

# Improving Language Understanding by Generative Pre-Training

---

**Name**

조병웅

NLP

2024/05/14



# Contents

---

- **Abstract**
- **Introduction**
- **Framework**
  - Unsupervised pre-training
  - Supervised fine-tuning
  - Task-specific input transformations
- **Experiments**
- **Supervised fine-tuning**
- **Analysis**
- **Conclusion**

# Abstract

---

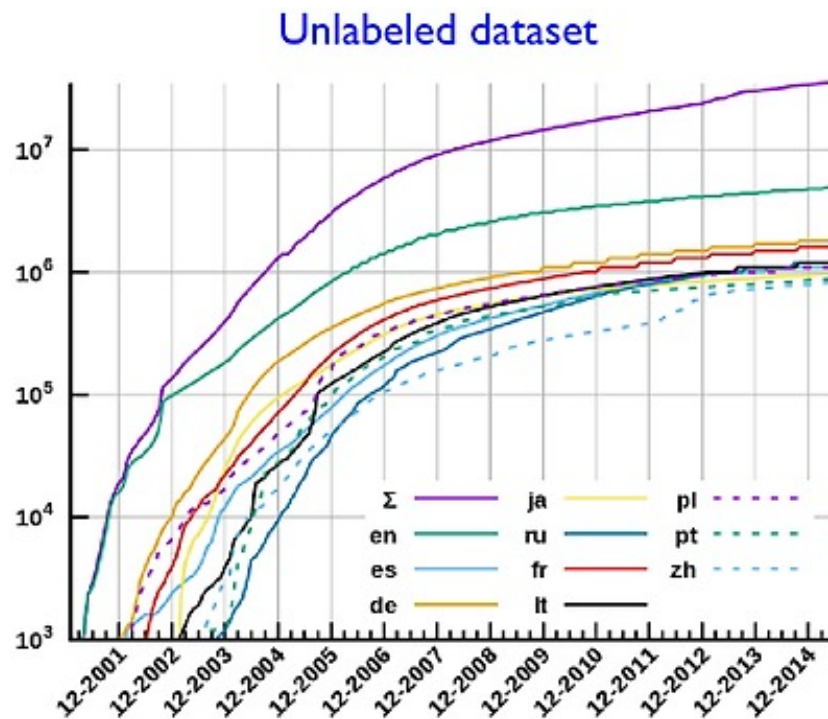
**NLP에는 다양한 Task 존재! 그러나,  
Labeled Data는 많이 없음.**

- Text Classification
- Text Entailment
- Similarity
- Multiple Choices

**Solution** -> 세상에 넘치는 Unlabeled Data를 활용해 LM 구축해보자!

# Abstract

## Unlabeled vs Labeled



### Labeled dataset

- STS Benchmark for sentence similarity: 8,628 sentences
- Quora question pairs: 404,290 question pairs
- CoLA dataset: 10,657 sentences

As of 24 February 2020, there are **6,020,081** articles in the [English Wikipedia](#) containing over **3.5 billion words**.

