

Decoupling Knowledge from Memorization: Retrieval-augmented Prompt Learning



Xiang Chen ^{1,2}, Lei Li ^{1,2}, Ningyu Zhang ^{1,2 *}, Xiaozhuan Liang ^{1,2}, Shumin Deng ^{1,2}

Chuanqi Tan ³, Fei Huang ³, Luo Si ³, Huajun Chen ^{1,2 *}

¹ Zhejiang University & AZFT Joint Lab for Knowledge Engine,

² Hangzhou Innovation Center, Zhejiang University, ³ Alibaba Group



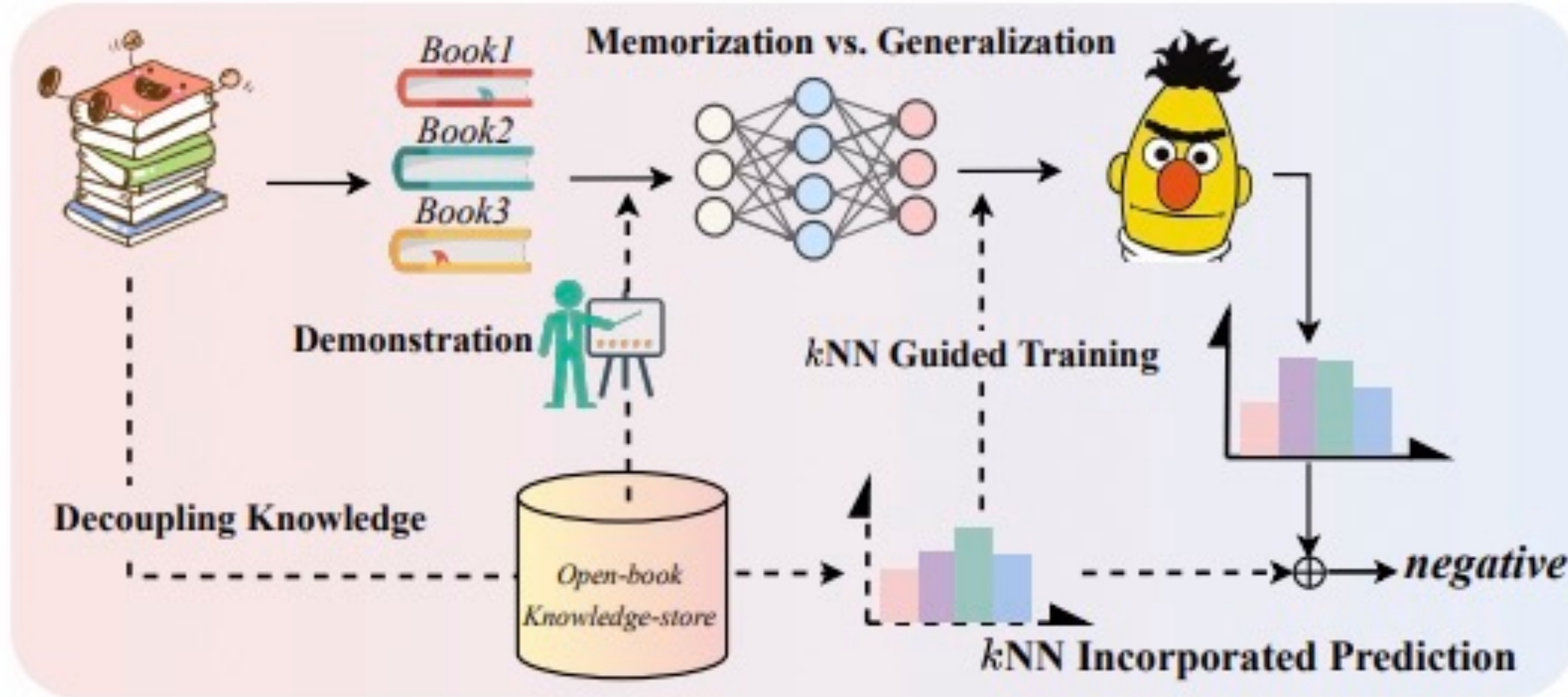
Introduction

• Limitations of Prompt Learning:

different Prompt learning with PLMs usually generalizes unstably in an extremely low-resource setting or emerging domains. One potential reason is that, it is non-trivial for parametric models to learn rare or hard patterns well with rote memorization, thus, resulting in inefficient generalizable performance.

