
Large Language Model as Attributed Training Data Generator: A Tale of Diversity and Bias

Yue Yu^{1*}, Yuchen Zhuang^{1*}, Jieyu Zhang^{2*}, Yu Meng³,
Alexander Ratner², Ranjay Krishna², Jiaming Shen⁴, Chao Zhang¹
¹ Georgia Institute of Technology ² University of Washington
³ University of Illinois at Urbana-Champaign ⁴ Google Research
{yueyu, yczhuang, chaozhang}@gatech.edu, yumeng5@illinois.edu
{jieyuz2, ajratner, ranjay}@cs.washington.edu, jmshen@google.com

Abstract

Large language models (LLMs) have been recently leveraged as training data generators for various natural language processing (NLP) tasks. While previous research has explored different approaches to training models using generated data, they generally rely on simple class-conditional prompts, which may limit the diversity of the generated data and inherit systematic biases of LLM. Thus, we investigate training data generation with diversely attributed prompts (e.g., specifying attributes like length and style), which have the potential to yield diverse and attributed generated data. Our investigation focuses on datasets with high cardinality and diverse domains, wherein we demonstrate that attributed prompts outperform simple class-conditional prompts in terms of the resulting model’s performance. Additionally, we present a comprehensive empirical study on data generation encompassing vital aspects like bias, diversity, and efficiency, and highlight three key observations: firstly, synthetic datasets generated by simple prompts exhibit significant biases, such as regional bias; secondly, attribute diversity plays a pivotal role in enhancing model performance; lastly, attributed prompts achieve the performance of simple class-conditional prompts while utilizing only 5% of the querying cost of ChatGPT associated with the latter. We release the generated dataset and used prompts to facilitate future research².

1 Introduction

Large language models (LLMs) have demonstrated exceptional performance across a broad range of NLP tasks [5, 38, 24, 36, 37, 62]. In recent research, LLMs have been proposed as *task-specific* training data generators, particularly for text classification, aiming to alleviate the need for task-specific data and annotations [55, 13, 56, 30, 59, 7]. While these efforts have showcased the effectiveness of LLMs as data generators, the focus has primarily been on advancing the training stage, where the generated data are utilized to train task-specific models, leaving the upstream data generation process relatively unexplored. Notably, the prevailing approach employs a simple class-conditional prompt for querying LLMs during data generation, potentially limiting the diversity of the generated data [7, 51, 60] and inheriting systematic biases inherent in LLMs [65, 21]. We refer to this simple class-conditional prompt as SimPrompt, providing an example in Table 1.

In this work, we ground the LLM to ChatGPT [37]³ for its ability to generate high-quality, human-like text [25], and consider four challenging topic classification tasks with high cardinality from various

*These authors contributed equally to this work.

²The data and code is available on <https://github.com/yueyu1030/AttrPrompt>.

³We use gpt-3.5-turbo in our main experiments.

Table 1: Prompt template for the NYT news dataset.

Method	Prompt
SimPrompt	Suppose you are a news writer. Please generate a {topic-class} news in NYT.
AttrPrompt	Suppose you are a news writer. Please generate a {topic-class} news in NYT following the requirements below: 1. Should focus on {subtopic}; 2. Should be in length between {length:min-words} and {length:max-words} words; 3. The writing style of the news should be {style}; 4. The location of the news should be in {location}.

domains. Our investigation primarily revolves around assessing the bias and diversity present within the generated training set through the lens of *data attributes*. In particular, data attributes encompass multiple attribute dimensions and their corresponding attribute values, where the latter represent possible instantiations of the former. For example, an attribute value such as “*shorter than 200 words*” could serve as an instantiation of the attribute dimension “*length*”.

On one hand, we employ a trained attribute classifier to examine the *attribute bias* present in the dataset generated using SimPrompt. When analyzing the “*location*” attribute in the NYT news dataset, we observe a striking bias towards “*North America*” in the predicted values of the generated data, accounting for a significant majority (68.01%). In contrast, instances associated with “*Africa*” are remarkably rare, comprising only 0.69% of the dataset (100 times less prevalent than “*North America*”). This regional bias exhibited in the generated dataset can pose substantial challenges when constructing reliable machine learning models [23, 6].

On the other hand, we explore the influence of *attribute diversity* on the downstream model performance. Specifically, we leverage ChatGPT to generate attributed data by incorporating desired attributes as constraints in the prompts. By comparing the performance of models trained on datasets generated using prompts with random attributes against those with fixed attributes, we observe a substantial underperformance of the latter, uncovering the importance of attribute diversity of the generated dataset.

To alleviate attribute biases and enhance the attribute diversity of the generated data, we propose to generate data with diversely attributed prompts. For a given classification task, we start by identifying attribute dimensions and their corresponding attribute values in an interactive, semi-automated process facilitated by the LLM. Subsequently, we generate diverse prompts by combining attributes randomly, replacing the simple class-conditional prompt typically used for querying data from the LLM. We refer to these diversely attributed prompts as AttrPrompt. An example of such prompts can be found in Table 1, where the LLM is instructed to generate training data based on attributes such as *location* and *style*.

On the four classification tasks, we empirically evaluate the generated datasets by measuring the performance of models trained using two scenarios: 1) solely on the generated dataset, and 2) on a merged dataset comprising the real training set and the generated set. In both scenarios, the dataset generated with AttrPrompt significantly outperforms its counterpart generated with SimPrompt. Furthermore, we demonstrate the superiority of AttrPrompt over SimPrompt in terms of data/budget efficiency and compatibility with different model sizes/various LLM-as-training-data-generator approaches. Notably, AttrPrompt achieves the performance of SimPrompt while utilizing only 5% of the querying cost of ChatGPT associated with SimPrompt. Lastly, we extend the LLM-as-training-data-generator paradigm to the more challenging *multi-label classification* tasks for the first time, and AttrPrompt outperforms SimPrompt across all evaluation metrics.

2 Related Work

LLMs as Training Data Generators. With the remarkable success of large language models (LLMs), researchers have recently attempted to leverage them as training data generators. Such applications include generating tabular data [4], relation triplets [8], sentence pairs [46], instruction data [40, 50, 53, 47], *etc.*. Among these applications, we anchor on training data generation for topic classification in a zero-shot setting where no labeled data is available. In this direction, existing approaches typically use simple class-conditional prompts while focusing on mitigating low-quality issues after generation. Initial explorations in this domain include SuperGen [30] and ZeroGen [55], which use LLMs for text classification and noise robust learning techniques [35, 52] to handle data

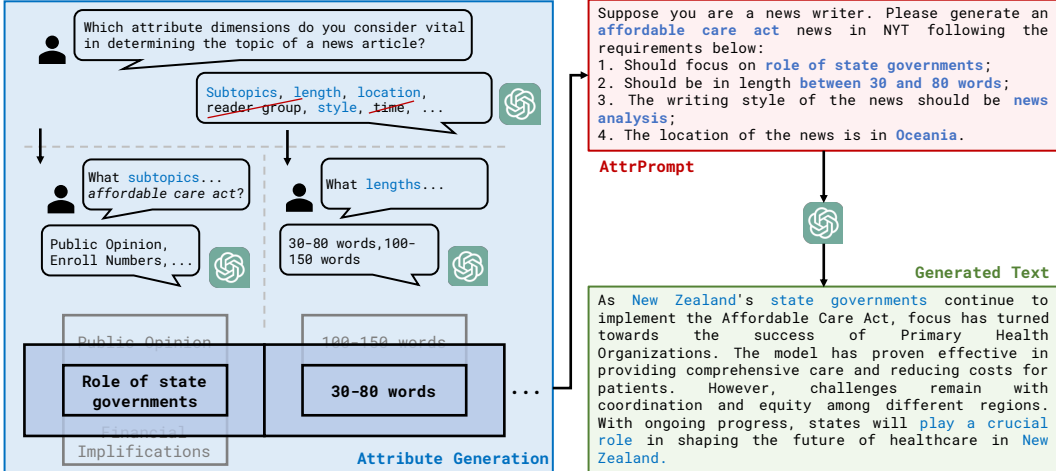


Figure 1: The overall workflow of AttrPrompt.

quality issues. SunGen [13] reweights the generated data during training with learned data quality weight, and ProGen [56] selects highly influential generated data via model feedback. In this work, we instead explore attributed prompts to reduce the issue of low informativeness and redundancy, which can be readily incorporated into the existing systems mentioned above. Notably, Chen et al. [7] also explore prompts to advance the data generation process, yet it adopts soft prompts and requires a white-box LLM and seed examples to tune them. In contrast, our method is applicable to black-box LLMs and even LLM APIs (e.g., ChatGPT) and does not rely on any labeled examples. A recent work WANLI [26] also considers human-AI collaboration for creating more challenging training data, but requires an initial dataset and a strong task model. Instead, we aim to generate training data without any initial dataset or a pre-existing task model, which allows us to effectively handle resource-limited scenarios.

Attribute-aware Text Generation. There are also several existing works [27, 44, 57] that incorporate attributes for controlled text generation, but these are concentrated on very different tasks like style transfer. Typically, these methods necessitate explicit provision of attributes. Differently, we introduce a semi-automated strategy that allows LLMs to propose attribute values autonomously.

Discrete Prompt Optimization. Several works attempt to optimize discrete prompts for querying LLMs with large language models [43, 64, 34]. More related to us, [33, 20] reframe prompts by decomposing a complex task instruction into multiple simple ones. However, these approaches mainly focus on the *inference* stage for directly predicting the answer and may rely on additional labeled examples for validation. Our focus is on an orthogonal setting, optimizing prompts for LLMs with attributes to diversify the generated training data. This approach improves the model’s overall performance without the need for additional labeled examples.

3 Large Language Model as Attributed Training Data Generator

In this section, we present the design of our proposed method, AttrPrompt. This technique employs class-conditional attributes as an enhancement to the query prompts employed in Large Language Models (LLMs). These augmented prompts enable more effective data generation for training purposes. A detailed workflow of the AttrPrompt can be referenced in Figure 1.

3.1 Datasets

While previous research has primarily focused on binary classification datasets [55, 30, 56] or datasets containing a maximum of 14 classes [13, 59], the performance of LLM as a data generator for topic

classification with high cardinality (*i.e.*, many topic classes) remains unclear. Thus, we consider the following datasets from various domains with the number of topics ranging from 23 to 50⁴:

- **NYT [31]**: The NYT dataset comprises news articles that were authored and published by *The New York Times*. These articles are categorized into 26 fine-grained categories.
- **Amazon [3]**: The Amazon dataset contains customer reviews on products from Amazon’s online store. It covers products from 23 different categories.
- **Reddit [15]**: The Reddit dataset consists of a vast collection of user-generated content from the popular social media platform Reddit. It encompasses a wide range of topics, discussions, and interactions among users across numerous communities.
- **StackExchange [15]**: The StackExchange dataset is a rich collection of structured data encompassing various online communities and knowledge-sharing platforms. It contains a vast array of questions, answers, comments, tags, and user interactions about specific technical problems.

We summarize the statistics of used dataset in Table 2, from which we can see that the involved datasets not only have high cardinality but also come with high imbalance ratio, *i.e.*, the ratio of the sample size of the majority class to that of the minority class, which reflects the long-tail class issue in real applications [1].

Table 2: Statistics of datasets.

Dataset	Domain	Task	# Train	# Valid	# Test	# Class	Imbalance Ratio
Amazon [3]	Reviews	Multi-class	15.0K	0.2K	1.2K	23	155.6
NYT [31]	News	Multi-class	9.0K	0.2K	1.2K	26	357.0
Reddit [15]	Web	Multi-class	26.6K	0.2K	2.3K	45	447.4
StackExchange [15]	Web	Multi-class	27.0K	0.3K	2.5K	50	1283.7

3.2 Interactive Attribute Generation

Different from the existing works [30, 55, 13] that directly use the simple class-conditional prompts for querying LLMs, our initial step involves identifying various types of data attributes (or metadata) that can be manipulated to generate attributed data samples. To facilitate this process, we employ ChatGPT to help establish both attribute dimensions and attribute values. Specifically, we begin by engaging ChatGPT in generating essential attribute dimensions. This is achieved by posing questions such as “Which attribute dimensions do you consider vital in determining the topic of a news article?” for the NYT dataset, resulting in responses like “*subtopics, length, location, reader group, style, time*”. Then, we adopt the *human-ai collaboration* scheme [26, 54, 61] to interactively select the attribute dimensions of the highest quality that best suit the dataset. Similarly, we prompt ChatGPT (the prompt format is listed in Appendix E) to suggest potential attribute values within each attribute dimension and choose high-quality candidates.

Attribute dimensions and values. There are two types of attribute dimensions: *class-independent* attributes and *class-dependent* attributes. Class-independent attributes, such as “*length*”, remain unchanged across different classes, while class-dependent attributes, like “*subtopic*”, have varying attribute values for each class. We list attribute dimensions and values for all datasets in Table 3. These data attributes provide a human-manipulatable interface for generating attributed data. In this study, we explore the potential of leveraging attributes to enhance the data generation process, while leaving the search for the optimal data attributes for a specific task to future work.

Class-dependent attribute value filtering. When dealing with class-dependent attributes, it is crucial to ensure that their attribute values are specifically associated with the corresponding class to avoid ambiguity and potential connections to multiple classes. For example, in the case of the “*economy*” class in the NYT dataset, a candidate attribute value generated by ChatGPT for the “*subtopic*” could be “*effect of trade tariffs on manufacturing companies*”, which is also relevant to the “*international business*” class in the NYT. This overlap may introduce ambiguity in the generated data. To address this issue, we employ a filtering process called Class-Dependent Attribute Value Filtering (CAF). First, we query ChatGPT for the top-5 similar classes and then check with ChatGPT whether

⁴We also conduct experiments on those datasets widely used in prior works [55, 56, 13, 59], the details can be found in Appendix C.4.

Table 3: Attribute dimensions and values. Attributes with an asterisk* are class-dependent attributes.

Dataset	# configurations / class	Attribute dimension	Attribute value
NYT	600	Subtopic*	Appendix G.1.1
		Location	Asia, North America, South America, Africa, Oceania, Europe
		Writing Style	Investigative journalism, Op-Eds, Feature writing, News analysis, Profiles and interviews
		Length	short (30-80 words); long (100-150 words)
Amazon	1000	Product Brands*	Appendix G.2.1
		Product Names*	Appendix G.2.2
		Usage Experience	Worst, Bad, Average, Good, Excellent
		Writing Style	Detailed Review; Comparative Review; Pros and Cons Review; Recommendation Review
Reddit	500	Length	short (30-80 words); long (100-150 words)
		Resources*	Appendix G.3.1
		Experience*	Appendix G.3.2
		Writing Style	Informative/Educational; Entertaining/Funny; Discussion; Storytelling; Help/Advice
StackExchange	400	Length	short (30-80 words); long (100-150 words)
		Scenario*	Appendix G.4.1
		Technical Depth	Beginner; Intermediate; Advanced; Expert
		Writing Style	Specific; Comparative; Problem-Solution; Troubleshooting; Tutorial
		Length	short (30-80 words); long (100-150 words)

each class-dependent attribute value is related to these top-5 similar classes. Then, if the answer is positive which indicates a potential ambiguity, we remove that attribute value for the specific class.

