

Deep Learning for Natural Language Processing

Lecture 4: Pre-trained Language Model

Quan Thanh Tho

Faculty of Computer Science and Engineering
Back Khoa University

Agenda

- GPT-2 and what we have done
- Practical Situation: Huge Data, Few Labeled
- Language Model: General Architecture
- Case Studies
- GPT, Transformer and BERT

**SYSTEM PROMPT
(HUMAN-WRITTEN)**

In a shocking finding, scientist discovered a herd of unicorns living in a remote, previously unexplored valley, in the Andes Mountains. Even more surprising to the researchers was the fact that the unicorns spoke perfect English.

**MODEL COMPLETION
(MACHINE-WRITTEN,
10 TRIES)**

The scientist named the population, after their distinctive horn, Ovid's Unicorn. These four-horned, silver-white unicorns were previously unknown to science.

Now, after almost two centuries, the mystery of what sparked this odd phenomenon is finally solved.

Dr. Jorge Pérez, an evolutionary biologist from the University of La Paz, and several companions, were exploring the Andes Mountains when they found a small valley, with no other animals or humans. Pérez noticed that the valley had what appeared to be a natural fountain, surrounded by two peaks of rock and silver snow.

trình chuyên môn ngành công nghệ này anh cho rằng thành phố ngàn thí sinh thực hiện việc nhận được các loại các chuyên g học đào tạo nhưng nhiều ngành nghề nghiệp trong thời gian từ

generate_text('một cô gái được các chàng trai yêu thích', 400)

'một cô gái được các chàng trai yêu thích với các con của nhữ ia sẽ dễ đánh giá vào ngày N cho thày thí sinh có thể nhận du hút nhiều người từng thi đấu chỉ theo một cách cho thày khách trước đó thông tin trong hành trình thi thpt quốc gia cũng đã

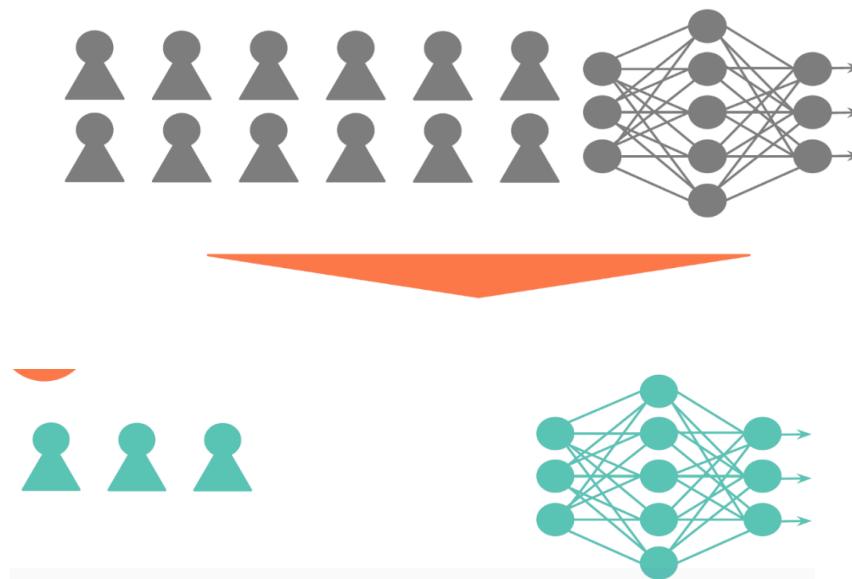
generate_text('andy murr', 400)

'andy murray chi tiết và có thể đánh chú ý cho những bức ảnh m ng bài học chuyên nghiệp ngoài ra chúng tôi đang phải thực hiệ bi và hiệu quả mới nhất giờ tay vợt người thái lan thế giới nh è tham gia công chúng cho thày công ty thiết kế ngày N trong t

Introduction

Common situation at “big-data-based companies”

Cheap large dataset vs. expensive **labeled data**





Nơi để bạn phát triển tương lai, Nhà giá 883
Thạnh Lộc 41, Phường Thạnh Lộc, Quận 12, TP.HCM

882 triệu 999 nghìn



Duy nhất chỉ có 1 căn, có thể chỉ dành cho bạn,
Đường TX 25, Phường Thạnh Xuân, Quận 12, TP.HCM

880 triệu



Đừng mua nhà quận 3 nếu chưa xem tin này -
Nguyễn Thượng Hiền, Phường 4, Quận 3, TP.HCM

7100 triệu

Bán nhà mặt phố Kim Mã, vị trí đẹp nhất phố (đang cho thuê 180 tr/th), 80m2, 8T, MT 4,5m, giá 33 tỷ

Giá : 33 tỷ

Diện tích : 80 m²

Quận/huyện : Ba Đình, Hà Nội

16/11/2018

Tôi cần bán CH Park 1 - Times City 458 Minh Khai, HN, 120m2, 3PN, căn góc, thiết kế đẹp, 35 tr/m2

Giá : 4.2 tỷ

Diện tích : 120 m²

Quận/huyện : Hai Bà Trưng, Hà Nội

16/11/2018

Asiana Capella Trần Văn Kiểu, Q6, chỉ 1.6 tỷ sở hữu ngay CH, duplex, smark office. LH: 0901 090 228

Giá : 1.6 tỷ

Diện tích : 50 m²

Quận/huyện : Quận 6, Hồ Chí Minh

16/11/2018



Nơi để bạn phát triển tương lai, Nhà giá 883
Thạnh Lộc 41, Phường Thạnh Lộc, Quận 12, TP.HCM

Xin chào, bạn cần tìm bất động sản như thế nào?

À tôi muốn mua một căn nhà mặt tiền ở
Quận 10 để kinh doanh

Bán nhà mặt phố
tr/th), 80m², 8T,
Giá :
Diện tích :
Quận/huyện :

Đây là 5 kết quả tìm được theo yêu cầu
của bạn

Cần bán nhà mặt tiền đường Sư Vạn
Hạnh, Phường 13, Quận 10. - Diện
tích: 4m x 27m, TDT: 115m². vuông
vức, không lô giới, lề rộng 7m, sổ
hồng chính chủ. - Kết cấu: 1 trệt, 2
lầu. - Vị trí rất đắc địa gần ngay và 3
Tháng 2, gần khu Trung tâm thương
mại An Phong đang xây dựng.khu
kinh doanh buôn bán đa ngành nghề
sầm uất. Tân trang nhiều trung tâm

Tôi cần bán CH Pa
3PN, căn góc, thiế
Giá :
Diện tích :
Quận/huyện :

Asiana Capella Tr
duplex, smark off

Giá : 1.6 ty
Diện tích : 50 m²
Quận/huyện : Quận 6, Hồ Chí Minh

CÓ

KHÔNG

Tư vấn

16/11/2018

chỉ dành cho bạn,
Xuân, Quận 12, TP.HCM

chưa xem tin này -
g 4, Quận 3, TP.HCM



Cần tiền bán gấp đất mặt tiền đường tam bửu tự - xã đa phước - bình chánh - tp . hcm

- diện tích : 2 , 000 m 2 (hai sổ đỏ , mỗi sổ là 1 , 000 m 2)
- đất trồng cây lâu năm , thích hợp làm nhà vườn , ao cá , trồng cây cảnh .
- đường bê tông xe hơi vào tận nơi ,
- đất cách đường quốc lộ 50 khoảng 200 m
- đất nằm gần đường dẫn cao tốc long thành dầu dây .
- cách bến xe quận 8 khoảng 6 km .

Predict

addr_street (S)	addr_district (D)	addr_city (C)	addr_ward (W)	position (P)	area (A)	price (M)	transaction_type (T)	realestate_type (E)
legal (L)	potential (O)	surrounding (U)	surrounding_characteristics (I)	surrounding_name (N)	interior_floor (F)	interior_room (R)	orientation (H)	
project (J)	normal ()	Reset						

Cần tiền bán gấp đất mặt tiền đường tam bửu tự - xã đa phước - bình chánh - tp . hcm

- diện tích : 2 , 000 m 2 (hai sổ đỏ , mỗi sổ là 1 , 000 m 2)
- đất trồng cây lâu năm , thích hợp làm nhà vườn , ao cá , trồng cây cảnh .
- đường bê tông xe hơi vào tận nơi ,
- đất cách đường quốc lộ 50 khoảng 200 m
- đất nằm gần đường dẫn cao tốc long thành dầu dây .
- cách bến xe quận 8 khoảng 6 km .
- cách khu dân cư phong phú khoảng 3 km

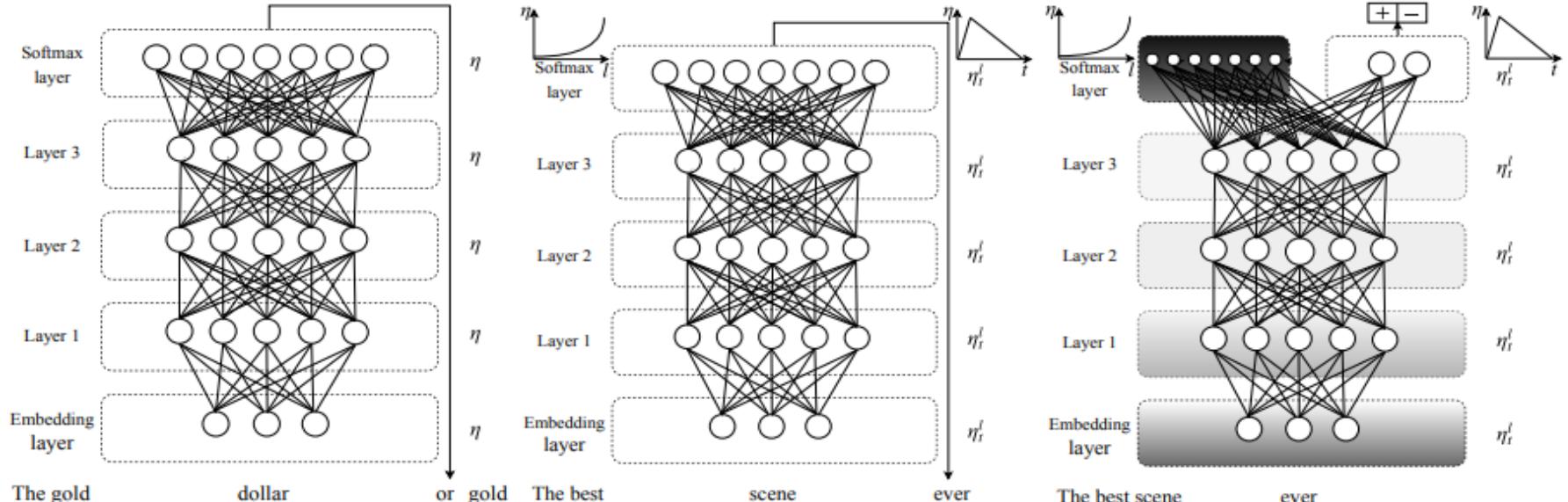
Neural-based Milestones for NLP

- 2001 - Neural language models
- 2008 - Multi-task learning
- 2013 - Word embeddings
- 2013 - Neural networks for NLP
- 2014 - Sequence-to-sequence models
- 2015 - Memory-based networks
- 2018 - Pretrained language models

**Character-based representations
Adversarial learning
Reinforcement learning**

**Conditional random fields
Bleu
Latent Dirichlet Allocation
Wikipedia exploitation**

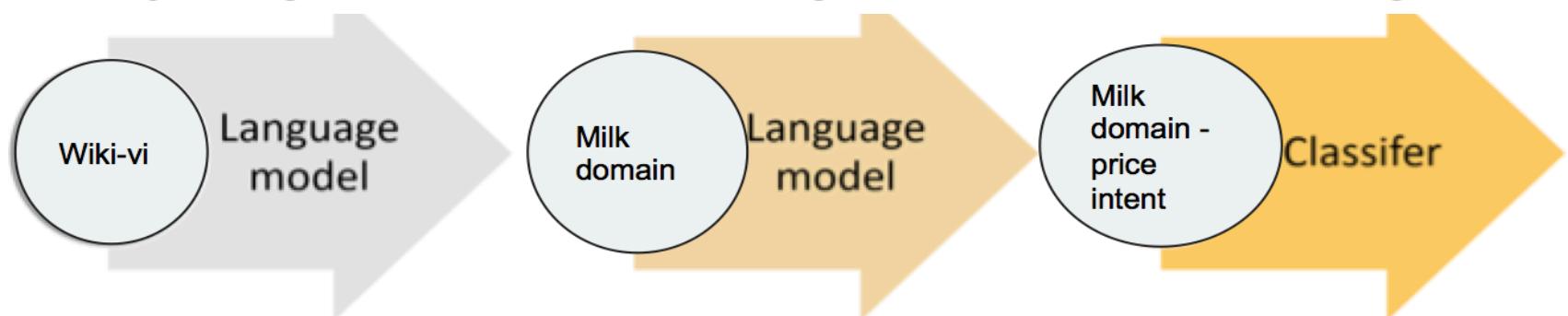
Pre-trained Neural Language Model



(a) LM pre-training

(b) LM fine-tuning

(c) Classifier fine-tuning



**Linh Nguyen**Comment on the wall **Bạch Hoàn**. Post: Cuối cùng thì dù có - Nov 14, 2018 08:19 AM GMT+07:00

0 0



Ôi lại hội sữa người, sữa thú huhu. Con tới tuổi đã lớn, thì uống sữa thú chứ biết sao. Sữa bắp, gạo, hạt... nguyên chất nào đủ hàm lượng canxi, mà mua loại hữu cơ có bổ sung canxi thì đắt gấp 3 lần sữa thú... Sữa thủy phân đắt và khó tìm, thường chỉ dùng cho trẻ bị chàm và dị ứng thôi, trẻ khỏe mạnh chẳng tội gì phải uống.



