

Using Synthetic Data for Domain Adaptation of Language Models

NLU Project Final Paper

Shubhankar Ranade

Akash Gupta

Abhinav Gupta

Victor Yu Cui

New York University

{shubhankar.r, aksq, gupta.abhinav, vyc8567}@nyu.edu

Abstract

In recent studies, Pretrained Language Models (PLMs) have demonstrated remarkable human level performance across multiple downstream tasks when fine-tuned on a large amount of task-specific training data. But getting large amount of labeled task data is often expensive and time consuming, and whatever annotated data is available may not be properly represented by the dataset the PLM is pre-trained on. In this paper, we utilize the few-shot learning capabilities of large PLMs (Radford et al., 2019, Keskar et al., 2019) and augment low-resource datasets using prompt-based techniques. We then fine-tune a domain adapted pretrained model DAPT (Gururangan et al.) on the augmented dataset. On comparing our model’s performance with state-of-the-art baseline models, we observe that our few-shot domain adaptation techniques benefit when the target domain highly overlaps with pre-training dataset domains of the generator PLM.

1 Introduction

Natural Language Processing (NLP) research these days has shifted from designing task-specific architectures to using large scale architectures pre-trained on multiple huge, general domain corpora. Pretrained language models (PLMs) (Brown et al., 2020, Devlin et al., 2019, Liu et al., 2019) have demonstrated remarkable human level performance across multiple downstream tasks when fine-tuned on a large amount of task-specific training data. (Wang et al., 2018, Howard and Ruder, 2018). However, procuring such a large amount of data is often expensive and time consuming. Therefore, it becomes important to develop PLMs that can extract maximum information from few labeled data samples.

In recent studies, it has been revealed that PLMs exhibit intriguing few-shot learning potential (Rad-

ford et al., 2019; Brown et al., 2020; Gao et al., 2021; Scao and Rush, 2021) and are able to leverage task-specific information. But when the target task is available in low quantity and has a different domain than the PLM’s source data (i.e. pre-training dataset), then performing well on downstream task becomes extremely difficult. Domain Adaptation as well as Data Augmentation techniques are widely used for tackling such issues.

Work done by Gururangan et al. and Lee et al., 2019 shows that additionally pretraining PLM on domain-specific corpus can help PLMs to adapt well on downstream tasks. However, these DAPT models require fine-tuning on moderately sized task-specific datasets, which can be scarce in many scenarios. This motivates us to look forward to several data augmentation methods and utilize these DAPT models. Unidirectional PLMs have demonstrated strong text generation power (Brown et al., 2020.) One of the ways to do controlled generation is prompt-based approach, where desired context is provided to the generator model through carefully designed prompts (Schick and Schütze, 2020). Building on this prompt-based approach, Meng et al., 2022 (SuperGen) explored the possibility of creating synthetic task-specific training data using a unidirectional PLM (e.g. CTRL Keskar et al., 2019 and Radford et al., 2019) as generator, then fine-tune a bidirectional PLM as classifier on the target task using the generated synthetic dataset.

Our work¹ relies directly on the SuperGen pipeline, but with the motivation to tackle the lack of task-specific data, we explore its potential in boosting performance of DAPT models (as in Gururangan et al.) and replacing expensive task-specific data. We decided to feed few-shot prompts to the generator because it can provide more contexts of the task domain to the generator thus steer the gen-

