Knowledgeable Prompt-tuning: Incorporating Knowledge into Prompt Verbalizer for Text Classification

Anonymous ACL submission

Abstract

Tuning pre-trained language models (PLMs) with task-specific prompts has been a promising approach for text classification. Particularly, previous studies suggest that prompt-tuning has remarkable superiority in the low-data scenario over the generic fine-tuning methods with extra classifiers. The core idea of prompt-tuning is to insert text pieces, i.e., template, to the input and transform a classification problem into a masked language modeling problem, where a crucial step is to construct a projection, i.e., verbalizer, between a label space and a label word space. A verbalizer is usually handcrafted or searched by gradient descent, which may lack coverage and bring considerable bias and high variances to the results. In this work, we focus on incorporating external knowledge into the verbalizer, forming a knowledgeable prompttuning (KPT), to improve and stabilize prompttuning. Specifically, we expand the label word space of the verbalizer using external knowledge bases (KBs) and refine the expanded label word space with the PLM itself before predicting with the expanded label word space. Extensive experiments on zero and few-shot text classification tasks demonstrate the effectiveness of knowledgeable prompt-tuning.

1 Introduction

001

004

006

011

012

014

017

023

027

028

034

040

Recent years have witnessed the prominence of Pretrained Language Models (PLMs) (Peters et al., 2018; Radford et al., 2018; Devlin et al., 2019; Raffel et al., 2020; Xu et al., 2021) due to their superior performance on a wide range of languagerelated downstream tasks such as text classification (Kowsari et al., 2019), question answering (Rajpurkar et al., 2016), and machine reading comprehension (Nguyen et al., 2016). To fathom the principles of such effectiveness of PLMs, researchers have conducted extensive studies and suggested that PLMs have obtained rich knowledge during pre-training (Petroni et al., 2019; Davison et al., 2019). Hence, how to stimulate and exploit such knowledge is receiving increasing attention.

043

045

047

048

050

051

054

059

060

061

062

063

064

065

066

067

068

069

070

071

072

075

076

077

078

079

081

One conventional approach to achieve that is fine-tuning (Devlin et al., 2019), where we add extra classifiers on the top of PLMs and further train the models under classification objectives. Fine-tuning has achieved satisfying results on supervised tasks. However, since the extra classifier requires adequate training instances to tune, it is still challenging to apply fine-tuning in few-shot learning (Brown et al., 2020) and zero-shot learning (Yin et al., 2019) scenarios. Originated from GPT-3 (Brown et al., 2020) and LAMA (Petroni et al., 2019, 2020), a series of studies using prompts (Schick and Schütze, 2020a; Liu et al., 2021) for model tuning bridge the gap between pre-training objective and down-stream tasks, and demonstrate that such discrete or continuous prompts induce better performances for PLMs on few-shot and zero-shot tasks.

A typical way to use prompts is to wrap the input sentence into a natural language template and let the PLM conduct masked language modeling. For instance, to classify the topic of a sentence x: "What's the relation between speed and acceleration?" into the "SCIENCE" category, we wrap it into a template: "A [MASK] question: x". The prediction is made based on the probability that the word "science" is filled in the "[MASK]" token. The mapping from *label words* (e.g., "science") to the specific class (e.g., class SCIENCE) is called the *verbalizer* (Schick and Schütze, 2020a), which bridges a projection between the vocabulary and the label space and has a great influence on the performance of classification (Gao et al., 2020).

Most existing works use manual verbalizers (Schick and Schütze, 2020a, 2021), in which the designers manually think up a single word to indicate each class. To ease the human effort of designing the class name, some works propose to learn the label words using discrete search (Schick et al.,

118

122

124

ter performance, only induces a few words or embeddings that are close to the class name in terms of word sense or embedding distance. Thus they are difficult to infer words across granularities (e.g. from "science" to "physics"). If we can expand the verbalizer of the above example into $\{\text{science, physics}\} \rightarrow \text{SCIENCE, the probability of}$ making correct predictions will be considerably enhanced. Therefore, to improve the coverage and reduce the bias of the manual verbalizer, we present to incorporate external knowledge into the verbaliz-117 ers to facilitate prompt-tuning, namely, knowledgeable prompt-tuning (KPT). Since our expansion 119 is not based on optimization, it will also be more 120 favorable for zero-shot learning. Specifically, KPT contains three steps: construc-

2020) or gradient descent (Liu et al., 2021; Ham-

bardzumyan et al., 2021). However, the learned-

from-scratch verbalizer, lack of human prior knowl-

edge, is still considerably inferior to the manual

verbalizers (see Appendix A for pilot experiments),

especially in few-shot setting, and even not appli-

cable in zero-shot setting, which leaves the manual

However, manual verbalizers usually determine

the predictions based on limited information. For

instance, in the above example, the mapping

 $\{\text{science}\} \rightarrow \text{SCIENCE}$ means that only predicting

the word "science" for the [MASK] token is re-

garded as correct during inference, regardless of

the predictions on other relevant words such as

"physics" and "maths", which are also informative.

Such handcrafted one-one mapping limits the coverage of label words, thus lacking enough infor-

mation for prediction and introducing bias into the

verbalizer. Therefore, manual verbalizers are hard

to be optimal in text classification, where the se-

mantics of label words are crucial for predictions.

The optimization-based expansion, though can

be combined with manual verbalizers to yield bet-

verbalizer a decent choice in many cases.

tion, refinement, and utilization. (1) Firstly, in the 123 construction stage, we use external KBs to generate a set of label words for each label (in § 3.2). 125 Note that the expanded label words are not sim-126 ply synonyms of each other, but cover different 127 granularities and perspectives, thus are more comprehensive and unbiased than the class name. (2) 129 Secondly, to cope with the noise in the unsuper-130 vised expansion of label words, we propose four 131 refinement methods, namely, frequency refinement, 132 relevance refinement, contextualized calibration, 133

and learnable refinement (in \S 3.3), whose effectiveness is studied thoroughly in § 4. (3) Finally, we apply either a vanilla average loss function or a weighted average loss function for the utilization of expanded verbalizers, which map the scores on a set of label words to the scores of the labels.

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

We conduct extensive experiments on zero-shot and few-shot text classification tasks. The empirical results show that KPT can reduce the error rate of classification by 16%, 18%, 10%, 7% on average in 0, 1, 5, 10 shot experiments, respectively, which shows the effectiveness of KPT. In addition to the performance boost, KPT also reduces the prediction variances consistently in few-shot experiments and yields more stable performances.¹

2 **Related Work**

Two groups of research are related to KPT: prompttuning, and the verbalizer construction.

Prompt-tuning. Since the emergence of GPT-3 (Brown et al., 2020), prompt-tuning has received considerable attention. GPT-3 (Brown et al., 2020) demonstrates that with prompt-tuning and incontext learning, the large-scale language models can achieve superior performance in the low-data regime. The following works (Schick and Schütze, 2020a,b) argue that small-scale language models (Radford et al., 2018; Devlin et al., 2019; Liu et al., 2019; Lan et al., 2019) can also achieve decent performance using prompt-tuning. Prompt-tuning has been applied to a large variety of tasks such as Text Classification (Schick and Schütze, 2020a), Natural Language Understanding (Schick and Schütze, 2020b; Liu et al., 2021), Relation Extraction (Han et al., 2021; Chen et al., 2021), and Knowledge Probing (Petroni et al., 2019; Liu et al., 2021), etc.

Verbalizer Construction. As introduced in § 1, the verbalizer is an important component in prompt-tuning and has a strong influence on the performance of prompt-tuning (Holtzman et al., 2021; Gao et al., 2020). Most works use humanwritten verbalizers (Schick and Schütze, 2020a), which are highly biased towards personal vocabulary and do not have enough coverage. Some other studies (Gao et al., 2020; Shin et al., 2020; Liu et al., 2021; Schick et al., 2020) design automatic verbalizer searching methods for better verbalizer choices, however, their methods require adequate training set and validation set for optimization. Moreover, the automatically determined

¹The source code will be available upon acceptance.



Figure 1: The illustration of KPT, the knowledgeable verbalizer maps the predictions over label words into labels. And the above part is the construction, refinement and utilization processes of KPT.

verbalizers are usually synonym of the class name, which differs from our intuition of expanding the verbalizer with a set of diverse and comprehensive label words using external KB. Schick et al. (2020); Shin et al. (2020) also try multiple label words for each class. The optimal size of their label words set for each class is generally less than 10, which lacks coverage when used in text classification tasks.

3 Knowledgeable Prompt-tuning

In this section, we present our methods to incorporate external knowledge into a prompt verbalizer. We first introduce the overall paradigm of prompttuning and then elucidate how to construct, refine and utilize the knowledgeable prompt.

