5B	40GB	2019
GPT-3	175B	≈2TB	2020

Openai [Accessed: 2023-03-01]. (2023). <https://openai.com/>

How it works?



XLNet

- Autoencoding (BERT):
 - [MASK] tokens do not appear during finetuning \Rightarrow pretrain-finetuning discrepancy.
 - Assume the predicted tokens are independent of each other given the unmasked tokens. Example: “New York is a city” \Rightarrow “[MASK] [MASK] is a city”
- Autoregressive (GPT):
 - Only trained to encode a unidirectional context (forward or backward).

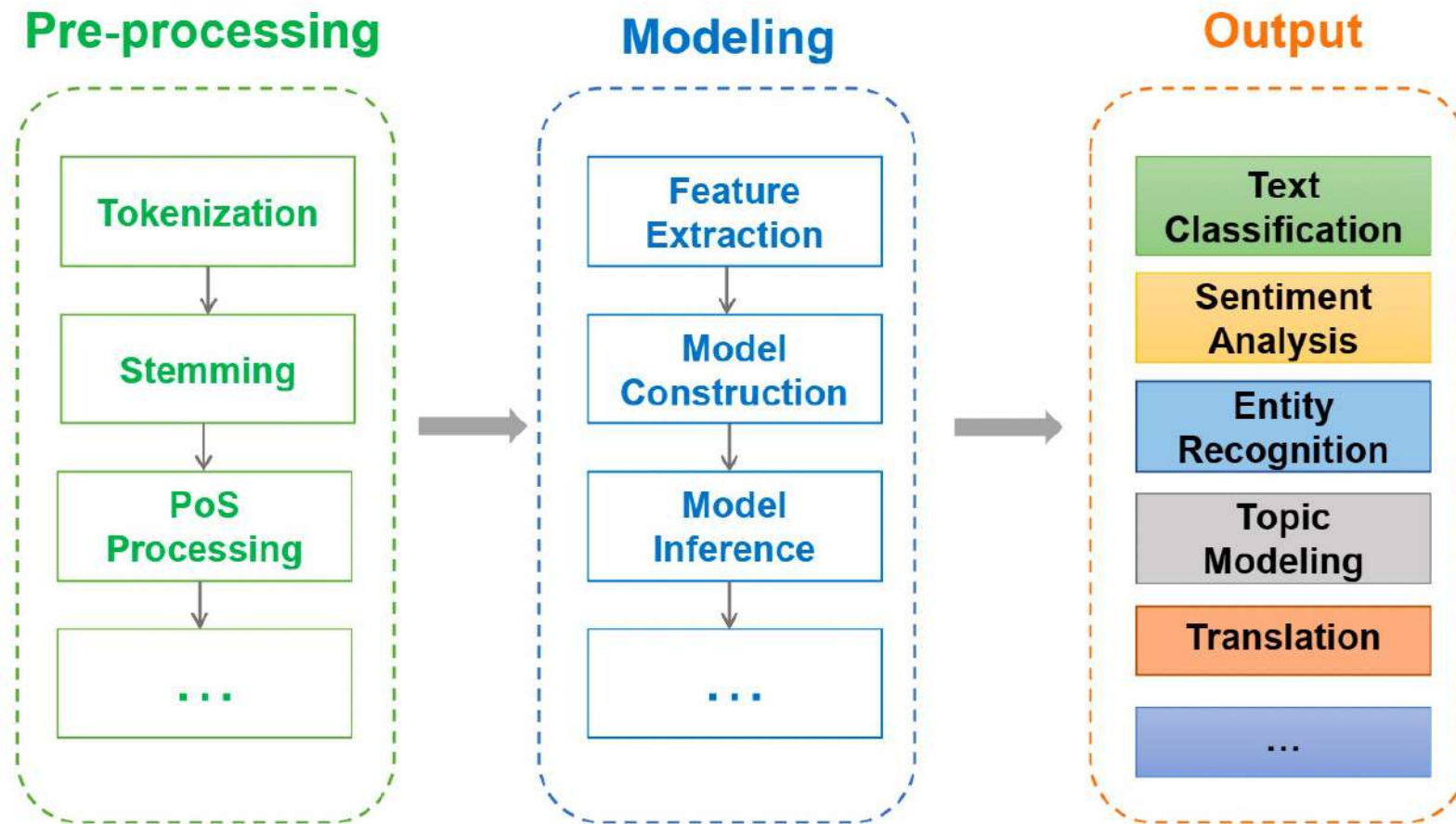
Yang, Z. et al. “XLNet: Generalized Autoregressive Pretraining for Language Understanding.” *NeurIPS* (2019)

