

# Paper Review

**[NAACL 2019] BERT: Pre-training of Deep Bidirectional Transformers  
For Language Understanding**

**Dept. of Computer Science & Engineering  
Artificial Intelligence and Data Mining Lab(AIDML)  
202122029 Meeyun Kim**

# Contents

**1. Introduction**

**2. About BERT**

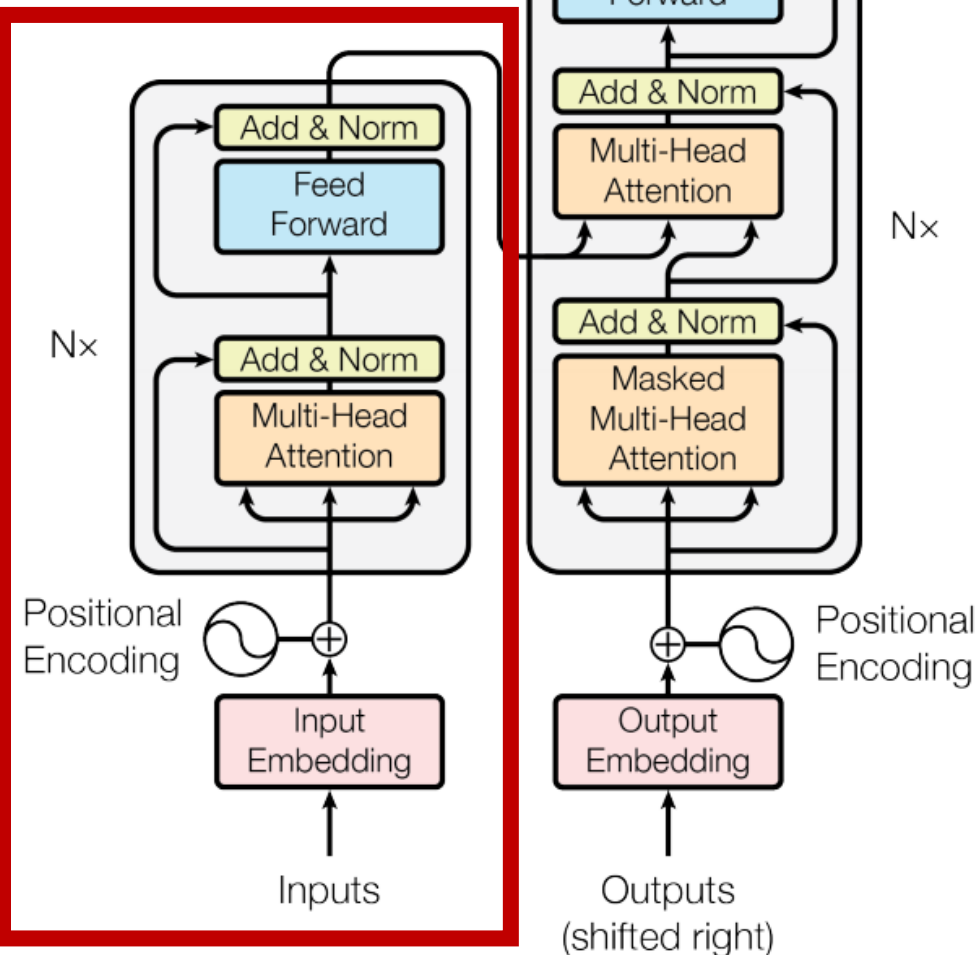
**3. Result**

# 1. Introduction

# 0) Transformer

- Encoder + Decoder 구조
- Transformer의 Encoder 부분을 발전시킨 것이 BERT!

**BERT!**



# 1) What is BERT?

- **Bidirectional Encoder Representations from Transformers**
- **Pre-Training, Fine-Tuning** 방식을 적용, 개선하여 transfer learning을 용이하게!
- **Masked Language Model(MLM)**을 이용한 성능 향상.

# \* Language Modeling

- 문장이 있을 때, 다음에 올 단어를 예측하는 task.



- 조건부 확률을 이용해 n번째 나올 단어의 확률 계산.

$$\prod_{i=1}^n P(w_i | w_1, \dots, w_{i-1})$$

- 확률값  $P(\text{books})$ ,  $P(\text{laptops})$ ,  $P(\text{minds})$ 를 비교.

# \* Language Modeling (Cont'd)



Q 인공지능 

Q 인공지능 - Google 검색

Q 인공지능 사례

Q 인공지능 스피커

Q 인공지능 수학

Q 인공지능 전문가

Q 인공지능 로봇

Q 인공지능 윤리

Q 인공지능 문제점

Q 인공지능 교육

Q 인공지능 기술

## 2) Why Bidirectional?

- **Unidirectional**한 구조는 성능을 제한할 수 있음.
- **left-to-right** 진행 구조 → **양쪽 문맥 이해 힘들.**
- ex) **OpenAI GPT** (next token만을 맞출 수 있는 language model 방식.)



