
Generating Training Data with Language Models: Towards Zero-Shot Language Understanding

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Pretrained language models (PLMs) have demonstrated remarkable performance
2 in various natural language processing tasks: Unidirectional PLMs (*e.g.*, GPT) are
3 well known for their superior text generation capabilities; bidirectional PLMs (*e.g.*,
4 BERT) have been the prominent choice for natural language understanding (NLU)
5 tasks. While both types of models have achieved promising few-shot learning
6 performance, their potential for zero-shot learning has been underexplored. In this
7 paper, we present a simple approach that uses both types of PLMs for fully zero-shot
8 learning of NLU tasks without requiring any task-specific data: A unidirectional
9 PLM generates class-conditioned texts guided by prompts, which are used as
10 the training data for fine-tuning a bidirectional PLM. With quality training data
11 selected based on the generation probability and regularization techniques (label
12 smoothing and temporal ensembling) applied to the fine-tuning stage for better
13 generalization and stability, our approach demonstrates strong performance across
14 seven classification tasks of the GLUE benchmark (*e.g.*, 72.3/73.8 on MNLI-m/mm
15 and 92.8 on SST-2), significantly outperforming zero-shot prompting methods and
16 achieving even comparable results to strong few-shot approaches using 32 training
17 samples per class¹.

18 **1 Introduction**

19 Pretrained language models (PLMs) [5, 8, 11, 19, 33, 37] have achieved human-level performance
20 on natural language understanding (NLU) tasks [62, 63] when fine-tuned on a large amount of
21 task-specific training data. However, such a supervised fine-tuning paradigm is drastically different
22 from how humans perform these tasks: We barely need to see many task-specific training samples
23 to perform well. Recently, many studies have revealed the intriguing few-shot learning potential of
24 PLMs: By converting task descriptions to natural language prompts and injecting them into PLMs,
25 prompt-based approaches [5, 13, 51, 52, 55] leverage task-specific information for better training
26 data efficiency and have achieved remarkable few-shot results.

27 When prompt-based methods are applied to the zero-shot setting, however, the PLMs' predictions
28 are much less accurate. For example, GPT-3's zero-shot performance is much degraded relative to
29 its few-shot performance [5], especially on challenging tasks like natural language inference (NLI).
30 Without any task-specific samples, it is indeed challenging for PLMs to effectively interpret the
31 prompts that come in different formats and are unseen in the pretraining data. To familiarize PLMs
32 with various prompts for zero-shot generalization to unseen tasks, a recent study proposes instruction
33 tuning [66], which fine-tunes PLMs on a large collection of different tasks described by instructions.
34 Despite its strong performance, its success is grounded in the large number of cross-task annotated
35 datasets (*e.g.*, train on many non-NLI tasks and transfer to NLI tasks) and the gigantic model size
36 (*e.g.*, hundreds of billions of parameters), posing great challenges for training and using them.

¹Code is shared in the supplementary material.

37 In this work, we study zero-shot learning of PLMs on NLU tasks without any task-specific or cross-
38 task data. Motivated by the strong text generation power of recent PLMs [5, 23, 29, 48], we propose
39 SuperGen, a **Sup**er**vi**sion **Ge**neration approach, wherein training data are created via a unidirectional
40 PLM (*i.e.*, the generator) which generates class-conditioned texts guided by label-descriptive prompts.
41 A bidirectional PLM (*i.e.*, the classifier) is then fine-tuned on the generated texts to perform the
42 corresponding task. Both PLMs can be of moderate size to fit in typical research hardware (*e.g.*, a
43 GPT-2-sized [47] generator and a RoBERTa_{Large}-sized [33] classifier). With supervision automatically
44 created by the generator, SuperGen eliminates the need for task-specific annotations and provides
45 the classifier PLM with a larger amount of training data than in few-shot scenarios. SuperGen is
46 compatible with any PLM as the classifier and any fine-tuning method. We note that the generator
47 creates synthetic samples in a zero-shot manner, and the classifier is fine-tuned on the synthetic data
48 (the classifier is thus not zero-shot, but there are no task-specific data required in such a process).

49 Across seven classification tasks of the GLUE benchmark [62], SuperGen significantly outperforms
50 the prompt-based zero-shot method and even achieves an overall better result in both average
51 performance and stability than strong few-shot approaches that use 32 annotated samples per class.
52 We identify several key factors to the strong performance of SuperGen through ablation studies: (1)
53 selecting quality training data based on their generated probability, and (2) using label smoothing and
54 temporal ensembling to regularize fine-tuning on generated data.

55 2 Related Work

56 2.1 Few-Shot and Zero-Shot Learning with PLMs

57 Instead of using a large amount of annotated training data for fine-tuning PLMs on downstream tasks,
58 few-shot learning studies how to better leverage only a small amount of task-specific training data,
59 a more realistic scenario in many applications. The most strict few-shot learning setting does not
60 assume access to any unlabeled data or large validation sets for hyperparameter tuning [44], where
61 prompt-based methods [5, 13, 32, 34, 51–53, 55, 59, 80] are prominently deployed to inject task
62 descriptions into PLMs and make effective use of their language modeling capability for improved
63 training data efficiency in low-data regimes. More broadly, semi-supervised learning additionally
64 leverages unlabeled task-specific data, where data augmentation [7, 69], regularization [39] and
65 bootstrapping [52] methods are commonly used.

66 Zero-shot learning, on the other hand, is a much more challenging setting with absolutely no access
67 to any task-specific data. When prompt-based methods are directly used to obtain predictions from
68 PLMs without any training, their zero-shot performance can be much worse [5, 13]—difficult NLU
69 tasks can be barely formulated as prompts that resemble the format of pretraining data, posing great
70 challenges for PLMs to accurately interpret and leverage the prompts without given any training
71 samples. The current mainstream of zero-shot learning is based on transfer learning: By converting
72 a set of tasks with abundant annotations into instruction templates [38, 50, 66, 70], entailment
73 pairs [75, 76] or question-answer formats [46, 82] and fine-tuning PLMs on them, the PLMs acquire
74 the cross-task transfer ability [74] to execute unseen tasks when they are formulated in a similar
75 format. Our work proposes a different approach from these studies: We use a unidirectional PLM to
76 generate training data for fine-tuning another PLM on the target task. This not only removes the need
77 for a large amount of cross-task annotations, but also eliminates the task difference in training and
78 inference. Moreover, different from previous studies [1, 72] that rely on labeled data to fine-tune the
79 generative PLM, we directly use prompts to guide data generation without fine-tuning.

80 2.2 Controlled Text Generation with PLMs

81 Controlled text generation [22] aims to steer the generated texts of language models towards desired
82 contents, styles or domains. Through fine-tuning PLMs on attribute-specific data, high-level control
83 (*e.g.*, generating certain topics or sentiments [84]), fine-grained control (*e.g.*, generating specific words
84 or phrases [6]) or both [24] can be achieved. Adapting PLMs to generate texts of specific attributes can
85 also be realized at inference time without any further training of the PLMs [10, 25, 26, 31, 43, 64, 71].
86 Different text attributes can also be represented during pretraining time as control codes [23] which
87 later can serve as explicit guidance for generating domain/attribute-specific texts.

88 Along another line of controlling text generation, the idea of using prompts as guidance has emerged
89 recently—Since natural language generation is largely based on contexts, providing certain prompts
90 as the beginning of a sequence can effectively steer the subsequent texts to be generated. The prompts

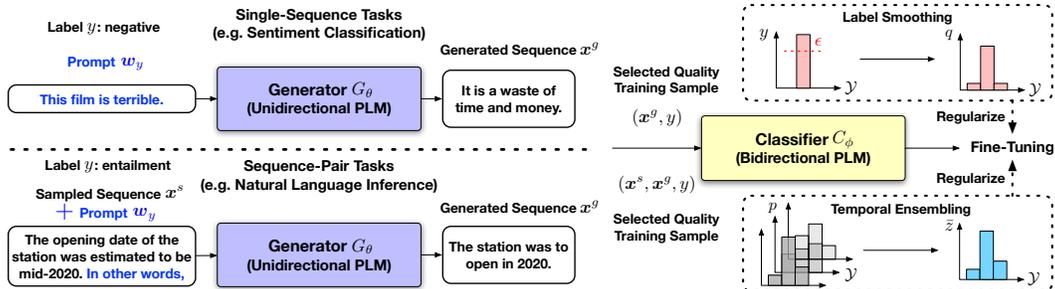


Figure 1: Overview of SuperGen for zero-shot learning of NLU tasks. A unidirectional PLM is used as the generator for creating training data guided by label-descriptive prompts. Quality training samples are selected based on average log generation probability. A bidirectional PLM is fine-tuned on the selected training set with label smoothing and temporal ensembling as regularization to perform the classification task.

91 can be either in natural language [53] or as learnable parameters [30]. In this work, we also guide
 92 text generation via prompts, but for the novel purpose of creating training data for NLU tasks. There
 93 have been studies with similar purposes, such as generating similar/dissimilar sentences for training
 94 sentence embeddings [54] and using labeled samples as demonstrations to prompt gigantic PLMs [77]
 95 for creating novel training data. In this work, we explore generating training data *without using any*
 96 *labeled samples* for a wide range of different NLU tasks. The similar setting is also explored in a
 97 concurrent study [73]. Compared to annotated task-specific data, the generated texts may contain
 98 noise and have domain difference from the downstream task. We will introduce several important
 99 strategies for effective fine-tuning of PLMs on generated data.

