

# Natural Language Processing (9)

## Transfer Learning

Daisuke Kawahara

Department of Communications and Computer Engineering,  
Waseda University

# Lecture Plan

1. Overview of Natural Language Processing
2. Formal Language Theory
3. Word Senses and Embeddings
4. Topic Models
5. Collocations, Language Models, and Recurrent Neural Networks
6. Sequence Labeling and Morphological Analysis
7. Parsing (1): Constituency Parsing
8. Parsing (2): Dependency Parsing
9. Transfer Learning
10. Knowledge Acquisition
11. Information Retrieval, Question Answering, and Machine Translation
12. Guest Talk (1)
13. Guest Talk (2)
14. Project: Survey or Programming
15. Project Presentation

The screenshot shows a web browser window with the Google AI Blog header. The main content is an article titled "Open Sourcing BERT: State-of-the-Art Pre-training for Natural Language Processing" posted on Friday, November 2, 2018. The article discusses the challenges of training data in NLP and introduces BERT, a new pre-training technique. It also covers what makes BERT different from previous models like ELMo and ULMFit, and why it's contextually aware. On the right side, there's a sidebar with a search bar, links for Labels, Archive, and Feed, a Twitter follow button, and a link to give feedback in Product Forums.

Google AI Blog

The latest news from Google AI

## Open Sourcing BERT: State-of-the-Art Pre-training for Natural Language Processing

Friday, November 2, 2018

Posted by Jacob Devlin and Ming-Wei Chang, Research Scientists, Google AI Language

One of the biggest challenges in [natural language processing](#) (NLP) is the shortage of training data. Because NLP is a diversified field with many distinct tasks, most task-specific datasets contain only a few thousand or a few hundred thousand human-labeled training examples. However, modern deep learning-based NLP models see benefits from much larger amounts of data, improving when trained on millions, or *billions*, of annotated training examples. To help close this gap in data, researchers have developed a variety of techniques for training general purpose language representation models using the enormous amount of unannotated text on the web (known as *pre-training*). The pre-trained model can then be *fine-tuned* on small-data NLP tasks like [question answering](#) and [sentiment analysis](#), resulting in substantial accuracy improvements compared to training on these datasets from scratch.

This week, we [open sourced](#) a new technique for NLP pre-training called Bidirectional Encoder Representations from [Transformers](#), or **BERT**. With this release, anyone in the world can train their own state-of-the-art question answering system (or a variety of other models) in about 30 minutes on a single [Cloud TPU](#), or in a few hours using a single GPU. The release includes source code built on top of [TensorFlow](#) and a number of pre-trained language representation models. In our [associated paper](#), we demonstrate state-of-the-art results on 11 NLP tasks, including the very competitive [Stanford Question Answering Dataset](#) (SQuAD v1.1).

**What Makes BERT Different?**

BERT builds upon recent work in pre-training contextual representations – including [Semi-supervised Sequence Learning](#), [Generative Pre-Training](#), [ELMo](#), and [ULMFit](#). However, unlike these previous models, BERT is the first *deeply bidirectional, unsupervised* language representation, pre-trained using only a plain text corpus (in this case, [Wikipedia](#)).

Why does this matter? Pre-trained representations can either be *context-free* or *contextual*, and *contextual* representations can further be *unidirectional* or *bidirectional*. Context-free models such as [word2vec](#) or [GloVe](#) generate a single [word embedding](#) representation for each word in the vocabulary. For example, the word "bank" would have the same context-free representation in "bank account" and "bank of the river." Contextual models instead generate a representation of each word that is based on the other words in the sentence. For example, in the sentence "*I accessed the bank account,*" a unidirectional contextual model would represent "bank" based on "*I accessed the*" but not "*account.*" However, BERT represents "bank" using both its previous and next context – "*I accessed the ... account*" – starting from the very bottom of a deep neural network, making it *deeply bidirectional*.

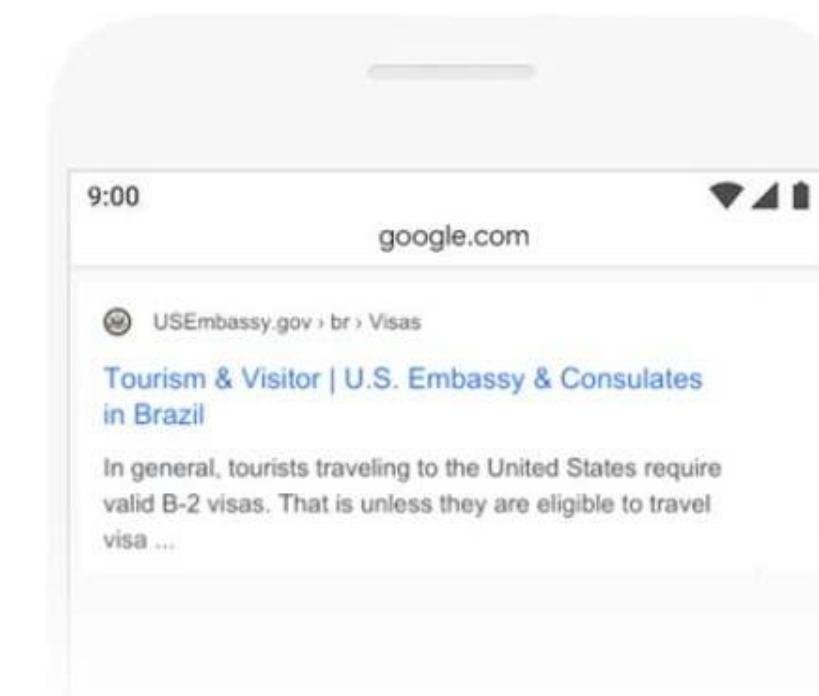
<https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html>

A screenshot of a web browser window. The address bar shows the URL [blog.google/products/search/search-language-understanding-bert/](https://blog.google/products/search/search-language-understanding-bert/). The page content is a Google search results page for the query "2019 brazil traveler to usa need a visa". The search bar at the top contains the same query. Below the search bar, there are two sections labeled "BEFORE" and "AFTER", each showing a smartphone screen displaying a news article from the Washington Post.

## BEFORE



## AFTER



This breakthrough was the result of Google research on transformers: models that process words in relation to all the other words in a sentence, rather than one-by-one in order. BERT models can therefore consider the full context of a word by looking at the words that come before and after it—particularly useful for understanding the intent behind search queries.

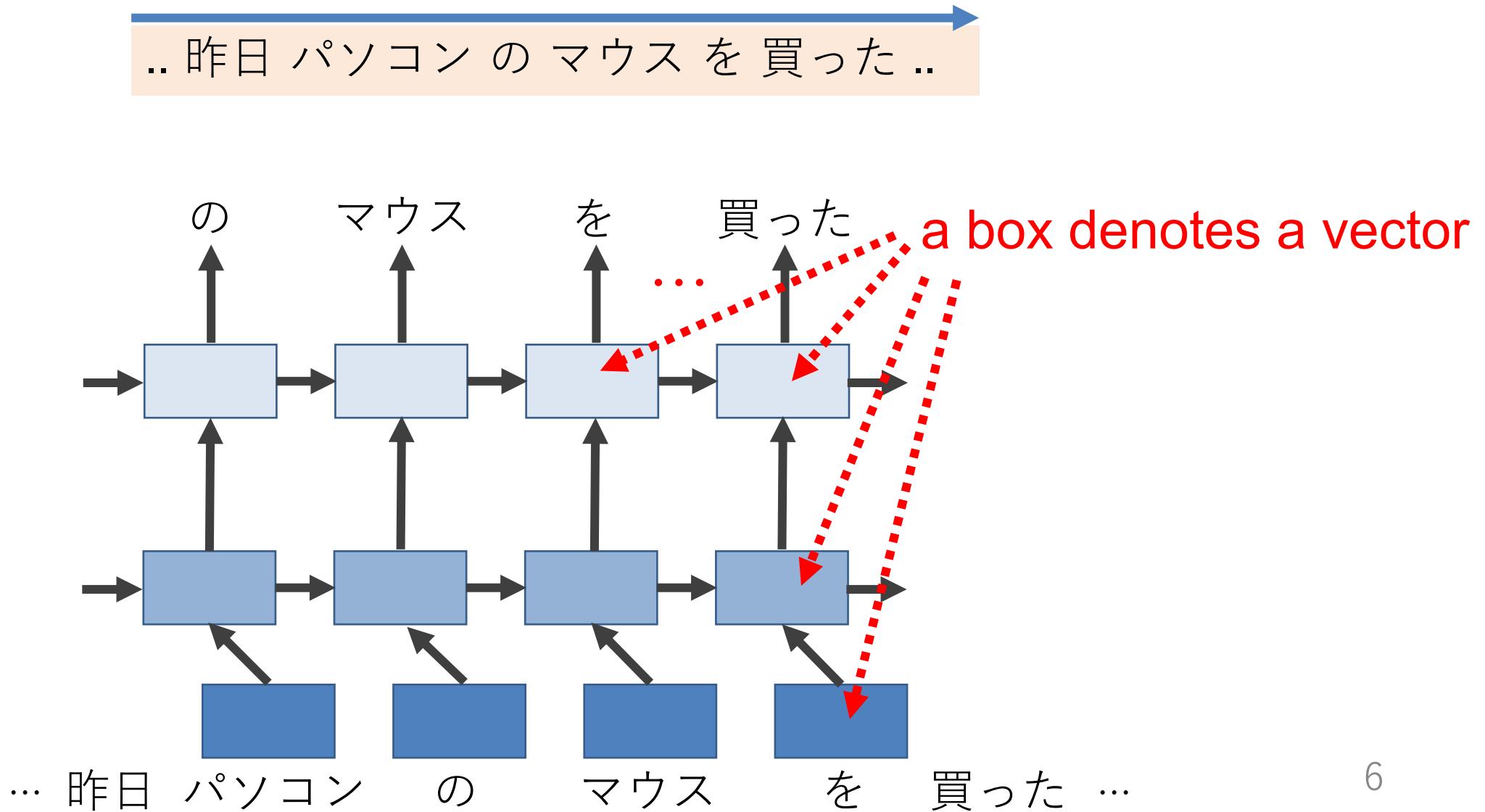
<https://www.blog.google/products/search/search-language-understanding-bert/>

# ELMo [Peters+ 2018]

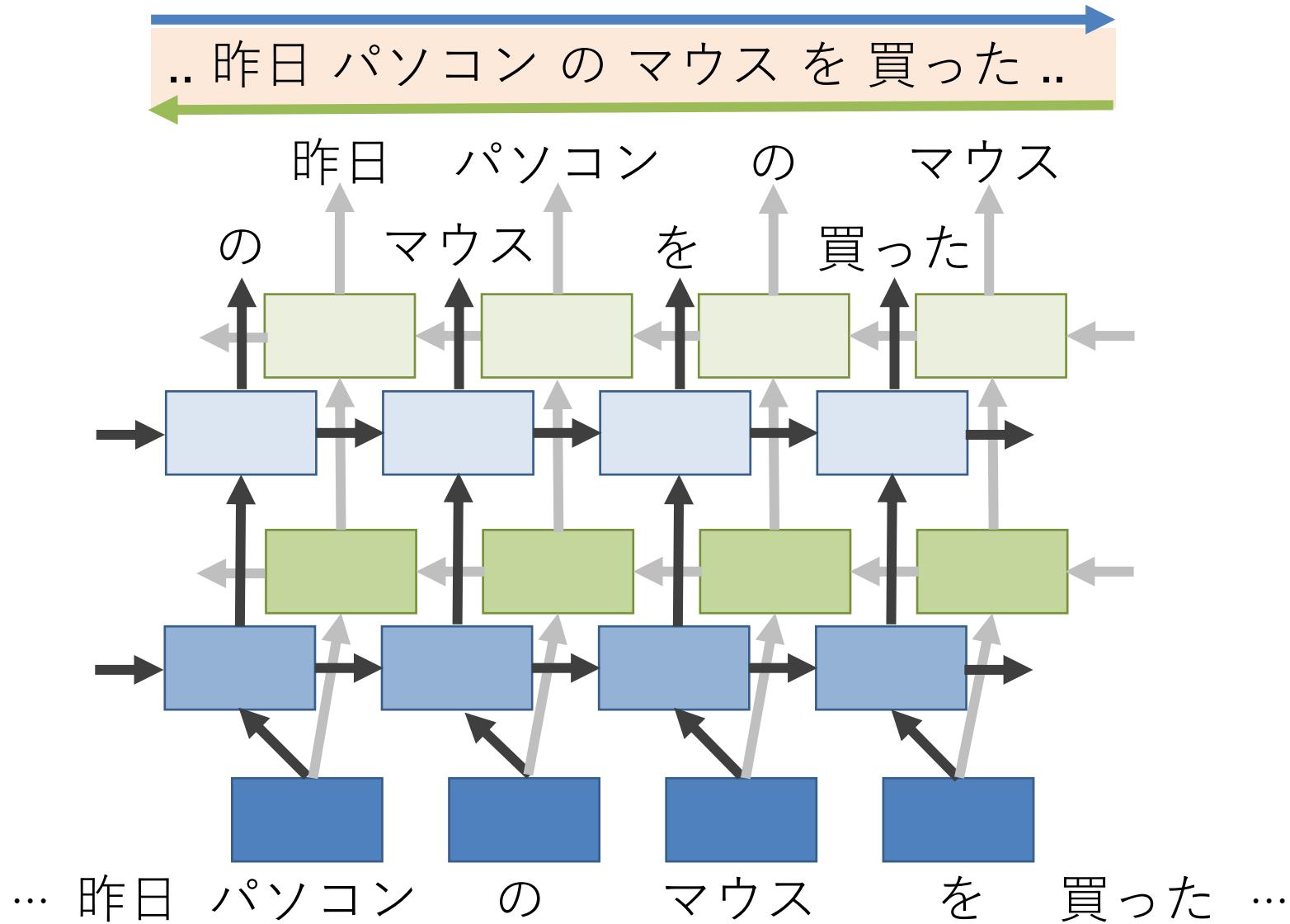
## Embeddings from Language Models

- Purpose
  1. To learn word meanings according to a context
  2. To learn language representations that consider syntax and semantics
- Method
  - Learn a bidirectional language model from a large corpus to obtain language representations (pre-training)
    - Note that language model learning is not a purpose
- Result
  - Achieved SOTA on six tasks

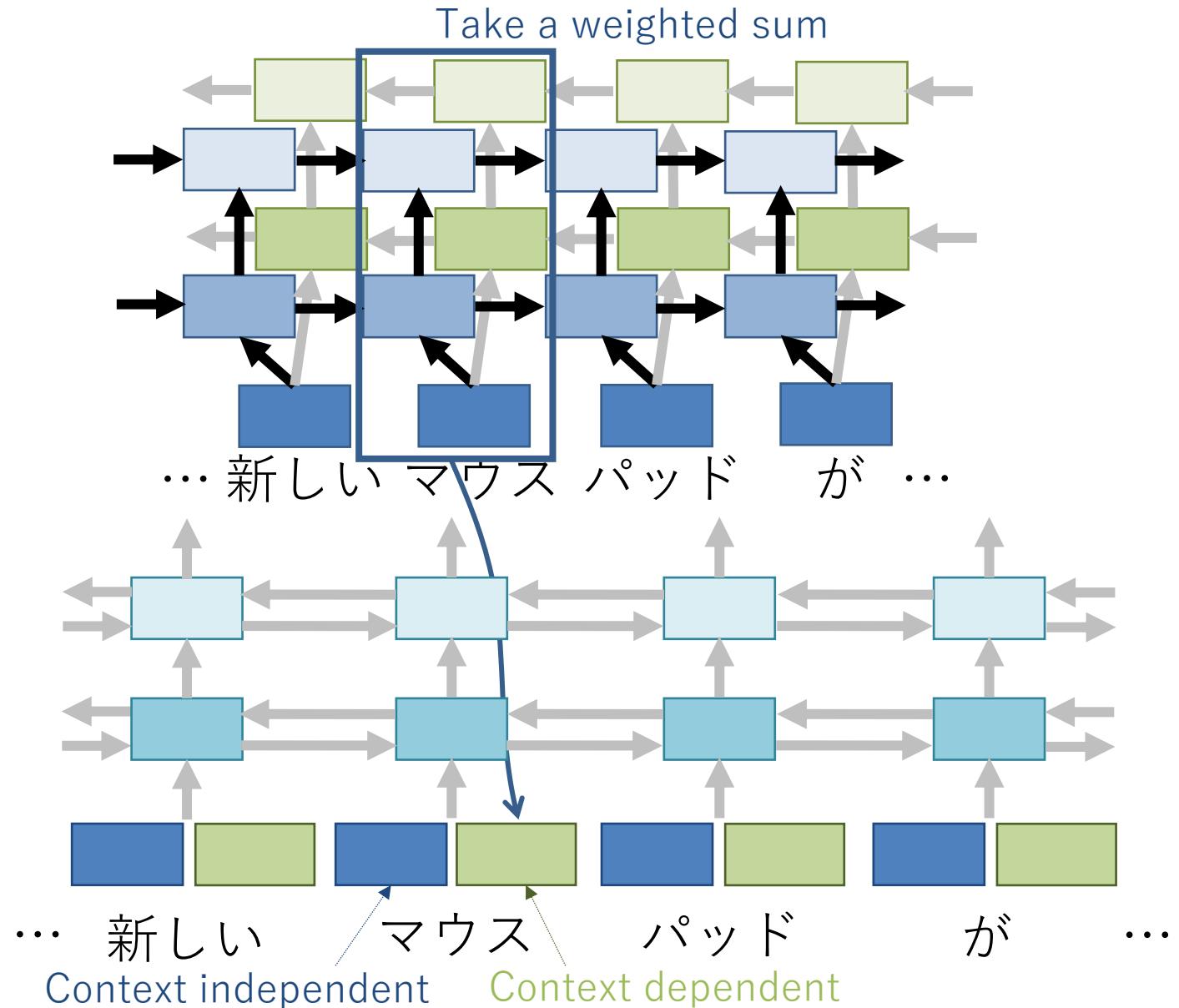
# Pre-training: Learn a bidirectional language model



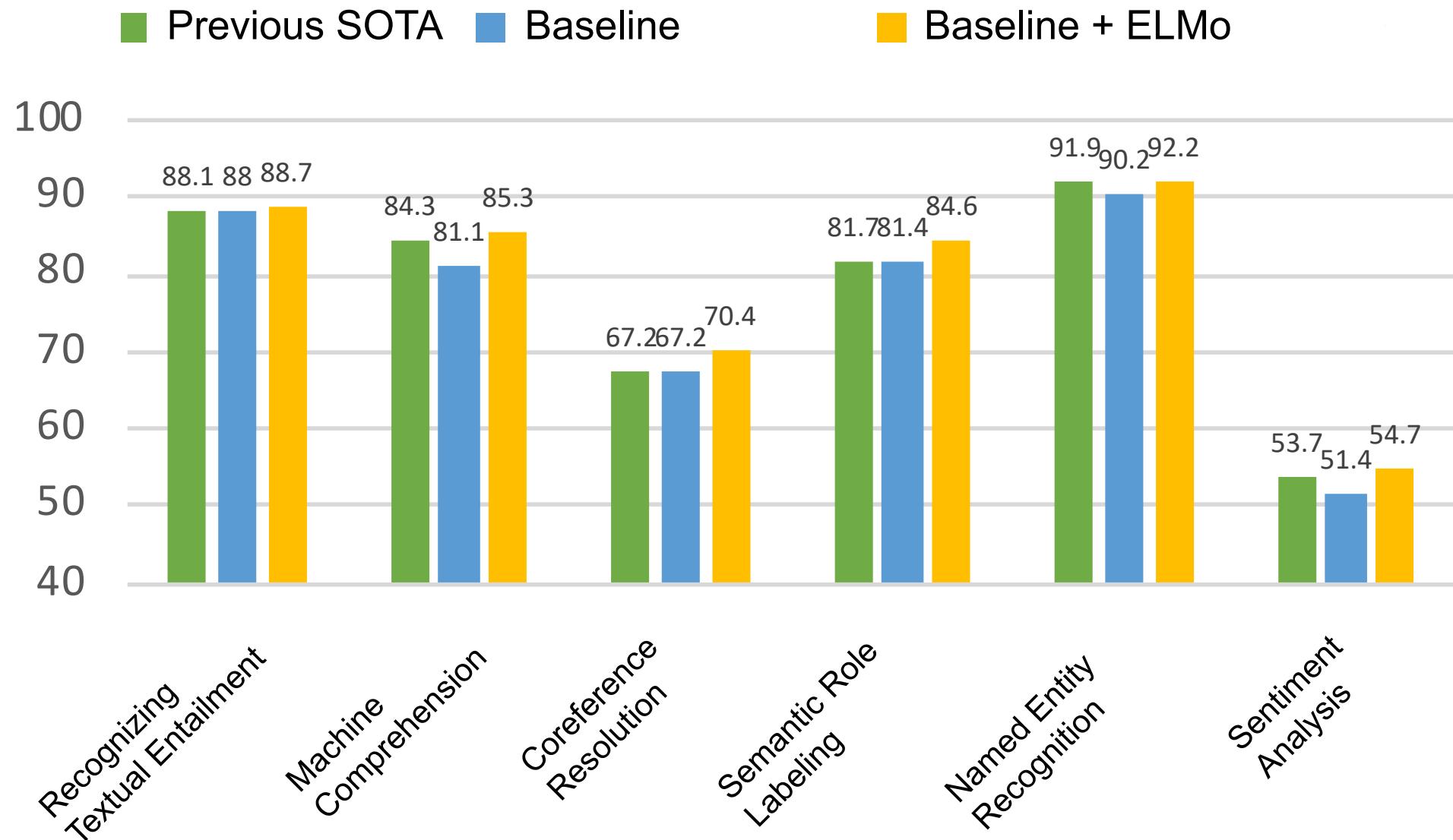
# Pre-training: Learn a bidirectional language model



