

# 지식 증류를 이용한 한국어 RoBERTa의 경량화

Knowledge distillation for lightweight RoBERTa of Korean

강동찬<sup>1</sup>, 나승훈<sup>1</sup>, 최윤수<sup>2</sup>, 이혜우<sup>2</sup>, 장두성<sup>2</sup>

<sup>1</sup>전북대학교, <sup>2</sup>KT

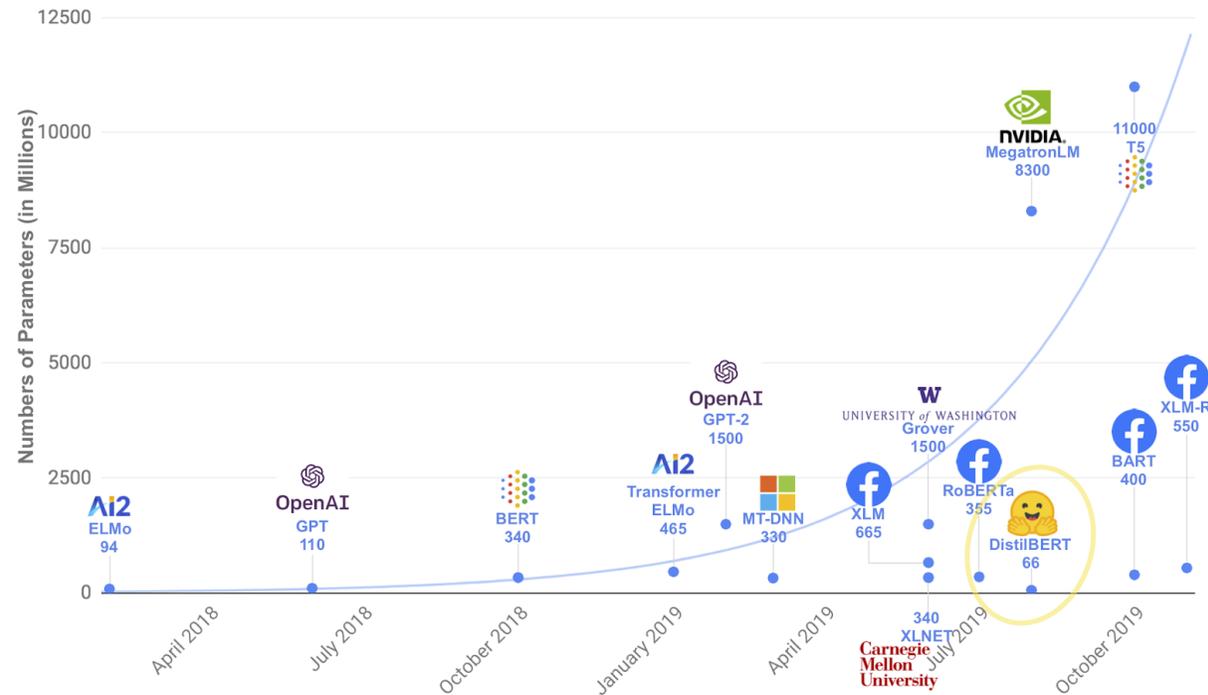
# Contents

1. Introduction: trends and problems
2. Related work: DistilBERT and TinyBERT
3. Experiment
4. Conclusion and future work

Introduction: trends and problems

# Trends

- 최근 NLP 분야에서 BERT와 같이 사전학습 된 언어모델들은 state-of-the-art 인코더
- 사전학습 된 언어모델의 크기는 점점 커지고 있는 추세



# Problems

## 1. Memory usage

- 단말기에 탑재하기 어려움
- 확장 되기 어려움

## 2. Inference speed

- 서비스 되기 어려운 속도

## 3. Environmental cost

- 연구에 사용되는 환경적인 비용

# Problems

## 1. Memory usage

- 단말기에 탑재하기 어려움
- 확장 되기 어려움

## 2. Inference speed

- 서비스 되기 어려운 속도

## 3. Environmental cost

- 연구에 사용되는 환경적인 비용

### Financial cost

연구/개발 단계에서 재정적인 비용과 직접적으로  
연관

# Problems

## 1. Memory usage

- 단말기에 탑재하기 어려움
- 확장 되기 어려움

## 2. Inference speed

- 서비스 되기 어려운 속도

## 3. Environmental cost

- 연구에 사용되는 환경적인 비용



<b>Consumption</b>	<b>CO<sub>2</sub>e (lbs)</b>
Air travel, 1 passenger, NY↔SF	1984
Human life, avg, 1 year	11,023
American life, avg, 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000