100 3 Method

101 3.1 Preliminaries

102 **Problem Formulation.** We consider solving a classification problem² where we are only given the
 103 label space \mathcal{Y} and a mapping $\mathcal{M} : \mathcal{Y} \rightarrow \mathcal{W}$ that converts each label $y \in \mathcal{Y}$ into a label-descriptive
 104 prompt (*i.e.*, a short phrase) $w_y \in \mathcal{W}$. We assume access to a unidirectional PLM G_θ as the generator
 105 and a bidirectional PLM C_ϕ which will be fine-tuned as the classifier³. We also assume the pretraining
 106 corpus \mathcal{D} (*e.g.*, Wikipedia) is available.

Text Generation with Unidirectional PLMs. A unidirectional PLM G_θ is pretrained to maximize the generation probability of each token in a sequence $\mathbf{x} = [x_1, x_2, \dots, x_n]$ conditioned on previous tokens:

$$\max_{\theta} \prod_{i=1}^n p_{\theta}(x_i | \mathbf{x}_{<i}), \quad \text{where} \quad p_{\theta}(x_i | \mathbf{x}_{<i}) = \frac{\exp(\mathbf{e}_i^{\top} \mathbf{h}_i)}{\sum_{j=1}^{|\mathcal{V}|} \exp(\mathbf{e}_j^{\top} \mathbf{h}_i)}.$$

107 Here, $p_{\theta}(\cdot)$ is usually parameterized using token embeddings e and contextualized embeddings \mathbf{h}
 108 given by a Transformer [61] encoder.

109 After pretraining, G_θ can be directly used to generate new texts by recursively sampling tokens from
 110 its output probability distribution. Typically, a temperature hyperparameter $\tau > 0$ is introduced
 111 during sampling [20] to adjust the sharpness of the probability distribution:

$$p_{\theta}(x_i | \mathbf{x}_{<i}) = \frac{\exp(\mathbf{e}_i^{\top} \mathbf{h}_i / \tau)}{\sum_{j=1}^{|\mathcal{V}|} \exp(\mathbf{e}_j^{\top} \mathbf{h}_i / \tau)}, \quad (1)$$

²We do not consider regression tasks in this work due to the difficulty of generating texts conditioned on a continuous label space. However, there exist approaches [14, 49] that solve regression tasks by training on classification tasks. We leave the integration of SuperGen with these methods as future work for solving regression tasks.

³We assume the classifier to be bidirectional PLMs since they generally work better than unidirectional PLMs in NLU tasks; we can in principle use any PLM as the classifier.

112 where $\tau \rightarrow 0$ approximates greedily picking the most probable next token; $\tau \rightarrow \infty$ induces a uniform
 113 distribution. Additionally, sampled tokens can be confined to the top- k most probable ones to avoid
 114 low-quality tokens. In this work, we find such top- k sampling with temperature is sufficient to
 115 produce coherent and meaningful texts as training data for NLU tasks. Exploring more sophisticated
 116 sampling strategies [21] is left for future work.

117 3.2 Training Data Generation

118 When given a label-descriptive prompt such
 119 as “Write a negative review:”, humans are
 120 able to produce texts pertaining to the cor-
 121 responding class. We aim to leverage the
 122 strong text generation power of a unidirec-
 123 tional PLM G_θ for the same purpose of cre-
 124 ating class-conditioned training data. We
 125 note that G_θ is directly used for generation
 126 without any parameter updates. The prompts
 127 used for different NLU tasks in GLUE are
 128 summarized in Table 1.

Generating Single Sequences. For single-
 sequence NLU tasks such as sentiment clas-
 sification (*e.g.*, SST-2), we simply use a
 prompt w_y corresponding to label y as the
 beginning of the sequence and let G_θ gener-
 ate the remaining sequence:

$$x^g \leftarrow G_\theta(w_y),$$

129 where $G_\theta(w_y)$ denotes using w_y as the in-
 130 put to G_θ and recursively sampling tokens
 131 from the distribution in Eq. (1) until a full
 132 sequence is generated; x^g denotes the *gen-*
 133 *erated* sequence (*i.e.*, excluding the prompt),
 134 which will be paired with y to form one train-
 135 ing sample (x^g, y) .

136 For syntactic tasks like linguistic acceptabil-
 137 ity classification (*e.g.*, CoLA) which requires generating both linguistically acceptable and unaccept-
 138 able sequences, we start the sequence with random stop words and use varying sampling temperatures
 139 for generating different sequences. A smaller temperature (*e.g.*, $\tau = 0.1$ in Equation (1)) sharpens the
 140 sampling probability distribution towards the most probable tokens, thus the resulting sequence will
 141 more likely to be linguistically acceptable. Using a larger temperature (*e.g.*, $\tau = 10$ in Equation (1))
 142 flattens the sampling probability distribution to be more uniform, and the generated tokens will be
 143 nearly random, which can create linguistically incorrect sequences.

Generating Sequence Pairs. Sequence-pair classification tasks require generating two sequences
 of specific relationships (*e.g.*, entailment, contradiction). We sample⁴ the first sequence x^s from the
 pretraining corpus \mathcal{D} , concatenate the prompt w_y with x^s , and let G_θ generate the second sequence
 x^g :

$$x^g \leftarrow G_\theta([x^s; w_y]), x^s \sim \mathcal{D}.$$

144 The sequence pair training sample will then be formed as (x^s, x^g, y) .

145 **Rewarding and Penalizing Repetitions for Sequence Pair Generation.** A common issue in text
 146 generation is degenerate repetition [21, 23, 47, 67] where generated texts can get stuck in repetition
 147 loops. To address this issue, one approach is to discourage repetition by reducing the logits of tokens
 148 that are already in the sequence before performing sampling [23]. In sequence pair generation,
 149 however, it is sometimes desirable to encourage the second sequence to repeat some words in the
 150 first sentence (*e.g.*, for generating an entailment or a paraphrase). Therefore, we propose a simple

⁴In principle, we can also generate the first sequence using G_θ , but we find sampling from \mathcal{D} improves the
 diversity of texts.

Table 1: Prompts used to generate class-conditioned texts for different GLUE tasks. SST-2 is a single-
 sequence classification task and the rest are sequence-
 pair classification tasks. Generation for CoLA does
 not use prompts but by varying sampling tempera-
 tures. x^s denotes a sequence randomly sampled from
 the pretraining corpus; x^g denotes the sequence to be
 generated by G_θ ; \dots denotes skipping at least one
 sequence. See Appendix B for more details.

| Task | Label | Prompt |
|-------|----------------|---|
| SST-2 | positive | Rating: 5.0 x^g |
| | negative | Rating: 1.0 x^g |
| MNLI | entailment | x^s . In other words, x^g |
| | neutral | x^s . Furthermore, x^g |
| | contradiction | There is a rumor that x^s . However, the truth is: x^g |
| QNLI | entailment | $x^s?$ x^g |
| | not entailment | $x^s?$... x^g |
| RTE | entailment | x^s . In other words, x^g |
| | not entailment | x^s . Furthermore, x^g |
| MRPC | equivalent | x^s . In other words, x^g |
| | not equivalent | x^s . Furthermore, x^g |
| QQP | equivalent | $x^s?$ In other words, x^g |
| | not equivalent | $x^s?$ Furthermore, x^g |

151 modification of Eq. (1) that rewards/penalizes repetition based on whether the token has appeared in
 152 $\mathbf{x}^s/\mathbf{x}^g$:

$$p_\theta(x_i|\mathbf{x}_{<i}) = \frac{\exp(\mathbf{e}_i^\top \mathbf{h}_i/\omega)}{\sum_{j=1}^{|\mathcal{V}|} \exp(\mathbf{e}_j^\top \mathbf{h}_i/\omega)}, \quad \text{where } \omega = \begin{cases} \tau\alpha & x_i \in \mathbf{x}^s \wedge x_i \notin \mathbf{x}^g \\ \tau\beta & x_i \in \mathbf{x}^g \\ \tau & \text{else} \end{cases}, \quad (2)$$

153 and $\alpha > 0, \beta > 0$ are hyperparameters. By setting $\alpha < 1$ and $\beta > 1$, we can promote tokens in
 154 \mathbf{x}^s that have not appeared in \mathbf{x}^g to have a higher chance of being generated, and discourage the
 155 generation of repetitive tokens in \mathbf{x}^g to mitigate degenerate repetition. The parameters used for
 156 different tasks are listed in Appendix C Table 8.

157 3.3 Effective Fine-Tuning on Generated Texts

158 With the generated training data, one can fine-tune a bidirectional PLM C_ϕ as the classifier to perform
 159 the NLU task. However, training C_ϕ via standard supervised training on all generated texts is likely
 160 to yield suboptimal performance on downstream tasks because (1) the generated texts may contain
 161 noise as G_θ may not always produce texts pertaining to the desired class, especially for challenging
 162 sequence pair tasks with subtle semantic relationships; and (2) the generated texts can be considered
 163 as originated from the domain of G_θ 's pretraining data, with a potentially different distribution from
 164 the downstream task; straightforward application of supervised training will result in overfitting
 165 to the pretraining domain and diminishing generalization ability, a common challenge in transfer
 166 learning [60, 83]. To address these challenges, we next introduce several simple and important
 167 strategies for more effective and stable fine-tuning on generated texts.