3.1 Overview

183

184

190

191

192

195

196

198

199

207

210

Let \mathcal{M} be a language model pre-trained on large scale corpora. In text classification task, an input sequence $\mathbf{x} = (x_0, x_1, ..., x_n)$ is classified into a class label $y \in \mathcal{Y}$. Prompt-tuning formalizes the classification task into a masked language modeling problem. Specifically, prompt-tuning wraps the input sequence with a *template*, which is a piece of natural language text. For example, assuming we need to classify the sentence $\mathbf{x} =$ "What's the relation between speed and acceleration?" into label SCIENCE (labeled as 1) or SPORTS (labeled as 2), we wrap it into

 $\mathbf{x}_{p} = [CLS] A [MASK] question : \mathbf{x}$

Then \mathcal{M} gives the probability of each word vin the vocabulary being filled in [MASK] token $\mathcal{P}_{\mathcal{M}}([MASK] = v | \mathbf{x}_p)$. To map the probabilities of words into the probabilities of labels, we define a *verbalizer* as a mapping f from a few words in the vocabulary, which form the *label word* set \mathcal{V} , to the label space \mathcal{Y} , i.e., $f: \mathcal{V} \mapsto \mathcal{Y}$. We use \mathcal{V}_y to denote the subset of \mathcal{V} that is mapped into a specific label $y, \cup_{y \in \mathcal{Y}} \mathcal{V}_y = \mathcal{V}$. Then the probability of label y, i.e., $P(y|\mathbf{x}_p)$, is calculated as 215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

233

234

235

236

237

239

240

241

242

243

244

245

246

247

$$P(y|\mathbf{x}_{p}) = g\left(P_{\mathcal{M}}([MASK] = v|\mathbf{x}_{p})|v \in \mathcal{V}_{y}\right), \quad (1)$$

where g is a function transforming the probability of label words into the probability of the label. In the above example, regular prompt-tuning may define $V_1 = \{\text{"science"}\}, V_2 = \{\text{"sports"}\}$ and g as an identity function, then if the probability of "science" is larger than "sports", we classify the instance into SCIENCE.

We propose KPT, which mainly focuses on using external knowledge to improve verbalizers in prompt-tuning. In KPT, we use KBs to generate multiple label words related to each class y, e.g., $\mathcal{V}_1 = \{$ "science", "physics", ... $\}$. And we propose four refinement methods to eliminate the noise in the expanded \mathcal{V} . Finally, we explore the vanilla average and weighted average approaches for the utilization of the expanded \mathcal{V} . The details are in the following sections.

3.2 Verbalizer Construction

The process of predicting masked words based on the context is not a single-choice procedure, that is, there is no standard correct answer, but abundant words may fit this context. Therefore, the label words mapped by a verbalizer should be equipped by two attributes: *wide coverage* and *little subjective bias*. Such a comprehensive projection is crucial to the imitation of pre-training, which is the essence of prompt-tuning. Fortunately, external structured knowledge could simultaneously meet both requirements. In this section, we introduce how we use external knowledge for two text classification tasks: topic classification and sentiment classification.

248

249

254

255

261

262

263

265

266

269

270

271

272

273

275

276

277

278

281

282

284

For topic classification, the core issue is to extract label words related to the topic from all aspects and granularities. From this perspective, we choose Related Words², a knowledge graph \mathcal{G} aggregated from multiple resources, including word embeddings, ConceptNet (Speer et al., 2017), WordNet (Pedersen et al., 2004), etc., as our external KB. The edges denote "relevance" relations and are annotated with relevance scores. We presume the the name of each class v_0 is correct and use them as the anchor node to get the neighborhood nodes $N_{\mathcal{G}}(v_0)$ whose scores are larger than a threshold η as the related words ³. Thus, each class is mapped into a set of label words $\mathcal{V}_{y} = N_{\mathcal{G}}(v_0) \cup \{v_0\}$. For binary sentiment classification, the primary goal is to extend the binary sentiment to sentiment of more granualities and aspects. We use the sentiment dictionary summarized by previous researchers ^{4,5}. Several examples of the label words in the KPT are in Table 1.

Dataset	Label	Label Words
AG's News	POLITICS SPORTS	politics, government, diplomatic, law sports, athletics, gymnastics, sportsman
IMDB	NEGATIVE POSITIVE	abysmal, adverse, alarming, angry, absolutely, accepted, acclaimed,

Table 1: Examples of the expanded label words.

3.3 Verbalizer Refinement

Although we have constructed a knowledgeable verbalizer that contains comprehensive label words, the collected label words can be very noisy since the vocabulary of the KB is not tailored for the PLM. Thus it is necessary to refine such verbalizer by retaining high-quality words. In this section, we propose four refinement methods addressing different problems of the noisy label words.

Frequency Refinement. The first problem is to handle the rare words. We assume that several words in the KB are rare to the PLM, thus the prediction probabilities on these words tend to be

```
<sup>2</sup>https://relatedwords.org
```

```
<sup>3</sup>We take \eta = 0 in the experiments
```

```
<sup>4</sup>https://www.enchantedlearning.com/
```

wordlist/positivewords.shtml

inaccurate. Instead of using a word-frequency dictionary, we propose to use *contextualized prior* of the label words to remove these words. Specifically, given a text classification task, we denote the distribution of the sentences x in the corpus as \mathcal{D} . For each sentence in the distribution, we wrap it into the template and calculate the predicted probability for each label word v in the masked position $P_{\mathcal{M}}([MASK]=v|\mathbf{x}_p)$. By taking the expectation of the probability over the entire distribution of sentences, we can get the prior distribution of the label words in the masked position. We formalize it as 287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

319

320

321

322

323

324

325

326

327

329

330

331

332

$$P_{\mathcal{D}}(v) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} P_{\mathcal{M}}([MASK] = v | \mathbf{x}_{\mathbf{p}}).$$
(2)

Empirically, we found that using a small-size *unlabeled support set* \tilde{C} sampled from the training set and with labels removed, will yield a satisfying estimate of the above expectation. Thus, assuming that the input samples $\{\mathbf{x} \in \tilde{C}\}$ have a uniform prior distribution, the contextualized prior is approximated by

$$P_{\mathcal{D}}(v) \approx \frac{1}{|\tilde{\mathcal{C}}|} \sum_{\mathbf{x} \in \tilde{\mathcal{C}}} P_{\mathcal{M}}([MASK] = v | \mathbf{x}_{p}).$$
(3)

Then we remove the label words whose prior probabilities are less than a threshold. Details can be found in Appendix C.

Relevance Refinement. As our construction of knowledgeable label words is fully unsupervised, some label words may be more relevant to their belonging class than the others. To measure the relevance of a label word to each class, we obtain the prediction probability of the label word on the support set \tilde{C} as the vector representation \mathbf{q}^v of the label words, i.e., \mathbf{q}^v 's *i*-th element is

$$\mathbf{q}_{i}^{v} = P_{\mathcal{M}}([\text{MASK}] = v | \mathbf{x}_{ip}), \mathbf{x}_{i} \in \tilde{\mathcal{C}}.$$
 (4)

To estimate the class's representation, we presume that the name of each class v_0 , such as "science" for SCIENCE, though lack of coverage, is very relevant to the class. Then we use the vector representation \mathbf{q}^{v_0} of the these names as the class's representation \mathbf{q}^y . Therefore the relevance score between a label word v and a class y is calculated as the cosine similarity between the two representation:

$$r(v,y) = \cos(\mathbf{q}^v, \mathbf{q}^y) = \cos(\mathbf{q}^v, \mathbf{q}^{v_0}).$$
 (5)

Moreover, some label words may contribute positively to multiple classes, resulting in confusion between classes. For example, the potential label

⁵https://www.enchantedlearning.com/

wordlist/negativewords.shtml

382

383

384

401

403

404

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

405

402

3.4 Verbalizer Utilization

approach is studied in Appendix E.3.

The final problem is how to map the predicted probability on each refined label word to the decision of the class label y.

Intuitively, in the training process, a small weight

is expected to be learned for a noisy label word

to minimize its influence on the prediction. Note

that in few-shot setting, calibration may not be

necessary because the probability of a label word

can be trained to the desired magnitude, i.e.,

In addition to these refinement methods, since

many label words are out-of-vocabulary for the

PLM and are split into multiple tokens by the tok-

enizer. For these words, we simply use the average

prediction score of each token as the prediction

score for the word. The influence of this simple

 $P_{\mathcal{M}}([MASK]=v|\mathbf{x}_{p}) = P_{\mathcal{M}}([MASK]=v|\mathbf{x}_{p}).$

Average. After refinement, we can assume that each label word of a class contributes equally to predicting the label. Therefore, we use the average of the predicted scores on \mathcal{V}_{y} as the predicted score for label y. The predicted label \hat{y} is

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \left(\frac{1}{|\mathcal{V}_y|} \sum_{v \in \mathcal{V}_y} \tilde{P}_{\mathcal{M}} ([MASK] = v | \mathbf{x}_p) \right).$$
(9)

We use this method in zero-shot learning since there is no parameter to be trained.

