

# GPT3Mix: Leveraging Large-scale Language Models for Text Augmentation

Kang Min Yoo<sup>1,2</sup>, Dongju Park<sup>2</sup>, Jaewook Kang<sup>2</sup>,  
Sang-Woo Lee<sup>1,2</sup>, and Woomyeong Park<sup>2</sup>

<sup>1</sup>NAVER AI Lab

<sup>2</sup>NAVER Clova AI

{kangmin.yoo, dongju.park, jaewook.kang}@navercorp.com

{sang.woo.lee, max.park}@navercorp.com

## Abstract

Large-scale language models such as GPT-3 are excellent few-shot learners, allowing them to be controlled via natural text prompts. Recent studies report that prompt-based direct classification eliminates the need for fine-tuning but lacks data and inference scalability. This paper proposes a novel data augmentation technique that leverages large-scale language models to generate realistic text samples from a mixture of real samples. We also propose utilizing soft-labels predicted by the language models, effectively distilling knowledge from the large-scale language models and creating textual perturbations simultaneously. We perform data augmentation experiments on diverse classification tasks and show that our method hugely outperforms existing text augmentation methods. Ablation studies and a qualitative analysis provide more insights into our approach.

## 1 Introduction

In the seminal work by [Brown et al. \(2020\)](#), a large-scale language model, specifically GPT-3, has been shown to achieve superior performance on zero-shot and few-shot learning tasks by prompt-based in-context learning. In-context learning utilizes a prompt, which usually consists of a task description and few examples, to solve unseen tasks without the hefty price of fine-tuning. Recognizing the potential research applications of in-context learning and prompt-based control, a part of the NLP community has shifted its focus on understanding and devising advanced methods for optimizing prompt-based approaches ([Schick and Schütze, 2020a](#); [Shin et al., 2020](#); [Zhao et al., 2021](#); [Reynolds and McDonnell, 2021](#)).

However, these prompt-based approaches with inference on a large-scale language model suffer from several drawbacks. First, the number of in-context training examples is hard limited by the

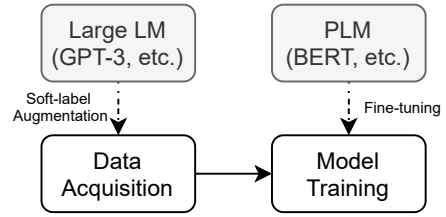


Figure 1: A conceptual diagram of text augmentation using large-scale language models.

maximum prompt length enabled by the inherent language model architecture. Second, prompt-based approaches require online inference on the expensive large-scale language models. The inference may not be scalable in real-world use cases, because it is slow and incurs huge memory overhead. Lastly, the prompt-based approaches do away with conventional machine learning techniques, making it mostly incompatible with existing established fine-tuning methods.

To overcome such limitations, we propose a more practical solution to utilize large-scale language models for downstream NLP tasks. In our proposed framework, as depicted in Figure 1, large-scale language models are not used as the pre-trained model for further domain-adaptive fine-tuning nor the backbone for prompt-based in-context learning but for imbuing the original training set with synthetic text data.

We propose GPT3Mix, a method for generating synthetic but hyper-realistic text samples from a mixture of real samples utilizing large-scale language models such as GPT-3<sup>1</sup>. GPT3Mix extracts few sample sentences from the task-specific training data, embed these samples in the prompt, and generates an augmented mixed sentence influenced by the sample sentences. GPT3Mix uses soft-

<sup>1</sup>Despite what the name suggests, we can apply GPT3Mix to any large-scale autoregressive language models.

labels predicted by the large-scale language model to transfer knowledge of probability as in knowledge distillation (Hinton et al., 2015). In short, our method achieves both (1) *data augmentation* via generating synthetic examples inspired by existing data samples and (2) *knowledge distillation* by training smaller classification models using soft-labels predicted by the large language model.

Our approach takes inspiration from the mix-based data augmentation methods in the vision domain (Zhang et al., 2017). Several mix-based data augmentation methods are suggested for NLP models. One of the notable methods is MixText (Chen et al., 2020), in which BERT is used to generate novel augmentation samples from interpolated embedding spaces. However, despite its great success in the vision domain, deep-mixing text augmentation methods have seen limited effectiveness in real-world cases due to the difficulty of interpolating language from latent spaces (Bowman et al., 2016). Synthetic language interpolated from a model’s hidden space such as the word embedding space of BERT may introduce noise, outweighing the benefit of novel sample discovery and causing deterioration in the training data distribution. Our work exploits the generative power of large-scale language models like GPT-3 to generate high-quality mixed samples from in-context examples.

We perform various data augmentation experiments on diverse classification tasks to verify our hypotheses and analyze our methodology. The contribution of our work is summarized as follows.<sup>2</sup>

1. We suggest employing prompt-based data augmentation using large-scale language models on top of the existing PLM fine-tuning paradigm to exploit the best of both worlds.
2. We propose GPT3Mix, a simple but effective text augmentation technique, that elicits knowledge and linguistic capability possessed by large-scale language models.
3. Our detailed analysis helps to understand the mechanism behind prompt-powered data augmentation, allowing us to better control the generation and the data augmentation behavior.

---

<sup>2</sup>The code to reproduce our results will be available at Github.

## 2 Related Work

**Knowledge Distillation** Knowledge distillation (Phuong and Lampert, 2019) is a technique that trains a smaller student classifier on the outputs of a larger teacher classifier. Knowledge distillation for language models in the context of model compression has been well-studied in the literature. There have been various distilled models and distillation methods proposed for pre-trained language models (Sanh et al., 2019; Tang et al., 2019). By utilizing soft-labels predicted by the large-scale language model, our approach helps to transfer knowledge to the downstream classifiers.