Selecting Quality Training Data. We aim to select generated texts \mathbf{x}^g that are most likely to
 pertain to the desired label y (i.e., with the highest $p(\mathbf{x}^g|y)$). The true probability $p(\mathbf{x}^g|y)$ is unknown
 and we estimate it via the generation probability given by G_θ conditioned on the prompt \mathbf{w}_y :

$$p(\mathbf{x}^g|y) \approx p_\theta(\mathbf{x}^g|\mathbf{w}_y) = \prod_{i=1}^n p_\theta(x_i|[\mathbf{w}_y; \mathbf{x}_{<i}^g]).$$

168 Since the above measure is biased towards shorter sequences, we instead use the geometric mean
 169 of the above conditional generation probability (or equivalently, the average log probability) of all
 170 tokens in \mathbf{x}^g as the ranking score, following [78]:

$$r = \frac{1}{n} \sum_{i=1}^n \log p_\theta(x_i|[\mathbf{w}_y; \mathbf{x}_{<i}^g]). \quad (3)$$

171 To construct a training set consisting of N samples per class, we will generate more samples (e.g.,
 172 $10N$), and select training data based on the score r in Eq. (3): For all tasks except CoLA, the
 173 top- N ones of each class are selected; for CoLA, the top- N ones are used as the training sample
 174 as linguistically acceptable sequences, and the bottom- N ones are as linguistically unacceptable
 175 sequences.

176 **Regularization for Better Generalization and Stability.** Even with the above training data
 177 selection procedure, the resulting training set may still contain noise and there exists domain difference
 178 from the downstream tasks. We apply two regularization techniques, *label smoothing* [58] and
 179 *temporal ensembling* [27] for better fine-tuning stability and generalization.

180 Given a training sample (\mathbf{x}^g, y) , *label smoothing* trains the classifier C_ϕ to minimize the standard
 181 cross-entropy loss between the label and the classifier's prediction $p_\phi(\mathbf{x}^g)$, except that the label is a
 182 weighted average of the one-hot vector and a uniform distribution over all labels:

$$\min_{\phi} - \sum_{j=1}^{|\mathcal{Y}|} q_j \log(p_\phi(\mathbf{x}^g)_j), \quad (4)$$

183 where $q_j = \mathbb{1}(j = y)(1 - \epsilon) + \epsilon/|\mathcal{Y}|$ and ϵ is the smoothing weight. By forcing the classifier to be
 184 less confident on training data, label smoothing improves robustness to label noise [35] and prevents
 185 overfitting to the training set [40], thus improving generalization to different domains.

186 The motivation for *temporal ensembling* is that neural networks usually first pick up easy and general
 187 patterns in the data before learning more sophisticated and dataset-specific features [79], and thus the
 188 earlier states of the network offer better generalizability to different domains. We therefore record the
 189 predictions $\mathbf{p}_\phi = p_\phi(\mathbf{x}^g)$ of C_ϕ on each training sample (\mathbf{x}^g, y) at different training steps, and use
 190 the accumulated moving-average predictions $\bar{\mathbf{z}}$ to regularize the latest model training. This also helps
 191 suppress the fluctuation in model predictions due to data noise, offering better noise-robustness [41].
 192 We update ensembled predictions $\bar{\mathbf{z}}$ once every B batches:

$$\hat{\mathbf{z}} \leftarrow \gamma \hat{\mathbf{z}} + (1 - \gamma) \mathbf{p}_\phi, \bar{\mathbf{z}} \leftarrow \hat{\mathbf{z}} / (1 - \gamma^t), \quad (5)$$

193 where $\hat{\mathbf{z}}$ has a zero initialization; γ is the momentum parameter; t is the number of updates $\bar{\mathbf{z}}$ has
 194 received; the division $(1 - \gamma^t)$ is for bias correction [27]. We also use the ensembled prediction $\bar{\mathbf{z}}$ as
 195 a reliable signal to filter out noisy training samples: Only those samples on which $\bar{\mathbf{z}}$ strongly agrees
 196 with the label y (i.e., $\bar{z}_y > \delta$ where $\delta > 0$ is a threshold parameter) will be used for training.

197 We regularize model training by extending Eq. (4) to add a KL divergence regularization term from
 198 the model prediction to the ensembled prediction weighed by λ :

$$\min_{\phi} - \sum_{j=1}^{|\mathcal{Y}|} q_j \log(p_\phi(\mathbf{x}^g)_j) - \lambda \sum_{j=1}^{|\mathcal{Y}|} \bar{z}_j \log \frac{p_\phi(\mathbf{x}^g)_j}{\bar{z}_j}. \quad (6)$$

199 We follow [27] to slowly ramp-up λ during training.

200 3.4 Overall Algorithm

201 We summarize SuperGen for single-
 202 sequence NLU tasks in Algorithm 1. Solv-
 203 ing sequence-pair problems follows the
 204 same algorithm except the pretraining cor-
 205 pus \mathcal{D} is needed for sampling the first se-
 206 quence \mathbf{x}^s .

207 4 Experimental Setup

208 **Downstream Tasks and Metrics.** We
 209 use all the tasks included in GLUE [62]
 210 except STS-B which is a regression task.
 211 Please refer to Appendix A for more de-
 212 tails about GLUE tasks. We follow the
 213 evaluation protocol of [13]: We use F1
 214 score as the metric for QQP and MRPC,
 215 Matthews correlation for CoLA, and accu-
 216 racy for the rest of the tasks. The original
 217 development sets of these tasks are used
 218 for testing. For all reported results, we in-
 219 clude the average and standard deviation
 220 over 5 different random seeds.

221 **Models.** Unless specified otherwise, we
 222 use CTRL (1.63B parameters) [23] as the
 223 generator G_θ and COCO-LM_{Large} (367M
 224 parameters) [37] as the classifier C_ϕ . We
 225 also show the results using similar-sized
 226 PLMs (GPT-2 [47]/RoBERTa [33]) as the
 227 generator/classifier in Appendix D.

228 **Fine-Tuning Settings and Hyperparam-**
 229 **eters.** We note that SuperGen is compat-
 230 ible with any fine-tuning method; while us-
 231 ing more sophisticated methods may grant
 232 further performance improvement, we use

Algorithm 1: SuperGen for Zero-Shot Learning.

Input: \mathcal{Y} : Label space; \mathcal{P} : Label-descriptive prompts; G_θ :
 Unidirectional PLM; C_ϕ : Bidirectional PLM.

Parameter: N : Number of training samples per class to
 generate; $M (\gg N)$: Number of total
 training samples to generate; T : Number of
 training steps; B : Ensemble prediction
 update interval; δ : Threshold parameter.

Output: C_ϕ^* : Classifier that classifies input texts into \mathcal{Y} .

```

for  $y \in \mathcal{Y}$  do
   $\mathcal{T}_y \leftarrow \{\}$ 
  // Class  $y$  train set init.
  for  $i \in [1, 2, \dots, M]$  do
     $\mathbf{x}^g \leftarrow G_\theta(\mathbf{w}_y)$ 
     $\mathcal{T}_y \leftarrow \mathcal{T}_y \cup \{(\mathbf{x}^g, y)\}$ 
  end
end
 $\mathcal{T} \leftarrow \{\}$ 
// Selected train set.
for  $y \in \mathcal{Y}$  do
  Sort  $\mathcal{T}_y$  in descending order by Eq. (3)
   $\mathcal{T} \leftarrow \mathcal{T} \cup \mathcal{T}_y[:N]$ 
end
 $\hat{\mathbf{z}} \leftarrow \mathbf{0}$ 
// Ensembled prediction init.
 $\mathcal{T}^* \leftarrow \mathcal{T}$ 
// Filtered train set.
for  $i \in [1, 2, \dots, T]$  do
  Fine-tune  $C_\phi$  via Eq. (6) on a minibatch of  $\mathcal{T}^*$ 
  if  $i \% B = 0$  then
    Update  $\hat{\mathbf{z}}, \bar{\mathbf{z}}$  via Eq. (5)
     $\mathcal{T}^* \leftarrow \{(\mathbf{x}^g, y) | \bar{z}_y > \delta, (\mathbf{x}^g, y) \in \mathcal{T}\}$ 
  end
end
return  $C_\phi^* = C_\phi$ 

```

the basic prompt-based fine-tuning with manual templates

Table 2: Results on seven GLUE classification tasks. We report average and standard deviation (as subscripts) performance over 5 different random seeds. †: Results from LM-BFF [13].