## 2. About BERT

# 1) Model Architecture

- **BERT<sub>BASE</sub>** : L=12, H=768, A=12, Total Parameters=110M

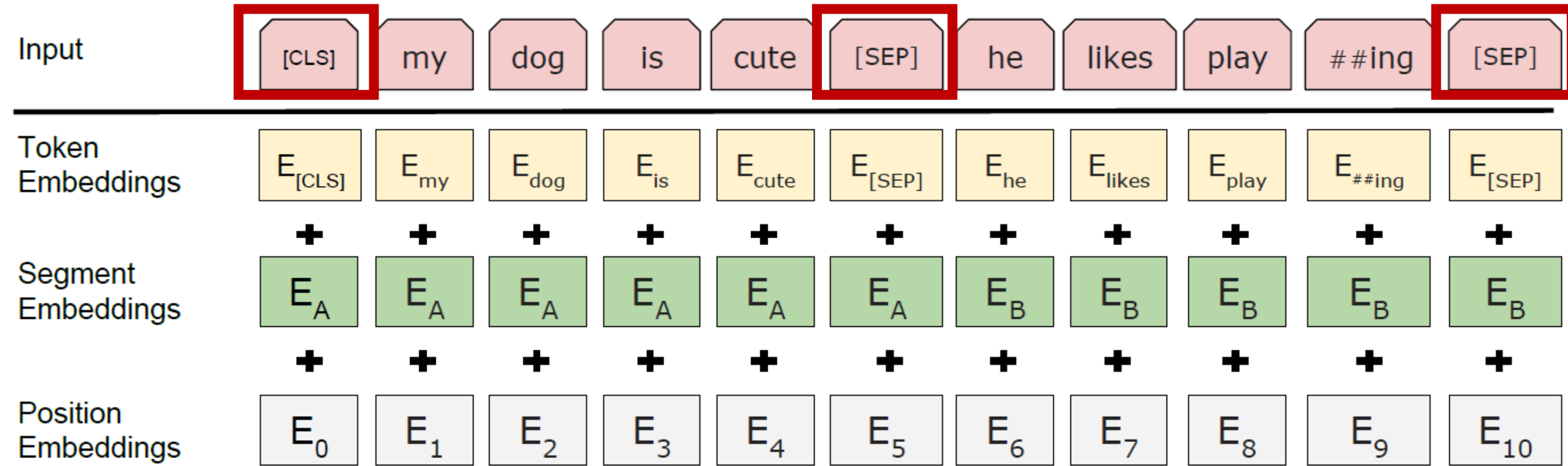
(OpenAI GPT와 hyper parameter가 동일!)

- **BERT<sub>LARGE</sub>** : L=24, H=1024, A=16, Total Parameters=340M

(number of layers - L, the hidden size - H, the number of self-attention heads - A)

- **BookCorpus**(800M words) + **English Wikipedia**(2,500M words)

