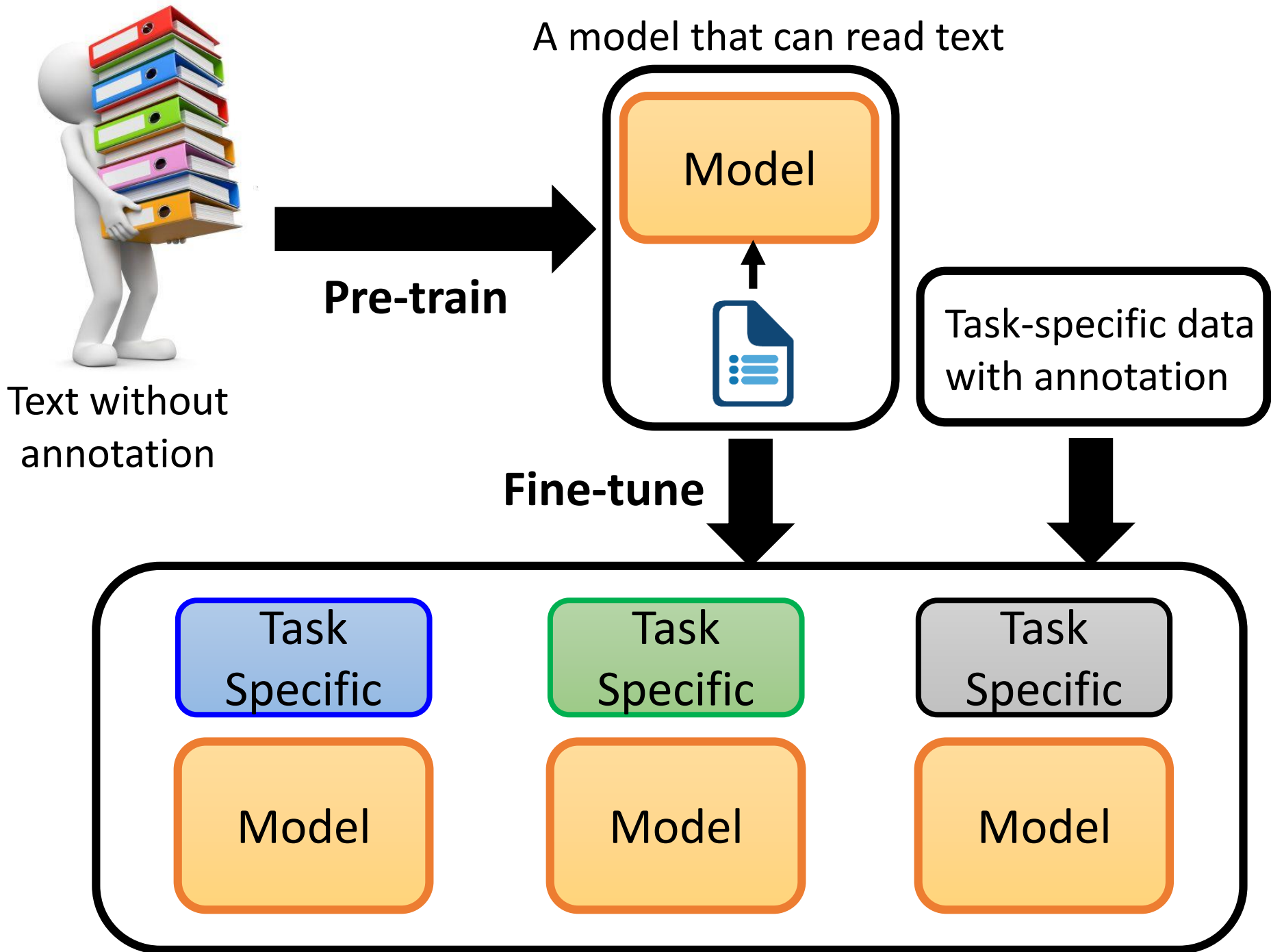


# BERT and its family

Hung-yi Lee 李宏毅



# Outline

What is pre-train model

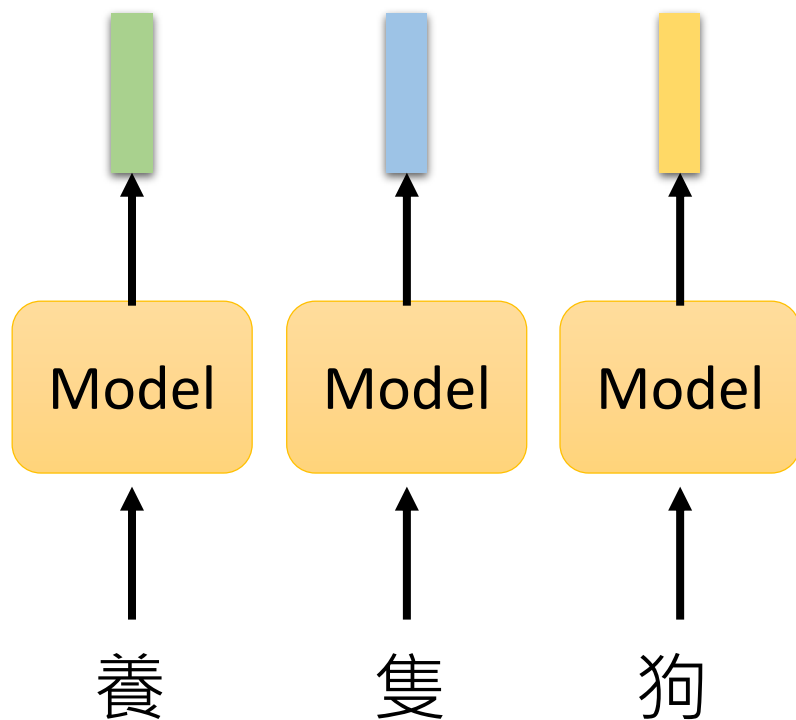
How to fine-tune

How to pre-train

Pre-train Model

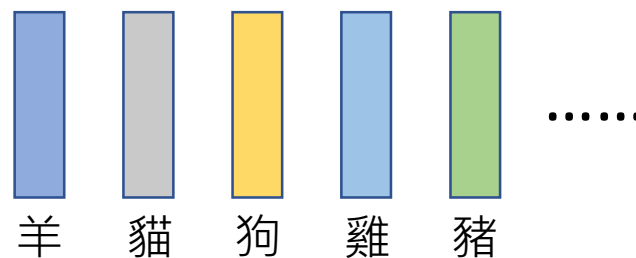
# Pre-train Model

Represent each token by a embedding vector



The token with the same type has the same embedding.

Simply a table look-up

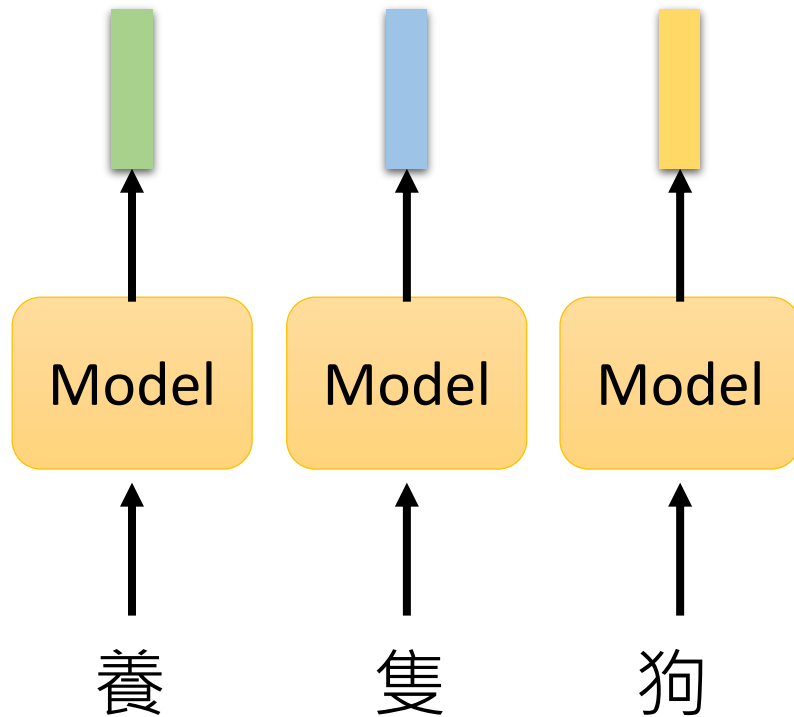


Word2vec [Mikolov, et al., NIPS'13]

Glove [Pennington, et al., EMNLP'14]

# Pre-train Model

Represent each token by a embedding vector

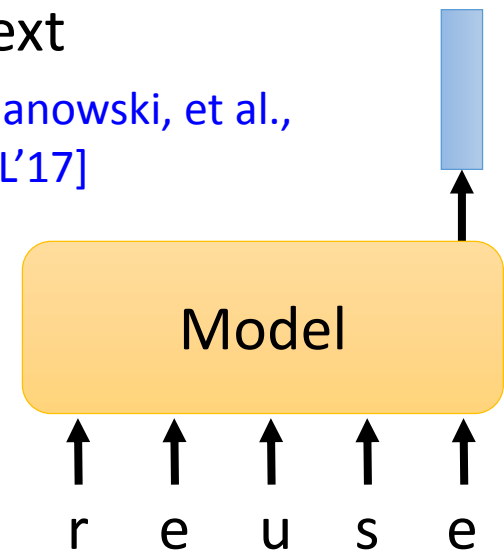


The token with the same type has the same embedding.

English word as token ...

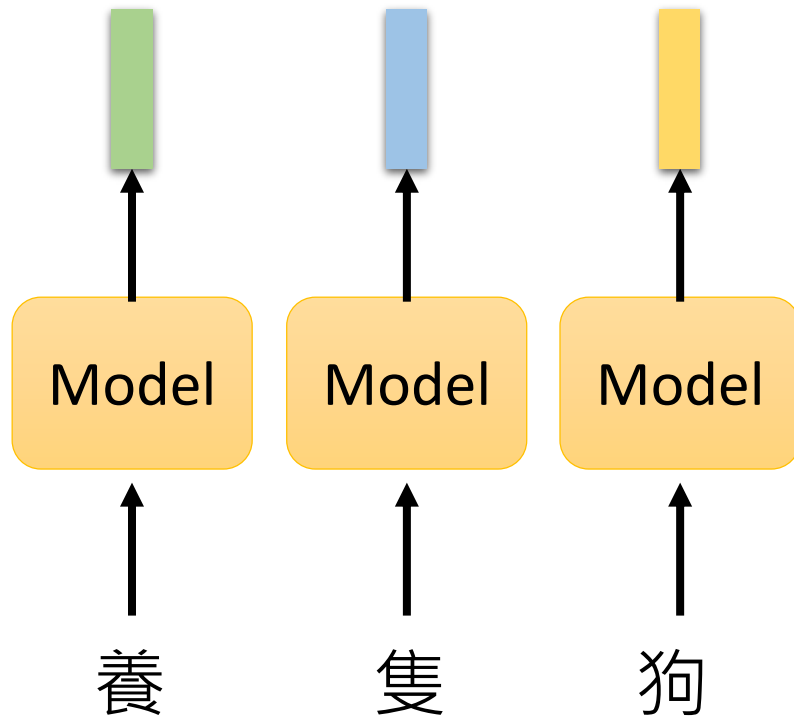
FastText

[Bojanowski, et al.,  
TACL'17]



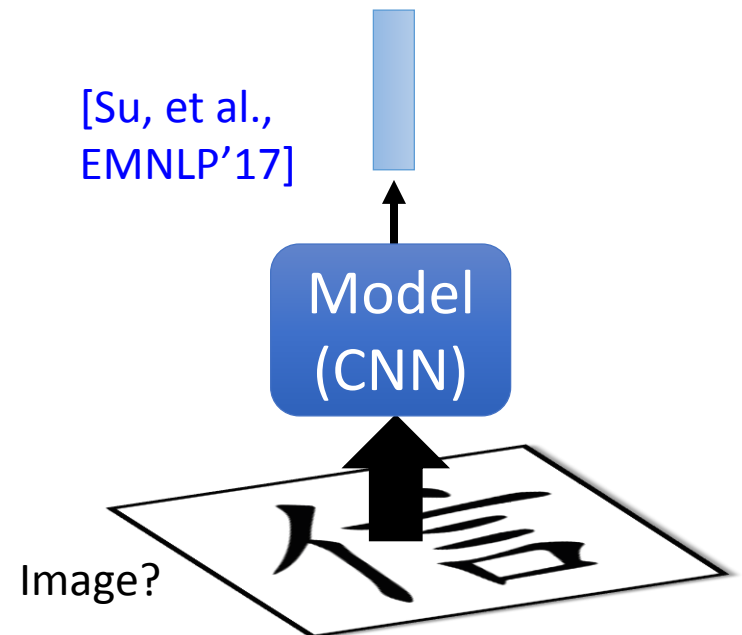
# Pre-train Model

Represent each token by a embedding vector



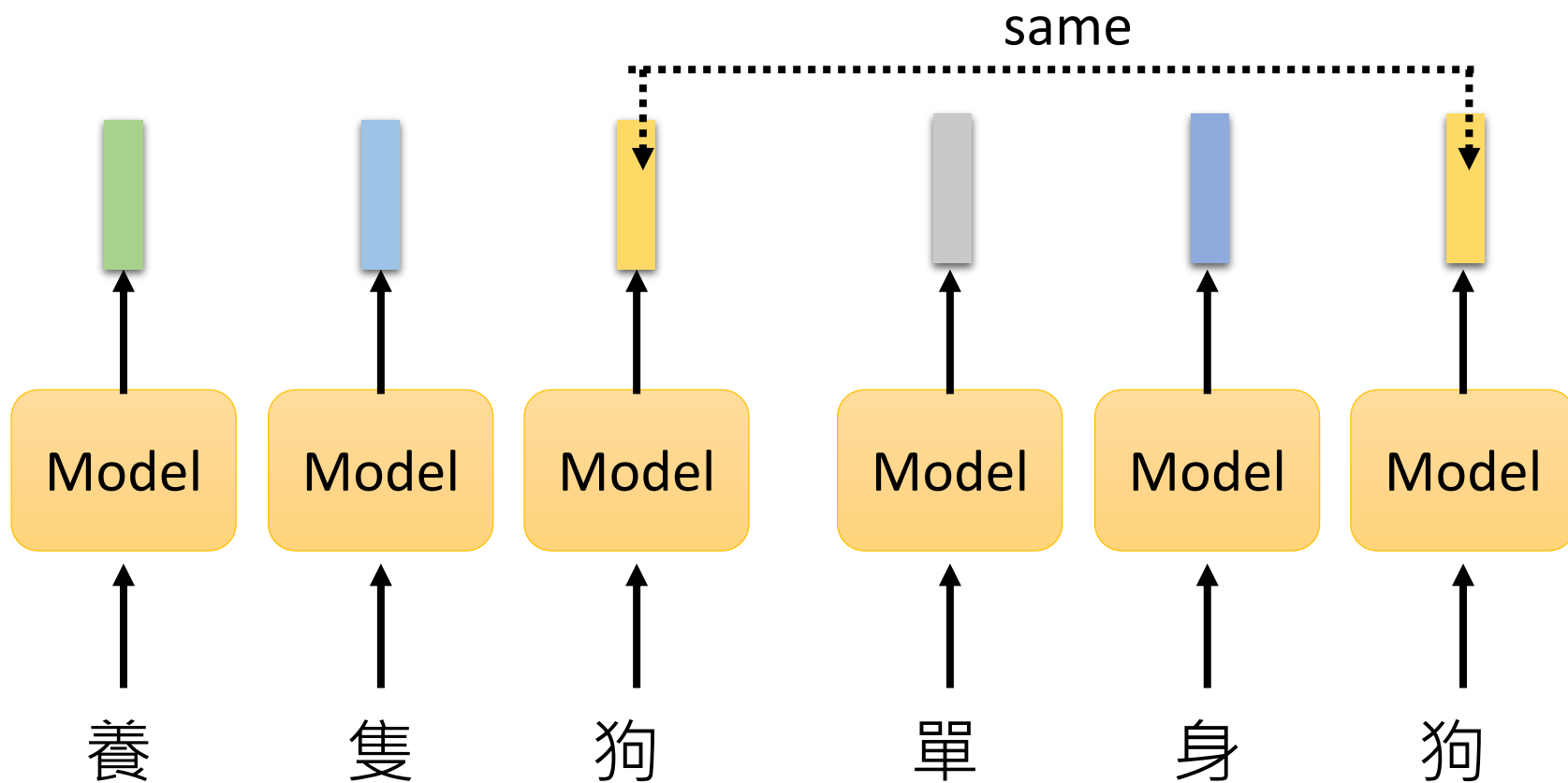
The token with the same type has the same embedding.

Chinese character as token ...