3

**Ngọc Khánh**Comment on the wall **Bạch Hoàn**. Post: Cuối cùng thì dù có - Nov 14, 2018 12:53 AM GMT+07:00

0 0

Các cháu có sữa uống là tốt rồi lại còn rẻ hơn thị trường nữa là có lợi cho túi tiền phụ huynh.
Còn tình trạng chung hoặc những nghi vấn về chất lượng sữa ở VN lại là một vấn đề rộng hơn



3

**Kim Thuan**Comment on the wall **Bạch Hoàn**. Post: Cuối cùng thì dù có - Nov 13, 2018 11:31 PM GMT+07:00

0 0

Cuối tuần này mới có kết quả chính thức mà, nhiều khi đấu thầu giá rẻ hơn chưa chắc đã thắng vì còn tiêu chí khác nữa. Hy vọng vụ này xong phẳng và Vinamilk thắng thầu.



3





Ngoc Linh

Comment on Hội Sữa Công Thức Sử Dụng Đúng Cách (Đọc Kỹ Nội Quy Ở Bài Ghim)'s group. Post: Bình thường e mua - Nov 08, 2018 08:16 PM GMT+07:00

都有自己家的牌子 , mua ở cửa hàng chính thức của th giá là nt đấy chị ơi . Cũng chả đắt đâu , chênh 20k mà dc cái đảm bảo cty , 2 nữa là hay có km nữa . Con e dùng vinamilk e cũng chỉ mua ở giắc mơ sữa việt thôi nè . Yên tâm hơn



2

1



Hoàng Diệu

0 0

Comment on Hội Sữa Công Thức Sử Dụng Đúng Cách (Đọc Kỹ Nội Quy Ở Bài Ghim)'s group. Post: Bình thường e mua - Nov 08, 2018 08:03 PM GMT+07:00

Giá chính xác dao động từ 215 đến 218. Bạn mua 240k là quá đắt. Còn khác nhau là do 2 nhà phân phối ak. Ko phải sữa giả đâu. Mình bán nên mình biết ak. Mọi yên tâm nha.



2

1



Hoan Nguyễn

0 0

Comment on Hội Sữa Công Thức Sử Dụng Đúng Cách (Đọc Kỹ Nội Quy Ở Bài Ghim)'s group. Post: Bình thường e mua - Nov 08, 2018 07:55 PM GMT+07:00

Tất cả các hãng sữa đều có 2 bao bì, do có 2 cty cung cấp bao bì đó ạ. Mn ko cần hoang mang đâu. Mua ở đại lý chuyên sữa nó sẽ bán đúng giá niêm yết, còn mua ở tạp hóa họ ăn số lượng nên bán rẻ hơn. Mua sữa lúc nào cũng gặp th như vậy đấy a.



2

1



Intuition

Sua pediasure giá bnhju vay

Intuition

Sua pediasure giá bnhju vay

Cái này giá **bao nhiêu** vậy
Giá nó bao nhiêu

Bao nhiêu một hộp sữa vậy
Giá hộp sữa **bao nhiêu**
Giá hộp này **bnhju**
Bnhju thì bán được

Case Studies

Language model
Intent classification
NER tagging

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Devlin et al., 2018 (Google AI Language)

NLP course

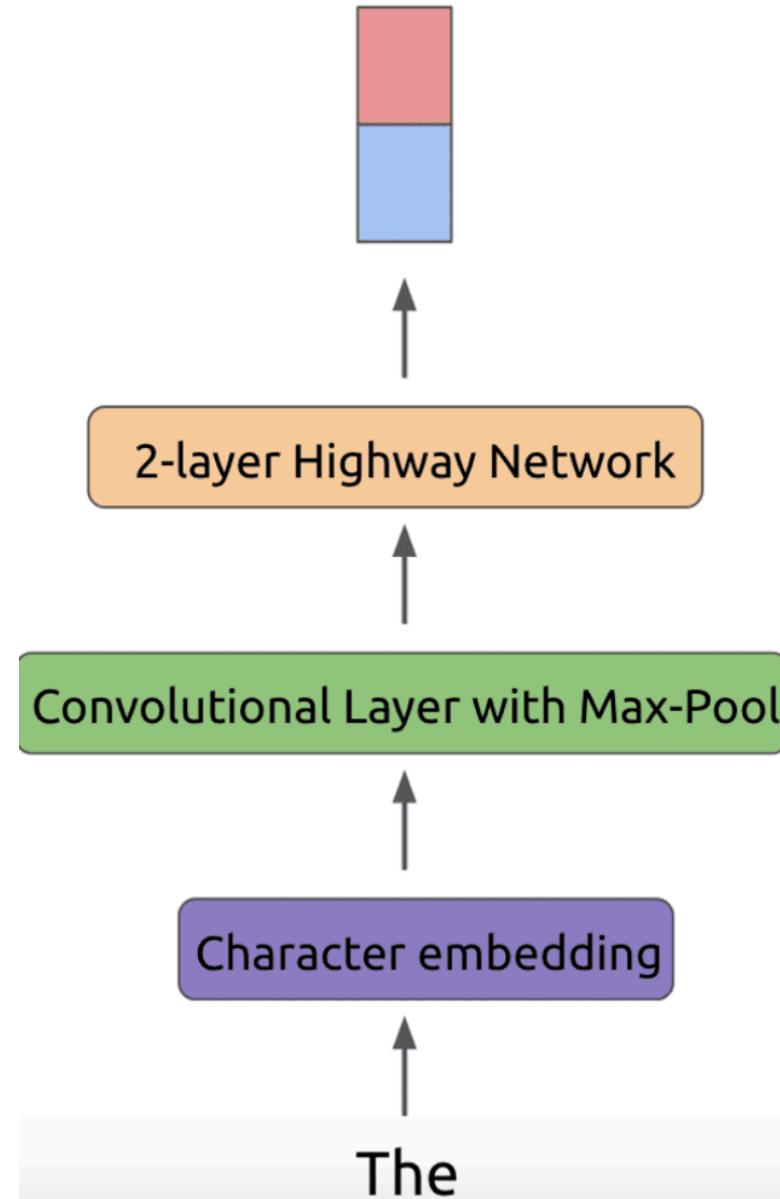
Outline

- Research context
- Main ideas
- BERT
- Experiments
- Conclusions

Research context

- Language model pre-training has been used to improve many NLP tasks
 - ELMo (Peters et al., 2018)
 - OpenAI GPT (Radford et al., 2018)
 - ULMFit (Howard and Ruder, 2018)
- Two existing strategies for applying pre-trained language representations to downstream tasks
 - *Feature-based*: include pre-trained representations as additional features (e.g., ELMo)
 - *Fine-tuning*: introduce task-specific parameters and fine-tune the pre-trained parameters (e.g., OpenAI GPT, ULMFit)

ELMO

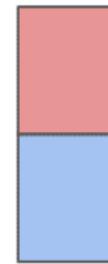
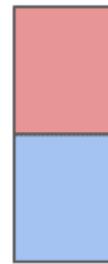
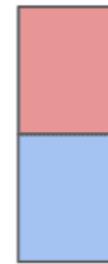
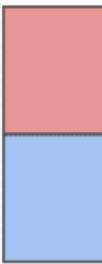
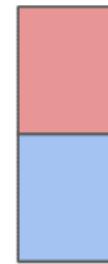
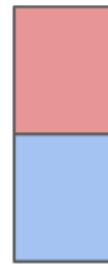
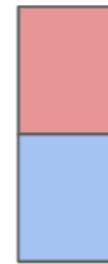
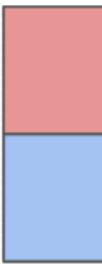


cat

is

happy

EOS



The

cat

is

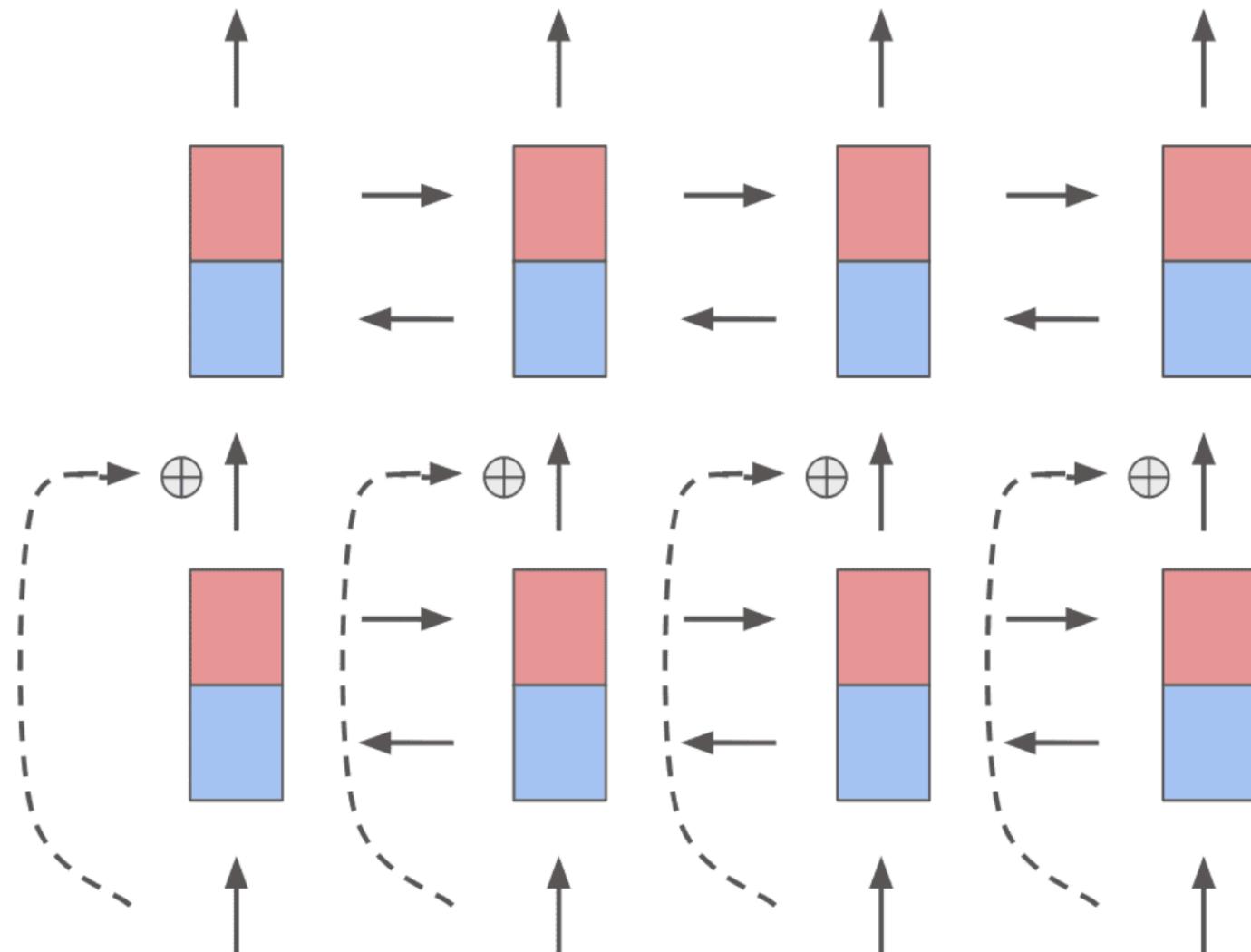
happy

cat

is

happy

EOS



The

cat

is

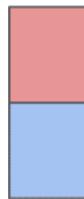
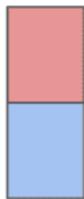
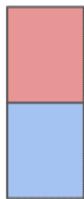
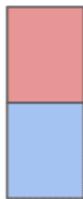
happy

cat

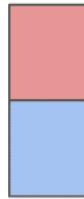
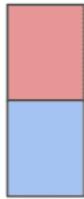
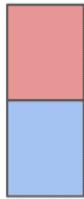
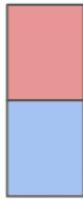
is

