

Text Grafting: Near-Distribution Weak Supervision for Minority Classes in Text Classification

Anonymous ACL submission

Abstract

For extremely weak-supervised text classification, pioneer research generates pseudo labels by mining texts similar to the class names from the raw corpus, which may end up with very limited or even no samples for the minority classes. Recent works have started to generate the relevant texts by prompting LLMs using the class names or definitions; however, there is a high risk that LLMs cannot generate in-distribution (i.e., similar to the corpus where the text classifier will be applied) data, leading to ungeneralizable classifiers. In this paper, we combine the advantages of these two approaches and propose to bridge the gap via a novel framework, *text grafting*, which aims to obtain clean and near-distribution weak supervision for minority classes. Specifically, we first use LLM-based logits to mine masked templates from the raw corpus, which have a high potential for data synthesis into the target minority class. Then, the templates are filled by state-of-the-art LLMs to synthesize near-distribution texts falling into minority classes. Text grafting shows significant improvement over direct mining or synthesis on minority classes. We also use analysis and case studies to comprehend the property of text grafting.

1 Introduction

Recent research has made rapid progress on extremely weak-supervised text classification (XWS-TC) (Wang et al., 2023), limiting the supervision to a brief natural-language description without any annotated samples. For example, text mining-based XWS-TC (Meng et al., 2020; Wang et al., 2021; Shen et al., 2021; Mekala et al., 2022; Zhao et al., 2023; Dong et al., 2023a) takes only class names or seed words from humans and discovers potential in-class texts following designated heuristics.

Minority classes are arguably the most challenging part of XWS-TC. The class distribution in real-world datasets is often a long-tailed distribution (Zhang et al., 2023), with a non-trivial number

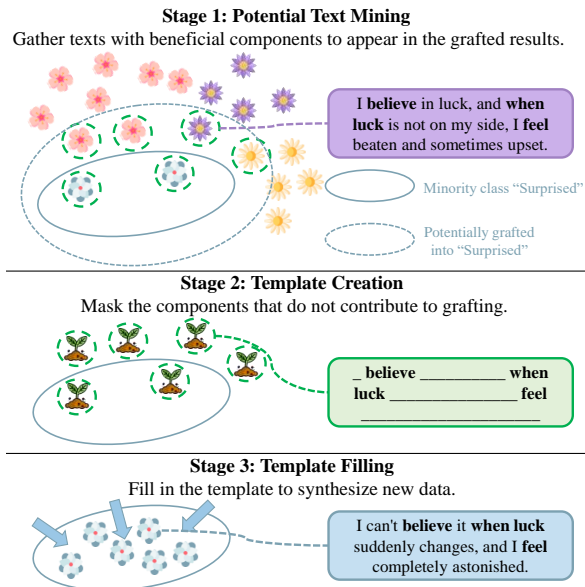


Figure 1: The framework of text grafting.

Framework	Mining?	Train Data	Data Quality	In-Distribution
Text Mining	Text	Raw	Noisy	Yes
Data Synthesis	None	Generated	Clean	Hardly
Text Grafting (ours)	Template	Grafted	Clean	Mostly

Table 1: High-level comparison among three discussed XWS-TC frameworks.

of minority classes. These minority classes have a very small number of documents in the raw corpus, therefore, it is difficult to locate the right documents by mining-based methods, leading to noisy pseudo-labels. Under extreme circumstances, the mining-based methods may end up with no sample for minority classes.

A potential way to address this issue is data synthesis-based XWS-TC (Ye et al., 2022a,b; Peng and Shang, 2024), which hopes to generate in-class texts by prompting large language models (LLM) (Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023a,b; Meta, 2024; Mesnard et al., 2024; OpenAI, 2024) with class names or definitions. However, such synthesized texts may follow

058 a distribution different from the corpus where the
059 text classifier will be later applied (Mitchell et al.,
060 2023), which makes the learned text classifier out-
061 of-distribution, leading to poor performance.

062 This paper combines the advantages of mining-
063 based and synthesis-based frameworks to propose
064 a new framework, *text grafting*, which aims to ob-
065 tain clean and near-distribution weak supervision
066 for minority classes. As specified in Figure 1, text
067 grafting incorporates three stages: (1) *Potential*
068 *Text Mining* gathers raw texts with beneficial com-
069 ponents to synthesize in-class texts for the target
070 minority class. (2) *Template Creation* forms tem-
071 plates by masking the components that do not con-
072 tribute to the in-class text synthesis. (3) *Template*
073 *Filling* synthesizes in-class texts by filling in the
074 blanks. Table 1 systematically compares the weak
075 supervision obtained by different frameworks.

076 To identify the words not contributing to the clas-
077 sification, we borrow the marginalization idea from
078 LLM reasoning (Holtzman et al., 2021). We get
079 the probability logit of each word in the raw text
080 by instructing LLMs (relatively small, specifically
081 Gemma (Mesnard et al., 2024)) to generate with or
082 without the in-class as a requirement. The differ-
083 ence between the two logits represents the potential
084 of each word to appear in the grafted text. As only
085 words with high potential will be left, we use the
086 average potential of top- $K\%$ words to represent
087 the text potential score. The bottom- $(100 - K)\%$
088 words will be masked to form the template for
089 data synthesis. We rank the templates by their po-
090 tential scores and select top- $T\%$ templates for the
091 last template-filling stage. Finally, these selected
092 templates are filled by prompting a state-of-the-art
093 LLM, GPT-4o (OpenAI, 2024).

094 We compare the three mentioned frameworks on
095 various raw corpora to classify different minority
096 classes. The experiment results show text grafting
097 can outperform state-of-the-art text mining and
098 dataset synthesis methods. The ablation study ver-
099 ifies that all stages and the intermediate template
100 contribute to the success of our proposed text graft-
101 ing. The mask-and-filling scenario also shows its
102 advantage over simple in-context generation, since
103 it forces the LLM to incorporate components from
104 the raw texts. We also involve an extreme situa-
105 tion where the target class does not appear in the
106 raw corpus completely. Remarkably, text grafting
107 shows its robustness to this extreme situation, indi-
108 cating its applicability does not require the target
109 class to appear in the raw corpus. This enables

110 text grafting to work on a very small corpus which
111 boosts efficiency.

112 Furthermore, we analyze and discuss the prop-
113 erty of text grafting. We apply principal component
114 analysis to visualize that the drafted texts are in-
115 deed near in-distribution. We also find the grafted
116 texts are near-distribution enough that we do not
117 need to synthesize negative samples as in tradi-
118 tional data synthesis, which reduces the cost. We
119 also conduct a comprehensive hyperparameter anal-
120 ysis of our method. Interestingly, we found that
121 The mask ratio is searched to be better set to a high
122 value like 0.75 and the mined template number can
123 be as small as 200. These case studies explore the
124 advantages of text grafting in distribution approxi-
125 mation and its failure when the raw texts are near
126 the distribution of LLM generation.