| Method | MNLI-(m/mm) (Acc.) | QQP (F1) | QNLI (Acc.) | SST-2 (Acc.) | CoLA (Matt.) | RTE (Acc.) | MRPC (F1) | AVG |
|---|---|----------------------------|----------------------------|----------------------------|-----------------------------|----------------------------|----------------------------|-------------|
| Zero-Shot Setting: No task-specific data (neither labeled nor unlabeled). | | | | | | | | |
| Prompting† | 50.8 _{0.0} /51.7 _{0.0} | 49.7 _{0.0} | 50.8 _{0.0} | 83.6 _{0.0} | 2.0 _{0.0} | 51.3 _{0.0} | 61.9 _{0.0} | 50.1 |
| SuperGen | 72.3 _{0.5} / 73.8 _{0.5} | 66.1 _{1.1} | 73.3 _{1.9} | 92.8 _{0.6} | 32.7 _{5.5} | 65.3 _{1.2} | 82.2 _{0.5} | 69.4 |
| - data selection | 63.7 _{1.5} /64.2 _{1.6} | 62.3 _{2.2} | 63.9 _{3.2} | 91.3 _{2.0} | 30.5 _{8.8} | 62.4 _{1.5} | 81.6 _{0.2} | 65.1 |
| - label smooth | 70.7 _{0.8} /72.1 _{0.7} | 65.1 _{0.9} | 71.4 _{2.5} | 91.0 _{0.9} | 9.5 _{1.0} | 64.8 _{1.1} | 83.0 _{0.7} | 65.2 |
| - temporal ensemble | 62.0 _{4.6} /63.6 _{4.8} | 63.9 _{0.3} | 72.4 _{2.0} | 92.5 _{0.9} | 23.5 _{7.0} | 63.5 _{1.0} | 78.8 _{2.2} | 65.3 |
| Few-Shot Setting: Use 32 labeled samples/class (half for training and half for development). | | | | | | | | |
| Fine-tuning† | 45.8 _{6.4} /47.8 _{6.8} | 60.7 _{4.3} | 60.2 _{6.5} | 81.4 _{3.8} | 33.9 _{14.3} | 54.4 _{3.9} | 76.6 _{2.5} | 59.1 |
| Manual prompt† | 68.3 _{2.3} /70.5 _{1.9} | 65.5 _{5.3} | 64.5 _{4.2} | 92.7 _{0.9} | 9.3 _{7.3} | 69.1 _{3.6} | 74.5 _{5.3} | 63.6 |
| + demonstration† | 70.7 _{1.3} / 72.0 _{1.2} | 69.8 _{1.8} | 69.2 _{1.9} | 92.6 _{0.5} | 18.7 _{8.8} | 68.7 _{2.3} | 77.8 _{2.0} | 66.9 |
| Auto prompt† | 68.3 _{2.5} /70.1 _{2.6} | 67.0 _{3.0} | 68.3 _{7.4} | 92.3 _{1.0} | 14.0 _{14.1} | 73.9 _{2.2} | 76.2 _{2.3} | 65.8 |
| + demonstration† | 70.0 _{3.6} /72.0 _{3.1} | 67.7 _{5.8} | 68.5 _{5.4} | 93.0 _{0.6} | 21.8 _{15.9} | 71.1 _{5.3} | 78.1 _{3.4} | 67.3 |
| Fully supervised† | 89.8/89.5 | 81.7 | 93.3 | 95.0 | 62.6 | 80.9 | 91.4 | 84.9 |

Table 3: Results with different groups of prompts. CoLA does not use prompts for generation. The number of prompt groups is equal to the number of the task labels.

| Prompt Group | MNLI-(m/mm) | QQP | QNLI | SST-2 | RTE | MRPC |
|----------------|---|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| # 0 (Original) | 72.3 _{0.5} / 73.8 _{0.5} | 66.1 _{1.1} | 73.3 _{1.9} | 92.8 _{0.6} | 65.3 _{1.2} | 82.2 _{0.5} |
| # 1 | 70.7 _{1.4} /72.4 _{1.2} | 65.5 _{1.4} | 71.9 _{1.7} | 92.2 _{0.9} | 64.4 _{1.6} | 81.9 _{0.4} |
| # 2 | 70.8 _{0.6} /72.1 _{0.8} | 65.6 _{1.1} | 72.2 _{2.2} | 92.4 _{0.8} | 64.7 _{1.8} | 81.8 _{0.8} |
| # 3 | 70.9 _{1.4} /72.2 _{1.4} | - | - | - | - | - |
| Mixed | 72.2 _{0.7} /73.4 _{0.6} | 66.9 _{1.5} | 73.0 _{1.7} | 92.8 _{0.9} | 66.3 _{1.0} | 81.3 _{2.0} |

233 approach for simplicity and clarity. For all tasks, we use the same templates and label words as in
 234 [13]. Under the zero-shot learning setting, it is not possible to tune hyperparameters due to the lack
 235 of validation sets. Therefore, we keep all fine-tuning hyperparameters (*e.g.*, learning rate, batch size,
 236 training epochs, number of generated training samples, label smoothing and temporal ensembling
 237 hyperparameters) same across all tasks. See Appendix C Table 9 for details.

238 **Compared Methods and Ablations.** We include the results of zero-shot prompting, standard
 239 few-shot fine-tuning and the four few-shot prompt-based fine-tuning methods proposed in [13]. We
 240 also conduct ablation studies by removing the following three techniques from SuperGen one at a
 241 time: (1) not using Eq. (3) for training data selection but randomly selecting the same amount of
 242 training data (- data selection); (2) not using label smoothing (- label smooth) but using one-hot
 243 labels; and (3) not using temporal ensembling (*i.e.*, using Eq. (4) instead of Eq. (6) as the training
 244 objective) (- temporal ensemble). Lastly, we report the fully supervised fine-tuning results trained on
 245 the entire training sets.

246 5 Evaluation

247 5.1 Main Results

248 We present the results of SuperGen, its ablations and compared methods in Table 2. Overall, SuperGen
 249 significantly outperforms zero-shot prompting and achieves an overall better result than all few-shot
 250 methods. Notably, SuperGen results in much smaller variance over different random seeds than
 251 few-shot approaches on most tasks—with access to more training data, fine-tuning of PLMs becomes
 252 much more stable. The ablation results demonstrate that all three strategies (*i.e.*, quality training
 253 data selection, label smoothing and temporal ensembling) play important roles in improving and
 254 stabilizing the final performance, especially on challenging tasks like MNLI.

255 5.2 Using Different Prompts

256 One important factor of SuperGen is the choice of label-descriptive prompts as they directly influence
 257 the quality of generated training samples. To study the impact of different prompt choices on the
 258 final model performance, we create different groups of prompts other than the original ones. We
 259 replace the prompt for one label used in Table 1 with a synonymous one and keep other prompts
 260 unchanged when forming a different prompt group (Please refer to Appendix B Table 7 for details).
 261 We also experiment with mixing the generated data by different prompt groups (mixed). The results
 262 are shown in Table 3. Overall, the model performance under different prompts is quite close, except
 263 on RTE whose test set is very small, potentially resulting in the higher variance. In this work, we
 264 manually choose simple prompts that make intuitive sense, and we leave the automatic searching of
 265 optimal prompts as future work.

266 5.3 Results with Different Amount of Generated Data

267 With training data automatically created by the generator, we can have
 268 a virtually infinite amount of training samples. We show the results
 269 of using different amount of generated data (after quality data selection)
 270 for fine-tuning the classifier C_ϕ in Fig. 2 on MNLI-m and SST-2.
 271 When the number of training data is small (*e.g.*, 100), the fine-tuning
 272 variance is high, resulting in the similar instability issue with few-shot
 273 settings. With more generated data used, both average performance and
 274 training stability improve, yielding
 275 comparable results (with smaller variance) to fine-tuning using few-shot task-specific data. However,
 276 when too many generated data (*e.g.*, 10,000) are used, the classifier’s performance slightly drops,
 277 probably due to increased label noise—recall that the training data are selected based on the ranking
 278 score in Eq. (3), so using more data results in the inclusion of more lower-ranking texts in the training
 279 set and reduced data quality. One way to address this issue is to use a fixed selection ratio and increase
 280 the total number of generated texts to obtain a larger number of high-quality training data. However,
 281 this comes at a greater computation cost in the generation step. An important future direction is thus
 282 to develop better data selection strategies.

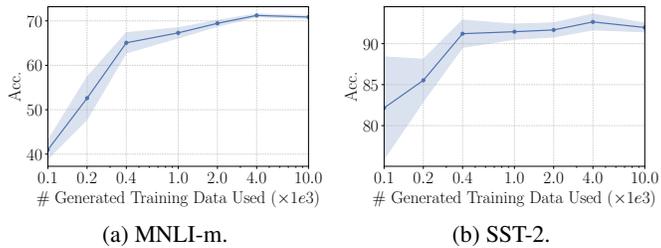


Figure 2: Results with different amount of generated training data used. Dots and error bars are the average performance and the standard deviation over 5 seeds, respectively.

290 5.4 Using Generators for Knowledge Distillation

291 Apart from using unidirectional PLMs G_θ for training
 292 data generation, one could also directly apply
 293 them to unlabeled data formulated as prompts to obtain
 294 zero-shot predictions (*i.e.*, prompting [5, 13]),
 295 which can then be used as soft labels to train the
 296 classifier C_ϕ . In Table 4, we show (1) the zero-
 297 shot prediction accuracy of CTRL (the best out of
 298 three different prompts, details in Appendix F) and
 299 (2) the classifier performance trained from CTRL’s
 300 predictions on the entire unlabeled training set as
 301 soft labels (*i.e.*, knowledge distillation). Similar
 302 to the observations in previous studies [5, 66, 81],
 303 the zero-shot predictions of unidirectional PLMs are quite inaccurate and directly using them as
 304 soft labels to train classifiers does not yield good results. We hypothesize that the advantages of
 305 using unidirectional PLMs for training data generation over using them for zero-shot predictions are
 306 twofold: (1) Better flexibility in prompt formats. When unidirectional PLMs are used for zero-shot
 307 predictions, the prompts have to be designed so that the label word is the last token in the sequence to
 308 be predicted, as unidirectional PLMs cannot attend to subsequent tokens. Such constraints may result
 309 in the prompt being dissimilar to the pretraining data distribution and worsen the prediction quality
 310 of the PLMs. On the contrary, using unidirectional PLMs for generation is not subject to any prompt

Table 4: Comparisons with using CTRL for zero-shot prompting and for knowledge distillation. \dagger : The entire training set is used as unlabeled data.

| Method | MNLI-(m/mm) | SST-2 |
|--------------------------------|---|----------------------------|
| SuperGen | 72.3 _{0.5} / 73.8 _{0.5} | 92.8 _{0.6} |
| CTRL Prompting | 38.5 _{0.0} /39.2 _{0.0} | 72.5 _{0.0} |
| Knowledge Distill [†] | 40.8 _{0.5} /41.5 _{0.6} | 73.6 _{0.8} |

311 format constraints. (2) More direct uses of PLMs’ language modeling ability. Using unidirectional
 312 PLMs for training data generation *directly* leverages the PLMs’ output token probability. Applying
 313 PLMs for zero-shot prediction, however, requires an additional step to convert token predictions to
 314 label predictions (*i.e.*, the verbalizer [52]), and such a mapping process usually necessitates manual
 315 curation and can hardly be optimal [13] especially without abundant task-specific data.