happy

EOS



$h_{2, happy}$



The

cat

is

happy

x_{happy}

$h_{1, happy}$

\rightarrow

f

ELMo
happy

ULMFiT

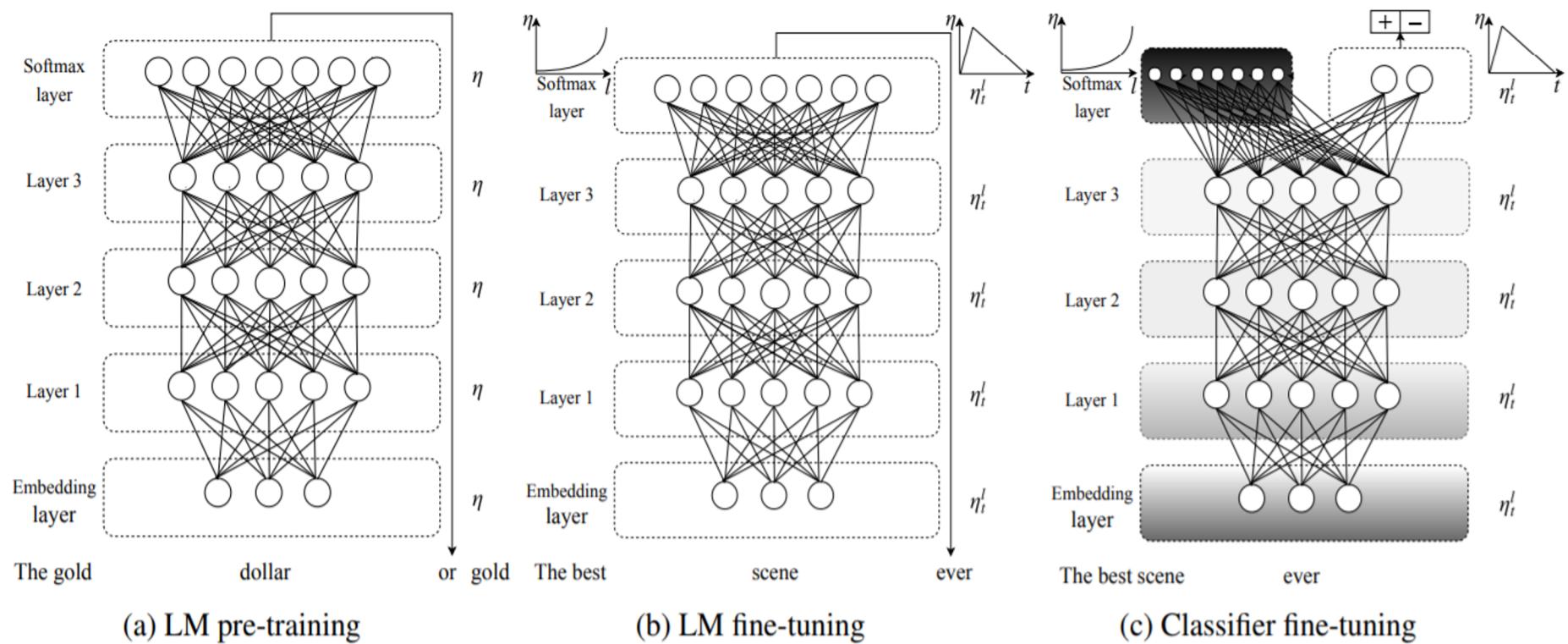


Figure 1: ULMFiT consists of three stages: a) The LM is trained on a general-domain corpus to capture general features of the language in different layers. b) The full LM is fine-tuned on target task data using discriminative fine-tuning ('*Discr*') and slanted triangular learning rates (STLR) to learn task-specific features. c) The classifier is fine-tuned on the target task using gradual unfreezing, '*Discr*', and STLR to preserve low-level representations and adapt high-level ones (shaded: unfreezing stages; black: frozen).

Limitations of current techniques

- Language models in pre-training are unidirectional, they restrict the power of the pre-trained representations
 - OpenAI GPT used left-to-right architecture
 - ELMo concatenates *forward* and *backward* language models
- Solution BERT: Bidirectional Encoder Representations from Transformers

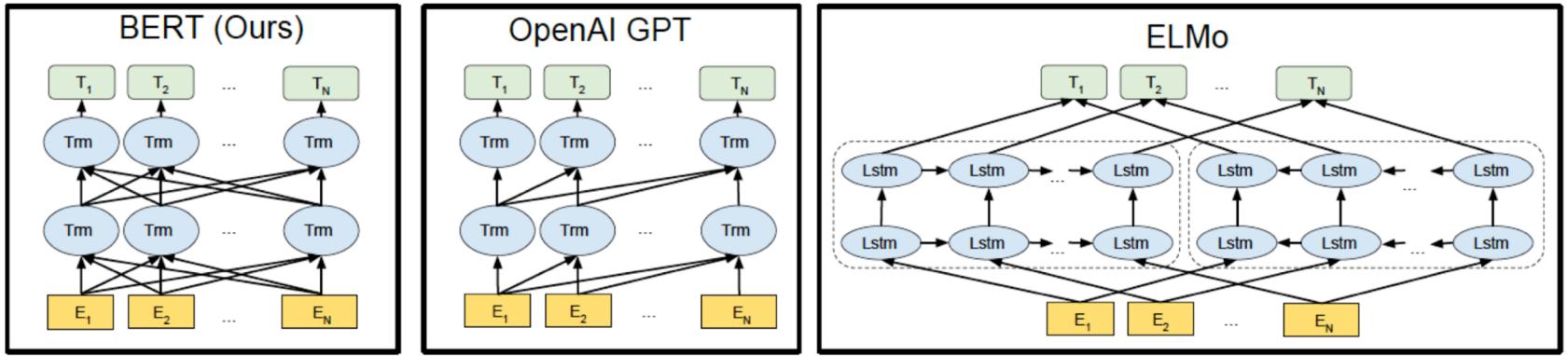
BERT: Bidirectional Encoder Representations from Transformers

- Main ideas
 - Propose a new pre-training objective so that a deep bidirectional Transformer can be trained
 - The “**masked language model**” (**MLM**): the objective is to predict the original word of a masked word based only on its context
 - “**Next sentence prediction**”
- Merits of BERT
 - Just fine-tune BERT model for specific tasks to achieve state-of-the-art performance
 - BERT advances the state-of-the-art for eleven NLP tasks

BERT: Bidirectional Encoder Representations from Transformers

Model architecture

- BERT's model architecture is a multi-layer **bidirectional Transformer encoder**
 - (Vaswani et al., 2017) "Attention is all you need"
- Two models with different sizes were investigated
 - BERT_{BASE}: L=12, H=768, A=12, Total Parameters=110M
 - (L : number of layers (Transformer blocks), H is the hidden size, A : the number of self-attention heads)
 - BERT_{LARGE}: L=24, H=1024, A=16, Total Parameters=340M

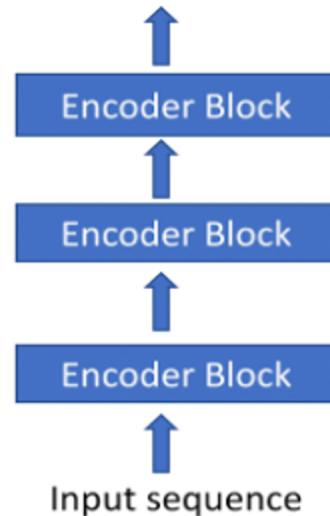


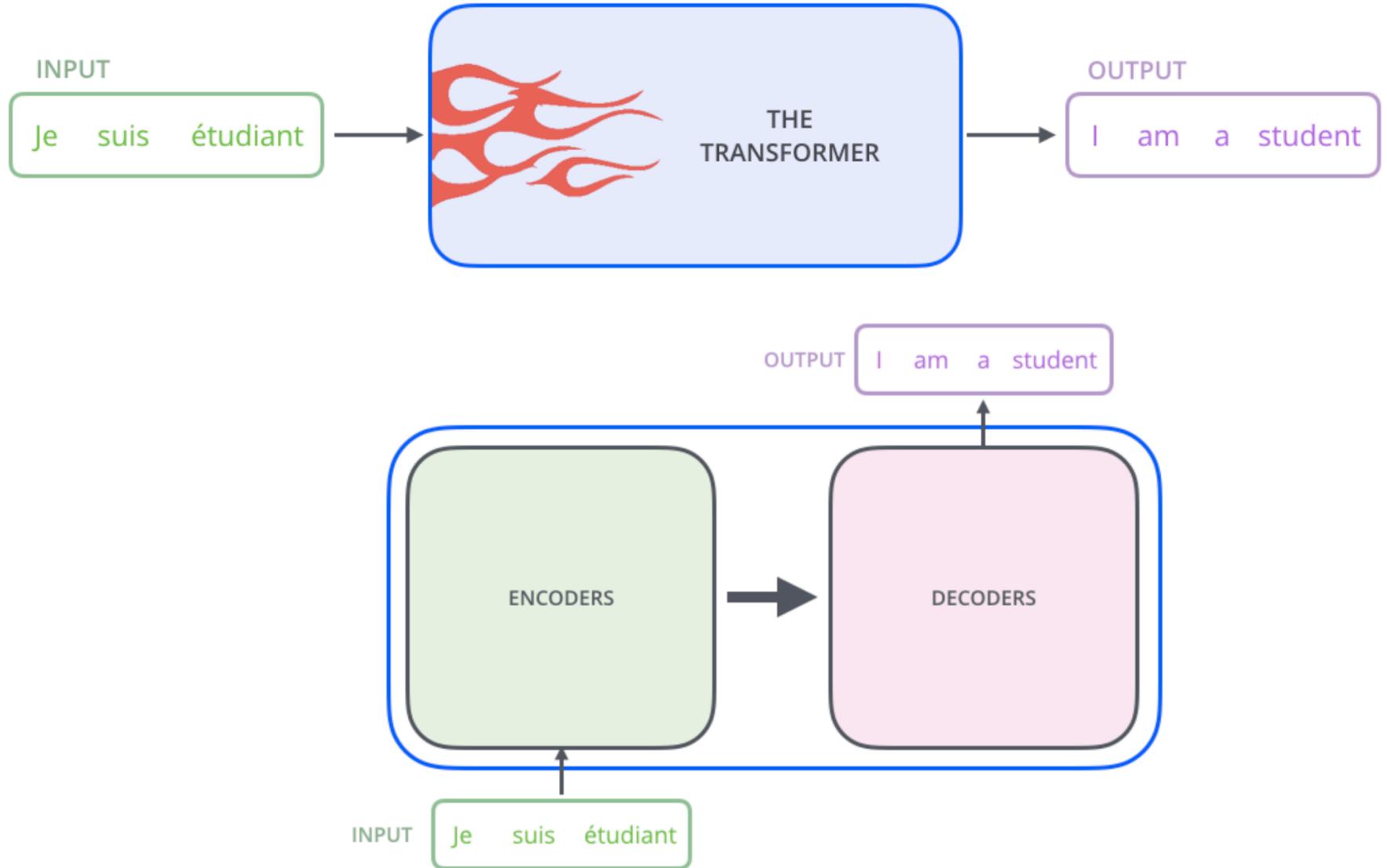
Differences in pre-training model architectures:
BERT, OpenAI GPT, and ELMo