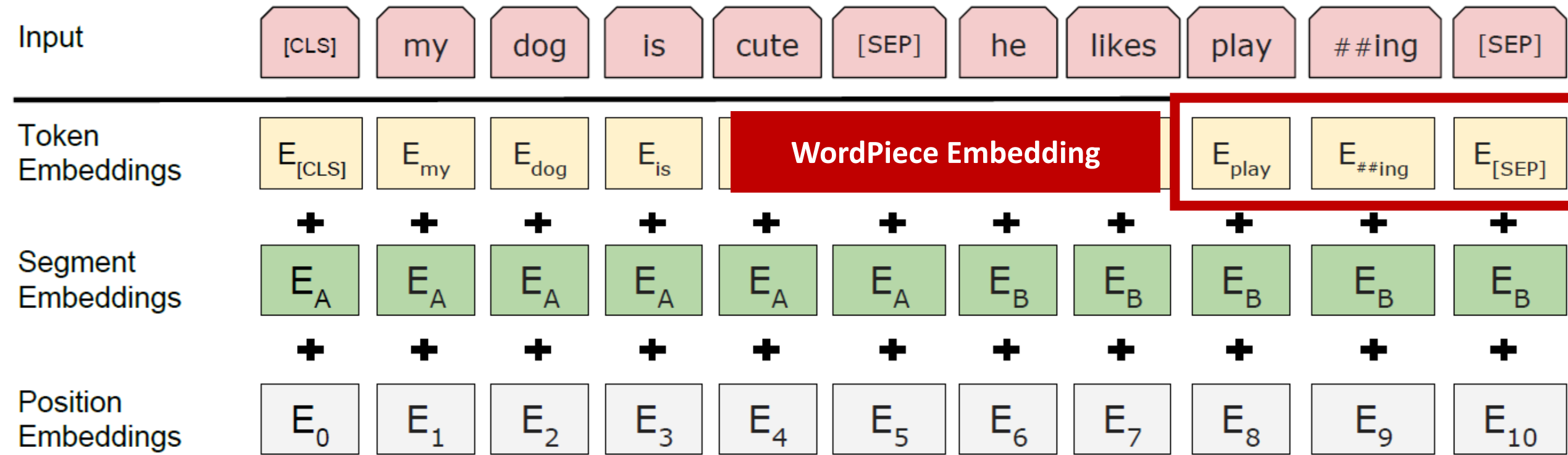
## 2) Input Representation



**[CLS]** : Classification task에 사용되기 위한 vector. 문장 전체가 하나의 vector로 표현된 special token.

**[SEP]** : 두 문장이 input으로 들어왔을 때, 이를 구분하기 위한 token.

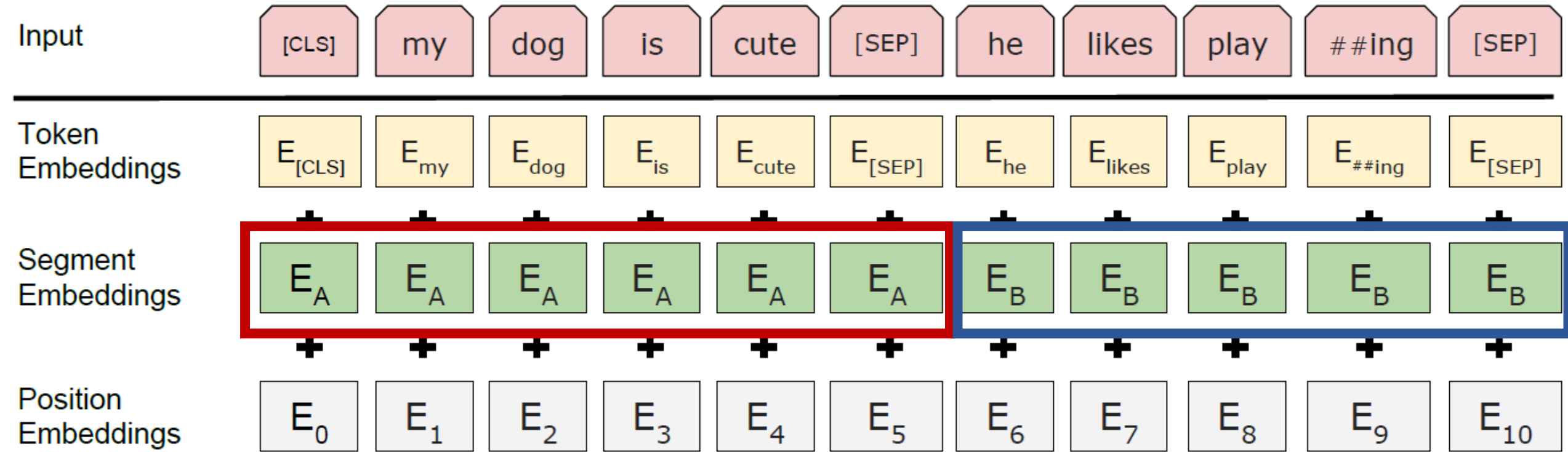
## 2) Input Representation (Cont'd)



### [WordPiece Embedding]

- 명확한 의미 전달이 가능해져서 model에 명확성 부여. (ex. play+ing)
- 신조어가 있는 입력에 대한 성능 향상. (ex. googling)

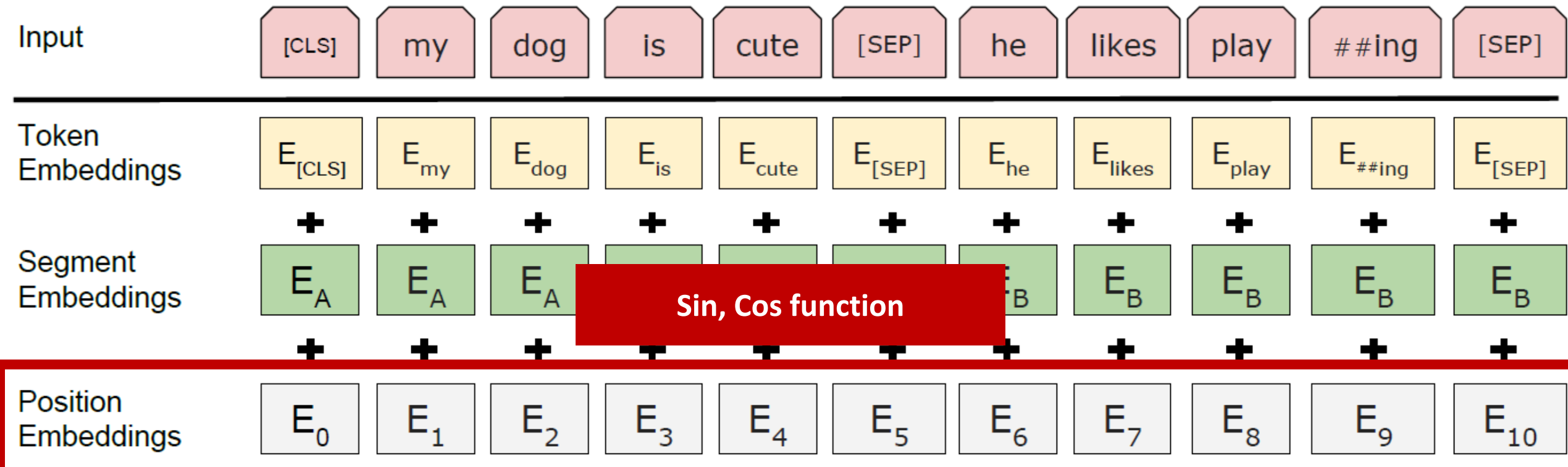
## 2) Input Representation (Cont'd)



### [Segment Embeddings]

- 두 개의 문장이 input일 때, 각 문장에 서로 다른 값을 더해 주어 문장 구분.
- $E_A$ 와  $E_B$ 는 고정된 값, 문장이 하나일 경우엔  $E_A$ 만 더해줌.

## 2) Input Representation (Cont'd)



### [Position Embeddings]

- sin과 cos의 출력 값은 입력에 따라 달라짐.
- 무한대 길이의 입력 값도 상대적인 위치로 출력이 가능.

### 3) Pre-training Tasks

**BERT = Pre-training + Fine-tuning**

A diagram showing the components of BERT. The word "Pre-training" in the equation above is enclosed in a red rectangular box. Two red lines originate from the bottom corners of this box and extend downwards to point at the two items listed below: "(1) Masked Language Model(MLM)" and "(2) Next sentence prediction".