316 5.5 Using SuperGen in Few-Shot Settings

317 We present a simple extension of Super-
 318 Gen to few-shot settings and show that
 319 the generated data of SuperGen may fur-
 320 ther improve the few-shot performance.

321 When few-shot samples are available, we
 322 first fine-tune the classifier on the few-
 323 shot training set (standard prompt-based
 324 fine-tuning without regularization), and
 325 then continue fine-tuning the classifier
 326 on the generated data by SuperGen as

327 described in Section 3.3. This allows the classifier to effectively leverage the knowledge from
 328 the few-shot training set to filter out noisy samples in the generated data, as temporal ensembling
 329 regularizes the classifier to remember the predictions learned previously and only keeps samples on
 330 which the model predictions agree with the label. As shown in Table 5, such a simple approach that
 331 applies SuperGen to few-shot settings improves both zero-shot SuperGen and few-shot prompt-based
 332 fine-tuning. We note that few-shot samples are not used in the training data generation stage, and
 333 we expect the results to be even better if they are leveraged to generate training data closer to the
 334 task-specific distribution. Possible ways to use few-shot samples for generation include using them
 335 as demonstrations [5], for creating augmentations [28] and for tuning the generators. We leave the
 336 explorations of generating higher quality data by leveraging few-shot samples for future work.

Table 5: Using SuperGen for few-shot learning. The few-shot setup follows [13].

| Method | MNLI-(m/mm) | QQP | CoLA |
|--------------------------------------|--|---------------------------|---------------------------|
| Zero-Shot SuperGen | 72.3 _{0.5} /73.8 _{0.5} | 66.1 _{1.1} | 32.7 _{5.5} |
| Few-Shot Manual Prompt + SuperGen | 68.3 _{2.3} /70.5 _{1.9} | 65.5 _{5.3} | 9.3 _{7.3} |
| | 73.1_{0.9}/74.3_{0.6} | 69.9_{2.5} | 45.5_{6.5} |

337 6 Discussions and Conclusions

338 **Ethical Considerations.** While PLMs have demonstrated remarkable text generation and un-
 339 derstanding capability, they can come with potential risks or harms [2, 3, 5] such as generating
 340 misinformation [42] or amplifying harmful biases [45]. The focus of our work is on utilizing exist-
 341 ing PLMs to generate training data for NLU tasks instead of developing new PLMs or generation
 342 methods. Therefore, our method can be used in company with any bias reduction and correction
 343 techniques [15, 36] to mitigate the risks of PLMs.

344 **Limitations.** One inherent limitation with zero-shot learning is the lack of access to task-specific
 345 samples for hyperparameter tuning, whereas the performance of neural networks is usually heavily
 346 dependent on the choice of hyperparameters even when the training algorithm and training set are
 347 fixed [44]. Also, without access to any labeled data, the generated training data quality may not be
 348 high enough to achieve good performance on challenging tasks, especially when the task distribution
 349 is significantly different from the pretraining data distribution (*e.g.*, the “linguistically incorrect” label
 350 of CoLA requires generating sequences with grammar mistakes – a different distribution from the one
 351 used to train PLMs). A promising direction to address the above limitations is extending SuperGen to
 352 few-shot settings and leveraging a small amount of labeled data for generating better quality data and
 353 for hyperparameter tuning.

354 **Conclusions.** We propose SuperGen, an automatic supervision generation approach for zero-shot
 355 learning of NLU tasks. By providing label-descriptive prompts as guidance to a unidirectional
 356 PLM, training data can be automatically created for fine-tuning a bidirectional PLM. Our framework
 357 differs from previous transfer-learning-based zero-shot methods in that SuperGen does not rely on
 358 cross-task annotations and eliminates the task difference in training and inference. We show that
 359 several strategies are important for effective and stable fine-tuning on generated data, including
 360 quality training data selection, label smoothing and temporal ensembling. SuperGen achieves strong
 361 performance on seven classification tasks of the GLUE benchmark, even yielding comparable or
 362 better results than sophisticated few-shot learning methods and offering better stability. There is large
 363 room for future work, including but not limited to: Extension to few-shot learning settings, exploring
 364 larger generator models, better fine-tuning techniques to leverage generated data and better strategies
 365 for selecting quality training data.

366 **References**

- 367 [1] Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev
368 Shlomov, Naama Tepper, and Naama Zwerdling. Do not have enough data? deep learning to
369 the rescue! In *AAAI*, 2020.
- 370 [2] Emily M. Bender and Alexander Koller. Climbing towards NLU: On meaning, form, and
371 understanding in the age of data. In *ACL*, 2020.
- 372 [3] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On
373 the dangers of stochastic parrots: Can language models be too big? In *ACM Conference on*
374 *Fairness, Accountability, and Transparency*, 2021.
- 375 [4] Luisa Bentivogli, Peter Clark, Ido Dagan, and Danilo Giampiccolo. The fifth pascal recognizing
376 textual entailment challenge. In *TAC*, 2009.
- 377 [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
378 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
379 Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M.
380 Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,
381 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,
382 Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020.
- 383 [6] Alvin Chan, Y. Ong, Bill Tuck Weng Pung, Aston Zhang, and Jie Fu. CoCon: A self-supervised
384 approach for controlled text generation. In *ICLR*, 2021.
- 385 [7] Jiaao Chen, Zichao Yang, and Diyi Yang. Mixtext: Linguistically-informed interpolation of
386 hidden space for semi-supervised text classification. In *ACL*, 2020.
- 387 [8] Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. ELECTRA:
388 Pre-training text encoders as discriminators rather than generators. In *ICLR*, 2020.
- 389 [9] Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment
390 challenge. In *Machine Learning Challenges Workshop*, 2005.
- 391 [10] Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason
392 Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled
393 text generation. In *ICLR*, 2020.
- 394 [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
395 deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.
- 396 [12] William B Dolan and Chris Brockett. Automatically constructing a corpus of sentential para-
397 phrases. In *International Workshop on Paraphrasing (IWP)*, 2005.
- 398 [13] Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot
399 learners. In *ACL*, 2021.
- 400 [14] Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence
401 embeddings. In *EMNLP*, 2021.
- 402 [15] Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. RealToxi-
403 cityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of EMNLP*,
404 2020.
- 405 [16] Danilo Giampiccolo, Bernardo Magnini, Ido Dagan, and Bill Dolan. The third pascal recog-
406 nizing textual entailment challenge. In *ACL-PASCAL workshop on textual entailment and*
407 *paraphrasing*, 2007.
- 408 [17] Aaron Gokaslan and Vanya Cohen. OpenWebText corpus. [http://Skylion007.github.
409 io/OpenWebTextCorpus](http://Skylion007.github.io/OpenWebTextCorpus), 2019.
- 410 [18] R Bar Haim, Ido Dagan, Bill Dolan, Lisa Ferro, Danilo Giampiccolo, Bernardo Magnini,
411 and Idan Szpektor. The second pascal recognising textual entailment challenge. In *PASCAL*
412 *Challenges Workshop on Recognising Textual Entailment*, 2006.

- 413 [19] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. DeBERTa: Decoding-enhanced
414 BERT with disentangled attention. In *ICLR*, 2021.
- 415 [20] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural
416 network. *ArXiv*, abs/1503.02531, 2015.
- 417 [21] Ari Holtzman, Jan Buys, Maxwell Forbes, and Yejin Choi. The curious case of neural text
418 degeneration. In *ICLR*, 2020.
- 419 [22] Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. Toward
420 controlled generation of text. In *ICML*, 2017.
- 421 [23] Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard
422 Socher. CTRL: A conditional transformer language model for controllable generation. *ArXiv*,
423 abs/1909.05858, 2019.
- 424 [24] Muhammad Khalifa, Hady ElSahar, and Marc Dymetman. A distributional approach to con-
425 trolled text generation. In *ICLR*, 2021.
- 426 [25] Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq R.
427 Joty, Richard Socher, and Nazneen Rajani. GeDi: Generative discriminator guided sequence
428 generation. In *EMNLP*, 2021.
- 429 [26] Sachin Kumar, Eric Malmi, Aliaksei Severyn, and Yulia Tsvetkov. Controlled text generation
430 as continuous optimization with multiple constraints. In *NeurIPS*, 2021.
- 431 [27] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. In *ICLR*,
432 2017.
- 433 [28] Kenton Lee, Kelvin Guu, Luheng He, Tim Dozat, and Hyung Won Chung. Neural data
434 augmentation via example extrapolation. *arXiv preprint arXiv:2102.01335*, 2021.
- 435 [29] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed,
436 Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. BART: Denoising sequence-to-sequence
437 pre-training for natural language generation, translation, and comprehension. In *ACL*, 2020.
- 438 [30] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.
439 In *ACL*, 2021.
- 440 [31] Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A.
441 Smith, and Yejin Choi. DExperts: Decoding-time controlled text generation with experts and
442 anti-experts. In *ACL*, 2021.
- 443 [32] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang.
444 GPT understands, too. *ArXiv*, abs/2103.10385, 2021.
- 445 [33] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy,
446 Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. RoBERTa: A robustly optimized BERT
447 pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- 448 [34] Robert L Logan IV, Ivana Balažević, Eric Wallace, Fabio Petroni, Sameer Singh, and Sebastian
449 Riedel. Cutting down on prompts and parameters: Simple few-shot learning with language
450 models. *arXiv preprint arXiv:2106.13353*, 2021.
- 451 [35] Michal Lukasik, Srinadh Bhojanapalli, Aditya Krishna Menon, and Surinder Kumar. Does label
452 smoothing mitigate label noise? In *ICML*, 2020.
- 453 [36] Xinyao Ma, Maarten Sap, Hannah Rashkin, and Yejin Choi. PowerTransformer: Unsupervised
454 controllable revision for biased language correction. In *EMNLP*, 2020.
- 455 [37] Yu Meng, Chenyan Xiong, Payal Bajaj, Saurabh Tiwary, Paul Bennett, Jiawei Han, and Xia
456 Song. COCO-LM: Correcting and contrasting text sequences for language model pretraining.
457 In *NeurIPS*, 2021.
- 458 [38] Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. Cross-task gener-
459 alization via natural language crowdsourcing instructions. *arXiv preprint arXiv:2104.08773*,
460 2021.