3.3 Data generation and model training

Given the data attributes, one could prompt LLMs to generate data samples with diverse *attribute configurations*. For example, an attribute configuration for the “*federal budget*” class of the NYT dataset could be {“*subtopic*”=“*defense spending*”, “*length*”=“*short:min-words=30, max-words=80*”, “*style*”=“*investigative journalism*”, “*location*”=“*North America*”}. In Table 3, we list the number of configurations per class, and one can further expand the number of configurations by adding more attribute dimensions and values. To generate attributed data samples, we prompt ChatGPT with random configurations. In particular, each time we generate a random configuration, complete a *prompt template* (see Table 1) with the generated configuration, and query ChatGPT with the completed prompt to collect generated data samples.

Throughout the experiments, we compare our method (AttrPrompt) against simple class-conditional prompt (SimPrompt, [55]) and the original training set of each dataset (Gold). For a fair comparison, we set the number of generated data the same as Gold for both AttrPrompt and SimPrompt. In principle, the generated dataset can be combined with any classifier (Sec. 6.4) and training techniques (Sec. 6.5); if not otherwise specified, we choose to fine-tune BERT-base-uncased [11] as the backbone and use the standard cross-entropy loss by default. For hyperparameter selection, we adhere to the recommendation in [41, 30] for SimPrompt and AttrPrompt, and *do not use the validation set* for model selection. Detailed hyperparameter configurations can be found in Appendix B.

A glimpse of the generated data. Here, we present examples of data generated by AttrPrompt and SimPrompt, and real data from the Gold set for the “*federal budget*” class in the NYT dataset (Table 4). It is evident that the data generated by ChatGPT exhibit high quality. Particularly, when comparing AttrPrompt to SimPrompt, we observe that AttrPrompt renders more diverse samples. This is because SimPrompt tends to generate news focused on the U.S., while AttrPrompt has the capability to generate news from various locations around the world.

4 Diversity Analysis of the Generated Data

Quantitative study of diversity. To quantify the diversity of the generated training data of SimPrompt and AttrPrompt, we first show the *vocabulary size* of the generated dataset and the Gold dataset, which is a natural way to check the lexical diversity of datasets (Table 5). From the table, we can see that AttrPrompt has higher lexical diversity than SimPrompt in terms of both vocabulary size of the whole dataset (All in the table) and the averaged vocabulary size across classes (Class Avg. in the table). Yet, both have much smaller vocabulary sizes than the Gold, indicating there is still room for improvement of the ChatGPT’s lexical diversity.

Table 4: Data examples of different datasets: the “*federal budget*” class of the NYT dataset.

Method	Example
Gold	<ol style="list-style-type: none"> 1. The emergency manager who was sent to reverse the fortunes of this financially troubled city asked some of its creditors on Friday to accept pennies on the dollar as he laid out his plan for tackling Detroit’s staggering debt, kick starting negotiations that could determine whether the city is headed to bankruptcy court... 2. Saying that its debt could reach 45 billion by 2017 if Congress does not act, the Postal Service on Wednesday called on lawmakers to give it the flexibility to change its business model to keep itself solvent. 3. Governmental Affairs Committee, Patrick R. Donahoe, the postmaster general, asked Congress to give the Postal Service permission to run its own health plan for employees and retirees, modify a Congressional mandate that requires the agency to pay...
SimPrompt	<ol style="list-style-type: none"> 1. Washington D.C. The United States government has recently passed its federal budget for the next fiscal year, setting aside a whopping 4.8 trillion for spending. This is the largest budget in the country’s history and reflects the government’s commitment to improving the country’s economy and infrastructural development. 2. WASHINGTON D.C. The Trump administration released its proposed budget for the federal government on Monday, calling for cuts to programs that provide assistance to low income Americans while increasing military spending. The proposed budget includes a 10 percent increase in defense spending, bringing the total military budget to 574 billion... 3. Washington D.C. The U.S. federal government has presented its budget for the fiscal year 2022, which includes a total of 6 trillion in spending and initiatives aimed at making the country more competitive and sustainable. The budget plan includes significant investments in infrastructure, healthcare, education, and renewable energy.
AttrPrompt	<ol style="list-style-type: none"> 1. In Germany, public opinion on federal budget issues has been influenced by The Federal Agency for Civic Education. Using in depth storytelling, this agency communicates budget policies to citizens, emphasizing their impact on society. Its efforts have led to greater engagement and awareness of budget issues among the public in recent years. 2. Op Eds The Government Spending Priorities in Mexico As the Mexican federal budget draws near, all eyes are on the government’s spending priorities. One area of concern is the role of the Banco de Mexico in managing the country’s finances. Many experts argue that the central bank should play a more active role in balancing the budget and stimulating economic growth... 3. TOKYO The recent federal budget proposal has triggered a wave of agency budget cuts and restructuring in Japan, with the Japan Foundation for Regional Development being one of the latest casualties. The foundation, which aims to promote regional development and revitalization, is set to have its budget slashed by 20 next year.

Table 5: Comparison of the vocabulary size of different datasets.

Method	NYT		Amazon		Reddit		StackExchange	
	All	Class Avg.	All	Class Avg.	All	Class Avg.	All	Class Avg.
Gold	70.8k	11.3k	44.7k	6.64k	50.8k	4.62k	52.3k	3.60k
SimPrompt	20.6k	3.13k	11.6k	2.50k	19.9k	3.06k	13.3k	2.20k
AttrPrompt	21.4k	3.50k	14.0k	2.76k	25.4k	3.64k	17.8k	2.93k

Table 6: Comparison of two quantitative metrics on diversity: the average pairwise sample similarity (APS) and inter-sample N-gram frequency (INGF) of different datasets. For APS, the *lower* stands for better diversity. For INGF, the *higher* stands for better diversity.

Method	NYT				Amazon			
	Inter-Class APS	Intra-Class APS	APS	INGF	Inter-Class APS	Intra-Class APS	APS	INGF
Gold	0.098	0.358	0.122	7618.1	0.101	0.251	0.114	4992.1
SimPrompt	0.101	0.568	0.135	5277.2	0.207	0.620	0.241	2266.5
AttrPrompt	0.159	0.474	0.182	6688.6	0.225	0.483	0.246	2605.5

Method	Reddit				StackExchange			
	Inter-Class APS	Intra-Class APS	APS	INGF	Inter-Class APS	Intra-Class APS	APS	INGF
Gold	0.044	0.261	0.054	9079.6	0.056	0.196	0.063	5492.4
SimPrompt	0.173	0.818	0.201	2697.8	0.282	0.804	0.302	2259.8
AttrPrompt	0.106	0.474	0.122	3994.5	0.105	0.375	0.114	2464.3

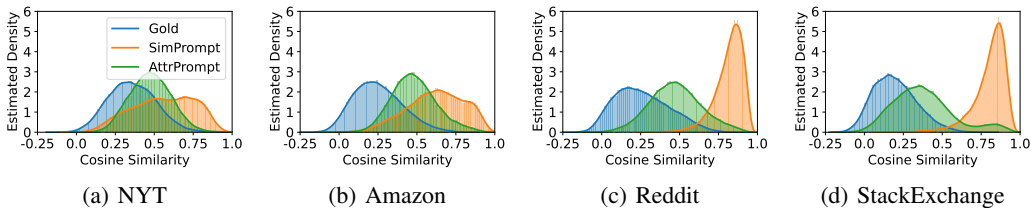


Figure 2: The distribution of cosine similarity of text pairs sampled from the same class.

We then visualize the diversity of datasets via the distribution of cosine similarity of same-class text pairs (Figure 2), where the cosine similarity is calculated based on the embedding of SentenceBERT [42], as well as including two additional metrics, namely *average pairwise sample similarity (APS)* and *inter-sample N-gram Frequency (INGF)* [32], as shown in table 6. We can see that the Gold dataset has the lowest cosine similarity, indicating that real data has the largest diversity. In contrast, the similarity between samples generated by SimPrompt is high. Compared to SimPrompt, dataset generated with AttrPrompt exhibits lower cosine similarity and the distribution is close to that of the Gold, which shows AttrPrompt could render more diverse data. Apart from the above automatic evaluation processes, we also conduct *human study* in Appendix D.1 to manually evaluate the quality of the generated training data.

The importance of the attribute diversity. We investigate the impact of attribute diversity within AttrPrompt on model performance. Specifically, we conduct experiments by fixing one attribute dimension to a candidate value while keeping other attribute values random. Then, we generate 50 data per class using such a one-fixed-others-random configuration to compose a dataset and evaluate the performance of the trained model. Note that for class-dependent attributes, we sample one value for each class and repeat it 5 times, since it is *computationally prohibitive* to enumerate all combinations of class-dependent attribute values. In Figure 3, each bar stands for a specific one-fixed-others-random configuration; compared to random configurations, most of one-fixed-others-random configurations result in a performance drop. To further reduce the attribute diversity, we pick the attribute value with the best performance for each attribute dimension (the highest bar within each attribute dimension) and compose them to a single configuration (the dashed blue line). We can see that the dashed blue line is significantly worse than the random configuration, even though it is composed of individually best attribute values. This illustrates the importance and necessity of designing prompts with diverse attributes.

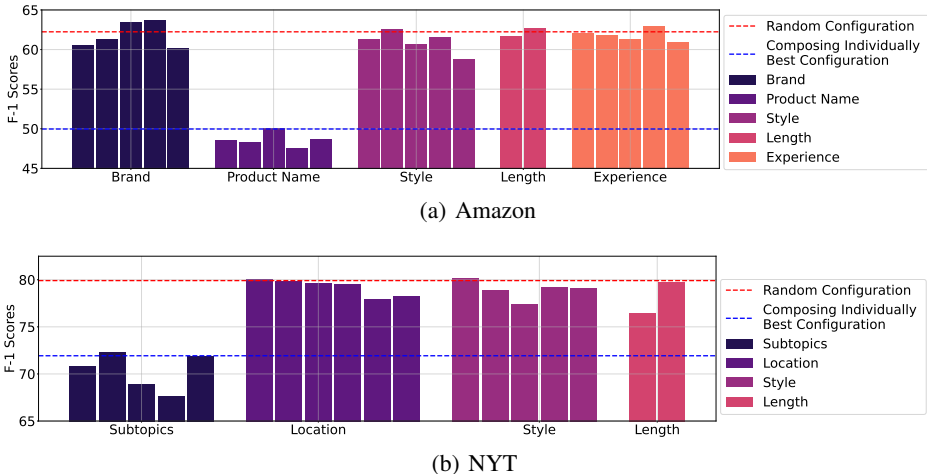


Figure 3: Bar charts of model performance with different attribute configurations of AttrPrompt.

5 Bias Analysis of the Generated Data

In this section, we study the attribute bias in both real and generated datasets. In particular, we pick the “location” attribute of the NYT data as a case study⁵. While existing works using LLMs as data generators usually overlook the bias embedded in the generated data, we hope that this preliminary analysis could raise the attention of the community to the attribute bias behind the generated data of LLMs such as ChatGPT.

We manually annotate the location for 100 samples from each of the Gold, SimPrompt, and AttrPrompt dataset. Note that we include “unknown” as an option in manual annotation to absorb text without clear location specifications. To visualize the distribution of annotated locations in datasets, we plot the pie charts in Figure 4. From the visualizations, one can see that both the Gold and SimPrompt dataset are largely biased towards “North America”, while the AttrPrompt datasets renders a relatively balanced “location” distribution.

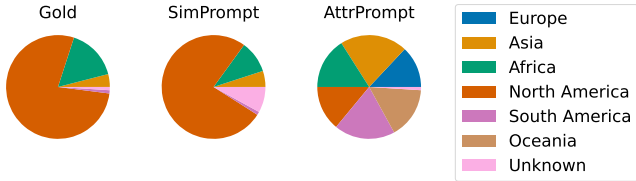


Figure 4: Pie charts of the distributions of the “location” attribute for the NYT dataset.

⁵Studies of attribute biases on other datasets can be found in Appendix D.2.

To scale up the study of attribute bias, we leverage the dataset generated by AttrPrompt as a probe. In particular, we employ the attributes associated with each data of AttrPrompt to train an *attribute classifier*, which is in turn used to make attribute predictions on Gold and SimPrompt dataset. Note that the attribute values associated with each data of AttrPrompt is not necessary the ground truth, yet since ChatGPT has shown remarkable performance in following instructions [38], the generated data could decently reflect the desired attributes and therefore the attribute classifier trained with them could partially reveal the underlying attribute distribution of tested dataset, *i.e.*, Gold and SimPrompt. In Appendix D.1, we justify the use of the attribute classifier by comparing the prediction of the attribute classifier and that of manual annotations.

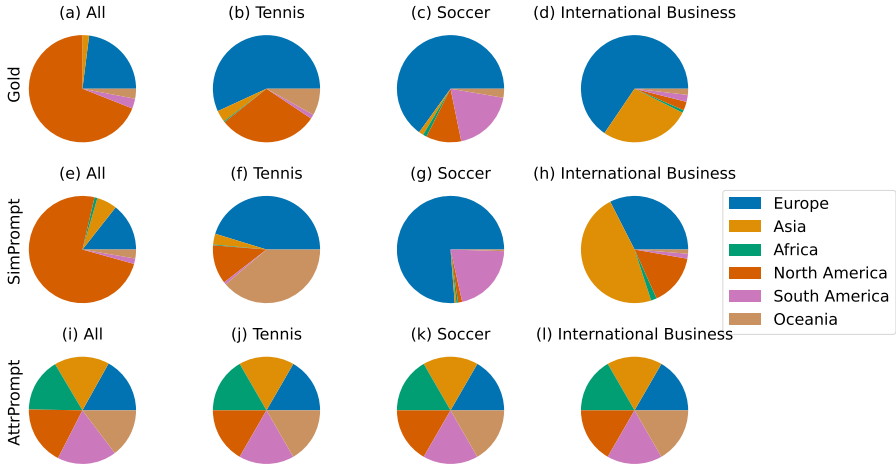


Figure 5: Pie charts of the distributions of “location” predicted by an attribute classifier for the NYT dataset. (a), (e), and (i) are “location” distributions over the whole dataset, while others are for specific classes.

We visualize the distributions of the predicted “location” in Figure 5. From the results, we can see that the “location” distribution of the whole dataset (the first column of Figure 5) is similar to that of manual annotations (Figure 4). Regarding the “location” distribution of specific classes, we can see that while the AttrPrompt still exhibits balanced distribution, the Gold and SimPrompt are biased towards continents other than “North America”. In addition, for the class “tennis”, the Gold dataset contains much more “North America” than “Oceania”, while the SimPrompt, in contrast, demonstrates an opposite trend with a higher representation of “Oceania” than “North America”. Such a noticeable disparity highlights the unpredictable nature of biases, potentially posing risks to models trained on such biased datasets.