**Text Augmentation** Text augmentation refers to methods for perturbing the linguistic space without altering class labels to improve the robustness and generalizability of the downstream models. Data augmentation has been studied extensively in the NLP scene. Text augmentation in the current literature comes with two flavors: shallow and deep augmentation. The shallow data augmentation techniques inject locally plausible small noises into the linguistic space (words or phrases), in the hopes that the perturbations produce linguistically acceptable samples while maintaining label consistency. Two examples are EDA (Wei and Zou, 2019) and synonym replacement (Zhang et al., 2016).

Another class of augmentation techniques employs external language models to improve global coherence and consistency. The back-translation approach exploits semantic consistency in translation language pairs to generate novel paraphrases (Fadaee et al., 2017). In the more recent line of work, pre-trained language models, such as BERT (Devlin et al., 2019) or the sequence-to-sequence variant BART (Lewis et al., 2020), are used to obtain more diverse and linguistically correct augmentation samples. For example, BART has been proven to be effective in populating text samples for data-scarce labels (Kumar et al., 2020). Ng et al. (2020) proposed using masked language models as a denoising autoencoder to generate synthetic texts. Some other researchers have taken the direction of perturbing the latent spaces, optionally by introducing variational inference in the architecture (Xia et al., 2020b,a; Hou et al., 2018; Yoo et al., 2019).

On the other hand, inspired by the mix-up technique (Zhang et al., 2017) proposed for the vision domain, there have also been works to mix existing text samples to produce realistic augmen-

tation texts using statistical methods (Guo et al., 2020; Sun et al., 2020; Chen et al., 2020). Furthermore, pseudo-labeling, the act of annotating unlabeled data with model predictions (Lee et al., 2013; Reed et al., 2014), has been actively used in semi-supervised learning settings (Chen et al., 2020; Xie et al., 2020; Berthelot et al., 2019).

**Large-scale Language Models** Pre-trained transformer-based language models (Devlin et al., 2019; Lewis et al., 2020) have initiated a new paradigm in the NLP scene, changing the way we design NLP pipelines. With the recent development of mega-scale language models (Shoeybi et al., 2019; Brown et al., 2020), we are witnessing another shift in the paradigm, namely prompt-based NLP. These large language models are essentially few-shot learners, allowing them to be controlled through natural text. There has been a steep rise in the community’s interest to better understand the prompt-based mechanisms (Reynolds and McDonell, 2021; Schick and Schütze, 2020a; Shin et al., 2020; Jiang et al., 2020; Zhao et al., 2021). Our work relies on the previous findings on prompt-based manipulation.

To the best of our knowledge, this work is the first to propose using the prompt-based approach to generate synthetic samples from large-scale language models for the purpose of text augmentation.

### 3 GPT3Mix

Mixup (Zhang et al., 2017) is a simple learning technique that has been shown to be effective in preventing memorization and improving generalizability for the vision domain. The technique has been very effective on image data, but it has been harder to establish a standard approach for texts due to the inherent sparse nature of linguistic distributions, which attributes to the challenges of identifying adversarial text examples (Li et al., 2017). Inspired by the technique, we propose GPT3Mix as a powerful yet simple method to generate highly fluent synthetic samples based on a data distribution.

The proposed method (Figure 2) consists of three steps: (1) selecting examples from the dataset, (2) constructing a GPT3Mix prompt from the selected examples and meta-information about the dataset, and finally (3) extracting augmentation from the language model generation. This section provides details about each step as follows.

**Example Selection** For simplicity, we confine the downstream task to text classification tasks. Given a classification task  $\mathcal{T}$ , the training dataset  $\mathcal{D}$  is a set of text  $\mathbf{x}$  and associated label  $y$  pairs:  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \mid 1 \leq i \leq N\}$ .

We randomly choose  $k$  examples from  $\mathcal{D}$  to be anchors. Large-scale language models are known to be highly sensitive to the choice and the order of examples in the prompt (Reynolds and McDonell, 2021; Zhao et al., 2021). We conjecture that by carefully choosing the examples, we are able to control the generated augmentation samples from the language model. We conduct qualitative analysis on the augmentation samples to confirm our hypothesis (§4.4.5).

In our implementation, we simply used uniform distribution to choose  $k$  examples:  $p_s(i) = 1/N$ . Otherwise stated, most experiments are carried out by setting  $k = 2$  to simulate Mix-up. As found in our ablation studies (§4.4.1),  $k = 2$  provides a good trade-off between cost and performance.

**Prompt Construction** Given a set of prompt examples  $\mathcal{D}_e = \{(\mathbf{x}_i, y_i) \mid 1 \leq i \leq k\}$  sampled from  $\mathcal{D}$ , we formulate the prompt as follows.

A GPT3Mix prompt consists of a description header, an enumeration of text-label pairs of  $\mathcal{D}_e$ , and the augmentation prefix. An example of the prompt is shown in the appendix (§A). Our prompt has been designed carefully with the current literature findings of GPT-3 prompts (Reynolds and McDonell, 2021) in mind.

Specifically, the prompt follows the general template shown in the appendix, but has task-specific information to allow the large-scale language models to generalize better about the data distribution. These task indicators are unique to each task and provide meta-information of the task.