- 461 [39] Takeru Miyato, Andrew M. Dai, and Ian J. Goodfellow. Adversarial training methods for
462 semi-supervised text classification. In *ICLR*, 2017.
- 463 [40] Rafael Müller, Simon Kornblith, and Geoffrey E. Hinton. When does label smoothing help? In
464 *NeurIPS*, 2019.
- 465 [41] Duc Tam Nguyen, Chaithanya Kumar Mummadi, Thi-Phuong-Nhung Ngo, Thi Hoai Phuong
466 Nguyen, Laura Beggel, and Thomas Brox. SELF: Learning to filter noisy labels with self-
467 ensembling. In *ICLR*, 2020.
- 468 [42] Artidoro Pagnoni, Vidhisha Balachandran, and Yulia Tsvetkov. Understanding factuality in
469 abstractive summarization with FRANK: A benchmark for factuality metrics. In *NAACL*, 2021.
- 470 [43] Damian Pascual, Béni Egressy, Clara Meister, Ryan Cotterell, and Roger Wattenhofer. A
471 plug-and-play method for controlled text generation. In *Findings of EMNLP*, 2021.
- 472 [44] Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models.
473 In *NeurIPS*, 2021.
- 474 [45] Shrimai Prabhumoye, Yulia Tsvetkov, Ruslan Salakhutdinov, and Alan W. Black. Style transfer
475 through back-translation. In *ACL*, 2018.
- 476 [46] Raul Puri and Bryan Catanzaro. Zero-shot text classification with generative language models.
477 *ArXiv*, abs/1912.10165, 2019.
- 478 [47] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever.
479 Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- 480 [48] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,
481 Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified
482 text-to-text transformer. *Journal of Machine Learning Research*, 2019.
- 483 [49] Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence embeddings using siamese
484 BERT-networks. In *EMNLP*, 2019.
- 485 [50] Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang A. Sutawika, Zaid Alyafeai,
486 Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, Manan Dey, M SAIFUL BARI,
487 Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan
488 Chhablani, Nihal V. Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang,
489 Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang,
490 Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Févry, Jason Alan
491 Fries, Ryan Teehan, Stella Rose Biderman, Leo Gao, T. G. Owe Bers, Thomas Wolf, and
492 Alexander M. Rush. Multitask prompted training enables zero-shot task generalization. *ArXiv*,
493 abs/2110.08207, 2021.
- 494 [51] Teven Le Scao and Alexander M. Rush. How many data points is a prompt worth? In *NAACL*,
495 2021.
- 496 [52] Timo Schick and Hinrich Schütze. Exploiting cloze-questions for few-shot text classification
497 and natural language inference. In *EACL*, 2021.
- 498 [53] Timo Schick and Hinrich Schütze. Few-shot text generation with natural language instructions.
499 In *EMNLP*, 2021.
- 500 [54] Timo Schick and Hinrich Schütze. Generating datasets with pretrained language models. In
501 *EMNLP*, 2021.
- 502 [55] Timo Schick and Hinrich Schütze. It’s not just size that matters: Small language models are
503 also few-shot learners. In *NAACL*, 2021.
- 504 [56] Iyer Shankar, Dandekar Nikhil, and Csernai Kornél. First Quora dataset re-
505 lease: Question pairs, 2017. URL [https://www.quora.com/q/quoradata/
506 First-Quora-Dataset-Release-Question-Pairs](https://www.quora.com/q/quoradata/First-Quora-Dataset-Release-Question-Pairs).
- 507 [57] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y
508 Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a
509 sentiment treebank. In *EMNLP*, 2013.

- 510 [58] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna.
511 Rethinking the inception architecture for computer vision. In *CVPR*, 2016.
- 512 [59] Derek Tam, Rakesh R Menon, Mohit Bansal, Shashank Srivastava, and Colin Raffel. Improving
513 and simplifying pattern exploiting training. In *EMNLP*, 2021.
- 514 [60] Lisa Torrey and Jude Shavlik. Transfer learning. In *Handbook of research on machine learning
515 applications and trends: algorithms, methods, and techniques*. 2010.
- 516 [61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
517 Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- 518 [62] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman.
519 GLUE: A multi-task benchmark and analysis platform for natural language understanding. In
520 *EMNLP Workshop BlackboxNLP*, 2018.
- 521 [63] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill,
522 Omer Levy, and Samuel R. Bowman. SuperGLUE: A stickier benchmark for general-purpose
523 language understanding systems. In *NeurIPS*, 2019.
- 524 [64] Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. Towards zero-label language learning.
525 *arXiv preprint arXiv:2109.09193*, 2021.
- 526 [65] Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability
527 judgments. In *TACL*, 2019.
- 528 [66] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
529 Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners. In *ICLR*,
530 2022.
- 531 [67] Sean Welleck, Ilya Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston.
532 Neural text generation with unlikelihood training. In *ICLR*, 2020.
- 533 [68] Adina Williams, Nikita Nangia, and Samuel Bowman. A broad-coverage challenge corpus for
534 sentence understanding through inference. In *NAACL-HLT*, 2018.
- 535 [69] Qizhe Xie, Zihang Dai, Eduard H. Hovy, Minh-Thang Luong, and Quoc V. Le. Unsupervised
536 data augmentation for consistency training. In *NeurIPS*, 2020.
- 537 [70] Hanwei Xu, Yujun Chen, Yulun Du, Nan Shao, Yanggang Wang, Haiyu Li, and Zhilin Yang.
538 ZeroPrompt: Scaling prompt-based pretraining to 1, 000 tasks improves zero-shot generalization.
539 *ArXiv*, abs/2201.06910, 2022.
- 540 [71] Kevin Yang and Dan Klein. FUDGE: Controlled text generation with future discriminators. In
541 *NAACL*, 2021.
- 542 [72] Yiben Yang, Chaitanya Malaviya, Jared Fernandez, Swabha Swayamdipta, Ronan Le Bras,
543 Ji Ping Wang, Chandra Bhagavatula, Yejin Choi, and Doug Downey. Generative data augmenta-
544 tion for commonsense reasoning. In *Findings of EMNLP*, 2020.
- 545 [73] Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Ling-
546 peng Kong. Zerogen: Efficient zero-shot learning via dataset generation. *ArXiv*, abs/2202.07922,
547 2022.
- 548 [74] Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. Crossfit: A few-shot learning challenge for
549 cross-task generalization in nlp. In *EMNLP*, 2021.
- 550 [75] Wenpeng Yin, Jamaal Hay, and Dan Roth. Benchmarking zero-shot text classification: Datasets,
551 evaluation and entailment approach. In *EMNLP*, 2019.
- 552 [76] Wenpeng Yin, Nazneen Rajani, Dragomir Radev, Richard Socher, and Caiming Xiong. Universal
553 natural language processing with limited annotations: Try few-shot textual entailment as a start.
554 In *EMNLP*, 2020.
- 555 [77] Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-Woo Lee, and Woomyeong Park. Gpt3mix:
556 Leveraging large-scale language models for text augmentation. In *EMNLP Findings*, 2021.

- 557 [78] Weizhe Yuan, Graham Neubig, and Pengfei Liu. BARTScore: Evaluating generated text as text
558 generation. In *NeurIPS*, 2021.
- 559 [79] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding
560 deep learning requires rethinking generalization. In *ICLR*, 2017.
- 561 [80] Ningyu Zhang, Luoqi Li, Xiang Chen, Shumin Deng, Zhen Bi, Chuanqi Tan, Fei Huang,
562 and Huajun Chen. Differentiable prompt makes pre-trained language models better few-shot
563 learners. In *ICLR*, 2022.
- 564 [81] Tony Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use:
565 Improving few-shot performance of language models. In *ICML*, 2021.
- 566 [82] Ruiqi Zhong, Kristy Lee, Zheng Zhang, and Dan Klein. Adapting language models for zero-shot
567 learning by meta-tuning on dataset and prompt collections. In *Findings of EMNLP*, 2021.
- 568 [83] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong,
569 and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 2020.
- 570 [84] Daniel M. Ziegler, Nisan Stiennon, Jeff Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul
571 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *ArXiv*,
572 abs/1909.08593, 2019.

