

Multi-CLS BERT:



An Efficient Alternative to Traditional Ensembling

Haw-Shiuan Chang*, Ruei-Yao Sun*, Kathryn Ricci*, and Andrew McCallum

Introduction

Background:

- Traditional ensembles of <u>multiple</u> BERT models boost performance on natural language understanding tasks over single models
- *However*, traditional ensembles are **expensive**
 - Computational cost, memory, space footprint

Research question:

Can we achieve the benefits of ensembling while minimizing the cost?

\rightarrow **Proposed method**:

Ensemble multiple CLS embeddings within a single BERT model

Main Result

				GLUE		SuperGLUE			
	Configuration \downarrow	Model Name ↓	Model Size \downarrow	100	1k	Full	100*	1k*	Full
-	BERT Base	Pretrained	109.5M	55.71	71.67	82.05	57.18	61.55	65.04
		MTL	109.5M	± 0.62 59.29	± 0.15 73.26 ± 0.13	$^{\pm 0.08}_{83.30\dagger}$	± 0.43 57.50	± 0.37 62.94	± 0.36 66.33 ± 0.33
		Ours (K=1)	111.3M	± 0.27 57.84	± 0.13 73.28	± 0.07 83.40	57.31	± 0.30 63.35	£ 0.33 66.29
				± 0.32	± 0.13	± 0.07	± 0.35	± 0.18	± 0.18
		Ours (K=5, $\lambda = 0$)	118.4M	61.54	74.14	83.41	58.29	63.71	66.80
		$O_{\rm resc}$ (V 5) 0.1)		± 0.32	± 0.12	± 0.07	± 0.33	± 0.26	± 0.25
		Ours (K=5, $\lambda = 0.1$)	118.4M	01.8U	/4.10	ð3.4 7	58.20	03.01	00./4
		Ours (K=5 $\lambda = 0.5$)	118 4M	± 0.33 60 49	± 0.13 74 02	± 0.03 83.47	58.41	± 0.27 63.78	£ 0.20
		Ours(R=3, N=0.0)	110.4141	± 0.35	± 0.12	± 0.08	± 0.38	± 0.25	± 0.24
		Ours (K=5, $\lambda = 1$)	118.4M	59.86	73.75	83.43	57.84	63.56	66.39
				± 0.34	± 0.14	± 0.07	± 0.40	± 0.22	± 0.22
-		MTL	335.2M	61.39	75.30	84.13	59.03	65.21	69.16
				± 0.37	± 0.27	± 0.11	± 0.54	± 0.38	± 0.37
	BERT Large	Ours $(K=1)$	338.3M	59.19	75.35	84.59	57.35	64.67	69.24
		$O (\mathbf{U} \mathbf{z})$	250 014	± 0.43	± 0.21	± 0.07	± 0.42	± 0.43	± 0.41
		Ours (K=5, $\lambda = 0$)	350.9M	63.19	15.13	84.51	59.46	65.43	69.56
		Ours $(K-5) = 0.1$	250 OM	± 0.49	± 0.26 76 77	± 0.05 94 61		± 0.38	± 0.31
		Ours (K=3, $\lambda = 0.1$)	550.9M	+ 0.40	+ 0.12	04.01 + 0.08	+ 0.43	03.30 ± 0.26	/ U.U + 0.25
		Ours (K=5 $\lambda = 0.5$)	350 9M	6302	75 95	84 49	59 42	65 84	69 79
		Outs (IX - 5, 7 - 0.0)	550.9141	± 0.42	± 0.10	± 0.08	± 0.34	± 0.25	± 0.25
		Ours (K=5, $\lambda = 1$)	350.9M	62.07	75.85	84.61	58.74	65.00	69.04
				± 0.45	± 0.17	± 0.07	± 0.50	± 0.29	± 0.27



Pretraining

- Learn multiple, diversified CLS embeddings:
 - → Adapt state-of-the-art pretraining objectives for BERT (Aroca-Ouellette and Rudzicz, 2020)
- Learn to represent fine-grained semantics:
 - \rightarrow Incorporate hard negatives into the objective

 Table 1: Results on GLUE and SuperGLUE for models derived from BERT Base and BERT Large.

Conclusion 1:

Efficient Multi-CLS BERT improves performance over baseline single BERT model and *K*=1 model with only small increase in model size.

<u>Conclusion 2</u>: Multi-CLS BERT is especially effective in the few-shot setting.



<u>Claim 1</u>: Increased performance is due to ensemble effects.

- Ensembling Multi-CLS BERT only slightly boosts the performance (Table 2)
- Expected calibration error (Table 4)
- Overlap of most uncertain dataset examples (Table 3)
- Qualitative analysis of nearest-neighbor embeddings (Paper appendix): Multiple CLS embeddings can learn to contribute in complementary ways to solving a task

<u>Claim 2</u>: Using Multi-CLS BERT vs. traditional ensemble reduces computational costs at the expense of a modest drop in performance.