• Decoupling knowledge from memorization :

with the motivation of decoupling knowledge from memorization to help the model strike a balance between generalization and memorization, we constructs an open-book knowledge-store from training instances and implements a retrieval mechanism during the process of input, training and inference, thus equipping the model with the ability to retrieve related contexts from the training corpus as cues for enhancement.



Retrieval-augmented Prompt Learning

➤ Open-book Knowledge-store

Given the i -th example (c_i, y_i) in the training data \mathcal{C} , we obtain the key-value pair (h_{c_i}, v_i) , in which $\hat{c}_i = \mathcal{T}(c_i)$, $h_{c_i} \in \mathbb{R}^d$ is the embedding of the [MASK] token in the last layer of the PLM, and $v_i = f(y_i)$ denotes the label word of the i -th example.

$$(\mathcal{K}, \mathcal{V}) = \{(h_{c_i}, v_i) \mid (c_i, y_i) \in \mathcal{C}\}$$

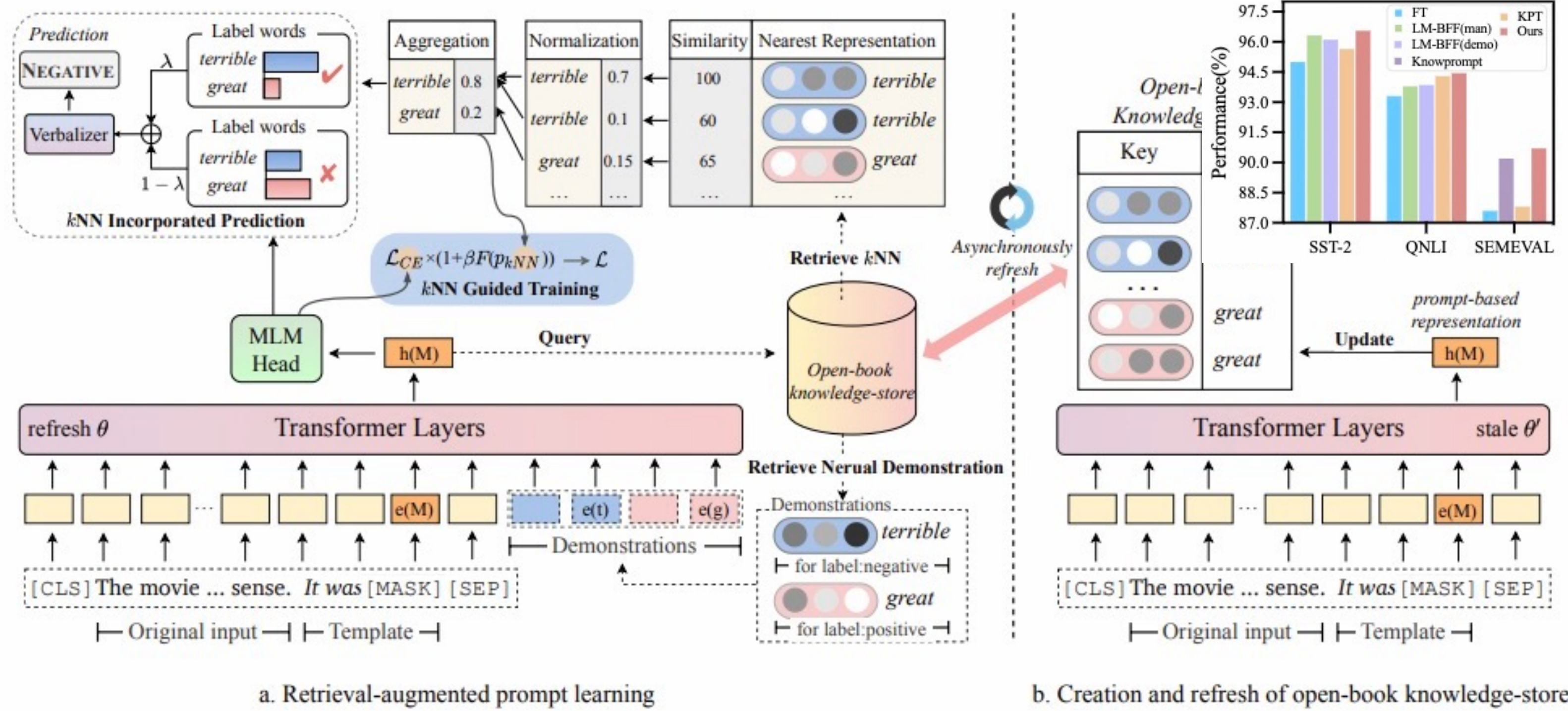


Figure 2: Overview of RETROPROMPT. Note that $e(\cdot)$ denotes word embedding function in the PLM \mathcal{M} , while “M”, “t” and “g” in $e(\cdot)$ specifically refers to “[MASK]”, “terrible” and “great”.

➤ Retrieval of Neural Demonstration