573 Checklist

- 574 1. For all authors...
- 575 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
576 contributions and scope? [Yes]
- 577 (b) Did you describe the limitations of your work? [Yes]
- 578 (c) Did you discuss any potential negative societal impacts of your work? [Yes]
- 579 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
580 them? [Yes]
- 581 2. If you are including theoretical results...
- 582 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 583 (b) Did you include complete proofs of all theoretical results? [N/A]
- 584 3. If you ran experiments...
- 585 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
586 mental results (either in the supplemental material or as a URL)? [Yes]
- 587 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
588 were chosen)? [Yes]
- 589 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
590 ments multiple times)? [Yes]
- 591 (d) Did you include the total amount of compute and the type of resources used (e.g., type
592 of GPUs, internal cluster, or cloud provider)? [Yes]
- 593 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 594 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 595 (b) Did you mention the license of the assets? [Yes]
- 596 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 597
- 598 (d) Did you discuss whether and how consent was obtained from people whose data you’re
599 using/curating? [N/A]
- 600 (e) Did you discuss whether the data you are using/curating contains personally identifiable
601 information or offensive content? [N/A]
- 602 5. If you used crowdsourcing or conducted research with human subjects...
- 603 (a) Did you include the full text of instructions given to participants and screenshots, if
604 applicable? [N/A]
- 605 (b) Did you describe any potential participant risks, with links to Institutional Review
606 Board (IRB) approvals, if applicable? [N/A]
- 607 (c) Did you include the estimated hourly wage paid to participants and the total amount
608 spent on participant compensation? [N/A]

609 **A GLUE Tasks**

610 We provide the details of the seven classification tasks included in the GLUE benchmark.

611 **MNLI:** Multi-genre Natural Language Inference [68] aims to predict whether a given premise
 612 sentence entails, contradicts or neutral with respect to a given hypothesis sentence.

613 **QQP:** Quora Question Pairs [56] aims to determine whether a pair of questions asked are semantically
 614 equivalent.

615 **QNLI:** Question Natural Language Inference aims to predict whether a given sentence contains the
 616 answer to a given question sentence.

617 **SST-2:** Stanford Sentiment Treebank [57] aims to determine if a movie review has positive or negative
 618 sentiment.

619 **CoLA:** Corpus of Linguistic Acceptability [65] aims to determine whether a given sentence is
 620 linguistically acceptable or not.

621 **RTE:** Recognizing Textual Entailment [4, 9, 16, 18] aims to predict whether a given premise sentence
 622 entails a given hypothesis sentence or not.

623 **MRPC:** Microsoft Research Paraphrase Corpus [12] aims to predict whether two sentences are
 624 semantically equivalent or not.

625 **B Details of Prompts Used for Different Tasks**

Table 6: Extensions of Table 1 with more details of prompts used to generate class-conditioned texts for different GLUE tasks. SST-2 and CoLA are single-sequence classification tasks and the rest are sequence-pair classification tasks. Generation for CoLA does not use prompts but by varying sampling temperatures. Text generation with CTRL [23] requires starting with control codes, and we use the ones that correspond to the pretraining corpus where the first sequence is sampled: For MNLI, RTE and MRPC, the first sequence is sampled from Wikipedia; for QNLI and QQP, the first sequence is sampled from OpenWebText [17]. x^s denotes a sequence randomly sampled from the pretraining corpus; x^g denotes the sequence to be generated by G_θ ; \dots denotes skipping at least one sequence. The prompts used for SST-2 are part of the CTRL [23] codes.

| Task | Task Type | Control Code | Label | Prompt |
|-------|-----------------|--------------|--|---|
| SST-2 | single-sequence | Reviews | positive negative | Rating: 5.0 x^g Rating: 1.0 x^g |
| CoLA | single-sequence | Links | grammatical not grammatical | x^g x^g |
| MNLI | sequence-pair | Wikipedia | entailment neutral contradiction | x^s . In other words, x^g x^s . Furthermore, x^g There is a rumor that x^s . However, the truth is: x^g |
| QNLI | sequence-pair | Links | entailment not entailment | $x^{s?} x^g$ $x^{s?} \dots x^g$ |
| RTE | sequence-pair | Wikipedia | entailment not entailment | x^s . In other words, x^g x^s . Furthermore, x^g |
| MRPC | sequence-pair | Wikipedia | equivalent not equivalent | x^s . In other words, x^g x^s . Furthermore, x^g |
| QQP | sequence-pair | Links | equivalent not equivalent | $x^{s?}$ In other words, x^g $x^{s?}$ Furthermore, x^g |

626 We present more details about the prompts used for different tasks in Table 6 which is an extended
 627 version of Table 1.

628 For SST-2, we fix the beginning of the generated sequence x^g to be “The/this film/movie” to make
 629 sure the generated texts are related to movie reviews.

630 For CoLA, we start the generated sequence x^g with a random stop word.

Table 7: Different prompt groups used in the experiments of Section 5.2. We replace the original prompt for each label with an alternative one and keep other prompts unchanged when forming a different prompt group.

| Task | Label | Original | Alternative |
|-------|----------------|--|--|
| SST-2 | positive | Rating: 5.0 x^g | Rating: 4.0 x^g |
| | negative | Rating: 1.0 x^g | Rating: 2.0 x^g |
| MNLI | entailment | x^s . In other words, x^g | x^s . To put it another way, x^g |
| | neutral | x^s . Furthermore, x^g | x^s . In addition, x^g |
| | contradiction | There is a rumor that x^s . However, the truth is: x^g | People believe that x^s . However, the truth is: x^g |
| QNLI | entailment | $x^s?$ x^g | Question: $x^s?$ Answer: x^g |
| | not entailment | $x^s? \dots x^g$ | Question: $x^s?$ Answer: $\dots x^g$ |
| RTE | entailment | x^s . In other words, x^g | x^s . To put it another way, x^g |
| | not entailment | x^s . Furthermore, x^g | x^s . In addition, x^g |
| MRPC | equivalent | x^s . In other words, x^g | x^s . To put it another way, x^g |
| | not equivalent | x^s . Furthermore, x^g | x^s . In addition, x^g |
| QQP | equivalent | $x^s?$ In other words, x^g | $x^s?$ To put it another way, x^g |
| | not equivalent | $x^s?$ Furthermore, x^g | $x^s?$ In addition, x^g |

631 For QNLI and QQP, the first sequence is always a question, and we require the sampled sequence
632 x^s to end with a question mark and begin with one of the following words: “how”, “what”, “why”,
633 “who”, “which”, “where”, “when”, “whom”, “whose”.

634 For QNLI, the generated sequence x^g for the “entailment” label is the one that immediately follows
635 the sampled sequence x^s ; the generated sequence x^g for the “not entailment” label is randomly
636 sampled from the paragraph following x^g excluding the first sequence that immediately follows x^g .

637 We also show the different prompt groups used in the experiments of Section 5.2 in Table 7.

638 C Hyperparameters and Reproducibility

639 Hyperparameters for Generating

640 **Training Data.** Table 8 lists the hyperparameters used in the training
641 data generation stage. For sequence-pair tasks, we use greedy sampling
642 for better reproducibility. For labels that require generating entailment,
643 paraphrase, or equivalent sequence pairs, we set $\alpha \leq 1$ to encourage
644 word overlapping between the second sequence and the first sequence;
645 otherwise, we set $\alpha = \beta > 1$ to discourage word repetition.
646

647 To construct a training set consisting
648 of N samples per class, we will generate M samples per class, and select
649 training data based on the score r in Eq. (3): For all tasks except CoLA
650 and the “neutral” label of MNLI, the top- N ones of each class are selected;
651 for CoLA, the top- N ones are used as the training sample as linguistically
652 acceptable sequences, and the bottom- N ones are as linguistically
653 unacceptable sequences; for the “neutral” label of MNLI, we find it better
654 to randomly select N samples from

Table 8: Hyperparameters for generating training data of different tasks. τ : Temperature during sampling ($\tau = 0$ means using greedy sampling); α and β : Repetition rewarding/penalizing parameters; M : Number of total generated texts per label. The top- k sampling (if $\tau > 0$) uses $k = 10$.

| Task | Label | τ | α | β | M |
|-------|-----------------|-----------|----------|---------|--------|
| SST-2 | positive | 0.2 | - | 1.2 | 25,000 |
| | negative | | - | 1.2 | 25,000 |
| CoLA | grammatical | [0.1, 10] | - | 1.2 | 10,000 |
| | not grammatical | | - | 1.2 | 10,000 |
| MNLI | entailment | 0 | 0.8 | 1.1 | 25,000 |
| | neutral | | 1.3 | 1.3 | 25,000 |
| | contradiction | | 1.1 | 1.1 | 25,000 |
| QNLI | entailment | 0 | 0.9 | 1.2 | 25,000 |
| | not entailment | | 0.9 | 1.2 | 25,000 |
| RTE | entailment | 0 | 0.8 | 1.1 | 30,000 |
| | not entailment | | 1.1 | 1.1 | 30,000 |
| MRPC | equivalent | 0 | 0.8 | 1.1 | 30,000 |
| | not equivalent | | 1.1 | 1.1 | 30,000 |
| QQP | equivalent | 0 | 1.0 | 1.2 | 25,000 |
| | not equivalent | | 1.2 | 1.2 | 25,000 |

Table 9: Hyperparameters used for fine-tuning on different tasks (they are kept same for all tasks). Fine-tuning-related hyperparameters (*e.g.*, learning rate, batch size) follow the default values (when the validation set is not available) in Appendix A of [13]; regularization-related hyperparameters follow the default values in label smoothing and temporal ensembling. lr : Learning rate; bs : Batch size; $N|\mathcal{Y}|$: Total number of selected generated data (*i.e.*, training set size); B : Ensemble prediction update interval; T : Number of training steps; ϵ : Label smoothing parameter; γ : Temporal ensembling momentum parameter; δ : Threshold for filtering out noisy data; λ_{\max} : Maximum weight (after ramp-up) of temporal ensembling regularization.

| lr | bs | $N \mathcal{Y} $ | B | T | ϵ | γ | δ | λ_{\max} |
|------|------|------------------|-----|-------|------------|----------|----------|------------------|
| 1e-5 | 16 | 6,000 | 100 | 1,125 | 0.15 | 0.8 | 0.8 | 10 |

Table 10: Results with different generator/classifier PLMs.

| PLMs | MNLI-(m/mm) | SST-2 |
|--|--|---------------------|
| G_θ : CTRL, C_ϕ : COCO-LM | 72.3 _{0.5} /73.8 _{0.5} | 92.8 _{0.6} |
| G_θ : CTRL, C_ϕ : RoBERTa | 69.0 _{0.8} /70.6 _{0.9} | 93.3 _{1.5} |
| G_θ : GPT-2, C_ϕ : COCO-LM | 69.5 _{1.2} /71.3 _{1.3} | 88.2 _{1.8} |
| G_θ : GPT-2, C_ϕ : RoBERTa | 68.3 _{0.9} /69.7 _{0.7} | 88.6 _{0.8} |

666 the total M samples instead of using the ranking score, probably because a neutral hypothesis with
667 respect to the premise has a wide range of possibilities (*i.e.*, any hypothesis that is not entailed by
668 or contradicts with the premise will be neutral), and random selection improves the diversity in
669 generated hypotheses of the neutral label.