127 We summarize our contributions as follows,

- We propose a novel XWS-TC framework for mi-
128 nority classes, text grafting, combining the in-
129 distribution advantage of text mining and the in-
130 class advantage of data synthesis.
- We implement text grafting following the
131 marginalization idea from LLM reasoning, uti-
132 lizing the probability logits for template mining
133 and masking.
- We provide comprehensive analysis and case
134 studies to show the strength, property, and possi-
135 ble failure of text grafting.¹

139 2 Related Works

140 Extremely Weak-Supervised Text Classification
141 (XWS-TC) needs only minimal human guidance to
142 label the text, such as a few rules by human experts
143 that match the text to the labels (Wang et al., 2023).
144 Mainstream XWS-TC methods can be divided into
145 two categories: **Text Mining** and **Data Synthesis**.

146 **Text Mining** is a fundamental task (Han and
147 Kamber, 2000) for natural language processing.
148 In XWS-TC, the text miner follows high-level
149 rules from humans to annotate raw texts, which
150 are used to train the text classifier. A mainstream
151 rule is whether a seed word appears in the raw
152 text (Mekala and Shang, 2020; Meng et al., 2020;
153 Wang et al., 2021), categorized as seed methods.
154 Another mining way is to prompt language models
155 for logits that reflect the probability of texts falling
156 in classes (Brown et al., 2020), which can be cali-
157 brated by several techniques (Holtzman et al., 2021;

¹The datasets and models used in the experiments will be released for reproducibility.

158 Zhao et al., 2021; Han et al., 2023). The strong per-
 159 formance of existing text mining methods is highly
 160 dependent on the precision of the class-indicative
 161 rules (Dong et al., 2023a), which is hard to main-
 162 tain for minority classes.

163 **Data Synthesis** (He et al., 2022) addresses the
 164 precision degradation in text mining by directly
 165 prompting LLMs with the label names to generate
 166 in-class texts (Ye et al., 2022a; Peng and Shang,
 167 2024). With the powerful generative ability of
 168 LLMs, the synthesized texts are generally clean
 169 (in-class) for training strong classifiers. However,
 170 synthesized texts hold LLM-specific patterns, dis-
 171 covered by LLM-generated text detectors (Mitchell
 172 et al., 2023; Wu et al., 2023). This pattern is hard to
 173 be eliminated even with in-context learning (Koike
 174 et al., 2024). Thus, synthesized texts are generally
 175 out-of-domain and consequently fine-tune a weaker
 176 classifier on the test set.

177 **Minority Classes** widely appear in classifica-
 178 tion datasets as a result of long-tailed distribu-
 179 tion (Zhang et al., 2023; Henning et al., 2023). For
 180 minority classes with supervised annotations, tech-
 181 niques like re-sampling (Shen et al., 2016; Pouyan-
 182 far et al., 2018; Tepper et al., 2020) and data aug-
 183 mentation (Wei and Zou, 2019; Juuti et al., 2020;
 184 Tian et al., 2021; Chen et al., 2021). However,
 185 these methods are applied to unbalanced annota-
 186 tions, which are unavailable under XWS.

187 **Counterfactual Augmentation** refers to generat-
 188 ing annotated data out of the dataset or raw corpus.
 189 Different from regular augmentation, counterfactual
 190 augmentation changes the reference, e.g., la-
 191 bel flipping (Zhou et al., 2022; Peng et al., 2023).
 192 Counterfactual augmentation is also applied for
 193 text-to-text tasks like translation (Liu et al., 2021)
 194 or summarization (Rajagopal et al., 2022). Coun-
 195 terfactual augmentation shares the same require-
 196 ment for known reference as regular augmentation.
 197 This paper explores a counterfactual augmentation
 198 method for unannotated raw text under XWS.

199 3 Text Grafting

200 3.1 Preliminary

201 **XWS Minority Class Classification** takes a raw
 202 corpus $\mathcal{D} = \{X_{(i)}\}_{i=1:|\mathcal{D}|}$ and the target minority
 203 class name c as the input to train a binary classifier
 204 $f(X)$ that discerns a text falling in c or not. We
 205 denote the j -th word in the i -th text of the raw
 206 corpus as $x_{(i,j)}$.

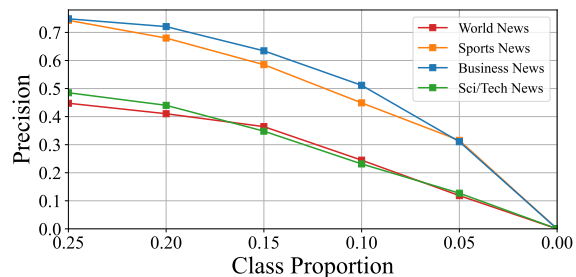


Figure 2: The precision of state-of-the-art text mining on same classes with different class proportions. “Precision” refers to the precision of the pseudo-labels. “Class Proportion” means the ratio of the texts of this class in the entire corpus after down-sampling.

207 **Text Mining** gathers in-class texts with high-
 208 level rules $g(X)$ that can precisely assign X to
 209 target class c . Example rules include whether X
 210 contains words indicating c (seed words) (Dong
 211 et al., 2023a) or X has top confidence to be in c
 212 by prompting LLMs (Brown et al., 2020) among
 213 \mathcal{D} . The mined $D^{(TM)} = \{X_{(i)}|g(X_{(i)})\}_{i=1:|\mathcal{D}|}$ is
 214 combined with some randomly sampled negative
 215 texts (due to the scarcity of c) to train $f(\cdot)$.

216 However, text miners fail in minority classes due
 217 to their low proportion in the raw corpus. By run-
 218 ning a state-of-the-art text mining method (Dong
 219 et al., 2023a) on AG-News (Zhang et al., 2015)
 220 with class name proportion modified by sampling,
 221 we observe the mining precision drops sharply with
 222 the decrease of proportion, presented in Figure 2.
 223 Another concern is the class might be too minor
 224 that even no ground truth can be mined from the
 225 raw corpus, limiting the precision to 0% no matter
 226 how intuitive the mining rule is.