# Use in a Downstream Task



# Experimental Results



# BERT [Devlin+ 2019]

- Model: Transformer [Vaswani+ 2017]
- Training: two steps (pre-training and fine-tuning)

## 1. Pre-training

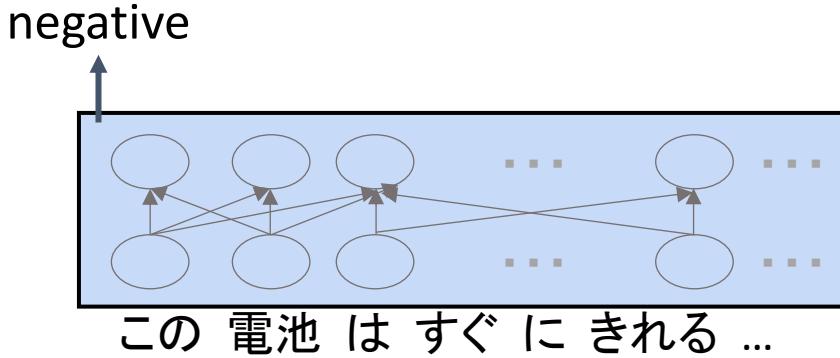
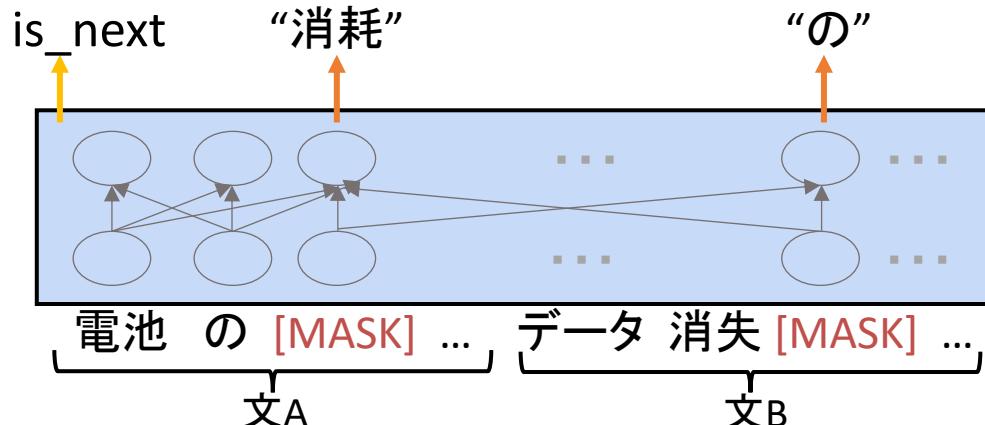
Learn vector representations using  
a large-scale raw corpus

## 2. Fine-tuning

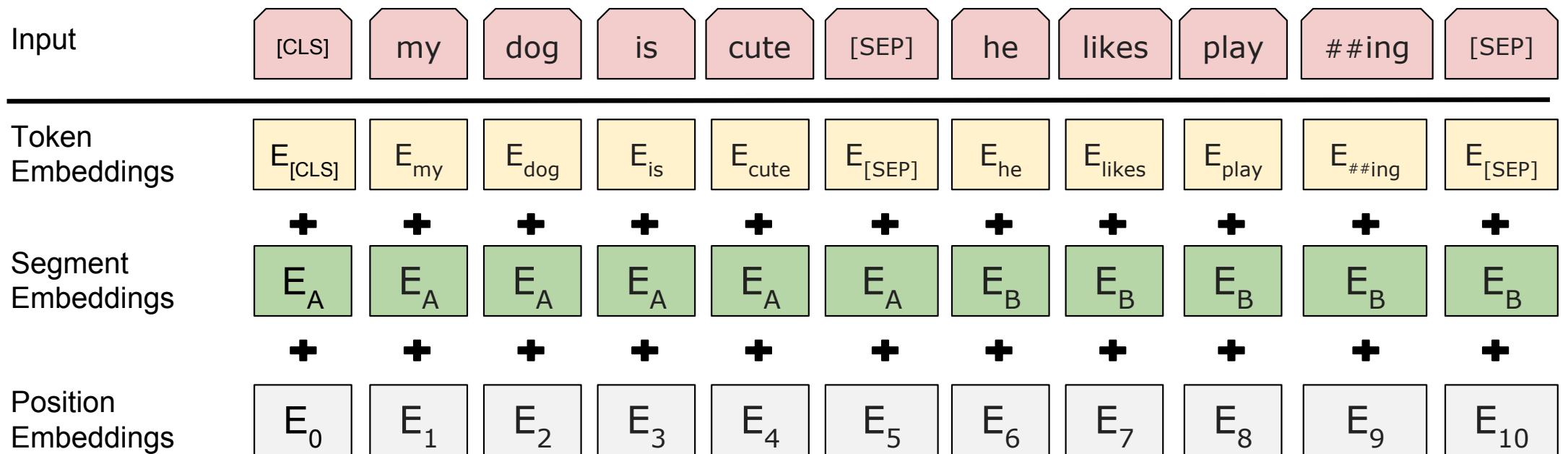
Fine-tune parameters for a specific task  
(e.g., sentiment analysis)

Next sentence prediction

Masked language model



# Input Representation



# Subwords

- A unit between a word and a character
- The purpose is to reduce unknown words
  - If subwords are not used, words other than the vocabulary (e.g., 30K words) become the [UNK] token
- Subwords have been introduced in neural machine translation
- An algorithm to induce subwords is Byte Pair Encoding (BPE)  
[Sennrich+ 2016]

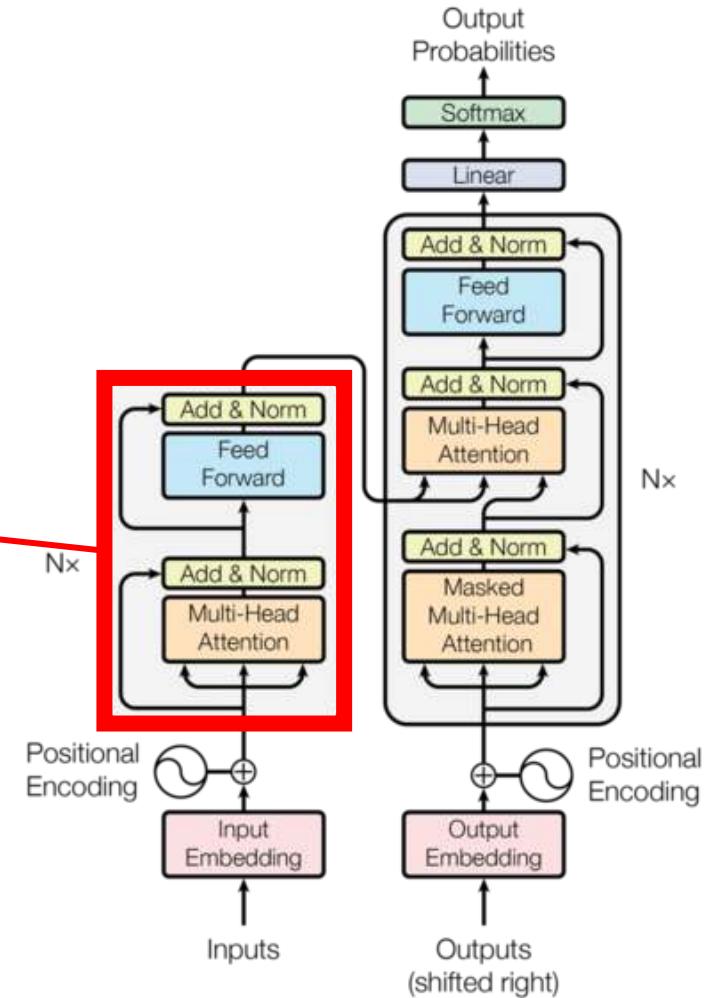
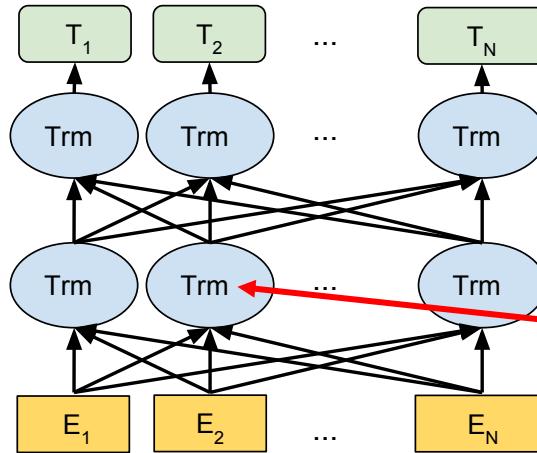
# Algorithm of BPE

1. Regard each character as a subword
2. Find a bigram (a pair of subword) with the highest frequency
3. Regard the bigram as a subword
4. Go to Step 2

Freq

4	low	low	<u>l</u> ow	lo w
2	lower	lower	<u>l</u> ower	lower
1	lowest	low <u>e</u> s t	<u>l</u> ow est	lo w est
6	newest	new <u>e</u> s t	new est	new est
3	widest	wid <u>e</u> s t	wid est	wid est

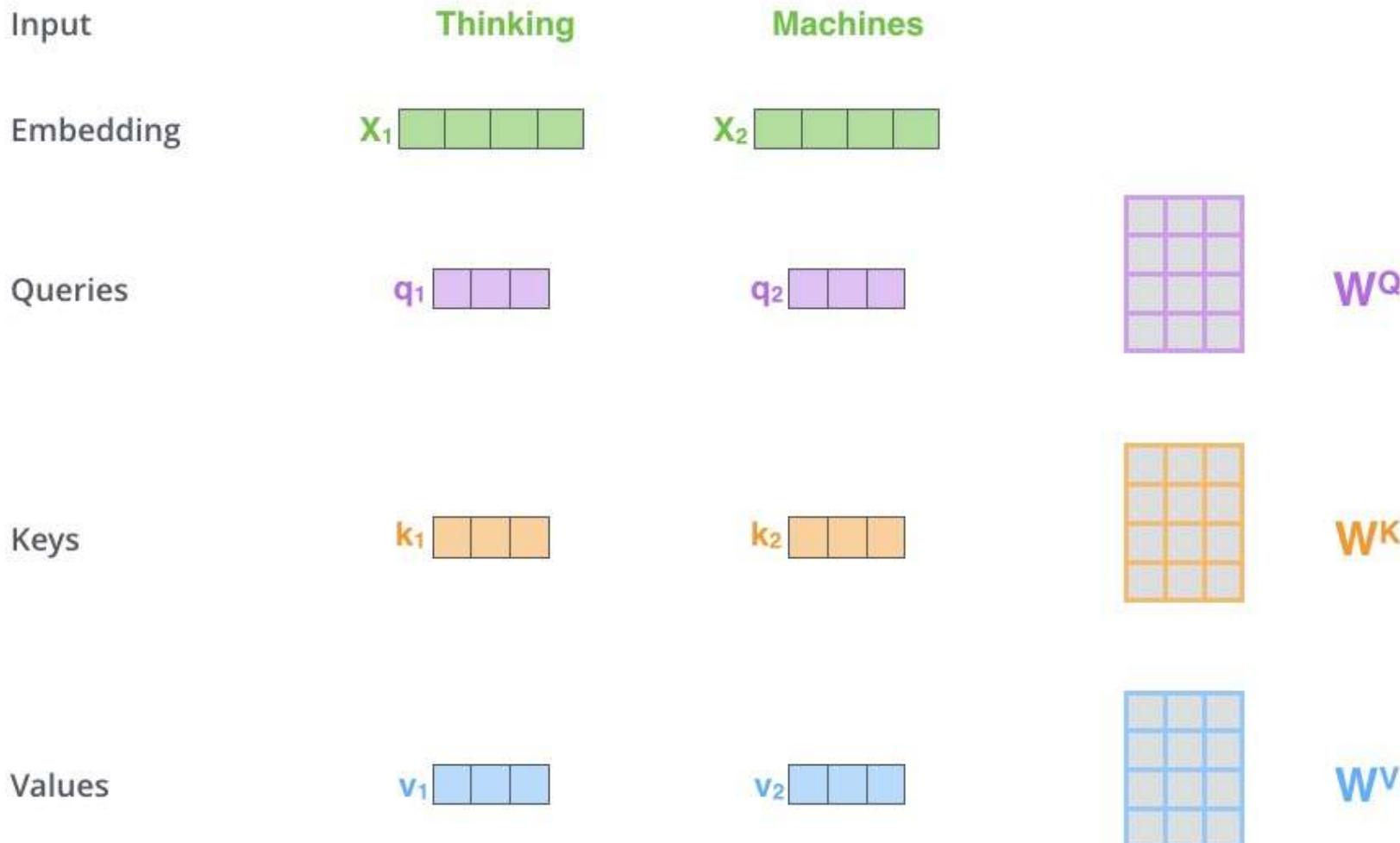
# Model Architecture



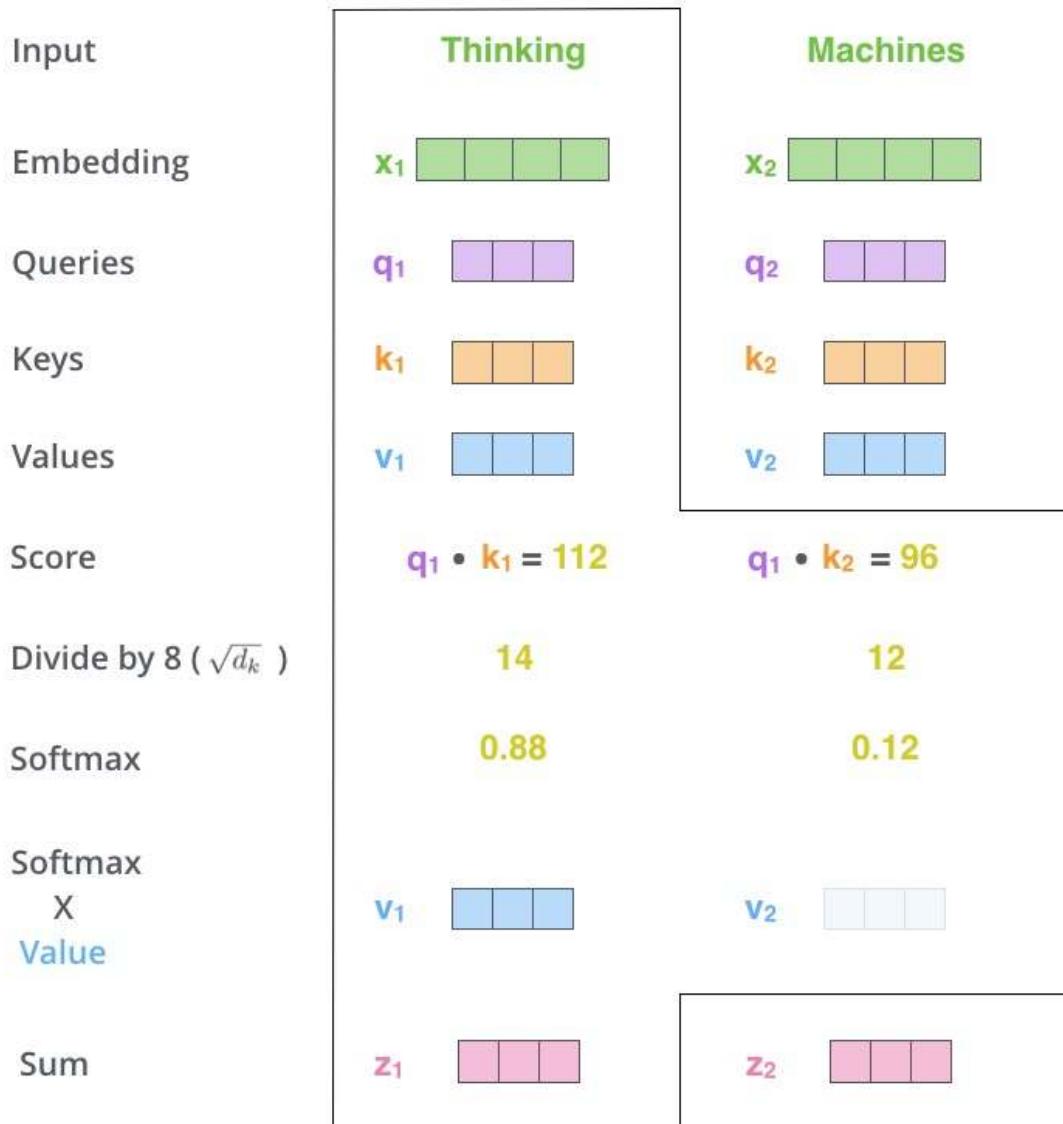
Transformer [Vaswani+ 2017]

The Annotated Transformer:  
<http://nlp.seas.harvard.edu/2018/04/03/attention.html>