# Pre-train Model

Represent each token by a embedding vector



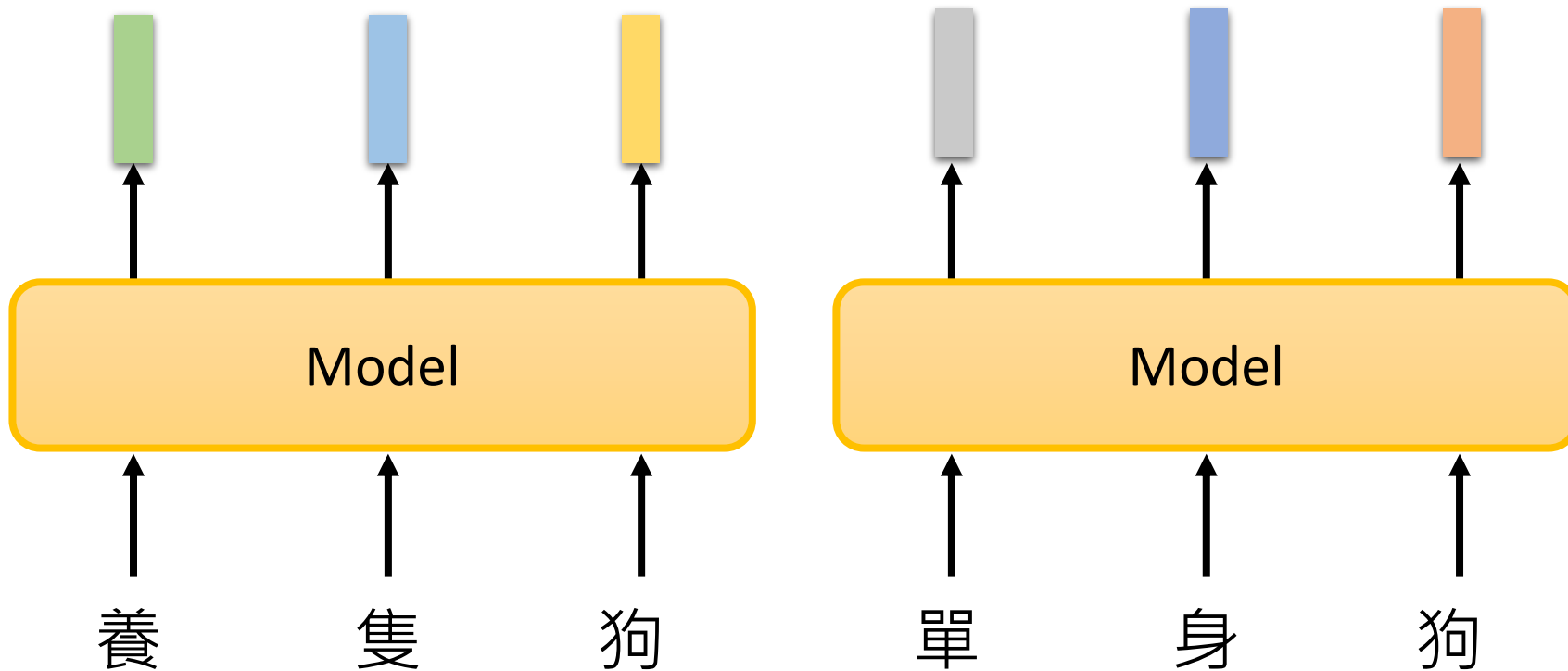


Pre-train  
Representation



# Pre-train Model

Contextualized Word Embedding



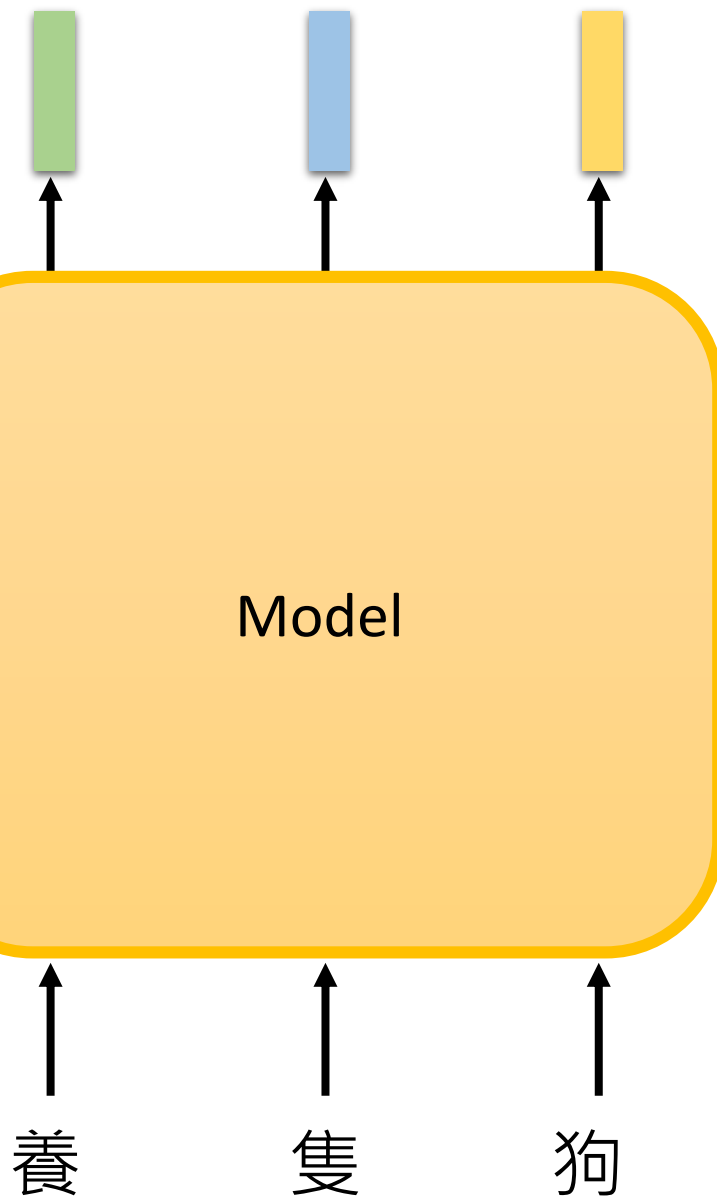
# Pre-train Model

Contextualized Word Embedding

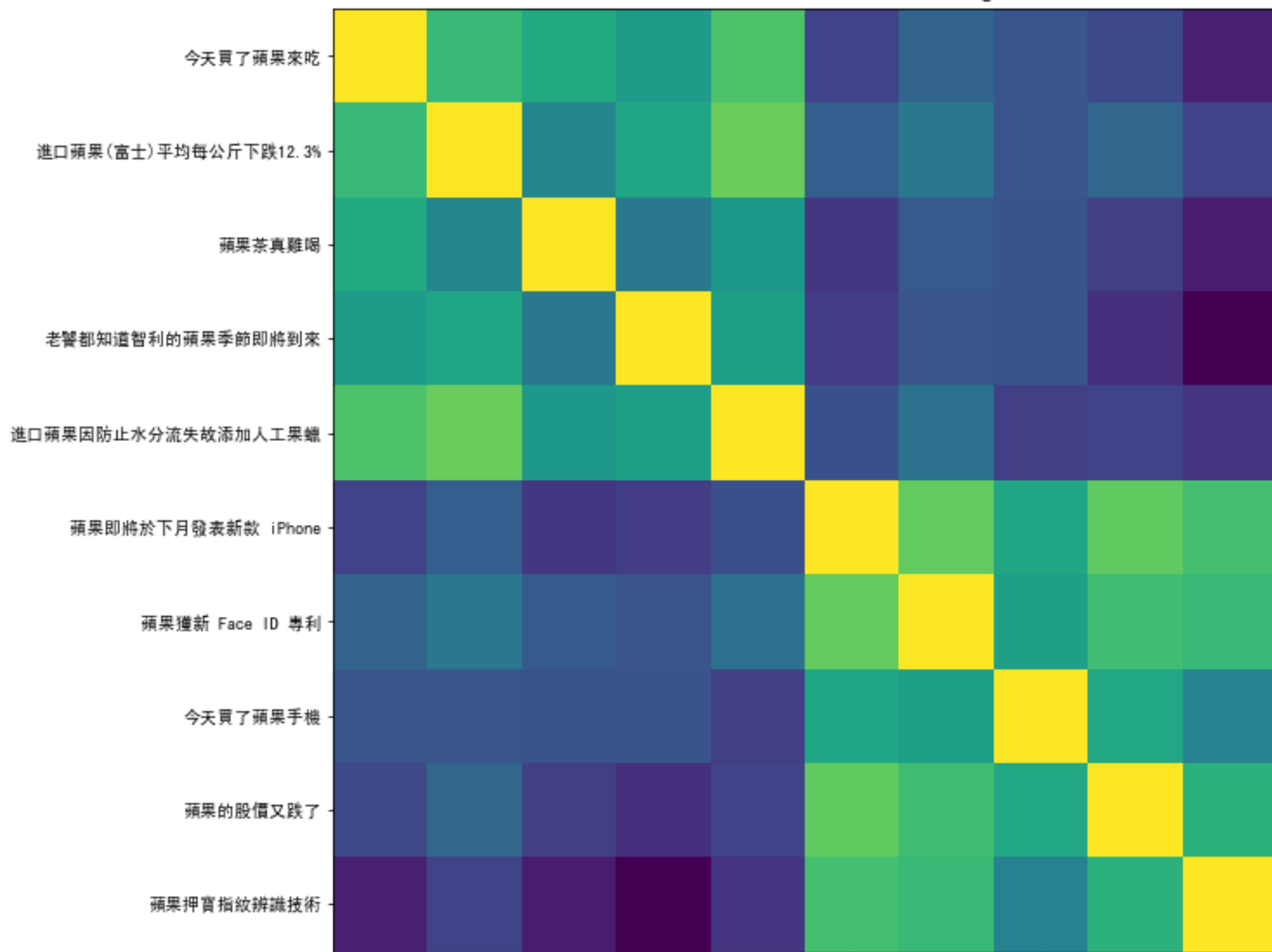
Many Layers

Model

- LSTM
- Self-attention layers
- Tree-based model (?)
  - Ref: <https://youtu.be/z0uOq2wEGcc>

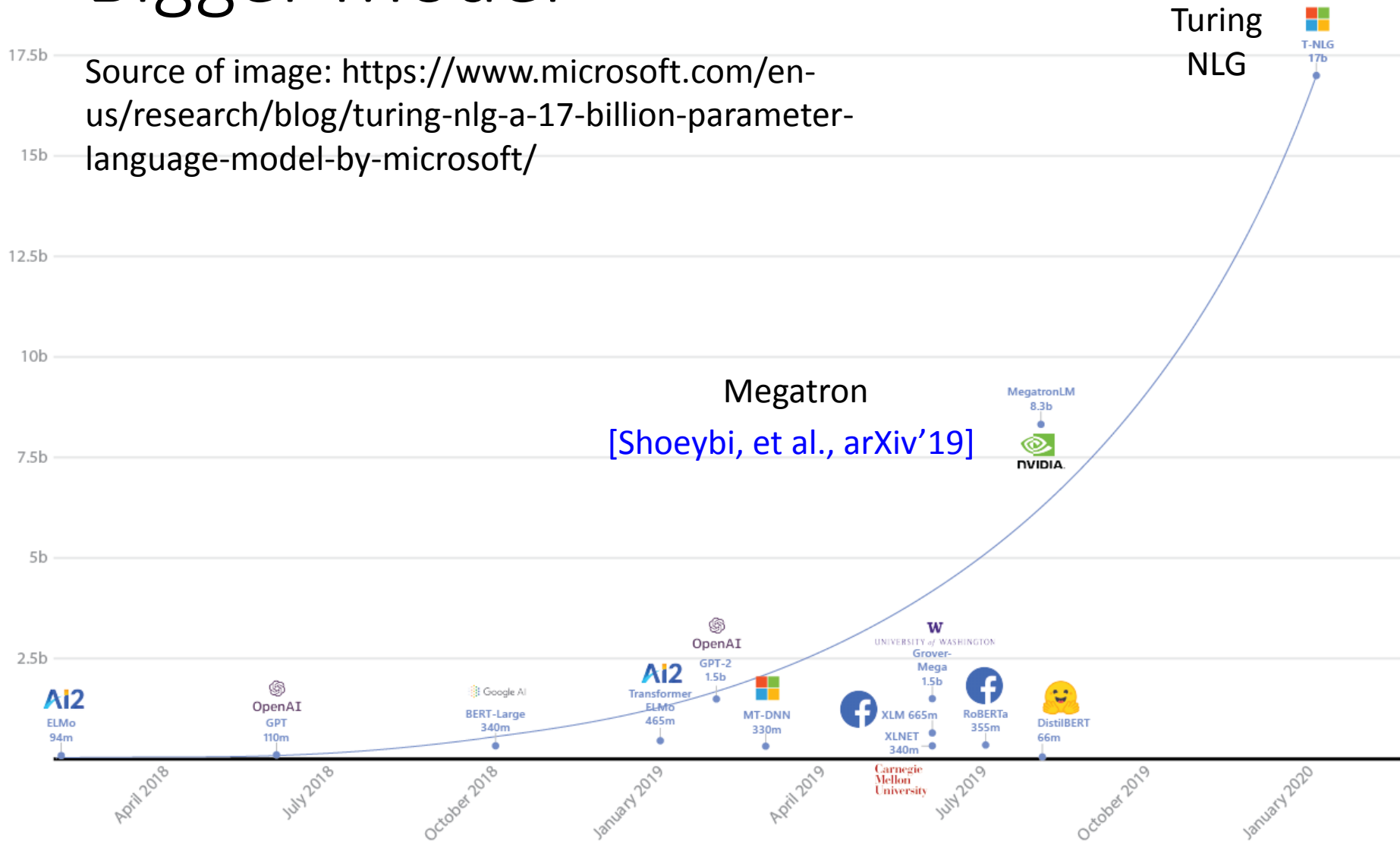


Cosine Similarities of BERT Embeddings



# Bigger Model

Source of image: <https://www.microsoft.com/en-us/research/blog/turing-nlg-a-17-billion-parameter-language-model-by-microsoft/>



# Smaller Model



Distill BERT

[Sanh, et al., NeurIPS workshop'19]

Tiny BERT [Jian, et al., arXiv'19]

Mobile BERT [Sun, et al., ACL'20]

Q8BERT

[Zafrir, et al., NeurIPS workshop 2019]

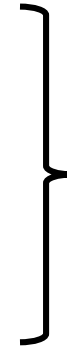
ALBERT [Lan, et al., ICLR'20]

# Smaller Model

- Network Compression

Ref: [https://youtu.be/dPp8rCAnU\\_A](https://youtu.be/dPp8rCAnU_A)

- Network Pruning
- Knowledge Distillation
- Parameter Quantization
- Architecture Design



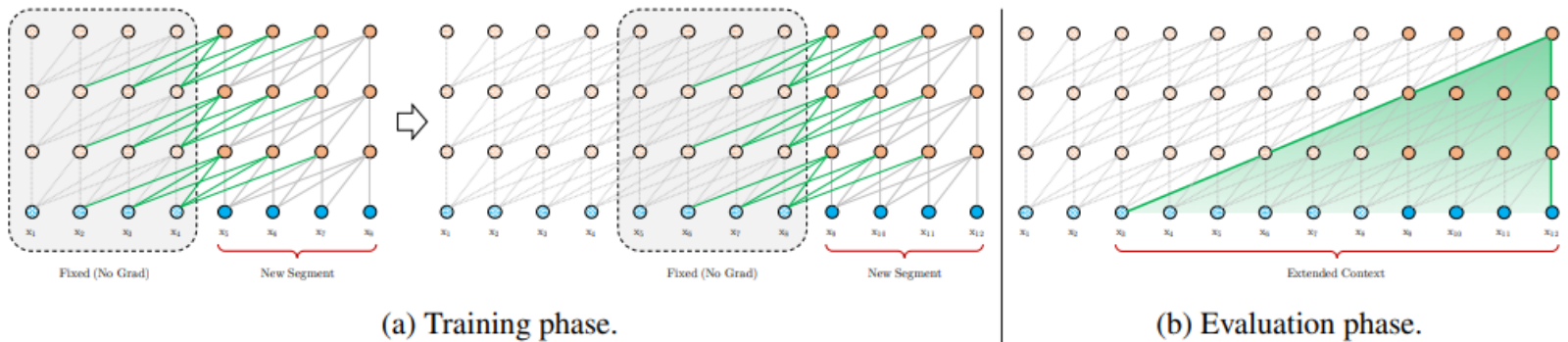
All of them have  
been tried.

Excellent reference:

<http://mitchgordon.me/machine/learning/2019/11/18/all-the-ways-to-compress-BERT.html>

# Network Architecture

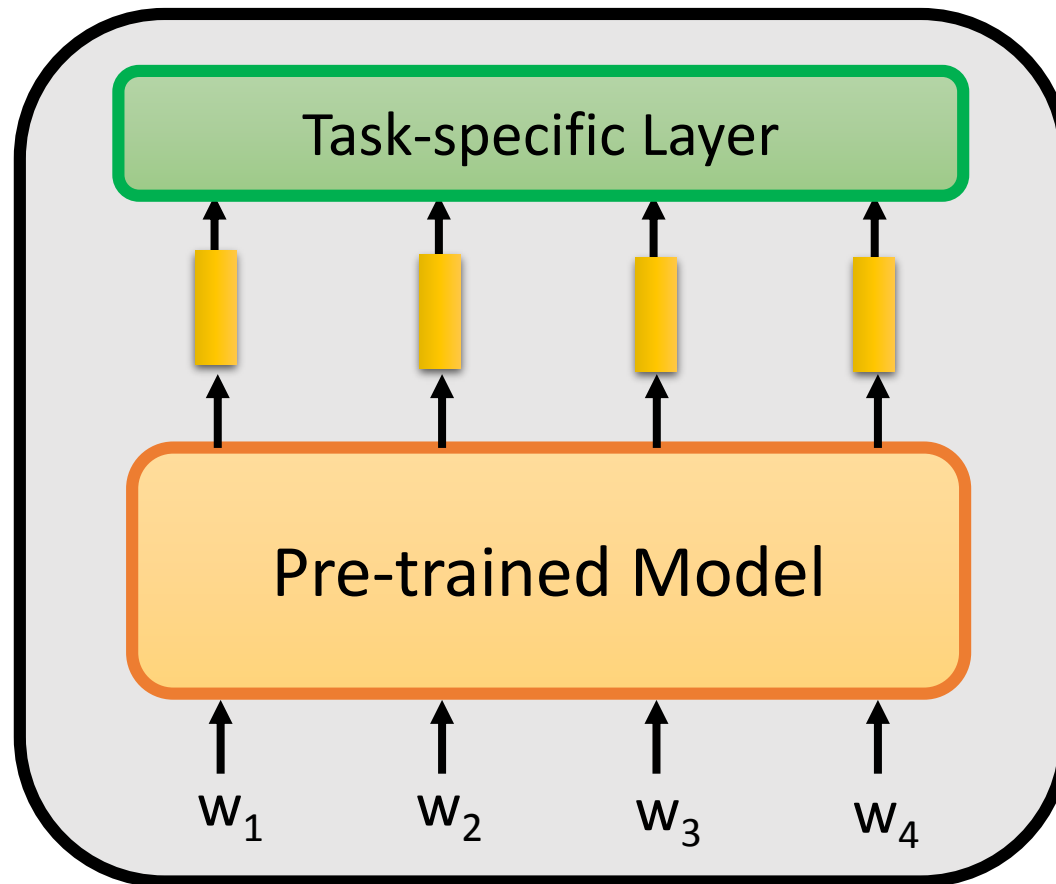
- Transformer-XL: Segment-Level Recurrence with State Reuse [Dai, et al., ACL'19]



- Reformer [Kitaev, et al., ICLR'20]
  - Longformer [Beltagy, et al., arXiv'20]
- } Reduce the complexity of self-attention