227 **Data Synthesis** does not annotate raw texts for
 228 classifier fine-tuning but directly prompts LLMs to
 229 generate in-class texts ($X' \sim \text{LLM}(I_c)$), where I_c
 230 is an instruction to write a text in class c . With the
 231 strong capability of state-of-the-art LLMs (OpenAI,
 232 2024; Meta, 2024), the generated X' are highly
 233 confident to fall in class. Another advantage of
 234 data synthesis is the ability of LLMs to generate
 235 negative samples (Ye et al., 2022a; Peng and Shang,
 236 2024). However, synthesized texts consist of pat-
 237 terns different from other sources (Mitchell et al.,
 238 2023), which indicates classifiers $f(\cdot)$ fine-tuned by
 239 synthesized texts are out-of-domain, consequently
 240 weaker in the classification task.

241 3.2 Overview of Text Grafting

242 As depicted in Figure 3, our text grafting is a hybrid
 243 method that combines the strengths of text mining

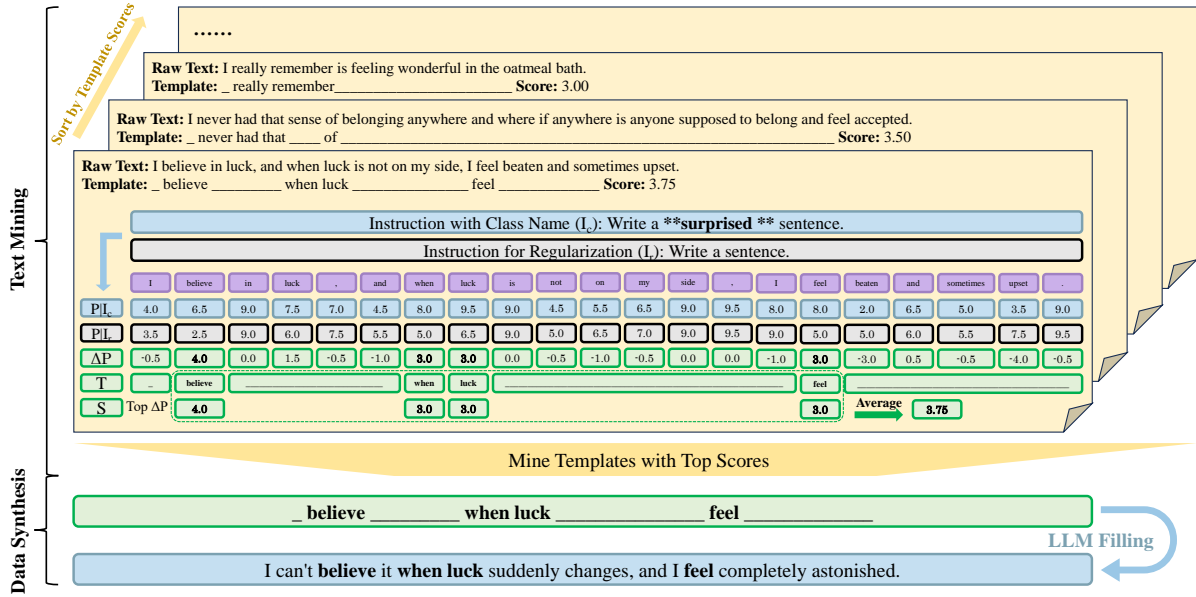


Figure 3: The overview of text grafting with the minority class “**Surprised**” in the Emotion dataset as an example. Text grafting includes two stages: 1) **Text (Template) Mining**: Create scored templates and select the ones with the top scores. 2) **Data Synthesis**: Prompt the LLM to fill in the templates to synthesize in-class texts.

and data synthesis. The core observation is that out-of-class texts can contain useful components for writing in-class texts. The text mining stage of text grafting aims to discover these potential components and formalize them as templates. In the data synthesis stage, the templates are filled by LLMs to produce in-class texts. With components from both raw texts and synthesis, the grafted texts are both in-class and near-distribution, which are supposed to fine-tune a better classifier than only text mining or data synthesis.

3.3 Implementation

In detail, the text mining stage includes **Potential Text Mining** and **Template Creation**, while in the data synthesis stage we conduct **Template Filling**. The text mining stage requires relatively small open-source LLMs with higher efficiency and accessible logits. Template Filling can utilize state-of-the-art LLMs even with API accessibility.

Potential Text Mining discovers texts with potential components to appear in the grafted texts. We evaluate the potential of each word $x_{(i,j)}$ in the raw text $X_{(i)}$ with regularized logits prompted from LLMs following the regularization idea in DC-PMI (Holtzman et al., 2021). The potential $\Delta p_{(i,j)}$ for $x_{(i,j)}$ is defined as the difference between the probability logit of $x_{(i,j)}$ prompted by an instruction with the class name (I_c) and an instruction for regularization (I_r). The difference can

also be viewed as the probability of $x_{(i,j)}$ raised by incorporating the class name c into the instruction.

$$\Delta p_{(i,j)} = \log P_{\text{LLM}}(x_{(i,j)}|I_c) - \log P_{\text{LLM}}(x_{(i,j)}|I_r) \quad (1)$$

The words with top- $K\%$ Δp among the words in text X_i will remain in the template. Thus, the average of their Δp represents the potential (ΔP_i) of the template created based on X_i . As we are mining potential templates rather than directly in-class texts, the mining rate $K\%$ can be much larger than text mining.

$$\Delta P_i = \left[\frac{1}{K\% \cdot |X_i|} \right] \sum_{\Delta p_i \in \text{Top-}K\%(\Delta p_{1:|X_i|})} \Delta p_i \quad (2)$$

Then the texts are ranked by their grafting potential ΔP and texts with top- $N\%$ potential are mined to create the templates.

Template Creation simply masks the words with bottom- $(100 - K)\%$ potential Δp by blank tokens “_” and uses the top- $K\%$ as template part. Text X_i is thus converted to template T_i , which is prepared for LLMs to fill in during the data synthesis stage. As the example in Figure 3, the components with the top potential to be in a grafted “Surprised” remain in the template such as “believe”, “when luck”, “feel”. These components support the data synthesis to better write an in-class text while keeping the style in distribution with the writing structure from the raw corpus.

Function	Prompt
TM (I_c)	“Please write a <label> <style>.”
TM (I_r)	“Please write a <style>.”
DS	“Fill in the blanks in the template to produce a <label> <style>.”

Table 2: The prompts used in text grafting. In prompts, <label> refers to the label names like “Surprised” while <style> represents the distribution like “Tweet”.

Template Filling prompts an LLM to fill in the blanks in T , which produces a grafted text that generally falls in the target class c . Referring to the example in Figure 3, the LLM well utilizes the writing structure in the template and fills in the blanks to produce the in-class text. As the template keeps the writing structure of the raw corpus, the grafted text is quite similar to the original one but flipped into the target minority class.

Specific prompts in these stages are shown in Table 2, where the label and distribution information is filled to support the text grafting.