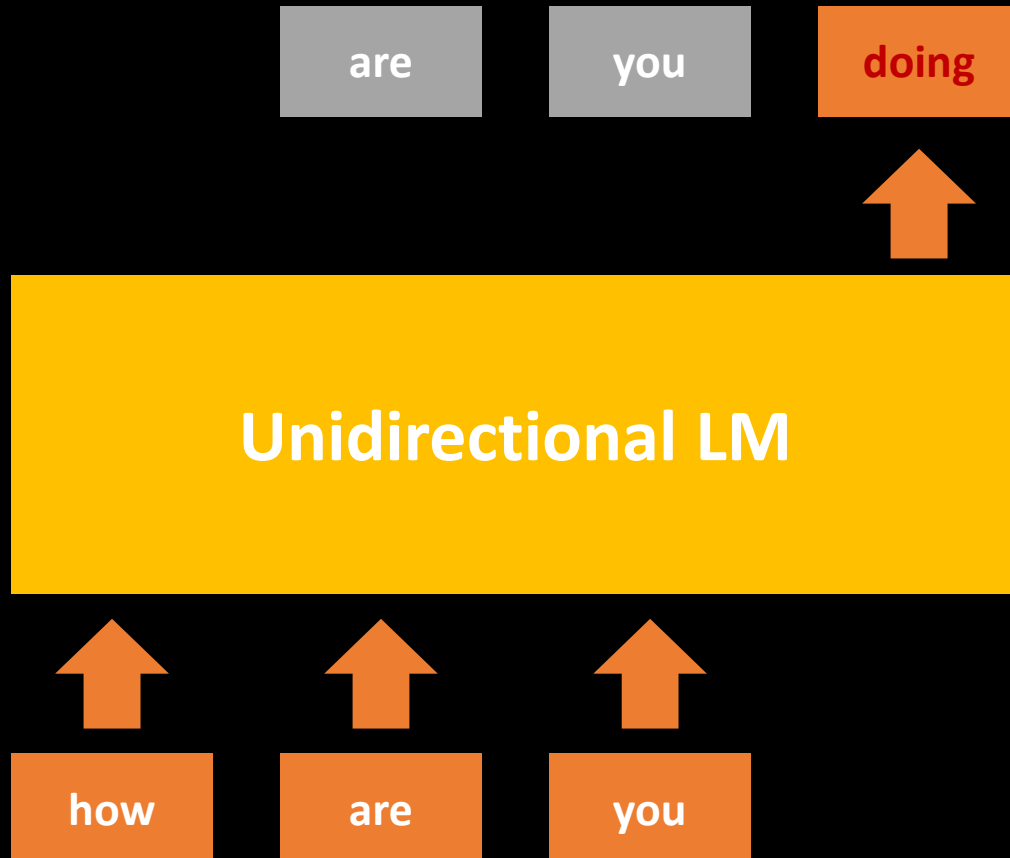
**(1) Masked Language Model(MLM)**

**(2) Next sentence prediction**

### 3) Pre-training Tasks

#### (1) Masked Language Model(MLM)

##### Unidirectional(Traditional) LM



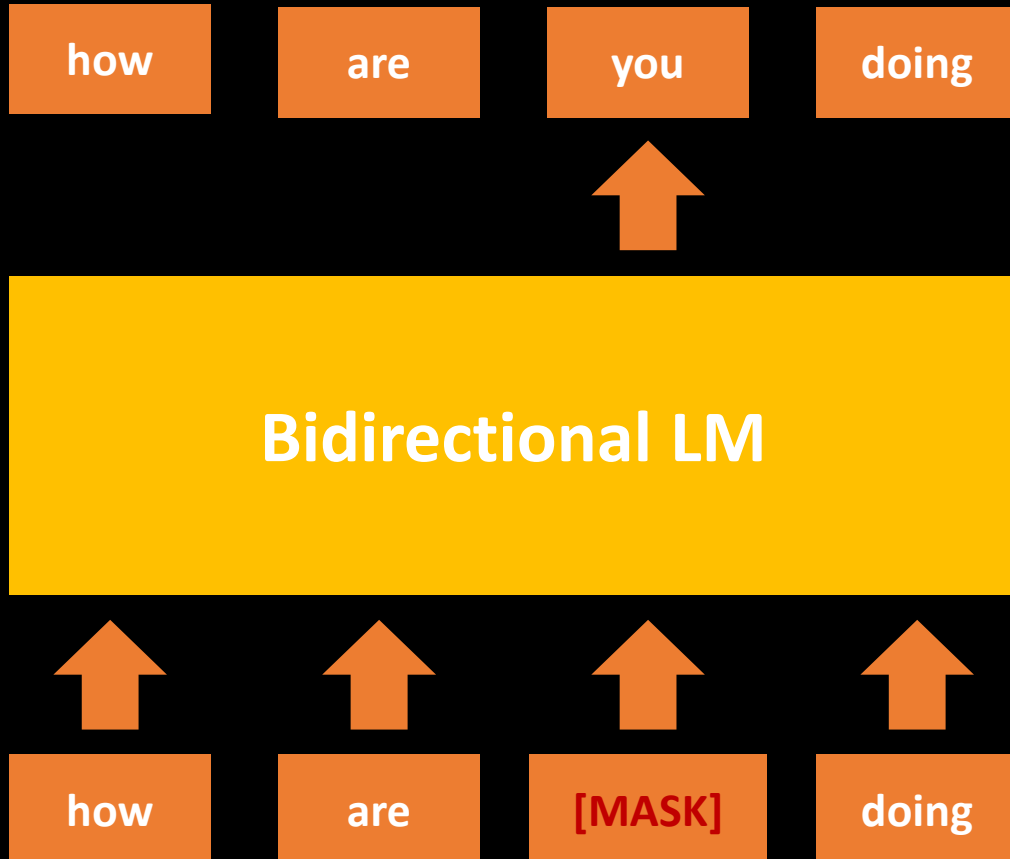
- Traditional한 Language model
- 단방향으로 학습됨. (ex. GPT)



### 3) Pre-training Tasks

#### (1) Masked Language Model(MLM) (Cont'd)

##### Bidirectional(Masked) LM(BERT)



- BERT에서 보여주는 MLM의 형태.

- 문장 전체를 학습하되, **[MASK]**로 가려진 단어를 예측하도록 학습됨.

### 3) Pre-training Tasks

#### (1) Masked Language Model(MLM) (Cont'd)

- Generator가 input 단어 중 random한 15%를 [MASK] token으로 바꿔줌.
  - 80% : [MASK] token
  - 10% : Random token
  - 10% : Unchanged token
- 이 과정을 통해 BERT는 context를 파악하는 능력을 기르게 됨.