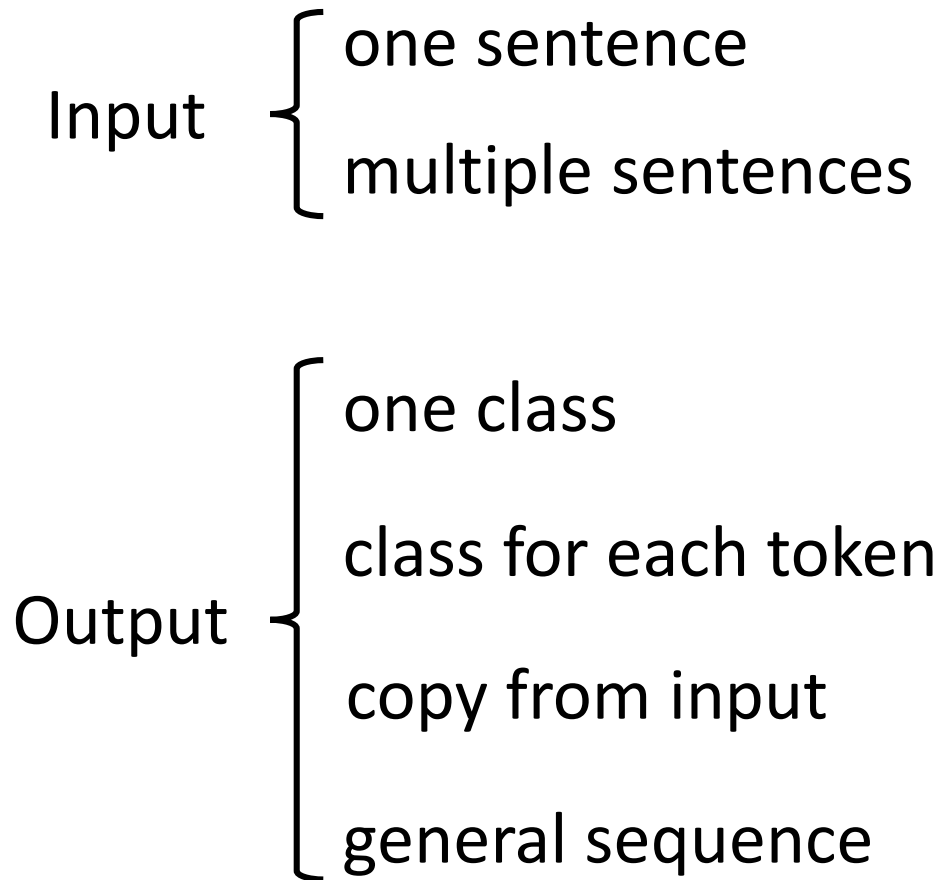


# How to fine-tune



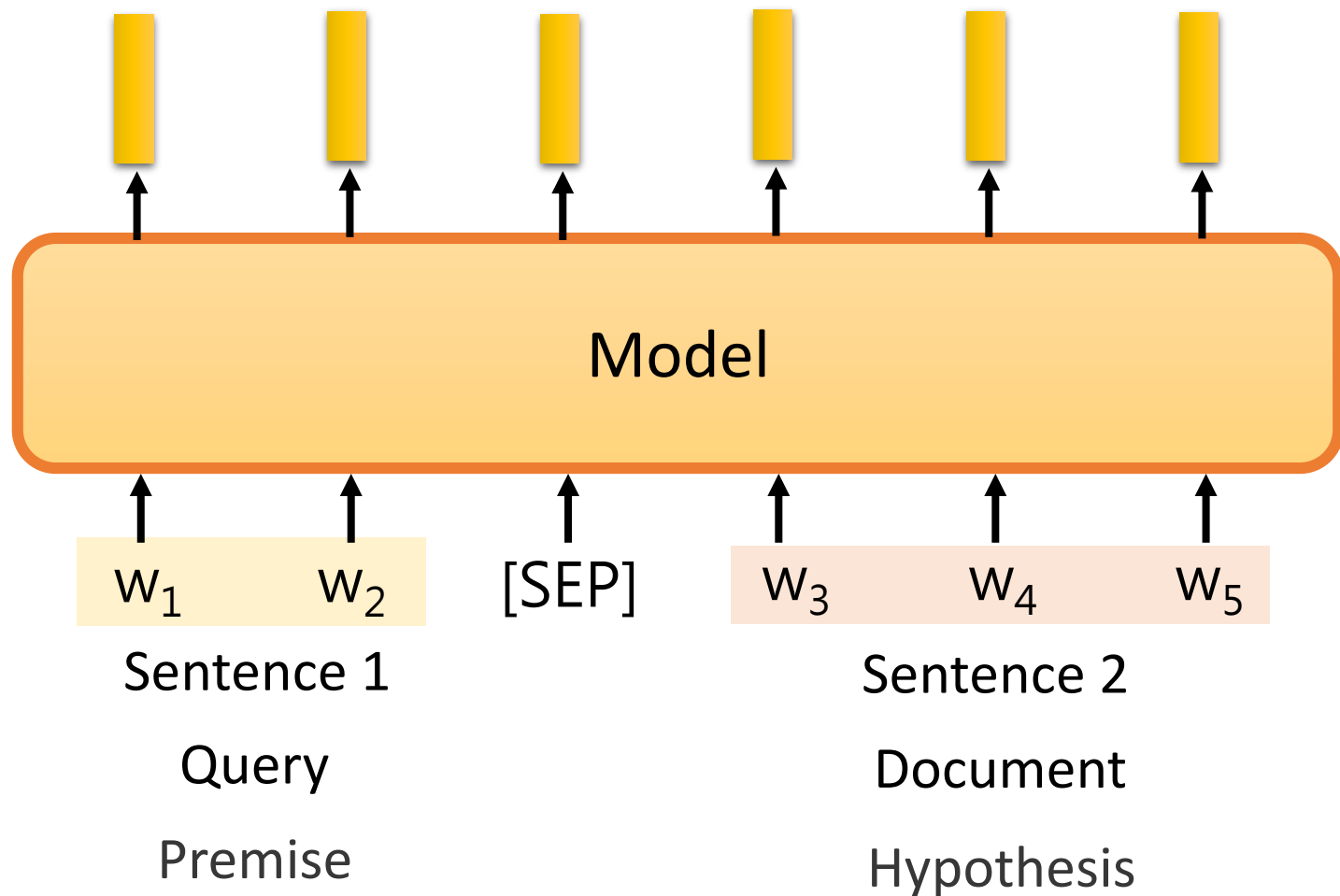
For a specific  
NLP task

# NLP tasks

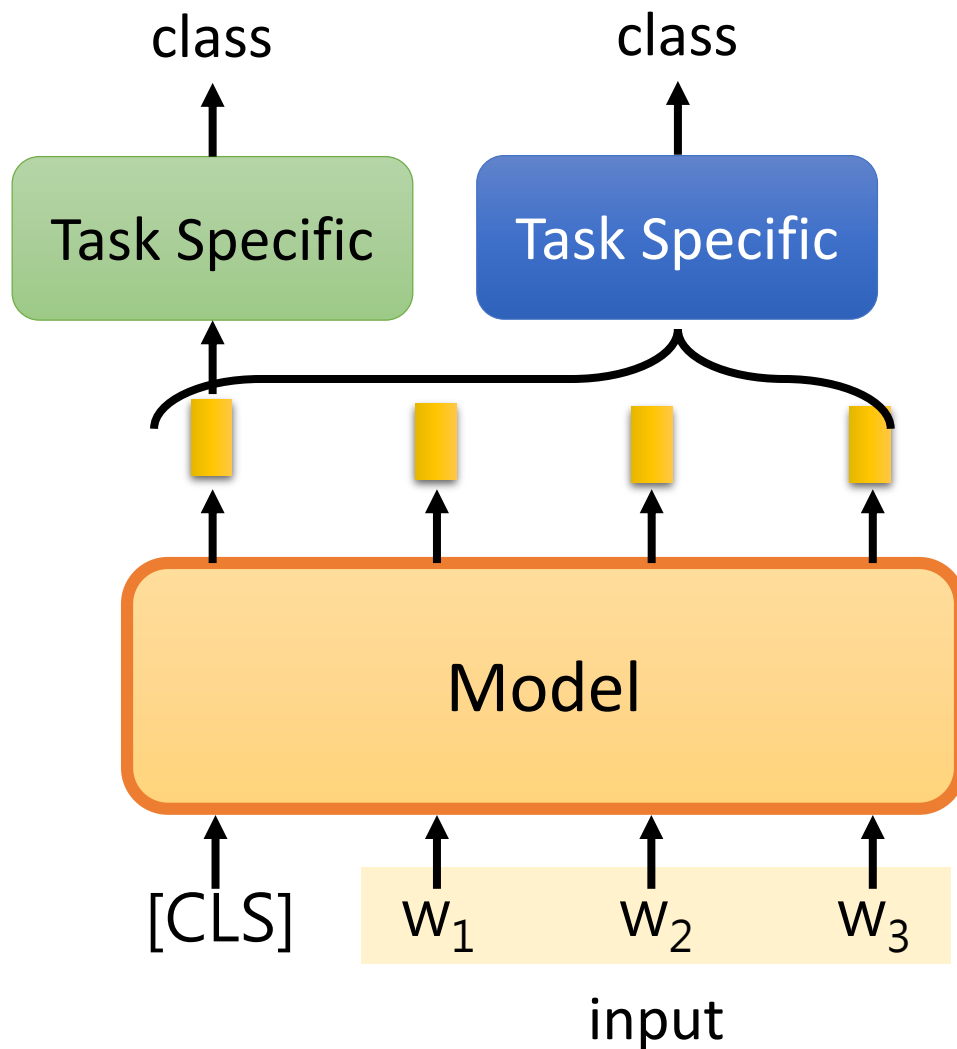


Input

one sentence  
multiple sentences



# Output



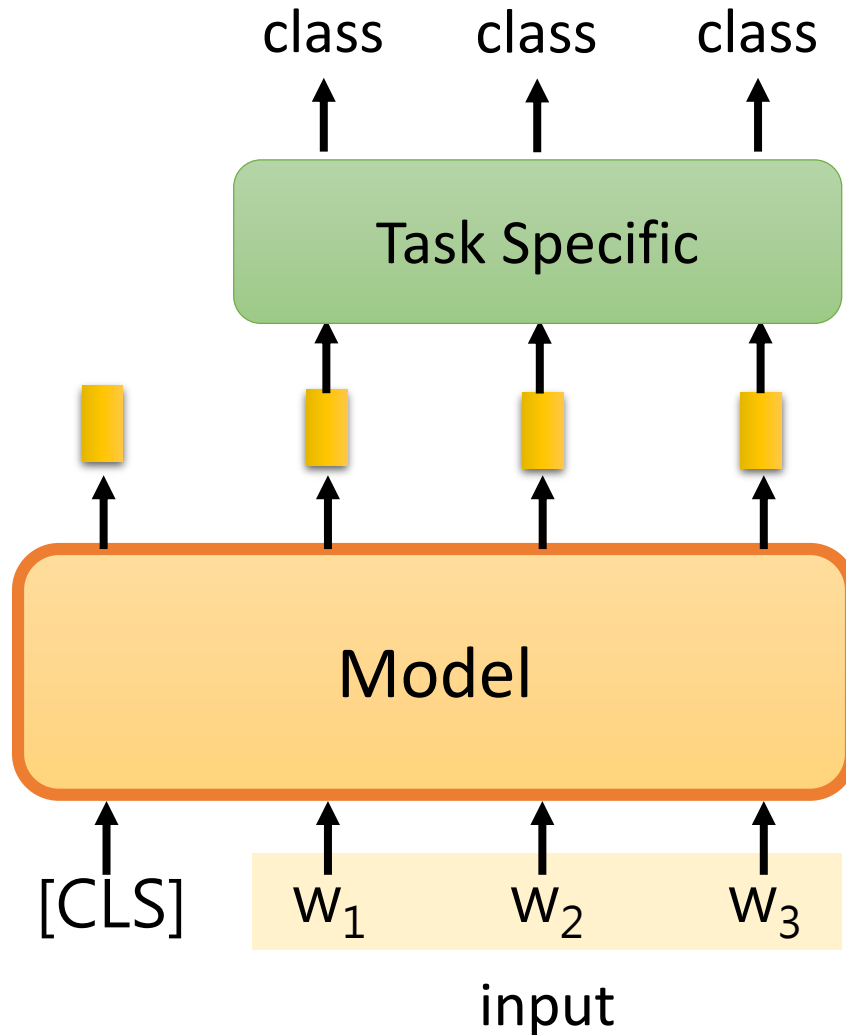
one class

class for each token

copy from input

general sequence

# Output



one class

class for each token

copy from input

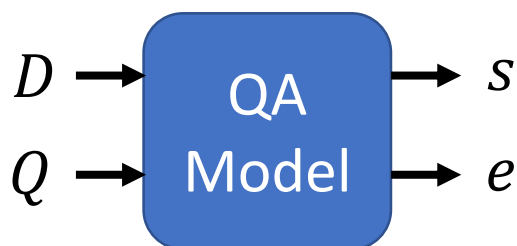
general sequence

# Output

- Extraction-based QA

**Document:**  $D = \{d_1, d_2, \dots, d_N\}$

**Query:**  $Q = \{q_1, q_2, \dots, q_M\}$



output: two integers ( $s, e$ )

**Answer:**  $A = \{d_s, \dots, d_e\}$

one class

class for each token

copy from input

general sequence

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **grau-pel** and hail... Precipitation occurs as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?  
**gravity**

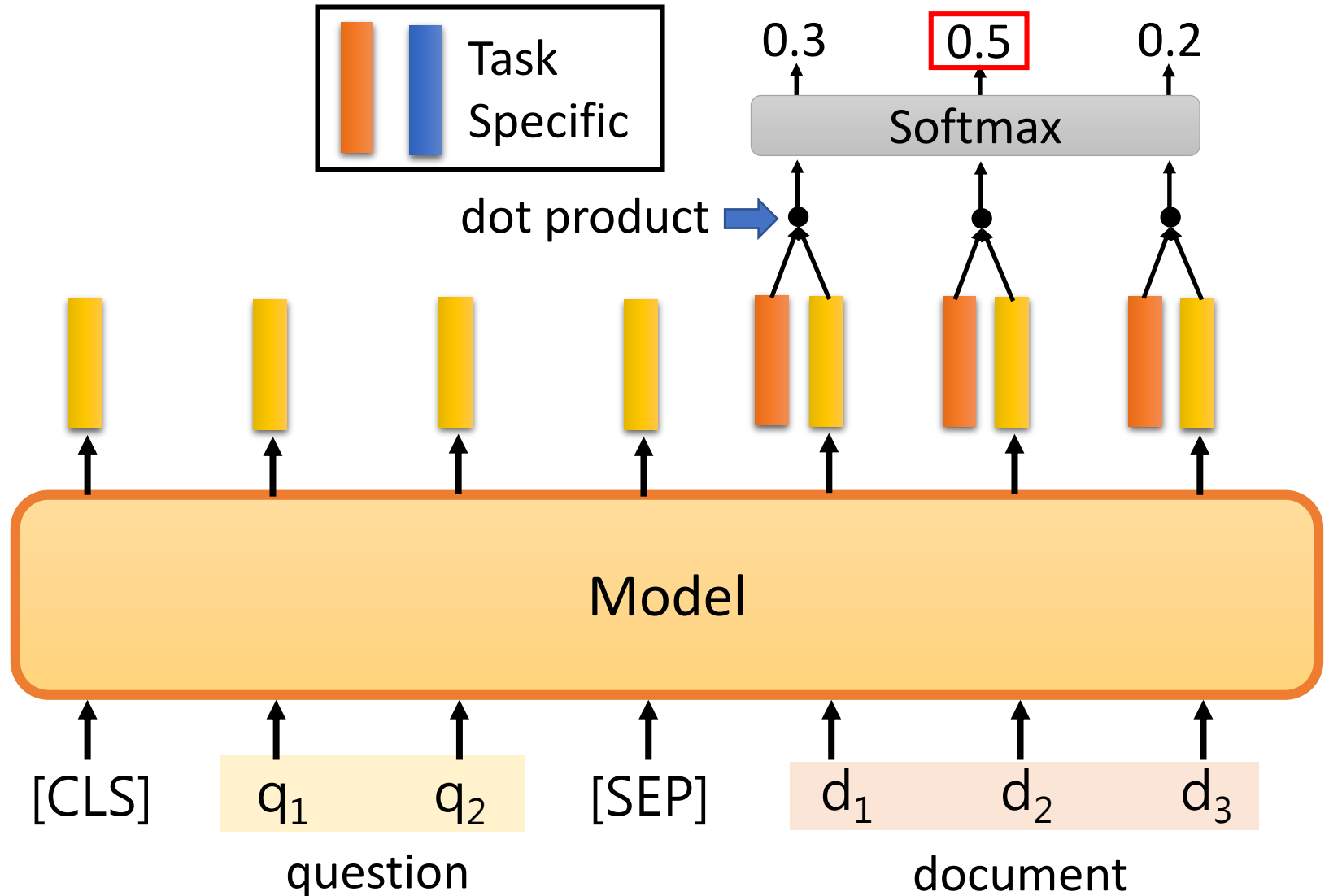
Where do water droplets collide with ice crystals to form precipitation?