670 **Hyperparameters for Fine-Tuning.** Table 9 lists the hyperparameters used in the fine-tuning stage.
671 We keep them same across all tasks except CoLA which uses $\delta = 0$ because half of the training data
672 for CoLA are intentionally made to be of low quality (*i.e.*, as linguistically unacceptable sequences)
673 and there is no need to filter them out. We follow [27] to slowly ramp-up λ in Equation (6) during the
674 first 10 ensembles: $\lambda(t) = \lambda_{\max} \exp(-5(1 - t/10)^2)$ where t is the number of prediction ensembles
675 performed.

676 **Computation Environment.** All experiments are conducted on NVIDIA GeForce RTX 3090
677 GPUs. SuperGen can be run on typical research hardware (*e.g.*, with > 10GB GPU memory). The
678 generator PLM G_θ does not need to be trained so a relatively large generator can be used (*e.g.*, a
679 1.63B-parameter CTRL model).

680 D Using Different PLMs

681 The final performance is also relevant to the choice of PLMs as the generator/classifier. Apart
682 from the default PLM choice, we report the results of using GPT-2_{XLarge} (1.54B parameters) [47]
683 as the generator and RoBERTa_{Large} (356M parameters) [33] as the classifier in Table 10 with every-
684 thing else unchanged. When using GPT-2, we change the prompt used for SST-2 to “The film is
685 bad/terrible/awful.” for the negative label and “The film is good/great/excellent.” for the positive
686 label, since the original prompts used for SST-2 in Table 1 are a part of the control codes of CTRL and
687 cannot be effectively leveraged by GPT-2. Overall, both CTRL and GPT-2 are able to generate quality
688 training data for good fine-tuned classifier performance; CTRL consistently yields better results
689 than GPT-2 regardless of the choice of the classifier PLM, probably because CTRL is pretrained
690 with control codes which provide explicit guidance for generating texts of certain domains and
691 attributes. We also observe that the generated text quality is strongly correlated to the generator’s
692 model size—using a smaller version of GPT-2 (*e.g.*, with 117M parameters) results in significantly
693 less coherent texts and can hardly serve as training data. An interesting future direction is to try larger
694 generator PLMs (*e.g.*, GPT-3) which may create training data of better quality.

Table 11: Example generated texts for SST-2, MNLI and QQP. *Sampled sequences* from pretraining corpus (x^s) are italicized; generated sequences (x^g) are underlined; **prompts** (w^y) are in bold.

| Task | Label | Generated Text |
|-------|----------------|---|
| SST-2 | positive | Rating: 5.0 <u>The film is a great example of the kind of movie that you can watch over and over.</u> |
| | negative | Rating: 1.0 <u>The film was a total waste of time. I would not recommend this movie to anyone.</u> |
| MNLI | entailment | <i>The construction of the station began in 2016, when the opening date was estimated to mid-2020.</i> In other words, <u>the station was to open in 2020.</u> |
| | neutral | <i>The construction of the station began in 2016, when the opening date was estimated to mid-2020.</i> Furthermore, <u>it is expected that a new bus terminal will be built at this station.</u> |
| | contradiction | There is a rumor that <i>The construction of the station began in 2016, when the opening date was estimated to mid-2020.</i> However, the truth is: <u>The construction started in 2017, and the official opening date was set for March 31, 2018.</u> |
| QQP | equivalent | <i>What are the most wear resistant steels?</i> In other words, <u>what are the most durable steels?</u> |
| | not equivalent | <i>What are the most wear resistant steels?</i> Furthermore, <u>what is the best way to clean them?</u> |

Table 12: Negative example generated texts for MNLI that do not pertain to the desired label. *Sampled sequences* from pretraining corpus (x^s) are italicized; generated sequences (x^g) are underlined; **prompts** (w^y) are in bold. In example #1, the generated text corresponding to the “entailment” label actually pertains to the “neutral” label. In example #2, the generated text corresponding to the “contradiction” label actually pertains to the “neutral” label. In example #3, the generated text corresponding to the “neutral” label actually pertains to the “contradiction” label.

| Example | Label | Generated Text |
|---------|---------------|--|
| #1 | entailment | <i>Directed by Daniel J. Sullivan, the cast features Kelly AuCoin, Annie Parisse, Lee Tergesen and Alex Wolff.</i> In other words, <u>the cast features a lot of people who are not in the original Broadway cast.</u> |
| | neutral | <i>Directed by Daniel J. Sullivan, the cast features Kelly AuCoin, Annie Parisse, Lee Tergesen and Alex Wolff.</i> Furthermore, <u>it was announced that a sequel to “The Last Airbender” will be released in 2019.</u> |
| | contradiction | <i>Directed by Daniel J. Sullivan, the cast features Kelly AuCoin, Annie Parisse, Lee Tergesen and Alex Wolff.</i> However, the truth is: <u>The cast includes only two actors from the original Broadway production.</u> |
| #2 | entailment | <i>As a result, minor problems begin to look serious and rare events begin to seem common.</i> In other words, <u>a minor problem can begin to look serious and rare events seem common.</u> |
| | neutral | <i>As a result, minor problems begin to look serious and rare events begin to seem common.</i> Furthermore, <u>the number of people who are affected by these changes is increasing.</u> |
| | contradiction | There is a rumor that <i>As a result, minor problems begin to look serious and rare events begin to seem common.</i> However, the truth is: <u>The only thing that can cause such an effect is a major change in the weather.</u> |
| #3 | entailment | <i>The character was announced on 20 September 2019, but casting wasn’t revealed until 10 December 2019.</i> In other words, <u>the casting wasn’t revealed until 10 December 2019.</u> |
| | neutral | <i>The character was announced on 20 September 2019, but casting wasn’t revealed until 10 December 2019.</i> Furthermore, <u>it is unknown whether the role will be played by a male or female actor.</u> |
| | contradiction | There is a rumor that <i>The character was announced on 20 September 2019, but casting wasn’t revealed until 10 December 2019.</i> However, the truth is: <u>The character was announced on 19 October 2018.</u> |

695 E Case Studies

696 We present concrete examples of generated texts guided by prompts of different labels in Table 11.
 697 The generated sequences are not only coherent, but also pertain to the corresponding labels. For
 698 easier tasks like SST-2, the generated texts almost always correctly reflect the desired sentiment
 699 polarity specified by the prompt. For more difficult tasks like MNLI, sometimes the generated texts
 700 are not of the correct label. Table 12 demonstrates more concrete examples where the generated texts
 701 may not correctly pertain to the label of the prompt. The existence of such label noise motivates our
 702 use of the regularization techniques in the fine-tuning stage.

703 We believe that larger generator PLMs (*e.g.*, GPT-3 [5]) can bring about better text generation quality
 704 and improve the accuracy in producing texts that pertain to the desired class. Furthermore, better
 705 filtering strategies can be developed in the future to select training data with the correct labels.

706 F Knowledge Distillation Baseline Details

707 We show the concrete prompts used for the knowledge distillation baseline in Tables 13 and 14
 708 on MNLI and SST-2, respectively. We use the best prompt (prompt # 1 in both tables) out of the
 709 three according to the zero-shot test set prediction accuracy for generating soft labels to train the
 710 classification model (*i.e.*, knowledge distillation). The classifier is trained with Kullback–Leibler

Table 13: Different prompts used on MNLi for CTRL zero-shot prompting and knowledge distillation baselines. x_1 and x_2 denote the first and second input sequence, respectively.

| Prompt | Template | Label name |
|--------|---|--|
| # 1 | Sentence 1: x_1 Sentence 2: x_2 Does Sentence 1 entail Sentence 2? The answer is: | entailment: Yes neutral: Maybe contradiction: No |
| # 2 | Premise: x_1 Hypothesis: x_2 Does the premise entail the hypothesis? Options: Yes. No. Maybe. The answer is: | entailment: Yes neutral: Maybe contradiction: No |
| # 3 | Premise: x_1 Hypothesis: x_2 What is the relation between the premise and the hypothesis? Options: Entailment. Neutral. Contradiction. The answer is: | entailment: Entailment neutral: Neutral contradiction: Contradiction |

Table 14: Different prompts used on SST-2 for CTRL zero-shot prompting and knowledge distillation baselines. x denotes the input sequence.

| Prompt | Template | Label name |
|--------|------------------------|--|
| # 1 | x This is | positive: good; negative: bad |
| # 2 | x It was | positive: good; negative: bad |
| # 3 | Review: x Sentiment: | positive: Positive; negative: Negative |

⁷¹¹ (KL) divergence as the objective to approximate the soft labels generated by CTRL on the entire
⁷¹² training set.