XLNet

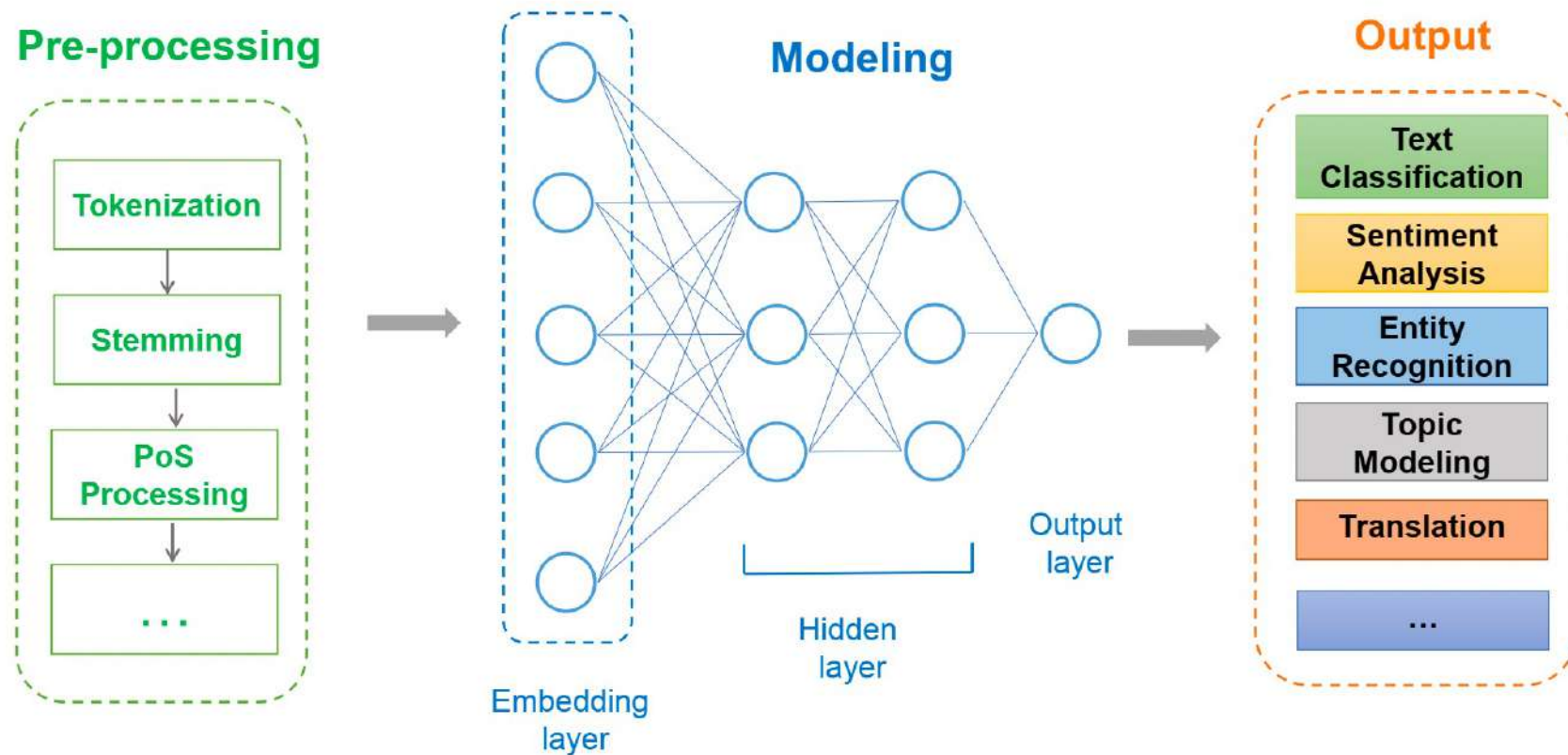
- XLNet combines pros from both while avoiding their cons.
- Techniques:
 - Permutation Language Modeling
 - Two-Stream Self-Attention for Target-Aware Representations
 - Incorporating Ideas from Transformer-XL
 - Modeling Multiple Segments