Weighted Average. In few-shot setting, supported by the Learnable Refinement, we adopt a weighted average of label words' scores as the prediction score. The refinement weights is used α_i as the weights for averaging. Thus, the predicted \hat{y} is

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \frac{\exp\left(s(y|\mathbf{x}_{\mathbf{p}})\right)}{\sum_{y'} \exp\left(s(y'|\mathbf{x}_{\mathbf{p}})\right)}, \quad (10)$$

where $s(y|\mathbf{x}_p)$ is

$$s(y|\mathbf{x}_{\mathbf{p}}) = \sum_{v \in \mathcal{V}_{y}} \alpha_{v} \log P_{\mathcal{M}}([MASK] = v|\mathbf{x}_{\mathbf{p}}).$$
(11)

This objective function is suitable for continuous optimization by applying a cross-entropy loss on the predicted probability.

3.5 Theoretical Illustration of KPT

We provide a theoretical illustration of the KPT framwork in Appendix B.

4 **Experiments**

We evaluate KPT on five text classification datasets to demonstrate the effectiveness of incorporating external knowledge into prompt-tuning.

word "physiology" of class SCIENCE may also be assigned with a high probability in a sentence of class SPORTS. To mitigate such confusion and filter the less relevant label words, we design a metric that favors the label word with high relevance merely to its belonging class and low relevance to other classes:

$$R(v) = r(v, f(v)) \frac{|\mathcal{Y}| - 1}{\sum_{y \in \mathcal{Y}, y \neq f(v)} (r(v, y))},$$
 (6)

where f(v) is the corresponding class of v.

333

334

335

338

339

341

342

344

345

349

351

354

355

361

363

369

374

375

379

Ideally, a good label word should at least has a higher relevance score for its belonging class than the average relevance score for the other classes. Therefore, we remove the label words with R(v) < v1. In practice, we have a slight modification to Equation 6, please refer to appendix C for details.

Essentially, this Relevance Refinement adopts the idea of the classical TFIDF (Jones, 1972) algorithm which estimates the relevance of a word to a document. It prefers to use a word that is relevant to a specific document while irrelevant to other documents as the keyword of the document. In KPT, a class is analogous to a document, while a label word is comparable to the word in the document. From this perspective, equation 6 is a variant of TFIDF metric.

Contextualized Calibration. The third problem is the drastic difference in the prior probabilities of label words. As previous works (Zhao et al., 2021; Holtzman et al., 2021) have shown, some label words are less likely to be predicted than the others, regardless of the label of input sentences, resulting in a biased prediction. In our setting, the label words in the KB tend to have more diverse prior probabilities, resulting in a severer problem (see Table 2). Therefore, we use the contextualized prior of label words to calibrate the predicted distribution, namely, contextualized calibration (CC):

$$\tilde{P}_{\mathcal{M}}(\text{[MASK]}=v|\mathbf{x}_{p}) \propto \frac{P_{\mathcal{M}}(\text{[MASK]}=v|\mathbf{x}_{p})}{P_{\mathcal{D}}(v)}, \quad (7)$$

where $P_{\mathcal{D}}(v)$ is the prior probability of the label word. The final probability is normalized to 1.

Learnable Refinement. In few-shot learning, the refinement can be strengthen by a learning process. Specifically we assign a learnable weight w_v to each label word v (may be already refined by the previous methods). The weights form a vector $\mathbf{w} \in \mathbb{R}^{|\mathcal{V}|}$, which is initialized to be a zero vector. The weights are normalized within each \mathcal{V}_{η} :

381
$$\alpha_v = \frac{\exp(w_v)}{\sum_{u \in \mathcal{V}_y} \exp(w_u)}.$$
 (8)

Method	AG's News	DBPedia	Yahoo	Amazon	IMDB
РТ	75.1 ± 6.2 (79.0)	$66.6 \pm 2.3 \ \text{(68.4)}$	$45.4 \pm 7.0 \ \text{(52.0)}$	80.2 ± 8.8 (87.8)	86.4 ± 4.0 (92.0)
PT+CC	$79.9 \pm 0.7 \ (81.0)$	$73.9 \pm 4.9 \ \text{(82.6)}$	$58.0 \pm 1.4 \ \text{(58.8)}$	$91.4 \pm 1.6 \ (93.5)$	$91.6 \pm 3.0 \ \text{(93.7)}$
KPT	84.8 ± 1.2 (86.7)	82.2 ± 5.4 (87.4)	61.6 ± 2.2 (63.8)	92.8 ± 1.2 (94.6)	91.6 ± 2.7 (94.0)
-FR	82.7 ± 1.5 (85.0)	81.8 ± 4.6 (86.2)	60.9 ± 1.5 (62.7)	92.8 ± 1.2 (94.6)	91.6 ± 2.8 (94.1)
-RR	$81.4 \pm 1.5 \ \text{(83.7)}$	$81.4 \pm 4.5 \ (85.8)$	$60.1 \pm 1.0 \ \text{(61.4)}$	92.8 ± 1.2 (94.6)	$91.6 \pm 2.8 \ (94.1)$
-CC	$55.5 \pm 2.8 \ {}^{(58.3)}$	$64.5 \pm 6.8 \ \text{(73.0)}$	$42.4 \pm 5.0 \ (46.8)$	$86.2 \pm 5.7 \ (92.5)$	$90.3 \pm 2.8 \ (\textbf{94.1})$

Table 2: Results of zero-shot text classification. The results of the best templates are shown in the brackets. Indentation means that the experimental configuration is a modification based on the up-level indentation.

4.1 Datasets and Templates

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

We carry out experiments on three topic classification datasets: AG's News (Zhang et al., 2015), DB-Pedia (Lehmann et al., 2015), and Yahoo (Zhang et al., 2015), and two sentiment classification datasets: IMDB (Maas et al., 2011) and Amazon (McAuley and Leskovec, 2013). The statistics of the datasets are shown in Table 8. The detailed information and the statistics of each dataset is in Appendix F.

We test all prompt-based methods using four manual templates and report both the average results (with standard error) of the four templates and the results of the best template (shown in (brackets)). The reasons for using manual templates and the specific templates for each dataset are in Appendix F.

4.2 Experiment Settings

Our experiments are based on OpenPrompt (Ding et al., 2021), which is an open-source toolkit to conduct prompt learning. For the PLM, we use RoBERTalarge (Liu et al., 2019) for all experiments. For test metrics, we use Micro-F1 in all experiments. For all zero-shot experiments, we repeat the experiments 3 times using different random seeds if randomness is introduced in the experiments, and for all few-shot experiments, we repeat 5 times. Note that considering the four templates and five/three random seeds, each reported score of prompt-based methods is the average of 20/12 experiments, which greatly reduces the randomness of the evaluation results. For the refinement based on the support set C, the size of the unlabeled support set |C| is 200. For few-shot learning, we conduct 1, 5, 10, and 20-shot experiments. For a k-shot experiment, we sample k instances of each class from the original training set to form the fewshot training set and sample another k instances per class to form the validation set. We tune the entire model for 5 epochs and choose the checkpoint with the best validation performance to test. Other

hyper-parameters can be found in Appendix G.

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

501

503

504

505

507

4.3 Baselines

In this subsection, we introduce the baselines we compare with. To better understand our proposed methods, we also compare within the performance of KPT using different configuration.

Fine-tuning (FT). Traditional fine-tuning method inputs the hidden embedding of [CLS] token of the PLM into the classification layer to make predictions. Note that fine-tuning can not be applied to the zero-shot setting, since the classification layer is randomly initialized.

Prompt-tuning (PT). The regular prompttuning method uses the class name as the only label word for each class, which is used in PET (Schick and Schütze, 2020a) and most existing works. For a fair comparison, we do not use the tricks in PET, such as self-training and prompt ensemble, which are orthogonal to our contributions.

Automatic Verbalizer (AUTO). The automatic verbalizer is proposed by PETAL (Schick et al., 2020), which uses *labeled* data to select the most informative label words *inside* a PLM's vocabulary. It is targeted at the situation when no manually defined class names are available. It's not obvious how to combine it with the manually defined class name to boost the performance, and how it can be applied in a zero-shot setting. Therefore we only compare it in the few-shot setting with no class name information given.

Soft Verbalizer (SOFT). The soft verbalizer is proposed by WARP (Hambardzumyan et al., 2021). They use a continuous vector for each class and use the dot product between the masked language model output and the class vector to produce the probability for each class. In our experiments, its class vectors are initialized with the class names' word embedding, since it is more effective with manual class names as the initial values (see Appendix A). As an optimization-based method, Soft Verbalizer is not applicable in the zero-shot setting.