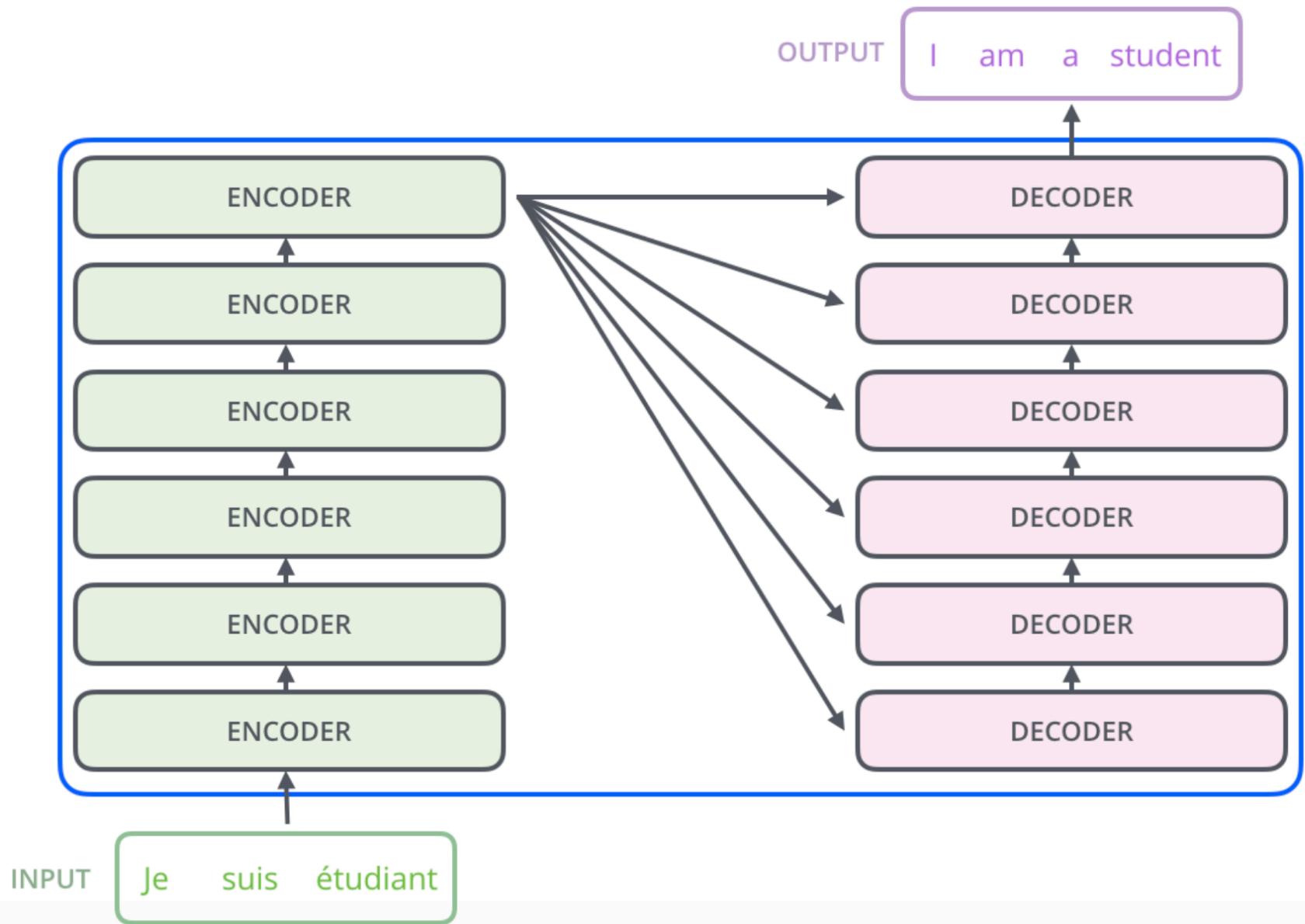
Transformer Encoders

Vaswani et al. (2017) *Attention is all you need*

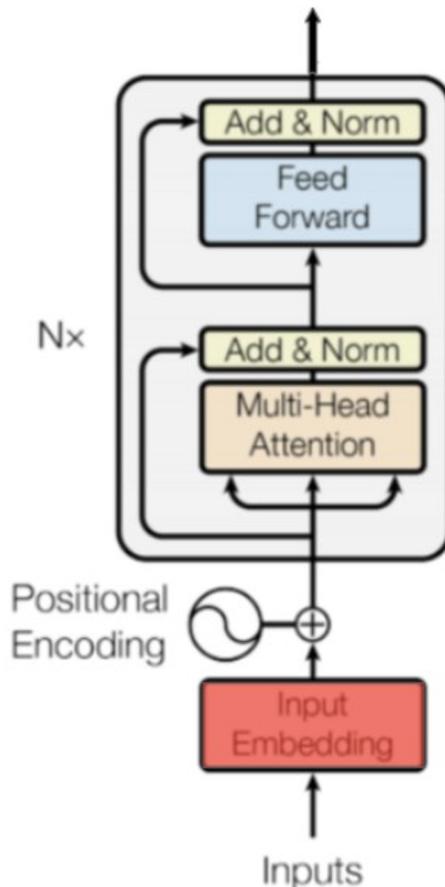
- Transformer is an attention-based architecture for NLP
- Transformer composed of two parts: **Encoding** component and **Decoding** component
- BERT is a multi-layer bidirectional Transformer **encoder**







Inside an Encoder Block

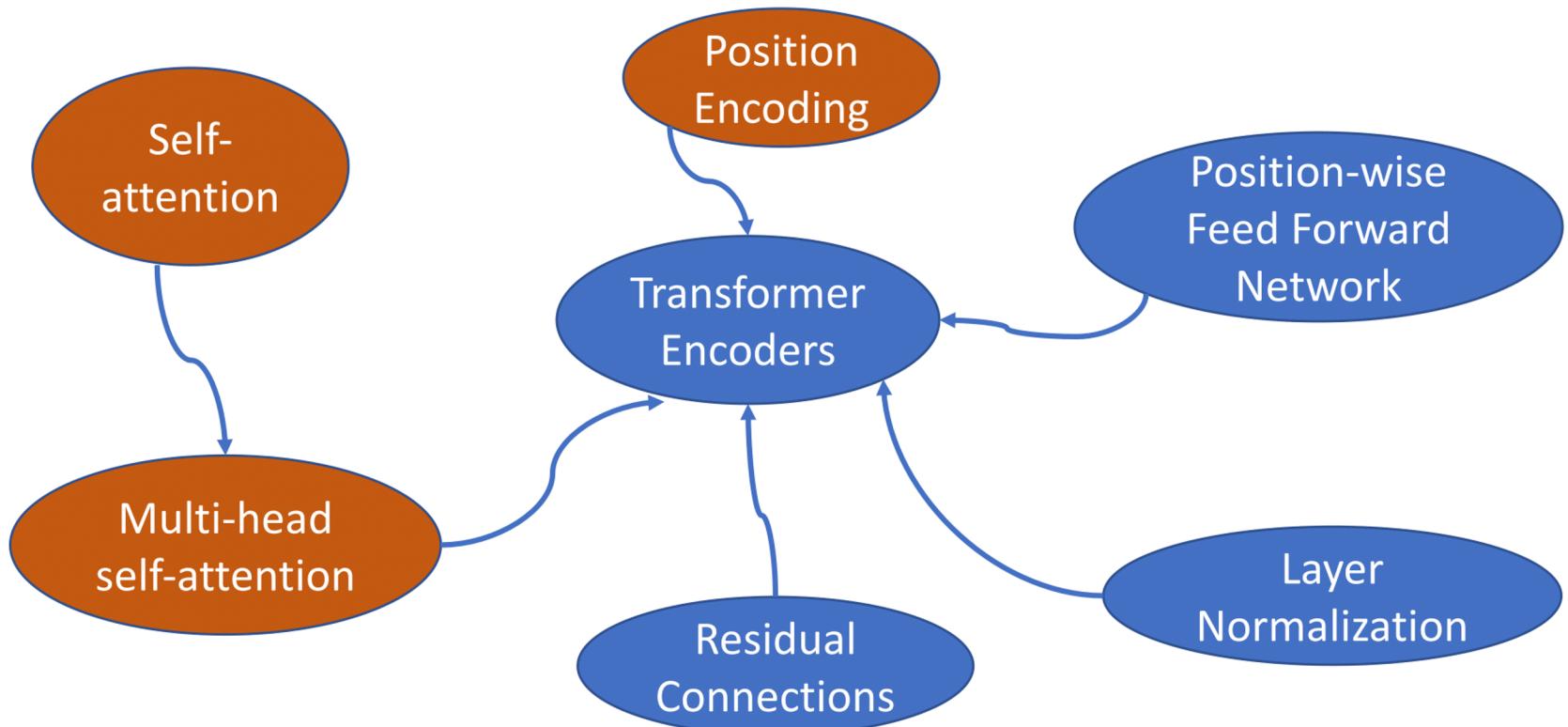


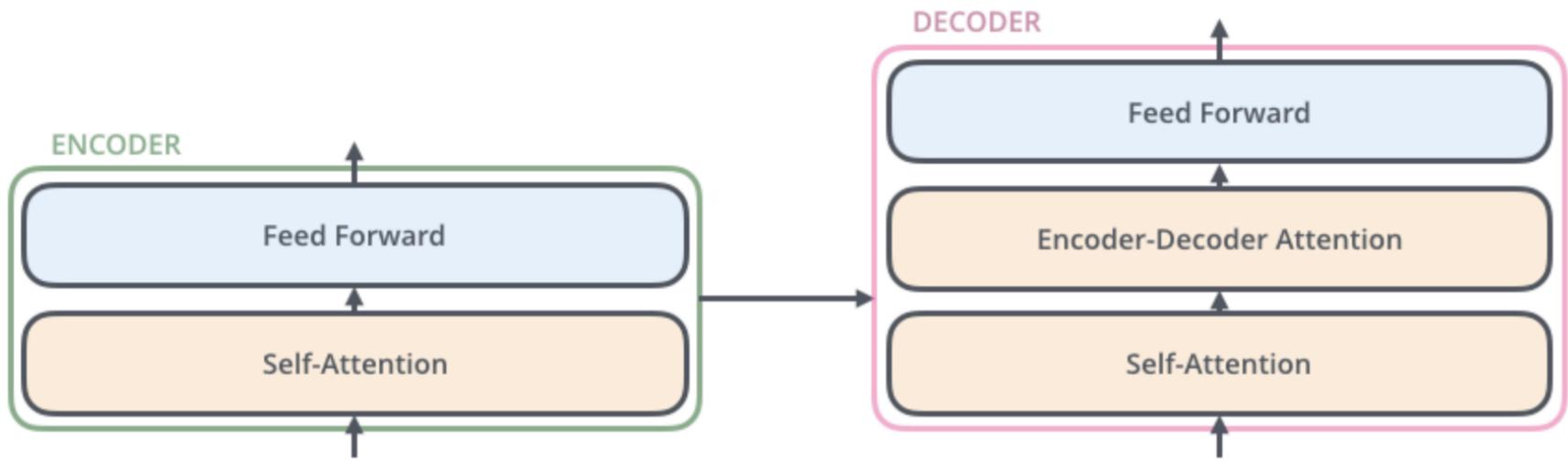
In BERT experiments, the number of blocks N was chosen to be 12 and 24.

Blocks do not share weights with each other

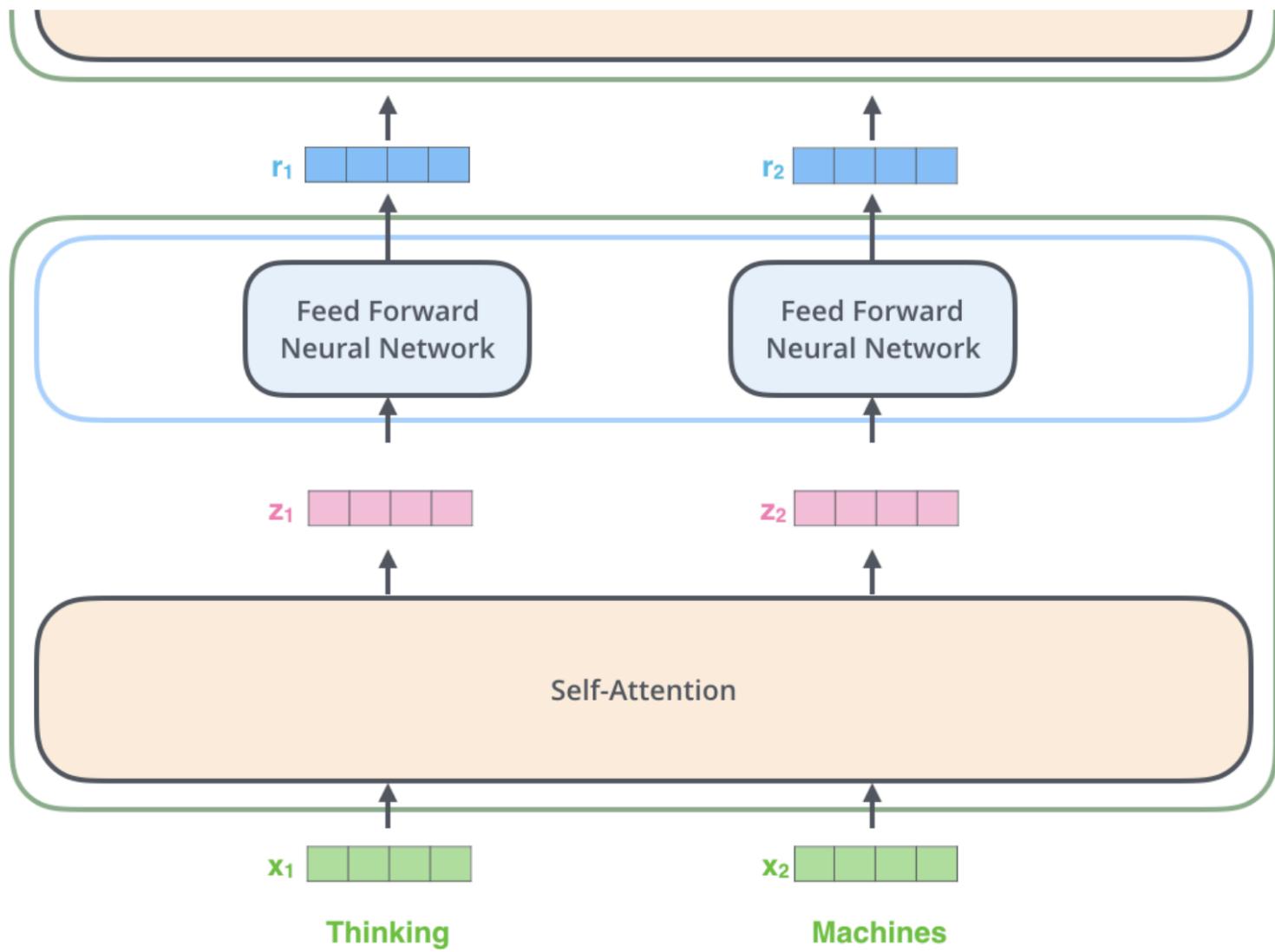
Source: <https://medium.com/dissecting-bert/dissecting-bert-part-1-d3c3d495cdb3>