6 Experiments on the Trained Models

6.1 Training with generated data

We quantitatively evaluate the quality of generated datasets via the test performance of models trained with them. Apart from the AttrPrompt and the direct baseline SimPrompt, we include an additional baseline MetaPrompt [43] which leverage LLM to generate additional guidance information for improving upon SimPrompt. The details for MetaPrompt are shown in Appendix J. In addition, we use ChatGPT as a zero-shot predictor for comparison. The results are in Table 7. Besides the test performance, we include the cost of querying the ChatGPT per 1000 data in the table.

From the results, we can draw the following conclusions. First, the AttrPrompt consistently renders better performance compared to the SimPrompt with a margin of 6–10 points⁶. Second, the class-dependent attribute value filter (CAF) is beneficial since the AttrPrompt outperforms its variant

⁶In Appendix C.2 we show that simply increasing the temperature t for SimPrompt does not significantly improve its performance.

without CAF⁷. Third, out of the four datasets, the AttrPrompt outperforms the LLM zero-shot method on three datasets in terms of accuracy, while for the F1 score, the AttrPrompt surpasses the LLM zero-shot on all the datasets; combined with the observation that the LLM zero-shot inference incurs much higher costs compared to data generation and the fact that the generated data is re-usable for training any model, we argue that for topic text classification generating training data could be a better practice of leveraging LLM than direct zero-shot inference. Lastly, in most cases, the generated data underperform the original training set, indicating that there is still room for future improvement. We conduct further studies in Appendix C.3 to illustrate the performance over different classes.

Table 7: Performance of the models trained with created datasets and the cost of constructing the datasets. The results are averaged over five runs. The gain of AttrPrompt has passed the statistical test with $p < 0.05$. We also include the performance and cost of using LLM as a zero-shot predictor.

Method	NYT			Amazon			Reddit			StackExchange		
	Acc.	F1	Price/1k	Acc.	F1	Price/1k	Acc.	F1	Price/1k	Acc.	F1	Price/1k
LLM Zero-Shot	74.16	69.84	5.44	59.55	54.56	2.11	67.00	56.66	2.89	44.70	43.80	3.12
Gold	83.80	81.02	—	82.23	81.12	—	84.22	83.38	—	67.56	63.28	—
SimPrompt	75.47	76.22	0.76	57.34	56.96	0.77	53.48	53.81	0.65	42.88	41.30	0.69
MetaPrompt	79.58	79.83	0.87	56.35	55.98	0.84	54.61	54.30	0.74	44.81	44.02	0.83
AttrPrompt w/o CAF	80.40	80.92	0.91	61.67	61.57	0.82	61.22	60.18	0.72	45.90	44.84	0.81
AttrPrompt	81.30	82.26	1.05	66.08	65.65	0.87	63.33	63.10	0.84	48.99	47.42	0.90

6.2 Augmenting existing dataset with generated data

Here, we merge the generated dataset and the original training set into a single training set, and then test the model performance when it is trained with the merged dataset to see whether the generated dataset can further improve model performance with the original training set available. We present the results in Table 8. From the table, we can see that the generated dataset is an effective complement to the original training set, since most of the generated datasets introduce performance gain when combined with the original training set, especially our AttrPrompt which leads to improvement for all the cases. This notable improvement with simple dataset merge may motivate future studies of more advanced ways of using the generated data as augmentations to boost existing dataset.

Table 8: Performance of the model trained with the original training set/augmented with the generated dataset. We present the performance gain/drop compared to using the original training set in green/red.

Method	NYT		Amazon		Reddit		StackExchange	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
SimPrompt	85.56 +1.76	86.34 +5.32	81.85 -0.38	80.23 -0.89	85.11 +0.89	84.88 +1.50	74.53 +6.97	74.23 +10.95
Metaprompt	87.14 +3.34	87.33 +6.31	82.12 -0.11	80.14 -0.98	84.71 +0.49	84.62 +1.24	76.02 +8.46	75.70 +12.42
AttrPrompt w/o CAF	85.71 +1.91	87.18 +6.16	82.24 +0.01	80.76 -0.36	85.86 +1.64	85.65 +2.27	75.16 +7.60	74.64 +11.36
AttrPrompt	87.47 +3.67	88.06 +7.04	83.95 +1.72	83.93 +2.81	86.08 +1.86	85.98 +2.60	76.86 +9.30	76.53 +13.25

6.3 The budget and sample efficiency of the generated data

Here, we aim to study two types of efficiency of the generated dataset, *i.e.*, budget efficiency and sample efficiency, on the model performance. First, in Figure 6, we compare the budget efficiency of AttrPrompt against that of SimPrompt. Surprisingly, AttrPrompt only requires 5% of budget to be on par with or outperform SimPrompt with 100% of budget across all the datasets. This observation highlights the significance of diverse prompts in the training data generation process.

Secondly, we examine the sample efficiency of Gold, SimPrompt, and AttrPrompt in Figure 7. While both SimPrompt and AttrPrompt exhibit better sample efficiency than Gold in the low-data regime, with superior performance when the dataset size is relatively small, Gold data shows better sample efficiency in the high-data regime. Overall, AttrPrompt renders better sample efficiency than SimPrompt, which suggests that increasing the diversity of the prompts could be an effective way to improve the unsatisfactory data scaling trend of using LLM as data generator [56].

⁷Examples of the filtered attributes are exhibited in Appendix H.

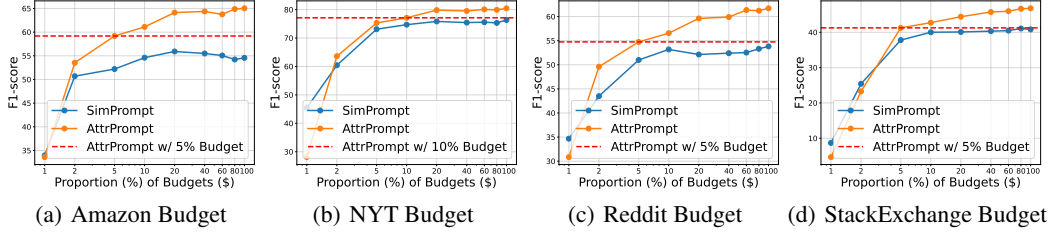


Figure 6: The comparisons on budget efficiency on four datasets.

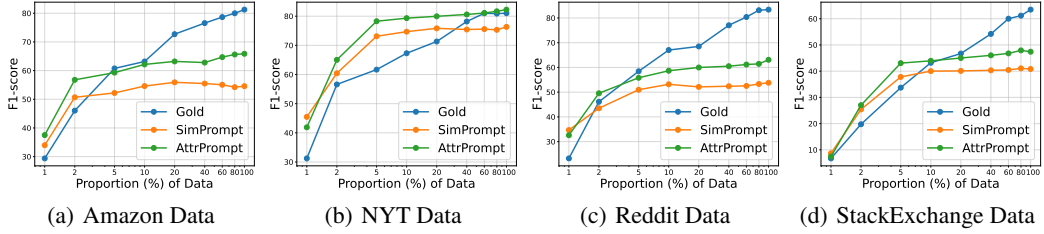


Figure 7: The comparisons on data efficiency on four datasets.

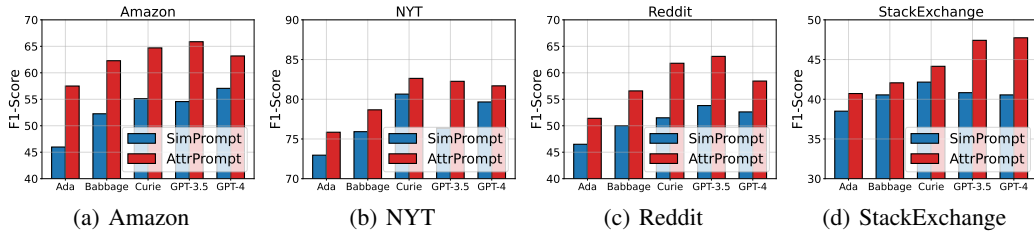


Figure 8: The barplot of performance with LLM Generators of different parameter sizes. Note that due to budget limit, for GPT-4 model, the size of the generated dataset is only 10% of the full set thus the result is not directly comparable with other models.

6.4 The performance with respect to model parameter size

Effect of the Model Size for LLM Generators. To study the effect of different LLMs on AttrPrompt, we use other instruction-finetuned GPT models as the generator, namely `text-ada-001` [38], `text-babbage-001` [38], `text-curie-001` [38], and GPT-4 [36] (due to budget constraints, we only generate a subset with 10% size of the original dataset). Under all settings, our model outperforms the direct baseline SimPrompt by a great margin. Besides, the performance is generally better with larger models, as they often have better instruction-following capabilities. In addition, an interesting finding is that for SimPrompt (but not for AttrPrompt), the average performance of using ChatGPT is worse than `text-curie-001`. This suggests that straightforward class-dependent prompts might not exploit the capabilities of LLMs as effectively as our proposed approaches.

Effect of the Model Size for Classifiers. We experiment with other model choices in addition to the BERT-base-uncased [11] used throughout the paper. They are TinyBERT [19], DistillBERT [45], DeBERTa-V3-base, and DeBERTa-V3-Large [18] with parameter size from 14M to 435M⁸. We visualize the results in Figure 9. Overall, AttrPrompt outperforms SimPrompt by a large margin yet underperforms the Gold across different model choices. With a light-weighted backbone of 66M parameters, AttrPrompt can often outperform the SimPrompt trained with the model containing 435M parameters. This indicates that diversely attributed prompts could help close the performance gap between the Gold and simple class-conditional prompts, and such an improvement is robust to model parameter size.

⁸Additionally, we also evaluate under the setting of linear probing [22] which trains a linear classifier over the frozen pretrained features. The results are deferred to Appendix C.1.

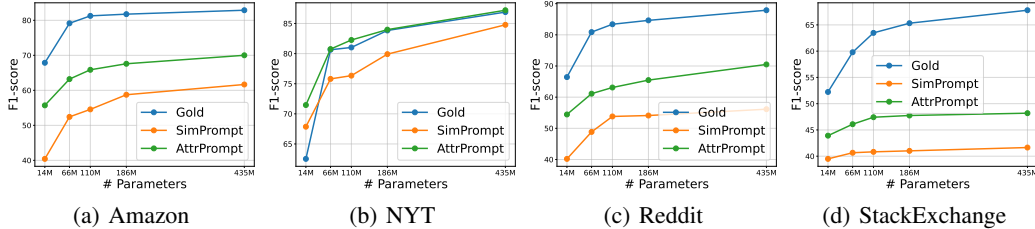


Figure 9: The performance curves with classifier models of different parameter sizes.

6.5 Plugging AttrPrompt in existing approaches

In this section, we demonstrate that AttrPrompt can be painlessly integrated with prior zero-shot training data generation techniques. Table 9 shows the results for several recently proposed methods, which design additional techniques based on the noisy-robust loss to further reduce the effect of noisy labeled data [13, 30, 56], and leverage in-context examples for data generation [56]. Despite these approaches achieving notable performance gains on simple binary classification tasks, their gains become more marginal for fine-grained classification: the performance gain is less than 2% for all methods on two datasets. Instead, using AttrPrompt lead to consistent performance boosts (more than 5% in all cases) for those approaches, indicating that compared with label noise, *data diversity* is a more crucial bottleneck for existing dataset generation methods.

More interestingly, AttrPrompt even benefits dataset generation approaches that do not use LLMs. To demonstrate this, we use the LLM-generated contents (subtopics for NYT and product name for Amazon) to enrich the label names used in ReGen [59], a retrieval-based approach for training data generation. With the expanded label names, AttrPrompt largely improves (14%–26% absolute gain) the performance of ReGen on fine-grained classification tasks. These results justify the advantage of AttrPrompt for serving as a generic plug-in module for existing training data generation approaches.

Table 9: Performance comparison when AttrPrompt serves as a plug-in for existing approaches.

		SuperGen	SuperGen w/ AttrPrompt	SunGen	SunGen w/ AttrPrompt	ProGen	ProGen w/ AttrPrompt	ReGen	ReGen w/ AttrPrompt
NYT	Acc.	76.11	82.05	75.82	81.65	77.05	80.93	70.01	82.18
	F1	76.80	82.62	76.52	82.70	76.70	81.32	68.14	82.50
Amazon	Acc.	58.17	66.76	54.30	63.89	58.40	66.43	34.70	58.40
	F1	56.06	66.33	53.50	63.76	56.95	66.02	30.93	56.00

6.6 Extension to multi-label classification

In this section, we take the first attempt to extend the paradigm of using the LLM as a training data generator to the more challenging multi-label classification setting. In particular, we adopt the arXiv dataset [9] consisting of 98 fine-grained classes, on which we apply both SimPrompt and AttrPrompt. Following [16, 48], we consider different evaluation metrics including Micro/Macro-F1, Precision@ k , Normalized Discounted Cumulative Gain (nDCG@ k), and Mean Reciprocal Rank (MRR). The experimental details are in Appendix A. We present the results in Table 10. Similar to our findings for single-label classification, AttrPrompt largely outperforms SimPrompt across all the metrics, which not only strengthens the superiority of AttrPrompt but also opens the door to using LLM as a training data generator for future research in multi-label classification.

Table 10: Multi-label classification performance (in %) on the arXiv dataset.

Method	Macro F1	Micro F1	Precision@1	Precision@5	NDCG@5	MRR	Price/1k
Gold	27.34	58.22	73.71	27.00	79.70	82.16	—
SimPrompt	21.03	26.75	37.00	15.08	42.49	49.60	1.41
AttrPrompt	27.10	37.88	49.27	18.79	54.74	61.23	1.53

7 Limitation, Social Impact, and Future Work

7.1 Limitation

One limitation of this research is that the investigation focuses primarily on text classification. While this provides valuable insights into the performance of attributed prompts compared to simple class-conditional prompts in these specific contexts, the generalizability of the findings to other data types and tasks remains to be explored. In addition, the AttrPrompt assumes the label name of the target task following prior works [30, 55] and requires some extent of human effort in creating data attributes.

In addition, one issue with LLM-based training data generation is the phenomenon of hallucination, wherein the model generates information that may sound plausible but is, in fact, not accurate or grounded in reality [2]. This can lead to the propagation of misinformation, particularly in sensitive areas such as news generation where accuracy is of utmost importance. To alleviate this issue, it is possible to leverage additional fact-checking mechanisms to cross-verify the generated text with a reliable knowledge base or dataset [39]. Furthermore, incorporating an additional layer of human review can also help in mitigating hallucinations and ensuring the faithfulness of generated outputs.

7.2 Social impact

The findings of this research have significant implications for the development of trustworthy machine learning models. By addressing the limitations of simple class-conditional prompts and exploring the potential of attributed prompts, the research aims to mitigate biases and enhance the diversity of generated training data. This has a positive social impact as it contributes to the development of more robust and fair machine learning models, reducing the potential for biased decision-making systems and promoting inclusivity in various domains. We do not involve human subjects research and do not contain any personally identifiable information. Possible misuse may lead to negative outcomes, such as intentionally introducing biases to the generated dataset via attributes.