Shot	Method	AG's News	DBPedia	Yahoo	Amazon	IMDB
1	FT PT AUTO SOFT	$19.8 \pm 10.4 \\ 80.0 \pm 6.0 (84.4) \\ 52.8 \pm 9.8 (57.6) \\ 80.0 \pm 5.6 (82.4) \\$	$\begin{array}{c} 8.6 \pm 4.5 \\ 92.2 \pm 2.5 & (94.3) \\ 63.0 \pm 8.9 & (68.3) \\ 92.3 \pm 2.3 & (93.3) \end{array}$	$\begin{array}{c} 11.1 \pm 4.0 \\ 54.2 \pm 3.1 \ (55.7) \\ 23.3 \pm 4.5 \ (25.0) \\ 54.3 \pm 2.7 \ (55.9) \end{array}$	$\begin{array}{c} 49.9 \pm 0.2 \\ 91.9 \pm 2.7 \ (93.2) \\ 66.6 \pm 12.5 \ (72.7) \\ 90.9 \pm 5.8 \ (93.6) \end{array}$	50.0 ± 0.0 91.2 \pm 3.7 (93.7) 75.5 \pm 15.5 (83.1) 89.4 \pm 8.9 (93.1)
	КРТ	83.7 ± 3.5 (84.6)	$\textbf{93.7} \pm \textbf{1.8} \hspace{0.1 in} (95.3)$	63.2 ± 2.5 (64.1)	$93.2 \pm 1.3 \ \textbf{(93.9)}$	$92.2 \pm 3.0 \ (93.6)$
	- LR - RR - RR - LR	$\begin{array}{c} 83.5 \pm 3.8 & (84.3) \\ 82.2 \pm 3.2 & (82.6) \\ 81.8 \pm 3.3 & (82.5) \end{array}$	$\begin{array}{c} 93.0 \pm 1.8 \hspace{0.1cm} (94.5) \\ 92.9 \pm 1.8 \hspace{0.1cm} (94.1) \\ 91.3 \pm 1.7 \hspace{0.1cm} (92.6) \end{array}$	$\begin{array}{c} 62.2 \pm 2.9 & (63.6) \\ 61.3 \pm 4.2 & (62.5) \\ 60.7 \pm 4.2 & (61.4) \end{array}$	$\begin{array}{c} \textbf{93.3} \pm \textbf{1.3} \hspace{0.1cm} (\textbf{93.9}) \\ \textbf{93.1} \pm \textbf{1.5} \hspace{0.1cm} (\textbf{93.7}) \\ \textbf{93.2} \pm \textbf{1.5} \hspace{0.1cm} (\textbf{93.9}) \end{array}$	$\begin{array}{c} 92.2\pm2.8 (93.6)\\ \textbf{92.6}\pm\textbf{1.7} (93.6)\\ 92.6\pm1.5 (93.5) \end{array}$
5	FT PT AUTO SOFT	$\begin{array}{c} 37.9 \pm 10.0 \\ 82.7 \pm 2.7 (84.0) \\ 72.2 \pm 10.1 (75.6) \\ 82.8 \pm 2.7 (84.3) \end{array}$	$\begin{array}{c} 95.8 \pm 1.3 \\ 97.0 \pm 0.6 \ (97.3) \\ 88.8 \pm 3.9 \ (91.5) \\ 97.0 \pm 0.6 \ (97.2) \end{array}$	$\begin{array}{c} 25.3 \pm 14.2 \\ 62.4 \pm 1.7 \ \ (63.9) \\ 49.6 \pm 4.3 \ \ (51.2) \\ 61.8 \pm 1.8 \ \ (63.1) \end{array}$	$52.1 \pm 1.3 \\92.2 \pm 3.3 (93.5) \\87.5 \pm 7.4 (90.8) \\93.2 \pm 1.6 (94.2)$	$\begin{array}{c} 51.4\pm1.4\\ 91.9\pm3.1 \ (92.7)\\ 86.8\pm10.1 \ (92.1)\\ 91.6\pm3.4 \ (93.9) \end{array}$
	KPT	85.0 ± 1.2 (85.9)	97.1 ± 0.4 (97.3)	$\textbf{67.2} \pm \textbf{0.8} \hspace{0.1in} \textbf{(67.8)}$	$93.4 \pm 1.9 \ (94.1)$	92.7 ± 1.5 (92.9)
	- LR - RR - RR - LR	$\begin{array}{c} \textbf{85.1} \pm \textbf{1.0} & (85.8) \\ 84.3 \pm 1.8 & (84.9) \\ 84.2 \pm 1.7 & (84.5) \end{array}$	$\begin{array}{c} 97.1 \pm 0.4 \hspace{0.1cm} (97.2) \\ \textbf{97.2} \pm \textbf{0.4} \hspace{0.1cm} (97.3) \\ 97.1 \pm 0.4 \hspace{0.1cm} (97.3) \end{array}$	$\begin{array}{c} 67.0 \pm 1.1 & (67.5) \\ \textbf{67.2} \pm \textbf{0.8} & (67.7) \\ \textbf{66.6} \pm 1.4 & (67.5) \end{array}$	$\begin{array}{c} 93.4 \pm 1.9 \hspace{0.1cm} (94.1) \\ \textbf{93.6} \pm \textbf{1.4} \hspace{0.1cm} (94.1) \\ 93.4 \pm 2.0 \hspace{0.1cm} (94.1) \end{array}$	$\begin{array}{c} 92.8 \pm 1.5 (93.0) \\ \textbf{93.0} \pm \textbf{2.0} (93.8) \\ 93.0 \pm 2.1 (93.8) \end{array}$
10	FT PT AUTO SOFT	$\begin{array}{c} 75.9 \pm 8.4 \\ 84.9 \pm 2.4 (86.1) \\ 81.4 \pm 3.8 (84.1) \\ 85.0 \pm 2.8 (86.7) \end{array}$	$\begin{array}{c} 93.8 \pm 2.2 \\ 97.6 \pm 0.4 \ (97.8) \\ 91.5 \pm 3.4 \ (95.1) \\ 97.6 \pm 0.4 \ (97.8) \end{array}$	$\begin{array}{c} 43.8 \pm 17.9 \\ 64.3 \pm 2.2 \ (64.8) \\ 58.7 \pm 3.1 \ (60.9) \\ 64.5 \pm 2.2 \ (65.0) \end{array}$	$\begin{array}{c} 83.0 \pm 7.0 \\ 93.9 \pm 1.3 \ (94.6) \\ 93.7 \pm 1.2 \ (94.5) \\ 93.9 \pm 1.7 \ (93.9) \end{array}$	76.2 ± 8.7 $93.0 \pm 1.7 (94.0)$ $91.1 \pm 5.1 (93.3)$ $91.8 \pm 2.6 (93.0)$
	KPT	86.3 ± 1.6 (87.0)	$98.0 \pm 0.2 \hspace{0.15cm} \textbf{(98.1)}$	$\textbf{68.0} \pm \textbf{0.6} \hspace{0.1in} \textbf{(68.2)}$	$93.8 \pm 1.2 \ (94.1)$	92.9 ± 1.8 (93.3)
	- LR - RR - RR - LR	$\begin{array}{c} 85.9 \pm 1.9 & (87.1) \\ 85.6 \pm 1.4 & (86.2) \\ 85.1 \pm 1.4 & (86.0) \end{array}$	$\begin{array}{c} 98.0 \pm 0.2 \hspace{0.1cm} \textbf{(98.1)} \\ 97.9 \pm 0.2 \hspace{0.1cm} \textbf{(98.0)} \\ 97.8 \pm 0.2 \hspace{0.1cm} \textbf{(97.8)} \end{array}$	$\begin{array}{c} 67.9 \pm 0.7 \ \textbf{(68.2)} \\ 67.5 \pm 1.1 \ \textbf{(68.1)} \\ 66.8 \pm 1.1 \ \textbf{(67.6)} \end{array}$	$\begin{array}{c} 93.9 \pm 1.1 \ (94.1) \\ 94.0 \pm 1.0 \ (94.7) \\ 94.1 \pm 0.9 \ (94.8) \end{array}$	$\begin{array}{c} \textbf{93.0} \pm \textbf{1.7} & (93.2) \\ \textbf{92.7} \pm \textbf{2.1} & (93.0) \\ \textbf{93.0} \pm \textbf{2.0} & (93.4) \end{array}$

Table 3: Results of 1/5/10-shot text classification. Indentation means that the experimental configuration is a modification based on the up-level indentation. For results of 20-shot experiments, please see Appendix D.

PT+CC. For zero-shot setting, we further introduce PT combined with our proposed contextualized calibration ⁶ as a baseline to see how much improvement is made by contextualized calibration instead of knowledgeable verbalizers.

For **KPT**, we experiment with different variants to better understand the proposed methods such as refinement. **-FR**, **-RR**, **-CC** and **-LR** is the variant that does not conduct Frequency Refinement, Relevance Refinement, Contextualized Calibration, and Learnable Refinement, respectively. In few-shot experiments, we presume that the supervised training data can train the output probability of each label word to the desired magnitude, thus we don't use CC and FR in the KPT . This decision is justified in Appendix E.2.