¹<https://github.com/ag2307/NLU-Project>

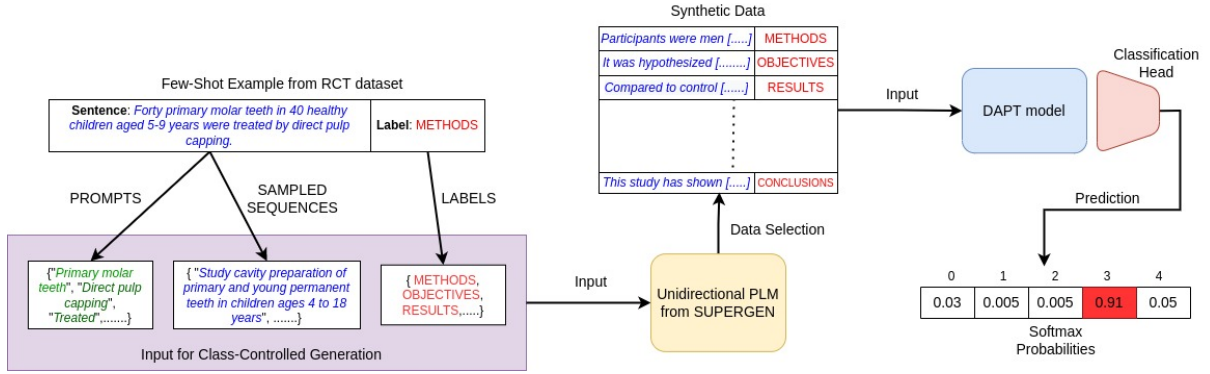


Figure 1: Pipeline for fine-tuning a DAPT model on synthetic data. Few-Shot examples from the RCT dataset are used to get prompts, sampled sequences and labels. Then a Unidirectional PLM from SuperGen generates synthetic examples which are used in the DAPT model fine-tuning. A feedforward layer is attached on to the DAPT for classification

erated texts closer to the domain of the task (Brown et al., 2020; Meng et al., 2022).

2 Related Work

2.1 Domain Adaptation

Domain adaptation can be thought of as a special kind of transfer learning where we want the model to learn domain-invariant features just like we want the model to learn task-invariant features in traditional transfer learning (Ganin et al., 2015, Tzeng et al., 2017, Khaddaj and Hajj, 2020). Domain Adaptation methods allow us to fine-tune models using easily accessible data that is not in the target domain of the task and allow models to generalize better on out-of-domain data. In recent years, many such adaptive techniques (Daumé III; Ramponi and Plank, 2020) with data selection strategies (Gururangan et al.) and representation learning (Ganin et al., 2015), have demonstrated good performance in fine-tuning PLMs without sufficient data in target domain. We closely follow the work done by Gururangan et al. and keep it as one of our baseline models, since they show that domain adaptive pretraining improves performance.

2.2 Few-shot Learning

Modern day PLMs depend on a large amount of labeled and unlabeled data for solving NLU tasks. In comparison, few-shot learning methods try to make the best use of a small amount of task-specific training data, which is a more realistic scenario. We generate data in a true few-shot setting (Perez et al., 2021), which means we don’t have access to any unlabeled data or large held-out datasets for prompt

and hyperparameter tuning. In this scenario, to improve data efficiency, prompt based methods are largely used (Brown et al., 2020, Gao et al., 2021, Liu et al., 2021, Schick and Schütze, 2021a). Many approaches have used a meta-learning strategy to tackle data sparsity (Finn et al., 2017, Mishra et al., 2018). These approaches assume access to some auxiliary tasks, which helps them extract some transferable knowledge to learn the target problem. We use prompt based methods in a true few-shot learning setting.

2.3 Controlled Data Generation

Controlled text generation (Hu et al., 2017) aims to generate text that follows the distribution of the task dataset. PLMs can be made to generate desired sentiments or topics, also called high-level control (Ziegler et al., 2019), or they can be made to generate specific words or phrases (Chan et al., 2021), also called low-level control. Both high-level and low-level control can also be achieved together (Khalifa et al., 2021). It is also possible to do controlled text generation at inference time, without further training the PLM (Dathathri et al., 2020; Krause et al., 2021; Kumar et al., 2021) Natural language prompts can be used as a guidance source for the model to generate text with desired attributes (Schick and Schütze, 2021b). We use prompts to guide text generation in this work.

3 Methods

In this work, we begin with synthetic data generation for fine-tuning on a task. We assume access to a unidirectional PLM \mathcal{G} (Meng et al., 2022). The technique uses carefully crafted prompts for

different downstream tasks to generate sequences. The PLM generates augmented task dataset \mathcal{T} which, after cleaning, is used to fine-tune a domain adapted PLM \mathcal{D} .

We compare the baselines ROBERTA model and DAPT, with a DAPT model (Gururangan et al.) fine tuned using SuperGen (Meng et al., 2022) inspired generator in a few-shot setting on 8 different tasks. The above helps to discern about how much does training on synthetic data improves performance of a domain adapted model. Another comparison of interest for this study is between ROBERTA model fine-tuned using SuperGen and the DAPT model. This tells us about how does training on only task-specific synthetic data compare to a domain adapted model.