1. **Text Type  $T$ :** Meta-type of the input text  $\mathbf{x}$ . For example, in movie review sentiment analysis, the text type corresponds to `movie review`.
2. **Label Type  $L$ :** Meta type of the label class  $y$ . For the example above, the label type corresponds to `sentiment`.
3. **Label-token Verbalizer  $v : \mathcal{Y} \rightarrow \mathcal{V}$ :** Similar to the concept of verbalizers in the work of Schick and Schütze (2020b), the one-to-one mapping between the label classes  $y \in \mathcal{Y}$  and

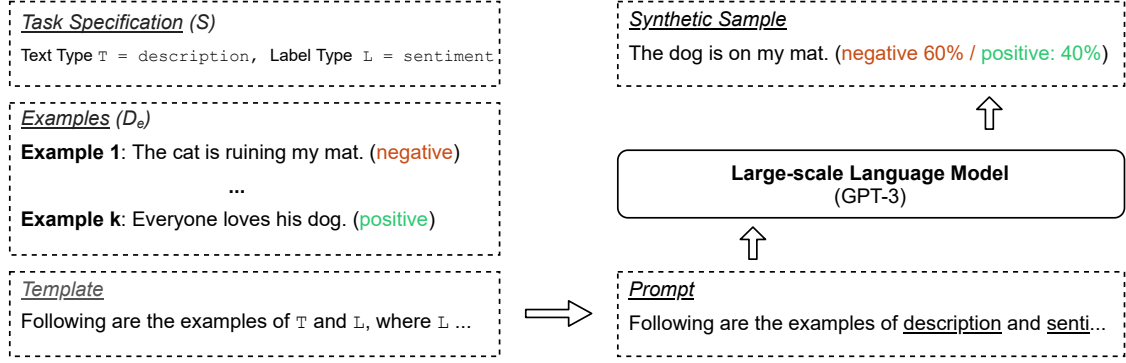


Figure 2: An illustration of GPT3Mix. The soft-labels of augmentation are extracted from the normalized label-token distributions predicted by the language model. Note that  $v$  has been omitted in the task specification  $S$  due to space limits.

word tokens in the language model’s vocabulary  $\mathcal{V}$ <sup>3</sup> is needed to formulate the prompt.

The triple of the meta information above forms the task specification  $S = (T, L, v)$ . Each task  $\mathcal{T}$  requires a task specification  $S_{\mathcal{T}}$  to be able to formulate a prompt for GPT3Mix. By default, the generic task specification  $S_{\text{generic}} = (\text{text}, \text{label}, I)$  is used to construct prompts, where  $I$  is the identity function assuming that the class label exists in the vocabulary  $\mathcal{V}$ .

**Augmentation Extraction** As with other prompt-based approaches, the augmentation text  $\mathbf{x}'$  and the label  $y'$  are generated in succession after the prompt as a natural text. A predefined prompt template in the examples signals the language model to generate  $(\mathbf{x}', y')$  with a structure, allowing us to extract respective values through pattern matching. Joint text and label generation also constraints the generated text to be associated with the correct label.

As illustrated in the prompt exhibit (§A), our particular prompt design ensures that the label token that corresponds to  $v(y')$  is generated after  $\mathbf{x}$ . This approach is inspired by the findings that, when inducing language models to come to a verdict, they require sufficient token lengths of “silent reasoning” prior to coming to a conclusion.

As large-scale language models are known to be few-shot learners (Brown et al., 2020), we also leverage GPT-3 to perform pseudo-labeling. The likelihood of generating the label-tokens is normalized to obtain the soft-label probability of the

augmentation text  $\mathbf{x}'$ . Concretely, the pseudo-label probability of an augmentation text  $\mathbf{x}'$  being labelled with label  $y'$  is as follows:

$$p(y' | \mathbf{x}') \propto p_{\text{LM}}(v_{\mathcal{T}}(y') | \mathcal{P}(\mathbf{x}', S_{\mathcal{T}})), \quad (1)$$

where  $p_{\text{LM}}$  is the language modeling likelihood and  $\mathcal{P} : \mathcal{S} \rightarrow \mathcal{X}$  is the function that constructs the prompt given a task specification.

Our approach effectively combines text perturbation, pseudo-labeling, and knowledge distillation in a single augmentation operation. In practice, augmentation samples with pseudo-labels are trained along with the real samples using the cross-entropy loss. This is in contrast to prior work, in which pseudo-labels are usually used for consistency regularization in the context of semi-supervised learning (Berthelot et al., 2019).

## 4 Experiments

We evaluate our augmentation approach on the following six classification benchmarks:

- **SST-2** (Socher et al., 2013): Stanford Sentiment Treebank is a sentiment classification dataset that contains movie reviews crawled from Rotten Tomatoes and their corresponding binary labels.
- **CR** (Hu and Liu, 2004): The Customer Review (CR) dataset is a set of Amazon product reviews and their binary labels.
- **COLA** (Warstadt et al., 2018): The Corpus of Linguistic Acceptability is a collection of sentences extracted from publications expertly annotated with grammaticality.

<sup>3</sup>In our implementation, we do not consider cases where a label class corresponds to multiple tokens. Regardless, expanding our work to incorporate multiple label tokens should be trivial.