### 3) Pre-training Tasks

#### (2) Next Sentence prediction (NSP)

- 두 문장 사이의 관계를 이해하기 위한 task를 위함. (like QA)
- Binarized next sentence prediction task



→ **50%** : Sentence B가 실제 input의 다음 문장(IsNext로 labeled)

→ **나머지 50%** : B가 Random sentence로 치환(NotNext로 labeled)

## 3) Pre-training Tasks

### (3) Pre-training Procedure

- Pre-training의 기본적인 절차는 LM에서 수행하는 것과 동일.
- **BERT\_English** : BookCorpus(800M words) + English Wikipedia(2,500M words)

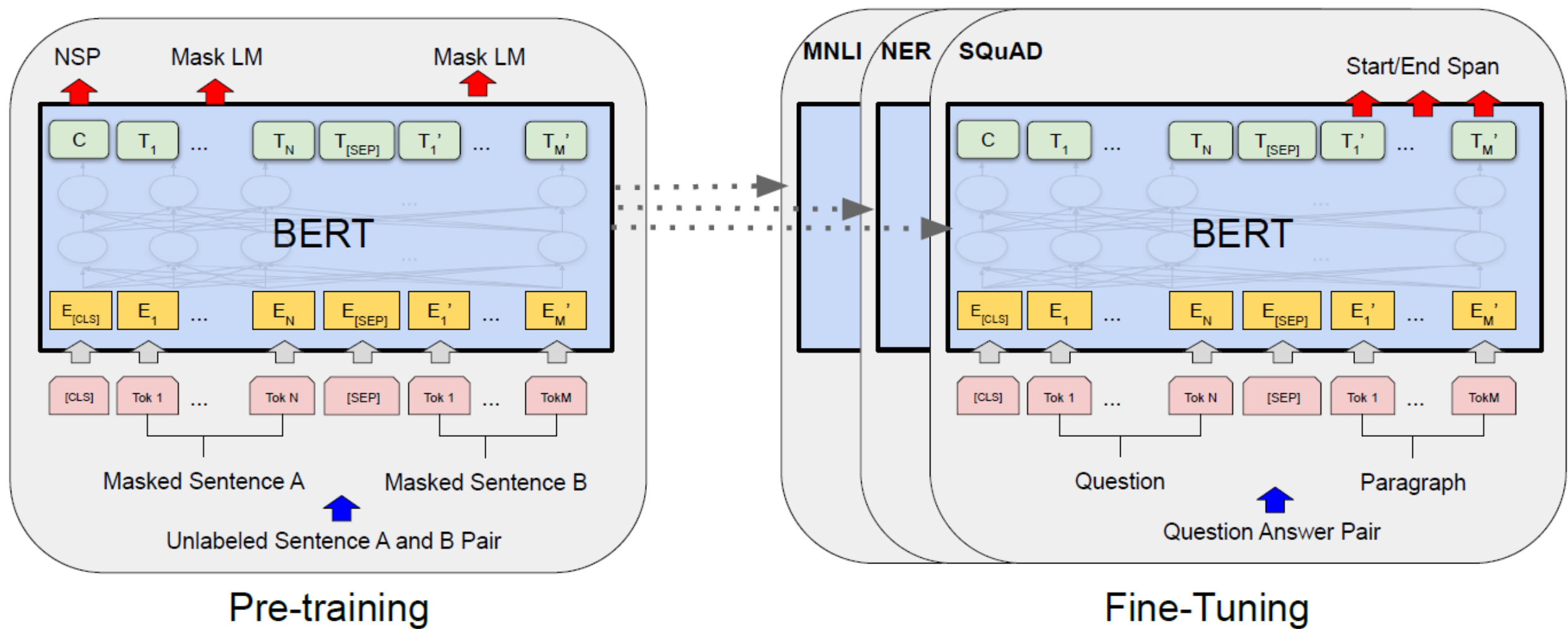
\* Wikipedia data - list, table, header등은 모두 제거. **Only text passage.**

→연속적인 시퀀스만을 추출, 학습시키기 위함.