**within a cloud**

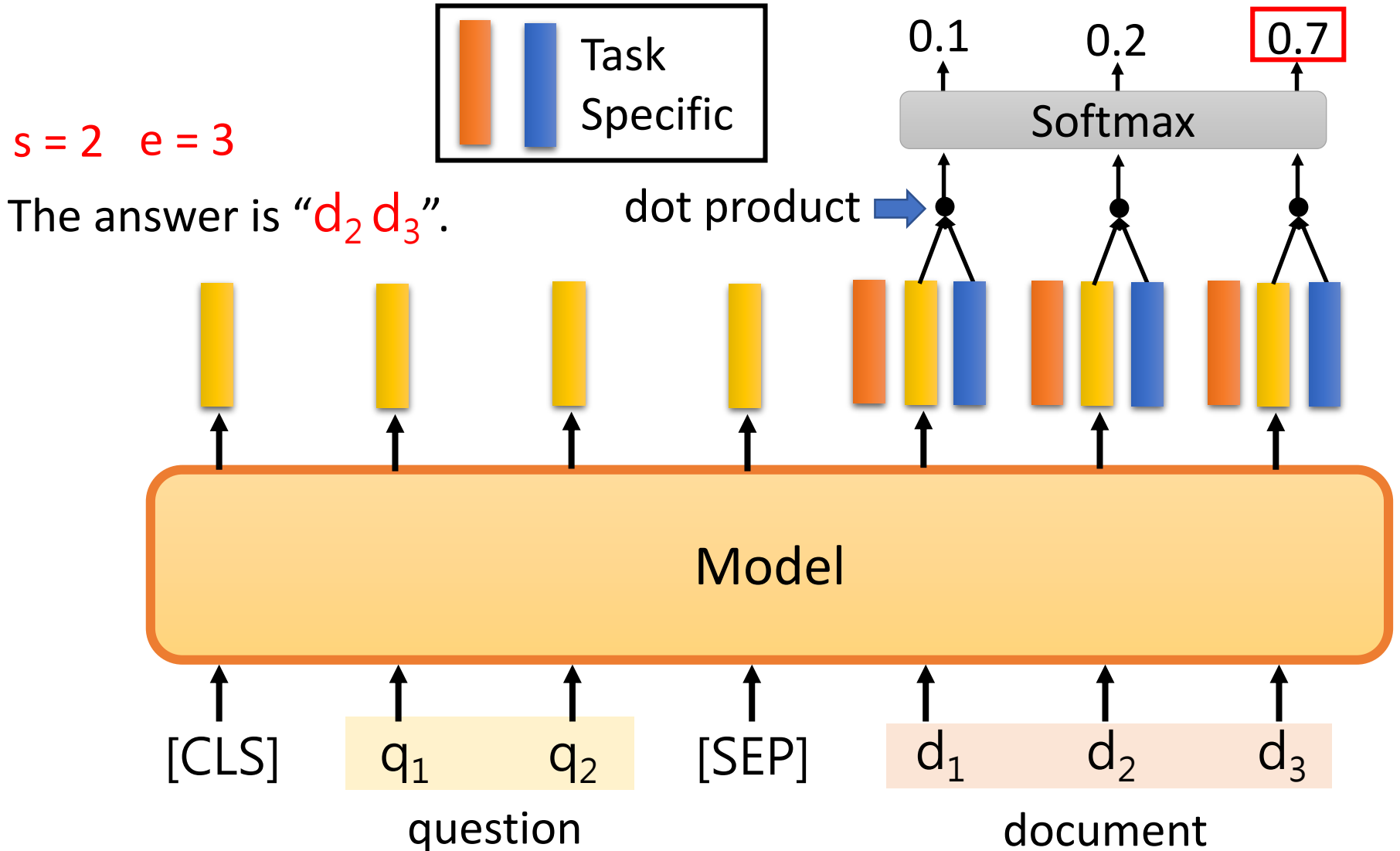
$s = 77, e = 79$

# Copy from Input (BERT)

$s = 2$



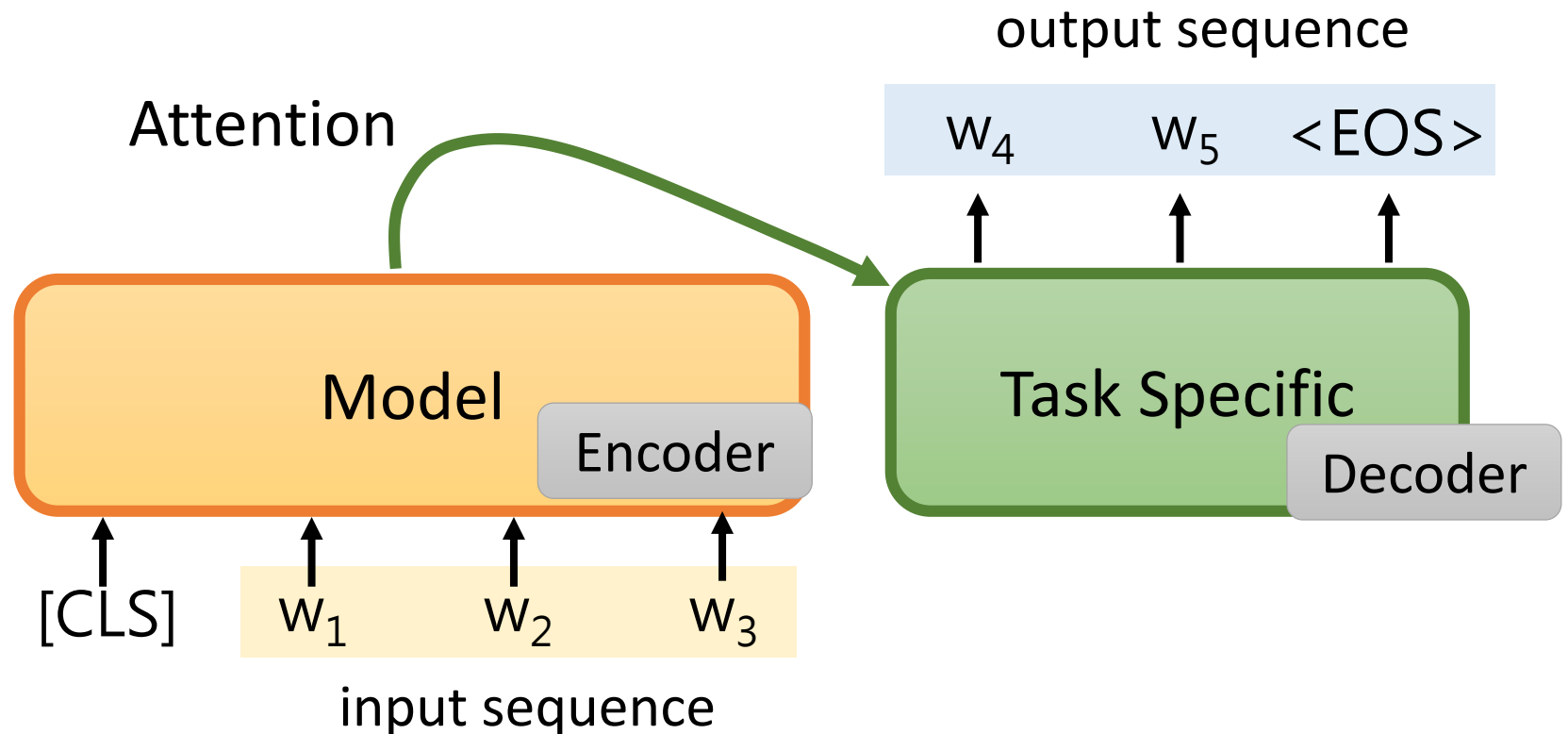
# Copy from Input (BERT)