7.3 Future work

Broader Coverage of Task and Modalities. In future work, it would be valuable to extend the investigation to other types of classification tasks beyond text classification. Exploring the effectiveness of attributed prompts in image classification, audio classification, or other modalities could provide insights into the broader applicability and effectiveness of AttrPrompt.

Automatic Prompt Optimization. In this work, we aim to demonstrate the usefulness of attributes and why they are useful through the lens of diversity and bias, and we believe this opens the door for future work to further explore the methodology for automatic attribute discovery and selection for optimal performance, such as exploring automated or semi-automated methods for identifying high-quality attribute dimensions and values. One promising direction is to incorporate a few labeled examples as demonstrations for generating better training data. Under this scenario, leveraging Language Models can aid in automatically detecting diverse attributes, modeling the high-order interactions among attributes as well as estimating attribute value distributions for each class, which further reduces the human efforts for selecting key attributes.

8 Conclusion

We delve into the realm of training data generation using complex, attributed prompts, which possess the potential to produce a wide range of diverse and attributed generated data. Specifically, we focus on datasets characterized by diverse domains and high cardinality and class-imbalance, and our results demonstrate the superior performance of attributed prompts compared to simple class-conditional prompts. Furthermore, we present a comprehensive empirical study on training data generation that covers essential aspects such as bias, diversity, and efficiency.

References

- [1] Neel Alex, Eli Lifland, Lewis Tunstall, Abhishek Thakur, Pegah Maham, C. Jess Riedel, Emmie Hine, Carolyn Ashurst, Paul Sedille, Alexis Carlier, Michael Noetel, and Andreas Stuhlmüller. RAFT: A real-world few-shot text classification benchmark. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021.
- [2] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, et al. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*, 2023.
- [3] John Blitzer, Mark Dredze, and Fernando Pereira. Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *Proceedings of the 45th annual meeting of the association of computational linguistics*, 2007.
- [4] Vadim Borisov, Kathrin Sessler, Tobias Leemann, Martin Pawelczyk, and Gjergji Kasneci. Language models are realistic tabular data generators. In *The Eleventh International Conference on Learning Representations*, 2023.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33, 2020.
- [6] Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. Assessing cross-cultural alignment between ChatGPT and human societies: An empirical study. In *Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)*, 2023.
- [7] Derek Chen, Celine Lee, Yun-Yun Lu, Domenic Rosati, and Zhou Yu. Mixture of soft prompts for controllable data generation. *ArXiv*, abs/2303.01580, 2023.
- [8] Yew Ken Chia, Lidong Bing, Soujanya Poria, and Luo Si. Relationprompt: Leveraging prompts to generate synthetic data for zero-shot relation triplet extraction. In *Findings of the Association for Computational Linguistics: ACL 2022*, 2022.
- [9] Colin B Clement, Matthew Bierbaum, Kevin P O’Keeffe, and Alexander A Alemi. On the use of arxiv as a dataset. *arXiv preprint arXiv:1905.00075*, 2019.
- [10] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2019.
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019.
- [12] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023.
- [13] Jiahui Gao, Renjie Pi, Lin Yong, Hang Xu, Jiacheng Ye, Zhiyong Wu, Weizhong Zhang, Xiaodan Liang, Zhenguo Li, and Lingpeng Kong. Self-guided noise-free data generation for efficient zero-shot learning. In *International Conference on Learning Representations*, 2023.
- [14] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021.
- [15] Gregor Geigle, Nils Reimers, Andreas Rücklé, and Iryna Gurevych. Tweac: transformer with extendable qa agent classifiers. *arXiv preprint arXiv:2104.07081*, 2021.
- [16] Chuan Guo, Ali Mousavi, Xiang Wu, Daniel N Holtmann-Rice, Satyen Kale, Sashank Reddi, and Sanjiv Kumar. Breaking the glass ceiling for embedding-based classifiers for large output spaces. *Advances in Neural Information Processing Systems*, 32, 2019.
- [17] Tatsunori B. Hashimoto, Hugh Zhang, and Percy Liang. Unifying human and statistical evaluation for natural language generation. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019.

- [18] Pengcheng He, Jianfeng Gao, and Weizhu Chen. DeBERTav3: Improving deBERTa using ELECTRA-style pre-training with gradient-disentangled embedding sharing. In *The Eleventh International Conference on Learning Representations*, 2023.
- [19] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. In *Findings of EMNLP*, 2020.
- [20] Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*, 2023.
- [21] Hannah Rose Kirk, Yennie Jun, Haider Iqbal, Elias Benussi, Filippo Volpin, Frédéric A. Dreyer, Aleksandar Shtedritski, and Yuki M. Asano. Bias out-of-the-box: An empirical analysis of intersectional occupational biases in popular generative language models. In *Neural Information Processing Systems*, 2021.
- [22] Ananya Kumar, Aditi Raghunathan, Robbie Matthew Jones, Tengyu Ma, and Percy Liang. Fine-tuning can distort pretrained features and underperform out-of-distribution. In *International Conference on Learning Representations*, 2022.
- [23] Yizhi Li, Ge Zhang, Bohao Yang, Chenghua Lin, Anton Ragni, Shi Wang, and Jie Fu. HERB: Measuring hierarchical regional bias in pre-trained language models. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2022*, 2022.
- [24] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2022.
- [25] Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng, Can Zheng, Junxiang Wang, Tanmoy Chowdhury, Yun Li, Hejie Cui, Tianjiao Zhao, et al. Beyond one-model-fits-all: A survey of domain specialization for large language models. *arXiv preprint arXiv:2305.18703*, 2023.
- [26] Alisa Liu, Swabha Swayamdipta, Noah A. Smith, and Yejin Choi. WANLI: Worker and AI collaboration for natural language inference dataset creation. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, 2022.
- [27] Lajanugen Logeswaran, Honglak Lee, and Samy Bengio. Content preserving text generation with attribute controls. *Advances in Neural Information Processing Systems*, 31, 2018.
- [28] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019.
- [29] Andrew 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*, 2011.
- [30] Yu Meng, Jiaxin Huang, Yu Zhang, and Jiawei Han. Generating training data with language models: Towards zero-shot language understanding. In *Advances in Neural Information Processing Systems*, 2022.
- [31] Yu Meng, Jiaming Shen, Chao Zhang, and Jiawei Han. Weakly-supervised hierarchical text classification. In *Proceedings of the AAAI conference on artificial intelligence*, 2019.
- [32] Swaroop Mishra, Anjana Arunkumar, Bhavdeep Sachdeva, Chris Bryan, and Chitta Baral. Dqi: Measuring data quality in nlp. *arXiv preprint arXiv:2005.00816*, 2020.
- [33] Swaroop Mishra, Daniel Khashabi, Chitta Baral, Yejin Choi, and Hannaneh Hajishirzi. Reframing instructional prompts to GPTk’s language. In *Findings of the Association for Computational Linguistics: ACL 2022*, 2022.
- [34] Swaroop Mishra and Elnaz Nouri. HELP ME THINK: A simple prompting strategy for non-experts to create customized content with models. In *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.
- [35] Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? *Advances in neural information processing systems*, 32, 2019.
- [36] OpenAI. Gpt-4 technical report. *arXiv*, 2023.
- [37] OpenAI. Introducing chatgpt, 2023.

- [38] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35, 2022.
- [39] Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, et al. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*, 2023.
- [40] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *ArXiv*, abs/2304.03277, 2023.
- [41] Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models. In *Advances in Neural Information Processing Systems*, 2021.
- [42] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019.
- [43] Laria Reynolds and Kyle McDonell. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, 2021.
- [44] Giuseppe Russo, Nora Hollenstein, Claudiu Musat, and Ce Zhang. Control, generate, augment: A scalable framework for multi-attribute text generation. *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020.
- [45] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [46] Timo Schick and Hinrich Schütze. Generating datasets with pretrained language models. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 2021.
- [47] Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie Huang, Nan Duan, and Weizhu Chen. Synthetic prompting: Generating chain-of-thought demonstrations for large language models. *arXiv preprint arXiv:2302.00618*, 2023.
- [48] Jiaming Shen, Wenda Qiu, Yu Meng, Jingbo Shang, Xiang Ren, and Jiawei Han. Taxoclass: Hierarchical multi-label text classification using only class names. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021.
- [49] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, 2013.
- [50] Zhiqing Sun, Yikang Shen, Qinlong Zhou, Hongxin Zhang, Zhenfang Chen, David Cox, Yiming Yang, and Chuang Gan. Principle-driven self-alignment of language models from scratch with minimal human supervision. *arXiv preprint arXiv:2305.03047*, 2023.
- [51] Guy Tevet and Jonathan Berant. Evaluating the evaluation of diversity in natural language generation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, 2021.
- [52] Yisen Wang, Xingjun Ma, Zaiyi Chen, Yuan Luo, Jinfeng Yi, and James Bailey. Symmetric cross entropy for robust learning with noisy labels. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019.
- [53] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- [54] Sarah Wiegreffe, Jack Hessel, Swabha Swayamdipta, Mark Riedl, and Yejin Choi. Reframing human-ai collaboration for generating free-text explanations. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2022.

- [55] Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong. ZeroGen: Efficient zero-shot learning via dataset generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022.
- [56] Jiacheng Ye, Jiahui Gao, Zhiyong Wu, Jiangtao Feng, Tao Yu, and Lingpeng Kong. ProGen: Progressive zero-shot dataset generation via in-context feedback. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, 2022.
- [57] Dian Yu, Zhou Yu, and Kenji Sagae. Attribute alignment: Controlling text generation from pre-trained language models. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 2021.
- [58] Yue Yu, Chenyan Xiong, Si Sun, Chao Zhang, and Arnold Overwijk. COCO-DR: Combating distribution shifts in zero-shot dense retrieval with contrastive and distributionally robust learning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, 2022.
- [59] Yue Yu, Yuchen Zhuang, Rongzhi Zhang, Yu Meng, Jiaming Shen, and Chao Zhang. Regen: Zero-shot text classification via training data generation with progressive dense retrieval. In *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.
- [60] Jieyu Zhang, Bohan Wang, Zhengyu Hu, Pang Wei Koh, and Alexander Ratner. On the trade-off of intra-/inter-class diversity for supervised pre-training. *arXiv preprint arXiv:2305.12224*, 2023.
- [61] Rongzhi Zhang, Yue Yu, Pranav Shetty, Le Song, and Chao Zhang. Prboost: Prompt-based rule discovery and boosting for interactive weakly-supervised learning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2022.
- [62] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
- [63] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. *Advances in neural information processing systems*, 28, 2015.
- [64] Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. In *The Eleventh International Conference on Learning Representations*, 2023.
- [65] Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. Exploring ai ethics of chatgpt: A diagnostic analysis. *ArXiv*, abs/2301.12867, 2023.

A Dataset Preprocessing

A.1 Multi-class Classification

For Reddit and StackOverflow, we select the classes with more than 65 examples from the original corpus as the target set of topics. For each dataset, we use 50 examples per class for the test set and no more than 10 examples for the validation set (10 for NYT/Amazon and 5 for Reddit/StackOverflow). The remaining data is used to compose the gold training set. It is worth noting that, some of the class names on Reddit may contain toxic information. To eliminate their effects, we filter label names with Detoxify⁹, a tool that leverages the multilingual XLM-RoBERTa [10] for toxic comment identification. We follow [12] to use a threshold of 0.1 to filter out potentially toxic topic classes.

A.2 Multi-label Classification

For the arXiv multi-label classification, we use the arXiv labels from the website¹⁰ as the label names. We select topics under the category of *Physics*, *Mathematics*, *Computer Science*, *Quantitative Biology*, and *Quantitative Finance*. Besides, we remove the topic of *General Physics*, *General Mathematics*, *General Literature*, *Other Computer Science*, *Other Quantitative Biology* as they only contain generic topics without concrete semantics. In total, there are 98 classes in general.

To generate examples using the SimPrompt approach, we follow a specific procedure. For each document, we begin by sampling from the distribution of label numbers, utilizing prior knowledge obtained from the original training set. This enables us to determine the precise number of labels, denoted as n , that should be associated with the document. Once we have determined n , we proceed to randomly select n classes from the original set of 98 labels. These chosen classes are then utilized as the input for prompting the LLM, forming a foundation for generating the examples.

To generate examples using the AttrPrompt approach, we employ a distinct methodology. For each document, we consider a set of attributes, namely `subtopics`, `technique`, `writing style`, and `length`. It is important to note that, unlike multi-class classification, where an additional step is typically taken to filter out irrelevant subtopics, our multi-label classification process follows a different approach. In our method, we introduce an additional step called *merging* to handle the different subtopics.

To determine the merging of subtopics, we utilize the fuzzy matching score¹¹, which calculates the similarity between the subtopics. Setting the threshold at $\gamma = 90$, we proceed to merge subtopics whenever their similarity score exceeds this threshold. Consequently, specific subtopics may be associated with multiple classes, providing a more comprehensive representation of the document content. When generating text utilizing these subtopics, the corresponding labels are determined by the classes associated with the specific subtopics. This ensures that the generated examples accurately reflect the relevant classes based on the subtopic associations.

B Implementation Details

B.1 Hardware information

All experiments are conducted on *CPU*: Intel(R) Core(TM) i7-5930K CPU @ 3.50GHz and *GPU*: NVIDIA GeForce RTX A5000 GPUs using python 3.8, Huggingface 4.6.0 and Pytorch 1.10.

B.2 Parameter configuration

We follow [40] to keep the parameter `top_p = 1.0` and temperature $t = 1.0$ for calling ChatGPT APIs [37] for the training data generation part. For finetuning the classifier, we optimize the model using AdamW [28] with a linear warmup of the first 5% steps and linear learning rate decay. The maximum number of tokens per sequence is 128. Table 11 lists the hyperparameters used for AttrPrompt and SimPrompt. For the generated synthetic dataset, we stick to the strict zero-shot

⁹<https://github.com/unitaryai/detoxify>

¹⁰<https://arxiv.org/>

¹¹<https://github.com/seatgeek/thefuzz>

learning setting [30], train all the models for 6 epochs and use the model from the last epoch *without using the validation set* for evaluation. For the original clean dataset, we train all models for 6 epochs and use the model with the best performance on the validation set for evaluation.

Backbone	Learning Rate lr	Batch Size	Training Epochs E	Weight Decay	Warmup Ratio
BERT-base-uncased [11]	5e-5	32	6	1e-4	6%
TinyBERT [19]	1e-4	32	6	1e-4	6%
DistilBERT-base-uncased [45]	5e-5	32	6	1e-4	6%
DeBERTa-v3-base [18]	5e-5	32	6	1e-4	6%
DeBERTa-v3-large [18]	2e-5	32	6	1e-4	6%

Table 11: Hyperparameters for fine-tuning classifiers on different tasks.