We intuitively aggregate the m neighbor vectors for each class according to their similarity and incorporate the demonstration into the input representation of \hat{x} after the word embedding layer of the \mathcal{M} as follows:

$$\mathcal{I} = e(\hat{x}) \oplus \left[\sum_{i \in [1:m]} \alpha_i^{(1)} h_{c_i}^{(1)}, e(v^{(1)}) \right] \oplus \dots \oplus \left[\sum_{i \in [1:m]} \alpha_i^{(L)} h_{c_i}^{(L)}, e(v^{(L)}) \right]; \alpha_i^{(l)} = \frac{e^{h_{c_i}^{(l)}}}{\sum_{i \in [1:m]} e^{h_{c_i}^{(l)}}}$$

➤ Retrieve kNN for Guiding Training

Our intuition is to differentiate between easy and hard examples according to the prediction of kNN.

$$P_{kNN}(y \mid q_t) \propto \sum_{(c_i, y_i) \in \mathcal{N}} \mathbf{1}_{y=y_i} \exp(d(h_{q_t}, h_{c_i}))$$

$$F(p_{kNN}) = -\log(p_{kNN}), \quad \mathcal{L} = (1 + \beta F(p_{kNN})) \mathcal{L}_{CE}$$

➤ kNN based probability for Cloze-style Prediction

we reformulate the $P(y \mid q_t)$ by interpolating the P_{kNN} with the already-trained base PLM's MLM prediction $P_{\mathcal{M}}$ using parameter λ to produce the final probability of the label:

$$P(y \mid q_t) = \lambda P_{kNN}(y \mid q_t) + (1 - \lambda) g(P_{\mathcal{M}}([MASK] = v \mid \mathcal{T}(q_t)))$$