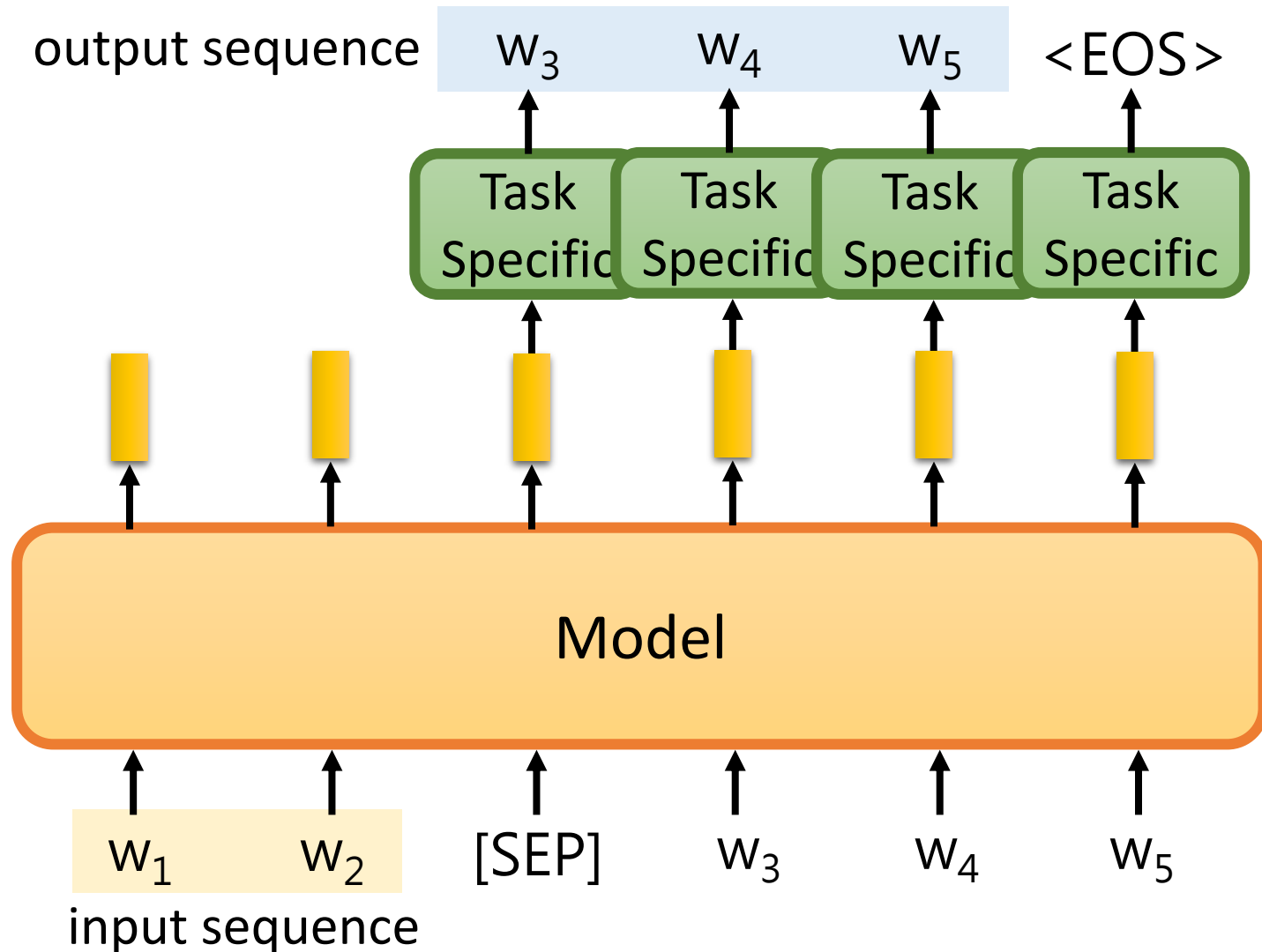


# Output – General Sequence (v1)

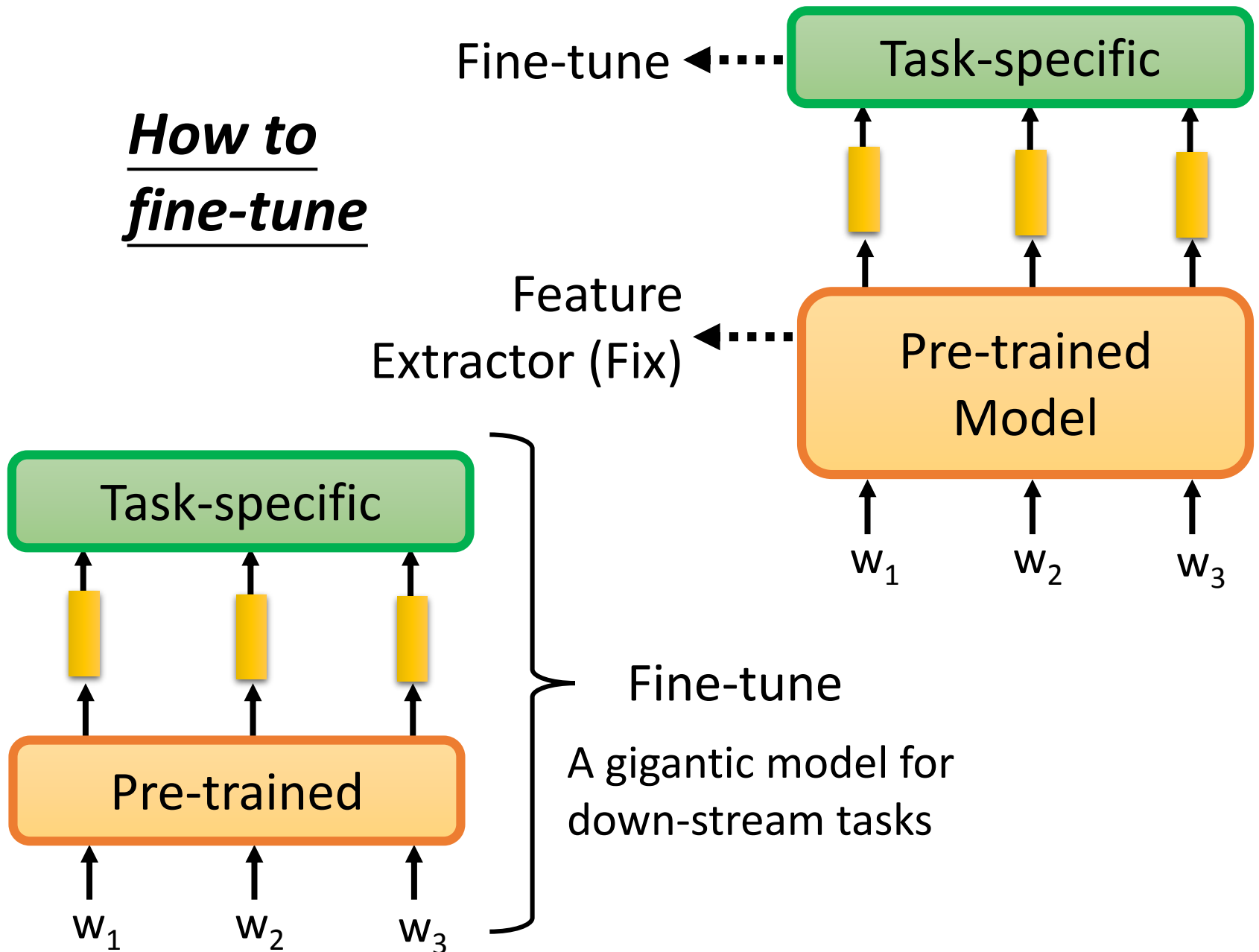
- Seq2seq model



# Output – General Sequence (v2)

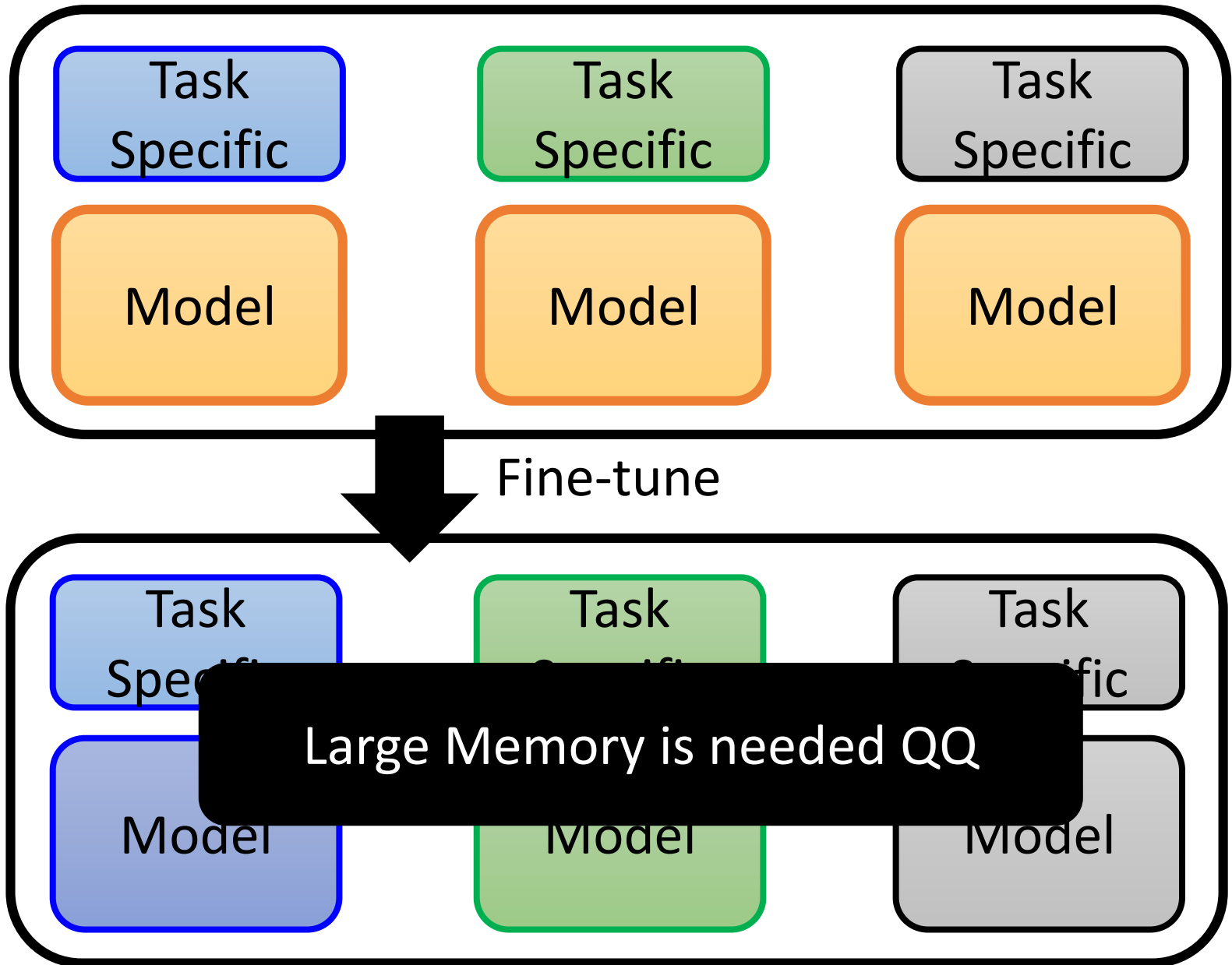


## How to *fine-tune*



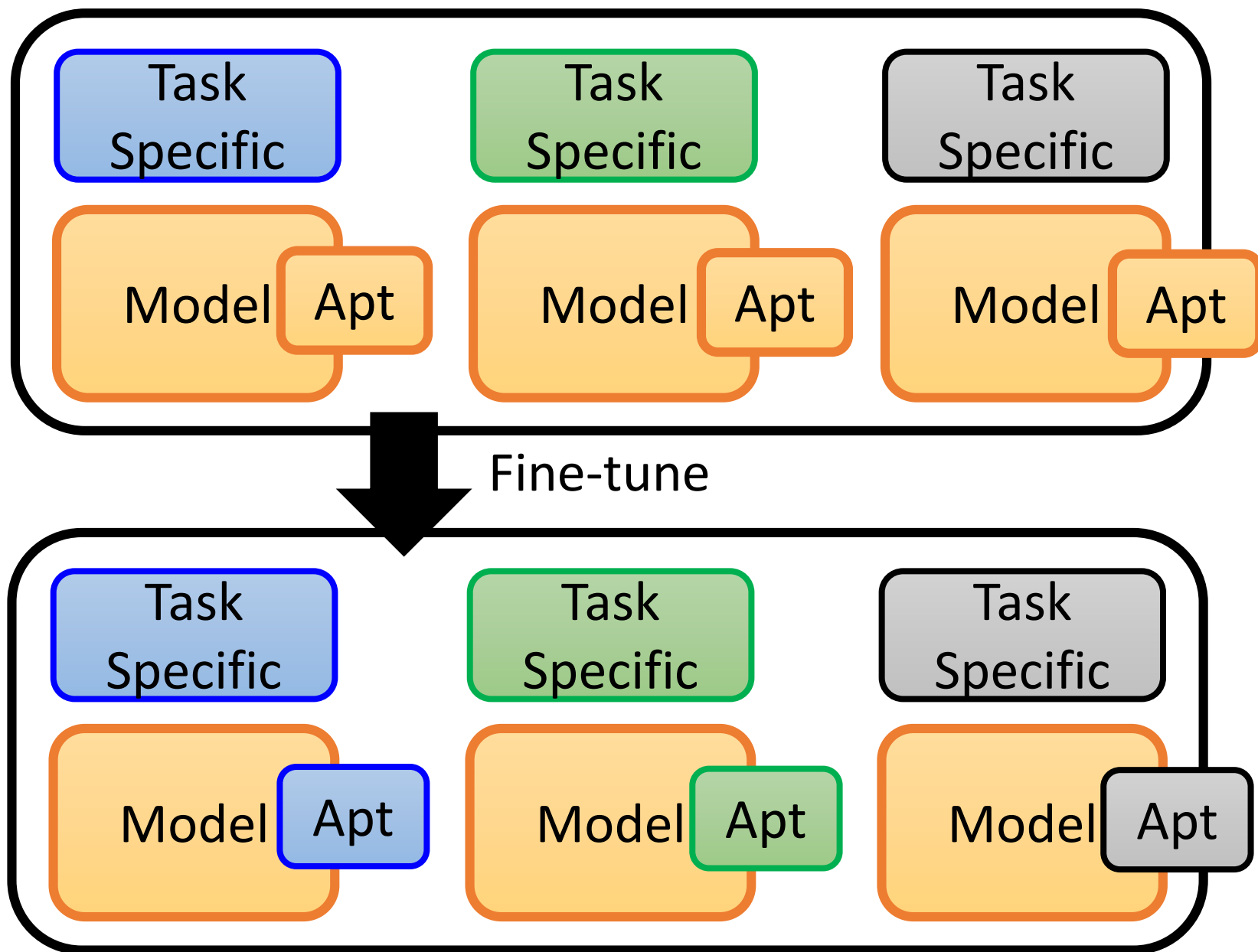
# Adaptor

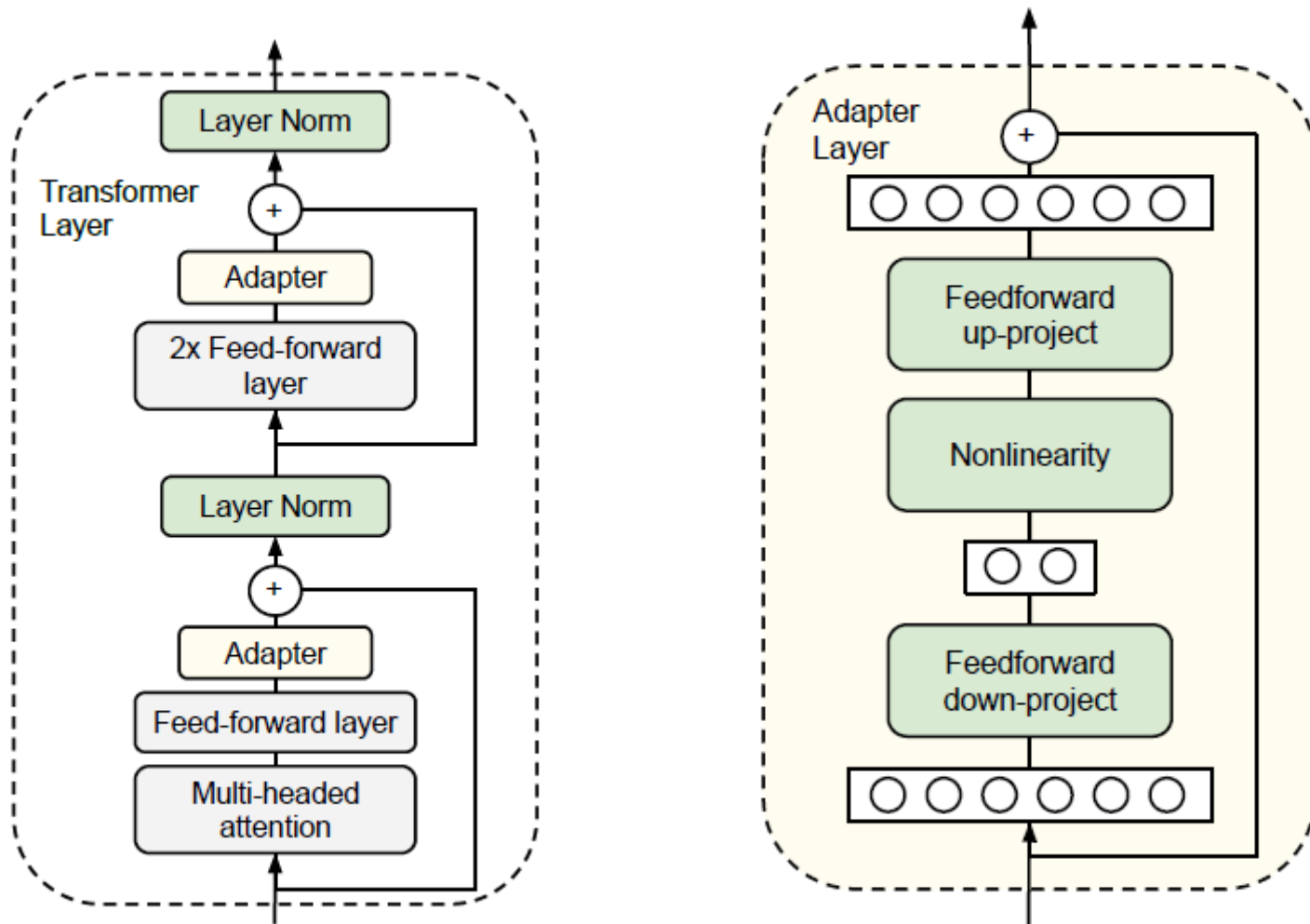
[Stickland, et al., ICML'19] [Houlsby, et al., ICML'19]



# **Adaptor**

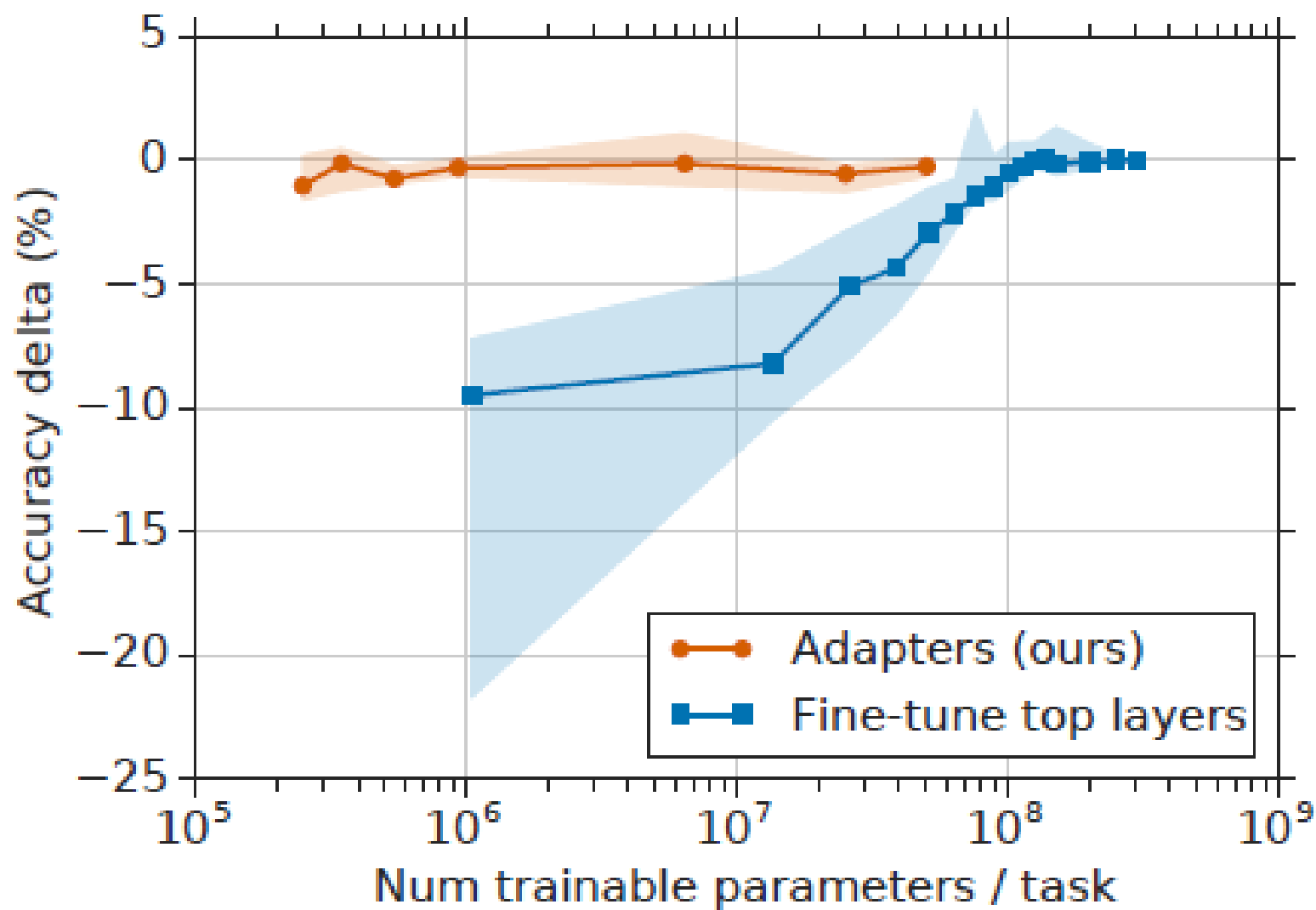
[Stickland, et al., ICML'19] [Houlsby, et al., ICML'19]





Source of image: <https://arxiv.org/abs/1902.00751>

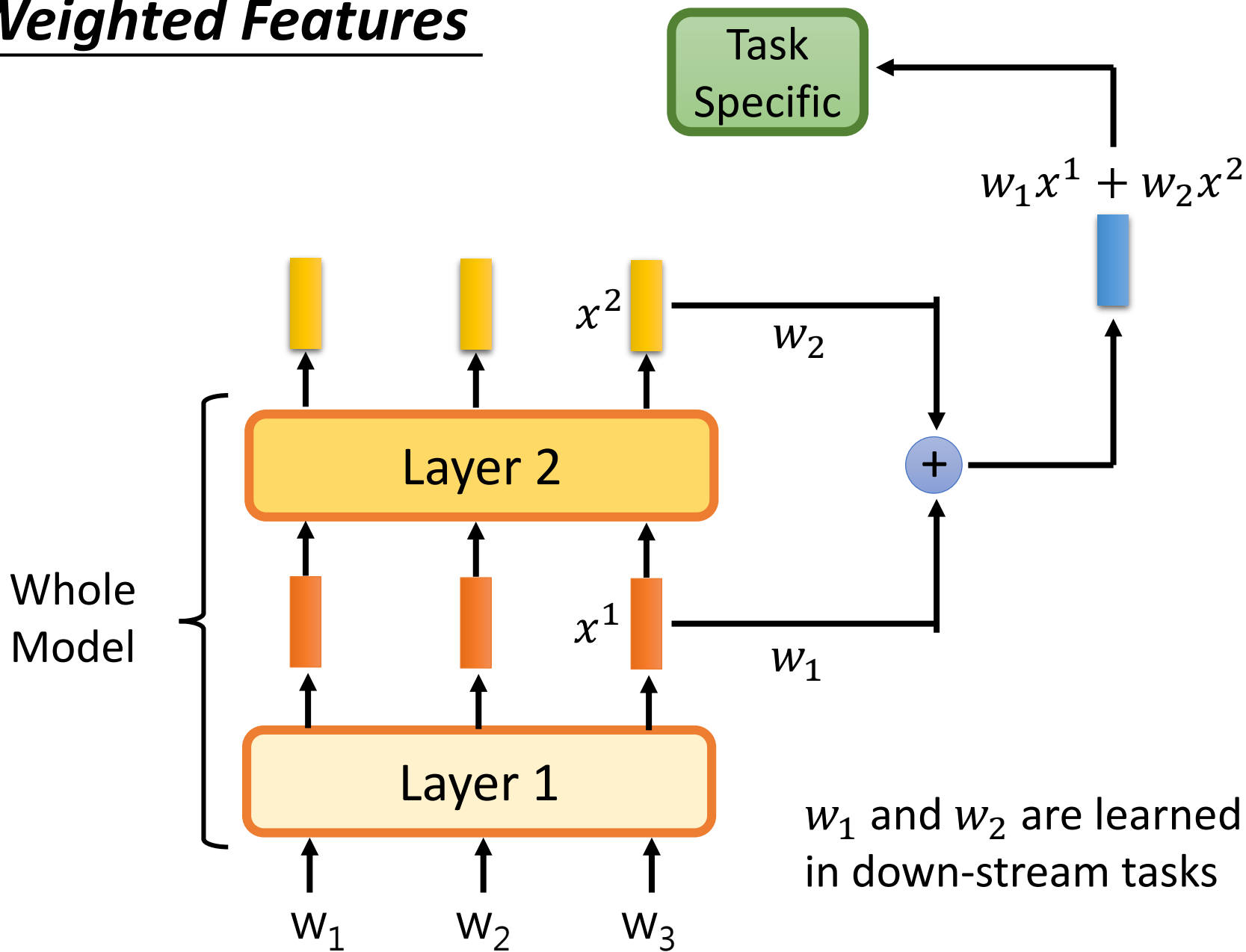
[Houlsby, et al., ICML'19]



Source of image: <https://arxiv.org/abs/1902.00751>

[Houlsby, et al., ICML'19]

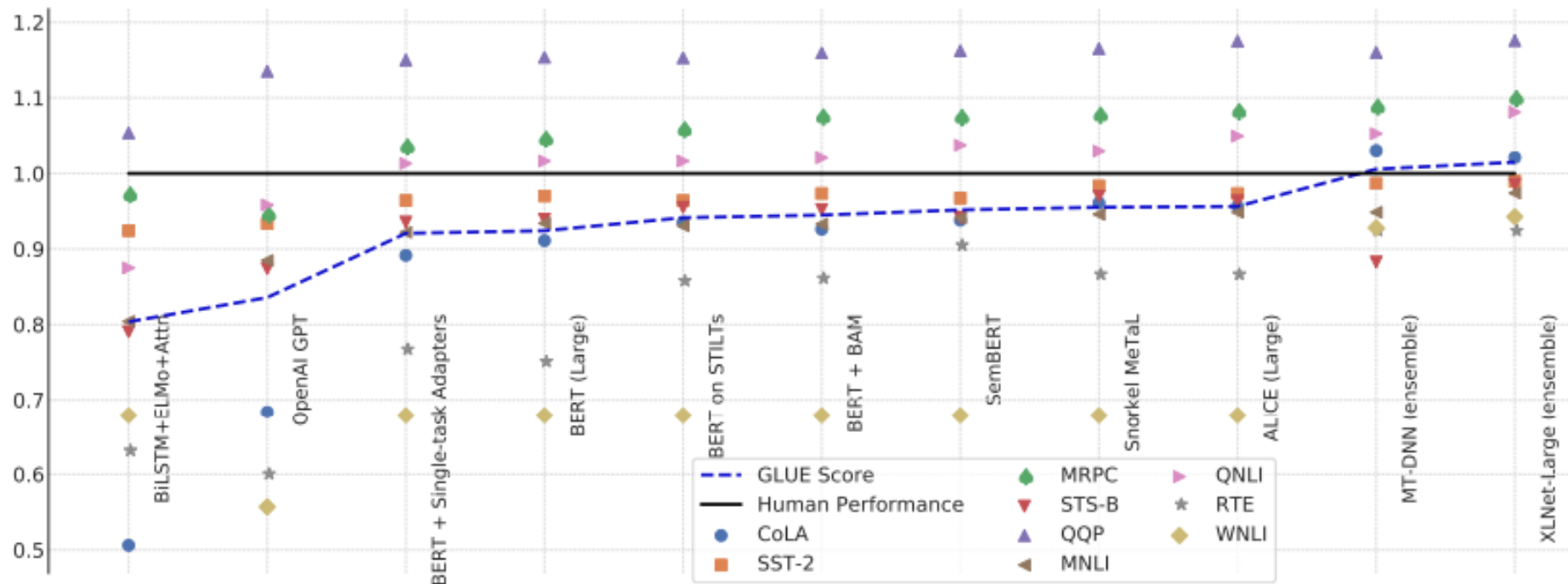
# Weighted Features





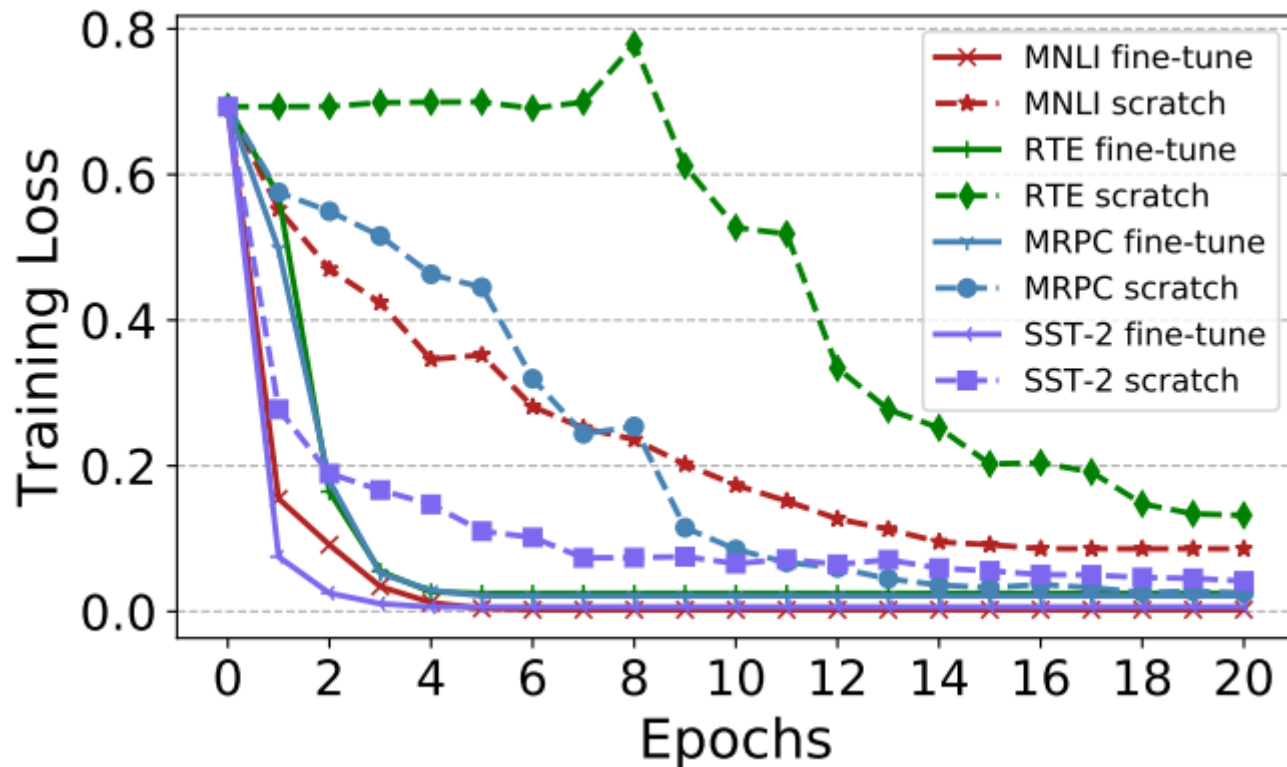
# Why Pre-train Models?