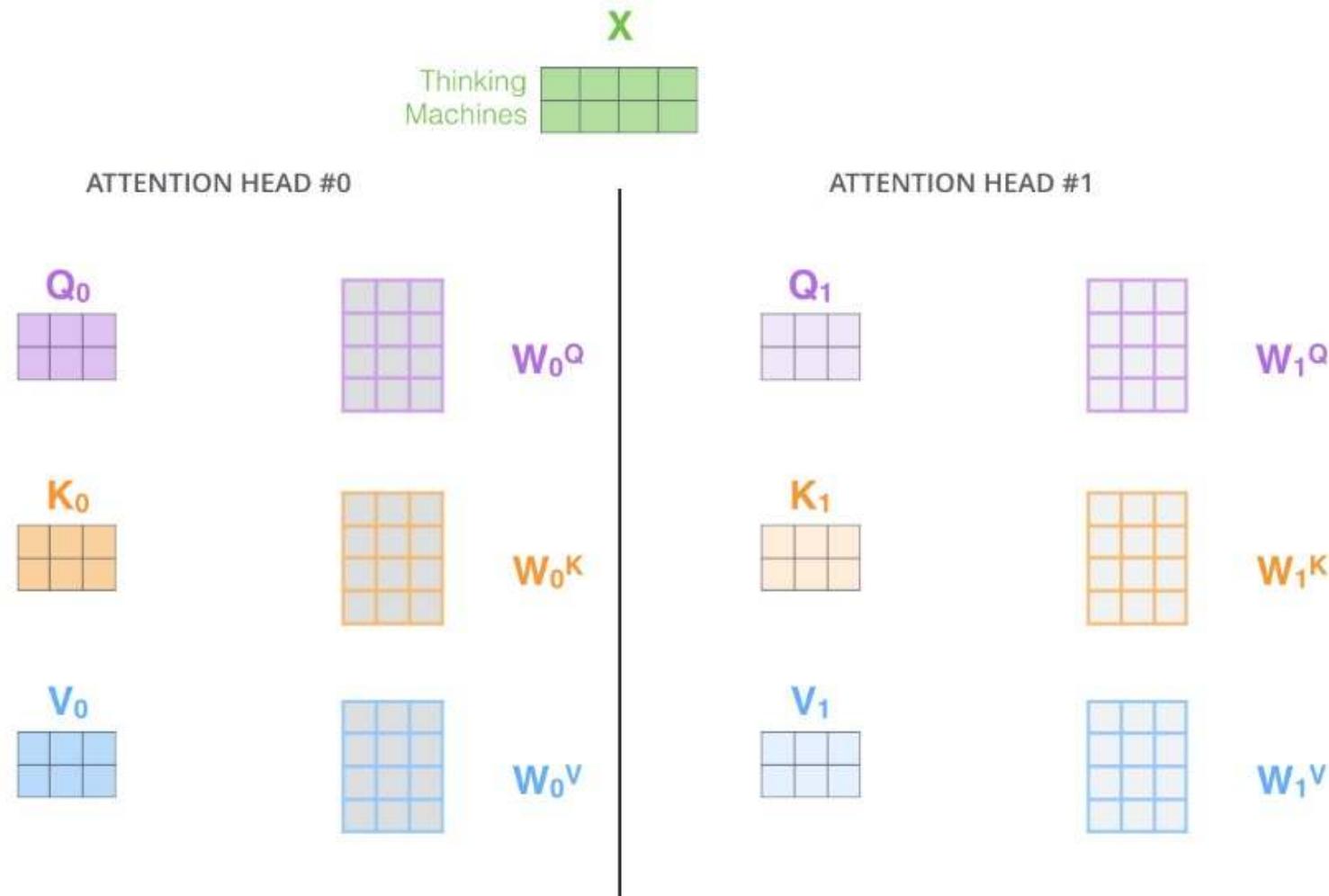
# Query, Key, Value



# Self-Attention

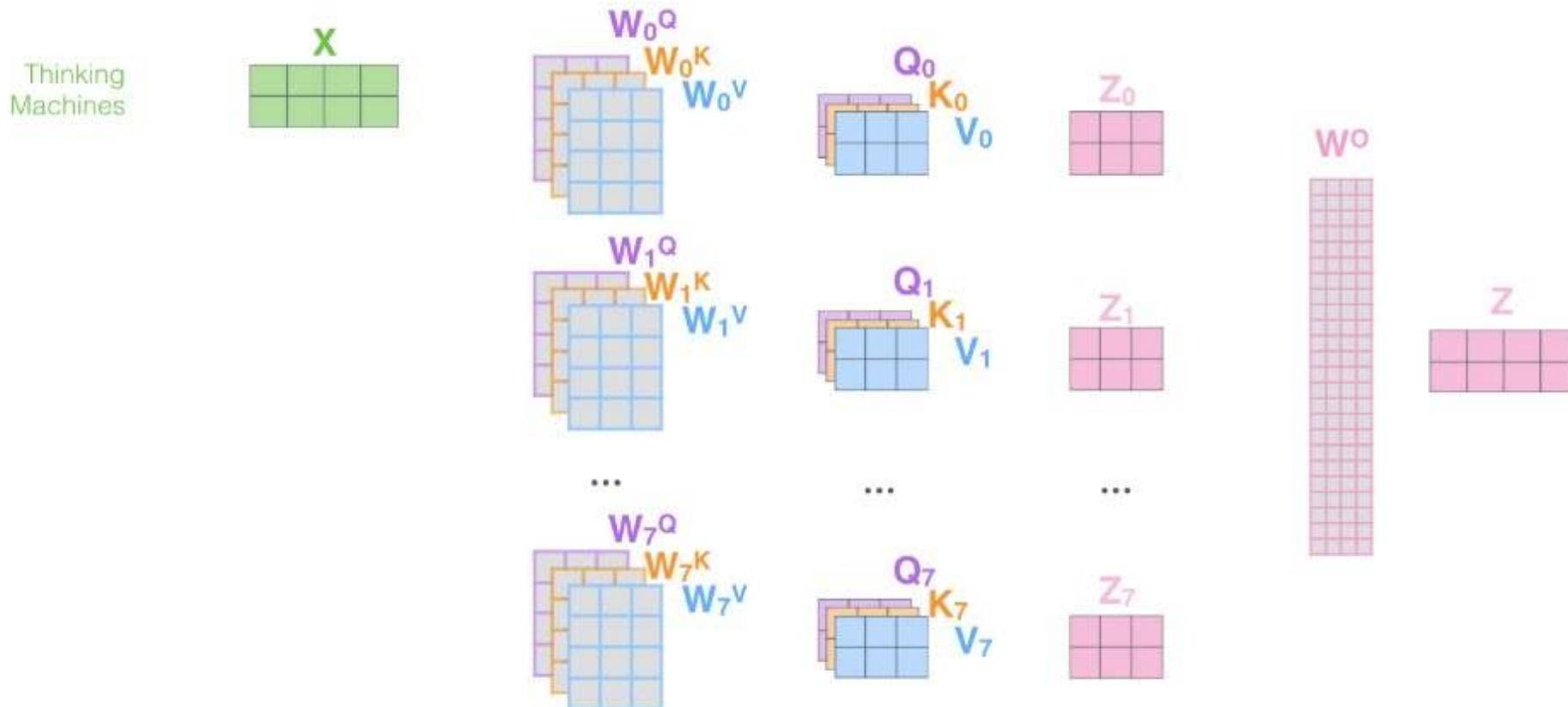


# Multiple Heads

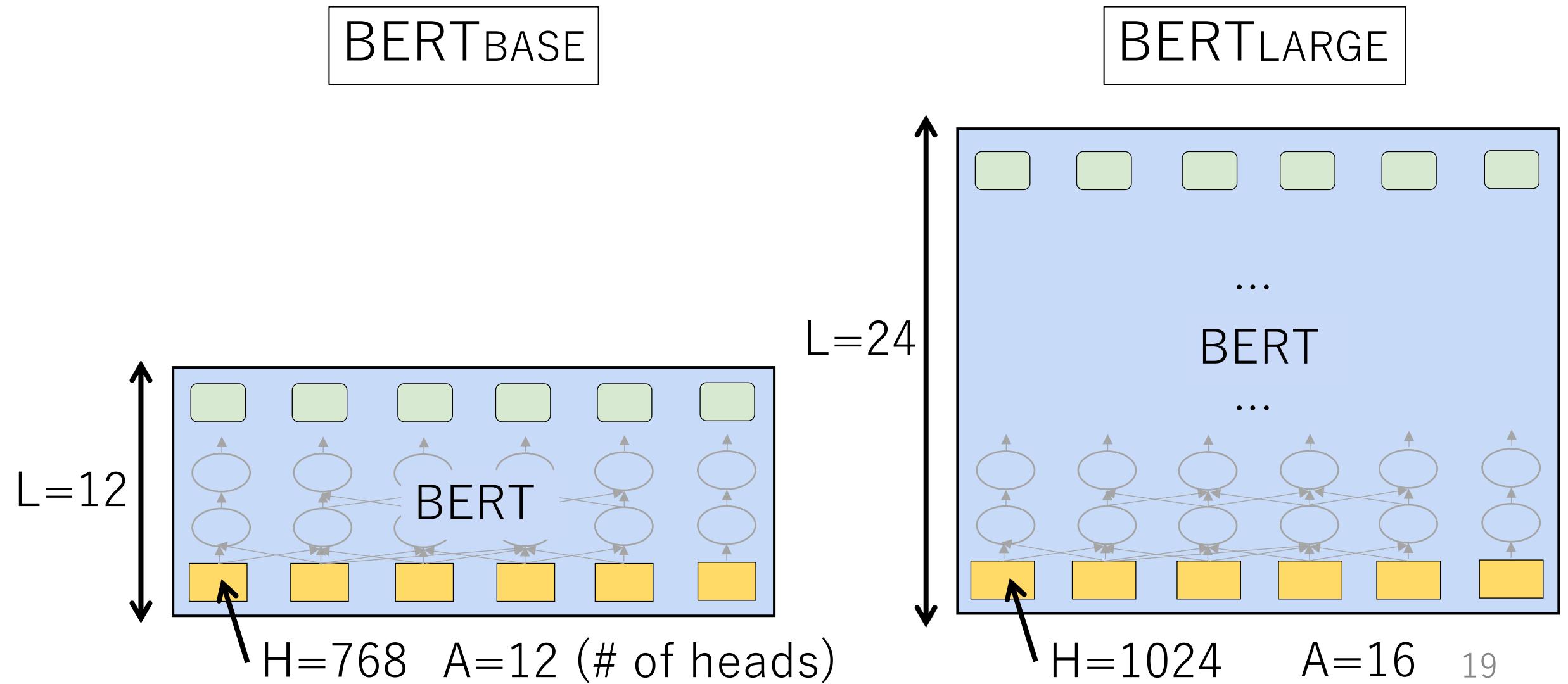


# Summary

- 1) This is our input sentence\*  
Thinking Machines
- 2) We embed each word\*  
 $X$
- 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

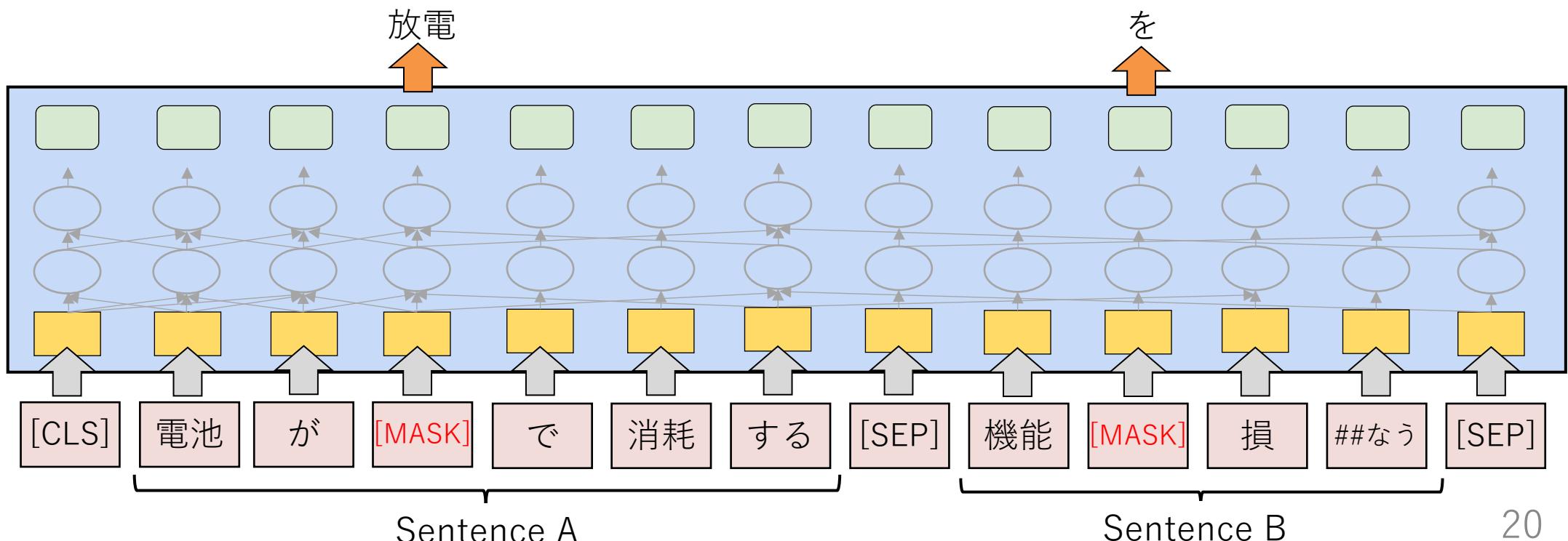


# Models



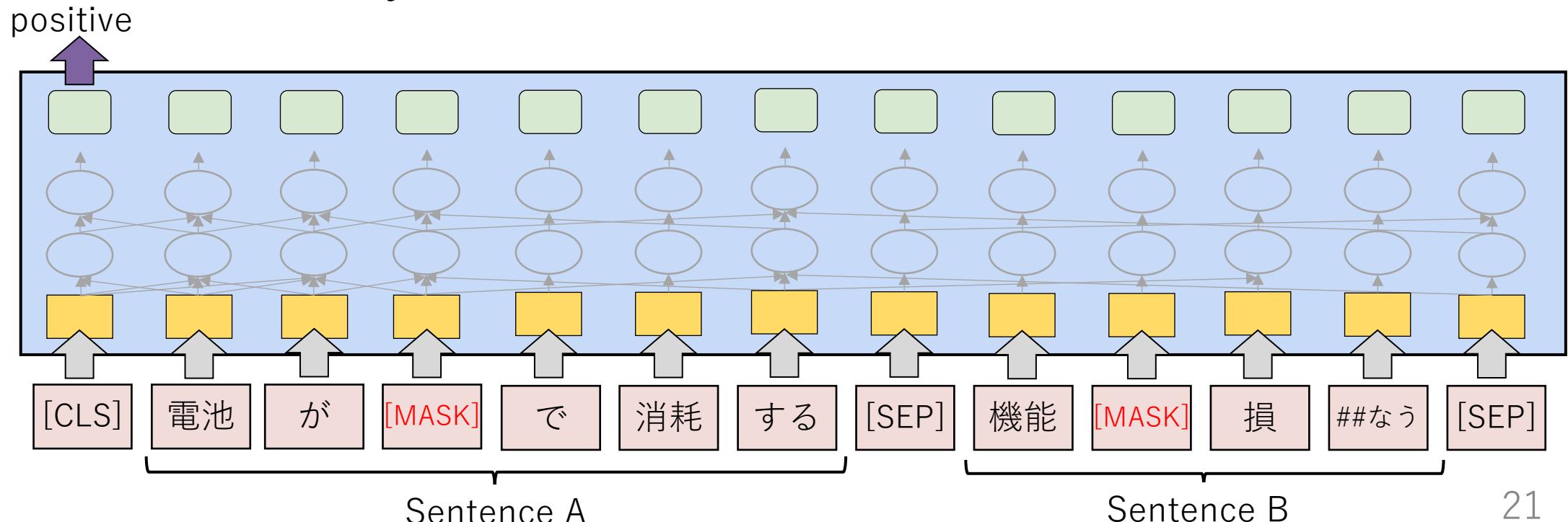
# Pre-training: 1. Masked Language Model

- Mask a randomly selected token in an input sentence
- Predict the masked token from its context
  - It is required to learn a syntactic and semantic representation



# Pre-training: 2. Next Sentence Prediction

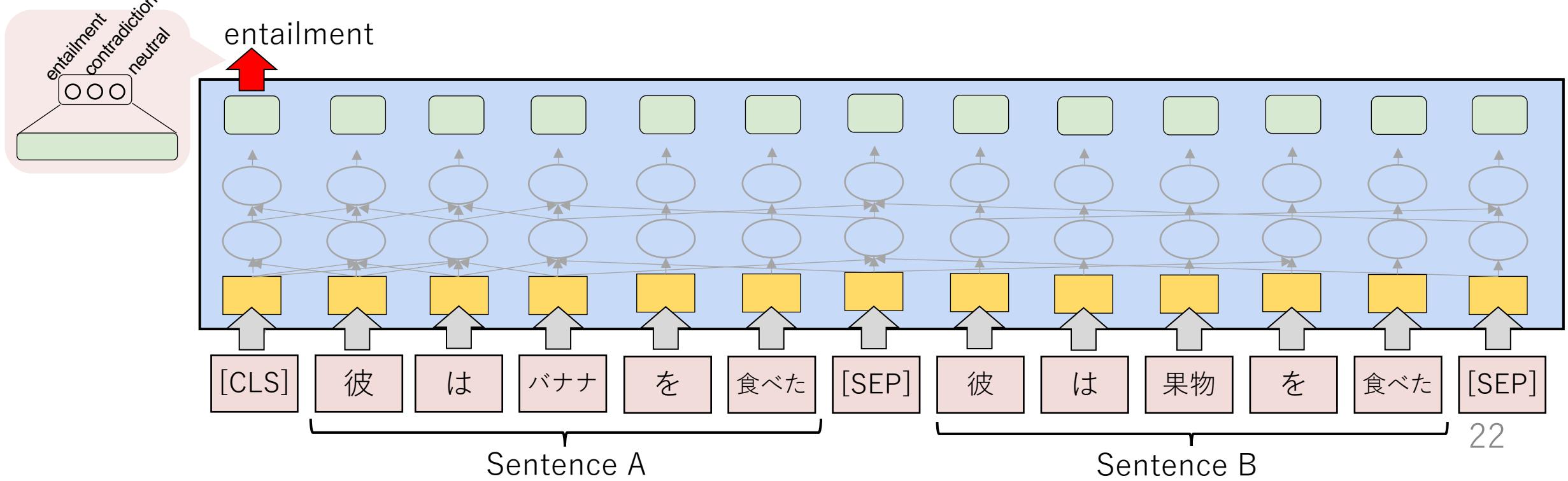
- It is necessary to capture relations between two sentences for entailment recognition and question answering
- Solve a problem to distinguish sentence B following sentence A from a randomly selected sentence



# Fine-Tuning (1/2)

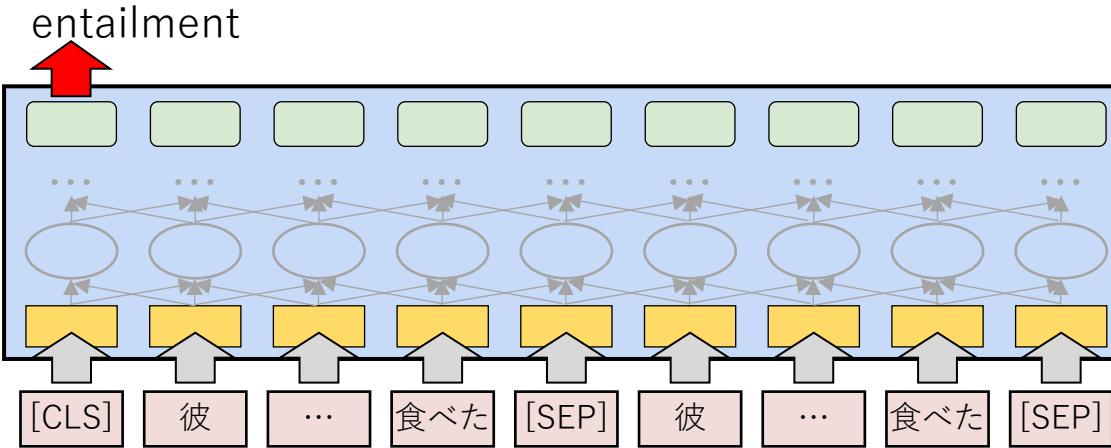
- Add a top layer according to a downstream task
- Update the parameters of the top layer and the Transformer

Sentence pair classification (e.g., recognizing textual entailment (RTE))

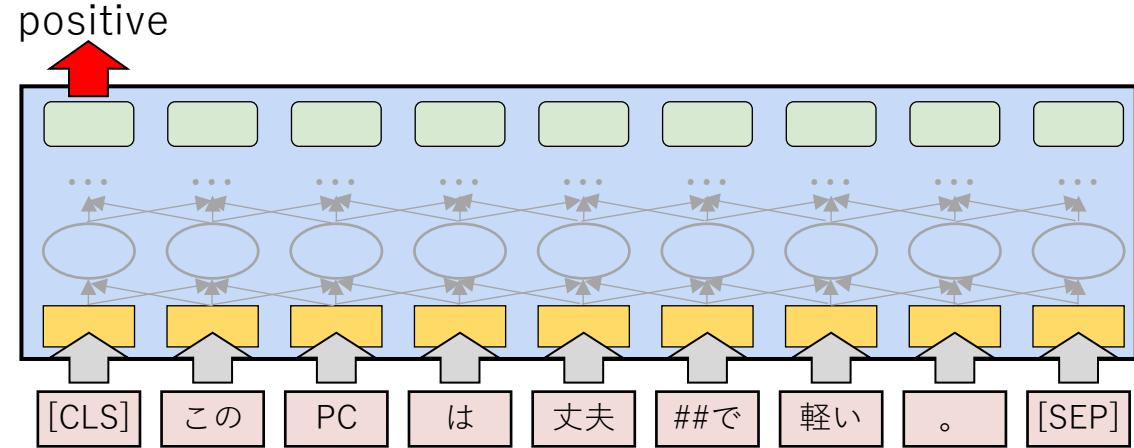


# Fine-Tuning (2/2)

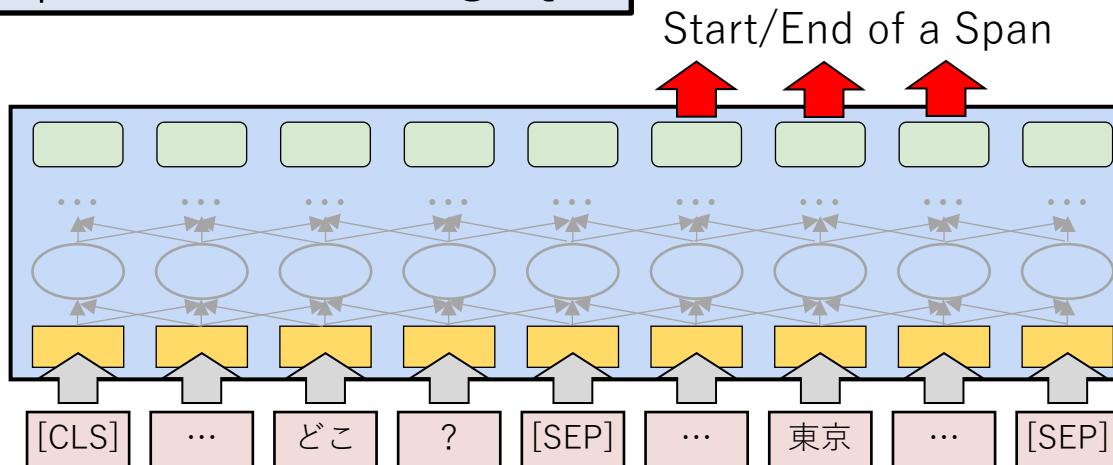
Sentence pair classification (e.g., RTE)