		DistilBERT <sub>base</sub>				BERT <sub>base</sub>			
Dataset	Sub-sample	NoAug	EDA	BT	GPT3Mix	NoAug	EDA	BT	GPT3Mix
SST2	0.1%	55.8 <sub>5.1</sub>	54.4 <sub>4.4</sub>	55.7 <sub>5.2</sub>	<b>61.0</b> <sub>6.0</sub>	56.5 <sub>24.3</sub>	56.9 <sub>6.1</sub>	58.0 <sub>4.1</sub>	<b>63.3</b> <sub>4.0</sub>
	0.3%	64.9 <sub>8.0</sub>	67.7 <sub>7.9</sub>	64.4 <sub>7.5</sub>	<b>71.7</b> <sub>8.7</sub>	70.7 <sub>2.9</sub>	70.0 <sub>5.4</sub>	64.8 <sub>7.5</sub>	<b>78.0</b> <sub>5.6</sub>
	1.0%	77.9 <sub>3.6</sub>	74.4 <sub>5.0</sub>	77.3 <sub>2.8</sub>	<b>80.3</b> <sub>3.0</sub>	<b>82.5</b> <sub>1.8</sub>	79.7 <sub>1.1</sub>	81.1 <sub>3.5</sub>	82.0 <sub>2.7</sub>
COLA	0.1%	64.9 <sub>4.7</sub>	59.0 <sub>9.7</sub>	51.4 <sub>6.8</sub>	<b>68.5</b> <sub>0.2</sub>	60.9 <sub>6.7</sub>	57.5 <sub>5.0</sub>	52.8 <sub>4.2</sub>	<b>67.8</b> <sub>1.5</sub>
	0.3%	62.2 <sub>7.2</sub>	66.1 <sub>3.0</sub>	57.3 <sub>5.6</sub>	<b>68.5</b> <sub>0.2</sub>	65.5 <sub>3.0</sub>	63.5 <sub>5.6</sub>	63.5 <sub>4.8</sub>	<b>68.7</b> <sub>0.9</sub>
	1.0%	67.8 <sub>1.6</sub>	64.0 <sub>2.3</sub>	61.8 <sub>3.6</sub>	<b>69.2</b> <sub>0.5</sub>	67.1 <sub>4.8</sub>	66.5 <sub>5.8</sub>	66.7 <sub>2.7</sub>	<b>71.0</b> <sub>1.8</sub>
TREC6	0.1%	29.3 <sub>5.0</sub>	29.0 <sub>4.6</sub>	26.3 <sub>5.4</sub>	<b>37.6</b> <sub>6.9</sub>	<b>35.7</b> <sub>8.2</sub>	28.0 <sub>9.4</sub>	27.4 <sub>5.1</sub>	31.8 <sub>5.4</sub>
	0.3%	37.9 <sub>5.9</sub>	34.7 <sub>7.1</sub>	37.3 <sub>8.0</sub>	<b>40.1</b> <sub>4.0</sub>	37.9 <sub>9.2</sub>	35.8 <sub>4.1</sub>	<b>38.8</b> <sub>8.7</sub>	37.1 <sub>5.4</sub>
	1.0%	67.2 <sub>4.9</sub>	66.8 <sub>4.9</sub>	<b>68.2</b> <sub>9.5</sub>	65.8 <sub>3.4</sub>	67.6 <sub>8.2</sub>	69.6 <sub>11.5</sub>	<b>72.6</b> <sub>8.0</sub>	58.8 <sub>2.9</sub>
CR	0.1%	57.2 <sub>5.7</sub>	58.1 <sub>9.8</sub>	59.3 <sub>3.3</sub>	<b>60.5</b> <sub>1.1</sub>	56.9 <sub>6.4</sub>	56.3 <sub>5.7</sub>	56.9 <sub>6.1</sub>	<b>61.7</b> <sub>1.5</sub>
	0.3%	<b>66.6</b> <sub>7.8</sub>	63.0 <sub>4.7</sub>	60.9 <sub>2.9</sub>	65.5 <sub>5.1</sub>	65.7 <sub>6.0</sub>	64.4 <sub>3.8</sub>	63.7 <sub>2.9</sub>	<b>66.9</b> <sub>3.0</sub>
	1.0%	69.6 <sub>6.3</sub>	70.0 <sub>3.0</sub>	74.2 <sub>6.9</sub>	<b>79.5</b> <sub>2.5</sub>	77.1 <sub>3.4</sub>	74.4 <sub>3.3</sub>	76.2 <sub>5.1</sub>	<b>80.7</b> <sub>1.6</sub>
SUBJ	0.1%	<b>83.7</b> <sub>5.1</sub>	83.3 <sub>2.9</sub>	82.8 <sub>5.4</sub>	82.4 <sub>2.8</sub>	82.6 <sub>2.8</sub>	83.3 <sub>5.4</sub>	<b>84.3</b> <sub>3.6</sub>	83.8 <sub>2.9</sub>
	0.3%	88.5 <sub>0.9</sub>	<b>89.2</b> <sub>0.6</sub>	88.2 <sub>1.0</sub>	88.8 <sub>1.7</sub>	89.8 <sub>0.8</sub>	89.4 <sub>1.3</sub>	89.4 <sub>2.0</sub>	<b>89.8</b> <sub>1.1</sub>
	1.0%	90.7 <sub>1.0</sub>	90.7 <sub>1.1</sub>	90.5 <sub>0.8</sub>	<b>91.4</b> <sub>0.4</sub>	91.5 <sub>1.1</sub>	91.2 <sub>1.2</sub>	91.7 <sub>0.5</sub>	<b>92.2</b> <sub>0.5</sub>
MPQA	0.1%	<b>69.7</b> <sub>1.3</sub>	67.5 <sub>3.7</sub>	65.7 <sub>7.0</sub>	66.9 <sub>2.5</sub>	58.2 <sub>4.9</sub>	<b>68.0</b> <sub>3.5</sub>	59.4 <sub>5.1</sub>	61.7 <sub>6.1</sub>
	0.3%	76.6 <sub>5.0</sub>	<b>78.0</b> <sub>2.3</sub>	76.2 <sub>3.9</sub>	77.7 <sub>5.8</sub>	71.1 <sub>2.6</sub>	77.1 <sub>3.1</sub>	74.6 <sub>4.2</sub>	<b>77.7</b> <sub>4.9</sub>
	1.0%	83.7 <sub>4.8</sub>	82.3 <sub>2.0</sub>	81.9 <sub>4.7</sub>	<b>85.3</b> <sub>3.2</sub>	83.3 <sub>2.1</sub>	82.7 <sub>2.1</sub>	82.4 <sub>3.2</sub>	<b>86.1</b> <sub>0.8</sub>
Average	0.1%	60.1	58.6	56.9	<b>62.8</b>	58.5	58.3	56.5	<b>61.7</b>
	0.3%	66.1	64.1	64.1	<b>68.7</b>	66.8	64.6	65.8	<b>69.7</b>
	1.0%	76.1	74.7	75.6	<b>78.6</b>	78.2	77.3	78.4	<b>78.5</b>