4 Experiments

4.1 Evaluation

Datasets We take several minority classes from popular text classification datasets to evaluate the performance of different XWS-TC methods on minority classes. We include 1) TweetEval (Barbieri et al., 2020) and Emotion (Saravia et al., 2018), which contain minority emotion classes “Optimism” (8.9%) and “Surprised” (3.6%); 2) 20 News (Lang, 1995), which contains minority news topic “Religion” (3.3%) and “Politics” (4.1%); 3) BigPatent (Sharma et al., 2019), which contains minority patent class “Mechanical Engineering” (7.0%). The raw corpus is down-sampled to 10,000 samples to improve experiment efficiency and save budget costs. We use the F1 score as the metric for evaluation.

Baselines We include various text mining and data synthesis methods as the baselines for comparison to illustrate the advantage of our text grafting.

Text mining methods include,

- **Prompting Confidence** (Brown et al., 2020), which is a prompting method that directly queries an LLM whether the text falls in the target minority class, and uses the probability logit of answering “yes” for ranking. Considering the class minority, the mining rate is set to 1%.

- **Debiased Seed Word** (Dong et al., 2023a), which is the current state-of-the-art XWS-TC method. This method uses a seed word (the same as the label name) to match the target minority class and then drops the seed word from the context to eliminate spurious correlation. Then the texts are filtered by text selection (Mekala et al., 2022) to produce the final mined texts.

Data synthesis methods include,

- **ZeroGen** (Ye et al., 2022a), which directly prompts the LLM to synthesize texts in or out of the target minority class.
- **In-Context Generation** (Dong et al., 2023b), which uses raw texts as the in-context examples to generate texts with a similar writing style as the raw corpus.
- **Incubator** (Peng and Shang, 2024), which uses instruction-tuned LLMs and in-context learning based on annotated instruction-to-dataset samples to generate data points for fine-tuning.

All text synthesis methods synthesize 1000 texts as positive (in the target minority class) or negative samples (out of the target minority class, 2000 in total).

The LLM used for text mining is a popular and advanced open-source LLM, Gemma (Mesnard et al., 2024) (gemma-1.1-7b-it) with accessible possibility logits. The LLM used for data synthesis is the state-of-the-art LLM, GPT-4o (OpenAI, 2024).

Grafting Hyperparameters The mining rates of our text grafter are set to 25% ($K\%$) for potential components in templates and 10% ($N\%$) for potential templates. Thus, the synthesized data number is less than 1000, not more than the data number from pure data synthesis.

Fine-tuning Hyperparameters We fine-tune a RoBERTa-Large (Liu et al., 2019) as the classifier with the AdamW (Loshchilov and Hutter, 2019) as the optimizer whose learning rate is initialized to 1×10^{-5} . The classifier is fine-tuned by 10 epochs with batch size 8 and 20% training data are split for validation to select the best-performing checkpoint. All the experiment results are achieved by an average of 5 runs. The two stages in text grafting apply the same LLM as text mining and data synthesis.

4.2 Main Result

The main results from our experiments are presented in Table 3. The comparison inside text

Dataset		TWEET	PATENT	EMOTION	20NEWS		Average
Distribution		Tweet	Patent	Tweet	News		
Minority Class	Class Proportion	Optimism 8.9%	Mechanical 7.0%	Surprised 3.6%	Religion 3.3%	Politics 4.1%	
Supervised		45.88	34.30	32.28	24.10	32.27	35.14
Text Mining (TM)	Prompting Confidence	17.93	14.59	7.00	6.50	15.77	12.81
	Debiased Seed Word	19.15	20.46	8.78	11.47	19.53	15.88
Data Synthesis (DS)	ZeroGen	10.82	24.17	7.19	6.97	17.60	13.35
	Incubator	22.46	20.86	7.44	23.96	24.48	19.84
	In-Context Generation	16.24	24.53	22.24	21.98	24.13	21.83
TM+DS	Text Grafting (Ours)	32.70	25.42	27.46	25.32	27.32	27.64
Ablation	w/o Mining	26.54	16.74	24.32	17.69	15.16	20.09
	w/o Synthesis (DC-PMI)	17.86	11.34	7.34	4.33	4.28	9.03
	w/ Random Masking	30.11	19.07	23.37	23.57	26.65	24.55
	w/ MF \rightarrow ICG	21.31	20.58	15.33	23.60	25.06	21.18
Zero-Occur	Debiased Seed Word	0.00	17.66	5.88	8.79	20.73	10.61
	In-Context Generation	18.84	23.15	19.50	20.63	24.11	21.25
	Text Grafting (Ours)	30.61	25.27	31.08	26.15	25.54	27.73

Table 3: Text mining performance (F1 Score) for minority classes among different datasets.

Method	EMOTION	TNEWS
Language	English	Chinese
Debiased Seed Word	19.14	22.84
+ Text Grafting	31.30	28.61

Table 4: Results (Macro F1 Score) on end-to-end XWS-TC for different languages. Emotion (English) contains minority classes ‘‘Surprised’’ and ‘‘Love’’ while TNEWS (Chinese) has a minority class ‘‘Stock’’.

mining methods shows the advantage of the seed method over the prompt method, consistent with the findings of Wang et al.. The comparison among text synthesis methods reflects the importance of knowledge about the distribution of the corpus, as in-context generation outperforms other baselines with raw texts as an example for synthesis. Finally, text grafting outperforms all the baselines, which verifies the benefit of text grafting to produce in-class and near-distribution texts.

However, there is still a significant gap between the performance of supervised classification and XWS-TC even with text grafting. This indicates the grafted texts still have differences with the raw corpus distribution for further improvement.

4.3 Ablation Study

Table 3 also includes the ablation study results for text grafting in the *Ablation* columns. The first comparison focuses on the necessity of text mining and data grafting in the pipelines of text grafting. **Without Mining** removes the template score-based

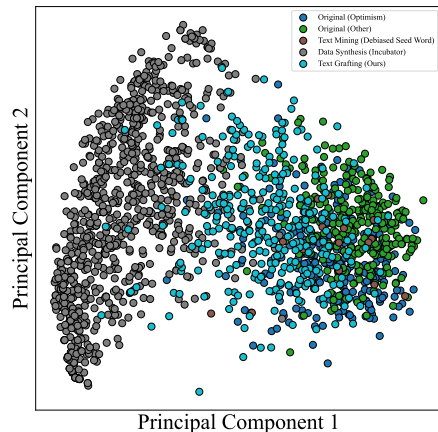


Figure 4: The visualization of text distributions from different methods.

