

Leveraging Variation Theory in Counterfactual Data Augmentation for Optimized Active Learning

Simret Araya Gebreegziabher

sgebreeg@nd.edu

University of Notre Dame

Kuangshi Ai

kai@nd.edu

University of Notre Dame

Zheng Zhang

zzhang37@nd.edu

University of Notre Dame

Elena L. Glassman*

glassman@seas.harvard.edu

Harvard University

Toby Jia-Jun Li*

toby.j.li@nd.edu

University of Notre Dame

Abstract

Active Learning (AL) allows models to learn interactively from user feedback. However, only annotating existing samples may hardly benefit the model’s generalization. Moreover, AL commonly faces a cold start problem due to insufficient annotated data for effective sample selection. To address this, we introduce a counterfactual data augmentation approach inspired by Variation Theory, a theory of *human concept learning* that emphasizes the essential features of a concept by focusing on what stays the same and what changes. We use a neuro-symbolic pipeline to pinpoint key conceptual dimensions and use a large language model (LLM) to generate targeted variations along those dimensions. Through a text classification experiment, we show that our approach achieves significantly higher performance when there are fewer annotated data, showing its capability to address the cold start problem in AL. We also find that as the annotated training data gets larger, the impact of the generated data starts to diminish. This work demonstrates the value of incorporating human learning theories into the design and optimization of AL.

1 Introduction

Active learning (AL) allows users to provide focused annotations to integrate human preferences and domain knowledge into machine learning models (Settles, 2009). It relies on a human’s iterative annotations to build and refine model performance (Budd et al., 2021). As a result, the model’s performance improvement with each annotation round depends on both the quality and quantity of annotated data. However, AL faces a cold start problem: in the early stages, when annotated data is limited, the model is often unstable and struggles to make informed decisions about which instances

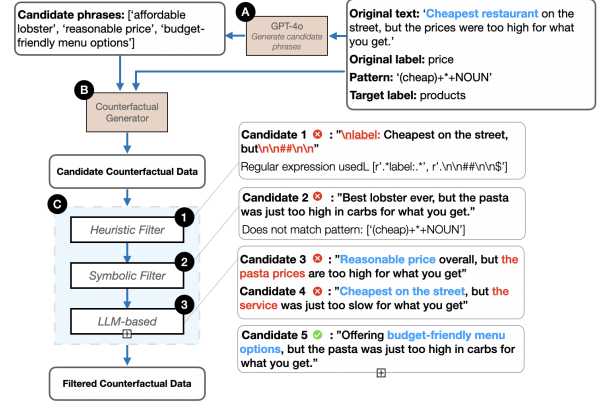


Figure 1: Our approach combines neuro-symbolic patterns with in-context learning to generate Variation Theory-based counterfactual examples for active learning.

to query for labeling, which hinders its initial performance (Yuan et al., 2020).

Counterfactual data augmentation techniques have been shown to enhance model performance (Yang et al., 2022a; Wang and Culotta, 2020; Reddy et al., 2023). Synthesized counterfactual data can be more effective in capturing meaningful variations than real data selected from the dataset. However, the scalable generation and selection of augmented data have been a consistent challenge (Liu et al., 2022; Li et al., 2023a). To address this, DISCO (Chen et al., 2023) proposed a method for automatically generating counterfactual data using task-agnostic models. Despite its robust approach to augmented data, DISCO’s use of a black-box pipeline makes debugging and improving the model difficult and does not allow meaningful presentation of variations that facilitate effective human annotation and sensemaking.

To address this, we propose a counterfactual generation pipeline that uses neuro-symbolic patterns to identify important features and uses them

*Co-senior authors contributed equally.

to guide the LLM’s counterfactual generation*. To motivate this approach, we draw on neuro-symbolic AI, which combines the representational power of neural networks with the interpretability and structure of symbolic reasoning (Hitzler and Sarker, 2022). Neuro-symbolic models integrate learned patterns with human-understandable rules, enabling systems to generalize in a transparent reasoning process. In the context of counterfactual generation, this hybrid approach allows us to generate examples that vary meaningfully along conceptually relevant dimensions while maintaining structural and semantic consistency. Specifically, we use a programming-by-example approach (Gulwani, 2011) to generate neuro-symbolic patterns (Gebregziabher et al., 2023). These patterns capture the syntactic and semantic similarities among similarly labeled examples. We then use the learned patterns to guide the LLM to generate counterfactual examples to be used in consecutive rounds of model re-training. The generated counterfactual examples change the assigned label into a different label while still keeping the original symbolic pattern in the data. In doing so, the generated examples introduce more meaningful variability in the data for subsequent model training. To further ensure the quality of the generated counterfactual examples, we design a three-step automatic filtering pipeline.