Experiments

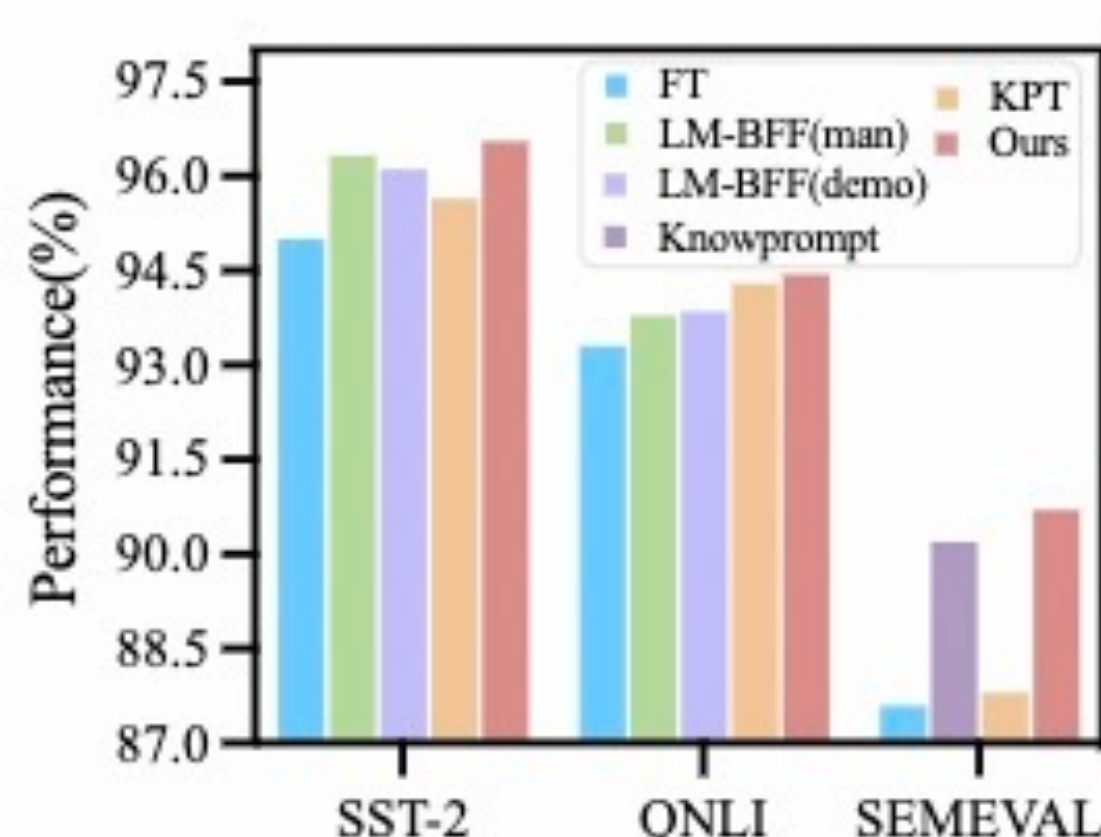
Few-shot/Zero-shot Results

St.	Model	Single Sentence					Model	Information Extraction			
		SST-2 (acc)	MR (acc)	CR (acc)	MNLI (acc)	QNLI (acc)		FewN (acc)	SemEval (acc)	TACRED (F1)	Avg.
16	FT	81.4 (3.8)	76.9 (5.9)	75.8 (3.2)	45.8 (6.4)	60.2 (6.5)	FT	52.7 (2.2)	66.1 (1.2)	25.8 (2.8)	60.6
	LM-BFF (man)	91.6 (1.2)	87.0 (2.0)	90.3 (1.6)	64.3 (2.5)	64.6 (5.4)	KnPr	65.3 (1.1)	80.9 (2.5)	33.2 (2.0)	71.4
	LM-BFF (D-demo)	91.8 (1.2)	86.6 (1.8)	90.2 (1.4)	64.8 (2.3)	69.2 (5.4)	KnPr (D-demo)	65.9 (1.5)	78.8 (2.1)	32.8 (1.7)	72.2*
	KPT †	90.3 (1.6)	86.8 (1.8)	88.8 (3.7)	61.4 (2.1)	61.5 (2.8)	KPT †	65.9 (1.5)	78.8 (2.1)	32.8 (1.7)	70.9
	Ours	93.9 (0.4)	88.0 (0.8)	91.9 (0.7)	71.1 (1.8)	71.6 (1.8)	Ours	67.3 (0.9)	81.5 (1.3)	40.7 (0.7)	75.6
4	FT	60.2 (2.8)	57.6 (1.4)	66.4 (5.5)	35.0 (0.3)	54.2 (3.9)	FT	32.7 (2.9)	38.8 (2.0)	14.7 (2.8)	45.8
	LM-BFF (man)	90.7 (0.8)	85.2 (2.8)	89.9 (1.8)	51.0 (2.5)	61.1 (6.1)	KnPr	52.5 (1.5)	58.4 (3.7)	28.8 (2.5)	62.8
	LM-BFF (D-demo)	90.2 (1.5)	85.5 (2.1)	89.7 (0.6)	56.1 (1.0)	61.7 (7.6)	KnPr (D-demo)	58.8 (2.2)	57.2 (3.2)	27.5 (2.2)	65.1*
	KPT †	88.2 (5.7)	83.4 (1.5)	87.2 (2.5)	53.7 (2.7)	59.2 (2.8)	KPT †	58.8 (2.2)	57.2 (3.2)	27.5 (2.2)	63.3
	Ours	91.5 (1.8)	87.4 (0.5)	91.4 (0.6)	57.6 (5.5)	62.2 (6.0)	Ours	60.9 (1.9)	59.2 (3.0)	32.1 (2.0)	67.6
0	LOTClass*	71.8	81.7	50.1	50.4	36.5	LOTClass*	11.5	9.8	2.5	41.1
	FT	49.1	50.0	49.8	34.4	49.5	FT	10.0	6.2	0.5	31.2
	LM-BFF (man)	83.5	80.3	78.4	49.7	50.5	KnPr	15.9	10.3	2.3	46.7