Model	Link
BERT-base-uncased [11]	https://huggingface.co/bert-base-uncased
TinyBERT [19]	https://huggingface.co/huawei-noah/TinyBERT_General_4L_312D
DistilBERT-base-uncased [45]	https://huggingface.co/distilbert-base-uncased
DeBERTa-v3-base [18]	https://huggingface.co/microsoft/deberta-v3-base
DeBERTa-v3-large [18]	https://huggingface.co/microsoft/deberta-v3-large
SentenceBERT [42]	https://huggingface.co/sentence-transformers/all-mpnet-base-v2
SimCSE [14]	https://huggingface.co/princeton-nlp/sup-simcse-bert-base-uncased
COCO-DR [58]	https://huggingface.co/OpenMatch/cocodr-base-msmarco

Table 12: The reference of model checkpoints used in this study.

B.3 Model checkpoint

All the checkpoint used in fine-tuning and linear probing experiments are listed in table 12.

C Additional Experimental Results

C.1 Linear probing with generated datasets

In this section, we evaluate the generated dataset under linear probing setting. We employ four text embedding methods and present the results in Table 13. We can see that AttrPrompt consistently outperforms SimPrompt and surpasses the Gold in the NYT dataset. This further indicates the datasets generated with AttrPrompt have higher quality than SimPrompt and this superiority is robust to the model choices.

C.2 Different temperature parameters for ChatGPT

Temperature (t) is one crucial hyperparameter of LLMs that controls the diversity of the generated text [17], while the studied attributed prompts are also for diversifying the generated data. We are then curious about the effectiveness of the temperature and how it compares to the AttrPrompt. We study different values of the temperature using the NYT dataset and present the results in Table 14. From the results, we can see that compared with the temperature, AttrPrompt brings more significant performance gain, demonstrating its superiority over temperature tuning.

C.3 The performance with respect to long-tail classes

As we have seen in Table 2, the original training sets of the involved datasets have severe long-tail classes issue since the imbalance ratio is high, yet the generated dataset are class-balanced, we are then curious how the class balance in the generated dataset benefits the model performance on long-tail classes. We take the NYT dataset as an example and plot the per-class F1 score of Gold, SimPrompt, and AttrPrompt in Figure 10, where the x-axis is classes sorted by their number of data in the Gold dataset in descending order. From the figure, we can see that out of 26 classes, AttrPrompt renders the best per-class F1 score on 10 classes, which is 13 for Gold and 3 for SimPrompt. Notably, for classes with few examples in the Gold set (the rightmost 4 classes in the figure), AttrPrompt is better than the Gold and SimPrompt, especially for the class “*abortion*” with the fewest examples. This suggests a data-centric way to handle the long-tail class issue in topic classification: one may use LLMs to

Table 13: Performance of linear probing with different pre-trained encoders.

Sequence Encoder	Method	NYT		Amazon		Reddit		StackExchange	
		Acc	F1	Acc	F1	Acc	F1	Acc	F1
BERT-CLS [11]	Gold	65.39	70.26	66.73	68.76	55.73	58.00	17.17	17.68
	SimPrompt	60.69	62.26	37.95	38.07	27.41	27.60	12.27	13.71
	MetaPrompt	58.43	57.05	27.52	24.43	23.93	23.99	11.40	12.60
	AttrPrompt	62.38	64.17	48.18	49.46	45.03	45.95	17.33	17.64
SentenceBERT [42]	Gold	81.16	84.43	80.95	83.30	80.76	81.37	55.13	54.25
	SimPrompt	84.32	85.91	67.55	68.29	69.94	70.62	45.37	47.79
	MetaPrompt	87.65	87.11	67.16	65.78	79.26	69.15	45.71	47.54
	AttrPrompt	88.77	88.36	68.10	69.11	71.62	71.73	46.80	48.35
SimCSE [14]	Gold	78.02	80.17	73.20	75.13	73.75	74.26	45.26	41.86
	SimPrompt	78.93	78.86	52.40	52.83	54.78	53.86	30.98	32.92
	MetaPrompt	78.17	78.27	53.53	52.55	51.30	48.93	30.54	31.50
	AttrPrompt	81.41	81.04	60.66	61.32	59.22	59.15	30.31	31.43
COCO-DR [58]	Gold	80.66	82.60	76.13	78.34	79.63	80.04	61.02	63.77
	SimPrompt	84.37	84.17	55.66	55.28	63.66	64.71	43.22	41.50
	MetaPrompt	84.51	84.26	57.96	56.58	65.36	65.91	44.30	42.97
	AttrPrompt	86.74	86.43	60.47	62.03	66.73	66.71	45.44	43.76
Average	Gold	76.31	79.37	74.25	76.38	72.47	73.42	44.65	44.39
	SimPrompt	77.08	77.80	53.39	53.62	53.95	54.20	32.96	33.98
	MetaPrompt	77.19	76.67	51.54	49.84	54.96	52.00	32.99	33.65
	AttrPrompt	79.83	80.00	59.35	60.48	60.65	60.89	34.97	35.30

Table 14: Study of the temperature.

Method	$t = 1.0$		$t = 1.5$		$t = 2.0$	
	Acc.	F1	Acc.	F1	Acc.	F1
SimPrompt	76.00	76.34	76.78	77.31	76.55	77.42
AttrPrompt	81.30	82.26	81.47	82.23	79.47	79.86

generate class-balanced training set or augment the existing training set with the LLM-generated data such that the augmented dataset is class-balanced, the in-depth study of which is left as future work.

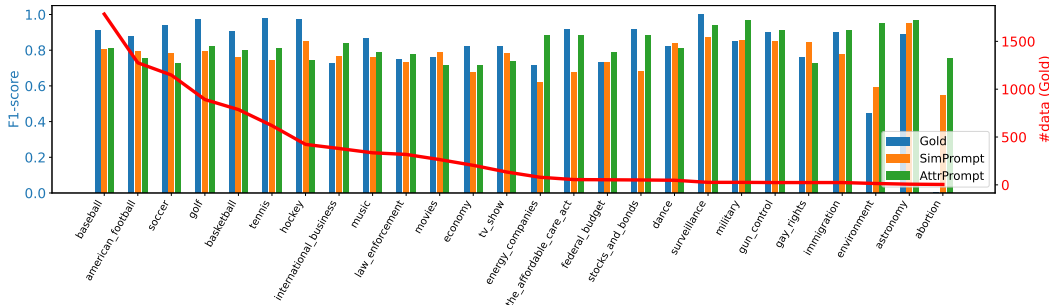


Figure 10: Per-class F1-score of the NYT dataset.

C.4 Experiment results on other datasets

To provide a comprehensive evaluation of AttrPrompt, we conducted additional empirical studies using datasets commonly employed in prior works, which typically exhibit *low cardinality*. Consistent with [56, 55, 59], we focused on four specific tasks: SST-2 [49], IMDB [29], Yelp [63], and AG News [63]. The IMDB, SST-2, and Yelp datasets involve binary sentiment classification, while the AG News dataset entails news classification with only four classes. This selection allows for diverse evaluation across different classification scenarios.

To ensure a fair comparison, we utilized DistilBERT [45] as the classifier, in line with [56, 55, 13]. For the purpose of fine-tuning the classifier, we generated a total of 6000 examples, evenly distributed among the classes. For fine-tuning, we use a learning rate of 2×10^{-5} with 5 epochs and use the model after the final epoch for evaluation.

Table 15: Results with recent baselines using DistilBERT [45]. *: Results are copied from the previous papers [56, 55, 13, 59].

Method/Dataset	IMDB	SST-2	Yelp	AG News	Avg.
<i>Simple fine-tuning with standard cross-entropy loss</i>					
SimPrompt	85.64	79.38	88.39	74.30	81.93
AttrPrompt	85.92	83.64	90.60	82.74	85.73
<i>Other results with dedicated training techniques (for reference only)</i>					
Zero-shot Prompting*	77.31	82.63	78.10	73.20	77.81
ZeroGen*	82.77	80.41	87.84	76.48	81.88
SuperGen*	85.05	84.28	89.88	77.40	84.15
ProGen*	87.20	84.12	89.39	80.81	85.38
SunGen*	89.45	89.38	89.19	80.49	87.12
ReGen*	85.32	87.84	89.00	80.18	85.58

The results of our evaluation are presented in table 15. Overall, we observe that AttrPrompt demonstrates competitive performance across all four tasks. It outperforms our direct baseline (SimPrompt) by an impressive margin of 3.8%. Moreover, when compared to other baselines incorporating dedicated noisy-aware learning techniques [52, 35], AttrPrompt still achieves promising performance, despite using the standard cross-entropy loss. This suggests that in addition to leveraging noisy-robust learning objectives, enhancing the diversity of input prompts serves as an effective alternative to improving downstream performance. It is interesting yet important to combine these different learning paradigms together to further boost the performance.

D Additional Bias Analysis

D.1 Regional Bias with Manual Annotations

In the main body of the paper, we manually annotate location for 100 samples from each of the Gold, SimPrompt, and AttrPrompt dataset. Here, we compare manual annotations with the prediction of the attribute classifier in Table 16 (the first three rows). Note that we include “*unkown*” as an option in manual annotation to absorb text without clear location specifications. As per the results, there is a significant alignment between the manual annotations and the attribute classifier predictions. The accuracy is greater than 0.95 for the Gold and AttrPrompt datasets, while it is 0.82 for SimPrompt. The lower accuracy in SimPrompt is attributed to its tendency to generate more data with unknown locations.

Furthermore, we use the manually annotated data to evaluate ChatGPT’s adherence to the location attribute stipulated in the prompts. By comparing the location attribute in the prompt associated with each AttrPrompt sample against the human annotations, we calculate an accuracy of 0.96 (as shown in the last row of Table 16). This high accuracy means that ChatGPT can effectively generate news adhering to the specified location attribute in the prompts.

Table 16: Accuracies based on manual annotation of the location attribute of the NYT dataset.

Comparison	Data Source	Accuracy
Attribute Classifier v.s. Human	Gold	0.96
	SimPrompt	0.82
	AttrPrompt	0.95
Attribute in Prompt v.s. Human	AttrPrompt	0.96

D.2 Biases other than regional bias

In addition to the regional bias we study in the main body of the paper, here we investigate bias in other datasets (Figure 11). From the results, we can see that the dataset generated with Gold and SimPrompt exhibit different level of biases, which shows that in general the biases exists in both

generated and real dataset, while dataset generated with AttrPrompt has relatively balanced attribute distribution.

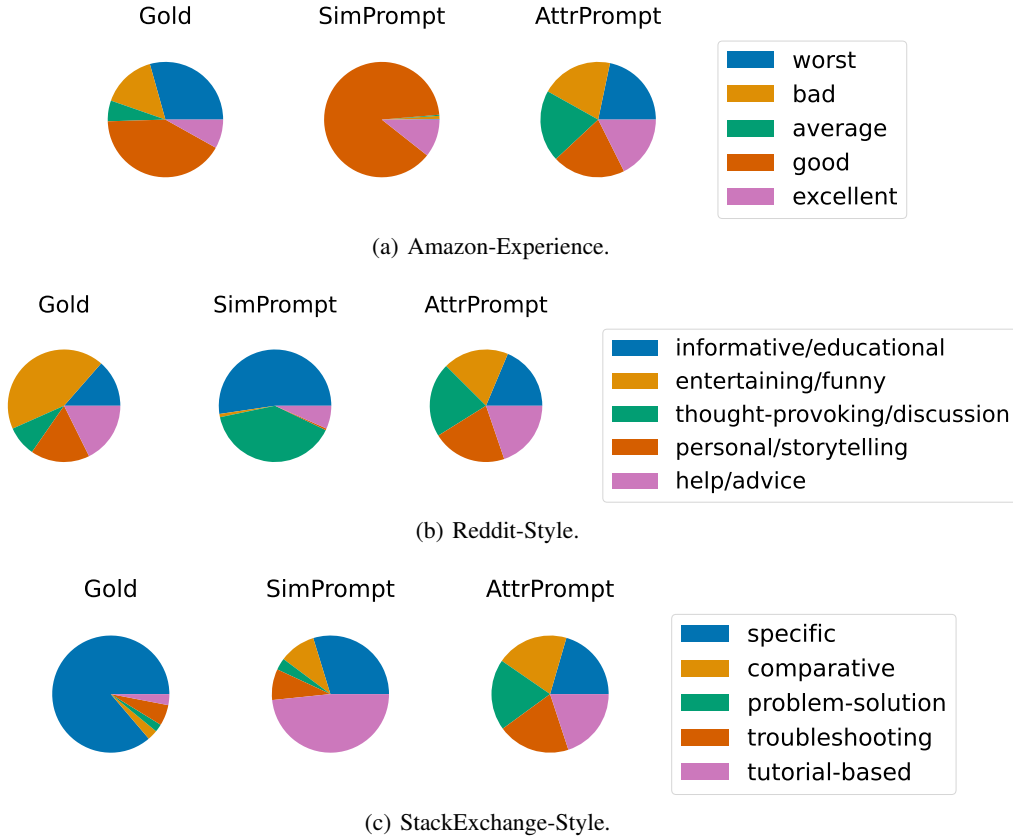


Figure 11: Pie charts of predicted attribute distributions on different datasets.

E The Prompt Format for LLM-Generated Attributes

E.1 Prompts for NYT Dataset

E.1.1 Key Attributes

Prompts to ChatGPT. What do you think are important attributes to generate some diverse news from NYT under a specific topic, such as Football, Federal Budget. Examples: subtopics, writing style...

Response from ChatGPT. Important attributes of news from the New York Times under a specific topic could include:

Subtopics: exploring different angles and subtopics within the main topic can lead to a wider range of news stories.

Writing style: varying the writing style can help to make the news stories more interesting and engaging for readers. This can include using different tones, perspectives, and sentence structures.

Sources: using a variety of sources, such as experts, eyewitnesses, and affected individuals, can provide different perspectives and add diversity to the news stories.

Geographical locations: covering the topic from different geographical locations can help to show how it affects different communities and provide a global perspective.

Timeliness: covering current events and breaking news related to the topic can add diversity and relevance to the news stories.

Multimedia: using a range of multimedia elements, such as images, videos, and infographics, can help to add variety and appeal to the news stories.

Length: varying the length of the news stories, from short and concise to longer and more in-depth, can help to provide a wider range of information and cater to different readers' preferences.

Selected Attributes

1. Subtopics
2. Writing Styles
3. Geographical locations (determined by six continents by default)
4. Length (determined by {30-80} and {100-150} by default)

E.1.2 Attributes Details

Prompts to ChatGPT for Similar Classes. List 5 similar classes for {class-name} news on NYT. The set of classes is listed as: {[a list of class-names]}.

Prompts to ChatGPT for Subtopics. List 10 diverse subtopics for {class_name} news on NYT.

Response from ChatGPT. See section G for details.

E.1.3 Writing Styles

Prompts to ChatGPT. List 5 different writing styles of NYT news.

Response from ChatGPT.

1. Investigative reporting
2. Opinion/editorial
3. Human interest/profiles
4. Long-form narrative
5. Breaking news

E.2 Prompts for Amazon Dataset

E.2.1 Key Attributes

Prompts to ChatGPT. What do you think are important attributes to generate some diverse product reviews. Examples: writing style...

Response from ChatGPT. When writing a review for a product, here are five important attributes to consider:

Product Type: Clearly mention the type of product you are reviewing, such as a smartphone, laptop, or fitness tracker. This helps readers understand the category and purpose of the product.