3.1 Task-specific Data

We evaluate the performance of our proposed approach in 8 different classification tasks across 4 domains, two tasks in each domain (Gururangan et al.). We use two datasets that are in the biomedical domain: ChemProt (Kringelum et al., 2016) and PubMed 200k RCT (Dernoncourt and Lee); Two in Computer Science domain: ACL articles on NLP (ACL-ARC) (Jurgens et al., 2018) and SciERC (Luan et al., 2018); Two in News articles domain: HyperPartisan News (Kiesel et al., 2019) and AGNews corpus (Zhang et al., 2015); And two in Movie Reviews domain: Helpfulness (McAuley et al., 2015) and IMDB reviews (Maas et al., 2011).

3.2 Data Generation

For each task, we randomly sample from the task-specific dataset, and use these sampled data to form prompts which are then fed into the generator \mathcal{G} to produce synthetic data. We take advantage of the few-shot examples which are used to sample these prompts, sequences and labels. Labels are passed along with prompts for class-controlled generation. The generated data then is selected based on the ranking score as stated in (Meng et al., 2022, Yuan et al., 2021) to form the synthetic training data, which later be used to fine-tune classification PLM. Due to space constraint, we do not restate SuperGen algorithms here, and treat it as black-box. Detailed descriptions can be accessed in (Meng et al., 2022). However, we still fine-tune hyperparameters including repetition penalties, temperature, and generating size. Repetition penalties prevent repetition loops in generation and controls beneficial token repetitions in sequence pair generation,

temperature. Temperature varies the sharpness of generating distribution. The generating size involves trade-off between generating consumption and quality selection, given the data-selection technique remains the same.

3.3 Models for Task-specific Fine-tuning

- **Baseline** For each task, we use the vanilla RoBERTa (Liu et al., 2019) and DAPT (Gururangan et al.) model and fine-tune them on the complete task specific datasets as our baseline models. The DAPT model serves as good baseline model as it was proven to have close to state-of-art performance on the tasks with or without further pretraining.
- **SuperGen Fine-tuned** Instead of fine-tuning DAPT model with just the complete task-specific dataset as in baseline, we separately fine-tune multiple DAPTs using datasets augmented using SuperGen. For low resource settings, we found that generating close to 100% of the size of existing data is best. This number is close to 20% for tasks that already have a large training corpus.

Domain	Tasks	RoB _A	DAPT	SuperGen + DAPT
BM	CHEMPROT	81.1	81.8	75.3
	RCT	86.6	87.5	86.1
CS	ACL-ARC	61.0	69.4	53.7
	SCIERC	72.3	78.1	59.3
NEWS	HYP.	85.0	86.0	81.8
	AGNEWS	92.9	93.1	92.5
REV.	HELPFUL.	65.0	65.7	66.3
	IMDB	93.4	95.3	89.2

Table 1: Comparison of results of the baseline models ROBERTA and DAPT against our method on 4 different domains each containing 2 tasks. The scores reported are macro-F₁ (except for CHEMPROT and RCT which follow micro-F₁, following Beltagy et al., 2019). Best task performance is boldfaced.

4 Experiments

As described in section 3.1, we evaluate our approach on task data used in Gururangan et al.. Following the implementation in Meng et al., 2022, we use CTRL (Keskar et al., 2019) as our generator PLM \mathcal{G} . However, we use the domain adapted model DAPT (Gururangan et al.) as the discriminator since that is our baseline.

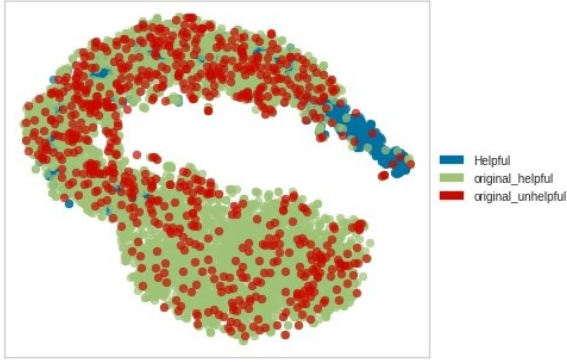


Figure 2: Domain overlap between the HELPFULNESS (McAuley et al., 2015) data and the data generated by CTRL for the same task.

To generate controlled text for each task, we use prompts starting with control codes as described in Keskar et al., 2019. We find 2.0 and 1.2 work best for temperature and repetition penalty respectively.