Table 1: Main data augmentation results on 0.1%, 0.3%, and 1.0% training set sub-sample levels. We compare different augmentation strategies by transformer architectures on the downstream classification performance. Experiments have been repeated 10 times and the statistics are presented in the mean<sub>std</sub> format.

- **TREC6** (Voorhees and Tice, 1999): The TREC dataset is a dataset for question classification consisting of open-domain, fact-based questions divided into broad semantic categories.
- **MPQA** (Wiebe et al., 2005): Multi-perspective Question Answering polarity datasets consists of opinions and their semantic polarity.
- **SUBJ** (Pang and Lee, 2004): Subjectivity dataset contains movie reviews annotated with binary labels of objectiveness.

#### 4.1 Experimental Settings

To showcase our approach, we conduct downstream classification experiments on artificially data-scarce tasks by sub-sampling the training set. For each experiment, we perform a class-balanced sub-sample on the training set. We account for statistical variance in our experiments by fixating the sub-samples on 5 different data seeds and repeating the augmentation procedure and downstream classification experiments 10 times on the sub-samples. The data seeds were not screened and chosen by

random<sup>4</sup>.

For the classifier architecture, we use the base size BERT (Devlin et al., 2019) and DistilBERT (Sanh et al., 2019) models. For each downstream classification trial, we initialize the classifier model with the pre-trained parameters provided by the Huggingface Transformers library (Wolf et al., 2019) and randomly initialize the classifier layers, which consist of two fully connected layers that predict the class labels from the output embeddings of the transformer architectures. The classifiers are trained automatically by employing early stopping against the validation score with patience of 20 training epochs.

#### 4.2 Implementation Details

For selecting the optimal task specification for each task in GPT3Mix augmentation, we evaluated the performance of few handcrafted task specification candidates on the validation set and chose the highest performing one. The details about the optimal task specifications are presented in §B. The inference on GPT-3 was carried out via the OpenAI API Beta Access program. We used the largest GPT-3 model available on (davinci) unless other-

<sup>4</sup>The data seeds were generated using a master seed.

	$k$			
Sub.	1	2	4	8
0.1%	65.5 <sub>3.3</sub>	71.2 <sub>6.5</sub>	74.6 <sub>3.9</sub>	72.0 <sub>6.7</sub>
0.3%	78.9 <sub>3.9</sub>	80.0 <sub>2.7</sub>	80.2 <sub>2.1</sub>	80.0 <sub>1.6</sub>
1.0%	85.2 <sub>0.6</sub>	84.3 <sub>0.7</sub>	84.3 <sub>0.7</sub>	84.2 <sub>1.2</sub>

Table 2: An ablation study on the number of examples  $k$  in GPT3Mix prompts. When  $k = 1$ , GPT-3 produces point-wise perturbed samples. Experiments are carried out on the SST-2 dataset.

wise stated. On average, a GPT3Mix augmentation roughly consumes 300 tokens in combined length (prompt and generation).

The augmentation ratio between the training set and the synthetic set was set to 1 unless otherwise stated.

During classifier training, we used the Adam optimizer with decoupled weight decay (Kingma and Ba, 2014; Loshchilov and Hutter, 2017) and a learning rate of  $3e-5$ . The learning rate had a warm-up period of 3 epochs. PyTorch and M40 GPUs were used to run the experiments.

### 4.3 Data Augmentation Experiments

We compare our approach to Easy Data Augmentation (EDA) (Wei and Zou, 2019) and back-translation (BT) (Fadaee et al., 2017). For the back-translation baseline, texts were translated to and from German using Transformer architectures trained on the WMT16 English-German corpus provided by Fairseq (Ott et al., 2019).

The results on data-scarce text augmentation are presented in Table 1. First, we notice that, in most cases, our approach outperforms other augmentation baselines by a large margin. Also, our approach offers much more performance stability in terms of the variance of repeated trials and inter-task fluctuations than any other augmentation methods. Although back-translation and EDA do outperform GPT3Mix at certain configurations (e.g. CR 0.3%, TREC6 1.0%), GPT3Mix offers the most stable and consistent performance boost for the downstream classifiers across all tasks. This is more evident when we observe the average downstream classification accuracies of all tasks, in which GPT3Mix improves the classification accuracies by 2-3% on average, whereas EDA and BT perform worse than the baseline without aug-

	Model Size			
Sub.	ada	babbage	curie	davinci
0.1%	61.9 <sub>4.1</sub>	65.2 <sub>6.9</sub>	65.9 <sub>5.3</sub>	67.6 <sub>7.2</sub>
0.3%	74.6 <sub>4.8</sub>	69.7 <sub>7.3</sub>	74.6 <sub>4.5</sub>	78.3 <sub>2.9</sub>
1.0%	81.6 <sub>1.0</sub>	82.5 <sub>1.1</sub>	83.4 <sub>1.8</sub>	84.3 <sub>1.1</sub>

Table 3: An ablation study on the size of the language model with the SST-2 dataset. Larger language models provide greater augmentation benefits in data-limited environment.

mentation<sup>5</sup>.