## 4) Fine-tuning Tasks

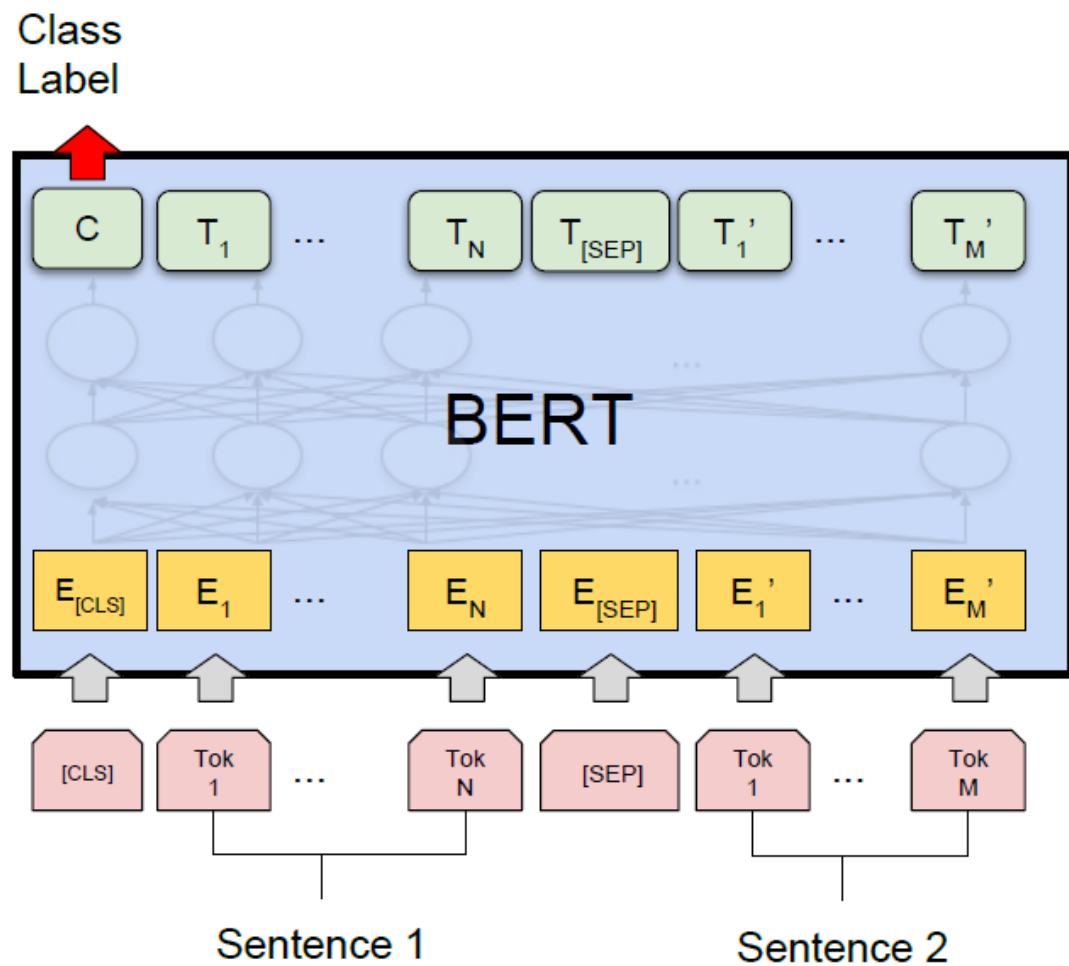
**BERT = Pre-training + Fine-tuning**

### 3) Fine-tuning Tasks

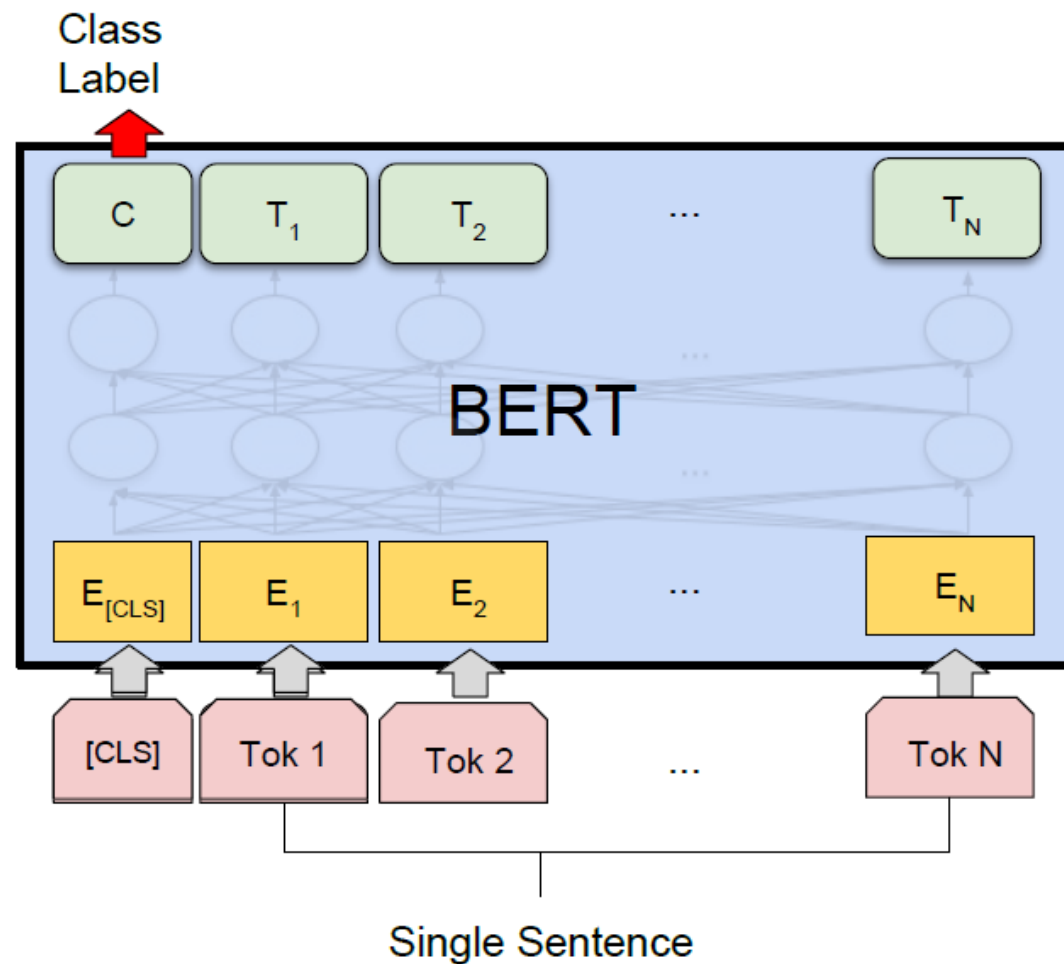


- 동일하게 pre-trained된 model들이 각자 다른 task를 수행하기 위해 모든 parameter를 **fine-tuning** 한다.