sorting and lets the LLM fill in randomly selected templates, which significantly underperforms the initial grafting. **Without Synthesis** does not create templates for data synthesis, but directly uses the Δp averaged over all words to mine texts for fine-tuning, equal to DC-PMI (Holtzman et al., 2021). The result is similar to the Prompting Confidence method, which shows the limitation of text mining for minority classes. Then we emphasize the necessity of intermediate templates. **With Random Masking** randomly masks the mined texts instead of following the word-level potential Δp , which also results in a performance drop. **With Mask Filling \rightarrow In-Context Generation** takes the mined

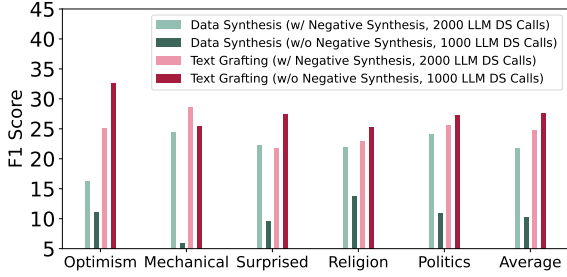


Figure 5: The analysis on the necessity of negative data synthesis.

423 texts as the in-context examples, which result in
 424 a similar performance as the one without mining,
 425 indicating the importance of template creation and
 426 filling. Based on these ablation results, our grafting
 427 framework is shown to be essential for achieving
 428 optimal performance by effectively combining data
 429 synthesis, text mining, and templates.

430 4.4 Further Analysis

431 **Q1: How does Text Grafting Benefit End-to-**
 432 **End XWS-TC?** Table 4 shows how text grafting
 433 can be integrated into end-to-end XWS-TC
 434 pipelines for different languages. We include the
 435 English Emotion dataset with “Surprised” and
 436 “Love” as the minority classes and the Chinese
 437 TNEWS dataset (Xu et al., 2020) with a minor-
 438 ity class “Stock”. For the minority classes, texts
 439 are synthesized by grafting while other classes
 440 apply the traditional debiased seed word method.
 441 The result shows text grafting improves end-to-end
 442 XWS-TC on different languages, which verifies
 443 the cross-lingual benefit of integrating text grafting
 444 into XWS-TC pipelines to handle minority classes.

445 **Q2: What if the class proportion is 0%?** In the
 446 *Zero-Occur* part of Table 3, we also include the dis-
 447 cussed extreme situation when the raw corpus does
 448 not contain any text falling in the target minority
 449 class. A dramatic drop appears in the performance
 450 of text mining as there is no ground truth that any
 451 miner can get. The data synthesis and text grafting
 452 methods are robust to this change as they do
 453 not require the existence of ground truth examples.
 454 Thus, text grafting is verified to be applicable to
 455 raw corpus without the target minority class. Thus,
 456 text grafting can be based on a small subset of the
 457 corpus which might not contain the target minority
 458 class to boost efficiency.

459 **Q3: How are grafted texts “near-distribution”?**
 460 In Figure 4, we apply semantic text embeddings

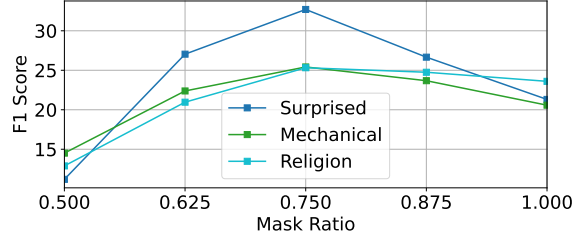


Figure 6: Analysis of the effect of mask ratio.

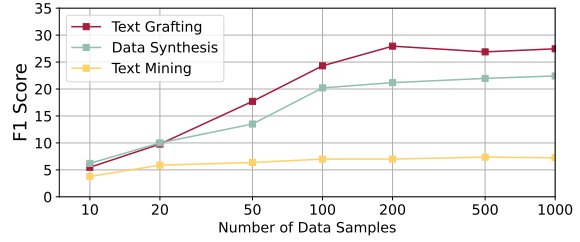


Figure 7: Analysis of the effect of data number.

(Gao et al., 2021) to represent the texts mined or
 461 synthesized by different methods. These embed-
 462 dings are then reduced to 2-dimension by principal
 463 component analysis (F.R.S., 1901) for visualiza-
 464 tion. We use the “Optimism” class of the TweetE-
 465 val benchmark and compare the most competitive
 466 methods (Debiased Seed Word, Incubator, Text
 467 Grafting) of different frameworks. We can observe
 468 that text mining only discovers a limited proportion
 469 of in-class texts. The synthesized texts fall into a
 470 very different domain from the raw corpus, which
 471 fine-tunes an out-of-domain classifier with limited
 472 generalizability. In contrast, the grafted texts are
 473 much more near-distribution, contributing to the
 474 performance of the fine-tuned classifier.
 475

476 **Q4: Is Negative Data Synthesis Necessary?** For
 477 data synthesis-based methods, the synthesis of neg-
 478 ative data is an essential stage in the pipeline, which
 479 doubles the calls for LLM to synthesize texts. In
 480 text grafting, we efficiently use the raw texts as the
 481 negative examples. Thus, we explore the necessity
 482 of negative synthesis by evaluating the performance
 483 of data synthesis (In-Context Generation) and text
 484 grafting with or without negative data synthesis
 485 with the results presented in Figure 5.

486 Based on the results, we observe negative data
 487 synthesis is very necessary to pure data synthesis as
 488 the performance drops dramatically by removing
 489 this stage. In contrast, text grafting without neg-
 490 ative data synthesis works even better, indicating
 491 that our text grafting can work more efficiently by

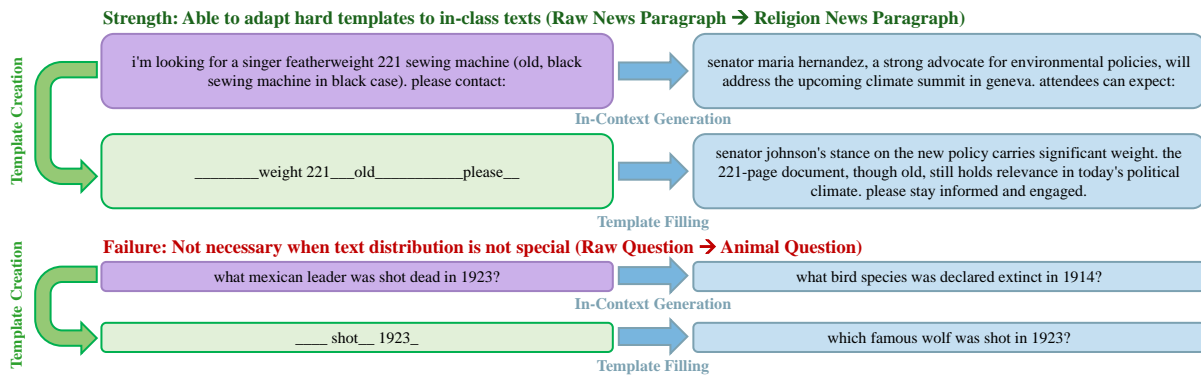


Figure 8: A case study on the strength and possible failure of text grafting.

reducing the effort to call LLM at double times. We attribute this efficiency to the near-distribution property of the grafted texts, which makes the discrimination between them and the original raw texts no longer degrade to the classifying of text sources (Mitchell et al., 2023).