4.4 Main Results

508

509

510

511

512

513

514

515

516

517

518

519

522

523

524

525

526

527

528

530

531

In this subsection, we introduce the specific results and provide possible insights of KPT.

Zero-shot. From Table 2, we see that all the variants of KPT, except for KPT-CC, consistently outperforms PT and PT+CC baselines, which indicates the effectiveness of our methods. Comparison between PT and PT+CC proves that Contextualized

Calibration is very effective in the zero-shot setting. The results of KPT-FR-RR-CC, which is the variant without any refinement, reveal the label noise is severe in the automatically constructed knowledgeable label words. The gap between KPT-FR-RR and KPT-FR-RR-CC is larger than the gap between PT+CC and PT, demonstrating the drastic difference in the prior probabilities of the knowledgeable label words as we hypothesized in § 3.3. Comparison between KPT, KPT-FR, KPT-FR-RR proves the effectiveness of the refinement methods.

For the analysis regarding each type of classification task, we observe that the performance boost compared to the baselines in topic classification is higher than sentiment classification, which we conjecture that topic classification requires more external knowledge than sentiment classification. While CC offers huge improvement (on average +13%) over PT baseline, the incorporation of external knowledge further improves over PT+CC up to 11% on DBPedia, and 6% on AG's News and Yahoo. We also observe that the improvement brought by the refinement methods is more noticeable for topic classification tasks. By looking at the fraction of label words maintained after the refinement process (See appendix E.4), we conjecture that the sentiment dictionary that we used in sentiment clas-

558

532

533

534

⁶The same support sets are used as KPT

sification tasks contains little noise. Moreover, the improvement brought by the refinement process justifies the resilience of our methods to recover from noisy label words.

559

560

561

564

565

566

569

570

572

574

577

578

582

585

586

589

591

594

595

596

599

601

606

Few-shot. From Table 3, we first find out that prompt-based methods win over fine-tuning by a dramatic margin under nearly all situations. The gap enlarges as the shot becomes fewer. Comparing the baseline methods, the Soft Verbalizer (SOFT) generally wins over the Manual Verbalizer(PT) by a slight margin. However, automatic verbalizer (AUTO), although free of manual effort, lags behind the other verbalizers especially in a low-shot setting. The reason is obvious since the selection of label words among the vocabulary becomes inaccurate when labeled data is limited.

When comparing KPT with the baseline methods, we find KPT or its variants consistently outperform all baseline methods. On average, 17.8%, 10.3%, and 7.4% error rate reduction from the best baseline methods are achieved on 1, 5, and 10 shot experiments, respectively. Comparing within the variants of KPT, we find that RR and LR are generally effective across shots on topic classification dataset, while in sentiment classification dataset, KPT works well without the refinements, which is consistent with our previous assumptions that the sentiment dictionary has little noise. Note that the KPT-RR variant does not utilize any unlabeled support set C since we do not conduct CC and FR by default in few-shot learning. This variant is still superior to the baseline methods in most cases. In terms of variance, we can see that KPT enjoys smaller variances than baseline methods in most cases, demonstrating that the better coverage of label words stabilizes the training.

5 Analysis

Ablation studies about our refinement methods have been shown in the previous section. In this section and Appendix E, we conduct more in-depth analyses on the proposed methods.

5.1 Diversity of Top Predicted Words

One advantage of KPT is that it can generate diverse label words across different granularities. To specifically quantify such diversity, we conduct a case study. For the correctly predicted sentences of a class y, we count the frequency of label words $v \in \mathcal{V}_y$ appearing in the top-5 predictions for the [MASK] position. Then we report the top-15 frequent label words in Figure 2. Due to space limit,

only the results of SPORTS and BUSINESS category of AG's News are shown. As shown in Figure 2, a diversity of label words, instead of mainly the original class names, are predicted. And the predicted label words cover various aspects of the corresponding topic. For example, for the topic SPORTS, the predicted "leagues", "football", "coach" are related to it from different angles.



Figure 2: Frequent words appearing in the top-5 predictions. The results for two classes: SPORTS (left) and BUSINESS (right) are drawn.

5.2 Other Analyses

In addition to the visualization, we study the influence of the support set's size on zero-shot text classification in Appendix E.1. Then we justify that few-shot learning via labeled data eases the need for calibration and frequency-based refinement in Appendix E.2. We also demonstrate that our approach to handling the out-of-vocabulary (OOV) words is reasonable in Appendix E.3. Moreover, we take a closer look at the refinement process by analyzing the fraction of label words maintained during refinement in Appendix E.4. Finally, we discuss the potential use of the proposed methods when knowledge bases resources are not readily available in Appendix E.5.

6 Conclusion

In this paper, we propose KPT, which expands the verbalizer in prompt-tuning using the external KB. To better utilize the KB, we propose refinement methods for the knowledgeable verbalizer. The experiments show the potential of KPT in both zero-shot settings and few-shot settings. For future work, there are open questions related to our research for investigation. (1) Better approaches for combining KB and prompt-tuning in terms of template construction and verbalizer design. (2) Incorporating external knowledge into prompt-tuning for other tasks such as text generation. 610 611 612

609

613 614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

References

645

646

647

659

670

671

672

673

674

675

686

687

693

- Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*.
 - Xiang Chen, Xin Xie, Ningyu Zhang, Jiahuan Yan, Shumin Deng, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2021. Adaprompt: Adaptive promptbased finetuning for relation extraction. *arXiv preprint arXiv:2104.07650*.
- Joe Davison, Joshua Feldman, and Alexander Rush. 2019. Commonsense knowledge mining from pretrained models. In *Proceedings of EMNLP*, Hong Kong, China. Association for Computational Linguistics.
 - Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL*.
 - Ning Ding, Shengding Hu, Weilin Zhao, Yulin Chen, Zhiyuan Liu, Hai-Tao Zheng, and Maosong Sun. 2021. Openprompt: An open-source framework for prompt-learning. arXiv preprint arXiv:2111.01998.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2020. Making pre-trained language models better few-shot learners. *arXiv preprint arXiv:2012.15723*.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. 2021. WARP: Word-level Adversarial ReProgramming. In *Proceedings of ACL*.
- Xu Han, Weilin Zhao, Ning Ding, Zhiyuan Liu, and Maosong Sun. 2021. Ptr: Prompt tuning with rules for text classification. *arXiv preprint arXiv:2105.11259*.
- Ari Holtzman, Peter West, Vered Schwartz, Yejin Choi, and Luke Zettlemoyer. 2021. Surface form competition: Why the highest probability answer isn't always right. *arXiv preprint arXiv:2104.08315*.
- Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*.
- Lingkai Kong, Haoming Jiang, Yuchen Zhuang, Jie Lyu, Tuo Zhao, and Chao Zhang. 2020. Calibrated language model fine-tuning for in- and outof-distribution data. In *Proceedings of EMNLP*.
- Kamran Kowsari, Kiana Jafari Meimandi, Mojtaba Heidarysafa, Sanjana Mendu, Laura Barnes, and Donald Brown. 2019. Text classification algorithms: A survey. *Information*, 10(4):150.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*.

Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia–a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic web*, 6(2):167–195. 699

700

702

703

705

706

707

708

709

711

713

714

715

716

717

719

720

721

722

723

724

725

726

727

729

730

732

733

734

737

738

739

740

741

742

743

744

745

746

747

748

749

750

- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021. Gpt understands, too. *arXiv preprint arXiv:2103.10385*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Andrew Maas, Raymond E Daly, Peter T Pham, Dan Huang, Andrew Y Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of ACL*.
- Julian McAuley and Jure Leskovec. 2013. Hidden factors and hidden topics: understanding rating dimensions with review text. In *Proceedings of RecSys*.
- Yu Meng, Yunyi Zhang, Jiaxin Huang, Chenyan Xiong, Heng Ji, Chao Zhang, and Jiawei Han. 2020. Text classification using label names only: A language model self-training approach. In *Proceedings of EMNLP*.
- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng.
 2016. Ms marco: A human generated machine reading comprehension dataset. In *Proceedings of CoCo@ NeurIPS*.
- Ted Pedersen, Siddharth Patwardhan, Jason Michelizzi, et al. 2004. Wordnet:: Similarity-measuring the relatedness of concepts. In *Proceedings of AAAI*, volume 4, pages 25–29.
- Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of NAACL*, pages 2227– 2237.
- Fabio Petroni, Patrick Lewis, Aleksandra Piktus, Tim Rocktäschel, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2020. How context affects language models' factual predictions. In *Automated Knowledge Base Construction*.
- Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. 2019. Language models as knowledge bases? In *Proceedings of EMNLP*, pages 2463– 2473.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.