<b>Training one model (GPU)</b>	
NLP pipeline (parsing, SRL)	39
w/ tuning & experimentation	78,468
Transformer (big)	192
w/ neural architecture search	626,155

Table 1: Estimated CO<sub>2</sub> emissions from training common NLP models, compared to familiar consumption.<sup>1</sup>

Related work: DistilBERT and TinyBERT

# DistilBERT

- **지식 증류 (Knowledge Distillation)**

큰 교사 모델의 아웃풋 분포를 작은 모델의 학습에 이용하여, 작은 학생 모델의 학습을 도움

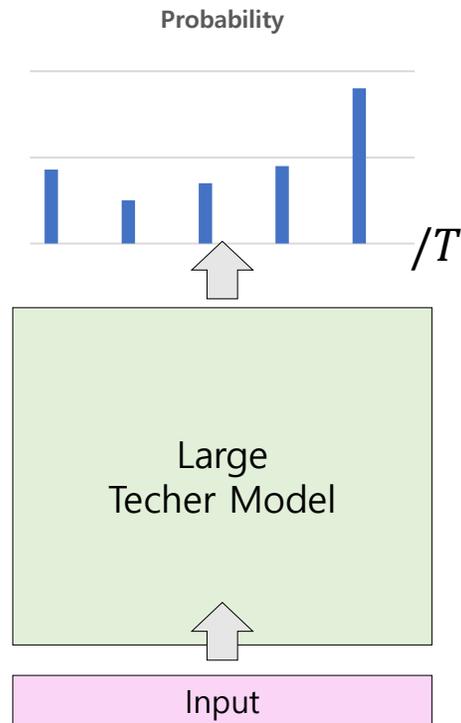
$$t^i = \frac{\exp(z_i^T / T)}{\sum_j \exp(z_j^T / T)}$$

$$s^i = \frac{\exp(z_i^S / T)}{\sum_j \exp(z_j^S / T)}$$

$$\mathcal{L}_{kd} = - \sum_i t^i * \log(s^i)$$

# DistilBERT

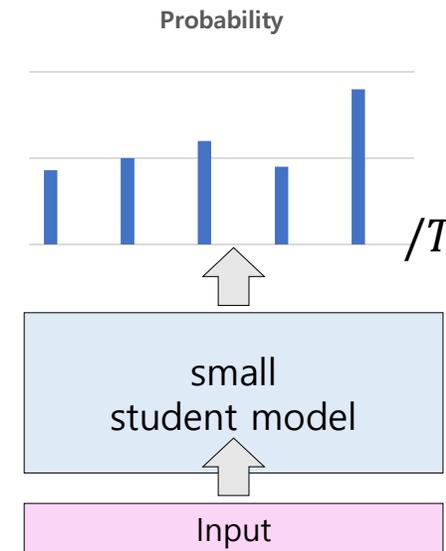
- 지식 증류 (Knowledge Distillation)



$$t^i = \frac{\exp(z_i^T / T)}{\sum_j \exp(z_j^T / T)}$$

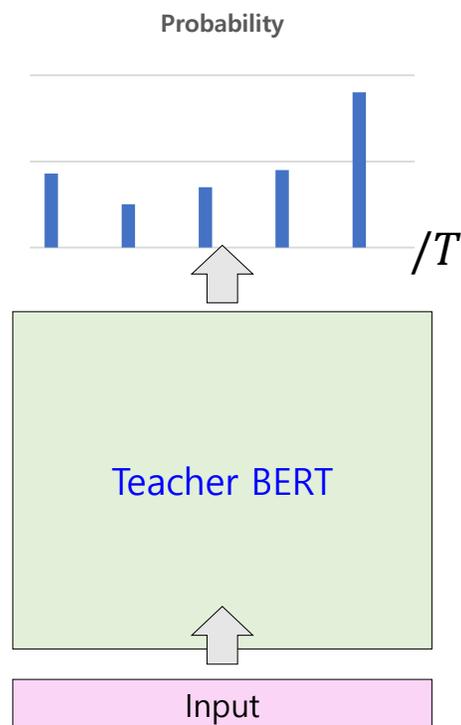
$$s^i = \frac{\exp(z_i^S / T)}{\sum_j \exp(z_j^S / T)}$$