Q5: What mask ratio to choose? In Figure 6, we analyze the mask ratio used in text grafting. Within the considered set of mask ratios, $\{0.5, 0.625, 0.75, 0.875, 1.0\}$, the best-performing ratio is 0.75 among different datasets, the same as the setup in our experiments. We can also observe a trend of performance decrease when the mask ratio becomes away from 0.75. This indicates a too-high masking ratio will make the synthesized text deviate from the domain of raw corpus (100% leads to in-context generation). On the other hand, a too-low mask ratio will limit the synthesizer to generate in-class texts, which might cause more severe performance drops.

Q6: How many templates to mine? In Figure 7, we further analyze the necessary number of templates to train a strong classifier, which can guide the efficient application of text grafting. The result of the “surprised” class shows about 200 samples can reach the best performance, which results in about \$0.2 budget for each class (OpenAI, 2024).

We also present how the efficiency of text mining (Debiased Seed Word) and data synthesis (In-Context Generation) is affected by sample numbers. Text mining cannot fine-tune a well-performing classifier due to severe noise in minority class mining. Data synthesis shows a similar scaling trend as text grafting but generally underperforms text grafting.

5 Case Study

In Figure 8, we depict workflows of text grafting in comparison with in-context generation to illustrate

the strength of grafting and possible failure.

Strength of text grafting is the ability of state-of-the-art LLMs to fill in hard templates as shown in the first case. While the template is not easy to be grafted into the target “Politics” class, the LLM comes up with the methodology to synthesize such a text. The text is also more similar in writing style to the original text than the in-context generation, which depicts the benefit from text grafting.

Failure of text grafting can happen when the corpus does not have a writing style very far from the way that LLMs can imitate. As shown in the second case, the LLM can synthesize the animal question without the intermediate template on the TREC corpus (Li and Roth, 2002), which reduces the necessity of text grafting. The XWS-TC of the minority class “Animal” on this corpus also shows a similar performance between data synthesis (F1 Score = 53.88) and text grafting (F1 Score = 53.46), which again emphasizes “near-distribution” to be an essential motivation to use text grafting.

6 Conclusion and Future Work

We introduced text grafting, a technique to generate in-distribution texts for minority classes using LLMs. By mining high-potential masked templates from the raw corpus and filling them with state-of-the-art LLMs, we achieve significant improvements in classifier performance on minority classes. Our analysis and case studies demonstrate the effectiveness of text grafting in enhancing text synthesis for minority classes. Future work will concentrate on improving the precision of template mining and the extension of text grafting to other tasks like information extraction.

564 Limitation

565 Despite the presented strengths in the paper, there
566 are still several limitations in the text grafting
567 pipeline. As a hybrid method, text grafting requires
568 a large raw corpus more than data synthesis and
569 LLM calls more than text mining. Other limitations
570 of text grafting also succeed from text mining and
571 data synthesis, such as the dependency on LLM
572 ability (for mining and synthesis). Thus, the appli-
573 cation scope for text grafting depends on how LLM
574 comprehends the class name semantics. The per-
575 formance of different classes might also be biased
576 to the LLM ability in different classes.

577 References

578 Francesco Barbieri, José Camacho-Collados, Luis Es-
579 pinosa Anke, and Leonardo Neves. 2020. [Tweeteval: Unified benchmark and comparative evaluation for tweet classification](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 1644–1650. Association for Computational Linguistics.

586 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
587 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
588 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
589 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
590 Gretchen Krueger, Tom Henighan, Rewon Child,
591 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
592 Clemens Winter, Christopher Hesse, Mark Chen, Eric
593 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
594 Jack Clark, Christopher Berner, Sam McCandlish,
595 Alec Radford, Ilya Sutskever, and Dario Amodei.
596 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

601 Junya Chen, Zidi Xiu, Benjamin Goldstein, Ricardo
602 Henao, Lawrence Carin, and Chenyang Tao. 2021.
603 [Supercharging imbalanced data learning with energy-based contrastive representation transfer](#). In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 21229–21243.

609 Chengyu Dong, Zihan Wang, and Jingbo Shang. 2023a.
610 [Debiasing made state-of-the-art: Revisiting the simple seed-based weak supervision for text classification](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 483–493. Association for Computational Linguistics.

617 Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong
618 Wu, Baobao Chang, Xu Sun, Jingjing Xu, Lei Li, and

Zhifang Sui. 2023b. [A survey on in-context learning](#). *Preprint*, arXiv:2301.00234.

Karl Pearson F.R.S. 1901. [Liii. on lines and planes of closest fit to systems of points in space](#). *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, 2(11):559–572.

Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. [Simcse: Simple contrastive learning of sentence embeddings](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 6894–6910. Association for Computational Linguistics.

Jiawei Han and Micheline Kamber. 2000. *Data Mining: Concepts and Techniques*. Morgan Kaufmann.

Zhixiong Han, Yaru Hao, Li Dong, Yutao Sun, and Furu Wei. 2023. [Prototypical calibration for few-shot learning of language models](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Xuanli He, Islam Nassar, Jamie Kiros, Gholamreza Haffari, and Mohammad Norouzi. 2022. [Generate, annotate, and learn: NLP with synthetic text](#). *Trans. Assoc. Comput. Linguistics*, 10:826–842.

Sophie Henning, William Beluch, Alexander Fraser, and Annemarie Friedrich. 2023. [A survey of methods for addressing class imbalance in deep-learning based natural language processing](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 523–540. Association for Computational Linguistics.

Ari Holtzman, Peter West, Vered Shwartz, Yejin Choi, and Luke Zettlemoyer. 2021. [Surface form competition: Why the highest probability answer isn’t always right](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 7038–7051. Association for Computational Linguistics.

Mika Juuti, Tommi Gröndahl, Adrian Flanagan, and N. Asokan. 2020. [A little goes a long way: Improving toxic language classification despite data scarcity](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020*, volume EMNLP 2020 of *Findings of ACL*, pages 2991–3009. Association for Computational Linguistics.