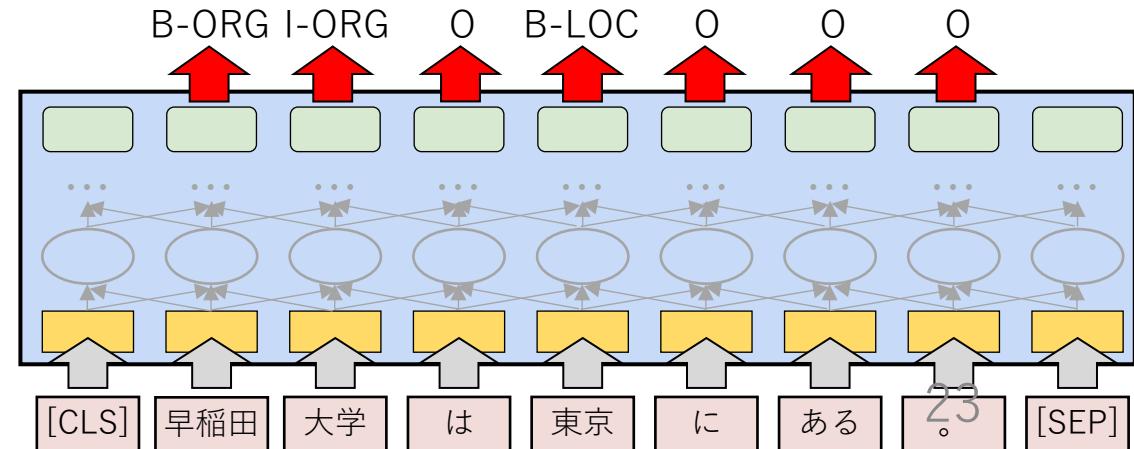
Sentence classification (e.g., sentiment analysis)



Span extraction (e.g., QA)



Sequence labeling (e.g., NER)



# GLUE (General Language Understanding Evaluation)

[Wang+ 18]

Task	Description
1s	SST-2 Polarity (positive/negative) of movie reviews
	CoLA Whether a sentence is acceptable
	MRPC Whether 2 sentences have the same meaning
	STS-B Similarity of 2 sentences (1-5)
2s	QQP Whether 2 question sentences have the same meaning
	MNLI 2 sentences have a relation of entailment/contradiction/neutral
	QNLI Whether a sentence contains an answer for a question (from SQuAD)
	RTE Whether 2 sentences have an entailment relation
	WNLI Whether 2 sentences have an entailment relation (from Winograd Schema Challenge)

# GLUE (Compositional Semantic Understanding Evaluation)

[Wang+ 18]

Task	Description
SST-2	Polarity (positive/negative) of movie reviews
CoLA	Whether a sentence is acceptable
MRPC	Whether 2 sentences have the same meaning
STS-B	Similarity of 2 sentences (1-5)

[same]

The top rate will go to 4.45 percent for all residents with taxable incomes above \$ 500,000 .  
For residents with incomes above \$ 500,000 , the income-tax rate will increase to 4.45 percent .

[not same]

While dioxin levels in the environment were up last year , they have dropped by 75 percent since the 1970s , said Caswell .  
The Institute said dioxin levels in the environment have fallen by as much as 76 percent since the 1970s .

# GLUE Leaderboard (as of November 2019)

Rank	Name	Model	URL	Scores	MMLU	SST-2	MRPC	STS-B	QQP	MNLI-m	MNLI-mm	QNLI	RTE	WNLI	AX					
1	T5 Team - Google	T5	<a href="#">🔗</a>	89.7	70.8	97.1	91.9/89.2	92.5/92.1	74.6/90.4	92.0	91.7	96.7	92.5	93.2	53.1					
2	ALBERT-Team Google Language	ALBERT (Ensemble)	<a href="#">🔗</a>	89.4	69.1	97.1	93.4/91.2	92.5/92.0	74.2/90.5	91.3	91.0	99.2	89.2	91.8	50.2					
+ 3	王玮	ALICE v2 large ensemble (Alibaba DAMO NLP)	<a href="#">🔗</a>	89.0	69.2	97.1	93.6/91.5	92.7/92.3	74.4/90.7	90.7	90.2	99.2	87.3	89.7	47.8					
4	Microsoft D365 AI & UMD	FreeLB-RoBERTa (ensemble)	<a href="#">🔗</a>	88.8	68.0	96.8	93.1/90.8	92.4/92.2	74.8/90.3	91.1	90.7	98.8	88.7	89.0	50.1					
5	Facebook AI	RoBERTa	<a href="#">🔗</a>										88.2	89.0	48.7					
6	XLNet Team	XLNet-Large (ensemble)	<a href="#">🔗</a>										86.3	90.4	47.5					
+ 7	Microsoft D365 AI & MSR AI	MT-DNN-ensemble	<a href="#">🔗</a>	87.6		96.5	92.7/90.3	91.1/90.7	73.7/89.9	87.9	87.4	96.0	86.3	89.0	42.8					
8	GLUE Human Baselines	GLUE Human Baselines	<a href="#">🔗</a>	87.1	66.4	97.8	86.3/80.8	92.7/92.6	59.5/80.4	92.0	92.8	91.2	93.6	95.9	-					
9	Stanford Hazy Research	Snorkel MeTaL	<a href="#">🔗</a>	83.2	63.8	96.2	91.5/88.5	90.1/89.7	73.1/89.9	87.6	87.2	93.9	80.9	65.1	39.9					
10	XLM Systems	XLM (English only)	<a href="#">🔗</a>	83.1	62.9	95.6	90.7/87.1	88.8/88.2	73.2/89.8	89.1	88.5	94.0	76.0	71.9	44.7					
11	Zhuosheng Zhang	SemBERT	<a href="#">🔗</a>	82.9	62.3	94.6	91.2/88.3	87.8/86.7	72.8/89.8	87.6	86.3	94.6	84.5	65.1	42.4					
12	Danqi Chen	SpanBERT (single-task training)	<a href="#">🔗</a>	82.8	64.3	94.8	90.9/87.9	89.9/89.1	71.9/89.5	88.1	87.7	94.3	79.0	65.1	45.1					
13	Kevin Clark	BERT + BAM	<a href="#">🔗</a>	82.3	61.5	95.2	91.3/88.3	88.6/87.9	72.5/89.7	86.6	85.8	93.1	80.4	65.1	40.7					
14	Nitish Shirish Keskar	Span-Extractive BERT on STILTs	<a href="#">🔗</a>						89.4/89.2	72.2/89.4	86.5	85.8	92.5	79.8	65.1	28.3				
15	Jason Phang	BERT on STILTs	<a href="#">🔗</a>						88.7/88.3	71.9/89.4	86.4	85.6	92.7	80.1	65.1	28.3				
16	廖亿	RGLM-Base (Huawei Noah's Ark Lab)	<a href="#">🔗</a>		56.9	94.2	90.7/87.7	89.7/89.1	72.2/89.4	86.1	85.4	92.1	78.5	65.1	40.0					
+ 17	Jacob Devlin	BERT: 24-layers, 16-heads, 1024-hidden	<a href="#">🔗</a>	80.5	60.5	94.9	89.3/85.4	87.6/86.5	72.1/89.3	86.7	85.9	92.7	70.1	65.1	39.6					
18	Neil Houlsby	BERT + Single-task Adapters	<a href="#">🔗</a>	80.2	59.2	94.3	88.7/84.3	87.3/86.1	71.5/89.4	85.4	85.0	92.4	71.6	65.1	9.2					
19	Zhuohan Li	Macaron Net-base	<a href="#">🔗</a>	79.7	57.6	94.0	88.4/84.4	87.5/86.3	70.8/89.0	85.4	84.5	91.6	70.5	65.1	38.7					
20	蘇大鈞	SesameBERT-Base	<a href="#">🔗</a>	78.6	52.7	94.2	88.9/84.8	86.5/85.5	70.8/88.8	83.7	83.6	91.0	67.6	65.1	35.8					
+ 21	MobileBERT Team	MobileBERT	<a href="#">🔗</a>	78.5	51.1	92.6	88.8/84.5	86.2/84.8	70.5/88.3	84.3	83.4	91.6	70.4	65.1	34.3					
22	Linyuan Gong	StackingBERT-Base	<a href="#">🔗</a>								84.4	84.2	90.1	67.0	65.1	36.6				
23	Huawei Noah's Ark Lab	TinyBERT (4-layers; 7.5x smaller and 9.4x faster than BERT-base)	<a href="#">🔗</a>								82.5	81.8	87.7	62.9	65.1	33.7				
24	shijing si	bert+pos6	<a href="#">🔗</a>						78.8	52.9	93.9	88.8/84.6	83.8/85.5	71.4/89.2	84.4	83.3	90.4	66.9	34.9	0.0
25	GLUE Baselines	BiLSTM+ELMo+Attn	<a href="#">🔗</a>		70.0	33.6	90.4	84.4/78.0	74.2/72.3	63.1/84.3	74.1	74.5	79.8	58.9	65.1	21.7				

Baseline (ELMo): 70.0

T5: 89.7

Human performance: 87.1

BERT: 80.5

# Pre-training for Japanese

Japanese Wikipedia  
(18M sentences)

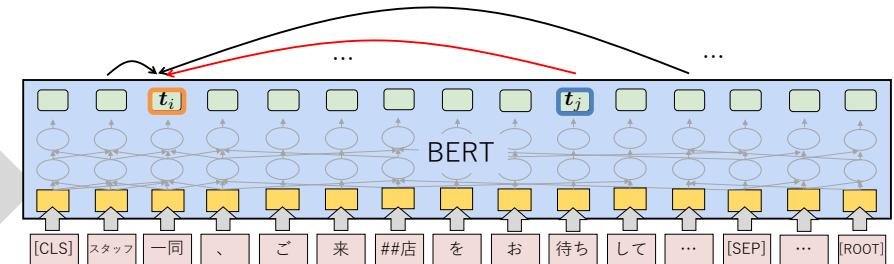
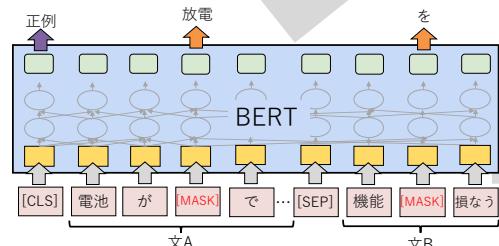
自動車 競技 は 四 横 の 自動車 あ  
るいは それ に 準ずる 車両 による  
競技 に 対して 主に 呼称 され、  
オートバイ や それ に 準ずる 車  
両 の 競技 に 対して は オートバ  
イ 競技 や モーター サイクル  
レース などと 呼ばれる。自動  
車 競技 は 操る 人 の …

PCM 音源（ピー ## シー ## エ  
ム おん ## げん）は、コンパクト  
ディスクなどで扱われるパ  
ルス 符号 変調 技術を用いたデ  
ジタル ## シンセサイザーの音源  
方式のひとつ。あらかじめメ  
モリに記録しておいたPCM  
波形（サンプル）を再生する  
こと で …

Data for  
pre-training

[CLS] 自動車 競技 は 四 横  
の 自動車 あるいは それ に  
準ずる 車両 による 競技 に  
対して 主に 呼称 され、  
オートバイ [MASK] … 呼ば  
れる。 [SEP] 自動車 競技  
は 操る [MASK] の … [SEP]  
Label: IsNext

Pre-training

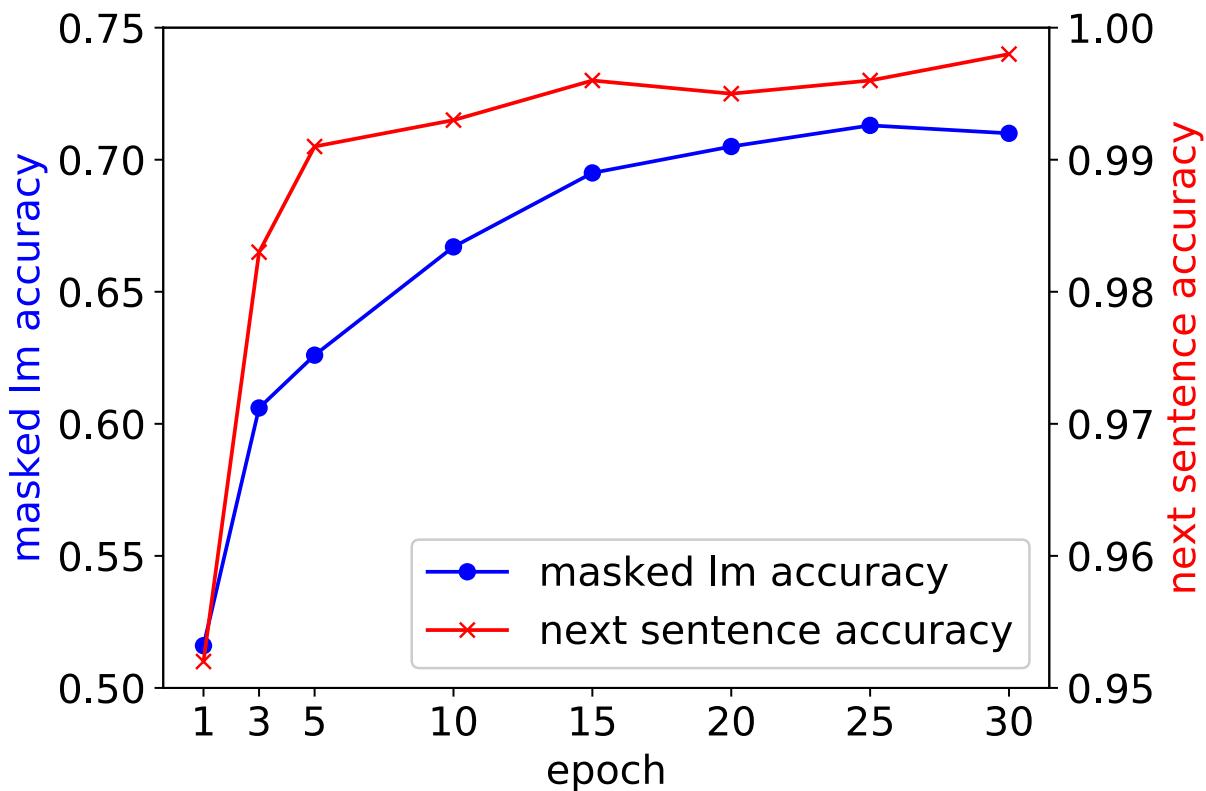


Using morphological analysis and BPE

Fine-tuning

# Pre-training for Japanese

- Vocab size: 32000, BERT<sub>BASE</sub>, max token number: 128



1 day for 1 epoch  
(using 1 GPU)

# BERT Masked LM

京都市左京区の京都大学では 17 日、試験について案内する看板の設置作業が行われた。

input	predictions
京都	京都, 同, 大阪, 福岡, 東京
市	市, 府, ・, 都, 大学
左京	左京, 東山, 中京, 右京, 伏見
区	区, 部, 地区, ##区, 地域
の	の, と, および, 、, ・
京都	京都, 立命館, 同志社, 近畿, 龍谷
大学	大学, 駅, 支所, キャンパス, 支社
で	で, に, へ, にて, と
は	は, 、, 翌, 同, 毎月
17	この, その, ある, 1, 24
日	日, 日間, 月, 年間, 年
、	に, から, 、, より, まで
試験	大学, キャンパス, 学生, これ, 授業
	に に, ##に, 結果, 詳細に, 的に
	について について, 向けて, おいて, 際して, 合わせて
	案内 説明, 告知, 案内, 解説, 揭示
	する する, 用, 表示, の, 説明
	看板 板, 看板, 装置, 揭示板, ため
	の の, を, と, や, が
	設置 設置, 撤去, 取り付け, 作成, 貼り
	作業 作業, 式, 工事, など, 実験
	が が, も, を, まで, で
	行わ 行わ, 行なわ, おこなわ, さ, 開か
	れた れた, れる, れ, られた, れて
	。 。, 、, に, も, ため

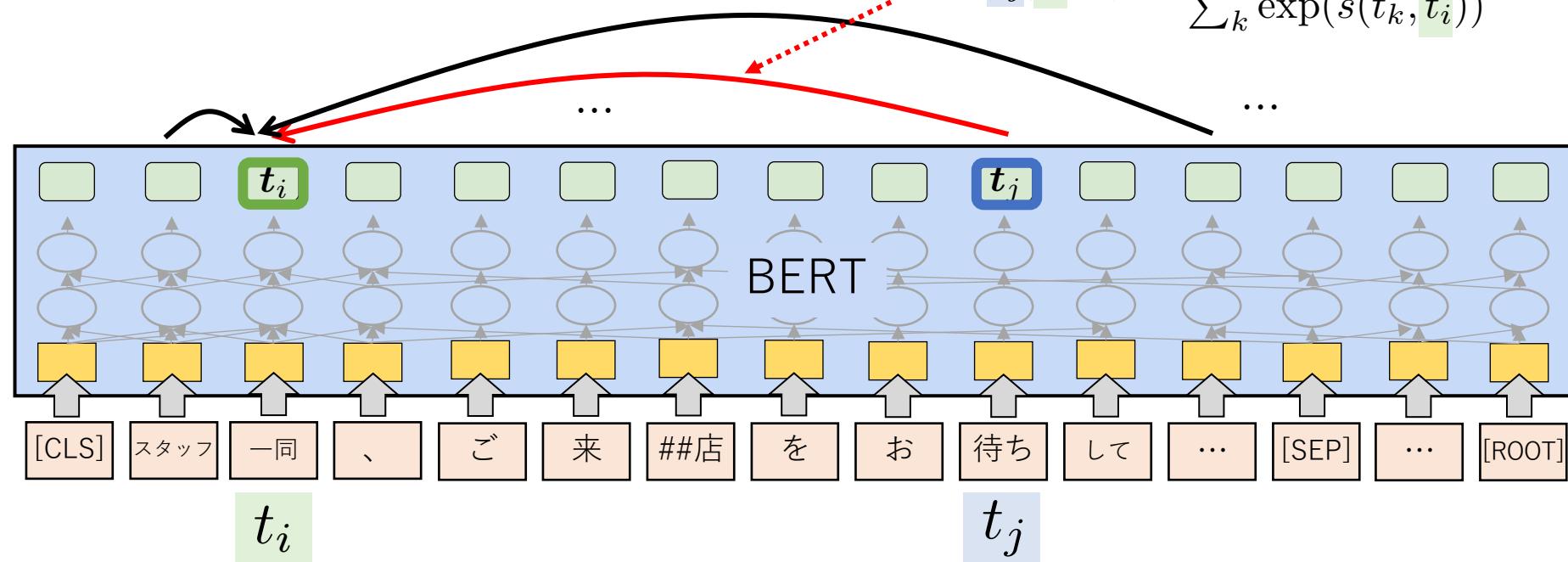
# BERT Fine-tuning for Dependency Parsing