# Abstract

---

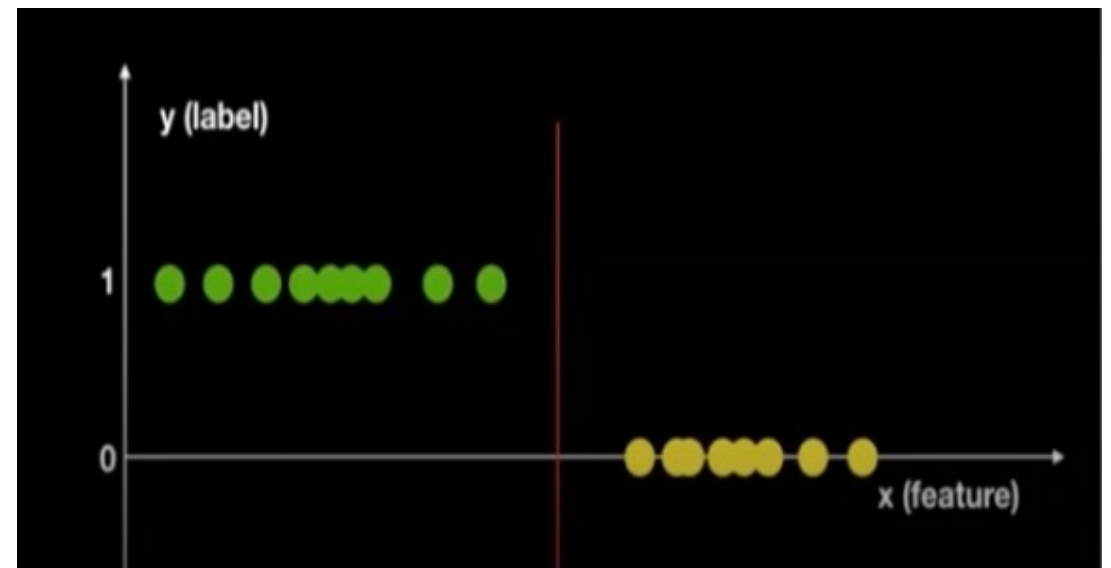
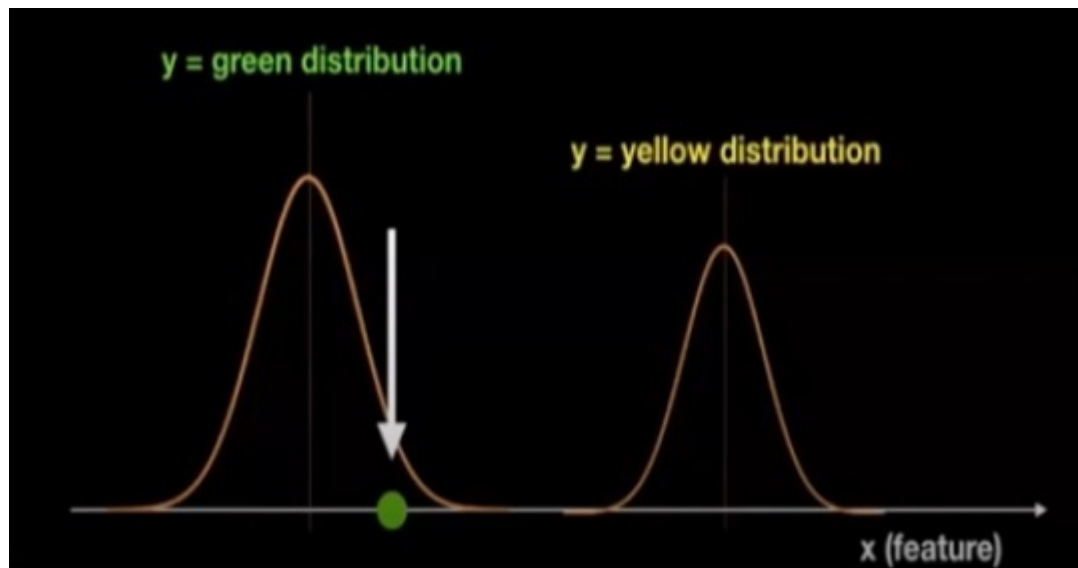
GPT is **Generative** **Pre-Training**

**Generative Learning vs Discriminative Learning**

**What is Pre-Training?**

# Abstract

## Generative Learning vs Discriminative Learning



# Abstract

---

## Generative Learning vs Discriminative Learning

### Generative Learning

라벨값의 분포를 활용, 데이터 '분포'를 모델링  
데이터가 더 많을 때 좋고, 데이터 과적합 확률이 적음  
So, 실제 분포를 나타낼 만큼 **충분한 데이터, 충분한 훈련 시간** 필요.