Ryuto Koike, Masahiro Kaneko, and Naoaki Okazaki. 2024. [OUTFOX: llm-generated essay detection through in-context learning with adversarially generated examples](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI*

676		2014, February 20-27, 2024, Vancouver, Canada, pages 21258–21266. AAAI Press.	
677			
678	Ken Lang. 1995. Newsweeder: Learning to filter net-news . In Armand Frieditis and Stuart Russell, editors, <i>Machine Learning Proceedings 1995</i> , pages 331–339. Morgan Kaufmann, San Francisco (CA).		
679			
680			
681			
682	Xin Li and Dan Roth. 2002. Learning question classifiers . In <i>COLING 2002: The 19th International Conference on Computational Linguistics</i> .		
683			
684			
685	Qi Liu, Matt J. Kusner, and Phil Blunsom. 2021. Counterfactual data augmentation for neural machine translation . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021</i> , pages 187–197. Association for Computational Linguistics.		
686			
687			
688			
689			
690			
691			
692			
693	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 . <i>CoRR</i> , abs/1907.11692.		
694			
695			
696			
697			
698	Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.		
699			
700			
701			
702			
703	Dheeraj Mekala, Chengyu Dong, and Jingbo Shang. 2022. LOPS: learning order inspired pseudo-label selection for weakly supervised text classification . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 4894–4908. Association for Computational Linguistics.		
704			
705			
706			
707			
708			
709			
710	Dheeraj Mekala and Jingbo Shang. 2020. Contextualized weak supervision for text classification . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 323–333. Association for Computational Linguistics.		
711			
712			
713			
714			
715			
716	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 <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 9006–9017. Association for Computational Linguistics.		
717			
718			
719			
720			
721			
722			
723			
724	Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer,	Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, and et al. 2024. Gemma: Open models based on gemini research and technology . <i>CoRR</i> , abs/2403.08295.	733 734 735 736 737 738 739 740 741
725			
726			
727			
728			
729			
730			
731			
732			
		Meta. 2024. Introducing meta llama 3: The most capable openly available llm to date .	742 743
		Eric Mitchell, Yoonho Lee, Alexander Khazatsky, Christopher D. Manning, and Chelsea Finn. 2023. Detectgpt: Zero-shot machine-generated text detection using probability curvature . In <i>International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA, volume 202 of Proceedings of Machine Learning Research</i> , pages 24950–24962. PMLR.	744 745 746 747 748 749 750 751
		OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> , abs/2303.08774.	752 753
		OpenAI. 2024. Hello gpt-4o . https://openai.com/index/hello-gpt-4o/ .	754 755
		Letian Peng and Jingbo Shang. 2024. Incubating text classifiers following user instruction with nothing but LLM . <i>CoRR</i> , abs/2404.10877.	756 757 758
		Letian Peng, Yuwei Zhang, and Jingbo Shang. 2023. Generating efficient training data via llm-based attribute manipulation . <i>CoRR</i> , abs/2307.07099.	759 760 761
		Samira Pouyanfar, Yudong Tao, Anup Mohan, Haiman Tian, Ahmed S. Kaseb, Kent Gauen, Ryan Dailey, Sarah Aghajanzadeh, Yung-Hsiang Lu, Shu-Ching Chen, and Mei-Ling Shyu. 2018. Dynamic sampling in convolutional neural networks for imbalanced data classification . In <i>IEEE 1st Conference on Multimedia Information Processing and Retrieval, MIPR 2018, Miami, FL, USA, April 10-12, 2018</i> , pages 112–117. IEEE.	762 763 764 765 766 767 768 769 770
		Dheeraj Rajagopal, Siamak Shakeri, Cícero Nogueira dos Santos, Eduard H. Hovy, and Chung-Ching Chang. 2022. Counterfactual data augmentation improves factuality of abstractive summarization . <i>CoRR</i> , abs/2205.12416.	771 772 773 774 775
		Elvis Saravia, Hsien-Chi Toby Liu, Yen-Hao Huang, Junlin Wu, and Yi-Shin Chen. 2018. CARER: Contextualized affect representations for emotion recognition . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3687–3697, Brussels, Belgium. Association for Computational Linguistics.	776 777 778 779 780 781 782
		Eva Sharma, Chen Li, and Lu Wang. 2019. BIG-PATENT: A large-scale dataset for abstractive and coherent summarization . In <i>Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers</i> , pages 2204–2213. Association for Computational Linguistics.	783 784 785 786 787 788 789