[柴田+ 19]

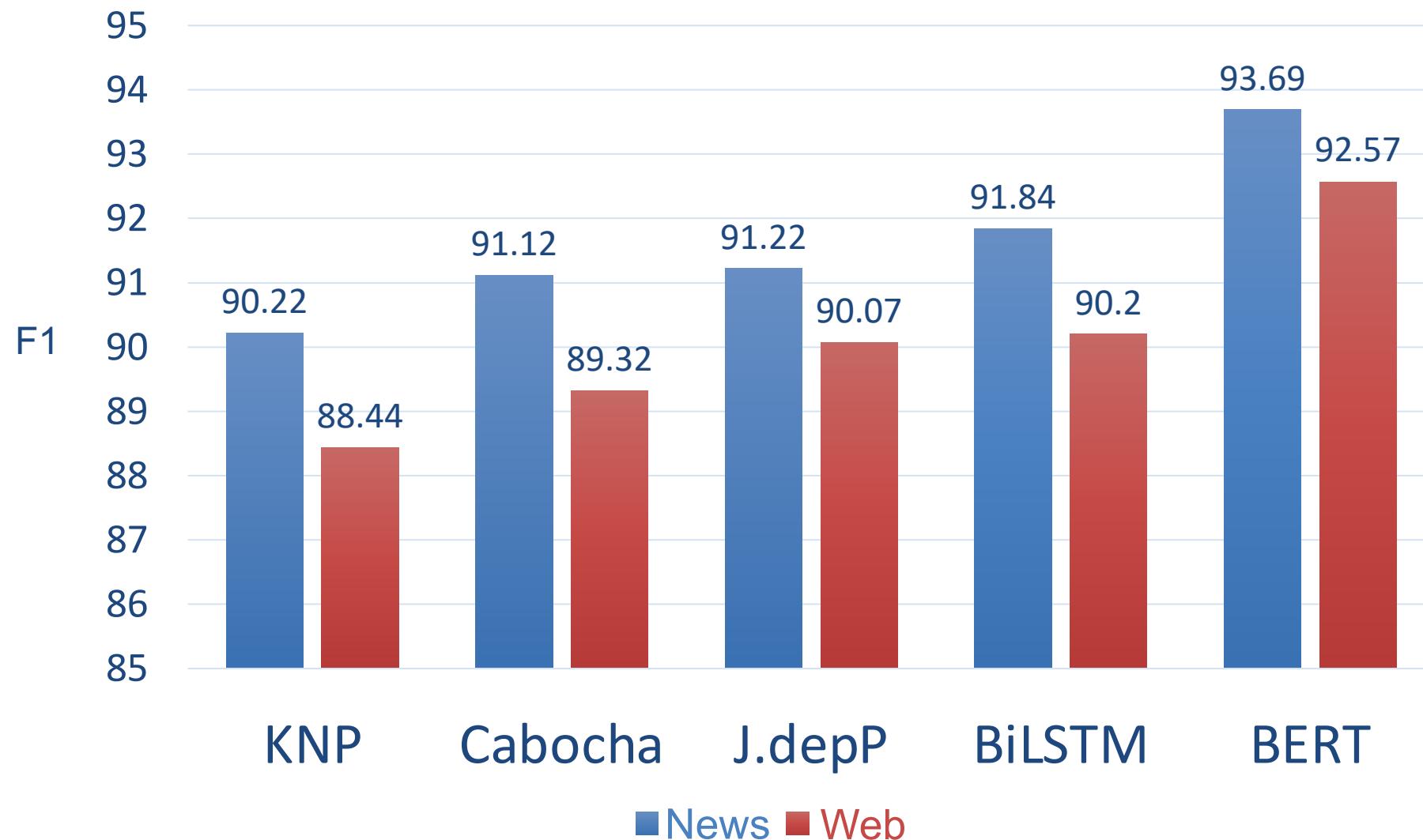
Based on  
head selection [Zhang+ 17]

$$s(t_j, t_i) = \mathbf{v}_h^T \tanh(U_h \mathbf{t}_j + W_h \mathbf{t}_i)$$

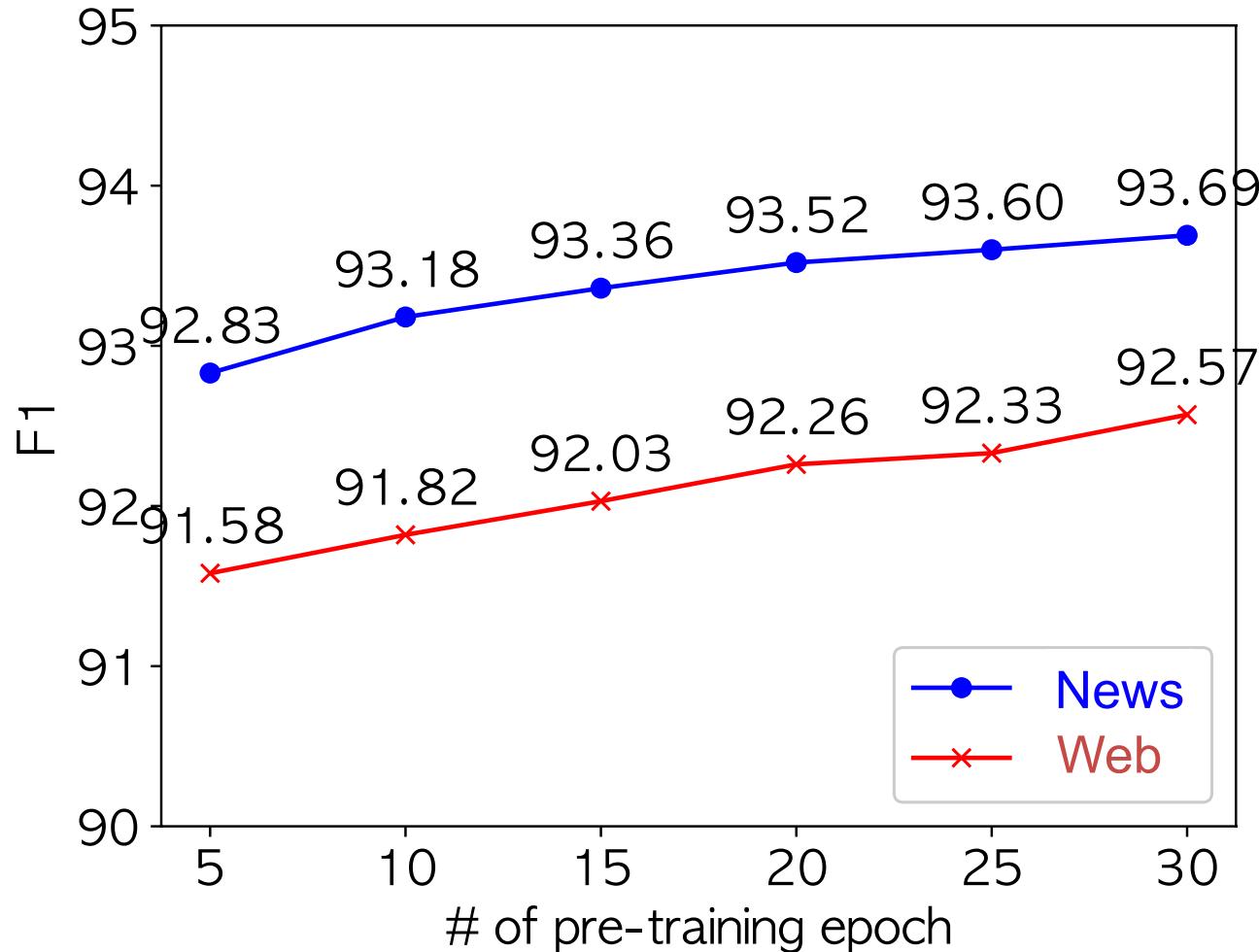
$$P_{\text{head}}(t_j | t_i, S) = \frac{\exp(s(t_j, t_i))}{\sum_k \exp(s(t_k, t_i))}$$



# Phrase-based Performance (F1)



# Learning Curve while Increasing Epochs of Pre-training



# Recent Topics around BERT

- Model improvements
- Use of knowledge
- Language models
- Generation models for translation, summarization, etc.
- Multilingualization
- Analysis of inside BERT
- Use of multimodal information

# RoBERTa [Liu+ 2019] (FAIR&UW)

A Robustly Optimized BERT Pretraining Approach

- Dynamic masking
- One-segment pre-training without NSP
- Larger batch size (32x BERT = 8K)
- Larger corpora (10x BERT = 160GB)
- Longer pre-training steps (16x BERT)
- Byte-level BPE
  - larger vocabulary (50K subwords)

# RoBERTa [Liu+ 2019]

Model	data	bsz	steps	SQuAD (v1.1/2.0)	MNLI-m	SST-2
RoBERTa						
with BOOKS + WIKI	16GB	8K	100K	93.6/87.3	89.0	95.3
+ additional data (§3.2)	160GB	8K	100K	94.0/87.7	89.3	95.6
+ pretrain longer	160GB	8K	300K	94.4/88.7	90.0	96.1
+ pretrain even longer	160GB	8K	500K	<b>94.6/89.4</b>	<b>90.2</b>	<b>96.4</b>
BERT <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	90.9/81.8	86.6	93.7
XLNet <sub>LARGE</sub>						
with BOOKS + WIKI	13GB	256	1M	94.0/87.8	88.4	94.4
+ additional data	126GB	2K	500K	94.5/88.8	89.8	95.6

# RoBERTa [Liu+ 2019]

GLUE

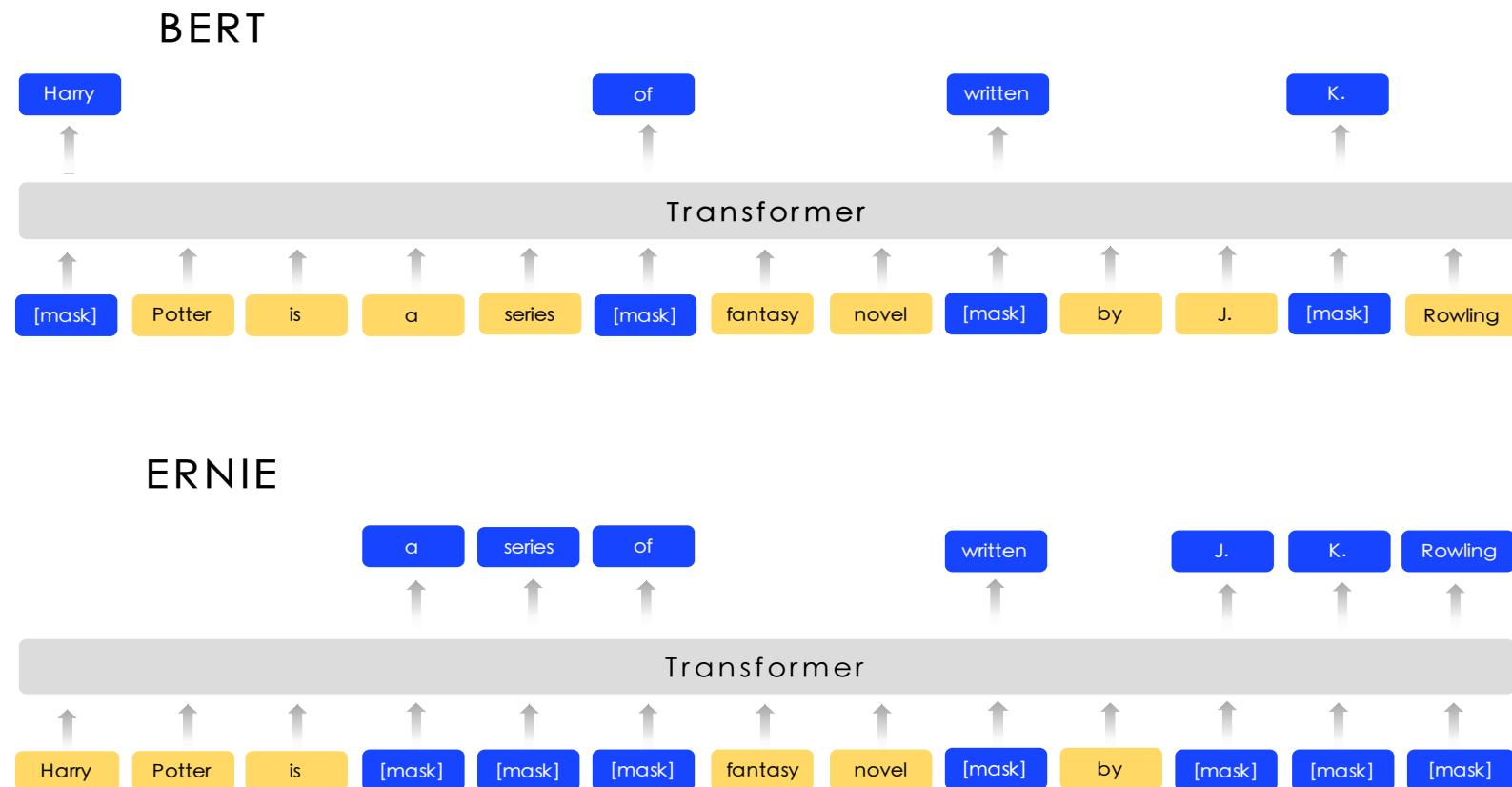
	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
<i>Single-task single models on dev</i>										
BERT <sub>LARGE</sub>	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet <sub>LARGE</sub>	89.8/-	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa	<b>90.2/90.2</b>	<b>94.7</b>	<b>92.2</b>	<b>86.6</b>	<b>96.4</b>	<b>90.9</b>	<b>68.0</b>	<b>92.4</b>	<b>91.3</b>	-
<i>Ensembles on test (from leaderboard as of July 25, 2019)</i>										
ALICE	88.2/87.9	95.7	<b>90.7</b>	83.5	95.2	92.6	<b>68.6</b>	91.1	80.8	86.3
MT-DNN	87.9/87.4	96.0	89.9	86.3	96.5	92.7	68.4	91.1	89.0	87.6
XLNet	90.2/89.8	98.6	90.3	86.3	<b>96.8</b>	<b>93.0</b>	67.8	91.6	<b>90.4</b>	88.4
RoBERTa	<b>90.8/90.2</b>	<b>98.9</b>	90.2	<b>88.2</b>	96.7	92.3	67.8	<b>92.2</b>	89.0	<b>88.5</b>

# Whole Word Masking (WWM)

- Input Text:  
the man jumped up , put his basket on phil ##am ##mon ' s head
- Original Masked Input:  
[MASK] man [MASK] up , put his [MASK] on phil [MASK] ##mon ' s head
- Whole Word Masked Input:  
the man [MASK] up , put his basket on [MASK] [MASK] [MASK] ' s head

# ERNIE [Sun+ 2019] (Baidu)

- Phrase-level masking (using chunking)
- Entity-level masking (using NER)



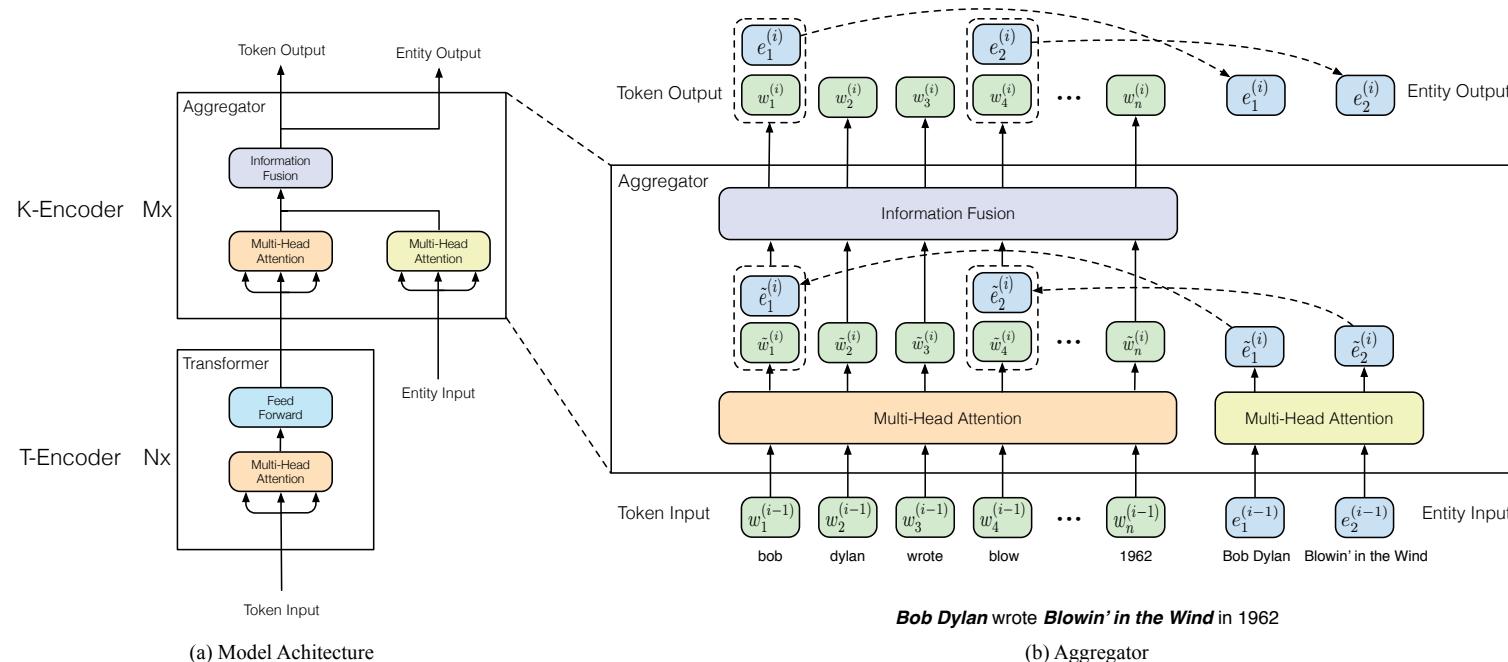
# ERNIE [Sun+ 2019]

Table 1: Results on 5 major Chinese NLP tasks

Task	Metrics	Bert		ERNIE	
		dev	test	dev	test
XNLI	accuracy	78.1	77.2	79.9 (+1.8)	78.4 (+1.2)
LCQMC	accuracy	88.8	87.0	89.7 (+0.9)	87.4 (+0.4)
MSRA-NER	F1	94.0	92.6	95.0 (+1.0)	93.8 (+1.2)
ChnSentiCorp	accuracy	94.6	94.3	95.2 (+0.6)	95.4 (+1.1)
nlpcc-dbqa	mrr	94.7	94.6	95.0 (+0.3)	95.1 (+0.5)
	F1	80.7	80.8	82.3 (+1.6)	82.7 (+1.9)

# ERNIE [Zhang+ 2019] (Tsinghua)

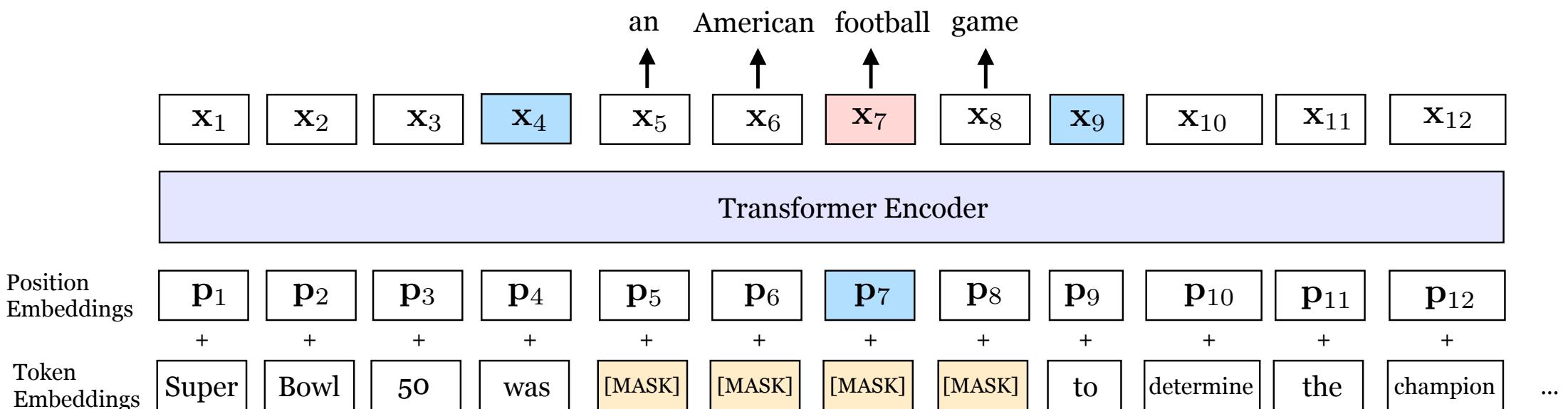
- Pre-training to predict entities
- Knowledgeable encoder
  - Entity embeddings (induced by TransE) are integrated