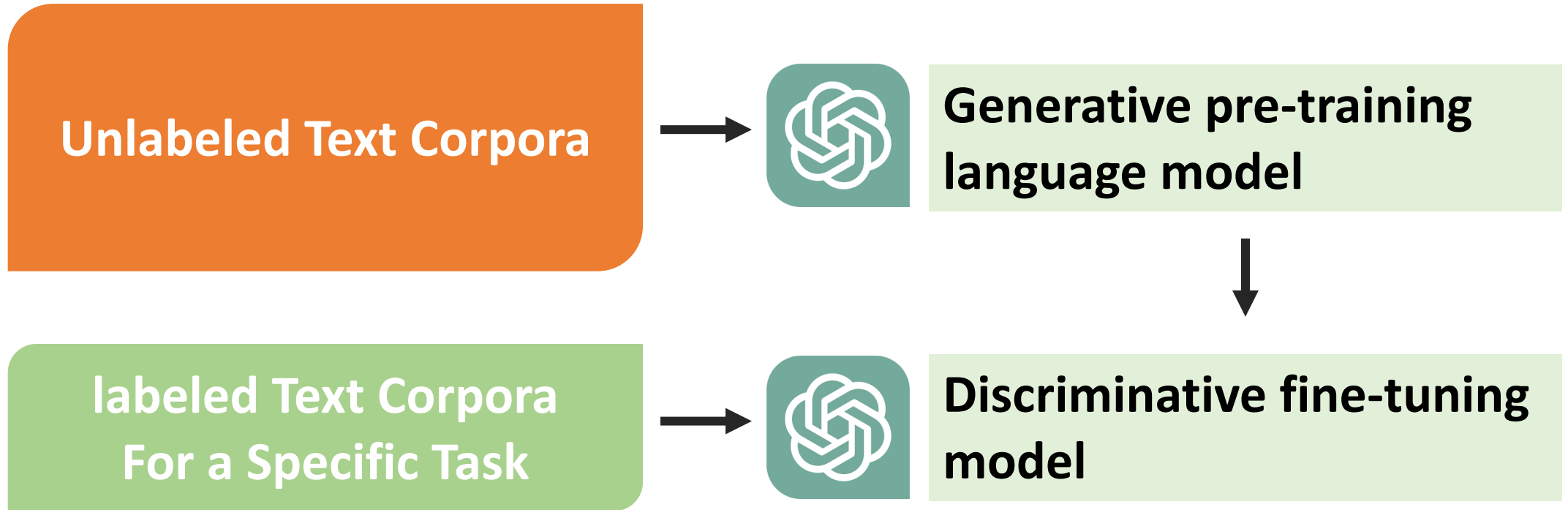
### Discriminative Learning

라벨값에 따른 데이터 간의 decision boundary를 정의  
데이터가 적어도 학습 가능  
샘플에 대한 과적합 확률이 높음

# Abstract

---

결국 GPT는..



***“Semi-Supervised Learning”***



# Introduction

---

## Why Semi-Supervised Learning?

➔ Leveraging more than word-level information from unlabeled text is challenging

- Transfer에 유용한 text 표현을 배우는 것에 어떤 형태의 최적화 목적(Optimization Objective)가 가장 좋은 지 불분명함 (정보부족, 목표의 불확실성, 부적절한 특성 등)
- 학습된 표현들을 목표 작업(target task)으로 transfer하는 가장 효과적인 방법이 불분명함.

이를 해결하고자 기존 연구에서는 Glove, discourse coherence 그리고 task-specific change를 적용하는 방식 등이 제안됨 -> 이번 논문에서는 semi-supervised를 그 해결책으로 제시.

# Introduction

---

## Why Semi-Supervised Learning?

➔ 최종적인 목표는 다양한 task에 작은 변화만으로 적용할 수 있는 representation을 학습하는 것.

### 과정과 의미

- K size context window 설정
- 특정한 i번째 단어를 예측하기 위해 i-1부터 i-k까지 단어를 보고, 그 가능성을 최대하는 방법  
(즉, 라벨이 없는 데이터에도 학습할 수 있도록 비지도 학습을 설계)
- 확률을 최대화하는 것이기 때문에 MLE(우도 최대화) 기법을 loss function으로 설정
- SGD를 이용해 역전파

# Introduction

---

## Why Transformer(decoder)?

이유

- RNN과 비교할 때, 멀리 떨어진 요소들 사이의 의존성을 학습하기에 좋음
- 전이 학습 시에 각 작업에 맞게 입력 데이터를 변환하는 방법을 사용하여 모델을 조정함.
- Encoder는 벡터로 출력되는 반면, Decoder는 확률값으로 표현되기에 예측이 가능.
- 답변, 응답, 해석, 번역 등의 분야에서 뛰어난 구조임.