Transformer Encoders: Key Concepts





ENCODER #2



Self-Attention

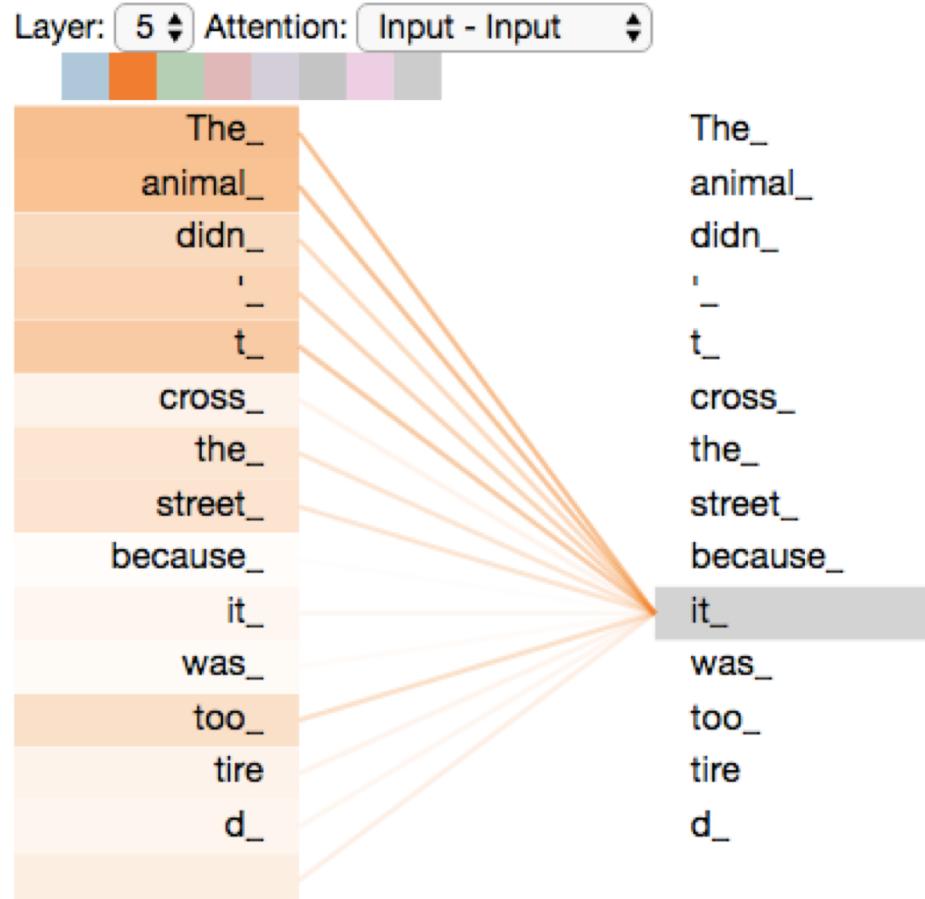
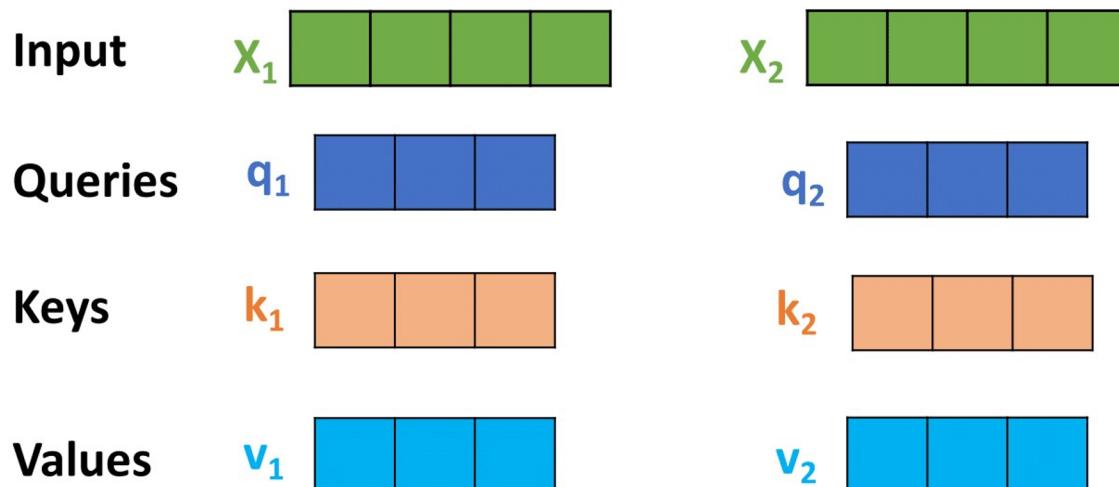


Image source: <https://jalammar.github.io/illustrated-transformer/>

Self-Attention in Detail

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



Use matrices W^Q , W^K and W^V to project input into query, key and value vectors

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

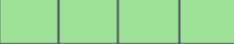
d_k is the dimension of key vectors

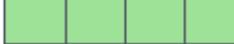
Input

Thinking

Machines

Embedding

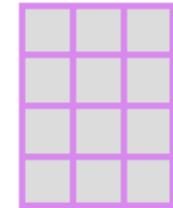
X_1 

X_2 

Queries

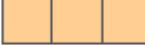
q_1 

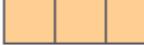
q_2 

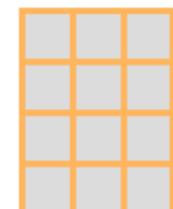


WQ

Keys

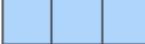
k_1 

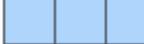
k_2 

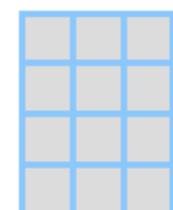


WK

Values

v_1 

v_2 

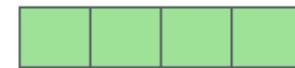


WV

Input

Thinking

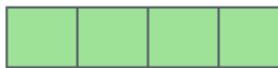
x_1



Embedding

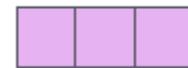
Machines

x_2

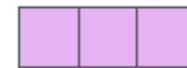


Queries

q_1

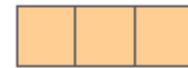


q_2

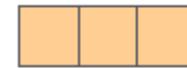


Keys

k_1

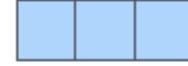


k_2

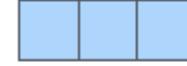


Values

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

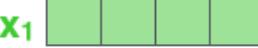
0.88

0.12

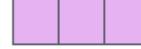
Input

Thinking

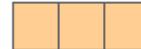
Embedding

x_1 

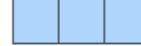
Queries

q_1 

Keys

k_1 

Values

v_1 

Score

$$q_1 \cdot k_1 = 112$$

$$q_1 \cdot k_2 = 96$$

Divide by 8 ($\sqrt{d_k}$)

$$14$$

$$12$$

Softmax

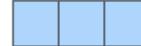
$$0.88$$

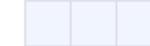
$$0.12$$

Softmax

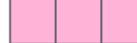
X

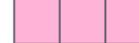
Value

v_1 

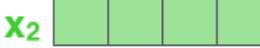
v_2 

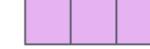
Sum

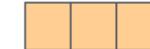
z_1 

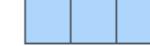
z_2 

Machines

x_2 

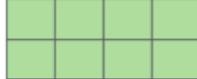
q_2 

k_2 

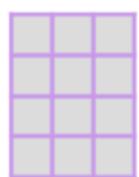
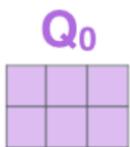
v_2 

X

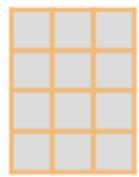
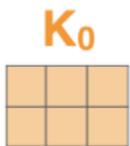
Thinking
Machines



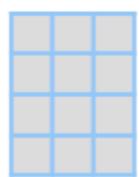
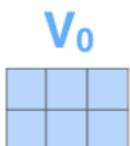
ATTENTION HEAD #0



W_0^Q

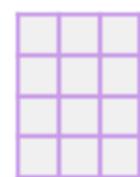
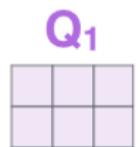


W_0^K

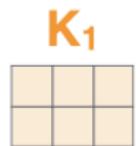


W_0^V

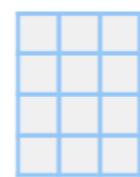
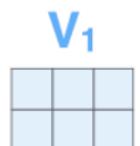
ATTENTION HEAD #1



W_1^Q

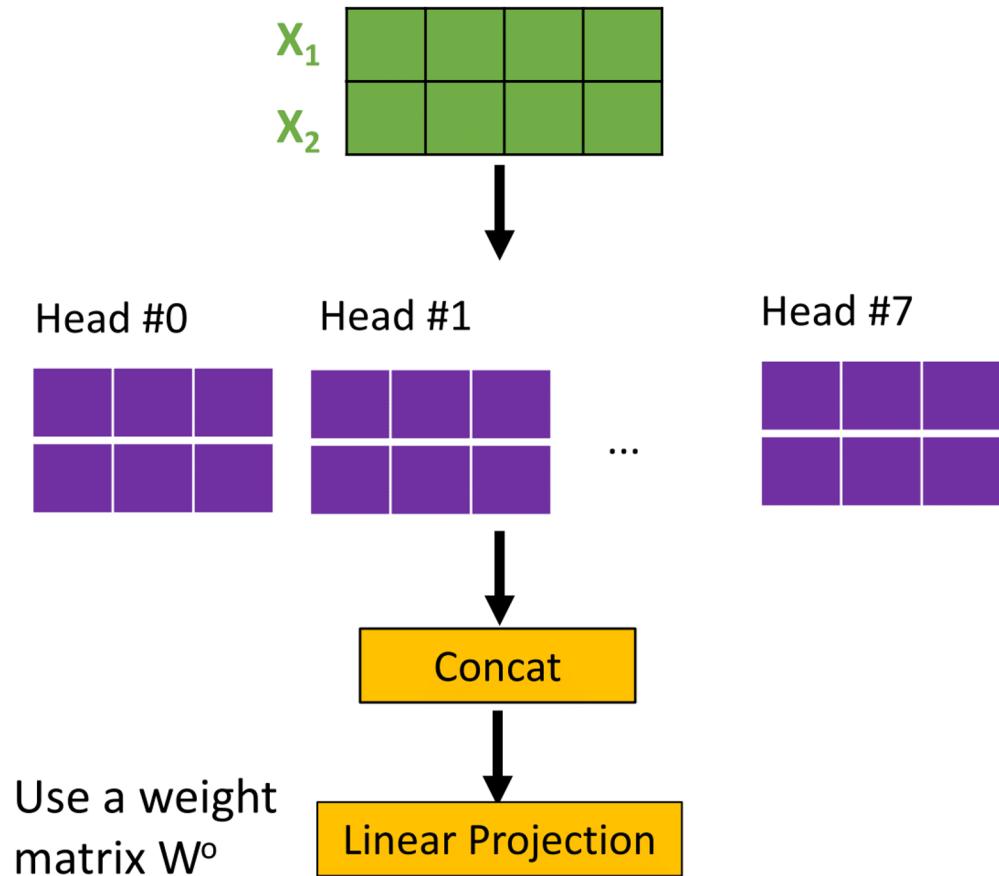
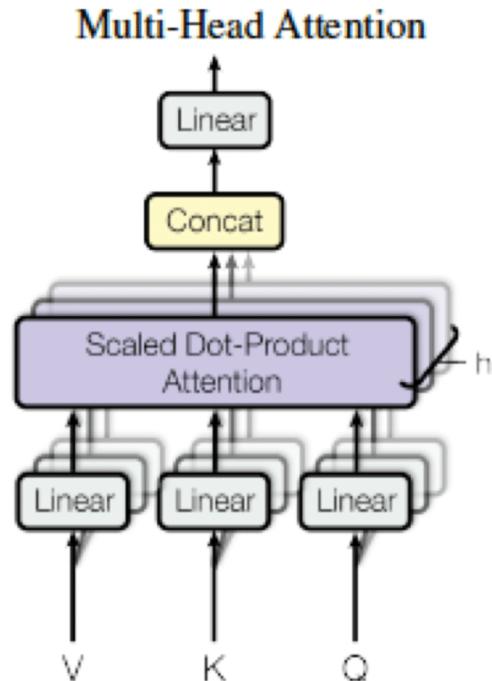