Applications in NLP

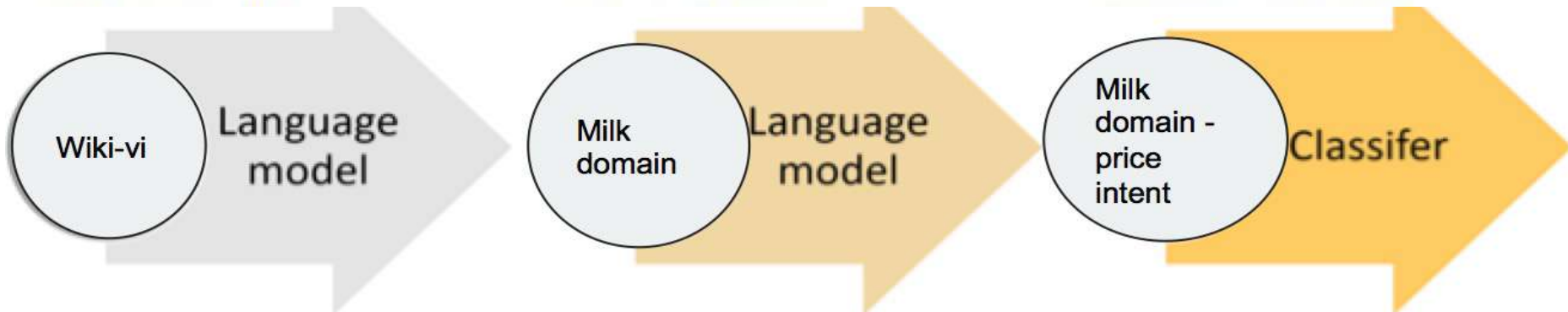
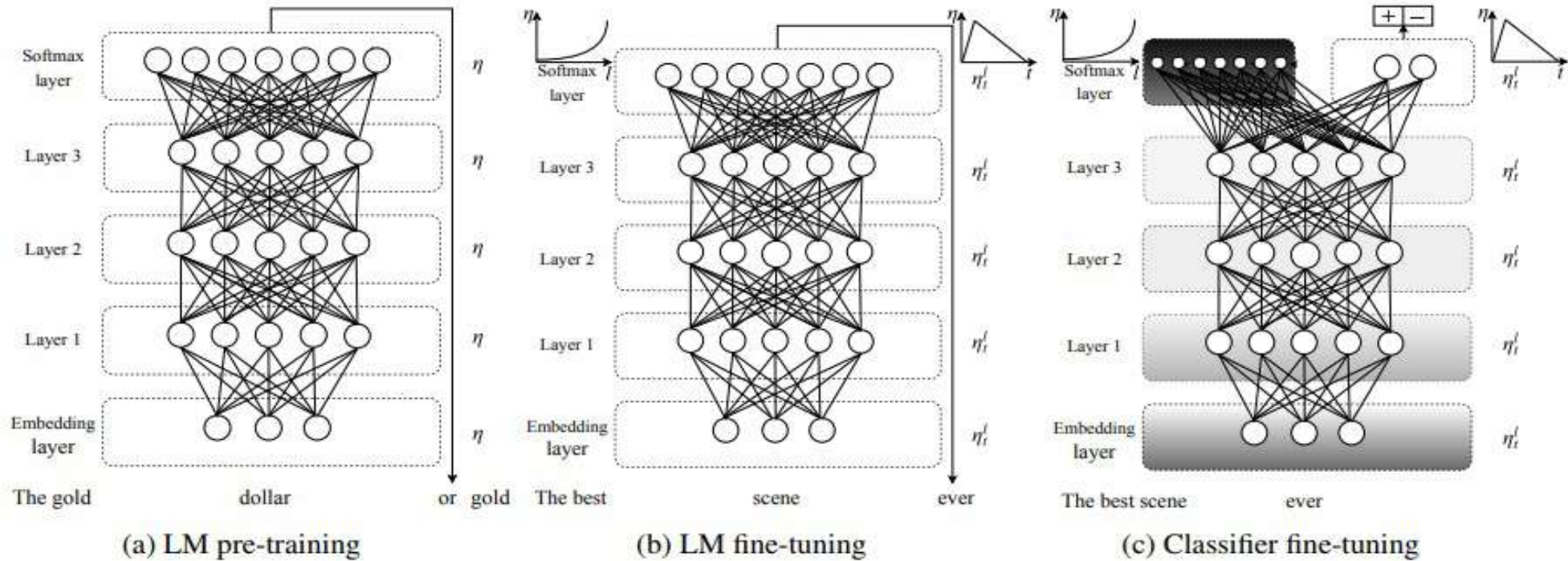
NLP typical pipeline