정리하면, 비지도 학습을 통해 대용량의 text 데이터를 Decoder 에 학습시켜서 LM으로 사전 학습한 것이 바로 GPT.

# Framework

---

학습은 두 가지 단계(STAGE)로 진행됨.

1. 대량의 말뭉치로 대용량의 언어 모델을 사전 학습
2. Labeled Data를 사용하여 목표 작업에 맞게 미세조정

# Framework

## Unsupervised pre-training

$$h_0 = UW_e + W_p$$

$$h_l = \text{transformer\_block}(h_{l-1}) \forall i \in [1, n]$$

$$P(u) = \text{softmax}(h_n W_e^T)$$

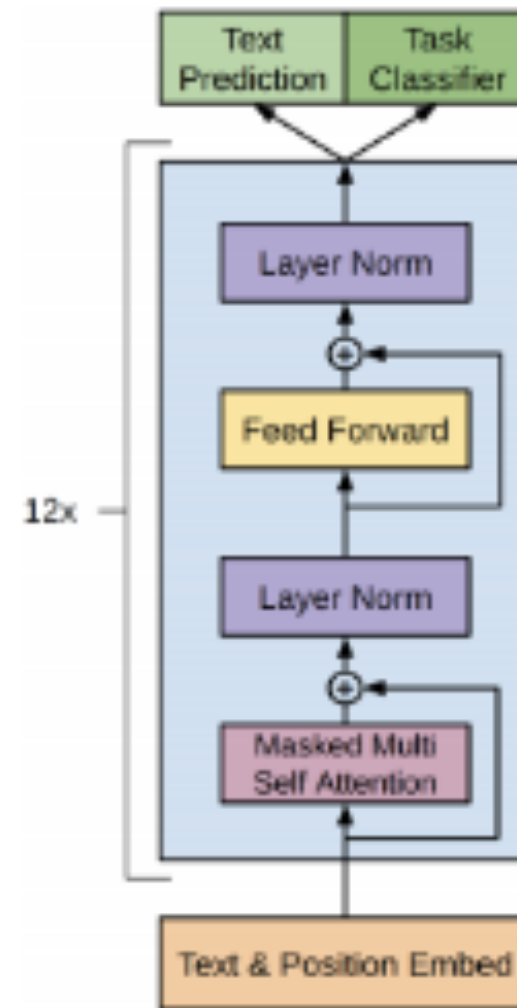
1. input sequence를 받아 word embedding, positional embedding 수행  
-> 이때,  $h_0$ 이 위와 같이 표현됨