- GLUE scores



Source of image: <https://arxiv.org/abs/1905.00537>

# Why Fine-tune?



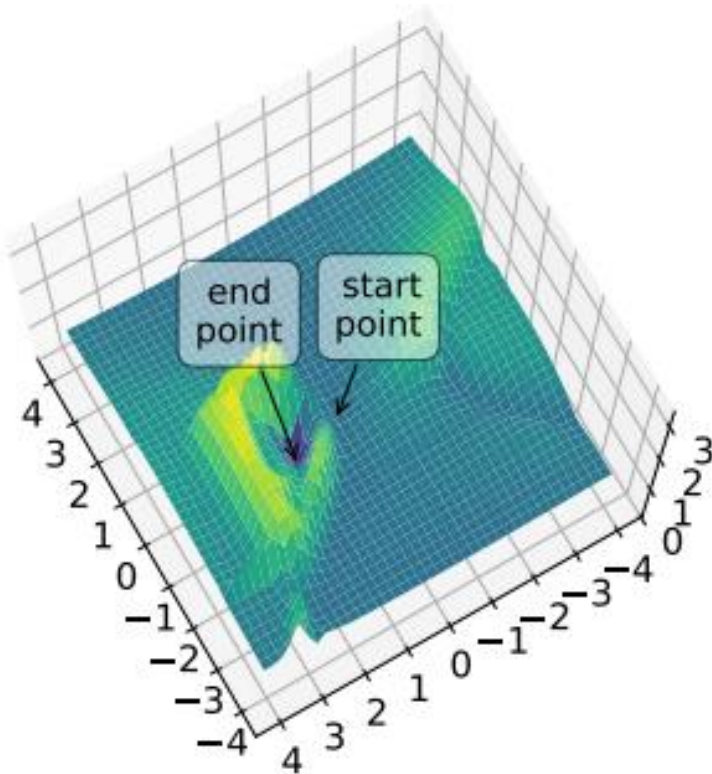
[Hao, et al., EMNLP'19] Source of image: <https://arxiv.org/abs/1908.05620>

# Why Fine-tune?

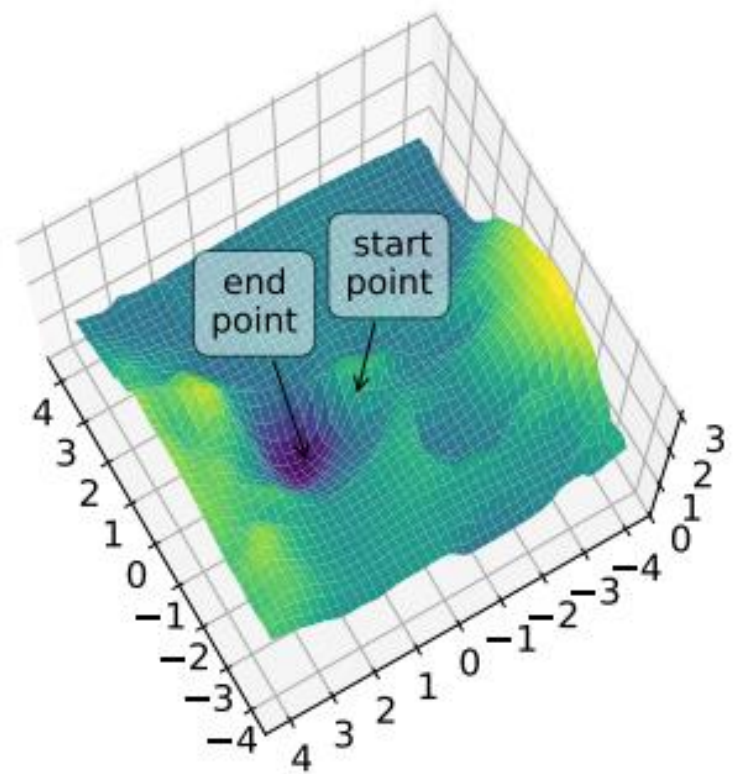
How to generate the figures below?

<https://youtu.be/XysGHdNOTbg>

Training from scratch



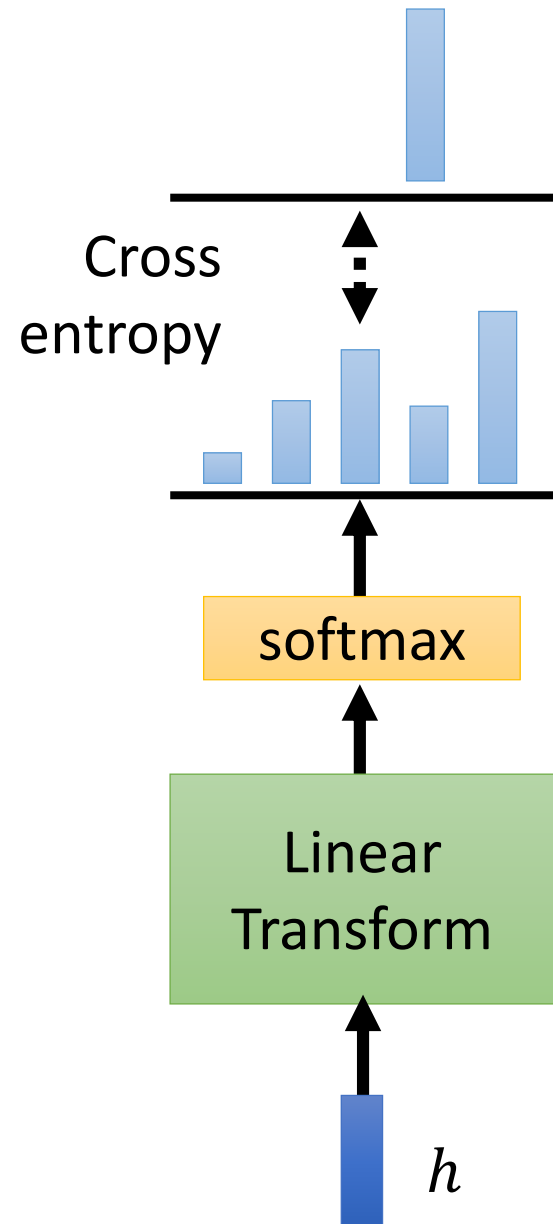
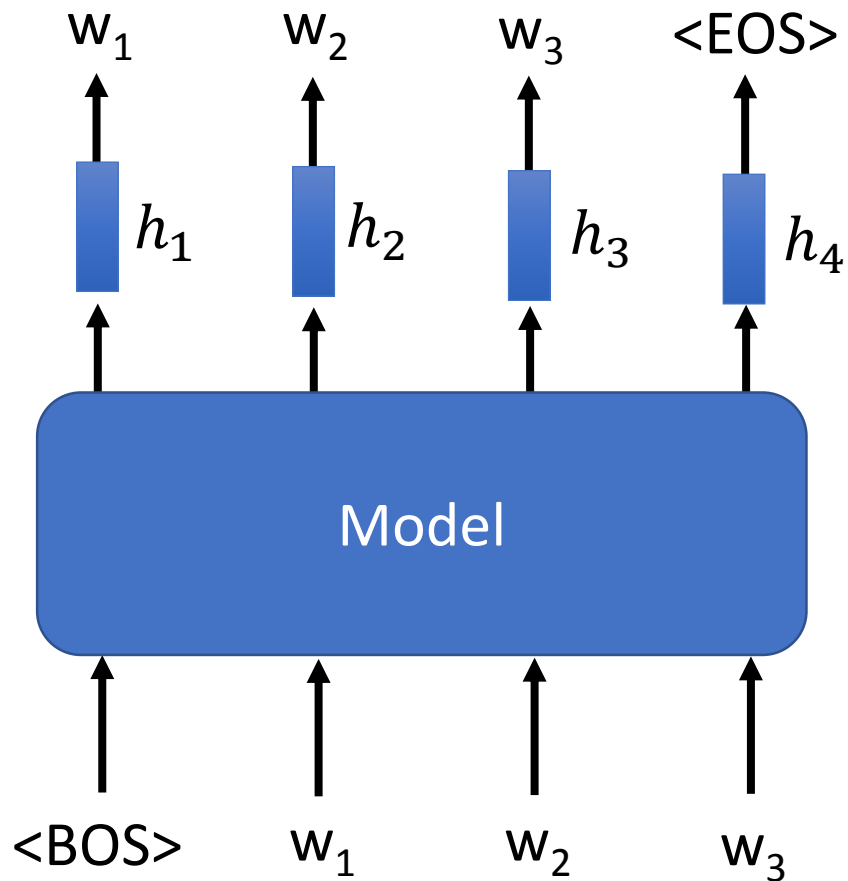
Fine-tuning BERT



[Hao, et al., EMNLP'19] Source of image: <https://arxiv.org/abs/1908.05620>

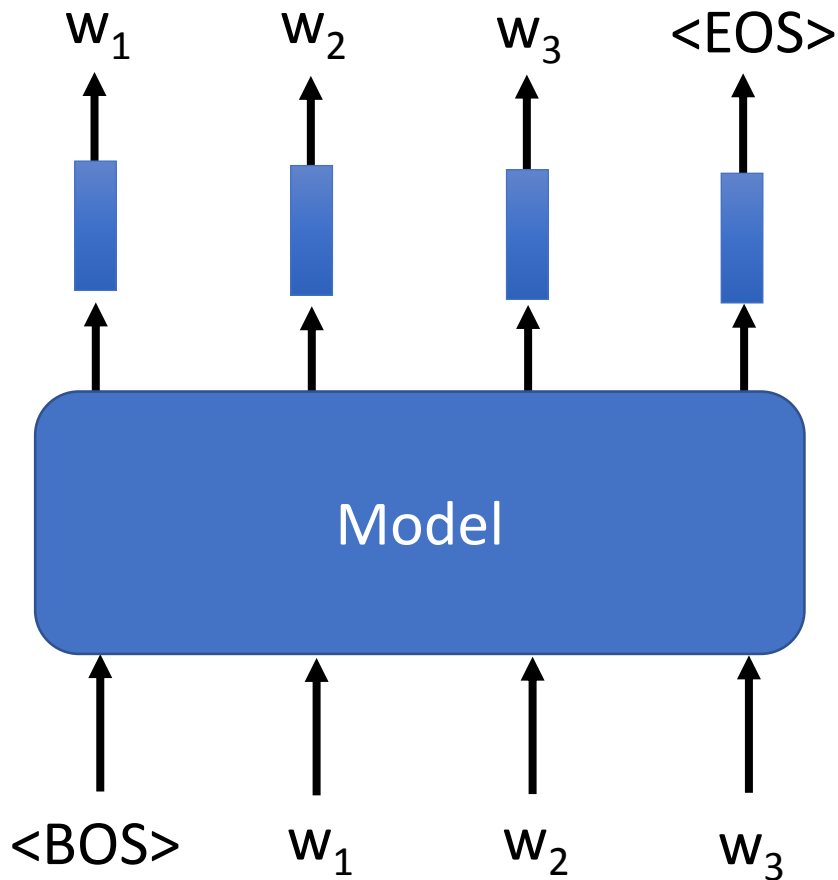
# Training Methods

# Predicting Next Token

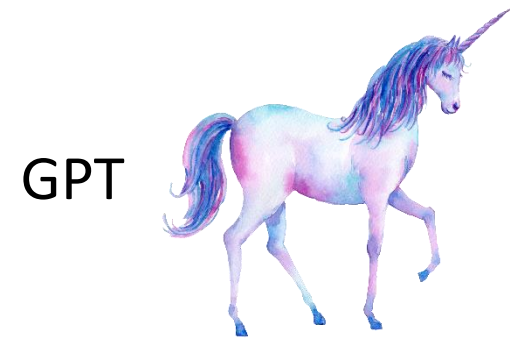


# Predicting Next Token

[Peters, et al., NAACL'18]



LSTM



Transformer-Decoder

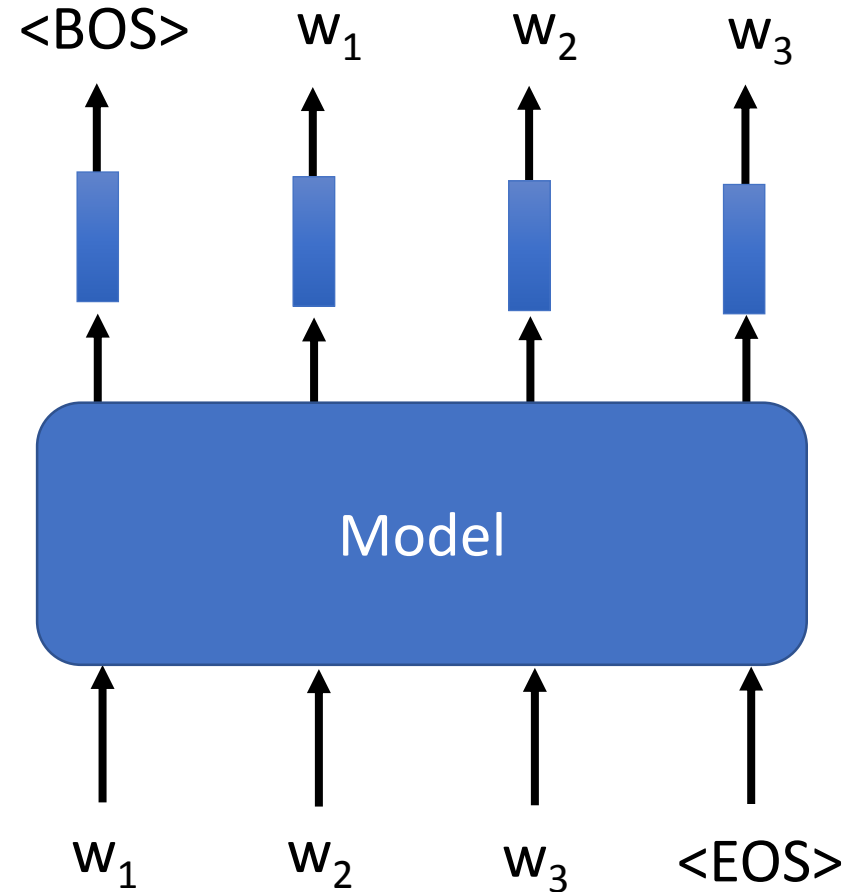
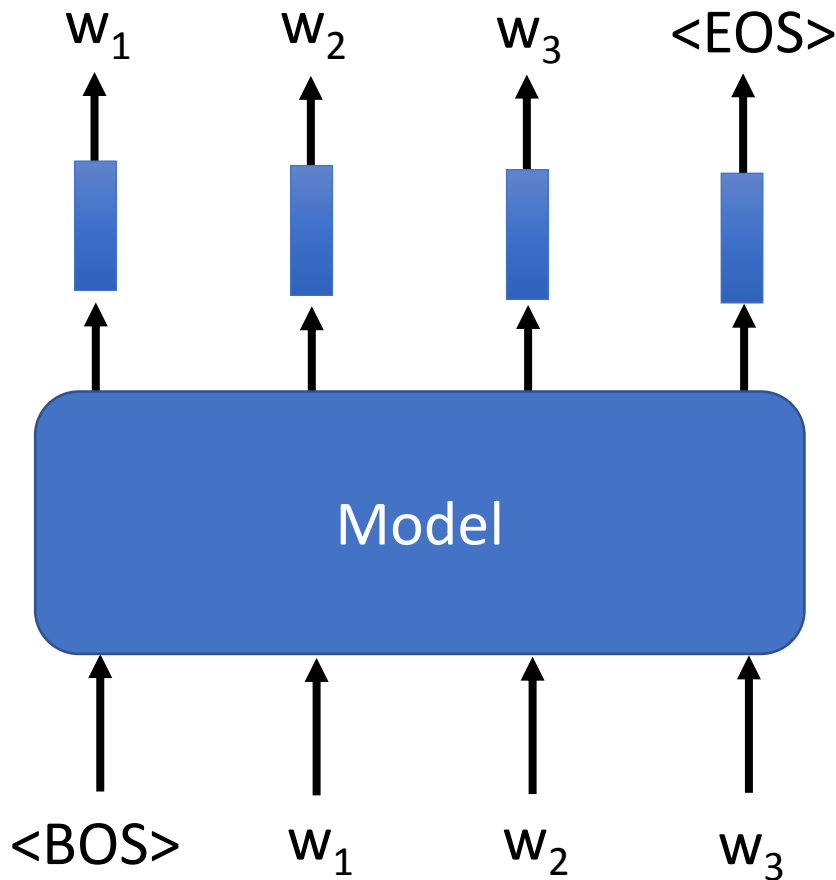
Don't use Transformer-Encoder