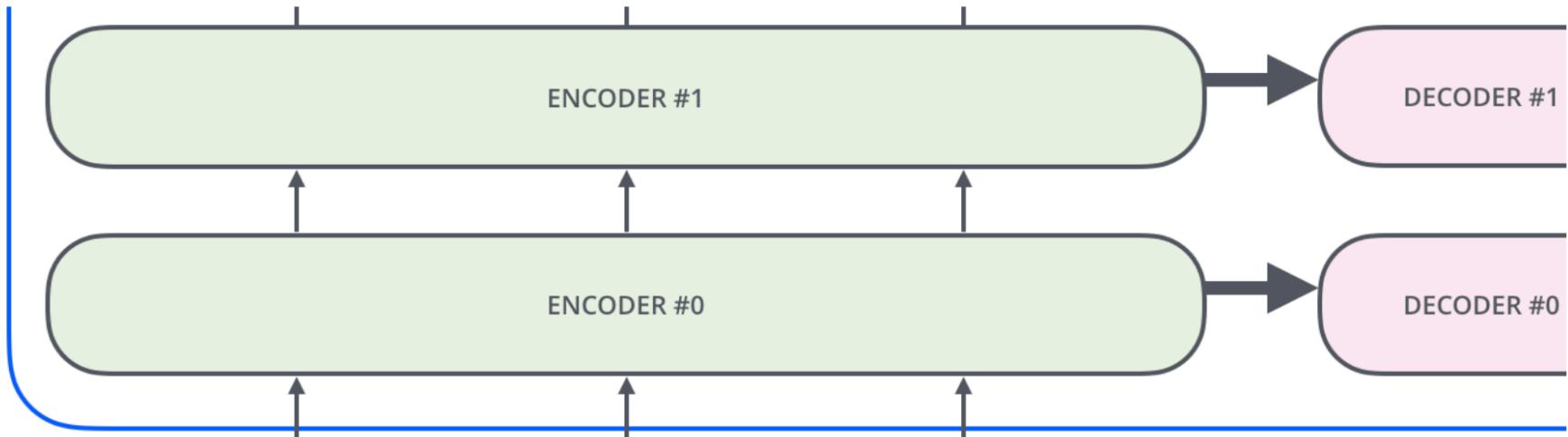
W_1^K



W_1^V

Multi-Head Attention





EMBEDDING
WITH TIME
SIGNAL

$$\mathbf{x}_1 \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}}$$

$$\mathbf{x}_2 \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}}$$

$$\mathbf{x}_3 \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}} \quad \boxed{\text{yellow}}$$

POSITIONAL
ENCODING

$$\mathbf{t}_1 \quad \boxed{\text{orange}} \quad \boxed{\text{orange}} \quad \boxed{\text{orange}}$$

$$\mathbf{t}_2 \quad \boxed{\text{orange}} \quad \boxed{\text{orange}} \quad \boxed{\text{orange}}$$

$$\mathbf{t}_3 \quad \boxed{\text{orange}} \quad \boxed{\text{orange}} \quad \boxed{\text{orange}}$$

EMBEDDINGS

$$\mathbf{x}_1 \quad \boxed{\text{green}} \quad \boxed{\text{green}} \quad \boxed{\text{green}}$$

$$\mathbf{x}_2 \quad \boxed{\text{green}} \quad \boxed{\text{green}} \quad \boxed{\text{green}}$$

$$\mathbf{x}_3 \quad \boxed{\text{green}} \quad \boxed{\text{green}} \quad \boxed{\text{green}}$$

INPUT

Je

suis

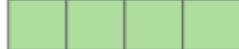
étudiant

POSITIONAL ENCODING

0	0	1	1
---	---	---	---

+

EMBEDDINGS

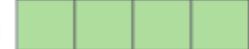
x_1 

0.84	0.0001	0.54	1
------	--------	------	---

+

0.91	0.0002	-0.42	1
------	--------	-------	---

+

x_2 

INPUT

Je

suis

étudiant

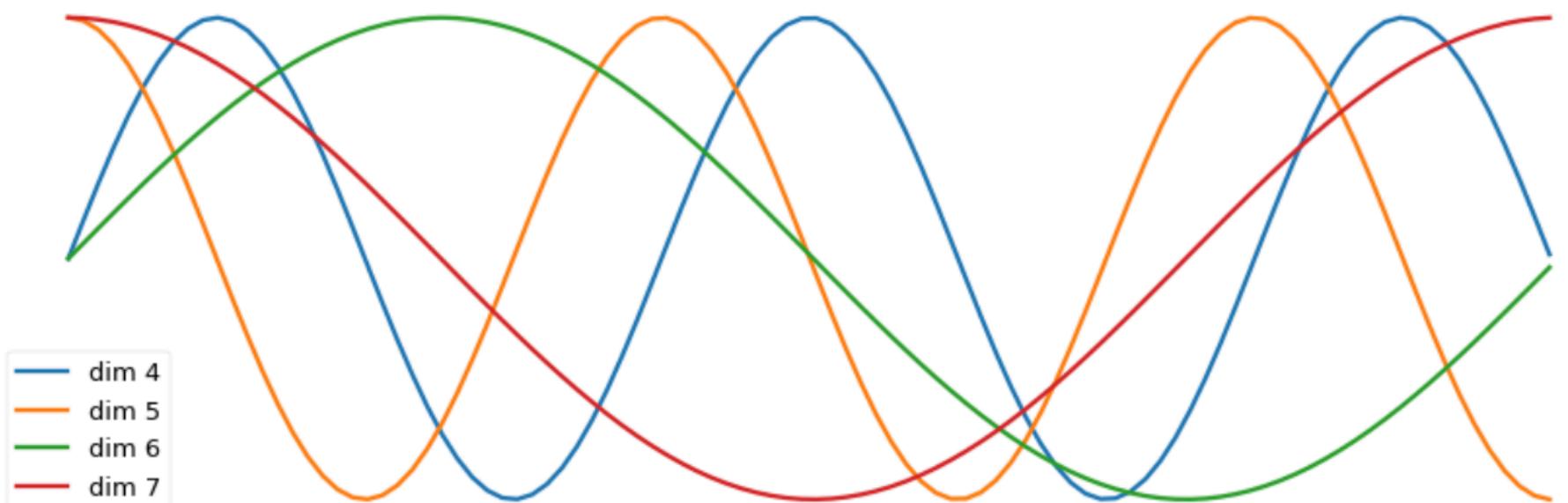
Position Encoding

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

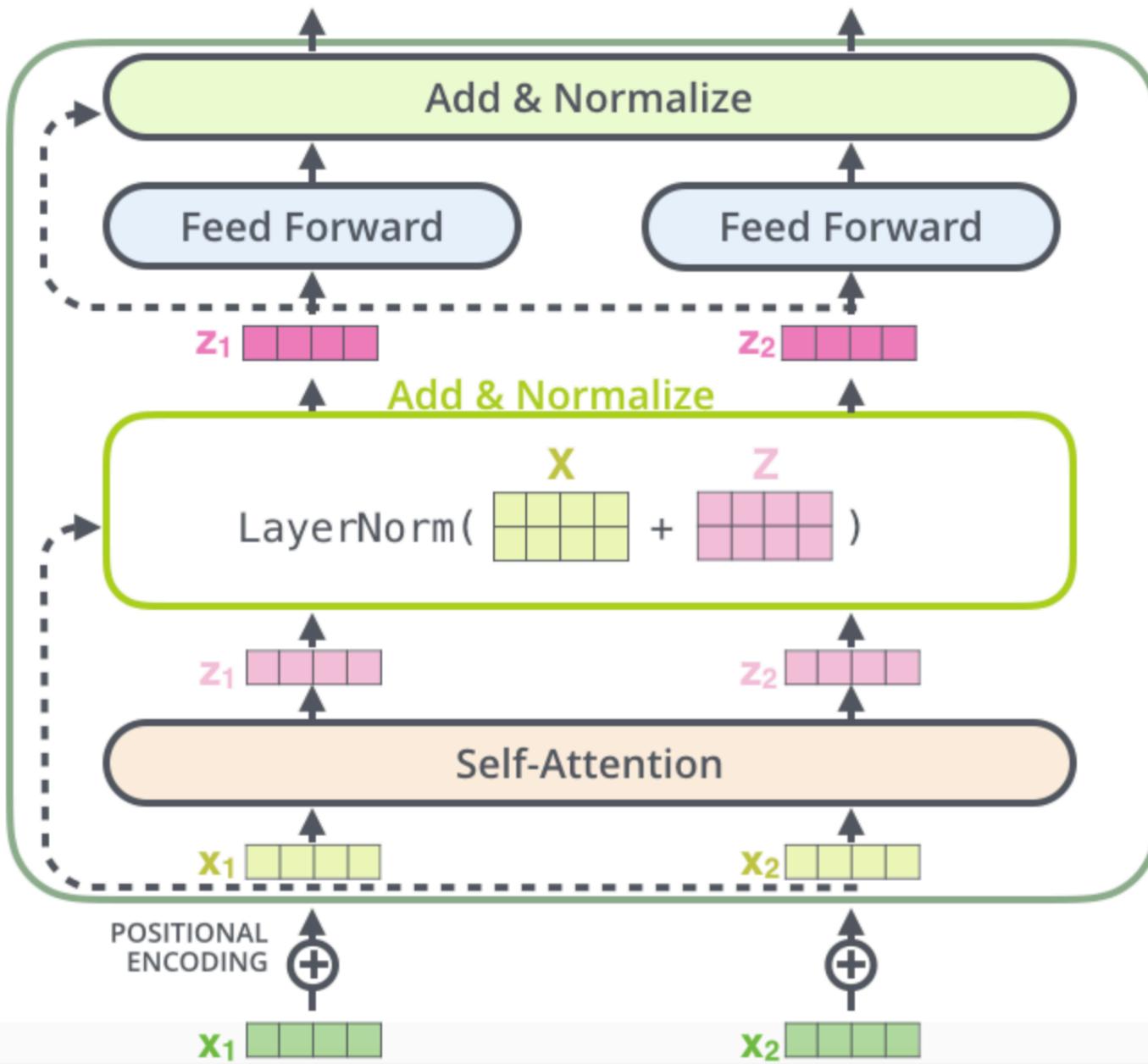
$$PE_{(pos,2i)} = \sin\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

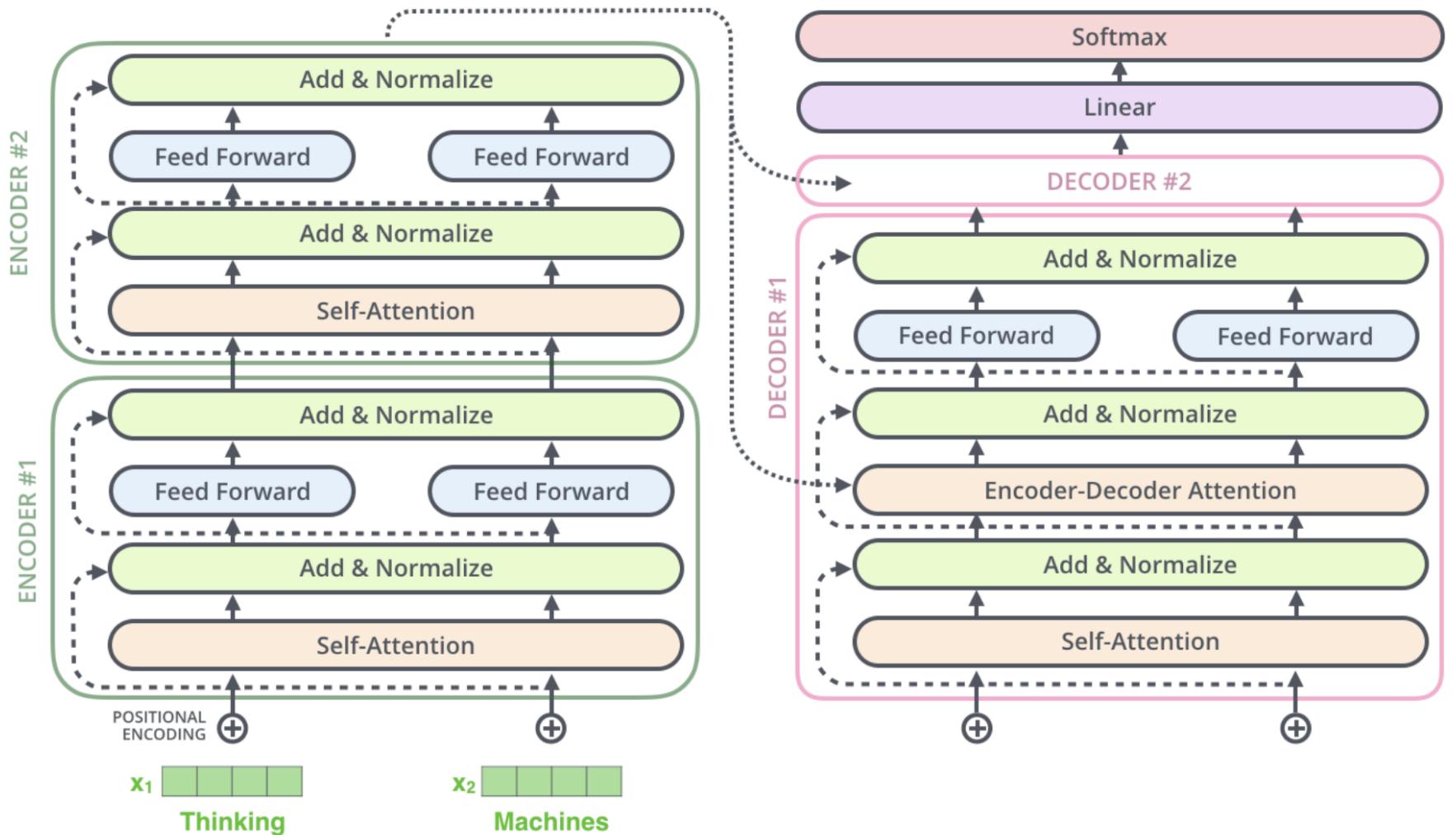
$$PE_{(pos,2i+1)} = \cos\left(\frac{pos}{10000^{\frac{2i}{d_{model}}}}\right)$$