$$\mathcal{L}_{kd} = - \sum_i t^i * \log(s^i)$$



# DistilBERT

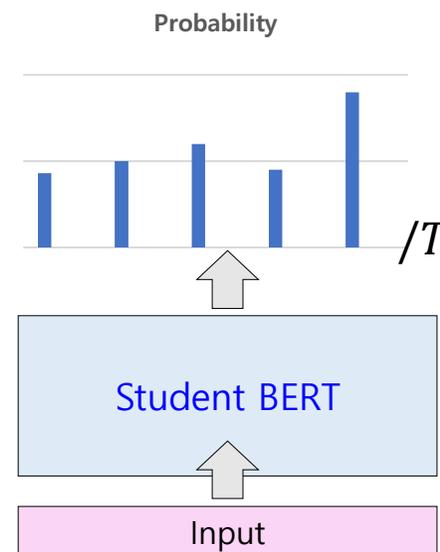
- Key idea



$$t^i = \frac{\exp(z_i^T / T)}{\sum_j \exp(z_j^T / T)}$$

$$s^i = \frac{\exp(z_i^S / T)}{\sum_j \exp(z_j^S / T)}$$

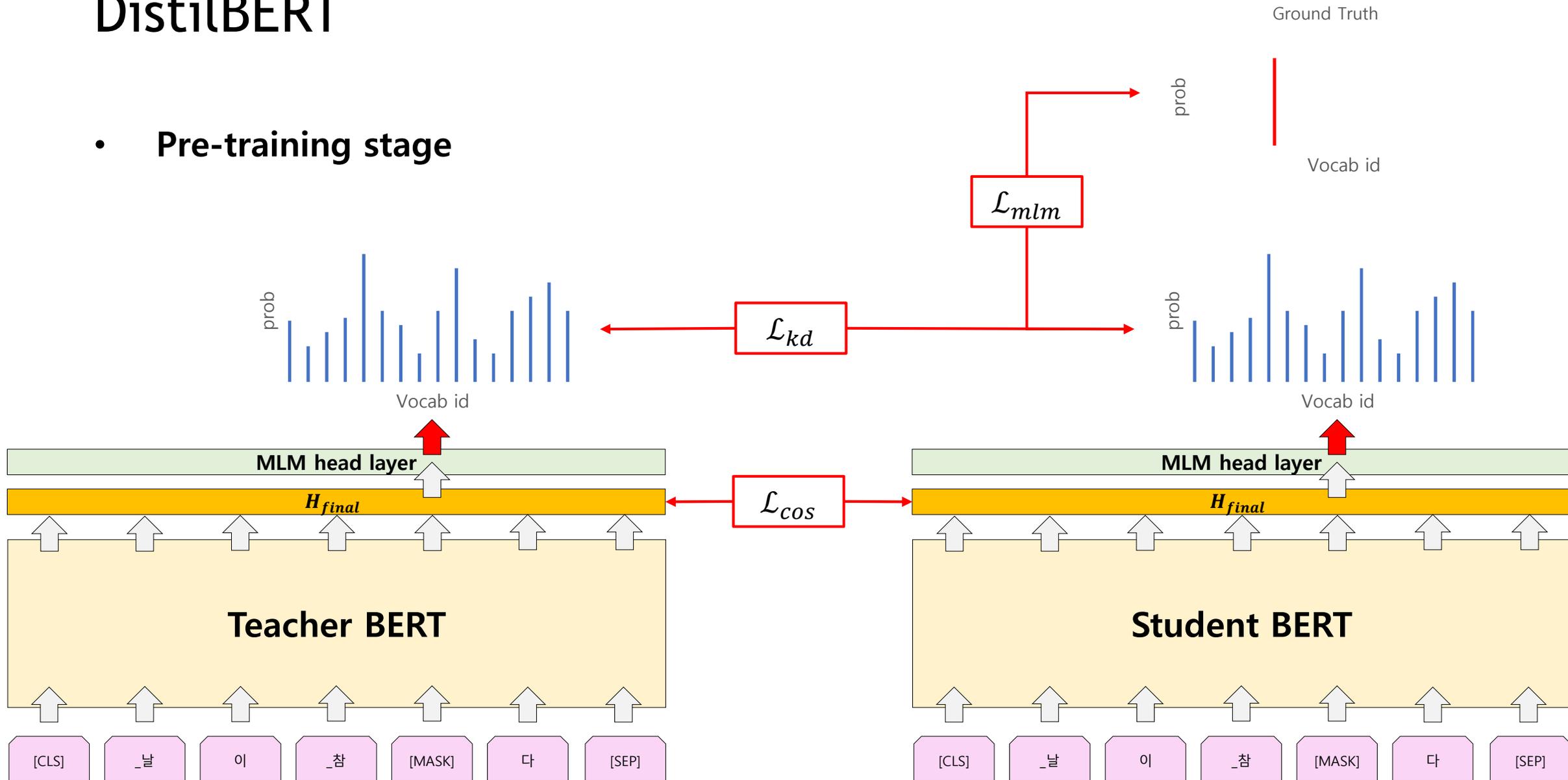
$$\mathcal{L}_{kd} = - \sum_i t^i * \log(s^i)$$