# SpanBERT [Joshi+ 2019] (FAIR&UW)

- Span-level masking
- Span Boundary Objective (SBO)
- One-segment pre-training without NSP

$$\mathcal{L}(\text{football}) = \mathcal{L}_{\text{MLM}}(\mathbf{x}_7) + \mathcal{L}_{\text{SBO}}(\mathbf{x}_4, \mathbf{x}_9, \mathbf{p}_7)$$



# SpanBERT [Joshi+ 2019]

	SQuAD 1.1		SQuAD 2.0	
	EM	F1	EM	F1
Human Perf.	82.3	91.2	86.8	89.4
Google BERT	84.3	91.3	80.0	83.3
Our BERT	86.5	92.6	82.8	85.9
Our BERT-1seq	87.5	93.3	83.8	86.6
SpanBERT	<b>88.8</b>	<b>94.6</b>	<b>85.7</b>	<b>88.7</b>

	NewsQA	TriviaQA	SearchQA	HotpotQA	NaturalQA	(Avg)
Google BERT	68.8	77.5	81.7	78.3	79.9	77.3
Our BERT	71.0	79.0	81.8	80.5	80.5	78.6
Our BERT-1seq	71.9	80.4	84.0	80.3	81.8	79.7
SpanBERT	<b>73.6</b>	<b>83.6</b>	<b>84.8</b>	<b>83.0</b>	<b>82.5</b>	<b>81.5</b>

# SpanBERT [Joshi+ 2019]

Coreference

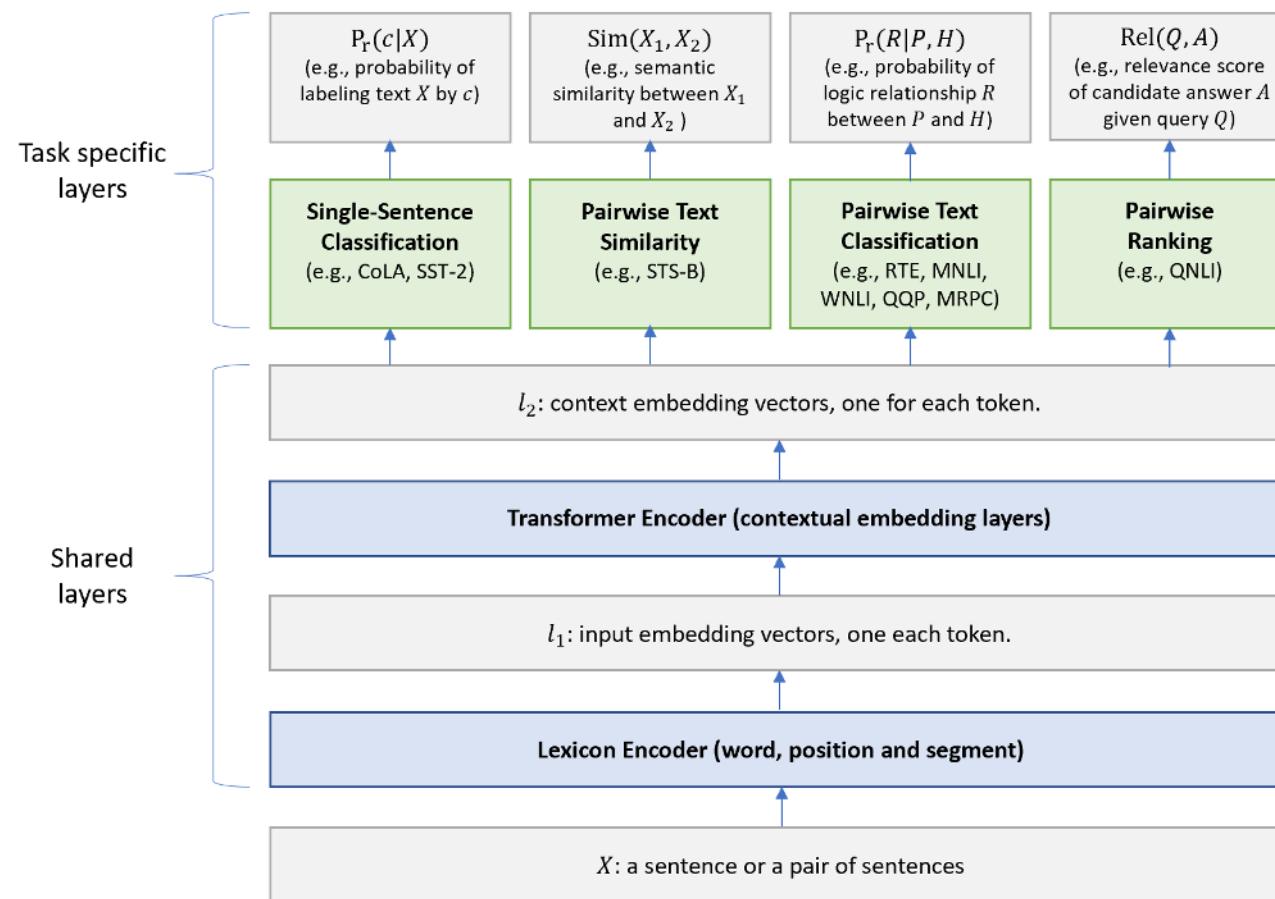
	MUC			B <sup>3</sup>			CEAF <sub>φ<sub>4</sub></sub>			Avg. F1
	P	R	F1	P	R	F1	P	R	F1	
Prev. SotA: (Lee et al., 2018)	81.4	79.5	80.4	72.2	69.5	70.8	68.2	67.1	67.6	73.0
Google BERT	84.9	82.5	83.7	76.7	74.2	75.4	74.6	70.1	72.3	77.1
Our BERT	85.1	83.5	84.3	77.3	75.5	76.4	75.0	71.9	73.9	78.3
Our BERT-1seq	85.5	84.1	84.8	77.8	76.7	77.2	75.3	73.5	74.4	78.8
SpanBERT	<b>85.8</b>	<b>84.8</b>	<b>85.3</b>	<b>78.3</b>	<b>77.9</b>	<b>78.1</b>	<b>76.4</b>	<b>74.2</b>	<b>75.3</b>	<b>79.6</b>

GLUE

	CoLA	SST-2	MRPC	STS-B	QQP	MNLI	QNLI	RTE	(Avg)
Google BERT	59.3	<b>95.2</b>	88.5/84.3	86.4/88.0	71.2/89.0	86.1/85.7	93.0	71.1	80.4
Our BERT	58.6	93.9	90.1/86.6	88.4/89.1	71.8/89.3	87.2/86.6	93.0	74.7	81.1
Our BERT-1seq	63.5	94.8	<b>91.2/87.8</b>	89.0/88.4	<b>72.1/89.5</b>	88.0/87.4	93.0	72.1	81.7
SpanBERT	<b>64.3</b>	94.8	<b>90.9/87.9</b>	<b>89.9/89.1</b>	<b>71.9/89.5</b>	<b>88.1/87.7</b>	<b>94.3</b>	<b>79.0</b>	<b>82.8</b>

# MT-DNN (BigBird) [Liu+ 2019]

- BERT + multi-task training



# MT-DNN (BigBird) [Liu+ 2019]

Model	CoLA 8.5k	SST-2 67k	MRPC 3.7k	STS-B 7k	QQP 364k	MNLI-m/mm 393k	QNLI 108k	RTE 2.5k	WNLI 634	AX	Score
BiLSTM+ELMo+Attn <sup>1</sup>	36.0	90.4	84.9/77.9	75.1/73.3	64.8/84.7	76.4/76.1	-	56.8	65.1	26.5	70.5
Singletask Pretrain Transformer <sup>2</sup>	45.4	91.3	82.3/75.7	82.0/80.0	70.3/88.5	82.1/81.4	-	56.0	53.4	29.8	72.8
GPT on STILTs <sup>3</sup>	47.2	93.1	87.7/83.7	85.3/84.8	70.1/88.1	80.8/80.6	-	69.1	65.1	29.4	76.9
BERT <sub>LARGE</sub> <sup>4</sup>	60.5	94.9	89.3/85.4	87.6/86.5	72.1/89.3	86.7/85.9	92.7	70.1	65.1	39.6	80.5
MT-DNN <sub>no-fine-tune</sub>	58.9	94.6	<b>90.1/86.4</b>	89.5/88.8	<b>72.7/89.6</b>	86.5/85.8	<b>93.1</b>	79.1	65.1	39.4	81.7
MT-DNN	<b>62.5</b>	<b>95.6</b>	<b>91.1/88.2</b>	<b>89.5/88.8</b>	<b>72.7/89.6</b>	<b>86.7/86.0</b>	<b>93.1</b>	<b>81.4</b>	65.1	<b>40.3</b>	<b>82.7</b>
Human Performance	66.4	97.8	86.3/80.8	92.7/92.6	59.5/80.4	92.0/92.8	91.2	93.6	95.9	-	87.1

# ERNIE 2.0 [Sun+ 2019] (Baidu)

a Continual Pre-training Framework for Language Understanding

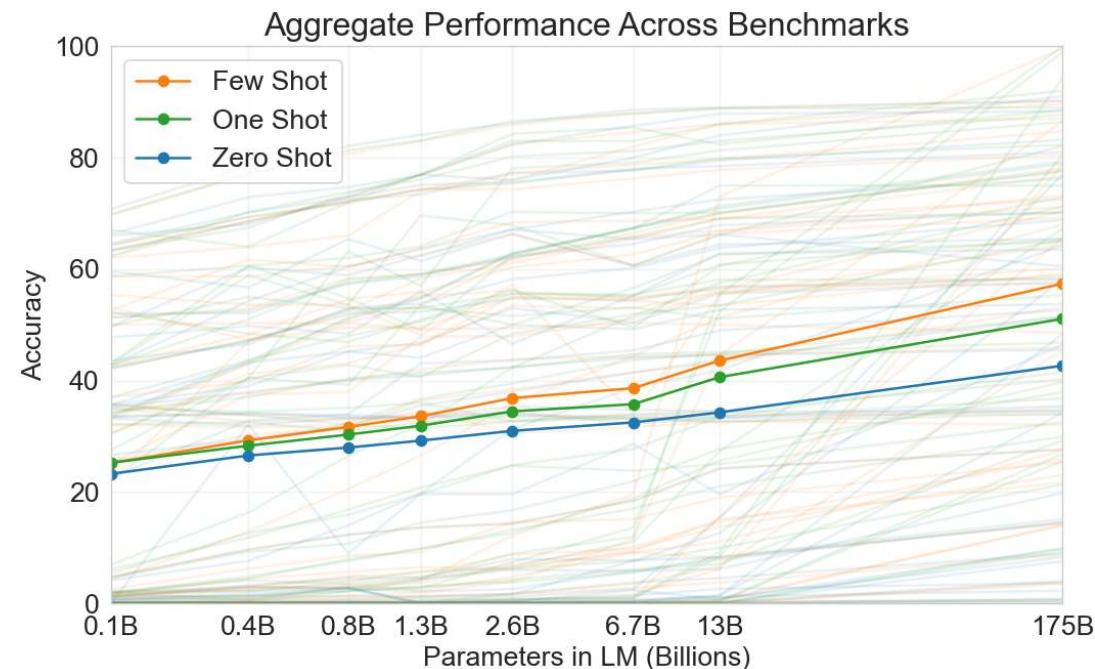
- Pre-training tasks are incrementally added
  - Task A → Tasks A + B → Tasks A + B + C → ...
- Pre-training tasks:
  - Word-aware tasks
    - Knowledge masking (phrase/entity masking; ERNIE 1.0)
    - Capitalization prediction
    - Token-document relation prediction
  - Structure-aware tasks
    - Sentence reordering
    - Sentence distance
  - Semantic-aware tasks
    - Discourse relation
    - IR relevance

# ERNIE 2.0 [Sun+ 2019]

Task(Metrics)	<i>BASE model</i>				<i>LARGE model</i>			
	Test		Dev			Test		
	BERT	ERNIE 2.0	BERT	XLNet	ERNIE 2.0	BERT	ERNIE 2.0	
CoLA (Matthew Corr.)	52.1	<b>55.2</b>	60.6	63.6	<b>65.4</b>	60.5	<b>63.5</b>	
SST-2 (Accuracy)	93.5	<b>95.0</b>	93.2	95.6	<b>96.0</b>	94.9	<b>95.6</b>	
MRPC (Accuracy/F1)	84.8/88.9	<b>86.1/89.9</b>	88.0/-	89.2/-	<b>89.7/-</b>	85.4/89.3	<b>87.4/90.2</b>	
STS-B (Pearson Corr./Spearman Corr.)	87.1/85.8	<b>87.6/86.5</b>	90.0/-	91.8/-	<b>92.3/-</b>	87.6/86.5	<b>91.2/90.6</b>	
QQP (Accuracy/F1)	89.2/71.2	<b>89.8/73.2</b>	91.3/-	91.8/-	<b>92.5/-</b>	89.3/72.1	<b>90.1/73.8</b>	
MNLI-m/mm (Accuracy)	84.6/83.4	<b>86.1/85.5</b>	86.6/-	<b>89.8/-</b>	89.1/-	86.7/85.9	<b>88.7/88.8</b>	
QNLI (Accuracy)	90.5	<b>92.9</b>	92.3	93.9	<b>94.3</b>	92.7	<b>94.6</b>	
RTE (Accuracy)	66.4	<b>74.8</b>	70.4	83.8	<b>85.2</b>	70.1	<b>80.2</b>	
WNLI (Accuracy)	<b>65.1</b>	<b>65.1</b>	-	-	-	65.1	<b>67.8</b>	
AX(Matthew Corr.)	34.2	<b>37.4</b>	-	-	-	39.6	<b>48.0</b>	
Score	78.3	<b>80.6</b>	-	-	-	80.5	<b>83.6</b>	

# Language Models

- GPT-2: zero-shot
- GPT-3: few-shot



The three settings we explore for in-context learning

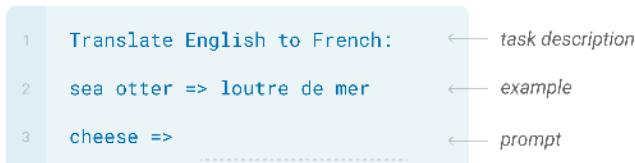
## Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.



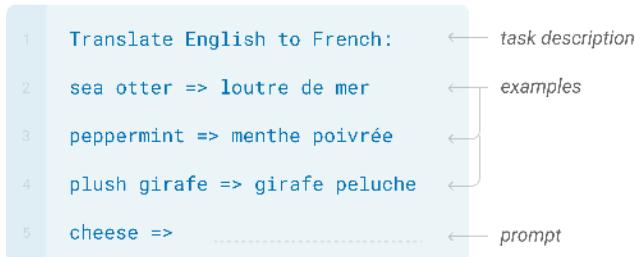
## One-shot

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.



## Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



Traditional fine-tuning (not used for GPT-3)

## Fine-tuning

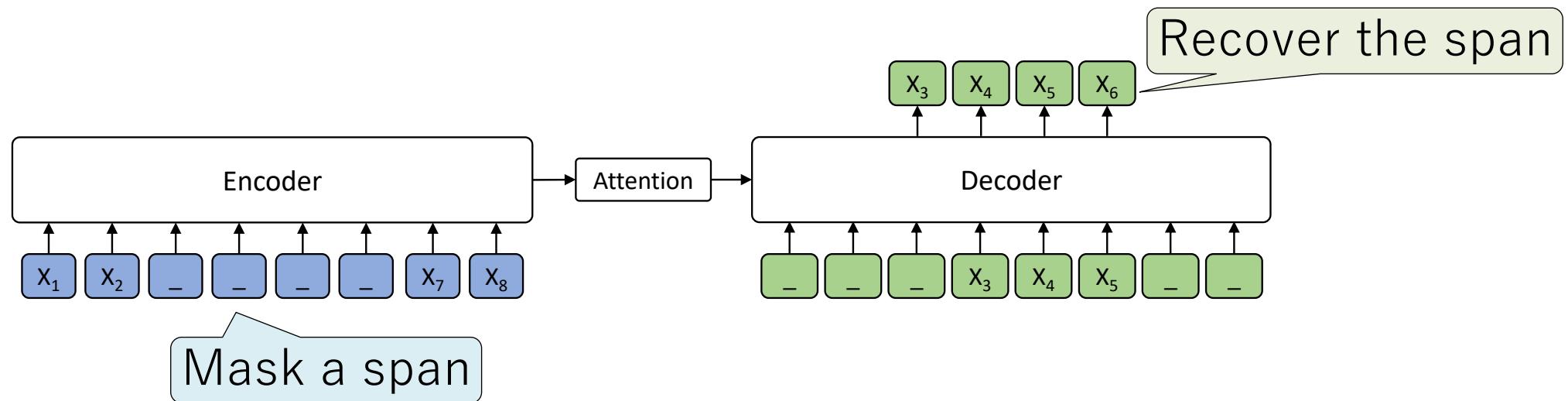
The model is trained via repeated gradient updates using a large corpus of example tasks.