### 3) Fine-tuning Tasks

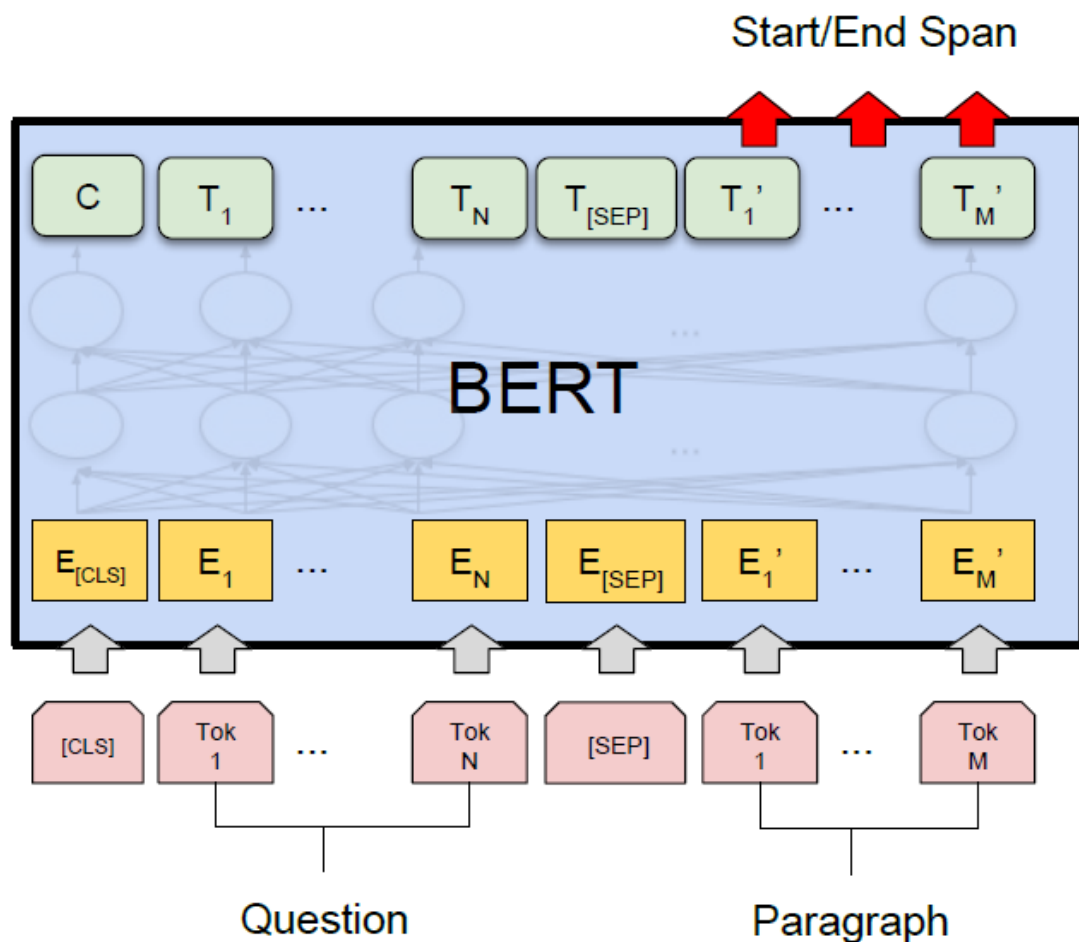


(a) Sentence Pair Classification Tasks:  
MNLI, QQP, QNLI, STS-B, MRPC,  
RTE, SWAG