# Predicting Next Token

ELMO



LSTM

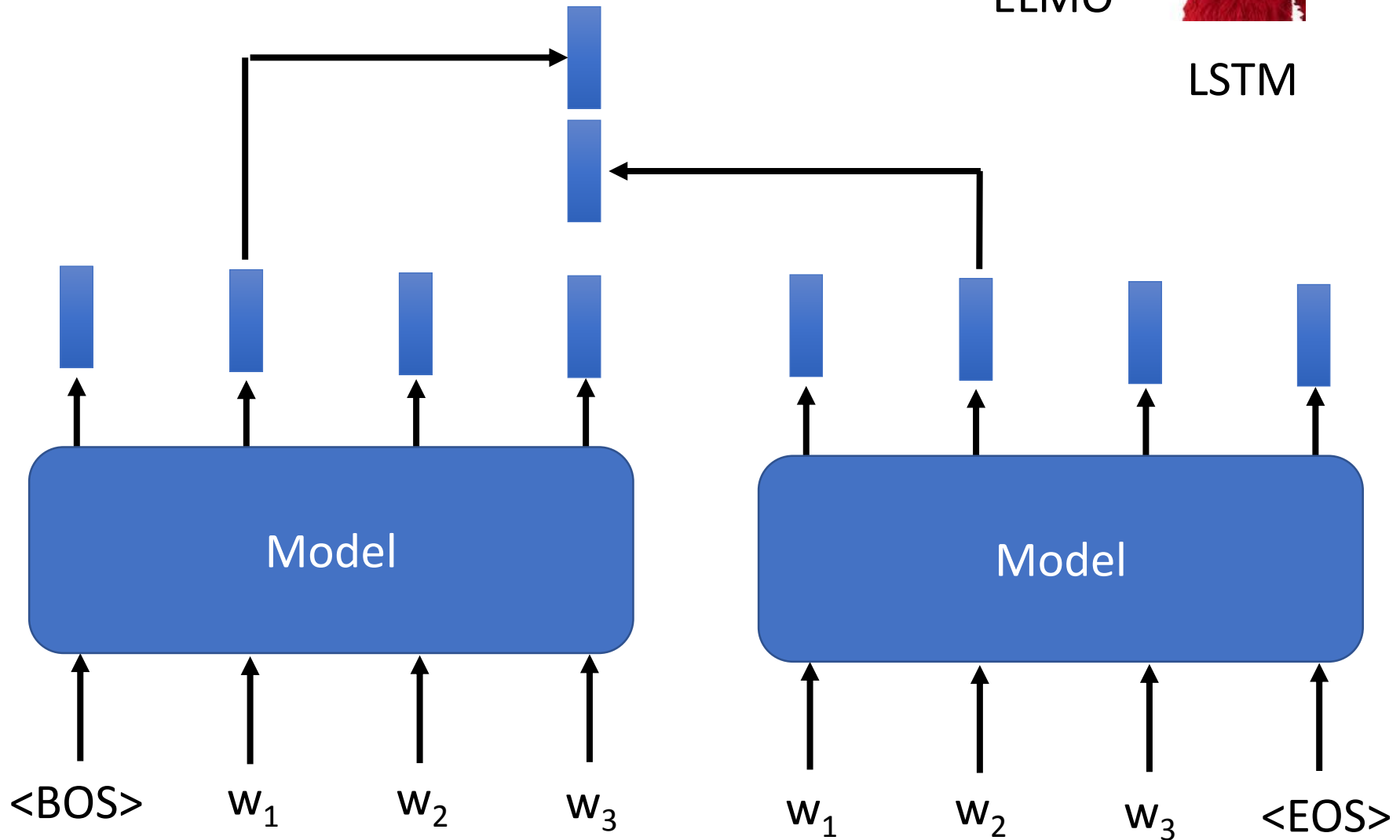


# Predicting Next Token

ELMO

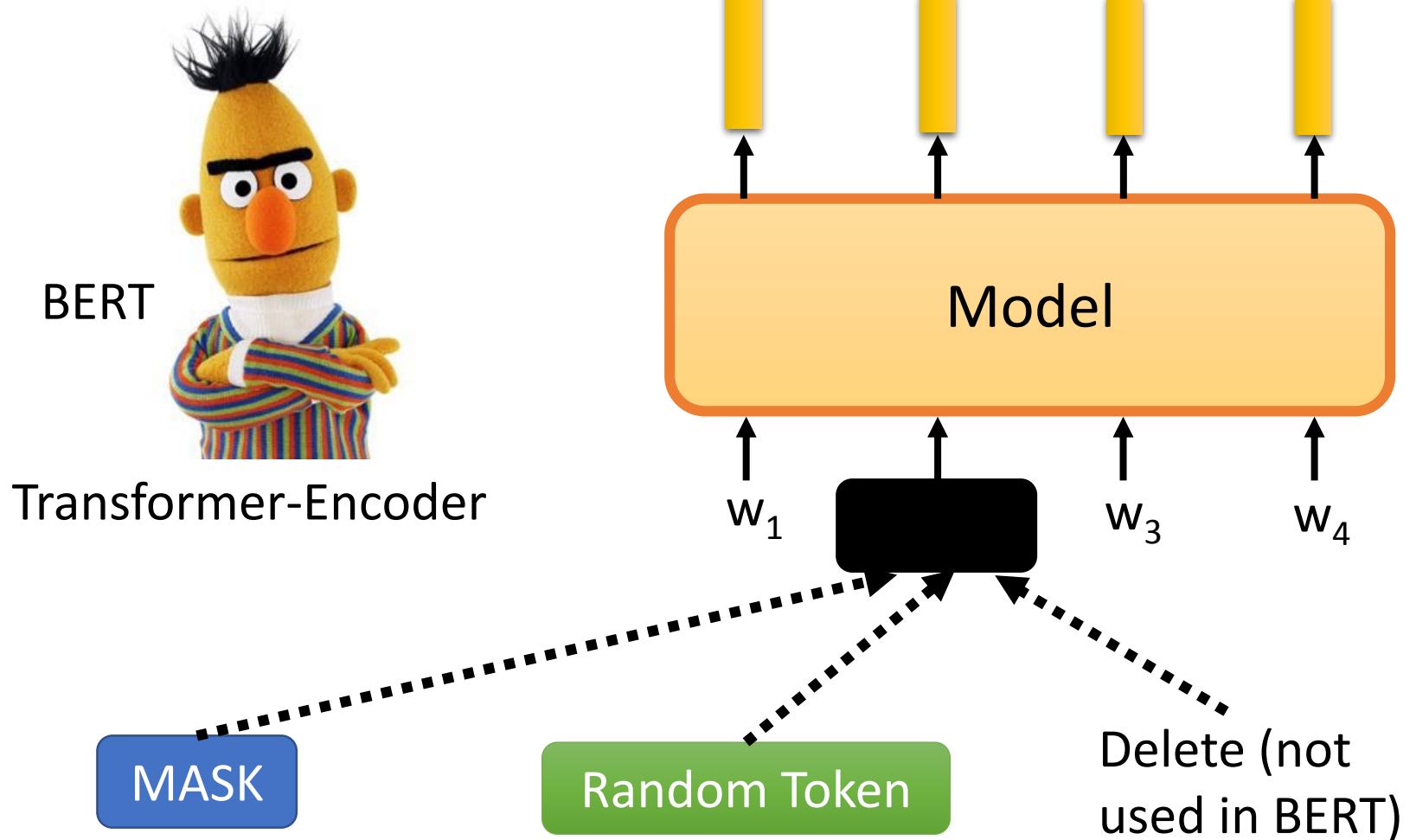


LSTM





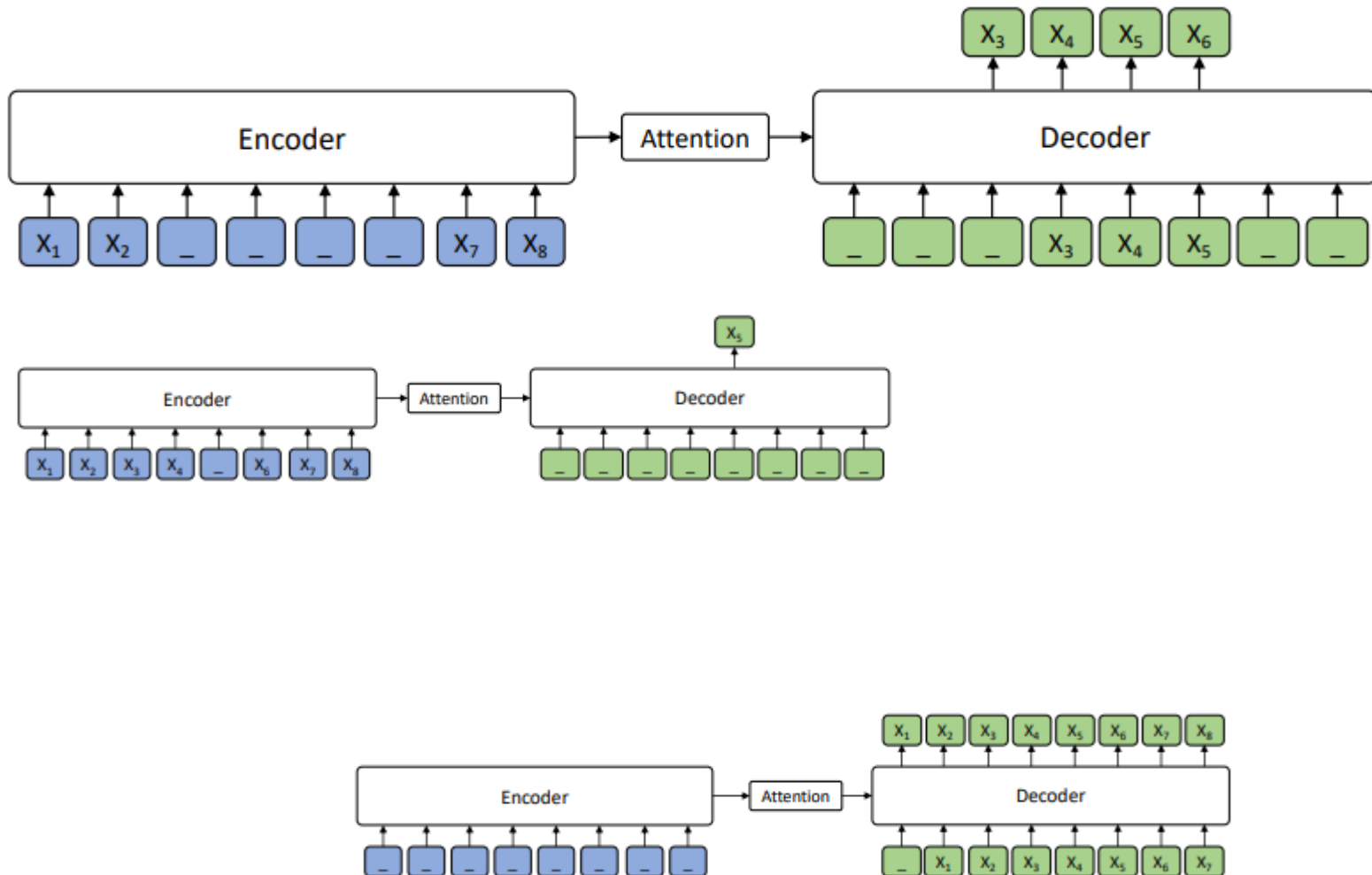
# Masking Input



# MASked Sequence to Sequence pre-training (MASS)

# MASS

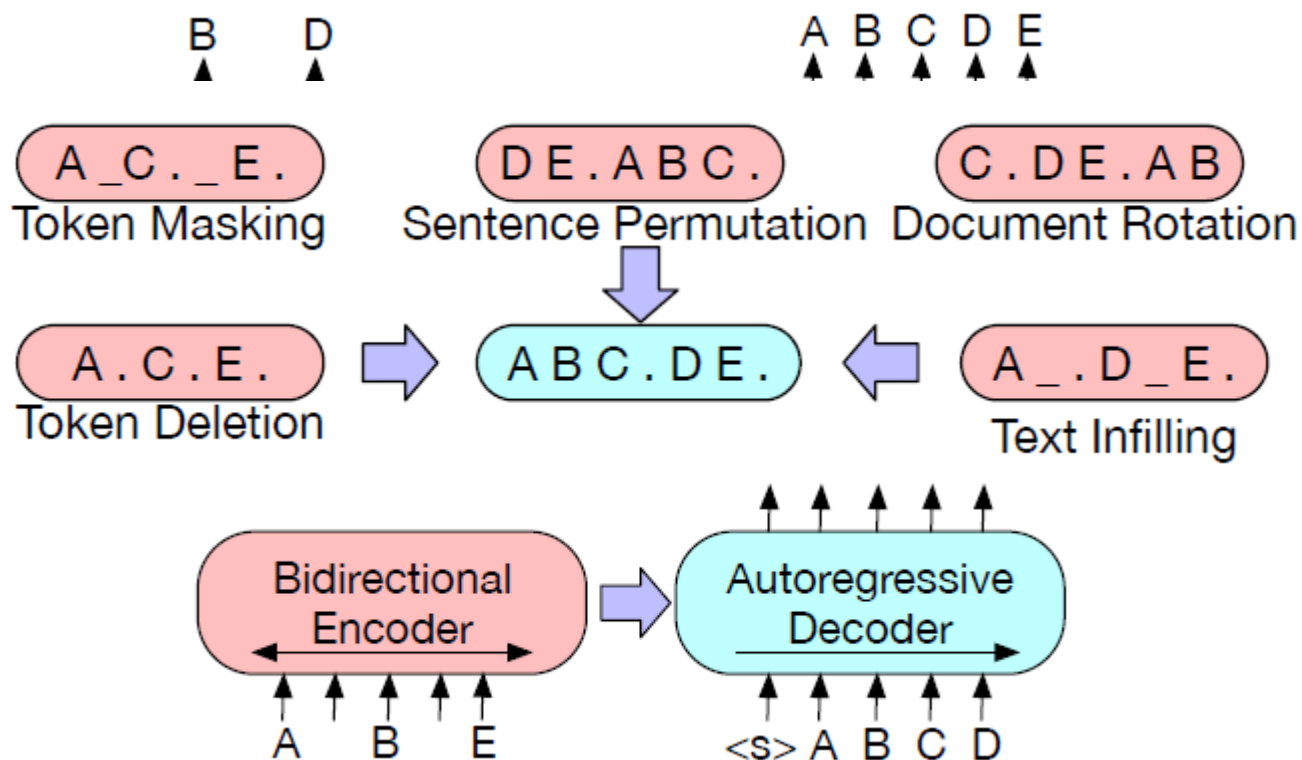
[Song, et al., ICML'19]



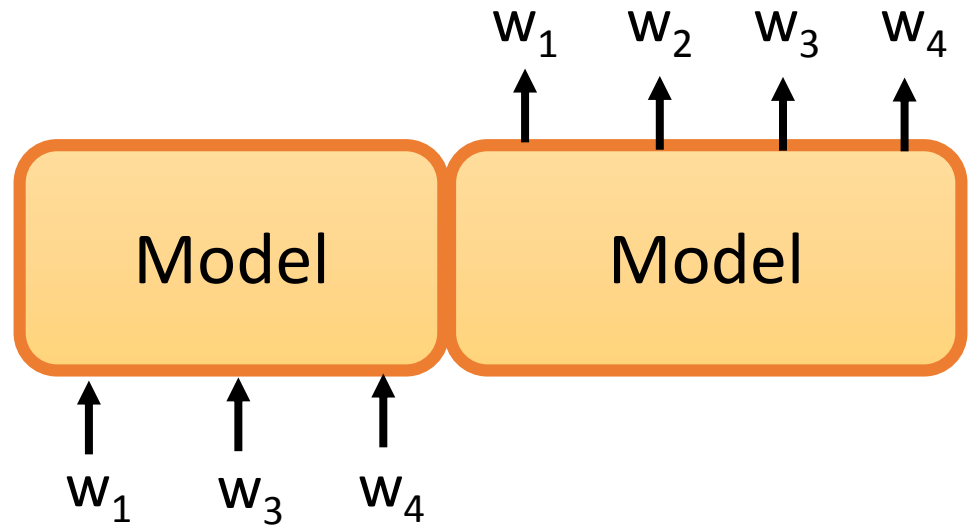
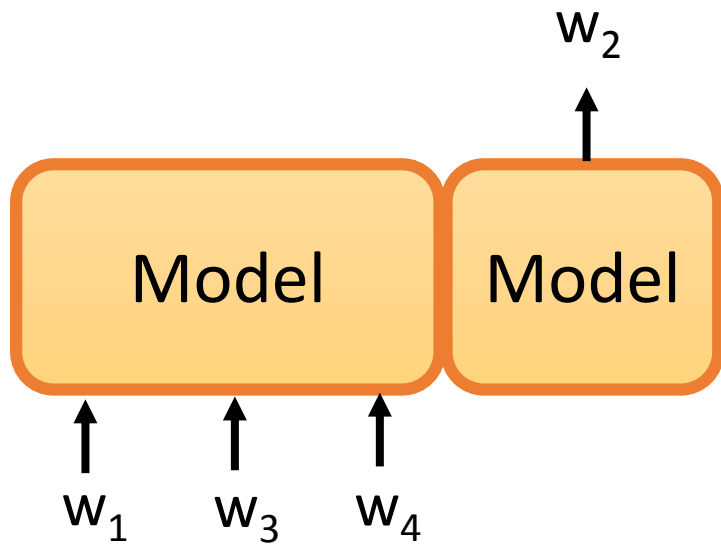
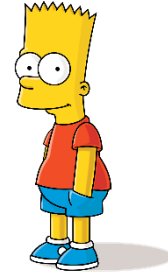
# BART

Bidirectional and Auto-Regressive  
Transformers

[Lewis, et al., arXiv'19]



# Masking Input



$w_2$  is deleted

# Masking Input

Is random masking  
good enough?