[Brown+ 2020]

# Generation Models based on Pre-training

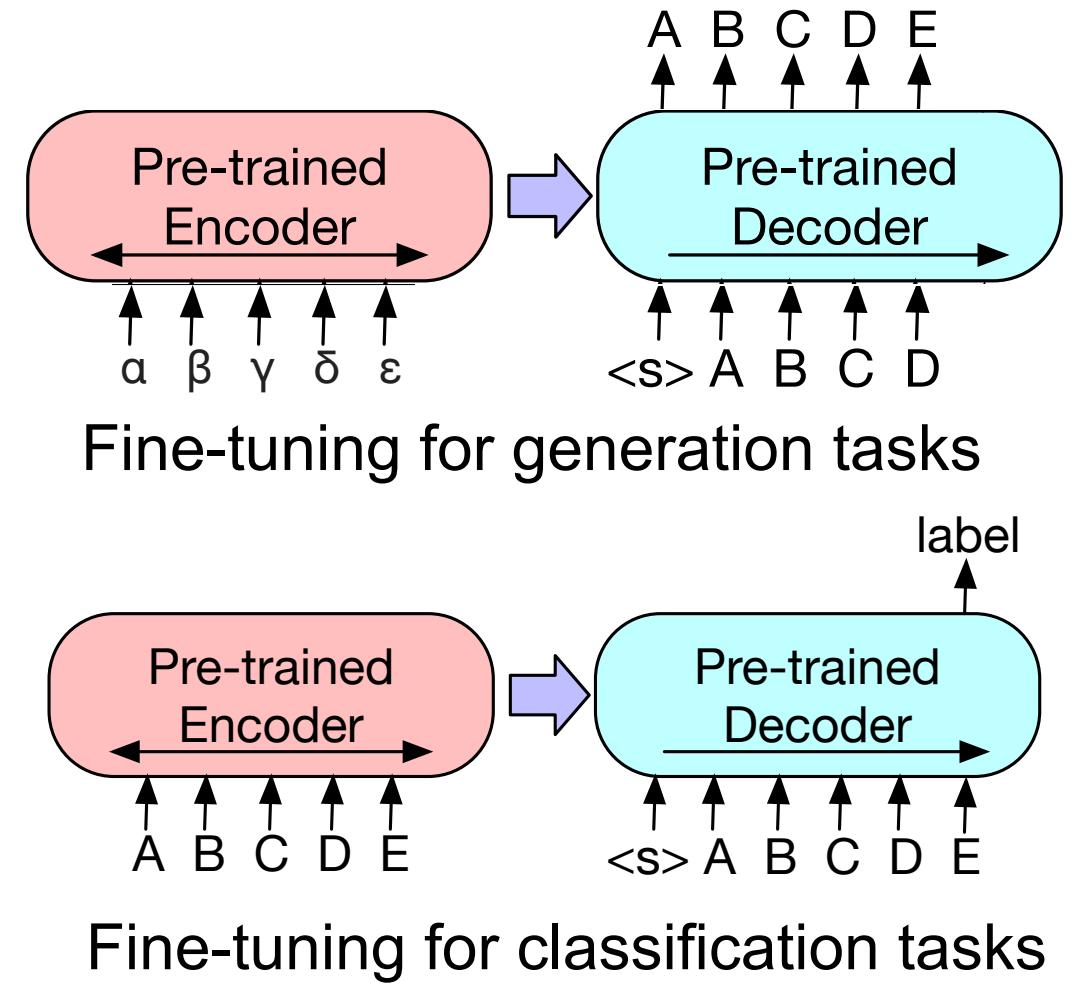
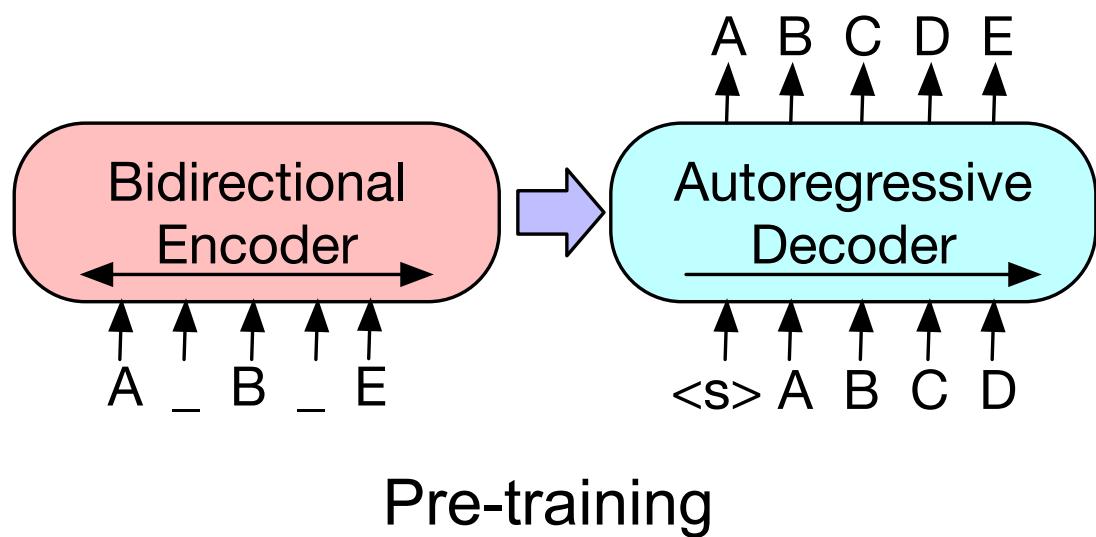
- MASS (Masked Sequence to Sequence Pre-training) [Song+ 19]



MASS was effective mainly for unsupervised machine translation

# BART (Bidirectional and Auto-Regressive Transformers)

[Lewis+ 19]



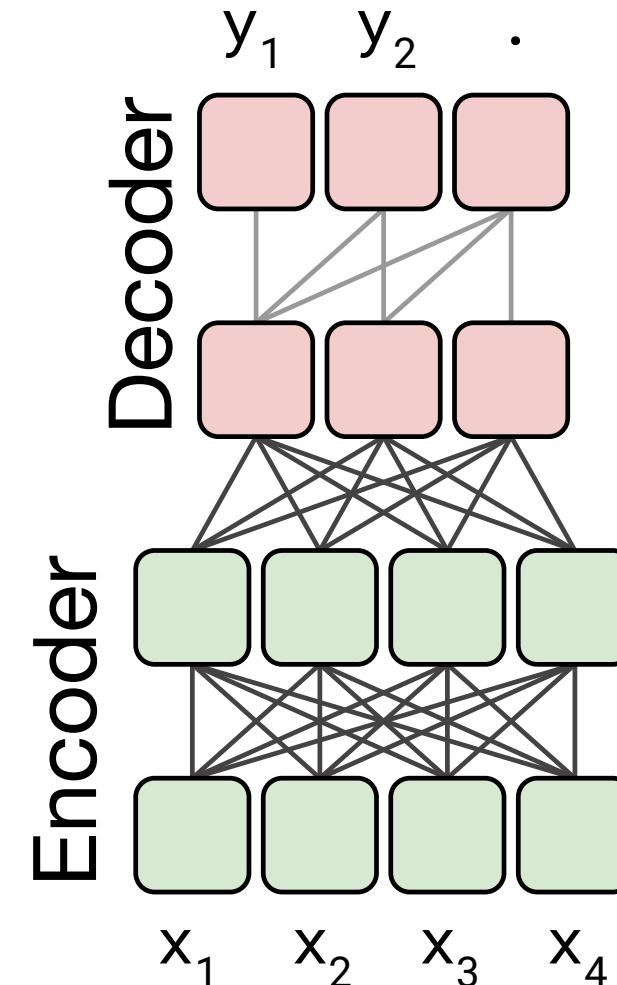
# Text-to-Text Transfer Transformer (T5) [Raffel+ 2019]

- All tasks are converted to a text-to-text format

task	source	target
COLA	cola sentence: The course is jumping well.	not acceptable
STS-B	stsbt sentence1: The rhino grazed on the grass. sentence2: A rhino is grazing in a field.	3.8
Summarization	summarize: state authorities dispatched emergency crews tuesday to survey the damage after an onslaught of severe weather in Mississippi ...	six people hospitalized after a storm in attala county.
MT	translate English to German: That is good.	Das ist gut.
WSC	wsc: The city councilmen refused the demonstrators a permit because *they* feared violence.	The city councilmen

# T5: Model

- An encoder-decoder model
  - Each consists of Transformer



# T5: Pre-training

Original text

Thank you ~~for inviting~~ me to your party ~~last~~ week.

Inputs

Thank you <X> me to your party <Y> week.

Targets

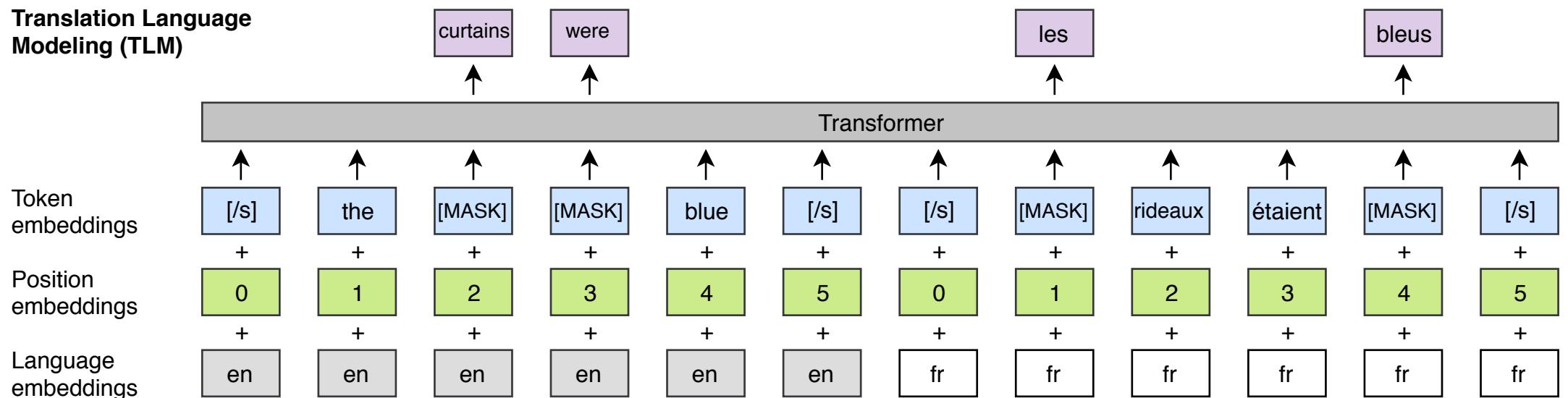
<X> for inviting <Y> last <Z>

# T5: Results

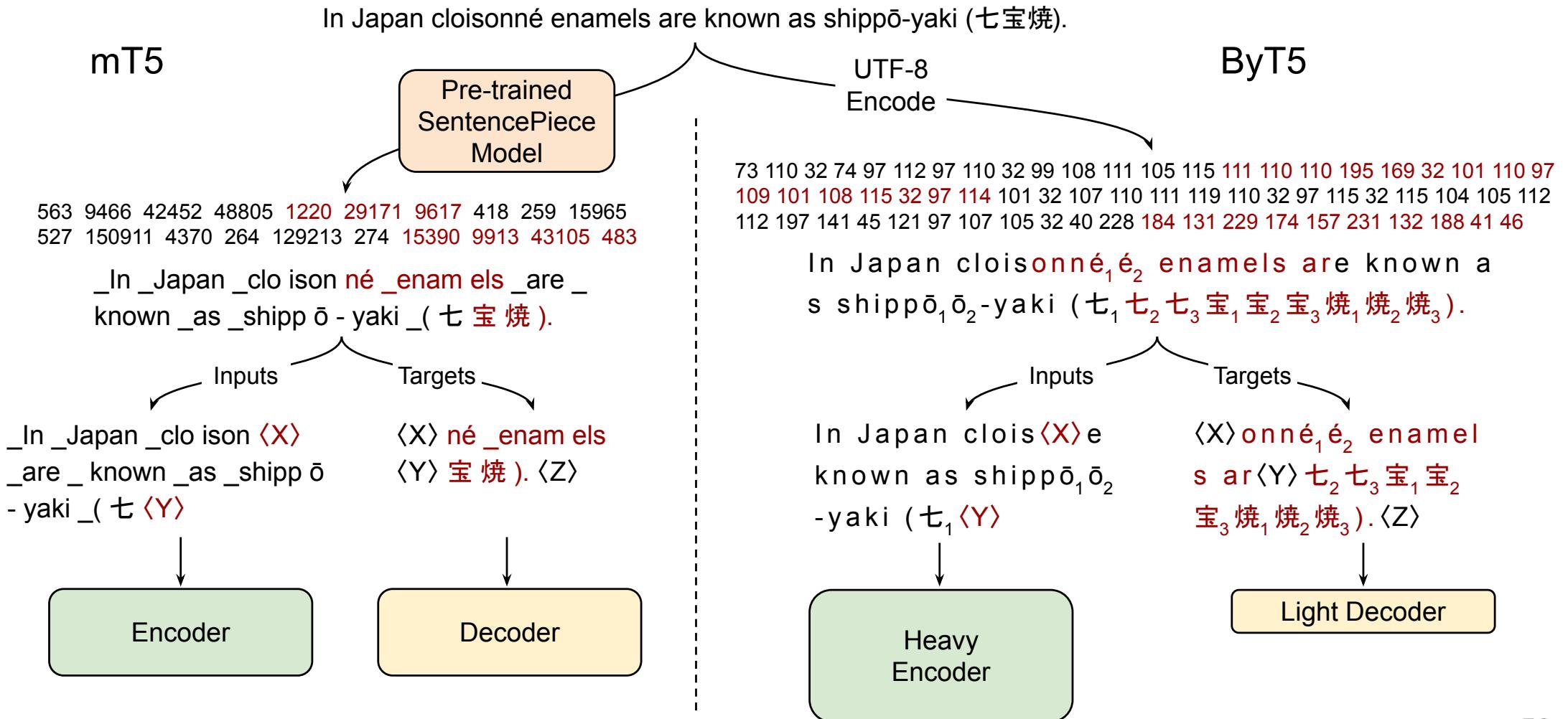
Model	GLUE Average	CoLA Matthew's	SST-2 Accuracy	MRPC F1	MRPC Accuracy	STS-B Pearson	STS-B Spearman
Previous best	89.4 <sup>a</sup>	69.2 <sup>b</sup>	<b>97.1<sup>a</sup></b>	<b>93.6<sup>b</sup></b>	<b>91.5<sup>b</sup></b>	<b>92.7<sup>b</sup></b>	<b>92.3<sup>b</sup></b>
T5-Small	77.4	41.0	91.8	89.7	86.6	85.6	85.0
T5-Base	82.7	51.1	95.2	90.7	87.5	89.4	88.6
T5-Large	86.4	61.2	96.3	92.4	89.9	89.9	89.2
T5-3B	88.5	67.1	97.4	92.5	90.0	90.6	89.8
T5-11B	<b>89.7</b>	<b>70.8</b>	<b>97.1</b>	91.9	89.2	92.5	92.1
Model	QQP F1	QQP Accuracy	MNLI-m Accuracy	MNLI-mm Accuracy	QNLI Accuracy	RTE Accuracy	WNLI Accuracy
Previous best	<b>74.8<sup>c</sup></b>	<b>90.7<sup>b</sup></b>	91.3 <sup>a</sup>	91.0 <sup>a</sup>	<b>99.2<sup>a</sup></b>	89.2 <sup>a</sup>	91.8 <sup>a</sup>
T5-Small	70.0	88.0	82.4	82.3	90.3	69.9	69.2
T5-Base	72.6	89.4	87.1	86.2	93.7	80.1	78.8
T5-Large	73.9	89.9	89.9	89.6	94.8	87.2	85.6
T5-3B	74.4	89.7	91.4	91.2	96.3	91.1	89.7
T5-11B	74.6	90.4	<b>92.0</b>	<b>91.7</b>	96.7	<b>92.5</b>	<b>93.2</b>
Model	SQuAD EM	SQuAD F1	SuperGLUE Average	BoolQ Accuracy	CB F1	CB Accuracy	COPA Accuracy
Previous best	88.95 <sup>d</sup>	94.52 <sup>d</sup>	84.6 <sup>e</sup>	87.1 <sup>e</sup>	90.5 <sup>e</sup>	95.2 <sup>e</sup>	90.6 <sup>e</sup>
T5-Small	79.10	87.24	63.3	76.4	56.9	81.6	46.0
T5-Base	85.44	92.08	76.2	81.4	86.2	94.0	71.2
T5-Large	86.66	93.79	82.3	85.4	91.6	94.8	83.4
T5-3B	88.53	94.95	86.4	89.9	90.3	94.4	92.0
T5-11B	<b>90.06</b>	<b>95.64</b>	<b>88.9</b>	<b>91.0</b>	<b>93.0</b>	<b>96.4</b>	<b>94.8</b>
Model	MultiRC F1a	MultiRC EM	ReCoRD F1	ReCoRD Accuracy	RTE Accuracy	WiC Accuracy	WSC Accuracy
Previous best	84.4 <sup>e</sup>	52.5 <sup>e</sup>	90.6 <sup>e</sup>	90.0 <sup>e</sup>	88.2 <sup>e</sup>	69.9 <sup>e</sup>	89.0 <sup>e</sup>
T5-Small	69.3	26.3	56.3	55.4	73.3	66.9	70.5
T5-Base	79.7	43.1	75.0	74.2	81.5	68.3	80.8
T5-Large	83.3	50.7	86.8	85.9	87.8	69.3	86.3
T5-3B	86.8	58.3	91.2	90.4	90.7	72.1	90.4
T5-11B	<b>88.2</b>	<b>62.3</b>	<b>93.3</b>	<b>92.5</b>	<b>92.5</b>	<b>76.1</b>	<b>93.8</b>
Model	WMT EnDe BLEU	WMT EnFr BLEU	WMT EnRo BLEU	CNN/DM ROUGE-1	CNN/DM ROUGE-2	CNN/DM ROUGE-L	
Previous best	<b>33.8<sup>f</sup></b>	<b>43.8<sup>f</sup></b>	<b>38.5<sup>g</sup></b>	43.47 <sup>h</sup>	20.30 <sup>h</sup>	40.63 <sup>h</sup>	
T5-Small	26.7	36.0	26.8	41.12	19.56	38.35	
T5-Base	30.9	41.2	28.0	42.05	20.34	39.40	
T5-Large	32.0	41.5	28.1	42.50	20.68	39.75	
T5-3B	31.8	42.6	28.2	42.72	21.02	39.94	
T5-11B	32.1	43.4	28.1	<b>43.52</b>	<b>21.55</b>	<b>40.69</b>	

# Multilingualization

- Pre-training on multilingual corpora (25-100 languages)
  - XLM [Lample+ 19], XLM-R [Conneau+ 19], mBART [Liu+ 20], mT5 [Xue+ 20], ByT5 [Xue+ 21]
- Pre-training on bilingual corpora (XLM [Lample+ 19])



# Token-free Model: ByT5 [Xue+ 21]

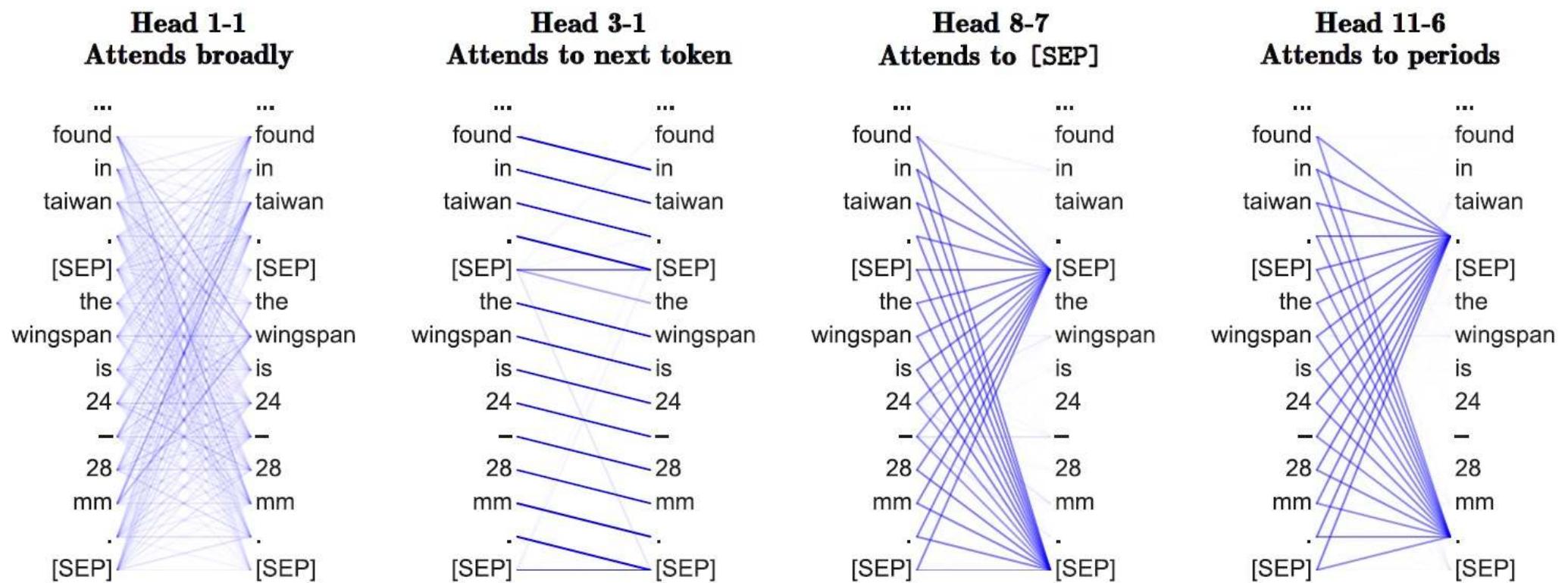


# Token-free Model: ByT5 [Xue+ 21]

	Small		Base		Large		XL		XXL	
	mT5	ByT5	mT5	ByT5	mT5	ByT5	mT5	ByT5	mT5	ByT5
<i>In-language multitask (models fine-tuned on gold data in all target languages)</i>										
WikiAnn NER	86.4	<b>90.6</b>	88.2	<b>91.6</b>	89.7	<b>91.8</b>	91.3	<b>92.6</b>	92.2	<b>93.7</b>
TyDiQA-GoldP	74.0 / 62.7	<b>82.6 / 73.6</b>	79.7 / 68.4	<b>86.4 / 78.0</b>	85.3 / 75.3	<b>87.7 / 79.2</b>	87.6 / 78.4	<b>88.0 / 79.3</b>	88.7 / 79.5	<b>89.4 / 81.4</b>
<i>Translate-train (models fine-tuned on English data plus translations in all target languages)</i>										
XNLI	72.0	<b>76.6</b>	79.8	<b>79.9</b>	<b>84.4</b>	82.8	<b>85.3</b>	85.0	<b>87.1</b>	85.7
PAWS-X	79.9	<b>88.6</b>	89.3	<b>89.8</b>	<b>91.2</b>	90.6	<b>91.0</b>	90.5	91.5	<b>91.7</b>
XQuAD	64.3 / 49.5	<b>74.0 / 59.9</b>	75.3 / 59.7	<b>78.5 / 64.6</b>	81.2 / 65.9	<b>81.4 / 67.4</b>	82.7 / 68.1	<b>83.7 / 69.5</b>	<b>85.2 / 71.3</b>	84.1 / 70.2
MLQA	56.6 / 38.8	<b>67.5 / 49.9</b>	67.6 / 48.5	<b>71.9 / 54.1</b>	73.9 / 55.2	<b>74.4 / 56.1</b>	75.1 / 56.6	<b>75.9 / 57.7</b>	<b>76.9 / 58.3</b>	<b>76.9 / 58.8</b>
TyDiQA-GoldP	49.8 / 35.6	<b>64.2 / 50.6</b>	66.4 / 51.0	<b>75.6 / 61.7</b>	75.7 / 60.1	<b>80.1 / 66.4</b>	80.1 / 65.0	<b>81.5 / 67.6</b>	<b>83.3 / 69.4</b>	83.2 / <b>69.6</b>
<i>Cross-lingual zero-shot transfer (models fine-tuned on English data only)</i>										
XNLI	67.5	<b>69.1</b>	<b>75.4</b>	<b>75.4</b>	<b>81.1</b>	79.7	<b>82.9</b>	82.2	<b>85.0</b>	83.7
PAWS-X	82.4	<b>84.0</b>	<b>86.4</b>	86.3	<b>88.9</b>	87.4	<b>89.6</b>	88.6	90.0	<b>90.1</b>
WikiAnn NER	50.5	<b>57.6</b>	55.7	<b>62.0</b>	58.5	<b>62.9</b>	<b>65.5</b>	61.6	<b>69.2</b>	67.7
XQuAD	58.1 / 42.5	<b>66.3 / 49.7</b>	<b>67.0 / 49.0</b>	66.6 / 48.1	<b>77.8 / 61.5</b>	61.5 / 43.9	<b>79.5 / 63.6</b>	57.7 / 43.0	<b>82.5 / 66.8</b>	79.7 / 63.6
MLQA	54.6 / 37.1	<b>60.9 / 43.3</b>	64.4 / 45.0	<b>66.6 / 47.3</b>	<b>71.2 / 51.7</b>	65.6 / 45.0	<b>73.5 / 54.4</b>	65.1 / 46.5	<b>76.0 / 57.4</b>	71.6 / 54.9
TyDiQA-GoldP	36.4 / 24.4	<b>54.9 / 39.9</b>	59.1 / 42.4	<b>69.6 / 54.2</b>	68.4 / 50.9	<b>75.4 / 59.4</b>	<b>77.8 / 61.8</b>	63.2 / 49.2	<b>82.0 / 67.3</b>	75.3 / 60.0