This paper makes the following contributions:

Evaluating the quality of generated counterfactual examples We assess the quality of generated counterfactual examples using a three-stage filtering mechanism. We define a high-quality counterfactual as a sentence that eliminates the original label (soft flip) while introducing the target label (label flip). The results show a high Soft Label Flip Rate (SLFR)—the rate of removal of original labels from counterfactual examples, and a high level of consistency in Label Flip Rate (LFR)—the rate of changing original labels into target labels in generated counterfactual examples. By evaluating how often new examples meaningfully alter the original label and capture valuable variations, we can assess the efficacy of the examples produced.

Evaluating the effectiveness of Variation Theory in active learning We investigate how incorporating Variation Theory into active learning

can improve robustness and address cold-start challenges (Yuan et al., 2020). Using a classification task, we compare our counterfactual-based method against four baselines—random, cluster-based, uncertainty-based selection, and counterfactuals without Variation Theory. Across three datasets and two models, our approach achieves up to 2× higher performance with fewer than 70 annotations. The benefits diminish as annotation volume grows, highlighting its effectiveness in low-data, cold-start settings. We also analyze the roles of annotation selection, syntactic diversity, and semantic diversity in driving this performance gain.

2 Related Work

2.1 Active Learning

Active Learning (AL) in machine learning is an approach in which the learning algorithm selectively chooses informative data points for model training. Although most sampling strategies rely on a pool of unlabeled data (Fu et al., 2013), there are strategies that synthesize data points in real time for annotation (Schumann and Rehbein, 2019). The second approach, also called Membership Query Synthesis (MQS) creates new examples that inform the model with more representative scenarios by either modifying existing instances (Wu et al., 2023, 2021) or generating new instances (Schumann and Rehbein, 2019).

In domains with scarce annotated data, data augmentation methods aim to enhance the quantity and quality of training data (Yang et al., 2022b). Traditional data augmentation techniques, such as geometric transformations and color space alterations, do not modify the fundamental causal generative process. As a result, they do not counteract biases like spurious correlations (Kaushik et al., 2021).

2.2 Data Generation and Augmentation

Counterfactual data augmentation has been widely used to counteract spurious correlations in data (Denton et al., 2020; ?; Yang et al., 2022a; Wang and Culotta, 2020). This approach employs counterfactual inference to control generative factors, facilitating the generation of samples that can address confounding biases. Many existing strategies use dataset-specific counterfactual augmentation methods in specific domains, such as sentiment analysis (Yang et al., 2022a; Kaushik et al., 2020), named entity recognition (Ghaddar et al., 2021), text classification (Wang and Culotta, 2020), and

* www.github.com/SimretA/Variation-Theory-in-Counterfactual-Data-Augmentation

neural machine translation (?). A popular approach to address spurious dependence in NLP datasets is to use human-guided counterfactual augmentation through crowd sourcing (Kaushik et al., 2021; Joshi and He, 2022). This approach presents individuals with data and preliminary labels, asking them to modify the data for an alternate label while avoiding unnecessary edits (Kaushik et al., 2020). This method depends on human efforts and expertise to overcome the challenge of automatically translating raw text into important features.

LLMs have been shown to possess extensive generative capacity, making them useful tools for counterfactual data generation. Li et al. (2023a) introduced a method utilizing LLMs to generate domain-specific counterfactual samples through prompt design, highlighting the alignment between the efficacy of LLMs in domain-specific counterfactual generation and their overall proficiency in that domain. Although in-context learning has been a promising direction to get LLMs to perform different tasks, Min et al. (2022) identified several key factors that influence its effectiveness, including the demonstration of the label space, the input text distribution, and the overall sequence format.

A consistent challenge in counterfactual generation has been the scalable generation and selection of augmented data (Liu et al., 2022; Li et al., 2023a). To address this, DISCO (Chen et al., 2023) introduced a method for automatically generating high-quality counterfactual data using task-agnostic “teacher and student” models to allow classifier models to learn causal representations. DISCO uses a neural syntactic parser to select the spans of the sentence to vary on to generate data using Large Language Models (LLMs). Although DISCO provides more robust models trained on augmented data, the use of black-box approaches to generate data could make model debugging and improvement harder. To address this, we adopt a neuro-symbolic approach to define the concept boundaries in user annotations (Gebreegziabher et al., 2023).