Brand: Specify the brand of the product as it often influences quality, reputation, and customer support. Discuss the brand's overall credibility and whether it aligns with your expectations.

User Experience: Evaluate the overall user experience of the product. Discuss its speed, accuracy, reliability, and efficiency in performing its intended tasks. Highlight any exceptional or lacking performance aspects.

Quality and Durability: Assess the quality of the product, including the materials used, construction, and overall durability. Discuss whether it feels well-made, solid, and likely to withstand regular use over time.

Features and Functionality: Describe the specific features and functions of the product. Highlight any unique or standout features that enhance its usability or set it apart from similar products in the market.

E.2.2 Selected Attributes

1. Product Name
2. Brands
3. User Experience
4. Writing Styles
5. Length

E.2.3 Attributes Details

Prompts to ChatGPT for Similar Classes. List 5 similar classes for `{product-class}` products on Amazon. The set of classes is listed as: `{[a list of class-names]}`.

Prompts to ChatGPT for Product Name/Brands. List 20 diverse product names/brands for `{product-class}` on Amazon.

Response from ChatGPT. See section [G.2](#) for details.

E.2.4 Usage Experience

Prompts to ChatGPT. List 5 different usage experiences of Amazon reviews.

Response from ChatGPT.

1. Worst: A highly unsatisfactory or problematic experience. Users may encounter multiple issues, such as usability problems, errors, and frustrations.
2. Bad: A below-average experience. Users may encounter some difficulties or frustrations while interacting with the product or service, but the overall experience falls short of meeting their expectations.
3. Average: A neutral or middling experience. Users may not encounter major issues, but the experience lacks any exceptional aspects or standout features.
4. Good: A positive and satisfactory experience. Users have a smooth and enjoyable experience with minimal issues or obstacles.
5. Excellent: The highest level of satisfaction. Users have a seamless, delightful, and memorable experience.

E.2.5 Writing Styles

Prompts to ChatGPT. List 5 different writing styles of Amazon reviews.

Response from ChatGPT.

1. Detailed Review: A detailed review provides a comprehensive analysis of the product, covering various aspects such as features, performance, durability, and user experience.
2. Comparative Review: A comparative review compares the product with similar options available on Amazon.
3. Pros and Cons Review: A pros and cons review presents a balanced assessment of the product by listing its positive and negative aspects.
4. Personal Experience Review: A personal experience review shares the reviewer's firsthand experience with the product.
5. Recommendation Review: A recommendation review offers a clear recommendation or endorsement of the product.

E.2.6 Prompt Template

See table [17](#) for details.

Table 17: Prompt template for the Amazon Product Review dataset.

Method	Prompt
SimPrompt	Suppose you are a review writer. Please write a review for <code>{product-class}</code> in Amazon.
AttrPrompt	Suppose you are a review writer. Please write a review for <code>{product-class}</code> product in Amazon following the requirements below: <ol style="list-style-type: none"> 1. The review should be about the product of <code>{subtopic}</code>; 2. The brand for the product should be <code>{brand}</code>; 3. Should be in length between <code>{length:min-words}</code> and <code>{length:max-words}</code> words; 4. Should describe the usage experience <code>{usage-experience}</code> 5. The writing style of the review should be <code>{style}</code>; 6. the review must be relevant to <code>{product-class}</code> and irrelevant to: <code>{similar-class}</code>.

E.3 Prompts for Reddit Dataset

E.3.1 Key Attributes

Prompts to ChatGPT. List 5 important attributes to generate some diverse reddit posts under a specific topic, like build_a_pc, premed, etc. Examples: writing style...

Response from ChatGPT. Informative Guides and Tutorials: Create posts that provide step-by-step guides, tutorials, or informative content. These can include detailed instructions, tips, or recommendations to help readers in their journey. For example, in "build_a_pc," you could create a guide on assembling a gaming PC or optimizing system performance.

Writing Style: Utilize different writing styles to create diverse posts. For instance, you can write informative guides or tutorials, opinion pieces, comparison reviews, troubleshooting tips, or personal stories related to the topic.

Discussion Starters: Pose thought-provoking questions or statements to initiate discussions and engage the Reddit community. These can be open-ended queries, asking for opinions or experiences, or presenting a controversial viewpoint to encourage debate and dialogue.

Personal Experiences: Share personal stories or experiences related to the topic. This allows you to connect with the community on a more personal level and provide insights based on your firsthand experiences. It can include successes, challenges, lessons learned, or interesting anecdotes.

Resources and Recommendations: Provide valuable resources such as articles, websites, tools, or recommended products/services related to the topic. These resources can help the community access additional information, useful tools, or trusted recommendations, enhancing their understanding and engagement with the subject.

E.3.2 Selected Attributes

1. Personal Experiences
2. Resources
3. Writing Styles
4. Length

E.3.3 Attributes Details

Prompts to ChatGPT for Similar Classes. List 5 similar classes for `{post-class}` products on Reddit. The set of classes is listed as: `{[a list of class-names]}`.

Prompts to ChatGPT for Personal Experiences/Resources. List 20 diverse Personal Experiences/Resources for posts about `{post-class}` in Reddit.

Response from ChatGPT. See section G.3 for details.

E.3.4 Writing Styles

Prompts to ChatGPT. List 5 different writing styles of Reddit Posts.

Table 18: Prompt template for the Reddit dataset.

Method	Prompt
SimPrompt	Suppose you are a writer for posts in Reddit. Please give a post for {post-class} in Reddit.
AttrPrompt	Suppose you are a review writer. Please give a post in Reddit on {post-class} community following the requirements below: 1. The post should focus on {experience}; 2. Should be in length between {length:min-words} and {length:max-words} words; 3. The writing style of the post should be {style}; 4. Should mention the resource of {resource}; 5. the post must be relevant to {post-class} and irrelevant to: {similar-class}.

Response from ChatGPT.

1. Informative/Educational: Create posts that aim to provide valuable information or educate the community on a particular topic. Share informative articles, research findings, guides, tutorials, or explain complex concepts in a simplified manner.
2. Entertaining/Funny: Craft posts that are meant to entertain or bring humor to the community. Share funny stories, jokes, memes, amusing videos, or any form of content that can elicit laughter and enjoyment.
3. Thought-Provoking/Discussion: Initiate posts that spark thoughtful discussions and invite users to share their opinions or experiences. Pose open-ended questions, present philosophical dilemmas, discuss current events, or explore societal issues to foster engaging conversations.
4. Personal/Storytelling: Share personal experiences, anecdotes, or narratives that others can relate to. Tell compelling stories, share achievements, challenges, or life lessons that can inspire or resonate with the community.
5. Help/Advice: Create posts where you seek help or offer advice to the community. Ask for guidance on a specific problem, seek recommendations, provide support, or share expertise in a particular field.

E.3.5 Prompt Template

See table 18 for details.

F Prompts for StackExchange Dataset

F.1 Selected Attributes

1. Depth of Knowledge
2. Writing Styles
3. Problem Scenarios
4. Length

F.2 Key Attributes

Prompts to ChatGPT. What do you think are important attributes to generate diverse posts under a specific topic on the stackexchange forum, like 3D meshes, server, etc. For example, one important attribute could be writing style...

Response from ChatGPT. When generating diverse posts under a specific topic on the Stack Exchange forum, several important attributes can contribute to the variety and richness of the discussions:

Depth of Knowledge: Participants with different levels of expertise can offer insights and perspectives from beginner to advanced, allowing for a broader range of questions and answers.

Diverse Perspectives: Encouraging participants from various backgrounds, industries, and experiences can bring fresh viewpoints and approaches to the topic.

Objectives: Encouraging participants to state their objectives or goals in their posts can help generate diverse content. Some posts may aim to seek advice, share experiences, provide step-by-step guides, analyze specific problems, or discuss broader concepts.

Problem Scenarios: Encourage users to discuss real-world problem scenarios related to the topic. This helps generate diverse posts by exploring different use cases, challenges, and approaches to problem-solving.

Problem Solving: Promoting problem-solving skills and encouraging users to share their challenges, roadblocks, and solutions can foster a collaborative and dynamic environment.

Clear Communication: Emphasizing clear and concise communication helps to ensure that questions are well-articulated and that answers are easily understood by a wide range of users.

F.3 Prompts for arXiv Dataset

F.3.1 Key Attributes

Prompts to ChatGPT. What do you think are important attributes to generate diverse papers under a specific topic on arXiv. For example, one important attribute could be writing style...

Response from ChatGPT. Generating diverse papers under a specific topic on platforms like arXiv involves considering various important attributes. Here are a few attributes that can contribute to generating diverse papers:

Writing Style: Encouraging papers with different writing styles, such as technical, expository, theoretical, or empirical, can bring diversity to the presentation and appeal to a wider range of readers.

Subtopics: Promoting papers that explore different subtopics within the broader topic can provide comprehensive coverage and delve into specific areas of interest.

Techniques: Encouraging papers that employ different research methodologies, such as experimental, computational, or analytical, can bring diverse approaches to studying the topic.

Data Sources: Promoting papers that utilize diverse data sources, such as surveys, simulations, real-world datasets, or case studies, can offer different perspectives and insights into the topic.

Interdisciplinary Perspectives: Encouraging papers that incorporate interdisciplinary perspectives, drawing insights from multiple fields or combining methodologies from different disciplines, can contribute to a richer understanding of the topic.

F.3.2 Selected Attributes

1. Subtopics
2. Techniques
3. Writing Styles
4. Length

Prompts to ChatGPT for Subtopics/Techniques. List 20 diverse subtopics/techniques for paper about `{paper-class}` in arXiv.

Response from ChatGPT. See section [G.5](#) for details.

F.3.3 Writing Styles

Prompts to ChatGPT. List 5 different writing styles of arXiv papers.

Response from ChatGPT.

1. Experimental: This style focuses on presenting empirical research with a strong emphasis on experimental design, data collection, and statistical analysis.
2. Theoretical: This style emphasizes the development of new theoretical frameworks or conceptual models.

Table 19: Prompt template for the arXiv multi-label dataset.

Method	Prompt
SimPrompt	Please give a synthetic paper abstract about <code>{paper-classes}</code> .
AttrPrompt	Please write an abstract of a <code>{paper-classes}</code> paper in arXiv following the requirements below: 1. Should focus on the <code>{subtopic}</code> ; 2. Should be in length between <code>{length:min-words}</code> and <code>{length:max-words}</code> words; 3. The paper should use the technique relevant to <code>{technique}</code> ; 4. The style of the paper should be in <code>{style}</code> ;

3. Review: This style involves synthesizing and analyzing existing literature on a specific topic.
4. Case Study: This style presents an in-depth examination of a specific case or situation, aiming to provide insights into complex phenomena or phenomena that are context-specific.

F.3.4 Prompt Template

See table 19 for details.

F.4 Prompts for AG News Dataset

The prompts are the same as the NYT dataset. The only difference is the subtopics for different labels. The subtopic information is listed in G.6 for details.

F.5 Prompts for Yelp/SST-2/IMDB Dataset

F.5.1 Key Attributes

Prompts to ChatGPT. What do you think are important attributes to generate some diverse reviews with sentiments? E.g. writing styles.

Response from ChatGPT. Writing Styles: Utilizing different writing styles such as descriptive, narrative, persuasive, or humorous can create a diverse range of tones and emotions in the reviews.

Language and Tone: Varying the language and tone used in the reviews can evoke different sentiments. This can include formal or casual language, positive or negative tone, or even incorporating elements of storytelling.

Personal Experience: Sharing personal experiences and anecdotes in the reviews can add authenticity and emotional depth, allowing readers to connect with the sentiments expressed.

Descriptive Details: Including specific and vivid descriptions of the ambiance, decor, flavors, smells, or textures of the food can create a sensory experience that evokes various sentiments.

F.5.2 Selected Attributes

1. Personal Experience
2. Descriptive Details
3. Writing Styles
4. Length

F.5.3 Personal Experiences/Descriptive Details

Prompts to ChatGPT. Details are listed in Appendix G.6 and G.8.

F.5.4 Writing Styles

Prompts to ChatGPT. List 5 different writing styles of Restaurant Reviews (for Yelp).

Response from ChatGPT.

1. Descriptive Style: Focuses on providing detailed descriptions of the restaurant’s ambiance, decor, and overall atmosphere. It aims to paint a vivid picture of the dining experience.

Table 20: Prompt template for the Yelp dataset.

Method	Prompt
SimPrompt	Suppose you are a review writer. Please give an example of a <code>{sentiment}</code> review in Yelp.
AttrPrompt	Please write a <code>{sentiment}</code> review for a restaurant, following the requirements below: <ol style="list-style-type: none"> 1. The overall review should be <code>{sentiment}</code>; 2. It should be a <code>{descriptive-details}</code> restaurant ; 3. Should focus on <code>{personal-experience}</code>; 4. Should be in length between <code>{length:min-words}</code> and <code>{length:max-words}</code> words; 5. The writing style of the review should be <code>{style}</code>;

Table 21: Prompt template for the SST-2/IMDB dataset.

Method	Prompt
SimPrompt	Suppose you are a review writer. Please give an example of a <code>{sentiment}</code> review for a movie.
AttrPrompt	Please write a <code>{sentiment}</code> review for a movie, following the requirements below: <ol style="list-style-type: none"> 1. The overall review should be <code>{sentiment}</code>; 2. It should be a <code>{descriptive-details}</code> movie ; 3. Should focus on <code>{personal-experience}</code>; 4. Should be in length between <code>{length:min-words}</code> and <code>{length:max-words}</code> words; 5. The writing style of the review should be <code>{style}</code>;

2. Personal Narrative Style: Incorporates personal anecdotes and experiences from the visit to the restaurant. It adds a personal touch to the review and allows the reader to connect with the reviewer’s perspective.
3. Analysis Style: Takes a more analytical approach by evaluating various aspects of the restaurant, such as food quality, service, cleanliness, and value for money.
4. Comparative Style: Compares the restaurant to similar establishments in terms of cuisine, price range, service quality, or other relevant factors.
5. Gastronomic Style: Focuses primarily on the food and beverage offerings of the restaurant.

Prompts to ChatGPT. List 5 different writing styles of Movie Reviews (for IMDB/SST-2).

Response from ChatGPT.

1. Descriptive: Focuses on providing a detailed and vivid description of the movie and its elements.
2. Analytical: Analyzes and evaluates various aspects of the movie, such as the plot, acting, cinematography, and direction.
3. Persuasive: Aims to convince the reader to watch or avoid the movie by presenting strong arguments and opinions.
4. Comparative: Compares the movie to other films within the same genre or by the same director.

F.5.5 Prompt Template

See table 20 and 21 for details.

G Details for LLM-Generated Attributes

In this section, we present a brief description of LLM-Generated Attributes. For a comprehensive list of attributes, please refer to the GitHub repository: <https://github.com/yueyu1030/attrprompt>.