2. 그 다음 hidden state를 decoder 블록에 계속 넣어 학습시킴

3. 마지막 최종 hidden state 값을 활용하여 확률값 출력

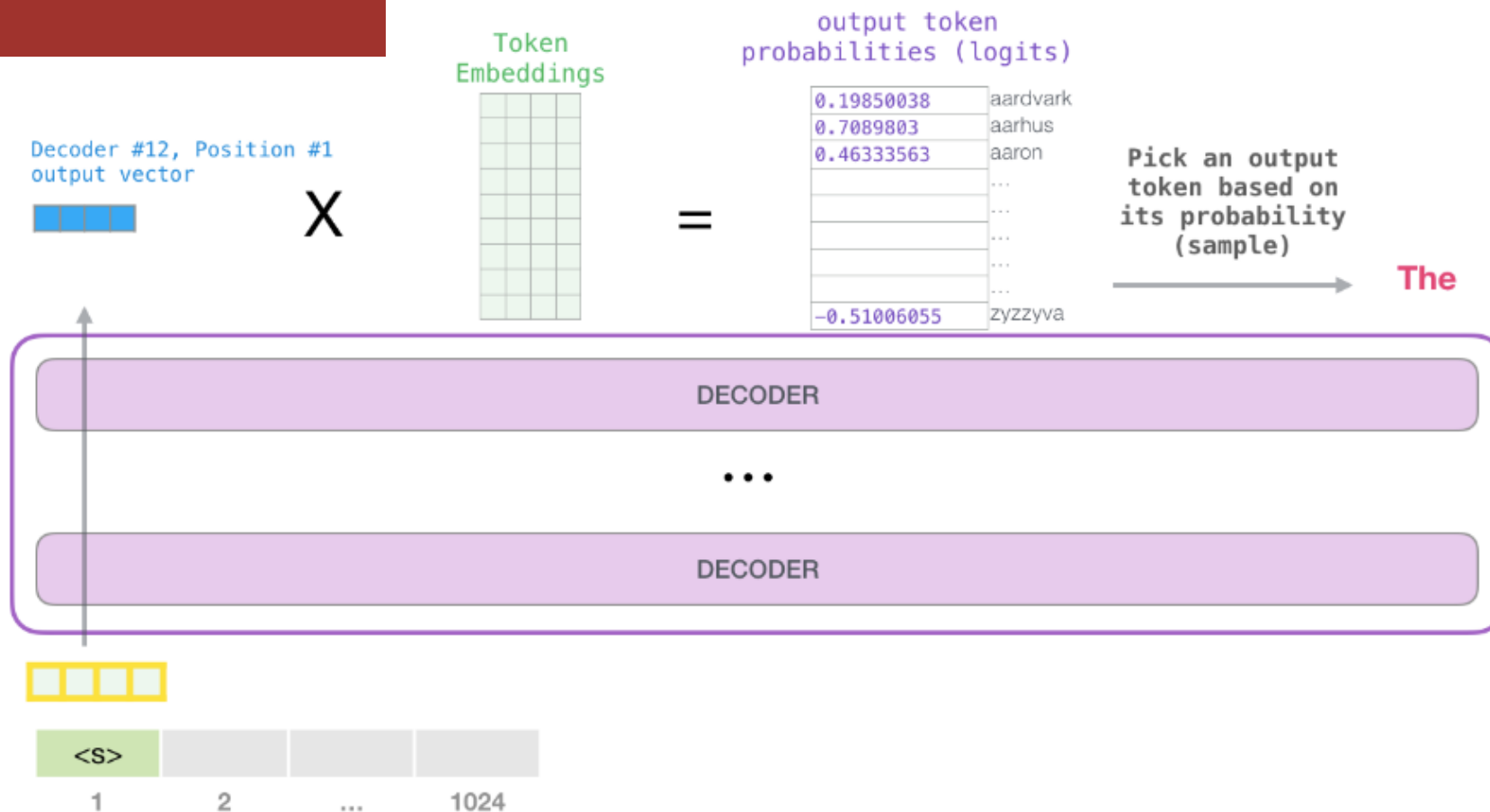
\*  $U$ 는 context vector,  $W_e$ 는 토큰 임베딩 matrix,  $W_p$ 는 positional matrix

\* No multihead attention, just masked multihead attention.



# Framework

## Unsupervised pre-training



# Framework

## Supervised fine-training

- 라벨 데이터를 활용하여, Task에 맞게 model을 fine tuning시키는 과정

the final transformer block's activation  $h_l^m$ , which is then fed into an added linear output layer with parameters  $W_y$  to predict  $y$ :

$$P(y|x^1, \dots, x^m) = \text{softmax}(h_l^m W_y).$$

$$L_2(\mathcal{C}) = \sum_{(x,y)} \log P(y|x^1, \dots, x^m).$$

- \* 데이터 셋에 label  $y$ 가 존재하기 때문에,  $y$ 에 대한 확률을 사용.
- \*  $\mathcal{C}$ 는 라벨링된 데이터 셋을 의미

# Framework

---

## Supervised fine-training

- 미세조정의 학습을 돕는 보조 목적 함수

$$L_3(\mathcal{C}) = L_2(\mathcal{C}) + \lambda * L_1(\mathcal{C})$$

\*L2 Function에 L1(C)를 추가적으로 사용할 때, 일반화와 학습속도 향상에 도움.



# Framework

## Task-specific input transformation

- 특정 task의 경우, 구조화된 입력을 제공해야 함.
- Task Specific하게 Tuning위해, Input transformation 진행
- linear+softmax의 layer만 추가하고 그 이전은 freeze.