2.3 Example-based Learning via Variation Theory

Based on previous studies on LLMs as counterfactual generators, our work seeks to learn from human cognition and example-based learning to better guide LLMs to generate higher quality data. *Will educational theories that work for human learners also work for AI?* Decades of research have

demonstrated that using example-based learning constitutes an effective instructional strategy for humans acquiring new skills (Gog and Rummel, 2010). Few-shot learning is an example-based learning method commonly used by LLMs.

How can we use human learning theories to support the annotation of data and training of LLM classifiers? Variation Theory (Marton, 2014), rooted in human learning research, gives us insights from human experience, e.g., (Cheng, 2016). The core concept of this theory involves presenting sets of examples that vary along specific dimensions, enabling learners to identify and conceptualize the dimensions as a useful coordinate space for describing instantiations of the underlying concept. This aligns with the foundational principle of counterfactual data augmentation in machine learning.

3 Approach

Our approach applies the Variation Theory of human learning to machine learning in the context of active learning (AL). We propose a new approach of counterfactual data generation by combining neuro-symbolic methods and LLMs. Specifically, we use domain-specific neuro-symbolic patterns to learn the syntactic representation of similarly labeled data that define a neuro-symbolic model’s learning space and concept boundaries. We then use the learned patterns to guide the generation of augmented data that helps a classification model learn important nuances about each label (Fig. 1-A,B).

Through this approach we generate counterfactual data that are *syntactically similar* to their original counterparts but semantically belong to a different label. To ensure the quality of the generated counterfactuals, we apply a three-level filtering mechanism (Fig. 1-C).

3.1 Using Neuro-symbolic Patterns to Define Concept Space

Variation Theory suggests that humans learn a concept most effectively when they are shown examples that vary in only one specific dimension at a time, while all other aspects stay the same. Therefore, an important aspect of Variation Theory is determining which features should vary to emphasize their effects in the learning process. We achieve this by learning critical features from labeled data by generating neuro-symbolic patterns and make small modifications on the original sentence while

maintaining consistency along the generated pattern.

3.1.1 Learning Neuro-symbolic Patterns

We use a programming-by-example (Lieberman, 2001) approach to establish the boundaries of concepts defined by data points and their associated ground truth labels. While our simulation study currently relies on ground truth labels, these will be substituted with human annotations in forthcoming interactive systems. After we randomly select a few annotations, we use PaTAT’s (Gebreegziabher et al., 2023) interactive program synthesis approach to generate domain-specific pattern rules that match the annotated examples. These pattern rules represent the lexical, syntactic, and semantic similarities of data under the same label. PaTAT’s pattern language includes the following components:

- Part-of-speech (POS) tags: VERB, PROP, NOUN, ADJ, ADV, AUX, PRON, NUM
- Word stemming: [WORD] (e.g., [have] will match all variants of have, such as *had*, *has*, and *having*)
- Soft match: (word) (e.g., (pricey) will match synonyms such as *expensive* and *costly*, etc.)
- Entity type: \$ENT-TYPE (e.g., \$LOCATION will match phrases of location type, such as *Houston, TX* and *California*; \$DATE will match dates; \$ORG will match names of organizations)
- Wildcard: * (will match any sequence of words)

Although the fundamental patterns are suitable for general domain text data, it is feasible to expand the pattern language to include specialized or domain-specific patterns.

This method generates a collection of regex-like patterns (but with semantically-enhanced tags) that match with the labeled positive examples while excluding the labeled negative examples. For example, if two data points in the domain of restaurant review “*Good food with great variety.*” and “*The food was amazing.*” have the same label “products”, PaTAT learns up to 5 patterns that collectively match the set of examples annotated with that label. In this case, two patterns match both sentences, i.e., “[food]+*+ADJ”, “(amazing)+*+”.

3.1.2 Using Neuro-symbolic Patterns for Counterfactual Data Generation

Using the learned neuro-symbolic patterns, we generate counterfactual examples by modifying the original text to be about a different label while still keeping the original pattern. To ensure minimal modifications and to make sure the reason for the original label is kept, we begin by generating candidate phrases for segments of the original sentence that matched the neuro-symbolic pattern (Fig. 1-A).

We use the generated candidate phrases as constraints to be included in the generated sentence. For example, in Fig. 1, the pattern (*cheap*)+*+NOUN has candidate phrases [‘*affordable lobster*’, ‘*reasonable price*’, ‘*budget-friendly menu*’]. When generating the counterfactual example, we instruct the LLM to always include one of these phrases in the modified sentence. This constraint ensures that counterfactual examples that vary in semantic content remain within the syntactic boundaries set by the pattern, which defines, at least in part, the particular label for which counterexamples are being generated (Fig. 1-B).