G.1 NYT Attributes

G.1.1 Subtopics

We randomly select 3 categories in NYT dataset and display the corresponding subtopic attributes for each category:

- astronomy:
 - Discoveries of exoplanets
 - Black holes and their role in shaping galaxies
 - The search for extraterrestrial life
 - Gravitational waves and the study of the universe’s origins
 - The use of telescopes to explore the universe
 - The mysteries of dark matter and dark energy
 - Solar flares and their impact on Earth
 - The history of the universe and its evolution over time
 - Exploring the possibility of space tourism
 - The exploration of our neighboring planets, such as Mars and Venus.
- baseball:
 - Recent controversy surrounding sign-stealing scandal in MLB
 - Breakdown of top prospects in minor league baseball
 - Analysis of new rule changes for upcoming baseball season
 - Coverage of recent World Series champions and their success
 - In-depth profile of influential baseball figures, such as managers or players
 - Updates on retired players and their post-baseball careers
 - Highlighting standout performances by individual players or teams in recent games
 - Coverage of international baseball leagues and their top players
- economy:
 - Job market and employment rates
 - Interest rates and monetary policy
 - Inflation and deflation
 - Economic growth and GDP
 - Consumer spending and retail sales
 - Income inequality and poverty
 - GDP growth and contraction
 - Labor market trends
 - Economic impacts of natural disasters and pandemics
 - Housing market and real estate

G.2 Amazon Attributes

G.2.1 Product Brands

We randomly select 3 categories in Amazon dataset and display the corresponding product brand attributes for each category:

- camera_photo.:
 - Canon
 - Nikon
 - Leica
 - Hasselblad
 - Fujifilm
 - Lumix (Panasonic)
 - GoPro
 - Polaroid
 - Yashica
 - Mamiya
- health_and_personal_care.:

- Johnson & Johnson
- Dove
- Colgate
- Head & Shoulders
- Merck
- Pfizer
- Merck
- Abbott Laboratories
- GlaxoSmithKline
- toys_games.:
 - Mattel
 - Fisher-Price
 - Hasbro
 - Lego
 - Nerf
 - Barbie
 - Hot Wheels
 - Playmobil
 - MGA Entertainment
 - Paw Patrol

G.2.2 Product Names

We randomly select 3 categories in Amazon dataset and display the corresponding product name attributes for each category:

- camera_photo.:
 - SnapShooter
 - FocusReady
 - ClickCapture
 - MemoriesMatter
 - FlashFinder
 - PicturePerfect
 - ShotSnap
 - VisionVibe
 - PixelPal
 - FreezeFrame
 - AngleAttack
 - SnapStash
 - FrameFlair
 - SmileSnaps
 - ImageImpact
 - ZoomZest
 - ClickCraze
 - PixelPassion
 - ShootSmart
 - CaptionChamp.
- baby.:
 - Baby Swing
 - Diaper Genie
 - Milk Warmer

- Baby Carrier
- Car Seat
- Baby Monitor
- Baby Food Maker
- Nursery Glider
- Changing Table
- Baby Bouncer
- Playpen
- Teething Rings
- Baby Wipes Dispenser
- Baby Bath Tub
- Bibs
- Baby Blankets
- Pacifier Clip
- Baby Sling
- Baby Napper
- Moses Basket
- outdoor_living.:
 - Sunbrella
 - Weber
 - Keter
 - Trex
 - Solaira
 - Tropitone
 - Bimini Solutions
 - La-Z-Boy Outdoor
 - Suncast
 - Beltwood
 - Quikrete
 - Cosco Outdoor Living
 - Anova Furnishings
 - Duramax
 - US Polymers
 - Ostrich Chairs
 - Carefree of Colorado
 - Tuff Coat
 - Fire Sense
 - Heritage Patios.

G.3 Reddit Attributes

G.3.1 Resources

We randomly select 3 categories in Reddit dataset and display the corresponding resource attributes for each category:

- buddy_crossing.:
 - Meetup.com - a website that allows you to find and join groups of individuals with similar interests in your local area, including hiking, book clubs, and social events.
 - The Buddy System: Understanding Mental Illness and Addiction - a book that explores the biology of addiction and provides a guide for friends and family members of individuals struggling with these issues.

- Lynda.com - a subscription-based online learning platform that provides courses on a variety of subjects including computer programming, business, web design, and more.
 - Codecademy.com - an interactive online platform that teaches coding skills for free or through a subscription.
 - Khan Academy - a nonprofit organization that provides free online courses in a wide range of subjects including math, science, and humanities to learners of all ages.
 - Duolingo - a language-learning app that is available for free on the App Store and Google Play, offering courses in a variety of languages including Spanish, French, and German.
 - MindBody App - a mobile app that helps users find and book local fitness, wellness, and beauty activities such as yoga classes, massages, and haircuts.
 - Headspace - a meditation app that offers guided meditation courses to help users reduce stress and improve focus.
 - The Knot - a website that provides tools and resources for wedding planning, including a Wedding Website Builder, guest list tracker, and registry management.
 - Khan Academy - a nonprofit organization that provides free online courses in a wide range of subjects including math, science, and humanities to learners of all ages.
 - Others resource for buddy_crossing.
- whats_that_book.:
 - Goodreads - A social platform for book lovers where users can search for books, create bookshelves, and write reviews.
 - LibraryThing - A community-driven cataloging website where users can create and share their personal book collections.
 - AbeBooks - An online marketplace for rare and out-of-print books, as well as other antique or collectible items.
 - Shelfari - An online book club where users can share book recommendations and read reviews from others.
 - Project Gutenberg - A digital library of freely available public domain books.
 - Paperback Swap - A book trading community where users can exchange books with others across the US.
 - Goodreads Librarians Group - A community of Goodreads users who help with book cataloging, including identifying books from incomplete information.
 - Book Riot - A website featuring book reviews and book-related news, with an emphasis on diverse and underrepresented voices.
 - The New York Times Book Review - A renowned weekly publication featuring book reviews, author interviews, and literary criticism.
 - Others resource for whats_that_book.
- pokemongo_friends.:
 - Pokemon GO Hub: A comprehensive website dedicated to news, guides, and analysis on Pokemon GO.
 - The Silph Road Subreddit: A community-run subreddit dedicated to research and analysis of Pokemon GO mechanics.
 - Poke Assistant: A website that offers a range of tools to help you optimize your Pokemon GO experience, including IV calculators and gym battle simulations.
 - The Trainer Club: A YouTube channel that provides daily updates, news, and tips for Pokemon GO trainers.
 - Gotta Catch 'Em All: A Facebook group where you can connect with other Pokemon GO players and coordinate raid battles and other activities.
 - Reddit's r/PokemonGOFriends Subreddit: A community of players looking for friends to exchange gifts and share invites for raids.
 - The PokeMap: A website that allows you to find nearby Pokemon on a map in real-time.
 - Poke Genie: An app that automatically calculates IVs and other stats for your Pokemon, saving you time and headaches.

- Pokemon GO Gamepress: A website that offers detailed breakdowns and analysis of Pokemon, movesets, and other game mechanics.
- The Go Ranger App: An app that helps you plan your raids and battles, with intuitive mapper tools and filters to help you find the Pokemon you're looking for.
- Others resource for pokemongo_friends.

G.3.2 Experience

We randomly select 3 categories in Reddit dataset and display the corresponding experience attributes for each category:

- build_a_pc.:
 - DIY PC Builds: Sharing personal experiences and success stories of building custom PCs, discussing component choices, troubleshooting, and performance optimizations.
 - Budget-Friendly Builds: Discussing experiences with building PCs on a tight budget, sharing cost-saving tips, and recommendations for budget-friendly components.
 - Cable Management: Sharing personal experiences and tips for effective cable management in PC builds, discussing cable routing techniques and showcasing clean build aesthetics.
 - RGB Lighting: Discussing experiences with RGB lighting setups in PC builds, sharing recommendations for RGB components, software customization, and lighting effects.
 - Troubleshooting Builds: Sharing experiences and tips for troubleshooting common issues in PC builds, helping fellow builders diagnose and solve hardware or software problems.
 - Silent and Quiet PC Builds: Discussing experiences and recommendations for building silent or quiet PCs, focusing on noise reduction techniques and quiet component choices.
 - Workstation Builds: Sharing experiences and insights into building PCs for professional workloads, such as video editing, 3D rendering, programming, and graphic design.
 - Water-Cooling Adventures: Sharing experiences and insights into custom water-cooling loops, discussing the challenges, benefits, and performance improvements achieved.
 - Unique and Custom Builds: Showcasing and discussing unique and custom PC builds, including themed builds, custom cases, or exotic cooling solutions.
 - Build Planning and Component Selection: Discussing experiences with planning PC builds, researching and selecting components, considering compatibility, and balancing performance and budget.
 - Modding and Case Customization: Sharing experiences with PC case modding and customization, discussing techniques, materials, and showcasing personal projects.
 - Compact and Small Form Factor Builds: Discussing experiences with building compact or small form factor PCs, sharing recommendations for mini-ITX cases, cooling solutions, and component choices.
 - Home Server and NAS Builds: Sharing experiences and insights into building home servers and network-attached storage (NAS) systems, discussing storage options, software, and data management.
 - Multimonitor Setups: Discussing experiences with multimonitor setups, sharing tips for optimizing productivity and gaming experiences across multiple displays.
 - PC Gaming Peripherals: Sharing experiences and recommendations for gaming peripherals, such as keyboards, mice, monitors, and headsets, discussing features and personal preferences.
- xbox.:
 - Xbox One exclusive games such as Halo 5, Forza Horizon 4, and Gears of War 4
 - Xbox One media and entertainment apps such as Netflix and Hulu
 - memorable gaming moments or achievements on the Xbox console.
 - Purchase Xbox One online.
 - Xbox Kinect motion sensor accessory

- Xbox Play Anywhere program
- Other Experience of Xbox
- metal_gear_solid.:
 - Tactical Weapon Customization: Experimenting with various weapons, attachments, and equipment to tailor loadouts to different mission objectives and playstyles.
 - Character Development: Witnessing the growth and development of iconic characters such as Solid Snake, Big Boss, or Raiden throughout their respective story arcs.
 - Stealthy Takedowns: Executing silent and non-lethal takedowns, utilizing tranquilizer darts, chokeholds, or sneaking up on enemies from behind.
 - Tactical Planning: Strategizing mission approaches, analyzing enemy patrols, setting traps, and utilizing distractions to gain the upper hand.
 - Memorable Characters: Developing connections with unique and memorable characters, such as Otacon, Meryl, Revolver Ocelot, or The Boss.
 - Stealthy Gadgets: Utilizing gadgets and tools, such as the cardboard box, thermal goggles, or the iconic Solid Eye, to gain advantages during missions.
 - Emotional Storytelling: Experiencing emotionally impactful moments within the narrative, exploring themes of loss, betrayal, loyalty, and the human cost of warfare.
 - Espionage Tactics: Participating in undercover missions, gathering intelligence, infiltrating enemy bases, and sabotaging their operations.
 - Lore and Mythology: Delving into the intricate lore, conspiracies, and historical events within the Metal Gear Solid universe, including topics like The Patriots or the Philosopher's Legacy.
 - Groundbreaking Game Design: Appreciating the innovative gameplay mechanics, cinematic presentation, and attention to detail that have made the Metal Gear Solid series a beloved and influential franchise in the gaming industry.
 - Others experience for metal_gear_solid.

G.4 StackExchange Attributes

G.4.1 Scenarios

We randomly select 3 categories in StackExchange dataset and display the corresponding scenario attributes for each category:

- multiplayer.:
 - Cheating/hacking in online games
 - Inappropriate player behavior
 - Unbalanced game mechanics
 - Difficulty connecting to multiplayer servers
 - Matchmaking errors
 - Unresponsive or laggy gameplay
 - Glitches in gameplay affecting online matches
 - Difficulty finding players to match with
 - Balancing player skills in matchmaking
 - Disconnects and dropped connections mid-game
 - Cross-platform compatibility issues
 - In-game communication problems
 - Difficulty managing and moderating game servers
 - Addressing griefing and trolling in multiplayer games
 - Managing player accounts and login systems
 - Implementing or improving anti-cheat measures
 - Community feedback and feature requests
 - Addressing game-breaking exploits
 - Ensuring fair and accurate reporting of player statistics

- Addressing server crashes and downtime.
- procedural-generation.:
 - Improving the efficiency of procedural generation algorithms in Python.
 - Troubleshooting issues with memory usage in large-scale procedural generation projects.
 - Debugging issues with randomized content generation in procedural levels.
 - Implementing procedural generation techniques in C++ for game development.
 - Exploring the potential of machine learning in procedural content generation.
 - Optimizing the generation of complex 3D models using procedural algorithms.
 - Managing complexity and maintaining consistency in procedurally generated game worlds.
 - Addressing issues with procedural generation of text-based content, such as dialogue or item descriptions.
 - Developing tools to aid in the creation and testing of procedural generation algorithms.
 - Balancing the need for randomness with player expectations for fairness and balance.
 - Addressing issues with the procedural generation of music and sound effects.
 - Improving the visual quality of procedurally generated game assets.
 - Exploring ethical concerns around the use of AI in procedural content generation.
 - Developing procedural generation techniques for non-linear narratives or branching storylines.
 - Improving the procedural generation of inhabited environments, such as procedurally generated NPCs.
 - Addressing issues with the procedural generation of terrain features such as rivers and mountains.
 - Implementing procedural generation techniques for user-generated content.
 - Supporting multithreaded execution in procedural generation algorithms.
 - Ensuring procedural generation techniques are compatible with various game engines and frameworks.
 - Improving the scalability of procedural generation algorithms for use in multiplayer games.
- networking.:
 - Difficulty in troubleshooting network connection issues on a Linux system
 - Configuring a wireless access point for a large office space
 - Implementing load balancing across multiple servers in a network
 - Optimizing network throughput to reduce latency in a gaming environment
 - Implementing firewall rules to block unauthorized access to a network
 - Troubleshooting DNS resolution issues on a Windows server
 - Designing and implementing a secure VPN connection
 - Setting up a network file server for shared access among multiple clients
 - Configuring SNMP to monitor network traffic and utilization
 - Designing a network topology for a large enterprise with multiple locations
 - Troubleshooting issues with Ethernet switches in a data center environment
 - Implementing QoS to prioritize network traffic for critical applications
 - Configuring NAT and PAT to enable internet access for multiple devices
 - Setting up and configuring VLANs to segment a network
 - Troubleshooting issues with network printers in an office environment
 - Configuring routing protocols in a large network
 - Securing wireless access points to prevent unauthorized access
 - Troubleshooting issues with VPN connection stability and speed
 - Implementing network virtualization with virtual LANs and virtual switches
 - Designing and implementing an effective network security strategy to prevent data breaches.