- Whole Word Masking (WWM)

<https://arxiv.org/abs/1906.08101>

---

**[Original BERT Input]**

使用语言 [MASK] 型来 [MASK] 测下一个词的 pro [MASK] ##lity 。

**[Whold Word Masking Input]**

使用语言 [MASK] [MASK] 来 [MASK] [MASK] 下一个词的 [MASK] [MASK] [MASK] 。

---

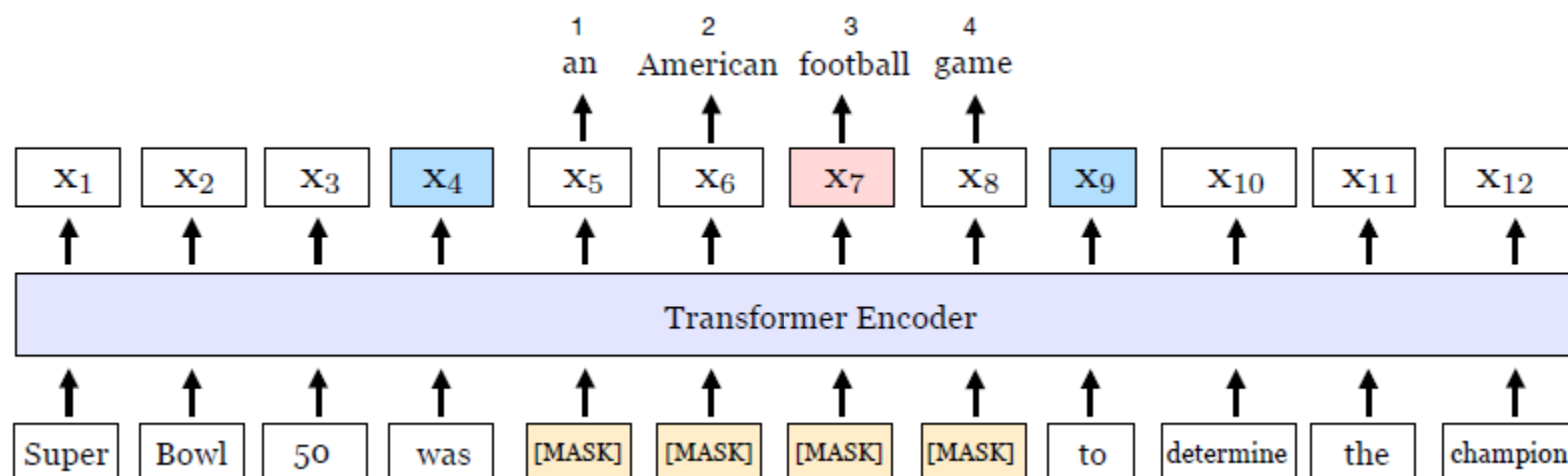
- ERNIE      Enhanced Representation through  
Knowledge Integration (ERNIE)



# SpanBert

NSP is not good

[Joshi, et al., TACL'20]



	SQuAD 2.0	NewsQA	TriviaQA	Coreference	MNLI-m	QNLI	GLUE (Avg)
Subword Tokens	83.8	72.0	76.3	<b>77.7</b>	86.7	92.5	83.2
Whole Words	84.3	72.8	77.1	76.6	86.3	92.8	82.9
Named Entities	84.8	72.7	78.7	75.6	86.0	93.1	83.2
Noun Phrases	85.0	<b>73.0</b>	77.7	76.7	86.5	93.2	<b>83.5</b>
Geometric Spans	<b>85.4</b>	<b>73.0</b>	<b>78.8</b>	76.4	<b>87.0</b>	<b>93.3</b>	83.4

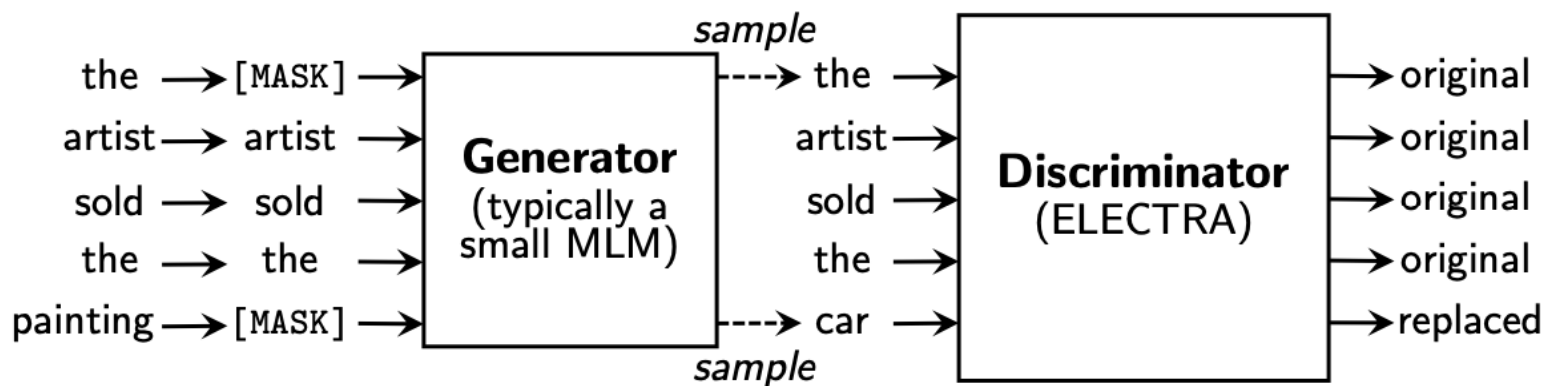
# Discriminator

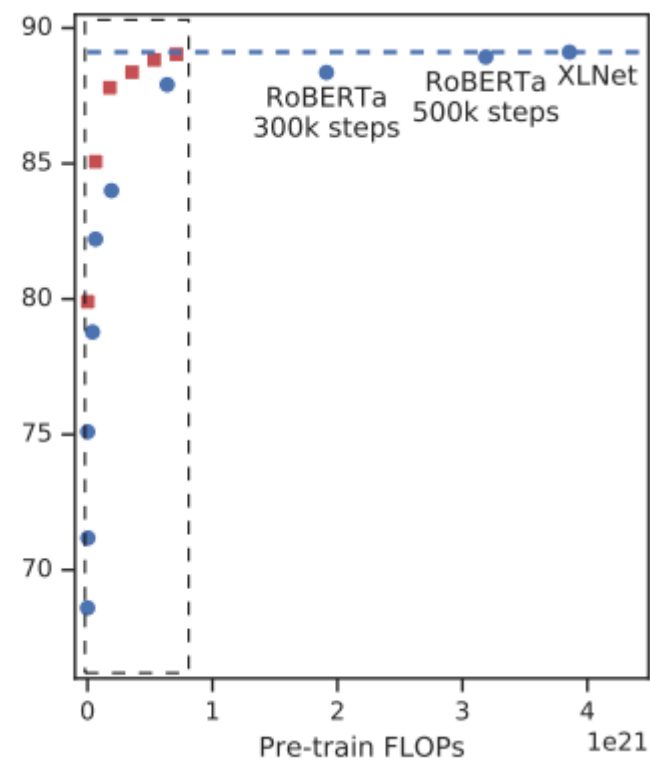
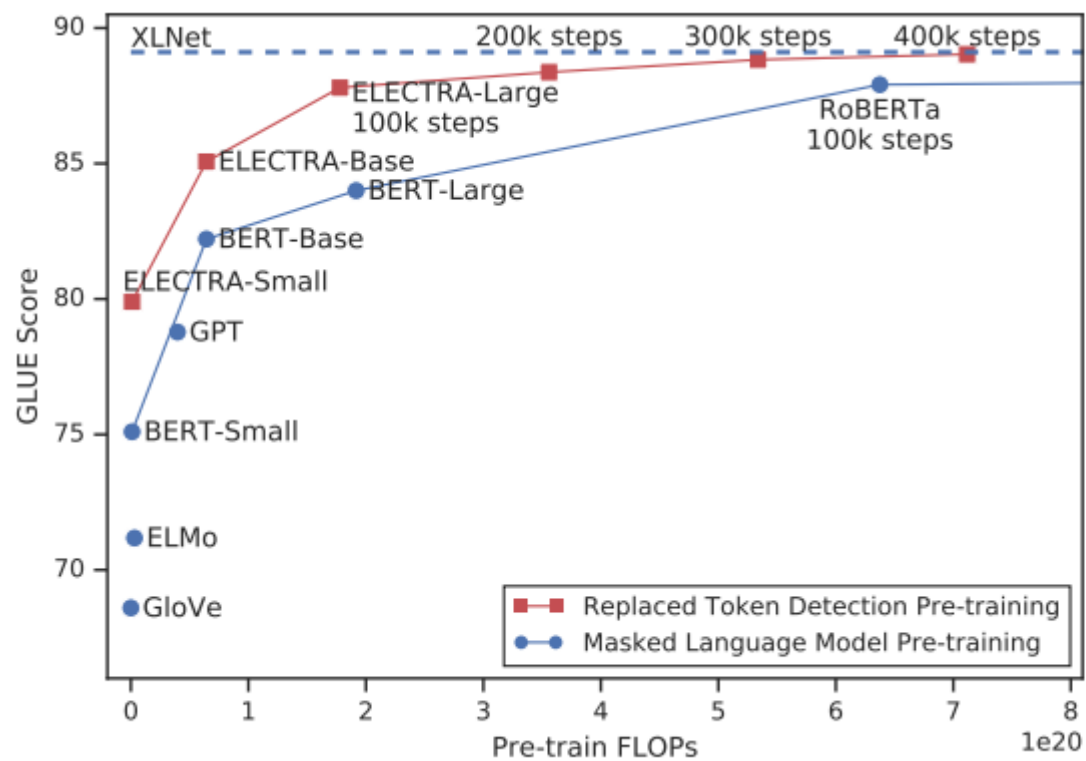
Efficiently Learning an Encoder that Classifies  
Token Replacements Accurately.”

<https://arxiv.org/abs/2003.10555>



ELECTRA



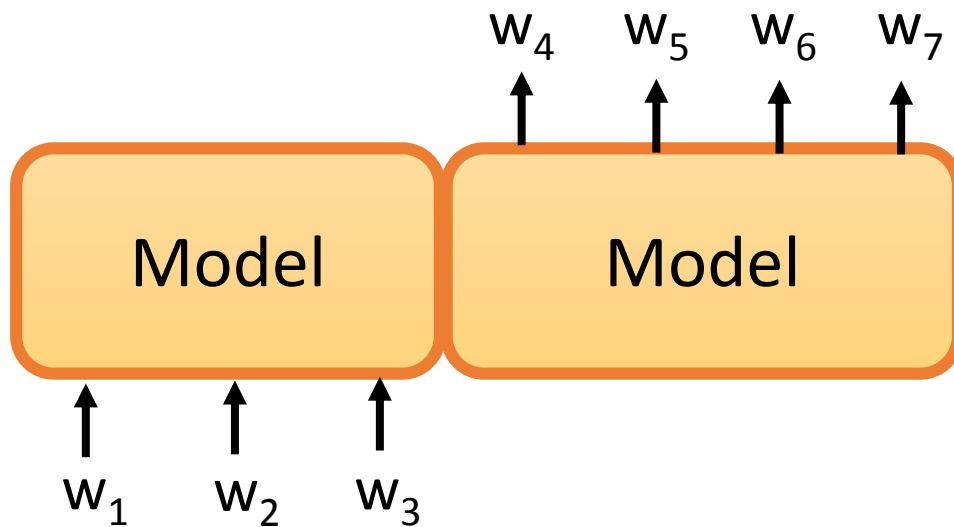






# Sentence Level

- Input one sentence, predict another



COVE:

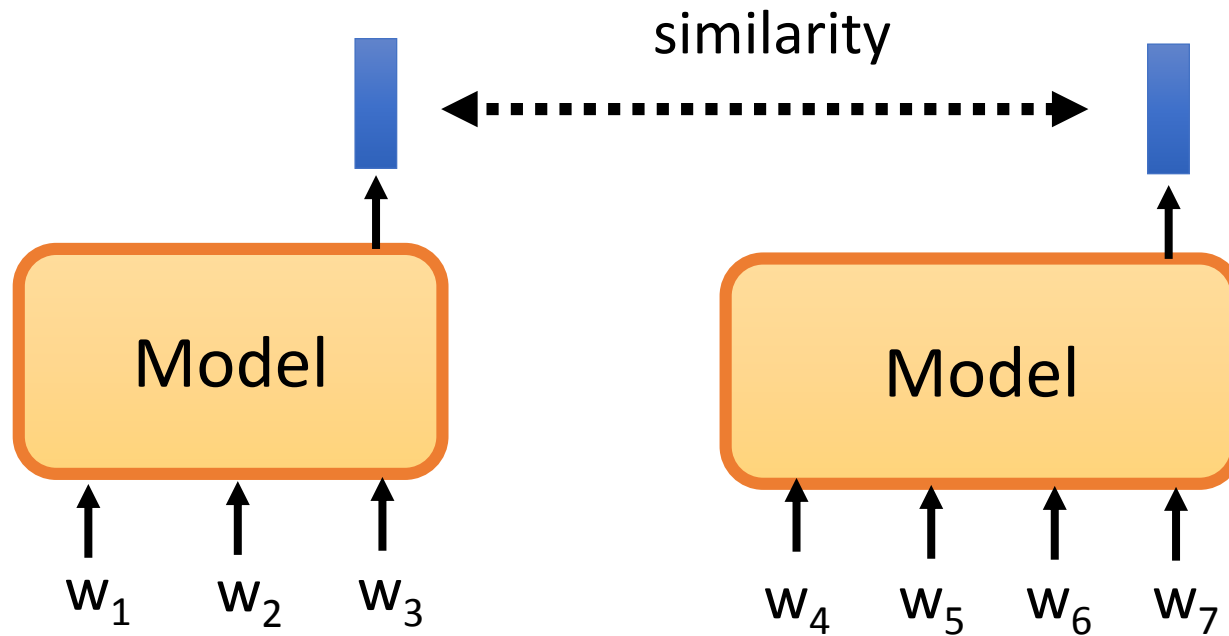
Input: A language  
output: B language

Skip-Thought

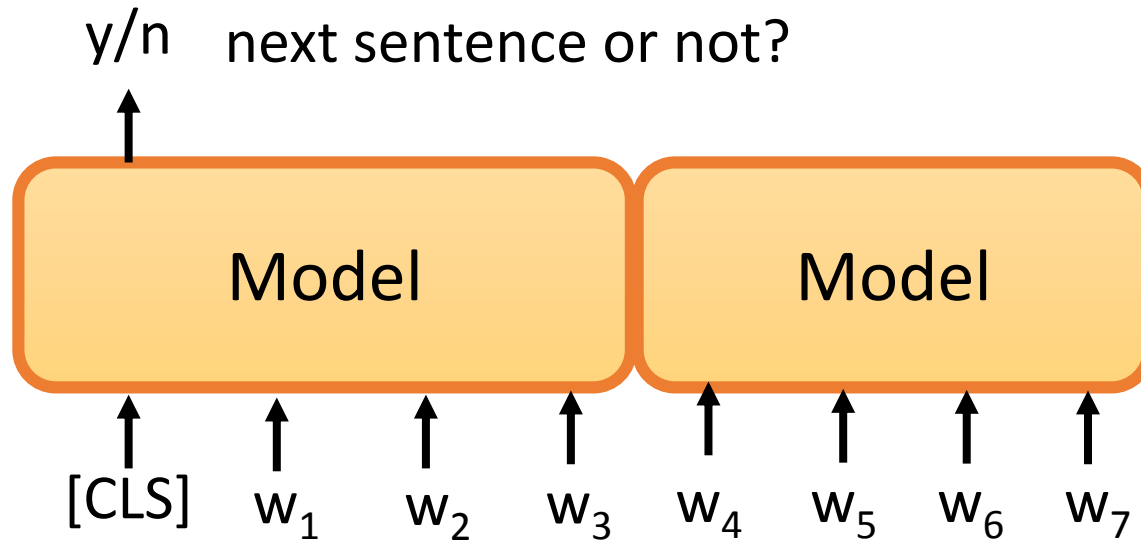
Input: previous sentence  
output: next sentence

# Sentence Level

- Quick Thought



# Sentence Level



*Robustly optimized BERT approach  
(ROBERTA)*

[Liu, et al., arXiv'19]



# Concluding Remarks

# Reference

- [Lewis, et al., arXiv'19] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, Luke Zettlemoyer, BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension, arXiv, 2019
- [Raffel, et al., arXiv'19] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, Peter J. Liu, Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer, arXiv, 2019
- [Joshi, et al., TACL'20] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, Omer Levy, SpanBERT: Improving Pre-training by Representing and Predicting Spans, TACL, 2020
- [Song, et al., ICML'19] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, Tie-Yan Liu, MASS: Masked Sequence to Sequence Pre-training for Language Generation, ICML, 2019
- [Zafrir, et al., NeurIPS workshop 2019] Ofir Zafrir, Guy Boudoukh, Peter Izsak, Moshe Wasserblat, Q8BERT: Quantized 8Bit BERT, NeurIPS workshop 2019

# Reference

- [Houlsby, et al., ICML'19] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin de Laroussilhe, Andrea Gesmundo, Mona Attariyan, Sylvain Gelly, Parameter-Efficient Transfer Learning for NLP, ICML, 2019
- [Hao, et al., EMNLP'19] Yaru Hao, Li Dong, Furu Wei, Ke Xu, Visualizing and Understanding the Effectiveness of BERT, EMNLP, 2019
- [Liu, et al., arXiv'19] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, Veselin Stoyanov, RoBERTa: A Robustly Optimized BERT Pretraining Approach, arXiv, 2019
- [Sanh, et al., NeurIPS workshop's] Victor Sanh, Lysandre Debut, Julien Chaumond, Thomas Wolf, DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter, NeurIPS workshop, 2019
- [Jian, et al., arXiv'19] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, Qun Liu, TinyBERT: Distilling BERT for Natural Language Understanding, arXiv, 19

# Reference

- [Shoeybi, et al., arXiv'19] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, Bryan Catanzaro, Megatron-LM: Training Multi-Billion Parameter Language Models Using Model Parallelism, arXiv, 19
- [Lan, et al., ICLR'20] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, Radu Soricut, ALBERT: A Lite BERT for Self-supervised Learning of Language Representations, ICLR, 2020
- [Kitaev, et al., ICLR'20] Nikita Kitaev, Lukasz Kaiser, Anselm Levskaya, Reformer: The Efficient Transformer, ICLR, 2020
- [Beltagy, et al., arXiv'20] Iz Beltagy, Matthew E. Peters, Arman Cohan, Longformer: The Long-Document Transformer, arXiv, 2020
- [Dai, et al., ACL'19] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V. Le, Ruslan Salakhutdinov, Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context, ACL, 2019
- [Peters, et al., NAACL'18] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, Luke Zettlemoyer, Deep contextualized word representations, NAACL, 2018

# Reference

- [Sanh, et al., NeurIPS workshop's] Victor Sanh, Lysandre Debut, Julien Chaumond, Thomas Wolf, DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter, NeurIPS workshop, 2019
- [Jian, et al., arXiv'19] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, Qun Liu, TinyBERT: Distilling BERT for Natural Language Understanding, arXiv, 19
- [Sun, et al., ACL'20] Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, Denny Zhou, MobileBERT: a Compact Task-Agnostic BERT for Resource-Limited Devices, ACL, 2020
- [Zafrir, et al., NeurIPS workshop 2019] Ofir Zafrir, Guy Boudoukh, Peter Izsak, Moshe Wasserblat, Q8BERT: Quantized 8Bit BERT, NeurIPS workshop 2019
- [Sun, et al., ACL'20] Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, Denny Zhou, MobileBERT: a Compact Task-Agnostic BERT for Resource-Limited Devices, ACL, 2020



# Reference

- [Pennington, et al., EMNLP'14] Jeffrey Pennington, Richard Socher, Christopher Manning, Glove: Global Vectors for Word Representation, EMNLP, 2014
- [Mikolov, et al., NIPS'13] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S. Corrado, Jeff Dean, Distributed Representations of Words and Phrases and their Compositionality, NIPS, 2013
- [Bojanowski, et al., TACL'17] Piotr Bojanowski, Edouard Grave, Armand Joulin, Tomas Mikolov, Enriching Word Vectors with Subword Information, TACL, 2017
- [Su, et al., EMNLP'17] Tzu-Ray Su, Hung-Yi Lee, Learning Chinese Word Representations From Glyphs Of Characters, EMNLP, 2017
- [Liu, et al., ACL'19] Xiaodong Liu, Pengcheng He, Weizhu Chen, Jianfeng Gao, Multi-Task Deep Neural Networks for Natural Language Understanding, ACL, 2019
- [Stickland, et al., ICML'19] Asa Cooper Stickland, Iain Murray, BERT and PALs: Projected Attention Layers for Efficient Adaptation in Multi-Task Learning, ICML, 2019