Note that the results for GPT3Mix can be improved further by increasing the augmentation ratio beyond 1 (refer to §4.4). As our API budget increases over time, we anticipate achieving better results as well in the future.

**Full Dataset Experiments** We also perform full dataset data augmentation experiments to confirm that GPT3Mix still offers benefits even when task-specific data are abundant. We augmented the full SST-2 dataset with 1:1 ratio of synthetic samples from GPT3Mix, and our experiments on the downstream classification show that GPT3Mix improves the accuracy of DistilBERT<sub>base</sub> from 90.28% to 90.70% (0.42% improvement) and the accuracy of BERT<sub>base</sub> from 90.33% to 93.25% (2.92% improvement). A larger improvement margin is observed in the architecture with the greater expressive power, and this observation is in agreement with prior work (Zhang et al., 2017; Shafahi et al., 2019).

### 4.4 Ablation Studies

We conduct a number of ablation experiments to study the underlying mechanism of GPT3Mix. Note that the augmentation results for GPT3mix in the following ablation studies outperform the results presented in §4.3 due to ablation studies having higher augmentation ratios<sup>6</sup>.

#### 4.4.1 Number of Prompt Examples

First, the effect of the number of examples in GPT3Mix prompts ( $k$ ) on the downstream augmentation performance is studied. GPT3Mix requires

<sup>5</sup>We employed the well-established hyperparameters and the best hyperparameters proposed by the authors for EDA and BT. Augmentation ratios for all methods were set to 1 for a fair comparison with GPT3Mix.

<sup>6</sup>We limit the ratio in the main data augmentation experiments to the lowest level a due to the constrained API budget.

<b>Example 1</b>	Laughably, irredeemably awful. ( <i>negative</i> )
<b>Example 2</b>	Well-made but mush-hearted. ( <i>positive</i> )
<b>GPT3Mix</b>	Groundbreaking, disturbing. ( <b>positive</b> : 75%, <i>negative</i> : 25%)
<b>Example 1</b>	Berry’s saucy, full-bodied performance give this aging series a much needed kick, making “Die Another Day” one of the most entertaining Bonds in years. ( <i>positive</i> )
<b>Example 2</b>	Moonlight Mile doesn’t quite go the distance but the cast is impressive and they all give life to these broken characters who are trying to make their way through this tragedy. ( <i>positive</i> )
<b>GPT3Mix</b>	“Gosford Park” takes a leisurely approach to story telling and along the way it mixes intrigue, comedy, and romance. ( <b>positive</b> : 98%, <i>negative</i> : 2%)
<b>Example 1</b>	It’s just not very smart. ( <i>negative</i> )
<b>Example 2</b>	It’s quite an achievement to set and shoot a movie at the Cannes Film Festival and yet fail to capture its visual appeal or its atmosphere. ( <i>negative</i> )
<b>GPT3Mix</b>	Excessively talky, occasionally absurd and much too long, Munich is a fascinating mess. ( <i>positive</i> : 21%, <b>negative</b> : 79%)

Table 4: SST-2 augmentation samples from GPT3Mix (davinci).

$k \geq 2$  to effectively mix existing samples and generate interpolated text samples. However, supplying one example ( $k = 1$ ) per prompt and expecting GPT-3 to introduce perturbations or paraphrases of the given example can be a viable strategy. We vary  $k$  on the SST-2 dataset and observe the downstream performances (Table 2). The second-largest GPT-3 model (curie) was used and the augmentation multiplier was set to 10.

From the results, we notice that when the data availability is severely limited (i.e. 0.1% and 0.3%), point-wise perturbation doesn’t offer the performance improvement as much as when  $k \geq 2$ . However, as data becomes more abundant, increasing the number of mixing samples offers marginally small benefits for data augmentation. Yet, increasing the number of examples incurs additional overhead to the GPT-3 inference cost.

Generally, over-providing prompt examples may constraint the degrees of freedom and causing the synthetic samples to overfit on the data, hurting the downstream performances. However, a significant improvement from  $k = 2$  to  $k = 4$  is observed for the 0.1% sub-sample level. In our data augmentation studies, we weigh in on  $k = 2$  as a reasonable balance between the trade-off between GPT-3 inference costs and performance gains.

#### 4.4.2 Language Model Capacity

Next, we study the influence of the model capacity of the augmenting language model on the quality

<b>Sub.</b>	NoAug	Hard Labels	soft-labels
0.1%	55.8 <sub>5.1</sub>	61.6 <sub>8.0</sub>	71.2 <sub>6.5</sub>
0.3%	64.9 <sub>8.0</sub>	67.7 <sub>5.9</sub>	80.0 <sub>2.7</sub>
1.0%	77.9 <sub>3.6</sub>	79.0 <sub>2.8</sub>	84.3 <sub>0.7</sub>

Table 5: An ablation study on the employment of pseudo-labels. Hard labels are obtained from the beam search of the entire sequence autoregressively generated by the language model.

of augmentations. OpenAI offers GPT-3 in four different capacities: ada, babbage, curie, and davinci<sup>7</sup>, listed in the increasing order of model complexity. In this study, the augmentation ratio is set to 5.

As expected, the results (Table 3) show that having larger and more expressive language models benefit data augmentation.