$$\mathcal{L} = \alpha_{mlm} \mathcal{L}_{mlm} + \alpha_{kd} \mathcal{L}_{kd} + \alpha_{cos} \mathcal{L}_{cos}$$

# DistilBERT

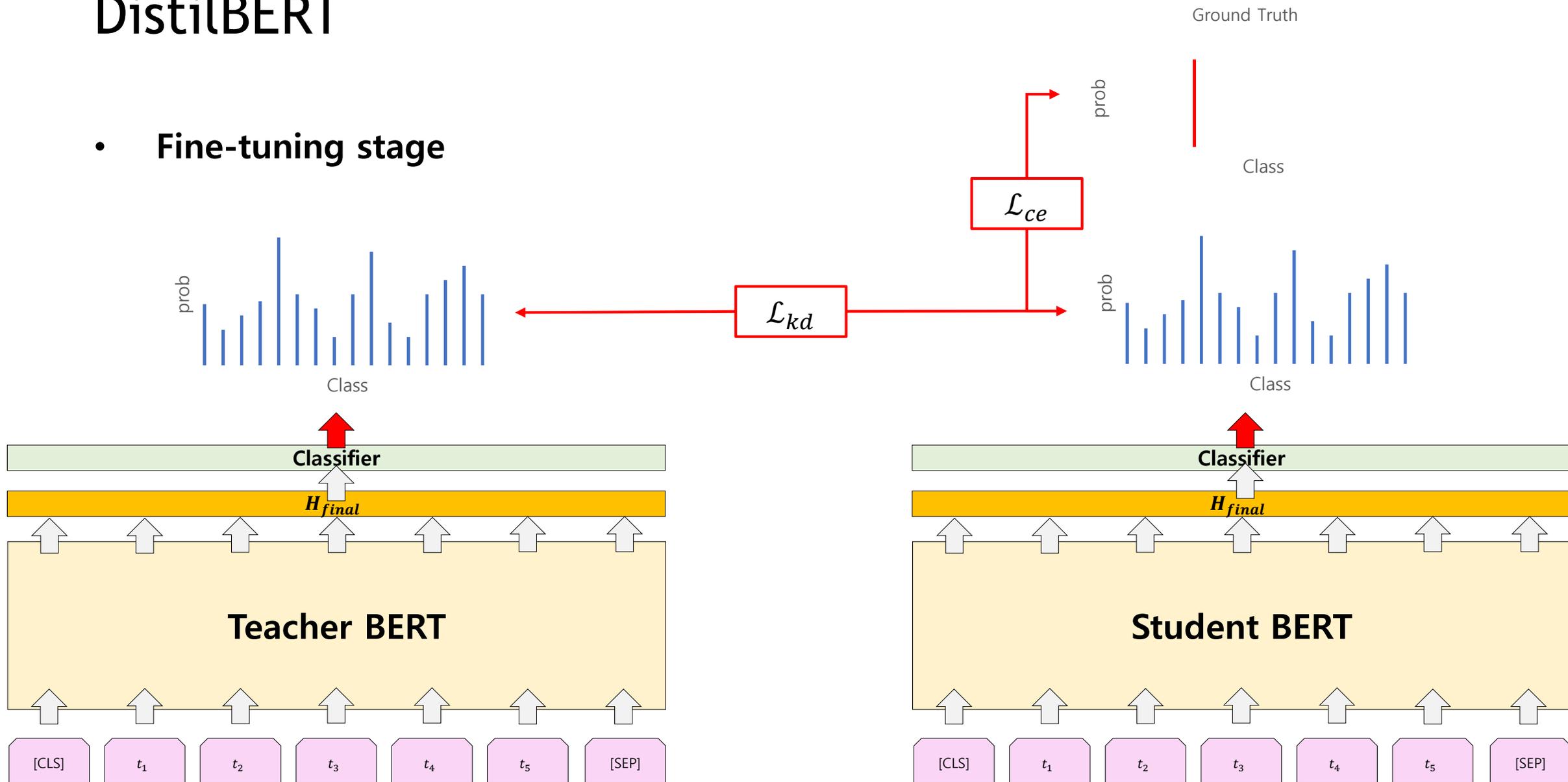
- Pre-training stage



$$\mathcal{L} = \alpha_{ce} \mathcal{L}_{ce} + \alpha_{kd} \mathcal{L}_{kd}$$