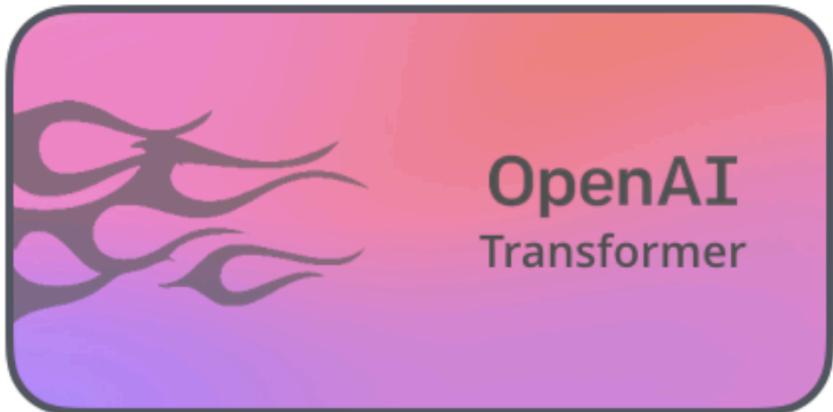
- pos is the position and i is the dimension



ENCODER #1

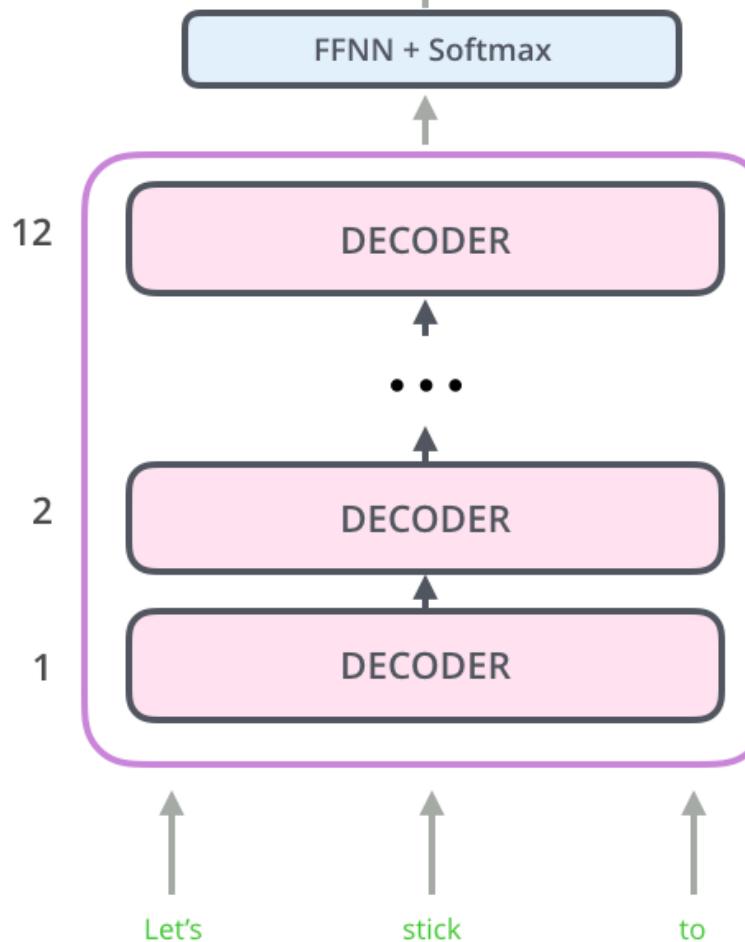




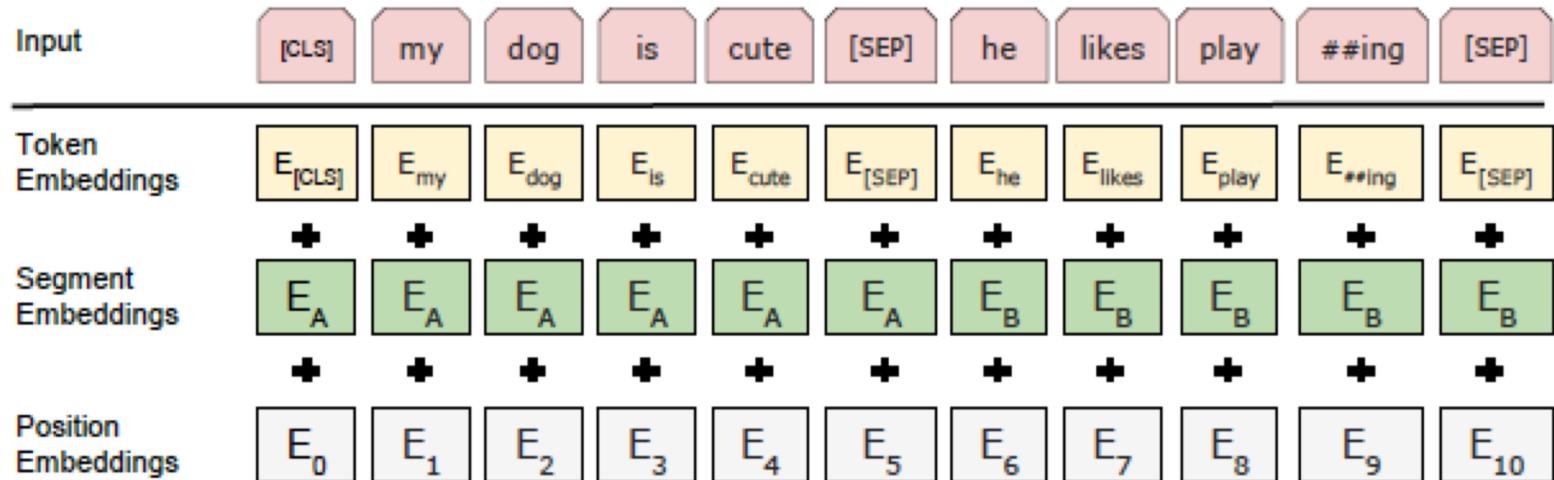


Possible classes:
All English words

0.1%	Aardvark
...	...
10%	Improvisation
...	...
0%	Zyzyva



Input Representation

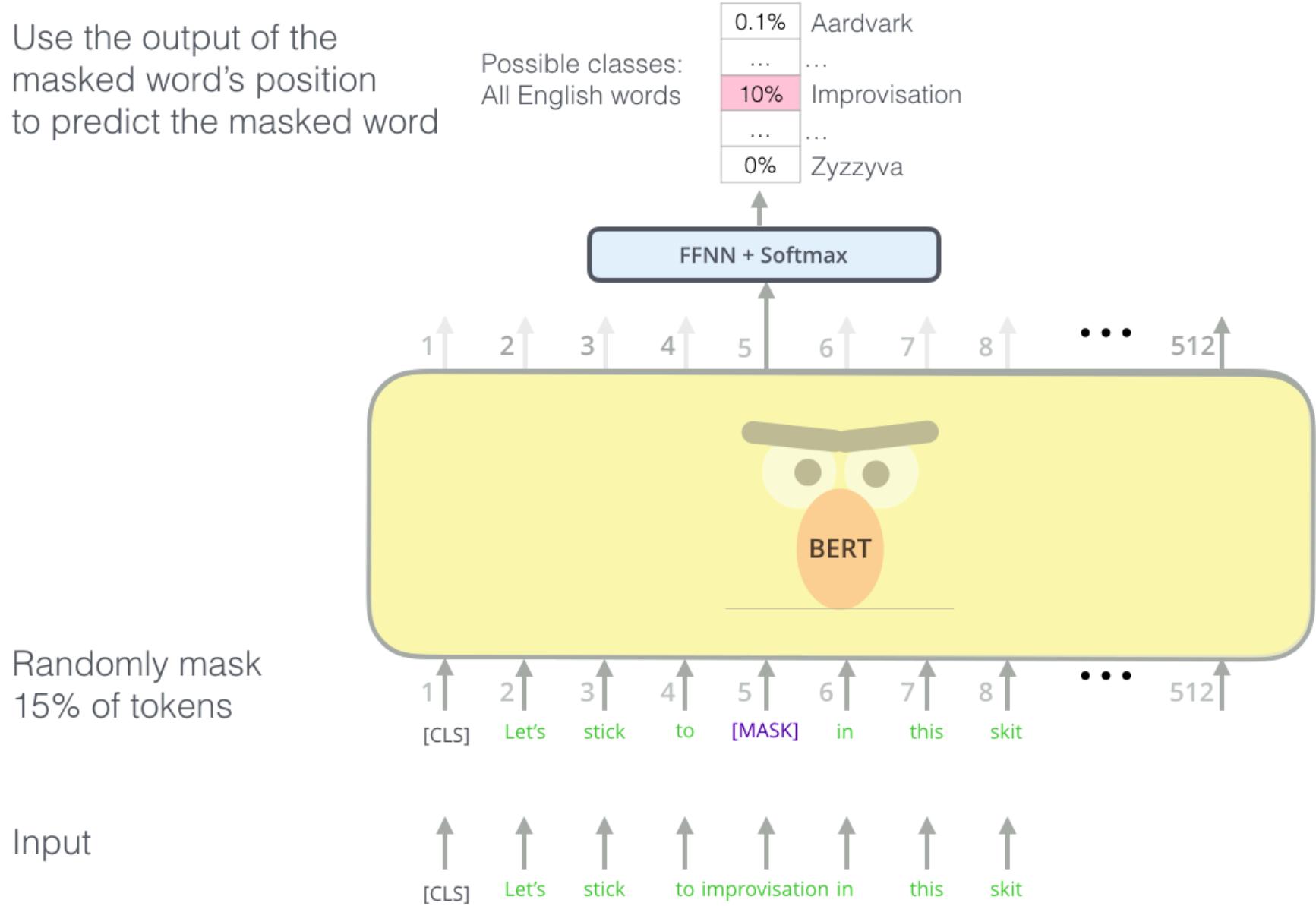


- Token Embeddings: Use pretrained WordPiece embeddings
- Position Embeddings: Use learned Position Embeddings
 - Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N. Dauphin. Convolutionalsequence to sequence learning. arXiv preprint arXiv:1705.03122v2, 2017.
- Added sentence embedding to every tokens of each sentence
- Use [CLS] for the classification tasks
- Separate sentences by using a special token [SEP]

Task#1: Masked LM

- 15% of the words are masked at random
 - and the task is to predict the masked words based on its left and right context
- Not all tokens were masked in the same way (example sentence "My dog is hairy")
 - 80% were replaced by the <MASK> token: "My dog is <MASK>"
 - 10% were replaced by a random token: "My dog is apple"
 - 10% were left intact: "My dog is hairy"

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



Task#2: Next Sentence Prediction

- Motivation
 - Many downstream tasks are based on understanding the relationship between two text sentences
 - Question Answering (QA) and Natural Language Inference (NLI)
 - Language modeling does not directly capture that relationship
- The task is pre-training binarized next sentence prediction task

Input = [CLS] the man went to [MASK] store [SEP] he bought a gallon
[MASK] milk [SEP]

Label = isNext

Input = [CLS] the man [MASK] to the store [SEP] penguin [MASK] are
flight ##less birds [SEP]

Label = NotNext

Predict likelihood
that sentence B
belongs after
sentence A



FFNN + Softmax

1

2

3

4

5

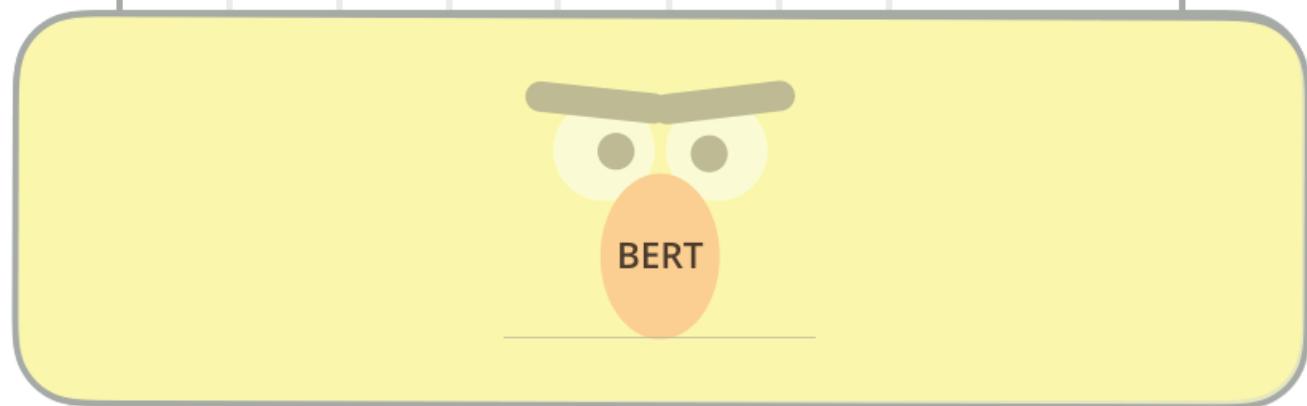
6

7

8

...