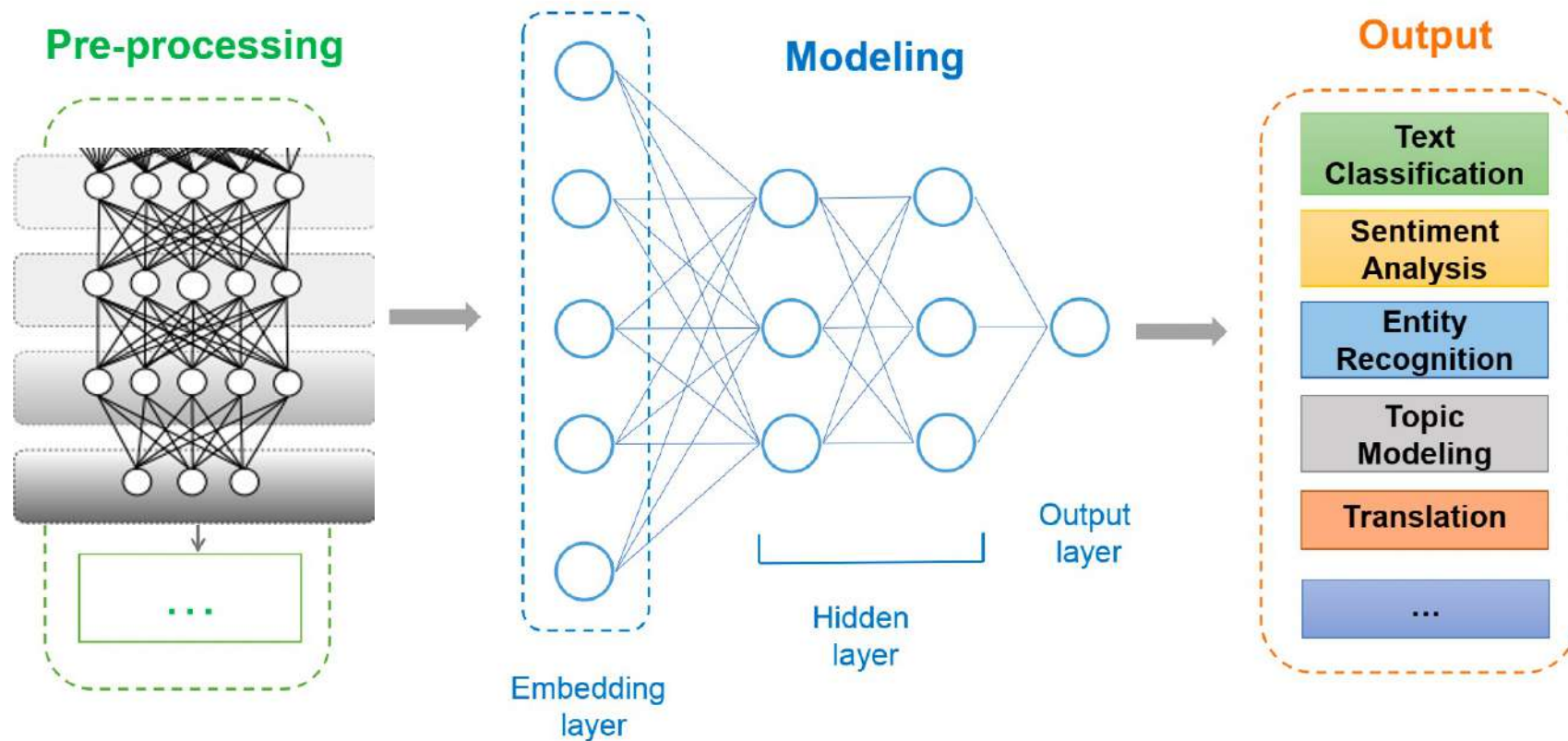
NLP DL-based pipeline



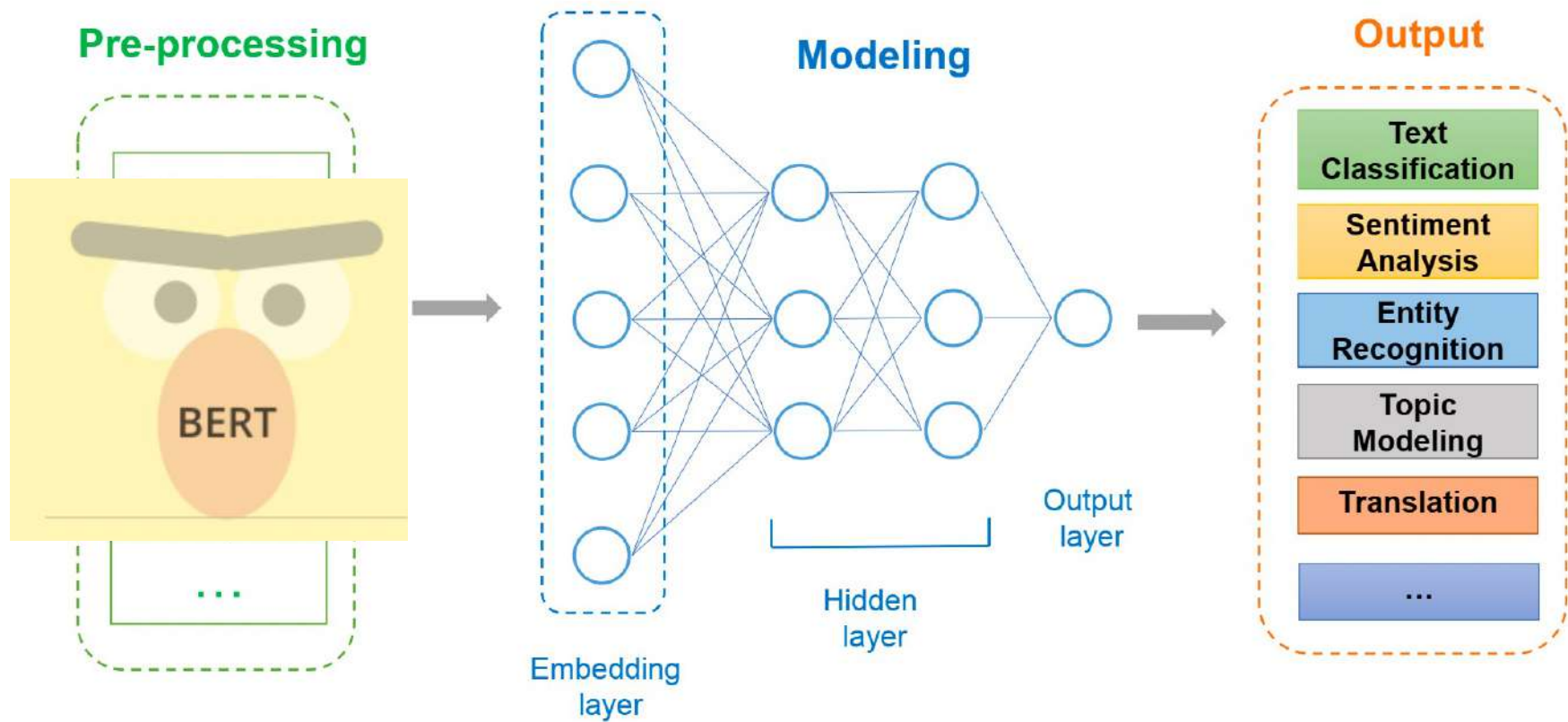
Pre-trained Neural Language Model



NLP LM-based pipeline

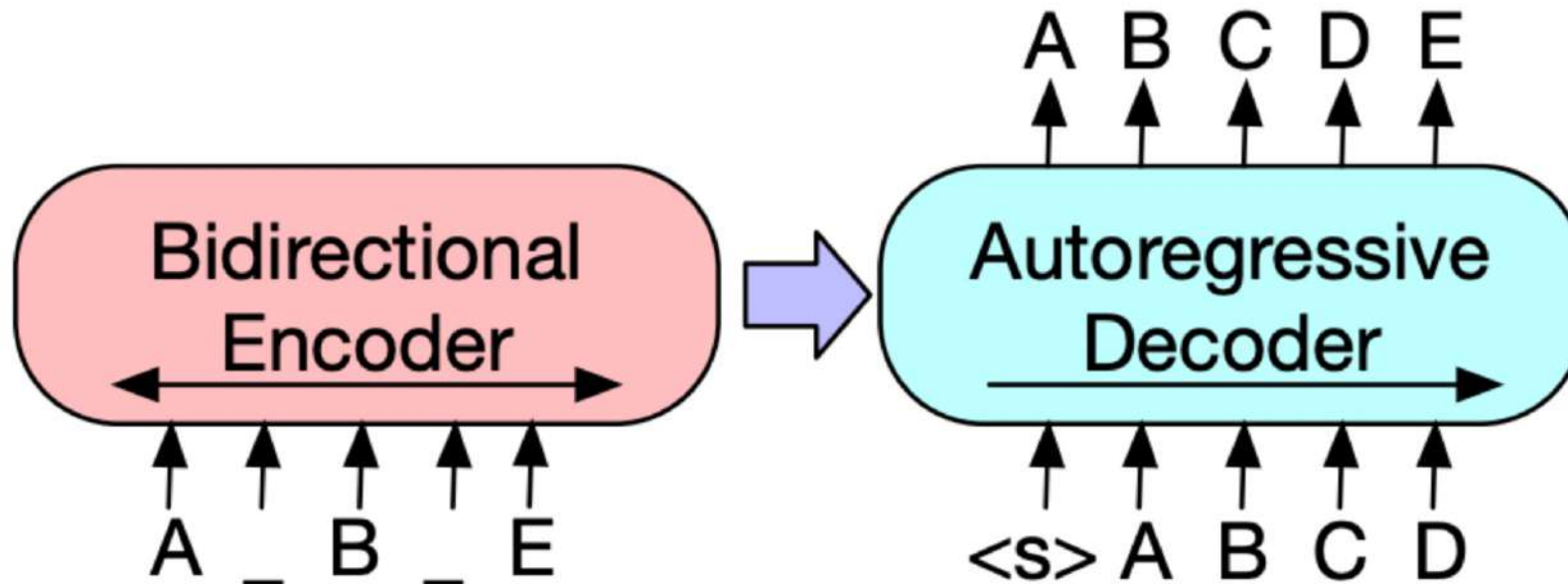


NLP LM-based pipeline



From BERT to BART

- BERT is not a fully Seq2Seq model (i.e. not a generative model)
- BART is introduced as an extended/complement



Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

From PhoBERT to BARTPho

VinAIResearch/ **BARTpho**



BARTpho: Pre-trained Sequence-to-Sequence
Models for Vietnamese (INTERSPEECH 2022)

 1

Contributor

 0

Issues

 75

Stars

 6

Forks



BARTPho for Vietnamese translation applications

- Pretrained with Vietnamese
- Implicitly processing “aligning” task
- More powerful if the target language has similar language to Vietnamese (Chinese, Bahnaric, etc.)

A Demo to be concluded



Trang chủ

Về hệ thống

Liên hệ

Đăng nhập

Xin chào, Khách

Mô hình Combined	Giọng Nam	Vùng Bình Định
Tiếng Việt xin chào, RADL 2023	Tiếng Bana Apinh chau, radl 2023.	ĐỌC

- <https://www.ura.hcmut.edu.vn/bahnar/nmt>

Thank you