# DistilBERT

- Fine-tuning stage



# TinyBERT

- 트랜스포머 증류 (Transformer Distillation)

TinyBERT에서 제안한 교사 BERT모델의 인코더들의 아웃풋 표현을 이용한 증류법

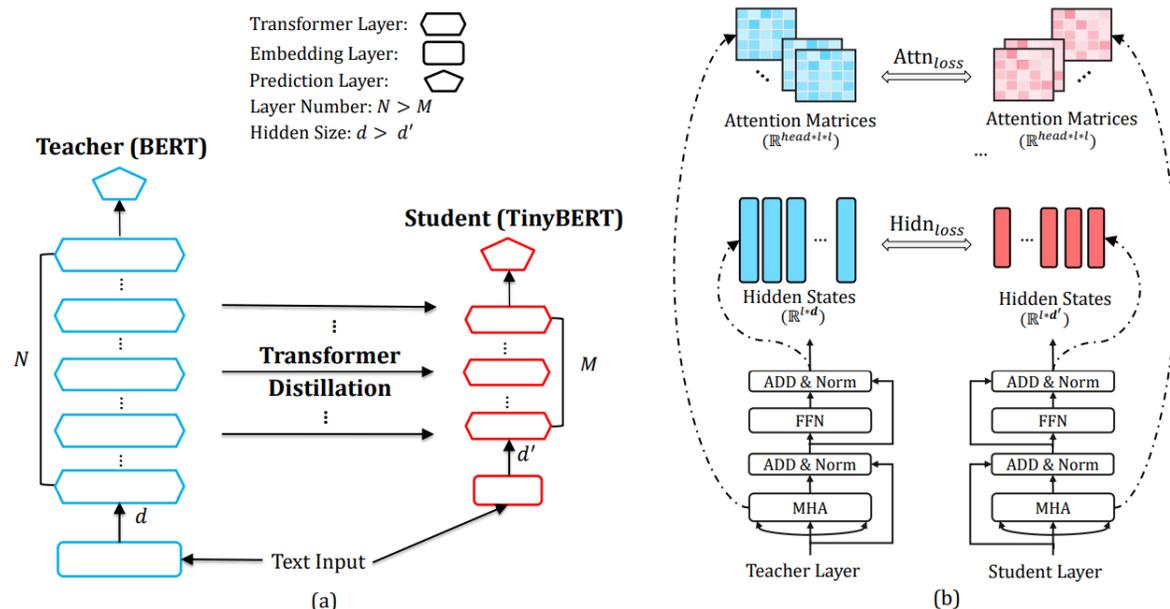
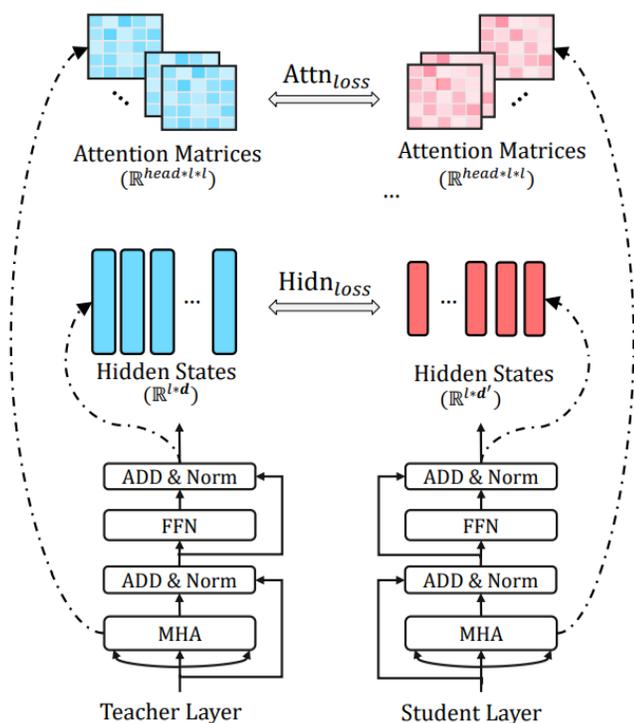