752

- 758 759 760 761 762
- 764 765 766 767 768 769
- 770 771 772
- 773 774

776

- 777 778 779 780 781 782
- 7 7 7 7 7
- 789 790 791 792 793 794
- 794 795
- 796 797
- 7

2

ð

- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1– 67.
 - Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of EMNLP*, pages 2383–2392.
- Timo Schick, Helmut Schmid, and Hinrich Schütze. 2020. Automatically identifying words that can serve as labels for few-shot text classification. In *Proceedings of COLING*, pages 5569–5578.
 - Timo Schick and Hinrich Schütze. 2020a. Exploiting cloze questions for few-shot text classification and natural language inference. *arXiv preprint arXiv:2001.07676*.
 - Timo Schick and Hinrich Schütze. 2020b. It's not just size that matters: Small language models are also few-shot learners. *arXiv preprint arXiv:2009.07118*.
 - Timo Schick and Hinrich Schütze. 2021. It's not just size that matters: Small language models are also few-shot learners. In *Proceedings of NAACL*, pages 2339–2352.
 - Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Autoprompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980.*
 - Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of AAAI*.
 - Han Xu, Zhang Zhengyan, Ding Ning, Gu Yuxian, Liu Xiao, Huo Yuqi, Qiu Jiezhong, Zhang Liang, Han Wentao, Huang Minlie, et al. 2021. Pre-trained models: Past, present and future. *arXiv preprint arXiv:2106.07139*.
 - Wenpeng Yin, Jamaal Hay, and Dan Roth. 2019. Benchmarking zero-shot text classification: Datasets, evaluation and entailment approach. In *Proceedings of EMNLP*, pages 3914–3923.
 - Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28.
 - Tony Z Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. *arXiv preprint arXiv:2102.09690*.

A Pilot Experiments

804

811

812

813

814

815

816

817

818

819

822

823

824

826

827

829

832

833

838

839

841

842

As pointed out by (Gao et al., 2020), manually defined verbalizer is competitive or even better than automatically searched/optimized verbalizers, which strengthens our motivation to improve over manual verbalizers by injecting more external human knowledge. To further illustrate the advantage of manual verbalizer, we conduct pilot experiments in soft verbalizer. Soft Verbalizer (Hambardzumyan et al., 2021) can be initialized with the predefined class name as the label words, which is adopted by us as a baseline in Table 3. It can also be randomly initialized without the manually defined class names. We test the performance of Soft Verbalizer with and without the manually defined class name in 5 and 10 shot experiments. From Table 4, we can see that the gaps between the variants are generally large. Therefore further improving the verbalizer with manually defined class name is a promising direction than the learned-from-scratch verbalizer without any human prior.

B A Theoretical Illustration of KPT

In this section we provide a theoretical analysis of the whole framework used by KPT. In prompt tuning, given a text \mathbf{x} , we wrap it into a template to form a wrapped sentence \mathbf{x}_p . we then predict the probability of the label word v using a PLM:

$$p([M]=v|\mathbf{x}) = P_{\mathcal{M}}([M]=v|\mathbf{x}_{p}), \qquad (12)$$

where [M] is short for [Mask], denoting the label word's prediction at the masked position of the template.

Then, if multiple label words are used to contribute to a single label, the predicted probability of the label is defined by marginalizing the probability of predicting all the label words, i.e.,

$$p(\mathbf{Y} = y | \mathbf{x}) = \sum_{v \in V_{\mathcal{Y}}} p(\mathbf{Y} = y, \text{[MASK]} = v | \mathbf{x}).$$
(13)

Since the prediction of Y is independent of x given v, we can write Equation 13 into

$$\sum_{v \in V_{\mathcal{Y}}} p(\mathbf{Y} = y | [M] = v) p([M] = v | \mathbf{x})$$

$$= \sum_{v \in V_{\mathcal{Y}}} p(\mathbf{Y} = y | [M] = v) P_{\mathcal{M}}([M] = v | \mathbf{x}_{p}).$$
(14)

Using the Bayes Theorem and assuming a balanced classification problem, Equation 15 can be transformed into

$$\sum_{v \in V_{\mathcal{Y}}} \frac{p([\mathbb{M}]=v|\mathbf{Y}=y)p(\mathbf{Y}=y)}{p([\mathbb{M}]=v)} P_{\mathcal{M}}([\mathbb{M}]=v|\mathbf{x}_{p})$$

$$\propto \sum_{v \in V_{\mathcal{Y}}} \frac{p([\mathbb{M}]=v|\mathbf{Y}=y)}{p([\mathbb{M}]=v)} P_{\mathcal{M}}([\mathbb{M}]=v|\mathbf{x}_{p}).$$
(15)

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

Now, the prediction probability of the label is composed of three parts.

(1) The first part p(v|Y = y) is the probability of predicting the specific label word v given the class label y. Intuitively, if a label word is relevant to label y, this term will be assigned a high probability. In KPT, the **Relevance Refinement** estimate this probability using a quantized objective, i.e., if a relevance score exceeds the threshold 1, it will be maintained, otherwise, it will be filtered. On the other hand, **Learnable Refinement** estimates this probability using continuous weights.

(2) The second part is p([M]=v) in the denominator. This term is actually the prior probability of label words v, which is estimated by our **Contextualized Calibration**. Previous works also try to approach this term using a context-free manner (Holtzman et al., 2021; Zhao et al., 2021).

(3) The last term $P_{\mathcal{M}}([MASK] = v | \mathbf{x}_p)$ is the probability of the label words v predicted by the PLM, which is the only component in most works such as Manual Verbalizers (Schick and Schütze, 2020a), yielding a sub-optimial solution compared to KPT.

Verbalizers with multiple label words for a class label can all be formalized into this framework once it uses Equation 13 as their backbone hypothesis. However, to the best of our knowledge, KPT is the first to combine all of the three components to form a powerful verbalizer.

C Practical Issues of Refinement

In this section, we detail the refinement process by making some practical modifications to the methods in § 3.3.

Frequency Refinement. For Frequency Refinement, since the absolute value distribution of the contextualized prior probability may be different for each task, determining a specific threshold of the contextualized prior probability may be tricky and elusive. We use a ranking-based threshold, i.e., we filter the label words that appear in the lower half of the contextualized prior probability.

Relevance Refinement. For Relevance Refinement, we observe that in the classification task with

Shot	Method	AG's News	DBPedia	Yahoo	Amazon	IMDB
5	SOFT SOFT w.o. M	$\begin{array}{c} 82.8 \pm 2.7 \hspace{0.1 cm} (84.3) \\ 63.4 \pm 11.3 \hspace{0.1 cm} (64.7) \end{array}$	$\begin{array}{c} 97.0 \pm 0.6 \hspace{0.2cm} (97.2) \\ 82.1 \pm 5.9 \hspace{0.2cm} (86.1) \end{array}$	$\begin{array}{c} 61.8 \pm 1.8 \hspace{0.1 cm} \tiny{(63.1)} \\ 24.5 \pm 6.2 \hspace{0.1 cm} \scriptstyle{(27.2)} \end{array}$	$\begin{array}{c} 93.2 \pm 1.6 \hspace{0.1 cm} (94.2) \\ 79.2 \pm 10.5 \hspace{0.1 cm} (85.5) \end{array}$	$\begin{array}{l} 91.6 \pm 3.4 \hspace{0.1 cm} (93.9) \\ 83.6 \pm 11.5 \hspace{0.1 cm} (93.4) \end{array}$
10	SOFT SOFT w.o. M	$\begin{array}{c} 85.0 \pm 2.8 \hspace{0.1in} (86.7) \\ 77.4 \pm 4.8 \hspace{0.1in} (79.1) \end{array}$	$\begin{array}{c} 97.6 \pm 0.4 \hspace{0.1 cm} (97.8) \\ 94.9 \pm 2.5 \hspace{0.1 cm} (95.9) \end{array}$	$\begin{array}{c} 64.5 \pm 2.2 \hspace{0.2cm} (65.0) \\ 42.6 \pm 8.3 \hspace{0.2cm} (48.1) \end{array}$	$\begin{array}{c} 93.9 \pm 1.7 \hspace{0.1 cm} (93.9) \\ 92.9 \pm 2.0 \hspace{0.1 cm} (94.0) \end{array}$	$\begin{array}{c} 91.8 \pm 2.6 \hspace{0.1 cm} (93.0) \\ 88.7 \pm 9.7 \hspace{0.1 cm} (93.8) \end{array}$

Table 4: Pilot experiment on soft verbalizer justifies the need of human (expert) knowledge into the verbalizer. SOFT is the soft verbalizer with class name and SOFT w.o. M is the variant without the manual verbalizer.