#### 4.4.3 Task Specification

We are also interested in how the design choice of task specification for prompt construction affects the downstream performance. To analyze the effect, we compare the optimal task specification  $S_{\mathcal{T}^*}$  to a generic one ( $S_{\text{generic}}$ ). For this study, we used curie as the augmenting language model with an augmentation ratio of 3. The results in Table

<sup>7</sup>The sizes of the language models are known to be 2.7B, 6.7B, 13B, and 175B respectively; however, OpenAI has not officially disclosed the exact numbers yet.

Dataset	Sub.	NoAug	$S_{\text{generic}}$	$S_{\mathcal{T}^*}$
SST2	0.1%	55.8 <sub>5.1</sub>	60.1 <sub>5.2</sub>	71.2 <sub>6.5</sub>
	0.3%	64.9 <sub>8.0</sub>	72.6 <sub>5.7</sub>	80.0 <sub>2.7</sub>
	1.0%	77.9 <sub>3.6</sub>	81.4 <sub>1.7</sub>	84.3 <sub>0.7</sub>
COLA	0.1%	64.9 <sub>4.7</sub>	68.4 <sub>0.4</sub>	68.6 <sub>0.0</sub>
	0.3%	62.2 <sub>7.2</sub>	65.7 <sub>2.7</sub>	68.7 <sub>0.2</sub>
	1.0%	67.8 <sub>1.6</sub>	68.7 <sub>0.3</sub>	69.1 <sub>1.1</sub>

Table 6: An ablation study on task specifications.  $S_{\text{generic}}$  denotes a generic task specification that does not hold task-specific meta-information (§3), and  $S_{\mathcal{T}^*}$  denotes the optimal specification for the corresponding task.

6 support our conjecture that the language model utilizes the meta-information about the dataset to generate better data samples, and thus prompt designs have a significant impact on the augmentation quality. However, the generic task specification outperforms other augmentation baselines, highlighting the effectiveness of employing large-scale language models as the augmentation source.

#### 4.4.4 Pseudo-labeling

Finally, we study the effect of employing pseudo-labels from the label token probabilities predicted by the large-scale language model. we compare the augmentation performance when the label tokens optimized from the sequence-wide beam search are used instead. Results on SST-2 (Table 5) show that employing soft-labels has a strong advantage over sequence-optimized labels. The performance gap between the hard and soft-labels can be considered as the benefit of utilizing the class distribution jointly predicted by the language model as a form of knowledge distillation for synthetic samples (Kim and Rush, 2016). *curie* was used as the GPT-3 model with the augmentation ratio of 5.

#### 4.4.5 Qualitative Analysis

Language models are known to be sensitive to the selection and the order of the examples presented in the prompt, causing the proceeding generation to exhibit bias (Zhao et al., 2021; Reynolds and McDonnell, 2021). Our proposed method hinges on this unique property of large-scale language models, hence we intend to qualitatively examine the augmentation samples to support our hypothesis.

The augmentation samples for the SST-2 dataset are presented in Table 4. First, we notice that the sentiment of the synthetic sample strongly depends

on the class labels of the examples. When both examples are either all *positive* or all *negative*, the sentiment of the augmentation sample is heavily influenced to have the corresponding bias. In the case of mixed class labels, the resulting class distribution is less peaky. Second, we also discover that the augmented sample follows the similar syntactic and semantic structure of the example texts. As demonstrated in the first case, the short and phrasal structure of the examples is well translated into the generated sample, supporting our claim that language models are able to learn from in-context examples, even for generation and pseudo-labeling tasks. Similar patterns are observed in the subsequent cases, in which even the entities and subjects of the examples are perturbed and carried over into the synthetic samples.

## 5 Conclusion

In this paper, we proposed a novel text augmentation technique called GPT3Mix that leverages large-scale language models and their abilities to perform controlled generation via prompts. Our extensive experiments on various classification tasks show that the augmentation samples can improve robustness and the classification performance of pre-trained transformer models, and the experiments suggest that GPT3-based data augmentation can become a competitive alternative to prompt-based task-solving (Brown et al., 2020) or direct fine-tuning (Liu et al., 2021). As future work, we are interested in achieving state-of-the-art results using existing pre-trained transformer architectures solely via GPT3Mix-enabled data augmentation. We are also working to improve data augmentation efficiency by example selection optimization and prompt construction optimization.

## References

- David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. 2019. Mixmatch: A holistic approach to semi-supervised learning. *arXiv preprint arXiv:1905.02249*.
- Samuel Bowman, Luke Vilnis, Oriol Vinyals, Andrew Dai, Rafal Jozefowicz, and Samy Bengio. 2016. Generating sentences from a continuous space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 10–21.
- 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*.
- Jiaao Chen, Zichao Yang, and Diyi Yang. 2020. Mix-text: Linguistically-informed interpolation of hidden space for semi-supervised text classification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2147–2157.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Marzieh Fadaee, Arianna Bisazza, and Christof Monz. 2017. Data augmentation for low-resource neural machine translation. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 567–573.
- Demi Guo, Yoon Kim, and Alexander M Rush. 2020. Sequence-level mixed sample data augmentation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5547–5552.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. *stat*, 1050:9.
- Yutai Hou, Yijia Liu, Wanxiang Che, and Ting Liu. 2018. Sequence-to-sequence data augmentation for dialogue language understanding. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1234–1245.
- Minqing Hu and Bing Liu. 2004. Mining and summarizing customer reviews. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 168–177.
- Zhengbao Jiang, Frank F Xu, Jun Araki, and Graham Neubig. 2020. How can we know what language models know? *Transactions of the Association for Computational Linguistics*, 8:423–438.
- Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1317–1327.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2020. Data augmentation using pre-trained transformer models. In *Proceedings of the 2nd Workshop on Life-long Learning for Spoken Language Systems*, pages 18–26.
- Dong-Hyun Lee et al. 2013. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.
- Yitong Li, Trevor Cohn, and Timothy Baldwin. 2017. Robust training under linguistic adversity. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 21–27.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. 2021. Gpt understands, too. *arXiv preprint arXiv:2103.10385*.
- Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.
- Nathan Ng, Kyunghyun Cho, and Marzyeh Ghassemi. 2020. Ssmba: Self-supervised manifold based data augmentation for improving out-of-domain robustness. *arXiv preprint arXiv:2009.10195*.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of NAACL-HLT 2019: Demonstrations*.
- Bo Pang and Lillian Lee. 2004. A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the 42nd Annual Meeting of the Association for Computational Linguistics (ACL-04)*, pages 271–278.
- Mary Phuong and Christoph Lampert. 2019. Towards understanding knowledge distillation. In *International Conference on Machine Learning*, pages 5142–5151. PMLR.
- Scott Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. 2014. Training deep neural networks on noisy labels with bootstrapping. *arXiv preprint arXiv:1412.6596*.
- Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. *arXiv preprint arXiv:2102.07350*.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.