Figure 1: An overview of Transformer distillation: (a) the framework of Transformer distillation, (b) the details of Transformer-layer distillation consisting of  $Attn_{loss}$  (attention based distillation) and  $Hidn_{loss}$  (hidden states based distillation).

# TinyBERT

- Key idea



$$\mathcal{L}_{\text{embd}} = \text{MSE}(E^S W_e, E^T)$$

$$\mathcal{L}_{\text{hidn}} = \text{MSE}(H^S W_h, H^T)$$

$$\mathcal{L}_{\text{attn}} = \frac{1}{h} \sum_{i=1}^h \text{MSE}(A_i^S, A_i^T)$$

$$\mathcal{L}_{\text{pred}} = - \sum_i t^i * \log(s^i)$$

$$\mathcal{L}_{\text{layer}}(S_m, T_{g(m)}) = \begin{cases} \mathcal{L}_{\text{embd}}(S_0, T_0), & m = 0 \\ \mathcal{L}_{\text{hidn}}(S_m, T_{g(m)}) + \mathcal{L}_{\text{attn}}(S_m, T_{g(m)}), & M \geq m > 0 \\ \mathcal{L}_{\text{pred}}(S_M, T_{N+1}), & m = M + 1 \end{cases}$$

$$\mathcal{L} = \sum_{m=0}^{M+1} \lambda_m \mathcal{L}_{\text{layer}}(S_m, T_{g(m)})$$

# TinyBERT

- Pre-training stage, Fine-tuning stage

KD Methods	KD at Pre-training Stage				KD at Fine-tuning Stage					
	INIT	Embd	Attn	Hidn	Pred	Embd	Attn	Hidn	Pred	DA
Distilled BiLSTM <sub>SOFT</sub>									✓	✓
BERT-PKD	✓							✓ <sup>3</sup>	✓	
DistilBERT	✓				✓ <sup>4</sup>				✓	
TinyBERT (our method)		✓	✓	✓		✓	✓	✓	✓	✓

Our Work: Two stage Distillation in fine-tuning stage

# Motivation

- **Stage-wise training in Fitnets**

교사 모델로부터 표현 학습 후, KD 적용.

---

**Algorithm 1** FitNet Stage-Wise Training.

The algorithm receives as input the trained parameters  $\mathbf{W}_T$  of a teacher, the randomly initialized parameters  $\mathbf{W}_S$  of a FitNet, and two indices  $h$  and  $g$  corresponding to hint/guided layers, respectively. Let  $\mathbf{W}_{\text{Hint}}$  be the teacher's parameters up to the hint layer  $h$ . Let  $\mathbf{W}_{\text{Guided}}$  be the FitNet's parameters up to the guided layer  $g$ . Let  $\mathbf{W}_r$  be the regressor's parameters. The first stage consists in pre-training the student network up to the guided layer, based on the prediction error of the teacher's hint layer (line 4). The second stage is a KD training of the whole network (line 6).

---

**Input:**  $\mathbf{W}_S, \mathbf{W}_T, g, h$

**Output:**  $\mathbf{W}_S^*$

- 1:  $\mathbf{W}_{\text{Hint}} \leftarrow \{\mathbf{W}_T^1, \dots, \mathbf{W}_T^h\}$
  - 2:  $\mathbf{W}_{\text{Guided}} \leftarrow \{\mathbf{W}_S^1, \dots, \mathbf{W}_S^g\}$
  - 3: Initialize  $\mathbf{W}_r$  to small random values
  - 4:  $\mathbf{W}_{\text{Guided}}^* \leftarrow \underset{\mathbf{W}_{\text{Guided}}}{\operatorname{argmin}} \mathcal{L}_{HT}(\mathbf{W}_{\text{Guided}}, \mathbf{W}_r)$
  - 5:  $\{\mathbf{W}_S^1, \dots, \mathbf{W}_S^g\} \leftarrow \{\mathbf{W}_{\text{Guided}}^{*1}, \dots, \mathbf{W}_{\text{Guided}}^{*g}\}$
  - 6:  $\mathbf{W}_S^* \leftarrow \underset{\mathbf{W}_S}{\operatorname{argmin}} \mathcal{L}_{KD}(\mathbf{W}_S)$
-