790	Jiaming Shen, Wenda Qiu, Yu Meng, Jingbo Shang,	Melanie Kambadur, Sharan Narang, Aurélien Ro-	850
791	Xiang Ren, and Jiawei Han. 2021. Taxoclass: Hi-	driguez, Robert Stojnic, Sergey Edunov, and Thomas	851
792	erarchical multi-label text classification using only	Scialom. 2023b. Llama 2: Open foundation and	852
793	class names . In <i>Proceedings of the 2021 Conference</i>	and fine-tuned chat models . <i>CoRR</i> , abs/2307.09288.	853
794	<i>of the North American Chapter of the Association</i>		
795	<i>for Computational Linguistics: Human Language</i>	Zihan Wang, Dheeraj Mekala, and Jingbo Shang. 2021.	854
796	<i>Technologies, NAACL-HLT 2021, Online, June 6-11,</i>	X-class: Text classification with extremely weak su-	855
797	<i>2021</i> , pages 4239–4249. Association for Computa-	pervision . In <i>Proceedings of the 2021 Conference</i>	856
798	tional Linguistics.	<i>of the North American Chapter of the Association</i>	857
		<i>for Computational Linguistics: Human Language</i>	858
		<i>Technologies, NAACL-HLT 2021, Online, June 6-11,</i>	859
799	Li Shen, Zhouchen Lin, and Qingming Huang. 2016.	<i>2021</i> , pages 3043–3053. Association for Computa-	860
800	Relay backpropagation for effective learning of deep	tional Linguistics.	861
801	convolutional neural networks . In <i>Computer Vision -</i>		
802	<i>ECCV 2016 - 14th European Conference, Amsterdam,</i>	Zihan Wang, Tianle Wang, Dheeraj Mekala, and Jingbo	862
803	<i>The Netherlands, October 11-14, 2016, Proceedings,</i>	Shang. 2023. A benchmark on extremely weakly su-	863
804	<i>Part VII</i> , volume 9911 of <i>Lecture Notes in Computer</i>	pervised text classification: Reconcile seed matching	864
805	<i>Science</i> , pages 467–482. Springer.	and prompting approaches . In <i>Findings of the As-</i>	865
		<i>sociation for Computational Linguistics: ACL 2023,</i>	866
		<i>Toronto, Canada, July 9-14, 2023</i> , pages 3944–3962.	867
		Association for Computational Linguistics.	868
806	Naama Tepper, Esther Goldbraich, Naama Zwerdling,	Jason W. Wei and Kai Zou. 2019. EDA: easy data	869
807	George Kour, Ateret Anaby-Tavor, and Boaz Carmeli.	augmentation techniques for boosting performance	870
808	2020. Balancing via generation for multi-class text	on text classification tasks . In <i>Proceedings of the</i>	871
809	classification improvement . In <i>Findings of the As-</i>	<i>2019 Conference on Empirical Methods in Natural</i>	872
810	<i>sociation for Computational Linguistics: EMNLP</i>	<i>Language Processing and the 9th International</i>	873
811	<i>2020, Online Event, 16-20 November 2020</i> , volume	<i>Joint Conference on Natural Language Processing,</i>	874
812	<i>EMNLP 2020 of Findings of ACL</i> , pages 1440–1452.	<i>EMNLP-IJCNLP 2019, Hong Kong, China, Novem-</i>	875
813	Association for Computational Linguistics.	<i>ber 3-7, 2019</i> , pages 6381–6387. Association for	876
		Computational Linguistics.	877
814	Jiachen Tian, Shizhan Chen, Xiaowang Zhang, Zhiyong	Junchao Wu, Shu Yang, Runzhe Zhan, Yulin Yuan,	878
815	Feng, Deyi Xiong, Shaojuan Wu, and Chunliu Dou.	Derek F. Wong, and Lidia S. Chao. 2023. A survey	879
816	2021. Re-embedding difficult samples via mutual in-	on llm-generated text detection: Necessity, methods,	880
817	formation constrained semantically oversampling for	and future directions . <i>CoRR</i> , abs/2310.14724.	881
818	imbalanced text classification . In <i>Proceedings of the</i>		
819	<i>2021 Conference on Empirical Methods in Natural</i>	Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao,	882
820	<i>Language Processing, EMNLP 2021, Virtual Event</i>	Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong	883
821	<i>/ Punta Cana, Dominican Republic, 7-11 November,</i>	Yu, Yin Tian, Qianqian Dong, Weitang Liu, Bo Shi,	884
822	<i>2021</i> , pages 3148–3161. Association for Computa-	Yiming Cui, Junyi Li, Jun Zeng, Rongzhao Wang,	885
823	tional Linguistics.	Wei Jian Xie, Yanting Li, Yina Patterson, Zuoyu Tian,	886
		Yiwen Zhang, He Zhou, Shaowei Hua Liu, Zhe Zhao,	887
		Qipeng Zhao, Cong Yue, Xinrui Zhang, Zhengliang	888
824	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	Yang, Kyle Richardson, and Zhenzhong Lan. 2020.	889
825	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	CLUE: A chinese language understanding evaluation	890
826	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	benchmark . In <i>Proceedings of the 28th International</i>	891
827	Azhar, Aurélien Rodriguez, Armand Joulin, Edouard	<i>Conference on Computational Linguistics, COLING</i>	892
828	Grave, and Guillaume Lample. 2023a. Llama: Open	<i>2020, Barcelona, Spain (Online), December 8-13,</i>	893
829	and efficient foundation language models . <i>CoRR</i> ,	<i>2020</i> , pages 4762–4772. International Committee on	894
830	abs/2302.13971.	Computational Linguistics.	895
831	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiang-	896
832	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	tao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong.	897
833	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	2022a. ZeroGen: Efficient zero-shot learning via	898
834	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-	dataset generation . In <i>Proceedings of the 2022 Con-</i>	899
835	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	<i>ference on Empirical Methods in Natural Language</i>	900
836	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	<i>Processing, EMNLP 2022, Abu Dhabi, United Arab</i>	901
837	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-	<i>Emirates, December 7-11, 2022</i> , pages 11653–11669.	902
838	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	Association for Computational Linguistics.	903
839	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,		
840	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	Jiacheng Ye, Jiahui Gao, Zhiyong Wu, Jiangtao Feng,	904
841	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	Tao Yu, and Lingpeng Kong. 2022b. Progen: Pro-	905
842	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	gressive zero-shot dataset generation via in-context	906
843	tinet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	feedback . In <i>Findings of the Association for Com-</i>	907
844	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-	<i>putational Linguistics: EMNLP 2022, Abu Dhabi,</i>	908
845	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,		
846	Ruan Silva, Eric Michael Smith, Ranjan Subrama-		
847	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-		
848	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,		
849	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,		

- 909 *United Arab Emirates, December 7-11, 2022*, pages
910 3671–3683. Association for Computational Linguistics.
911
- 912 Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. 2015.
913 [Character-level convolutional networks for text clas-](#)
914 [sification](#). In *Advances in Neural Information Pro-*
915 *cessing Systems 28: Annual Conference on Neural In-*
916 *formation Processing Systems 2015, December 7-12,*
917 *2015, Montreal, Quebec, Canada*, pages 649–657.
- 918 Yifan Zhang, Bingyi Kang, Bryan Hooi, Shuicheng Yan,
919 and Jiashi Feng. 2023. [Deep long-tailed learning:](#)
920 [A survey](#). *IEEE Trans. Pattern Anal. Mach. Intell.*,
921 45(9):10795–10816.
- 922 Xuandong Zhao, Siqi Ouyang, Zhiguo Yu, Ming Wu,
923 and Lei Li. 2023. [Pre-trained language models can](#)
924 [be fully zero-shot learners](#). In *Proceedings of the 61st*
925 *Annual Meeting of the Association for Computational*
926 *Linguistics (Volume 1: Long Papers), ACL 2023,*
927 *Toronto, Canada, July 9-14, 2023*, pages 15590–
928 15606. Association for Computational Linguistics.
- 929 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and
930 Sameer Singh. 2021. [Calibrate before use: Improv-](#)
931 [ing few-shot performance of language models](#). In
932 *Proceedings of the 38th International Conference on*
933 *Machine Learning, ICML 2021, 18-24 July 2021, Vir-*
934 *tual Event*, volume 139 of *Proceedings of Machine*
935 *Learning Research*, pages 12697–12706. PMLR.
- 936 Jing Zhou, Yanan Zheng, Jie Tang, Li Jian, and Zhilin
937 Yang. 2022. [Flipda: Effective and robust data aug-](#)
938 [mentation for few-shot learning](#). In *Proceedings of*
939 *the 60th Annual Meeting of the Association for Com-*
940 *putational Linguistics (Volume 1: Long Papers), ACL*
941 *2022, Dublin, Ireland, May 22-27, 2022*, pages 8646–
942 8665. Association for Computational Linguistics.