- 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*.
- Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. 2019. Adversarial training for free! *arXiv preprint arXiv:1904.12843*.
- 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*.
- Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. 2019. Megatron-lm: Training multi-billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. 2013. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1631–1642.
- Lichao Sun, Congying Xia, Wenpeng Yin, Tingting Liang, S Yu Philip, and Lifang He. 2020. Mixup-transformer: Dynamic data augmentation for nlp tasks. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3436–3440.
- Raphael Tang, Yao Lu, Linqing Liu, Lili Mou, Olga Vechtomova, and Jimmy Lin. 2019. Distilling task-specific knowledge from bert into simple neural networks. *arXiv preprint arXiv:1903.12136*.
- Ellen M Voorhees and Dawn M Tice. 1999. The trec-8 question answering track evaluation. In *TREC*, volume 1999, page 82. Citeseer.
- Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. 2018. Neural network acceptability judgments. *arXiv preprint arXiv:1805.12471*.
- Jason Wei and Kai Zou. 2019. Eda: Easy data augmentation techniques for boosting performance on text classification tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6383–6389.
- Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2):165–210.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface’s transformers: State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*.
- Congying Xia, Caiming Xiong, Philip Yu, and Richard Socher. 2020a. Composed variational natural language generation for few-shot intents. *arXiv preprint arXiv:2009.10056*.
- Congying Xia, Chenwei Zhang, Hoang Nguyen, Jiawei Zhang, and Philip Yu. 2020b. Cg-bert: Conditional text generation with bert for generalized few-shot intent detection. *arXiv preprint arXiv:2004.01881*.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Thang Luong, and Quoc Le. 2020. Unsupervised data augmentation for consistency training. *Advances in Neural Information Processing Systems*, 33.
- Kang Min Yoo, Youhyun Shin, and Sang-goo Lee. 2019. Data augmentation for spoken language understanding via joint variational generation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 7402–7409.
- Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. 2017. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2016. [Character-level convolutional networks for text classification](#).
- 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*.

Dataset	$T$	$L$	$v$
Generic	text	label	$\cdot \rightarrow \cdot$
SST-2	movie review	sentiment	pos $\rightarrow$ positive, neg $\rightarrow$ negative
CR	customer review	sentiment	pos $\rightarrow$ positive, neg $\rightarrow$ negative
SUBJ	text	objective	subjective $\rightarrow$ no, objective $\rightarrow$ yes
COLA	text	grammar	acceptable $\rightarrow$ correct, unacceptable $\rightarrow$ incorrect
TREC6	question	type	ABBR $\rightarrow$ abbreviation, LOC $\rightarrow$ location, DESC $\rightarrow$ description, NUM $\rightarrow$ numeric ENTY $\rightarrow$ entity, HUM $\rightarrow$ human
MPQA	text	sentiment	pos $\rightarrow$ positive, neg $\rightarrow$ negative

Table 7: Optimal task specifications.

## A Prompts

The GPT3Mix prompt uses the following template. The template corresponds to the prompt-constructing function  $\mathcal{P}$ , which require a task specification  $S_{\mathcal{T}} = (T, L, v)$ .

Each item in the following list contains a `<text type>` and the respective `<label type>`. `<label type>` is one of '`<label token 1>`', ..., or '`<label token N>`'.

```
<text type>: <example text 1> (<label type>: <example label 1>)
...
<text type>: <example text k> (<label type>: <example label k>)
<text type>:
```

For example, given  $S_{\text{SST2}} = (\text{movie review}, \text{sentiment}, I)$ , the constructed GPT3Mix prompt is as follows.

Each item in the following list contains a movie review and the respective sentiment. The sentiment is one of 'positive' or 'negative'.

```
Movie review: Despite its Hawaiian setting, the science-fiction
trimmings and some moments of rowdy slapstick, the basic plot of
``Lilo`` could have been pulled from a tear-stained vintage Shirley
Temple script. (Sentiment: Negative)
Movie review: And people make fun of me for liking Showgirls.
(Sentiment: Negative)
Movie review:
```

## B Task Specifications

After validating candidate task specifications for each task, we have selected the following for conducting our experiments (Table 7).

Providing incorrect or suboptimal specifications to the prompt may cause a large drop in augmentation qualities. For example, in the case of designing task specifications for the COLA dataset, when “linguistic acceptability” is used as the label type (instead of the optimal “grammar”), the downstream performance on the 0.1% sub-dataset drops to 38.8%, resulting in performance worse than the non-augmented baseline of 68.80%.