# Our Work

- **Inspired by stage-wise training in Fitnets**

교사 BERT로부터 표현 학습( $\mathcal{L}_{trm}$ ) 후, KD 학습( $\mathcal{L}_{pred}$ ). 각 파인튜닝 데이터 셋의 Train set 의 10%를 Validation set 으로 사용

$$\mathcal{L}_{layer}(S_m, T_{g(m)}) = \begin{cases} \mathcal{L}_{embd}(S_0, T_0), & m = 0 \\ \mathcal{L}_{hidn}(S_m, T_{g(m)}) + \mathcal{L}_{attn}(S_m, T_{g(m)}), & M \geq m > 0 \end{cases}$$

$$\mathcal{L}_{trm} = \sum_{m=0}^M \lambda_m \mathcal{L}_{layer}(S_m, T_{g(m)})$$

$$\mathcal{L}_{pred} = f(Z^S, Z^T)$$

Experiment

# Experiment setup (pre-training stage)

- **Corpus**

한국어 위키피디아 + 뉴스 + 백과사전 = 15GB

- **Configuration**

	Hidden	Embd	Intermediate	#Attention heads	#Hidden layers	Vocab	#params
Teacher RoBERTa	768	768	3072	12	12	31331	110M
DistilRoBERTa	768	768	3072	12	6	31331	66M
TinyRoBERTa	312	312	1200	12	4	31331	14M

# Experiment setup (fine-tuning stage)

- **Dataset**

**NSMC, KorQuAD, ETRI NER**

- **KD function**

<b>NSMC</b>	$\mathcal{L}_{\text{pred}} = \text{MSE}(Z^{\mathcal{S}}, Z^{\mathcal{T}})$
<b>NER</b>	$\mathcal{L}_{\text{pred}} = \text{MSE}(Z^{\mathcal{S}}, Z^{\mathcal{T}})$
<b>KorQuAD</b>	$\mathcal{L}_{\text{pred}} = \{\text{MSE}(Z_{\text{start}}^{\mathcal{S}}, Z_{\text{start}}^{\mathcal{T}}) + \text{MSE}(Z_{\text{end}}^{\mathcal{S}}, Z_{\text{end}}^{\mathcal{T}})\}/2$

# Experiment result

- Result**

	corpus size	#params	NSMC (Acc)	KorQuAD 1.0 (EM / F1)	ETRI NER (F1)
BERT-Multilingual [8, 9]	-	110M	87.43	77.69 / 89.98	- / 91.92
BERT-ETRI [9]	23.5G	110M	-	84.72 / 92.74	-
BERT-형태소 태그 [8]	540M	110M	86.57	-	- / 91.58
RoBERTa-형태소 태그 [10]	15G	110M	89.88	-	- / 94.79
RoBERTa-자소 (Teacher) [11]	18G	110M	91.14	87.11 / 94.47	92.80 / 93.47
DistilRoBERTa	15G	66M	89.93	82.99 / 91.42	92.13 / 91.75
+ pred			90.02	83.97 / 92.47	92.22 / 92.09
+ trm + pred			89.89	83.13 / 91.66	90.73 / 92.31
TinyRoBERTa	15G	14M	87.20	72.01 / 83.22	86.88 / 85.70
+ pred			87.87	75.16 / 86.31	87.15 / 84.58
+ trm + pred			88.20	75.90 / 87.05	87.69 / 85.05

Conclusion and future work

# Conclusion and future work

- **Conclusion**

본 연구에서는 DistilBERT와 TinyBERT의 방법론을 사전학습 단계에 적용하여 기존 연구와 비교하고, 파인튜닝 단계에서 2단계 증류를 통해 효과를 확인하였다. 실험결과 파인튜닝 단계에서 예측 레이어의 증류는 성능의 향상에 도움이 되었던 한편, 트랜스포머 증류의 경우는 예측 레이어의 성능을 저하시키는 요인으로 작용하기도 하였다.

- **Future work**

추후연구에서는 더 다양한 경량화 방법을 시도하고, 교사 모델의 성능과 차이가 많이 나는 태스크에서 성능차이를 줄이고자 한다.

Thanks for your attention !