# Analysis of Inside BERT

## Analysis of attention [Clark+ 19]



# Multimodal Tasks of Language and Vision

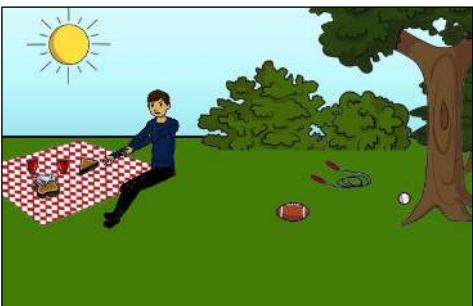
- Visual Question Answering (VQA) [Agrawal+ 15]



What color are her eyes?  
What is the mustache made of?



How many slices of pizza are there?  
Is this a vegetarian pizza?

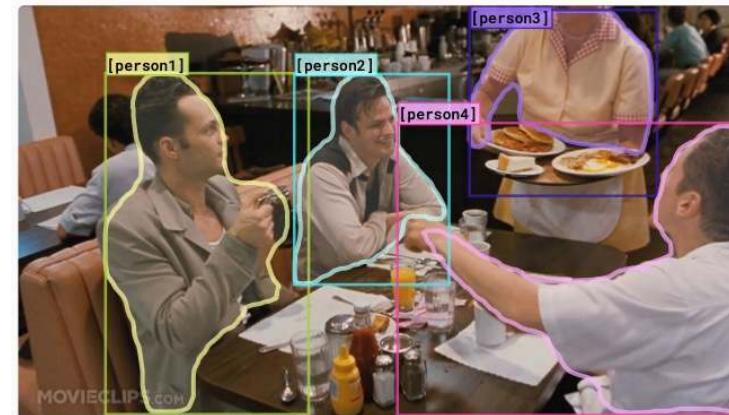


Is this person expecting company?  
What is just under the tree?



Does it appear to be rainy?  
Does this person have 20/20 vision?

- Visual Commonsense Reasoning (VCR) [Zellers+ 18]



hide all show all [person1] [person2] [person3] [person4]  
more objects »

Why is [person4] pointing at [person1]?

- a) He is telling [person3] that [person1] ordered the pancakes.
- b) He just told a joke.
- c) He is feeling accusatory towards [person1].
- d) He is giving [person1] directions.

Rationale: I think so because...

- a) [person1] has the pancakes in front of him.
- b) [person4] is taking everyone's order and asked for clarification.
- c) [person3] is looking at the pancakes both she and [person2] are smiling slightly.
- d) [person3] is delivering food to the table, and she might not know whose order is whose.

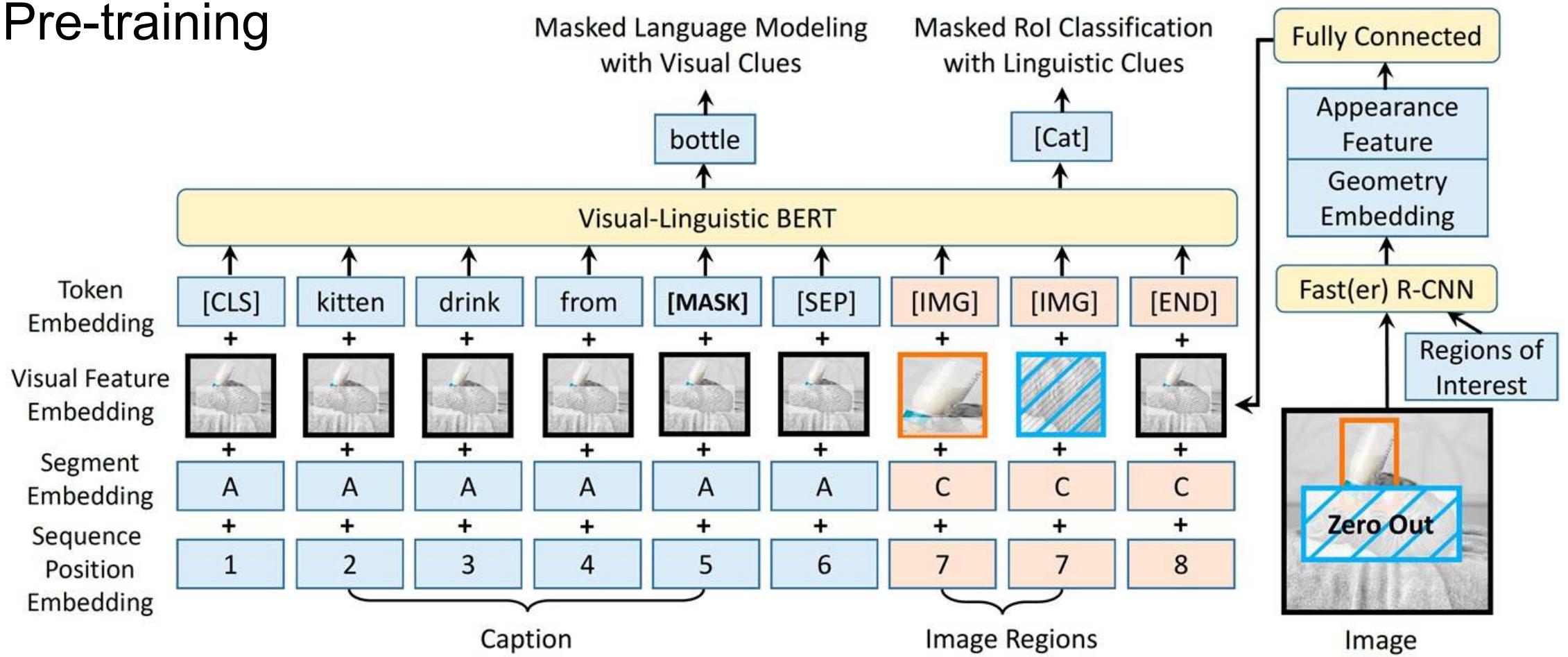
<https://visualcommonsense.com/>

# Models for Multimodal Tasks

- Input: 1 sequence
  - Encode a sequence that concatenates text and images
    - VideoBERT [Sun+ 19], VisualBERT [Li+ 19], Unicoder-VL [Li+ 19], VL-BERT [Su+ 19], UNITER [Chen+ 19], ImageBERT [Qi+ 20]
- Input: 2 sequences
  - Encode text and images respectively and then fuse them
    - CBT [Sun+ 19], ViLBERT [Lu+ 19], LXMERT [Tan+ 19], ERNIE-ViL [Yu+ 20]

# 1-sequence Model: VL-BERT [Su+ 19]

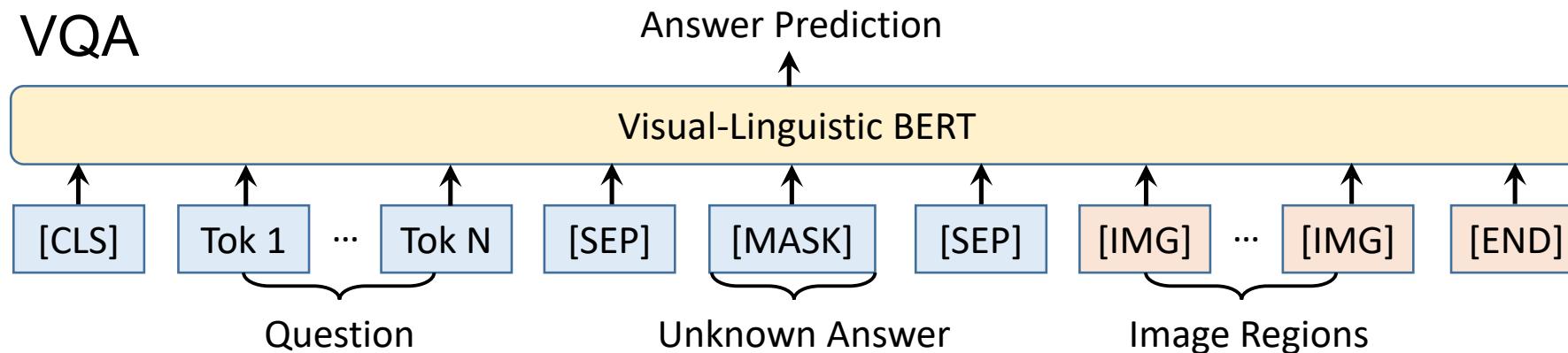
- Pre-training



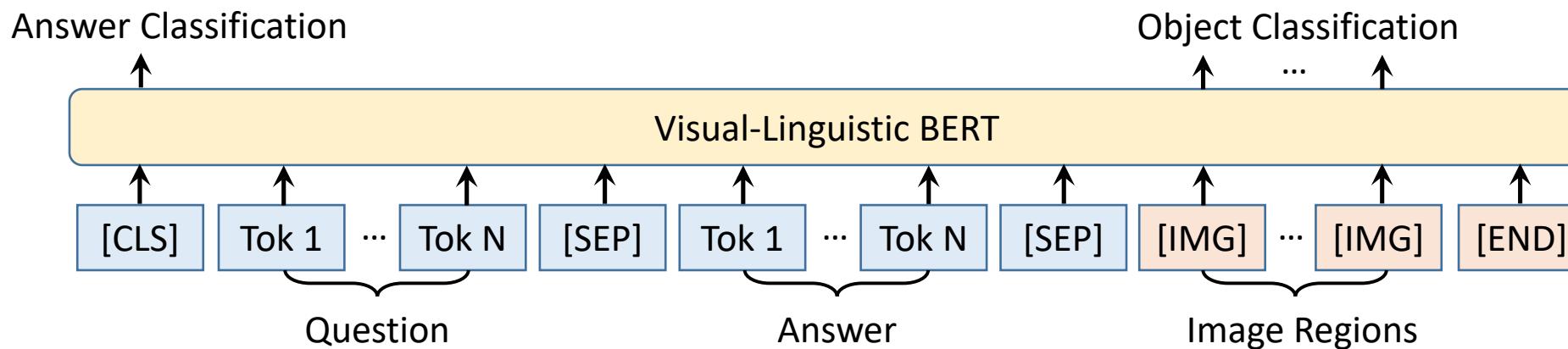
# 1-sequence Model: VL-BERT [Su+ 19]

- Fine-tuning

- VQA

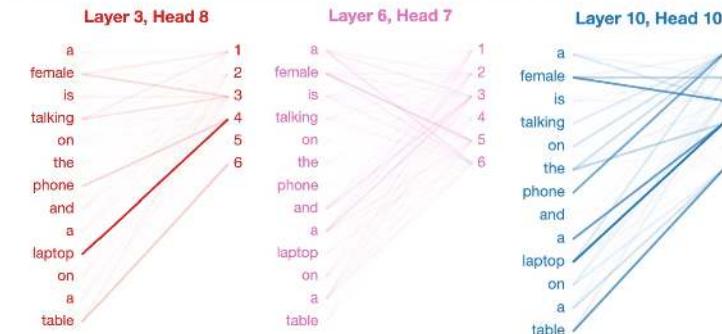
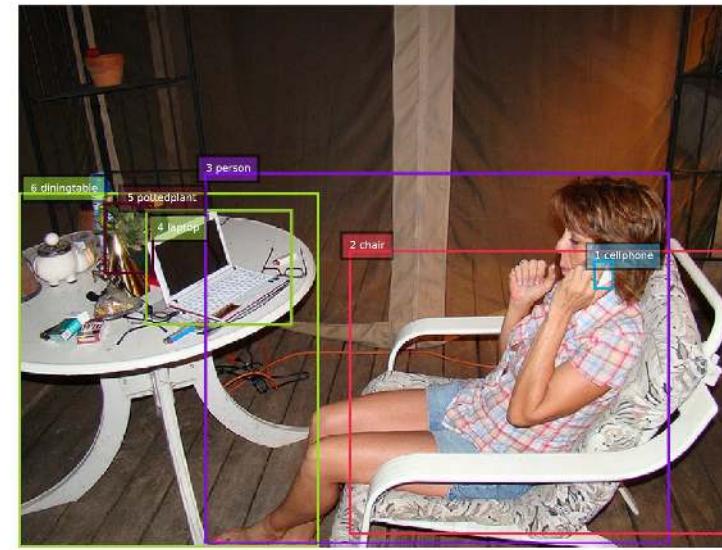
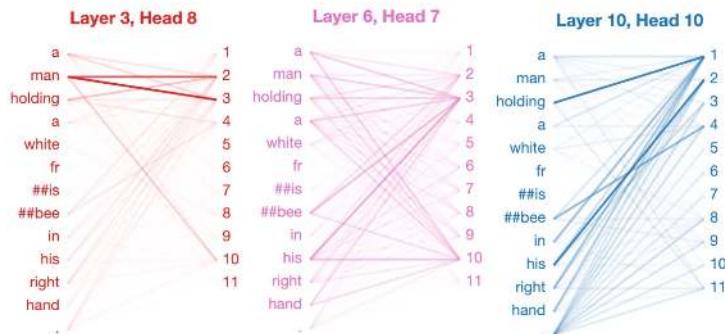
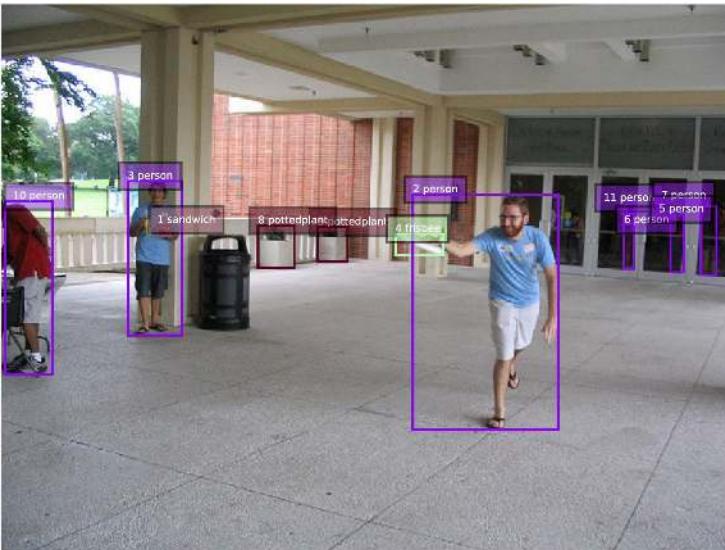


- VCR



# 1-sequence Model: VL-BERT [Su+ 19]

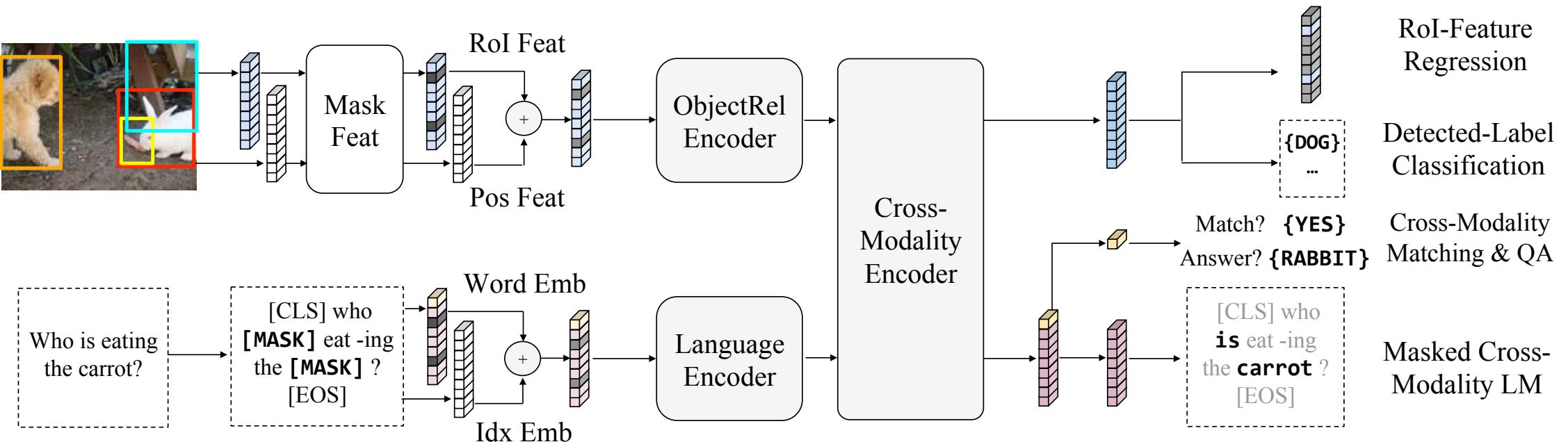
- Visualization



# 2-sequence Model: LXMERT

[Tan+ 19]

- Pre-training



Model	test-dev	test-std
BUTD (Anderson et al., 2018)	65.32	65.67
ViLBERT (Lu et al., 2019) <sup>†</sup>	70.55	70.92
VisualBERT (Li et al., 2019b) <sup>†</sup>	70.80	71.00
LXMERT (Tan & Bansal, 2019) <sup>†</sup>	72.42	72.54
VL-BERT <sub>BASE</sub> w/o pre-training	69.58	-
VL-BERT <sub>BASE</sub>	71.16	-
VL-BERT <sub>LARGE</sub>	71.79	72.22

Performance on VQA (from [Su+ 19])

# Libraries

- Pre-training
  - TensorFlow (by Google)
    - <https://github.com/google-research/bert> (TF 1.0)
    - <https://github.com/tensorflow/models/tree/master/official/nlp/bert> (TF 2.0)
      - Multi-GPU is supported
- Fine-tuning
  - TensorFlow (by Google)
  - PyTorch: Transformers (by Hugging Face)
    - <https://github.com/huggingface/transformers>
    - Latest models, such as RoBERTa and T5, are also supported