G.5 ArXiv Attributes

G.5.1 Subtopics

We randomly select 3 categories in arXiv dataset and display the corresponding product brand attributes for each category:

- machine_learning.:
 - Text generation
 - Natural language understanding for chatbots
 - Sentiment analysis and opinion mining
 - Text summarization and keyword extraction
 - Machine translation
 - Named entity recognition and entity linking
 - Dialogue systems and conversational agents
 - Cross-lingual and Multilingual NLP
 - Text-to-speech systems
 - Phonetics and phonology in computational linguistics
 - Grammatical error detection and correction
 - Speech recognition and acoustic modeling
 - Semantic role labeling
 - Discourse analysis and coherence modeling
 - Lexical semantics and word sense disambiguation
 - Computational lexicography and machine-readable dictionaries
 - Language Modeling
 - question answering
 - Language resources and corpora
 - Computational sociolinguistics and dialectology.
- number_theory.:
 - Prime numbers
 - Diophantine equations
 - Modular arithmetic
 - Cryptography
 - Continued Fractions
 - Pell’s Equation
 - Fermat’s Last Theorem
 - Algebraic Number Theory
 - Riemann Hypothesis
 - Arithmetic Geometry
 - Quadratic Forms
 - L-Functions
 - Automorphic Forms
 - Galois Theory
 - Ramsey Theory
 - Distribution of Prime Numbers
 - Number Theory in Cryptography
 - Summation Formulas
 - Gaussian Integers
 - The Goldbach Conjecture
- geophysics.:
 - Seismic imaging

- Earthquake prediction
- Geothermal energy
- Volcanic eruptions
- Plate tectonics
- Geomagnetism
- Paleomagnetism
- Geophysical surveying
- Geophysical fluid dynamics
- Gravity measurements
- Rock physics
- Crustal deformation
- Geomorphology
- Mineral exploration
- Earth structure modeling
- Geodetic techniques
- Hydrogeophysics
- Earth modeling
- Electrical geophysics
- Remote sensing geophysics

G.5.2 Techniques

We randomly select 3 categories in the arXiv dataset and display the corresponding attributes for each category:

- genomics.:
 - Genome assembly and annotation using hybrid approaches.
 - Comparative genomics for analyzing evolutionary relationships between genomes.
 - Differential gene expression analysis using RNA sequencing data.
 - Metagenomics for studying the microbial communities in different environments.
 - Epigenetic analysis for understanding gene regulation.
 - Network analysis for identifying gene interactions and pathways.
 - Structural variation analysis for detecting genomic rearrangements.
 - Functional genomics for studying gene function and pathway regulation.
 - Genome-wide association studies for identifying genetic variants associated with complex traits.
 - High-throughput screening methods for identifying genes involved in specific biological processes.
- number_theory:
 - Primality testing using elliptic curves
 - Continued fraction factorization method
 - Algorithm for solving Diophantine equations
 - Quadratic sieve algorithm for integer factorization
 - Pollard rho algorithm for integer factorization
 - Digital sum subtraction method for computing discrete logarithm
 - Fermat’s method for factorization of primes
 - Chinese remainder algorithm for solving modular equations
 - Exponential-sum algorithm for computing in algebraic number fields
 - Generalized Ramanujan-Selberg formula for counting integer points on algebraic varieties.
- geophysics.:
 - Seismic attribute interpretation

- Full waveform inversion
- Gravity inversion
- Spherical geometries
- Ground penetrating radar imaging
- Time-lapse reservoir monitoring
- Electrical resistivity tomography
- Joint inversion of geophysical data
- Radiometric dating
- Geomagnetic field modeling

G.6 AG News Attributes

G.6.1 Subtopics

The corresponding subtopic attributes for each category are shown as follows:

- business:
 - Corporate earnings and financial reports
 - Stock market updates and analysis
 - Mergers and acquisitions
 - Business regulations and policies
 - Startups and entrepreneurship
 - Industry trends and forecasts
 - Economic indicators and market trends
 - Business strategies and management practices
 - Corporate governance and ethics
 - Consumer behavior and market research
 - Business leadership and executive profiles
 - Banking and finance industry updates
 - Energy and sustainability in business
 - Retail and e-commerce trends
 - Real estate and property market updates
 - Business disruptions and crisis management
 - Corporate social responsibility and sustainability initiatives
- sci_tech:
 - Artificial intelligence
 - Robotics
 - Quantum computing
 - Biotechnology
 - Nanotechnology
 - Internet of Things
 - Renewable energy
 - Virtual reality
 - Augmented reality
 - Cybersecurity
 - Genetic engineering
 - Big data
 - Autonomous vehicles
 - 3D printing
 - Blockchain technology
 - Bioinformatics
 - Machine learning

- Biomedical engineering
- Clean technology
- sports:
 - Soccer
 - Basketball
 - Baseball
 - Tennis
 - Golf
 - Cricket
 - Rugby
 - Athletics
 - Formula 1
 - Olympics
 - Boxing
 - Swimming
 - Volleyball
 - Ice hockey
 - American football
 - Cycling
 - Motorsports
 - Martial arts
 - Horse racing
 - Surfing
- world:
 - International politics and diplomacy
 - Global conflicts and war
 - Terrorism and security threats
 - Human rights issues and social justice movements
 - Migration and refugee crises
 - Climate change and environmental policies
 - Global health crises and pandemics
 - Natural disasters and emergencies
 - Cross-border crime and corruption
 - Cultural and social developments worldwide
 - Geopolitical tensions and territorial disputes
 - International aid and development efforts
 - Humanitarian crises and relief efforts
 - Cultural heritage preservation and promotion
 - International collaborations and partnerships

G.7 SST-2 Attributes

G.7.1 Subtopics

We display the corresponding subtopic attributes for each category as follows:

- positive:
 - Compelling Storyline: A strong and engaging narrative that captures the audience’s attention from beginning to end.
 - Well-Developed Characters: Memorable and relatable characters that evoke emotions and drive the story forward.

- Skillful Direction: Effective direction that showcases the filmmaker’s vision, ensuring cohesive storytelling and engaging visual elements.
- Excellent Acting: Convincing performances from the cast that bring the characters to life and immerse the audience in the story.
- Cinematography: Expertly captured visuals, including the use of framing, lighting, and camera movements, to enhance the storytelling and create a visually appealing experience.
- Engaging Dialogue: Well-written dialogue that is natural, meaningful, and contributes to character development and plot progression.
- Sound Design and Music: Thoughtful and immersive sound design, including sound effects and a well-curated soundtrack or original score, that enhances the overall cinematic experience.
- Production Design: Attention to detail in creating visually appealing and authentic sets, costumes, and overall aesthetics that contribute to the film’s atmosphere and world-building.
- Editing: Skillful editing that maintains a good pace, effectively transitions between scenes, and enhances the overall flow and impact of the story.
- Emotional Impact: A movie that evokes emotions, whether it be through humor, drama, suspense, or other means, leaving a lasting impression on the audience.
- negative:
 - Weak Plot: A poorly developed or uninteresting storyline that fails to engage the audience.
 - Lackluster Performances: Unconvincing or uninspired performances by the actors that fail to bring the characters to life.
 - Poor Production Quality: Subpar production values, including low-quality visuals, amateurish cinematography, and weak special effects.
 - Incoherent Storytelling: Confusing or disjointed narrative structure that makes it difficult to follow or understand the plot.
 - Unmemorable Characters: Underdeveloped or forgettable characters that fail to resonate with the audience.
 - Weak Soundtrack: A forgettable or poorly composed soundtrack that fails to enhance the mood or add depth to the movie.
 - Poor Dialogue: Uninteresting or poorly written dialogues that fail to engage or resonate with the audience.
 - Disjointed Atmosphere: A lack of coherence or consistency in creating an immersive and believable world for the viewers.
 - Unresolved Plotlines: Loose ends or unresolved plotlines that leave the audience feeling unsatisfied or confused.
 - Lack of Entertainment Value: A movie that fails to deliver an enjoyable or engaging experience for the audience, leaving them feeling bored or uninterested.

G.7.2 Descriptive Details

We use movie genres as the characteristics of movies, and the attributes are listed as follows:

- Action
- Drama
- Comedy
- Thriller
- Romance
- Horror
- Adventure
- Science Fiction
- Fantasy
- Animation

G.8 Yelp Attributes

G.8.1 Subtopics

We display the corresponding subtopic attributes for each category as follows:

- positive:
 - Quality of Food: The taste, flavor, and presentation of the dishes.
 - Fresh Ingredients: The use of fresh and high-quality ingredients in the preparation of the food.
 - Menu Variety: A diverse range of options catering to different dietary preferences and restrictions.
 - Presentation: The visually appealing presentation of the dishes.
 - Service: Attentive, friendly, and prompt service from the restaurant staff.
 - Value for Money: Offering good quality and portion sizes at reasonable prices.
 - Cleanliness: A clean and well-maintained dining area, including tables, utensils, and restrooms.
 - Special Dietary Accommodations: Catering to specific dietary needs such as vegetarian, vegan, gluten-free, etc.
 - Unique and Creative Dishes: Offering innovative and creative dishes that stand out.
 - Efficient Operations: Smooth and well-coordinated operations to minimize waiting times and delays.
- negative:
 - Poor Service: Slow or inattentive service from the restaurant staff. Unfriendly Staff: Rude or unhelpful behavior from the restaurant staff.
 - Long Waiting Times: Excessive waiting times for a table or food.
 - Incorrect Orders: Receiving incorrect or poorly prepared food orders.
 - Unappetizing Presentation: Dishes that are poorly presented or lack visual appeal.
 - Unpleasant Ambience: Uncomfortable or uninviting atmosphere in the restaurant.
 - Dirty or Unhygienic Conditions: Lack of cleanliness in the dining area, restrooms, or utensils.
 - Limited Menu Options: A limited selection of dishes or lack of variety.
 - Poor Food Quality: Dishes that are poorly cooked, tasteless, or of low quality.
 - Overpriced: Excessive prices for the quality and portion sizes of the food.

G.8.2 Descriptive Details

We use cuisine types as the characteristics of restaurants, and the attributes are listed as follows:

- Turkish
- Spanish
- Greek
- Italian
- French
- American
- Mexican
- Canadian
- Cajun
- Tex-Mex
- Brazilian
- Peruvian
- Argentinean

- Colombian
- Venezuelan
- Ethiopian
- Moroccan
- South African
- Nigerian
- Egyptian
- Chinese
- Japanese
- Indian
- Thai
- Korean
- Australian
- New Zealand
- Polynesian
- Hawaiian
- Singaporean

H Examples for Filtered Attribute Values

Here we give some examples of the filtered attributes.

For the Amazon product review dataset, some filtered attributes are listed as follows.

- beauty:
 - Hair Dryer (close to health and personal care)
 - Hair Straightener (close to health and personal care)
- electronics:
 - Car dashcam (close to automotive)
 - Wireless earbuds (close to cell_phones_service)
- office_products:
 - Mouse pad (close to electronics)

For NYT dataset, some filtered attributes are listed as follows:

- american_football:
 - The economic impact of football on local communities and businesses.
 - The role of nutrition and hydration in optimal football performance
 - The history of football’s integration with America’s armed forces
- stocks_and_bonds:
 - The role of artificial intelligence in stock market prediction
 - The effects of climate change on bond performance
- movies:
 - The future of cinema in a post-COVID world
 - The role of international markets in Hollywood’s global success

For Reddit dataset, some filtered attributes are listed as follows:

- amc_stock:

- Working in the stock market industry.
- Impact of COVID-19 on amc_stock.
- Having a background in finance or business.
- what_car_should_i_buy:
 - Towing capacity and trailer options
 - Financing and leasing options
 - Road trip and travel experiences
- pittsburgh:
 - Attending a game at PNC Park
 - Watch a Steelers football game at Heinz Field

For StackExchange dataset, some filtered attributes are listed as follows:

- algorithm:
 - Debugging an optimization algorithm for a complex logistics problem.
 - Improving the robustness of an encryption algorithm against attacks.
 - Finding a way to better handle overflow errors in a cryptographic algorithm.
- iPhone:
 - Unable to use certain features in certain regions/countries
 - Trouble with GPS not working properly
- opengl:
 - Difficulty setting up a custom camera in OpenGL
 - Inconsistent rendering of 3D models in OpenGL
 - Difficulty integrating OpenGL with a cross-platform game engine

I Examples for Merged Attribute on the arXiv dataset

Here we give some examples of the filtered attributes.

- optimization techniques for computer vision:
 - computer vision and pattern recognition
 - machine learning
- black holes:
 - classical physics
 - popular physics
- evolving networks for financial trading:
 - neural and evolutionary computing
 - computational finance
- multi-lingual information retrieval:
 - information retrieval
 - computation and language
- iot operator systems and their security challenges:
 - operating systems
 - cryptography and security
- complexity of social networks:
 - computational complexity
 - social and information networks

J The Prompt Format for MetaPrompt

In this section, we provide more details about the implementation of MetaPrompt [43]. It first leverages LLMs to provide additional task-specific information, and then augments the class-dependent prompt with the task-specific knowledge.

Prompts to ChatGPT: [Original Prompt]. What does this task ask us to do?

Response from ChatGPT for NYT dataset: The task of generating an example of a NYT news asks the model to create high-quality, informative, and engaging news articles that mimic the writing style and structure commonly found in The New York Times (NYT).

Response from ChatGPT for Amazon dataset: The task of generating an example of a product review asks the model to create a detailed and informative review for a particular product, highlighting its features, performance, pros and cons, and overall user experience. This review should be helpful for potential buyers, providing them with valuable insights to make an informed purchasing decision.

Response from ChatGPT for Reddit dataset: The task of generating an example of a post similar to Reddit asks the model to craft a well-written and engaging post that adheres to the guidelines and community standards of the specific subreddit it's intended for. The post should spark discussions, share valuable information, or showcase interesting content that resonates with the subreddit's audience.

Response from ChatGPT for StackExchange dataset: The task of generating an example of a post similar to StackExchange asks the model to compose a clear and concise question or answer related to a specific topic within the given StackExchange community. The post should be well-researched and supported by relevant sources or personal expertise, aiming to provide valuable insights or solutions to other users' queries.

K Dataset Information

K.1 Dataset Documentations

The dataset is provided in *json* format; there are three json files corresponding to the original `train`, `validation` and `test` split. We also include two additional file, `attrprompt` and `simprompt`, which is generated by AttrPrompt and SimPrompt as the synthetic training data.

Each data point contains the following fields:

- `label`: the label for the example. For multi-class classification, the `label` field is an integer, while for multi-label classification, the `label` field is a `list[int]` containing one or multiple integers as each example may refer to multiple classes;
- `text`: a content of each example.

K.2 Intended Uses

AttrPrompt and SimPrompt are intended for researchers in machine learning, natural language processing, and related fields to innovate novel methods for training data generation problems.

K.3 Hosting and Maintenance Plan

The codebase is hosted and version-tracked via GitHub. It will be available under the link <https://github.com/yueyu1030/attrprompt>. The download link of all the datasets can be found in the Github repository.

Note that it is a community-driven and open-source initiative. We are committed and have the resources to maintain and actively develop it for at minimum the next five years. We plan to grow the GitHub repo by including new tasks and datasets and warmly welcome external contributors.

K.4 Licensing

We license our work using Apache 2.0¹². All of the `train/validation/test` data are publicly released by previous work [3, 31, 15].

K.5 Author Statement

We the authors will bear all responsibility in case of violation of rights.

¹²<https://www.apache.org/licenses/LICENSE-2.0>