512



Tokenized
Input

1
[CLS]

2
the

3
man

4
[MASK]

5
to

6
the

7
store

8
[SEP]

...

512

Input

[CLS] the man [MASK] to the store [SEP] penguin [MASK] are flightless birds [SEP]
Sentence A Sentence B

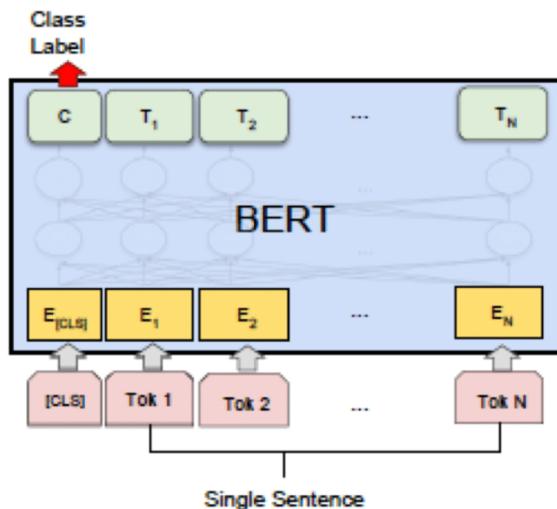
Pre-training procedure

- Training data: BooksCorpus (800M words) + English Wikipedia (2,500M words)
- To generate each training input sequences: sample two spans of text (A and B) from the corpus
 - The combined length is ≤ 500 tokens
 - 50% B is the actual next sentence that follows A and 50% of the time it is a random sentence from the corpus
- The training loss is the sum of the mean masked LM likelihood and the mean next sentence prediction likelihood

Fine-tuning procedure

- For sequence-level classification task
 - Obtain the representation of the input sequence by using the final hidden state (hidden state at the position of the special token [CLS]) $C \in R^H$
 - Just add a classification layer and use softmax to calculate label probabilities. Parameters $W \in R^{K \times H}$

$$P = \text{softmax}(CW^T)$$



Fine-tuning procedure

- For sequence-level classification task
 - All of the parameters of BERT and W are fine-tuned jointly
- Most model hyperparameters are the same as in pre-training
 - except the batch size, learning rate, and number of training epochs

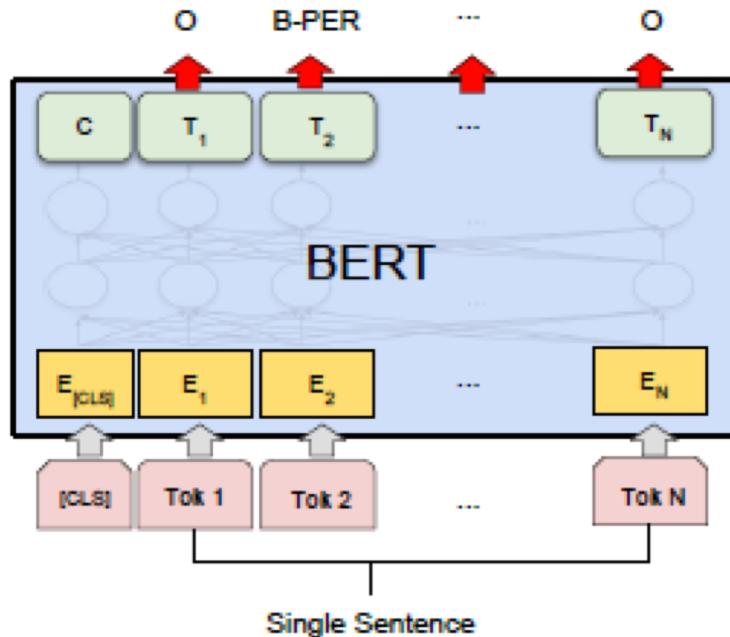
Fine-tuning procedure

- Token tagging task (e.g., Named Entity Recognition)
 - Feed the final hidden representation $T_i \in R^H$ for each token i into a classification layer for the tagset (NER label set)
 - To make the task compatible with WordPiece tokenization

Jim	Hen	##son	was	a	puppet	##eer
I-PER	I-PER	X	O	O	O	X

- Predict the tag for the first sub-token of a word
- No prediction is made for X

Fine-tuning procedure



Single Sentence Tagging Tasks: CoNLL-2003 NER

Fine-tuning procedure

- Span-level task: SQuAD v1.1
 - Input Question:
Where do water droplets collide with ice crystals to form precipitation?
 - Input Paragraph:
.... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud.
 - Output Answer:
within a cloud

Fine-tuning procedure

- Span-level task: SQuAD v1.1
 - Represent the input question and paragraph as a single packed sequence
 - The question uses the A embedding and the paragraph uses the B embedding
 - New parameters to be learned in fine-tuning are start vector $S \in \mathbb{R}^H$ and end vector $E \in \mathbb{R}^H$
 - Calculate the probability of word i being the start of the answer span

$$P_i = \frac{e^{S \cdot T_i}}{\sum_j e^{S \cdot T_j}}$$

- The training objective is the log-likelihood the correct and end positions

Comparison of BERT and OpenAI GPT

OpenAI GPT	BERT
Trained on BooksCorpus (800M)	Trained on BooksCorpus (800M) + Wikipedia (2,500M)
Use sentence separator ([SEP]) and classifier token ([CLS]) only at fine-tuning time	BERT learns [SEP], [CLS] and sentence A/B embeddings during pre-training
Trained for 1M steps with a batch-size of 32,000 words	Trained for 1M steps with a batch-size of 128,000 words
Use the same learning rate of 5e-5 for all fine-tuning experiments	BERT choose a task-specific learning rate which performs the best on the development set

Outline

- Research context
- Main ideas
- BERT
- **Experiments**
- Conclusions

Experiments

- GLUE (General Language Understanding Evaluation) benchmark
 - Distribute canonical Train, Dev and Test splits
 - Labels for Test set are not provided
- Datasets in GLUE:
 - MNLI: Multi-Genre Natural Language Inference
 - QQP: Quora Question Pairs
 - QNLI: Question Natural Language Inference
 - SST-2: Stanford Sentiment Treebank
 - CoLA: The corpus of Linguistic Acceptability
 - STS-B: The Semantic Textual Similarity Benchmark
 - MRPC: Microsoft Research Paraphrase Corpus
 - RTE: Recognizing Textual Entailment
 - WNLI: Winograd NLI

GLUE Results

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. OpenAI GPT = (L=12, H=768, A=12); BERT_{BASE} = (L=12, H=768, A=12); BERT_{LARGE} = (L=24, H=1024, A=16). BERT and OpenAI GPT are single-model, single task. All results obtained from <https://gluebenchmark.com/leaderboard> and <https://blog.openai.com/language-unsupervised/>.

SQuAD v1.1

System	Dev		Test	
	EM	F1	EM	F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
#1 Single - nlnet	-	-	83.5	90.1
#2 Single - QANet	-	-	82.5	89.3
Published				
BiDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

Table 2: SQuAD results. The BERT ensemble is 7x systems which use different pre-training checkpoints and fine-tuning seeds.

Reference: <https://rajpurkar.github.io/SQuAD-explorer>

Named Entity Recognition

System	Dev F1	Test F1
ELMo+BiLSTM+CRF	95.7	92.2
CVT+Multi (Clark et al., 2018)	-	92.6
BERT _{BASE}	96.4	92.4
BERT _{LARGE}	96.6	92.8

Table 3: CoNLL-2003 Named Entity Recognition results. The hyperparameters were selected using the Dev set, and the reported Dev and Test scores are averaged over 5 random restarts using those hyperparameters.

SWAG

- The Situations with Adversarial Generations (SWAG)

A girl is going across a set of monkey bars. She

- (i) jumps up across the monkey bars.
- (ii) struggles onto the bars to grab her head.
- (iii) gets to the end and stands on a wooden plank.
- (iv) jumps up and does a back flip.

- The only task-specific parameters is a vector $V \in \mathbb{R}^H$
- The probability distribution is the softmax over the four choices

$$P_i = \frac{e^{V \cdot C_i}}{\sum_{j=1}^4 e^{V \cdot C_i}}$$

SWAG Result

System	Dev	Test
ESIM+GloVe	51.9	52.7
ESIM+ELMo	59.1	59.2
BERT _{BASE}	81.6	-
BERT _{LARGE}	86.6	86.3
Human (expert) [†]	-	85.0
Human (5 annotations) [†]	-	88.0

Table 4: SWAG Dev and Test accuracies. Test results were scored against the hidden labels by the SWAG authors. [†]Human performance is measured with 100 samples, as reported in the SWAG paper.

Ablation Studies

- To understand
 - Effect of Pre-training Tasks
 - Effect of model sizes
 - Effect of number of training steps
 - Feature-based approach with BERT

Ablation Studies

- Main findings:
 - “Next sentence prediction” (NSP) pre-training is important for sentence-pair classification task
 - Bidirectionality (using MLM pre-training task) contributes significantly to the performance improvement

Tasks	Dev Set				
	MNLI-m (Acc)	QNLI (Acc)	MRPC (Acc)	SST-2 (Acc)	SQuAD (F1)
BERT _{BASE}	84.4	88.4	86.7	92.7	88.5
No NSP	83.9	84.9	86.5	92.6	87.9
LTR & No NSP	82.1	84.3	77.5	92.1	77.8
+ BiLSTM	82.1	84.1	75.7	91.6	84.9

Table 5: Ablation over the pre-training tasks using the BERT_{BASE} architecture. “No NSP” is trained without the next sentence prediction task. “LTR & No NSP” is trained as a left-to-right LM without the next sentence prediction, like OpenAI GPT. “+ BiLSTM” adds a randomly initialized BiLSTM on top of the “LTR + No NSP” model during fine-tuning.

Ablation Studies

- Main findings
 - Bigger model sizes are better even for small-scale tasks

#L	#H	#A	Hyperparams				Dev Set Accuracy		
			LM (ppl)	MNLI-m	MRPC	SST-2			
3	768	12	5.84	77.9	79.8	88.4			
6	768	3	5.24	80.6	82.2	90.7			
6	768	12	4.68	81.9	84.8	91.3			
12	768	12	3.99	84.4	86.7	92.9			
12	1024	16	3.54	85.7	86.9	93.3			
24	1024	16	3.23	86.6	87.8	93.7			

Table 6: Ablation over BERT model size. #L = the number of layers; #H = hidden size; #A = number of attention heads. “LM (ppl)” is the masked LM perplexity of held-out training data.

Conclusions

- Unsupervised pre-training (pre-training language model) is increasingly adopted in many NLP tasks
- Major contribution of the paper is to propose a deep *bidirectional* architecture from Transformer
 - Advance state-of-the-art for many important NLP tasks

Links

- TensorFlow code and pre-trained models for BERT:
<https://github.com/google-research/bert>
- PyTorch Pretrained Bert:
<https://github.com/huggingface/pytorch-pretrained-BERT>
- BERT-pytorch: <https://github.com/codertimo/BERT-pytorch>
- BERT-keras: <https://github.com/Separius/BERT-keras>

Remark: Applying BERT for non-English languages

- Pre-trained BERT models are provided for more than 100 languages (including Vietnamese)
 - <https://github.com/google-research/bert/blob/master/multilingual.md>
- Be careful with tokenization!!
- For Japanese (and Chinese): “spaces were added around every character in the CJK Unicode range before applying WordPiece” => Not a good way to do
 - Use SentencePiece:
<https://github.com/google/sentencepiece>
 - We may need to pre-train BERT model

References

1. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv preprint arXiv:1810.04805.
2. Vaswani et al. (2017). Attention Is All You Need. arXiv preprint arXiv:1706.03762. <https://arxiv.org/abs/1706.03762>
3. The Annotated Transformer:
<http://nlp.seas.harvard.edu/2018/04/03/attention.html>, by harvardnlp.
4. The Illustrated Transformer: <http://jalammar.github.io/illustrated-transformer/>
5. ELMo explained: <https://www.mihaileric.com/posts/deep-contextualized-word-representations-elmo/>
6. ULMFit explained: <https://yashuseth.blog/2018/06/17/understanding-universal-language-model-fine-tuning-ulmfit/>
7. Dissecting BERT: <https://medium.com/dissecting-bert>
8. The Illustrated BERT, ELMo, and co. (How NLP Cracked Transfer Learning):
<http://jalammar.github.io/illustrated-bert/>