We further fine-tune the DAPT and pretrained ROBERTA models on the above task specific synthetic data generated using CTRL. We use the same technique as Gururangan et al. and fine-tune for 10 epochs with a learning rate of $2e - 5$ and a dropout of 0.1. A more detailed description of hyperparameters and prompts used can be found in appendix A. All experiments were conducted on NVIDIA V100 GPUs available on the Greene cluster.

5 Results

We evaluate our methodology for the data generation and classification tasks.

5.1 Analyzing Generated Data

We first visualize the difference between the data generated by CTRL (Keskar et al., 2019) and original data. For each sentence, we use CountVectorizer to output an embedding vector. Once we have all the embeddings we further use t-SNE (van der Maaten and Hinton, 2008) to visualise the respective clusters. The t-SNE plot has been shown in 2. We note that there is a significant overlap between the two distributions thereby generating valid task-specific examples. For further use for the classification tasks, we filter the generated dataset by removing out-of-distribution examples.

5.2 Downstream Performance Comparison

Once we have the clean synthetic dataset we proceed with the fine-tuning on downstream tasks. We perform external evaluation using macro- F_1 as the

measure. As can be seen in table 1, augmented set seems to hurt performance on most classification tasks. We hypothesize that this is due to the gap in the generator PLM’s pretraining distribution and the target distribution, since the accuracy doesn’t drop but slightly increases in instances where there is some overlap between the source and the target domains (figure 2). Also, this overlap is expected since HELPFULNESS (McAuley et al., 2015) belongs to the pretraining corpus of CTRL.

6 Conclusion and Future Work

We investigate a generative variation to adapt pre-trained language models for domains which have low amount of labeled data available. We propose a model which augments the downstream task data using zero-shot and few-shot prompt based data generation techniques on the CTRL (Keskar et al., 2019) model. Our experiments reveal that baseline models like RoBERTa (Liu et al., 2019) which are already built upon million of parameters still struggle to encapsulate the distribution of a particular textual domain let alone all of the natural language. Our domain adaptive model shows strong performance in cases where the generator PLM can match the target domain as closely as possible. With large generative models becoming more accessible, this scenario is quite likely possible because of the vast and diverse nature of pretraining datasets. Hence, we show that training data augmentation can be marginally beneficial in certain situations.

In future, we would like to explore more ways of doing controlled labeled data generation by closely following the works of Khalifa et al., 2021 and Schick and Schütze, 2021b. We can also compare other medium sized auto-regressive language models like GPT-2 (Radford et al., 2019) to get a better picture. Expert guided seed selection from the input data for few-shot learning and proper pre-processing can help in improving performance.

7 Ethical Considerations

As PLMs exhibit intriguing few-shot learning potential and are able to leverage task-specific information, they are used for wide range of applications like data generation/augmentation (Radford et al., 2019; Keskar et al., 2019), text classification (Liu et al., 2019) and many others. We take this property into consideration and propose a model which employs domain adaptation methods to achieve good downstream task performance even if the available

data is very less.

In our work, we use CTRL (Keskar et al., 2019) auto-regressive model which use way less parameters than the huge GPT-3 (Brown et al., 2020) model. This lines up with our objective to reduce overall computation needs and essentially being a bit more environment friendly (Rolnick et al., 2019) compared to GPT-3 in terms of electrical power consumption and carbon emission. Our work also helps researchers and firms who have low availability of resources to transfer better performance to their natural language models.

But there are potential risks involved in using large scale PLMs and consequently our domain adaptation technique which is based on one. These generative models are susceptible to personal data extraction attack (Carlini et al., 2020) and can also lead to gender, racial and caste biases. Hence, as these large PLMs are based on deep neural networks which have low level of interpretability, controllable text generation is becoming a challenging task.

We see various opportunities for research to apply our domain adapted technique to improve upon a niche downstream task domain. But to mitigate the risks of using a generative model, we encourage researchers to understand the limitations of pre-trained autoregressive models (Solaiman et al., 2019) as well as not put undue trust on such large automated models and keep monitoring and testing them at various stages.

8 Collaboration Statement

All members worked together on brainstorming the idea, held weekly discussions, drafted the report and wrote the code scripts. Akash worked on getting the results for IMDB and RCT datasets. Shubhankar worked on CHEMPROT and SCIERC. Abhinav worked on ACL-ARC and AGNEWS and Victor worked on HELPFULNESS and HYPERPARTISAN.

References

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.

Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. Ctrl: A conditional transformer language model for controllable generation, 2019. URL <https://arxiv.org/abs/1909.05858>.

Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, et al. Don’t stop pretraining: Adapt language models to domains and tasks.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, et al. Language models are few-shot learners. *ArXiv*, abs/2005.14165, 2020.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding. *arXiv e-prints*, art. arXiv:1804.07461, April 2018.

Jeremy Howard and Sebastian Ruder. Universal language model fine-tuning for text classification. In *ACL*, 2018.

Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot learners. *ArXiv*, abs/2012.15723, 2021.

Teven Le Scao and Alexander M. Rush. How many data points is a prompt worth? In *NAACL*, 2021.

Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, et al. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics*, 36(4):1234–1240, 09 2019. ISSN 1367-4803. doi: 10.1093/bioinformatics/btz682. URL <https://doi.org/10.1093/bioinformatics/btz682>.

Timo Schick and Hinrich Schütze. It’s not just size that matters: Small language models are also few-shot learners. *CoRR*, abs/2009.07118, 2020. URL <https://arxiv.org/abs/2009.07118>.

Yu Meng, Jiaxin Huang, Yu Zhang, and Jiawei Han. Generating training data with language models: Towards zero-shot language understanding. 02 2022.

Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. 2015. doi: 10.48550/ARXIV.1505.07818. URL <https://arxiv.org/abs/1505.07818>.

Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2962–2971, 2017.

- Alaa Khaddaj and Hazem M. Hajj. Representation learning for improved generalization of adversarial domain adaptation with text classification. *2020 IEEE International Conference on Informatics, IoT, and Enabling Technologies (ICIoT)*, pages 525–531, 2020.
- Hal Daumé III. Frustratingly easy domain adaptation.
- Alan Ramponi and Barbara Plank. Neural unsupervised domain adaptation in nlp—a survey, 2020.
- Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models. *ArXiv*, abs/2105.11447, 2021.
- Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too. *ArXiv*, abs/2103.10385, 2021.
- Timo Schick and Hinrich Schütze. Exploiting cloze-questions for few-shot text classification and natural language inference. In *EACL*, 2021a.
- Chelsea Finn, P. Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.
- Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and P. Abbeel. A simple neural attentive meta-learner. In *ICLR*, 2018.
- Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P. Xing. Toward controlled generation of text. In *ICML*, 2017.
- Daniel M. Ziegler, Nisan Stiennon, Jeff Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *ArXiv*, abs/1909.08593, 2019.
- Alvin Chan, Y. Ong, Bill Tuck Weng Pung, Aston Zhang, and Jie Fu. Cocon: A self-supervised approach for controlled text generation. *ArXiv*, abs/2006.03535, 2021.
- Muhammad Khalifa, Hady ElSahar, and Marc Dymetman. A distributional approach to controlled text generation. *ArXiv*, abs/2012.11635, 2021.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text generation. *ArXiv*, abs/1912.02164, 2020.
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq R. Joty, Richard Socher, and Nazneen Rajani. Gedi: Generative discriminator guided sequence generation. *ArXiv*, abs/2009.06367, 2021.
- Sachin Kumar, Eric Malmi, Aliaksei Severyn, and Yulia Tsvetkov. Controlled text generation as continuous optimization with multiple constraints. *ArXiv*, abs/2108.01850, 2021.
- Timo Schick and Hinrich Schütze. Few-shot text generation with natural language instructions. In *EMNLP*, 2021b.
- Jens Kringelum, Sonny Kim Kjaerulff, Søren Brunak, Ole Lund, Tudor I. Oprea, and Olivier Taboureau. Chemprot-3.0: a global chemical biology diseases mapping. *Database : the journal of biological databases and curation*, 2016:bav123, Feb 2016. ISSN 1758-0463. doi: 10.1093/database/bav123. URL <https://pubmed.ncbi.nlm.nih.gov/26876982>. 26876982[pmid].
- Franck Dernoncourt and Ji Young Lee. PubMed 200k RCT: a dataset for sequential sentence classification in medical abstracts.
- David Jurgens, Srijan Kumar, Raine Hoover, Dan McFarland, and Dan Jurafsky. Measuring the evolution of a scientific field through citation frames. *Transactions of the Association for Computational Linguistics*, 6:391–406, 2018. doi: 10.1162/tacl.a.00028. URL <https://aclanthology.org/Q18-1028>.
- Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3219–3232, Brussels, Belgium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/D18-1360. URL <https://aclanthology.org/D18-1360>.
- Johannes Kiesel, Maria Mestre, Rishabh Shukla, Emmanuel Vincent, Payam Adineh, David Corney, Benno Stein, and Martin Potthast. SemEval-2019 task 4: Hyperpartisan news detection. In *Proceedings of the 13th International Workshop on Semantic Evaluation*, pages 829–839, Minneapolis, Minnesota, USA, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/S19-2145. URL <https://aclanthology.org/S19-2145>.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification, 2015. URL <https://arxiv.org/abs/1509.01626>.
- Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton van den Hengel. Image-based recommendations on styles and substitutes, 2015. URL <https://arxiv.org/abs/1506.04757>.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL <https://aclanthology.org/P11-1015>.