Shot	Method	AG's News	DBPedia	Yahoo	Amazon	IMDB
20	FT PT AUTO SOFT	$\begin{array}{l} 85.4 \pm 1.8 \\ 86.5 \pm 1.6 (87.0) \\ 85.7 \pm 1.4 (86.1) \\ 86.4 \pm 1.7 (87.1) \end{array}$	$\begin{array}{c} 97.9 \pm 0.2 \\ 97.9 \pm 0.3 & (98.1) \\ 92.2 \pm 2.7 & (94.9) \\ 98.0 \pm 0.3 & (98.1) \end{array}$	$\begin{array}{c} 54.2 \pm 18.1 \\ 67.2 \pm 1.1 \ \ (67.5) \\ 65.0 \pm 1.8 \ \ (66.9) \\ 67.4 \pm 0.7 \ \ (67.5) \end{array}$	$\begin{array}{l} 71.4 \pm 4.3 \\ 93.5 \pm 1.0 (94.4) \\ \textbf{93.9} \pm \textbf{1.1} (94.1) \\ 93.8 \pm 1.6 (94.2) \end{array}$	$\begin{array}{l} 78.5 \pm 10.1 \\ 93.0 \pm 1.1 \ (93.6) \\ 92.8 \pm 2.0 \ (94.0) \\ \textbf{93.5 \pm 0.9} \ (94.0) \end{array}$
	KPT	$87.2 \pm 0.8 \ (87.5)$	$98.1\pm0.3~(98.2)$	$68.9 \pm 0.8 \ \text{(69.3)}$	$93.7 \pm 1.6 \ (94.4)$	$93.1 \pm 1.1 \ \text{(93.5)}$
	- LR	87.7 ± 0.6 (87.8)	98.1 ± 0.3 (98.2)	68.8 ± 0.9 (69.8)	93.4 ± 2.3 (94.3)	93.4 ± 0.9 (93.6)
	- RR - RR - LR	$\begin{array}{l} 87.3 \pm 0.8 \hspace{0.1 cm} (87.5) \\ 87.1 \pm 0.9 \hspace{0.1 cm} (87.4) \end{array}$	$\begin{array}{c} 98.1 \pm 0.3 \hspace{0.1 cm} (98.2) \\ 98.1 \pm 0.3 \hspace{0.1 cm} (98.2) \end{array}$	$\begin{array}{l} 68.8 \pm 0.9 & \tiny{(68.9)} \\ \textbf{69.0} \pm \textbf{0.7} & \tiny{(69.3)} \end{array}$	$\begin{array}{c} 93.6 \pm 1.3 \hspace{0.2cm} (94.2) \\ 93.7 \pm 0.9 \hspace{0.2cm} (94.5) \end{array}$	$\begin{array}{l} 93.1 \pm 0.8 \hspace{0.1 cm} (93.6) \\ 93.1 \pm 0.8 \hspace{0.1 cm} (93.7) \end{array}$

Table 5: Results of 20-shot text classification. Average Micro-F1 scores and variances using four templates are shown. The Micro-F1 scores of the best templates are shown in the brackets.Indentation means that the experimental configuration is a modification based on the up-level indentation.

only a few classes, it's better to provide a stricter criterion to ensure that the relevance scores of a label word to *any other* class is lower than the score to the belonging class, i.e., *maximum* in the term of IDF-score is preferred. To keep a unified criterion, we use a norm-based IDF-score.

$$R^{d}(v) = r(v, f(v)) \left(\frac{|\mathcal{Y}| - 1}{\sum_{y \in \mathcal{Y}, y \neq f(v)} (r(v, y)^{d})}\right)^{1/d}$$
(16)

where

891

893

897

900

901

902

903

904

905

$$d = \frac{C}{|\mathcal{Y}| - 2 + \epsilon} + 1, C > 0.$$
(17)

This criterion will approximate the maximum value in $\{r(v, y)|y \in |Y|, y \neq f(v)\}$ in classification with only a few labels, and revert to the mean score in Equation 6 when conducting classification with many labels. We take C = 10 (without trial and error) in the experiments. And $0 < \epsilon \ll 1$ is a small number to prevent numerical error.

D Results of 20 Shot Setting

Due to the space limitation, we report the perfor-906 mance of 20-shot classification in this section. As 907 we can see in Table 5, the gap between different 908 methods narrows as the training data becomes suf-909 ficient. However, KPT and its variants still win 910 by a consistent margin over the baseline methods. 911 Surprisingly, with more training data, Learnable 912 Refinement does not become more powerful as we 913

may hypothesize. We conjecture that it is because all label words, even with some noise, *can* serve as training objectives for prompt tuning. This perspective is similar to Gao et al. (2020) that using "bad" as a label word for the class "positive" can still do classification even though the performance decreases.

E Further Analyses and Ablation Studies



Figure 3: Support set size w.r.t. zero-shot performance. The points at size=0 are the performances of PMI_{DC}.

E.1 Calibration and Contextualized Calibration.

In fine-tuning, calibration has been studied under the topic of prediction confidence and out-ofdistribution detection (Kong et al., 2020). Recently, it got renascent attention in prompt learning (Zhao et al., 2021; Holtzman et al., 2021). In prompt learning, the PLM has a natural tendency to predict one word over another word regardless of the real sentence input. For example, GPT-3 prefers to 914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

predict "positive" over "negative" given "N/A" as 932 the input sentence (Zhao et al., 2021). Therefore 933 the calibration is crucial (see Table 2) when no 934 posterior optimization is conducted, i.e., in zeroshot learning. Existing methods such as PMI_{DC} 936 propose only using the empty template without fill-937 ing the template with the instances in the corpus, 938 for example, "A [Mask] question :", to produce the calibration logits. Our proposed Contextualized Calibration utilizes a limited amount of unlabeled 941 support data to yield significantly better results. 942 However, since we target the data-scarce scenario, 943 we study in detail the amount of unlabeled data necessary to produce a satisfying calibration result. 945 In Figure 3, we draw the performance of KPT -946 RR with different support set sizes |C|. We also draw the performance of PMI_{DC} on the $|\mathcal{C}| = 0$ for comparison.

> From Figure 3, we find that $|\tilde{C}| \sim 50$ is enough to yield a satisfying calibration. Contextualized calibrate is more effective in classification with many classes, while calibrate without the context is effective in classification with few classes.

951

953

956

958

959

960

961

In addition, we must point out that if we have a set of sentences to classify in real-world scenarios, we can use these sentences themselves as the support set to conduct more accurate Contextualized Calibration.

E.2 Supervised Data Ease the Need for Calibration.

Although calibration is crucial for the zero-shot setting, we do not perform calibration for the few-963 shot setting because we assume that the posterior 964 probability of the label words can be trained to 965 the desired magnitude with only a few training instances. We also do not perform Frequency Re-967 finement for few-shot learning due to the same assumption. To verify the assumption empirically, we add both Contextualized Calibration and Frequency Refinement to KPT and test the performance under 971 different settings. The results are shown in Table 6. 972 The performance comparison to KPT without CC 973 and FR in Table 3 and Table 5 are shown using up 974 arrows and down arrows. We can see that except 975 in Yahoo, the improvement is not consistent for 976 even negative, which supports our assumption that 977 the need for calibration is greatly eased with the 978 supervised input data. 979

E.3 How to Handle the OOV Label Words?

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1003

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

Since the knowledgeable verbalizer is expanded using external resources which may not be tailored for the vocabulary of PLM. Thus, many label words are out-of-vocabulary (OOV) and are split into multiple tokens by the tokenizer. For these words, as mentioned in § 3.3, we average the prediction probability of each token in the *single* [MASK] position, which may not be very reasonable at the first glance. Therefore, we conduct an ablation study that whether forcing the label words to be a single token in the vocabulary of the PLM leads to better performance. The results under different shots are shown in Table 7. Surprisingly, making the singletoken restriction does not yield stable improvement, instead, in many cases, the performance degrades by minor margins. Therefore we conclude that our method to handle OOV label words that are split by the tokenizer into multiple tokens is simple yet reasonable. More importantly, the label words expanded by the knowledge bases but not in the 1000 PLM's vocabulary can serves as good label words 1001 in prompt tuning as well. 1002

E.4 Visualization of the Refinement Process.

In this section, we report the number of label words that remained after Frequency Refinement and Relevance Refinement process. As we can see, these refinement methods remove a large fraction of label words while retaining the ones that are most informative. However, even the fewest number of remaining label words exceeds 100, which is far more than the number of label words in the previous works (Schick et al., 2020). The broad coverage of label words contributes to the success of KPT.



Figure 4: The number of remaining label words after Frequency Refinement and Relevance Refinement.

Shot	Method	AG's News	DBPedia	Yahoo	Amazon	IMDB
1	KPT + CC + FR	83.4 ± 4.0 (84.6)	94.0 ↑± 2.0 (95.7)	63.3 ↑± 2.0 (64.9)	$93.2 \pm 1.2 \ (94.0)$	92.1 ± 3.2 (93.8)
5	KPT + CC + FR	84.6 ± 1.3 (85.1)	97.3 [↑] ± 0.3 (97.4)	67.3 [↑] ± 1.1 (67.7)	94.0 ↑± 1.2 (94.7)	92.7 ± 1.6 (93.1)
10	KPT + CC + FR	85.9 ± 1.7 (86.7)	98.1 ↑± 0.2 (98.2)	$68.0 \pm 1.1 \ (68.6)$	93.3 ± 1.8 (93.7)	92.9 ± 1.8 (93.6)
20	KPT + CC + FR	$87.3\uparrow\pm0.8$ (87.6)	$98.0 \downarrow \pm 0.4 (98.2)$	$69.1 \uparrow \pm 0.7$ (69.5)	93.5 ↓± 1.1 (93.9)	$93.1 \ \pm 1.3 \ \text{(93.5)}$

Table 6: Results of Contextualized Calibration and Frequency Refinement on few-shot experiments. The green up arrow \uparrow means the results is higher than KPT in Table 3 and Table 5, and the red down arrow \downarrow means the results is lower than KPT in Table 3 and Table 5.