• Model size (Table 1)



Architecture

- Instead of traditional single BERT CLS embedding, leverage a fixed number of multiple CLS embeddings:
 - \rightarrow Insert *K* special CLS tokens
 - \rightarrow Use the K final-layer CLS hidden states to represent the input text

• Inference time (Table 4)

			GLUE		SuperGLUE*	
$Model \downarrow$	Model Description \downarrow	$K\downarrow$	100	1k	100	1k
	Pretrained	1	56.85	71.68	57.90	62.14
Baselines	MLM only	1	55.38	70.74	57.39	61.77
(BERT	CMTL+	1	58.65	72.57	56.88	62.63
Base)	MLM + SO + TFIDF	1	60.35	72.65	57.88	62.60
	MTL	1	59.53	73.12	57.51	62.95
	No Inserted	1	58.06	73.18	57.97	63.34
	Layers	5	60.12	73.35	56.46	62.00
	No Hard	1	58.44	73.30	57.19	63.33
	Negative	5	61.77	74.18	58.89	63.86
Ours	Sum Aggregation	5	58.87	73.94	57.41	63.82
(BERT	Default	1	57.76	73.30	57.53	63.22
Base)		3	61.09	73.95	57.85	63.31
	Default	5	62.62	74.49	58.82	63.86
		10	60.99	73.59	58.25	62.82
	SWA	1	57.31	72.91	-	-
	Ensemble on Dropouts	1	58.45	72.86	-	-
		1	60.07	75.20	-	-
	Ensemble on FT Seeds	5	63.34	75.35	-	-
0	No Hard	1	60.36	75.69	58.47	65.04
Ours	Negative	5	63.23	75.77	60.33	65.75
(BEKI	Default	1	60.01	76.03	57.38	65.10
Large)		5	64.33	76.38	59.99	65.51

Table 2: Ablation studies for the few-shot setting. Conclusion: Architecture and pretraining methods improve GLUE scores (and SuperGLUE almost as consistently). Traditional ensembles can achieve somewhat higher performance.



	GLUE* 100	GLUE* 1k
Multi-CLS vs ENS	32.57	41.35
Dropout vs ENS	37.17	45.53
Least vs ENS	39.57	48.85
ENS vs ENS	38.67	50.14

Table 3: Overlap ratio of the 20% most uncertain dataset examples as predicted by the two given models ("ENS" = traditional ensemble). Conclusion: "Multi-CLS vs. ENS" overlap ratio approaches "ENS vs. ENS".

	Inference	GLUE* (ECE	
	Time (s)	100	1k
Ours (K=1)	0.2918	25.22	19.32
Ours (K=5, $\lambda = 0.1$)	${ \overset{\pm 0.0002}{\textbf{0.3119}} \atop { \pm 0.0004 } }$	± 1.99 15.46 ± 1.79	${ \pm 1.64 \atop \pm 1.64 \atop \pm 1.64 }$
Ensemble of Ours (K=1)	$\begin{array}{c}1.4590\\\pm0.0012\end{array}$	$\begin{array}{c}13.85\\\pm0.97\end{array}$	$\begin{array}{c} 10.80 \\ \scriptstyle \pm \ 0.88 \end{array}$

Table 4: Inference time vs. expected calibration error (ECE).

Conclusion: Multi-CLS BERT with K=5 achieves significantly lower ECE than the same model with K=1, but with very little increase in inference time.

- \rightarrow Aggregate the K CLS embeddings during fine-tuning/inference
- Prevent collapse of the multiple CLS embeddings:
 - → Insert linear layers in between selected BERT layers at multi-CLS input positions
 - \rightarrow Add novel parameterization



- Using multiple CLS embeddings in a single BERT model with our architecture and pretraining methodologies results in almost free performance gain
- Evidence suggests that performance gains are due to ensemble effects without the cost of a traditional ensemble of multiple models
- Our methods for successfully implementing diversified multiple CLS embeddings may be extensible in future studies of efficient ensembling using other types of architectures

Related Work and References



[1] Aroca-Ouellette, S. and Rudzicz, F., On Losses for Modern Language Models. EMNLP 2020
[2] H.-S. Chang, "Modeling the Multi-mode Distribution in Self-Supervised Language Models, "PhD Thesis 2022

[3] H.-S. Chang, and A. McCallum, "Softmax Bottleneck Makes Language Models Unable to Represent Multi-mode Word Distributions," ACL 2022

[4] H.-S. Chang*, Z. Yao*, A. Gon, H. Yu, and A. McCallum, "Revisiting the Architectures like
Pointer Networks to Efficiently Improve the
Next Word Distribution, Summarization
Factuality, and Beyond" ACL Findings 2023

Figure 2: Architecture.