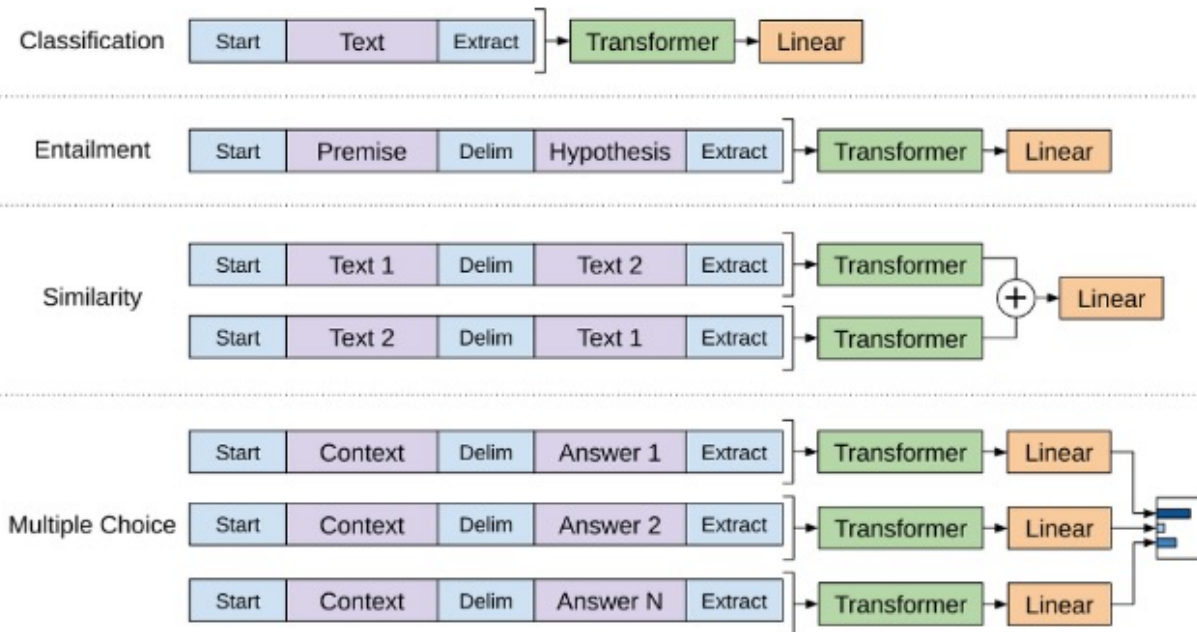
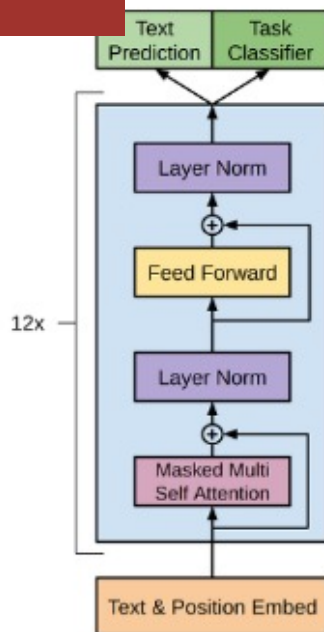


Figure 1: **(left)** Transformer architecture and training objectives used in this work. **(right)** Input transformations for fine-tuning on different tasks. We convert all structured inputs into token sequences to be processed by our pre-trained model, followed by a linear+softmax layer.

# Experiments

---

## Setup

- 7000개의 다양한 종류의 서적 데이터 – BookCorpus Dataset  
-> 상대적으로 긴 길이의 연속적 텍스트를 포함하여 long – range info를 습득
- 12개의 Transcormer decoder block + adam optimizer  
-> 전반적으로 Transformer의 구조를 따름.
- Bytepair Encoding(BPE)사용  
-> 데이터의 효율적 표현 가능