Shot	Method	AG's News	DBPedia	Yahoo	Amazon	IMDB
0	KPT + ST	$84.9 \uparrow \pm 1.0 (86.3)$	81.0 ± 4.3 (85.2)	62.7 ↑± 1.1 (64.4)	92.8 ± 1.2 (94.7)	91.5 ± 2.8 (94.1)
1	KPT + ST	83.4 ± 3.9 (84.2)	$94.0 \uparrow \pm 1.8 (95.8)$	$62.5 \pm 2.3 (63.5)$	93.3 ↑± 1.4 (94.1)	92.1 \pm 3.5 (93.6)
5	KPT + ST	$84.7 \pm 1.8 (85.4)$	$97.1 \pm 0.5 \ (97.2)$	66.8 ± 1.0 (67.3)	93.3 ± 2.1 (93.8)	93.1 ↑± 1.4 (93.3)
10	KPT + ST	$86.3 \pm 1.5 \ (86.8)$	$98.0 \pm 0.2 \ (98.1)$	$67.6 \pm 0.9 (67.9)$	94.0 ↑± 1.0 (94.1)	92.7 ± 1.8 (93.6)
20	KPT + ST	$87.2 \downarrow \pm 1.1 \ (87.6)$	$97.9 \downarrow \pm 0.4 (98.1)$	$68.6 \downarrow \pm 0.7 (69.1)$	$93.5 \uparrow \pm 1.8 \ (94.0)$	92.9 ± 1.2 (93.4)

Table 7: Results of restricting the expanded label word to be a single token in the PLM's vocabulary, where ST denotes "single token". The green up arrow \uparrow means the results is higher than KPT in Table 3 and Table 5, and the red down arrow \downarrow means the results is lower than KPT in Table 3 and Table 5.

E.5 Potential Usage without External KB.

Although KBs are ubiquitous in natural language processing, there are cases that no readily available KBs can be found for specific tasks. For these tasks, if we have enough unlabeled corpus, we can use the methods proposed by LOTClass (Meng et al., 2020) to mine potential label words from the corpus. More specifically, LOTClass (Meng et al., 2020) uses a self-supervised objective to train the PLM to extract the topic-related words from the whole unlabeled training corpus. Experiments that combine KPT with LOTClass are beyond the scope of our work, but we believe the combination of the two can be very effective.

F Datasets and Templates

We carry out experiments on three topic classification datasets: AG's News (Zhang et al., 2015), DB-Pedia (Lehmann et al., 2015), and Yahoo (Zhang et al., 2015), and two sentiment classification datasets: IMDB (Maas et al., 2011) and Amazon (McAuley and Leskovec, 2013). The statistics of the datasets are shown in Table 8.

Name	Туре	# Class	Test Size
AG's News	Topic	4	7600
DBPedia	Topic	14	70000
Yahoo	Topic	10	60000
Amazon	Sentiment	2	10000
IMDB	Sentiment	2	25000

Table 8: The statistics of each dataset.

Due to the rich expert knowledge contained, the manual templates are proven to be competitive with or better than auto-generated templates (Gao et al., 2020) even though they are simpler to be constructed. Therefore we use manual templates in our experiments. Manual templates are also more applicable than auto-generated templates in the zero-shot setting. To mitigate the influence of different templates, we test KPT under four manual templates that are either introduced by (Schick and Schütze, 2020a) or tailored to fit the dataset for each experimental configuration. The templates for each dataset is listed below.

AG's News. AG's News is a news' topic classification dataset. In this dataset, we follow PET (Schick and Schütze, 2020a) to design the templates. However, their best performance pattern $T_1(\mathbf{x}) = \text{``[MASK]}$ news : **x**'' requires the [MASK] token to be capitalized, which is not suitable for the label words in KB. And some of their templates are not informative and yield low performances. Therefore, we define four slightly changed templates:

$T_1(\mathbf{x}) = \mathbf{A}$ [MASK] news: \mathbf{x}	
$T_2(\mathbf{x}) = \mathbf{x}$ This topic is about [MASK].	
$T_3(\mathbf{x})=$ [Category : [MASK]] \mathbf{x}	
$T_4(\mathbf{x})=$ [Topic : [MASK]] \mathbf{x}	

DBPedia. In a DBPedia sample, we are given a paragraph b paired with a title a, in which the title is the subject of paragraph. The task is to determine the topic (or the type) of the subject. Different from other topic classifications, the paragraph can emphasize topics that are different from the title. For example, in a paragraph about an audio company, the main 1060 1061

1040

1042

1043

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1058

1062 1063 1064

1065

1067

1037 1038 1039

1015

1016

1017

1019

1020

1021

1022

1023

1025

1026

1027

1028

1030

1031

1032

1033

1034 1035

paragraph talks about music, albums, etc., but the correct label is "company" rather than "music". Therefore, we define the following templates:

> $T_1(\mathbf{a}, \mathbf{b}) = \mathbf{a} \mathbf{b} \,\tilde{\mathbf{a}} \text{ is a [MASK]} .$ $T_2(\mathbf{a}, \mathbf{b}) = \mathbf{a} \mathbf{b} \text{ In this sentence, } \tilde{\mathbf{a}} \text{ is a [MASK]} .$ $T_3(\mathbf{a}, \mathbf{b}) = \mathbf{a} \mathbf{b} \text{ The type of } \tilde{\mathbf{a}} \text{ is [MASK]} .$ $T_4(\mathbf{a}, \mathbf{b}) = \mathbf{a} \mathbf{b} \text{ The category of } \tilde{\mathbf{a}} \text{ is [MASK]} .$

where $\tilde{\mathbf{a}}$ means removing the last punctuate in the title.

Yahoo. Yahoo is a topic classification dataset about the questions raised in yahoo website (Zhang et al., 2015). We use the same templates as AG's News, except that we change the word "news" into "question" in the $T_1(\mathbf{x})$:

 $T_1(\mathbf{x}) = \mathbf{A} \text{ [MASK] question : } \mathbf{x}$ $T_2(\mathbf{x}) = \mathbf{x} \text{ This topic is about [MASK].}$ $T_3(\mathbf{x}) = \text{ [Category : [MASK]] } \mathbf{x}$ $T_4(\mathbf{x}) = \text{ [Topic : [MASK]] } \mathbf{x}$

IMDB. IMDB is a sentiment classification dataset about movie reviews. Similar to the template defined in (Schick and Schütze, 2020a) for sentiment classification, we define the following template:

 $T_1(\mathbf{x}) = \text{It was [MASK] .} \mathbf{x}$ $T_2(\mathbf{x}) = \text{Just [MASK] !} \mathbf{x}$ $T_3(\mathbf{x}) = \mathbf{x} \text{ All in all, it was [MASK].}$ $T_4(\mathbf{x}) = \mathbf{x} \text{ In summary, the film was [MASK].}$

Amazon. Amazon is another sentiment classification dataset , we define the following template:

 $T_1(\mathbf{x}) = \text{ It was [MASK] .x}$ $T_2(\mathbf{x}) = \text{ Just [MASK] ! x}$ $T_3(\mathbf{x}) = \mathbf{x} \text{ All in all, it was [MASK].}$ $T_4(\mathbf{x}) = \mathbf{x} \text{ In summary, it was [MASK]".}$

Since the test set of amazon is unnecessarily large for efficient testing, we randomly sample 10,000 samples from the 400,000 test samples to test, which is proven to have tiny influence on the performance in our pilot experiments.

G Experimental Settings

We list the hyper-parameters in Table 9. Most of the hyper-parameters are the default parameters from Huggingface Transformers⁷.

Hyper-parameter	Dataset	Value
truncate length	AG's News, DB- Pedia, Yahoo	128
truncate length	Amazon, Imdb	512
warmup steps	All	0
learning rate	All	3e-5
maximum epochs	All	5
adam epsilon	All	1e-8

Table 9: Hyper-parameter settings.

For soft verbalizer, we use a learning rate of 3e -10984 to its soft label words' embeddings to encourage1099a faster convergence.1100

1085

1086

1087

1068

1069

1070

1071

1072

1074

1075

1076

1077

1078

1079

1080

1081

1083

1090 1091 1092 1093

1088

1089

⁷https://huggingface.co/transformers/