	LM-BFF (D-demo)	82.9	80.7	81.4	52.2	53.5	KnPr (D-demo)	—	—	—	47.0*
	KPT †	78.4	81.9	71.4	37.1	55.3	KPT †	24.6	11.6	0.8	45.7
	Ours	86.8	83.5	79.7	53.7	56.2	Ours	41.3	12.2	2.8	52.5

Cross-domain Results

Model	Source	Target Domain
	16-shot MR	SST-2 CR
FT	76.9	71.4
LM-BFF (man)	87.0	88.9
LM-BFF (D-demo)	86.6	89.3
KPT	86.8	86.8
RETROPROMPT	88.0	91.4
	16-shot QQP	MRPC RTE
FT	60.7	43.7
LM-BFF (man)	65.4	20.9
LM-BFF (D-demo)	68.2	38.8
KPT	71.6	42.3
RETROPROMPT	74.0	49.4

Full-data Results



Analysis of Memorization

➤ Definition of Memorization Measurement

we define memorization measures as to how the classification varies when a training instance z is deleted from the trainset. We define and derive the memorization score for a training instance z as follows :

$$S_{\text{delete}}(z) \stackrel{\text{def}}{=} - \frac{dP(y \mid x; \hat{\theta}_{\xi, -z})}{d\xi} \bigg|_{\xi=0} = - \nabla_{\theta} P(y \mid x; \hat{\theta})^{\top} \frac{d\hat{\theta}_{\xi, -z}}{d\xi} \bigg|_{\xi=0} = - \nabla_{\theta} P(y \mid x; \hat{\theta})^{\top} H_{\hat{\theta}}^{-1} \nabla_{\theta} \mathcal{L}(z, \hat{\theta})$$

➤ Top-memorized Instances: Typical or Atypical?

we adopt SST-2 to analyze the memorization by judging the atypical of an instance by checking the percentage of positive phrases.

Mem Group	Negative			Positive		
	FT	LM-BFF	OURS	FT	LM-BFF	OURS
Top-10%	34.29	32.78	30.23	68.75	69.71	75.67
ALL	23.40	23.40	23.40	86.39	86.39	86.39
Bottom-10%	17.63	16.25	14.42	95.92	95.08	94.53
	FT	LM-BFF	OURS			
MEM SCORE	4.597	0.121	0.032			