# Supervised fine-tuning

## Natural Language Inference(NLI)

-> 두 문장 간의 관계를 파악하는 Task

| Method                              | MNLI-m      | MNLI-mm     | SNLI        | SciTail     | QNLI        | RTE         |
|-------------------------------------|-------------|-------------|-------------|-------------|-------------|-------------|
| ESIM + ELMo [44] (5x)               | -           | -           | <u>89.3</u> | -           | -           | -           |
| CAFE [58] (5x)                      | 80.2        | 79.0        | <u>89.3</u> | -           | -           | -           |
| Stochastic Answer Network [35] (3x) | <u>80.6</u> | <u>80.1</u> | -           | -           | -           | -           |
| CAFE [58]                           | 78.7        | 77.9        | 88.5        | <u>83.3</u> |             |             |
| GenSen [64]                         | 71.4        | 71.3        | -           | -           | <u>82.3</u> | 59.2        |
| Multi-task BiLSTM + Attn [64]       | 72.2        | 72.1        | -           | -           | 82.1        | <b>61.7</b> |
| Finetuned Transformer LM (ours)     | <b>82.1</b> | <b>81.4</b> | <b>89.9</b> | <b>88.3</b> | <b>88.1</b> | 56.0        |

[표-2]

# Supervised fine-tuning

## Q n A and Commonsense reasoning

-> Multiple Choice Task

| Method                          | Story Cloze | RACE-m      | RACE-h      | RACE        |
|---------------------------------|-------------|-------------|-------------|-------------|
| val-LS-skip [55]                | 76.5        | -           | -           | -           |
| Hidden Coherence Model [7]      | <u>77.6</u> | -           | -           | -           |
| Dynamic Fusion Net [67] (9x)    | -           | 55.6        | 49.4        | 51.2        |
| BiAttention MRU [59] (9x)       | -           | <u>60.2</u> | <u>50.3</u> | <u>53.3</u> |
| Finetuned Transformer LM (ours) | <b>86.5</b> | <b>62.9</b> | <b>57.4</b> | <b>59.0</b> |

[H-3]