- Weizhe Yuan, Graham Neubig, and Pengfei Liu. Bartscore: Evaluating generated text as text generation. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 27263–27277. Curran Associates, Inc., 2021. URL <https://proceedings.neurips.cc/paper/2021/file/e4d2b6e6fdeca3e60e0f1a62fee3d9dd-Paper.pdf>.
- Iz Beltagy, Kyle Lo, and Arman Cohan. Scibert: A pretrained language model for scientific text. 2019. doi: 10.48550/ARXIV.1903.10676. URL <https://arxiv.org/abs/1903.10676>.
- Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008. URL <http://jmlr.org/papers/v9/vandermaaten08a.html>.
- David Rolnick, Priya L. Donti, Lynn H. Kaack, Kelly Kochanski, Alexandre Lacoste, Kris Sankaran, Andrew Slavin Ross, Nikola Milojevic-Dupont, Natasha Jaques, Anna Waldman-Brown, Alexandra Luccioni, Tegan Maharaj, Evan D. Sherwin, S. Karthik Mukkavilli, Konrad P. Körding, Carla P. Gomes, Andrew Y. Ng, Demis Hassabis, John C. Platt, Felix Creutzig, Jennifer T. Chayes, and Yoshua Bengio. Tackling climate change with machine learning. *CoRR*, abs/1906.05433, 2019. URL <http://arxiv.org/abs/1906.05433>.
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. Extracting training data from large language models. *CoRR*, abs/2012.07805, 2020. URL <https://arxiv.org/abs/2012.07805>.
- Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, and Jasmine Wang. Release strategies and the social impacts of language models. *CoRR*, abs/1908.09203, 2019. URL <http://arxiv.org/abs/1908.09203>.

A Hyperparameters

As stated before, we use CTRL (Keskar et al., 2019) as the generator PLM for controlled text generation and ROBERTA as well as DAPT (Gururangan et al.) as discriminator models for the downstream classification task. Here, we describe the hyperparameters used in our experiments.

A.1 Controlled Text Generation

We found that a temperature of 2.0 and a repetition penalty of 1.2 gave the best generation results across all tasks. Table 2 describes the prompts used for each task. For the tasks not included in the table, we found it extremely difficult to generate coherent class conditioned texts and simply used the class labels and the domain as prompts (domain was used as a control code (Keskar et al., 2019)).

Task	Control Code(s)	Label	Prompt(s)
IMDB Reviews	Reviews	1 0	Rating: 5.0 Rating: 1.0
Hyperpartisan News	Politics, Conspiracy	true false	It is true that; In other news There is a rumor that
Amazon	Reviews	Helpful Unhelpful	Helpful review Unhelpful review
RCT	Science	CONCLUSIONS OBJECTIVE RESULTS METHODS BACKGROUND	The conclusion drawn from this medical research is that The goal of this experiment is The result of this experiment This research uses the following method The background of this reasearch is

Table 2: A sample of prompts used to generate class-conditioned training samples. All tasks are single-sequence classification tasks. IMDB prompts are the same used in CTRL (Keskar et al., 2019). We tried having more than alternative prompts for some tasks. Such prompts are separated by a semicolon.

A.2 Text Classification

Following the hyperparameter seeds in Gururangan et al., we fine-tune discriminator models ROBERTA and DAPT, for 5 – 10 epochs with a learning rate of $2e - 5$. We set the batch size to 16, and the dropout to 0.1. Since we use saved pretrained models as our starting point, we do not perform the pretraining step for ROBERTA and DAPT models.