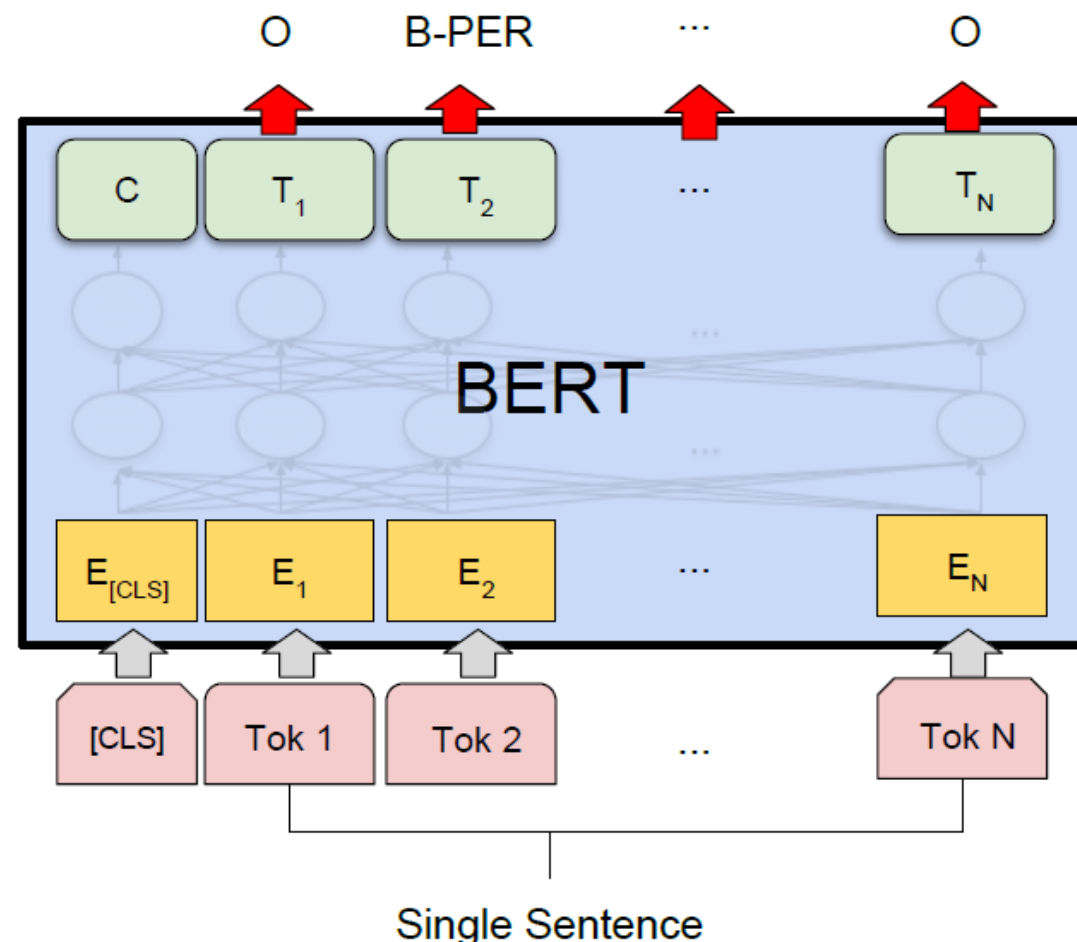


(b) Single Sentence Classification Tasks:  
SST-2, CoLA

### 3) Fine-tuning Tasks



(c) Question Answering Tasks:  
SQuAD v1.1



(d) Single Sentence Tagging Tasks:  
CoNLL-2003 NER



# 3. Result

# 1) GLUE Results

— : Single Sequence task  
 — : Sequence Pair task

Dataset	Description	Data example	Metric
CoLA	Is the sentence grammatical or ungrammatical?	"This building is than that one." = <b>Ungrammatical</b>	Matthews
SST-2	Is the movie review positive, negative, or neutral?	"The movie is funny , smart , visually inventive , and most of all , alive ." = <b>.93056 (Very Positive)</b>	Accuracy
MRPC	Is the sentence B a paraphrase of sentence A?	A) "Yesterday , Taiwan reported 35 new infections , bringing the total number of cases to 418 ." B) "The island reported another 35 probable cases yesterday , taking its total to 418 ." = <b>A Paraphrase</b>	Accuracy / F1
STS-B	How similar are sentences A and B?	A) "Elephants are walking down a trail." B) "A herd of elephants are walking along a trail." = <b>4.6 (Very Similar)</b>	Pearson / Spearman
QQP	Are the two questions similar?	A) "How can I increase the speed of my internet connection while using a VPN?" B) "How can Internet speed be increased by hacking through DNS?" = <b>Not Similar</b>	Accuracy / F1
MNLI-mm	Does sentence A entail or contradict sentence B?	A) "Tourist Information offices can be very helpful." B) "Tourist Information offices are never of any help." = <b>Contradiction</b>	Accuracy
QNLI	Does sentence B contain the answer to the question in sentence A?	A) "What is essential for the mating of the elements that create radio waves?" B) "Antennas are required by any radio receiver or transmitter to couple its electrical connection to the electromagnetic field." = <b>Answerable</b>	Accuracy
RTE	Does sentence A entail sentence B?	A) "In 2003, Yunus brought the microcredit revolution to the streets of Bangladesh to support more than 50,000 beggars, whom the Grameen Bank respectfully calls Struggling Members." B) "Yunus supported more than 50,000 Struggling Members." = <b>Entailed</b>	Accuracy
WNLI	Sentence B replaces sentence A's ambiguous pronoun with one of the nouns - is this the correct noun?	A) "Lily spoke to Donna, breaking her concentration." B) "Lily spoke to Donna, breaking Lily's concentration." = <b>Incorrect Referent</b>	Accuracy

# 1) GLUE Results

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average -
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	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>92.7</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>82.1</b>

- **GLUE** → **sequence classification task** (한 문장, 혹은 두 문장을 이용한 task)
- **BERT** model이 모든 task에서 **SOTA**를 달성.

## 2) SQuAD v1.1 Results

- SQuAD → Question Answering task
- BERT<sub>LARGE</sub>가 SOTA를 달성. (with wide margin)

System	Dev		Test	
	EM	F1	EM	F1
Top Leaderboard Systems (Dec 10th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
Published				
BiDAF+ELMo (Single)	-	85.6	-	85.8
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT <sub>BASE</sub> (Single)	80.8	88.5	-	-
BERT <sub>LARGE</sub> (Single)	84.1	90.9	-	-
BERT <sub>LARGE</sub> (Ensemble)	85.8	91.8	-	-
BERT <sub>LARGE</sub> (Sgl.+TriviaQA)	<b>84.2</b>	<b>91.1</b>	<b>85.1</b>	<b>91.8</b>
BERT <sub>LARGE</sub> (Ens.+TriviaQA)	<b>86.2</b>	<b>92.2</b>	<b>87.4</b>	<b>93.2</b>

### 3) SWAG Results

- **SWAG** → grounded common-sense inference task
- **BERT<sub>LARGE</sub>**가 **SOTA**를 달성. (with wide margin)

System	Dev	Test
ESIM+GloVe	51.9	52.7
ESIM+ELMo	59.1	59.2
OpenAI GPT	-	78.0
BERT <sub>BASE</sub>	81.6	-
BERT <sub>LARGE</sub>	<b>86.6</b>	<b>86.3</b>
Human (expert) <sup>†</sup>	-	85.0
Human (5 annotations) <sup>†</sup>	-	88.0

## 4) Ablation Studies

### (1) Effect of Pre-training Tasks

Tasks	Dev Set				
	MNLI-m (Acc)	QNLI (Acc)	MRPC (Acc)	SST-2 (Acc)	SQuAD (F1)
BERT <sub>BASE</sub>	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

- 논문에서 중요하다고 언급한 요소들을 제거해보며 중요성을 파악.
- **NSP**(Next Sentence Prediction)만 제거한 것, **MLM**도 제거한 것들과 결과 비교.

## 4) Ablation Studies

### (1) Effect of Pre-training Tasks (Cont'd)

Tasks	Dev Set				
	MNLI-m (Acc)	QNLI (Acc)	MRPC (Acc)	SST-2 (Acc)	SQuAD (F1)
BERT <sub>BASE</sub>	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

- **NSP**를 빼면, 두 문장 간의 구조 파악을 해야 하는 **NLI** 쪽에서 성능 크게 감소.
- **MLM** 대신 **LTR**을 쓰면 성능 하락은 더욱 심각해짐.

## 4) Ablation Studies

### (2) Effect of Model Size

Hyperparams				Dev Set Accuracy		
#L	#H	#A	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



- 모델이 커질수록(hyperparameter 수가 늘어날수록), 정확도가 증가한다.



**QnA**