# Supervised fine-tuning

## Classification and Semantic similarity

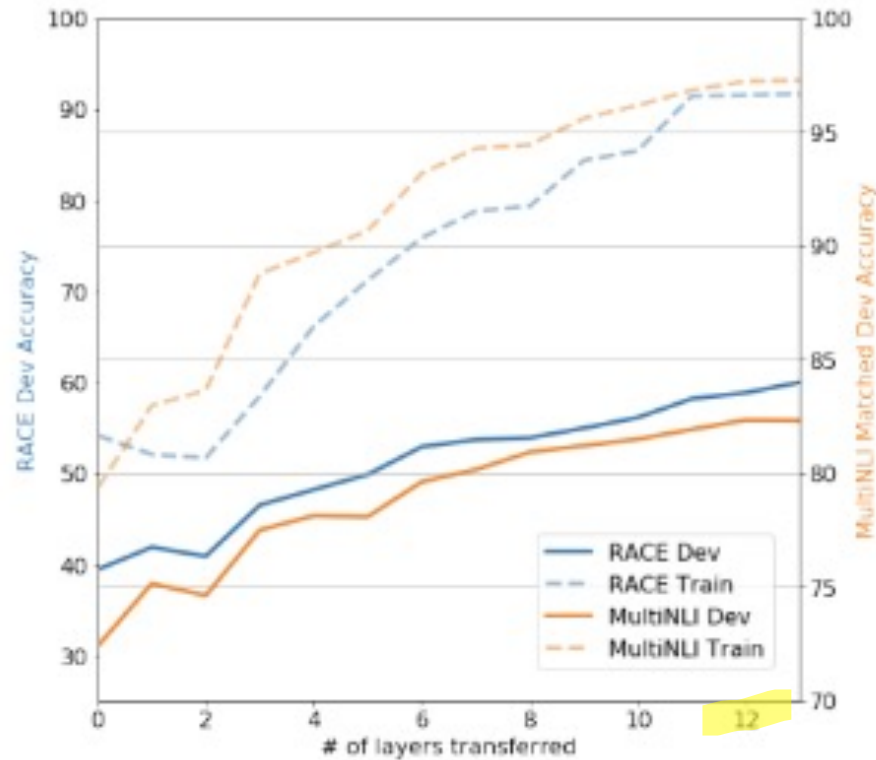
-> 분류와 문장간의 유사도 측정 task

| Method                                | Classification |               | Semantic Similarity |              |             | GLUE        |
|---------------------------------------|----------------|---------------|---------------------|--------------|-------------|-------------|
|                                       | CoLA<br>(mc)   | SST2<br>(acc) | MRPC<br>(F1)        | STSB<br>(pc) | QQP<br>(F1) |             |
| Sparse byte mLSTM [16]                | -              | <b>93.2</b>   | -                   | -            | -           | -           |
| TF-KLD [23]                           | -              | -             | <b>86.0</b>         | -            | -           | -           |
| ECNU (mixed ensemble) [60]            | -              | -             | -                   | <u>81.0</u>  | -           | -           |
| Single-task BiLSTM + ELMo + Attn [64] | <u>35.0</u>    | 90.2          | 80.2                | 55.5         | <u>66.1</u> | 64.8        |
| Multi-task BiLSTM + ELMo + Attn [64]  | 18.9           | 91.6          | 83.5                | 72.8         | 63.3        | <u>68.9</u> |
| Finetuned Transformer LM (ours)       | <b>45.4</b>    | 91.3          | 82.3                | <b>82.0</b>  | <b>70.3</b> | <b>72.8</b> |

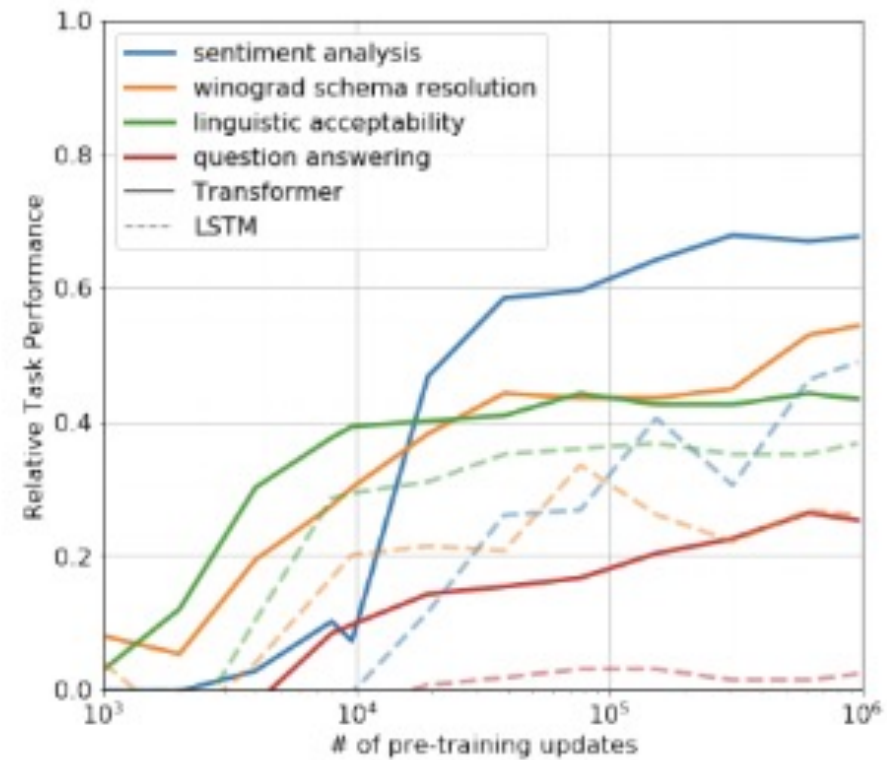
[표-4]

# Analysis

## Impact of number of layers & Zero shot Behaviors



\* 디코더 블록을 쌓을 수록 성능이 올라감



\* 사전훈련할수록 성능이 올라감

# Conclusion

---

- GPT는 semi-supervised learning을 사용하여 task specific한 LM을 구축한 것.
- 광대한 정보를 담은 연속된 텍스트로 이루어진 다양한 말뭉치들로 사전학습된 모델은 상당한 일반지식과 넓은 범위의 정보 처리 가능한 기능을 얻음
- 모델을 구성할 때, unlabeled data를 활용하여 unsupervised pre-training이 가능함을 보여줌.
- 다양한 벤치마크에서 SOTA를 보여줌.



TRAIN AND TEST