3.2 Filtering Generated Counterfactual Data

The ideal counterfactual example is a complete and coherent sentence that should keep the patterns of the original text, and successfully flip the original label to the target label. To ensure the quality of the fine-tuning dataset, we implement a three-stage filtering mechanism:

3.2.1 Regex Heuristic Filtering

We use a heuristic-based filter to identify and remove counterfactual data with common generation flaws. This filter ensures that the generated sentences are coherent and complete. This method uses regular expressions to detect common generation errors observed during our experimentation (Fig. 1-C1). We define rules to identify error patterns such as repetition of the prompt, inaccurate formatting, and incomplete generation, which were some common pitfalls we observed during generation.

3.2.2 Neuro-symbolic Filtering

The neuro-symbolic filter ensures that the generated counterfactual examples retain the original learned pattern. The original patterns represent features the model learns as useful conceptual boundaries. Therefore, keeping them in the counterfactually generated examples challenges the model’s cur-

rent boundary. To achieve this, we implement the filter using executable neuro-symbolic patterns defined in § 3.1. Specifically, we check whether each generated counterfactual example matches its original counterpart’s neuro-symbolic pattern (Fig. 1-C2). This filter excludes generated counterfactual examples that do not match with the provided pattern from being used in the consecutive training pipeline. To quantify this over the generated counterfactual examples, we calculate the pattern keeping rate (PKR) as defined below.

$$PKR = \frac{1}{N} \sum_{n=1}^N \mathbb{1}(\hat{p}_n = p_n)$$

where p_n is the original pattern, \hat{p}_n is the pattern given to the counterfactual data, and N is the size of the counterfactual data.

3.2.3 LLM-based Discriminator Filtering

Finally, we apply a filter using a GPT-4o discriminator. This filter removes counterfactuals that still keep their original label and all counterfactuals that do not change the label to the target label (Fig. 1-C3). This filter makes sure that the generated counterfactual examples have enough semantic changes that changes the original label to the target label. We adopt two matrices (Chen et al., 2023) to quantify this: the Label Flip Rate (LFR), and the Soft Label Flip Rate (SLFR) as defined below:

$$LFR = \frac{1}{N} \sum_{n=1}^N \mathbb{1}(\hat{l}_n = L_n)$$

$$SLFR = \frac{1}{N} \sum_{n=1}^N \mathbb{1}(\hat{l}_n \neq l_n)$$

where \hat{l}_n is the label given by GPT-4o discriminator, L_n is the target label, l_n is the original label.

SLFR measures the rate at which the generated counterfactual remove their original label, and LFR evaluates how often the counterfactual examples successfully adopt the target label.

4 Experiments

We evaluate the generated counterfactuals using two experiments*. First, we evaluate the quality of

*We spent approximately 400USD in total on API calls to OpenAI for running Experiments 1 and 2. Since the running of the experiments, the cost of GPT-4o has decreased by 79%.

generated counterfactual examples using the PKP, LFR, and SLFR metrics in § 3.2.

In the second experiment, we compare our proposed approach to other example selection techniques in a standard classification task, using two pre-trained models. We use five different data selection techniques in interactive AL: random selection, cluster-based selection, uncertainty-based selection, counterfactual examples generated without Variation Theory, and our proposed counterfactual based example selection. We use each dataset’s original label as ground truth and use GPT-4o and a BERT model as the target classification models.

To further understand the impact of each component of our filtering pipeline, we conduct an ablation study. In this study, we aim to understand the impact of each individual filter on the pipeline’s performance in downstream model training. Additional details can be found in Appendix B.

4.1 Datasets

- **YELP:** The YELP dataset (Asghar, 2016) consists of user reviews of different businesses and services. The dataset itself provides 4 ground-truth categories (i.e. service, price, environment and products), we randomly sampled 495 examples for this experiment.
- **MASSIVE:** The MASSIVE (FitzGerald et al., 2022) virtual assistant utterances with 18 labeled intents as ground-truth (e.g. audio, cooking, weather, recommendation etc). For this experiment we randomly selected 30 examples from each category, making up a total of 540 examples.
- **Emotions:** Includes a collection of English Twitter messages annotated with 6 emotions: anger, fear, joy, love, sadness, and surprise (Elgiriye-withana, 2024). For this experiment we randomly selected 500 examples while balancing the number of labels.

4.2 Experiment 1: Generated Counterfactual Quality

We evaluate the generated counterfactuals using two experiments. First, we evaluate the quality of generated counterfactual examples using the PKP, LFR, and SLFR metrics in § 3.2.

4.2.1 Results

	YELP	MASSIVE	Emotions
Pattern Keeping Rate	0.94	0.88	0.81
Soft Label Flip Rate	0.45	0.71	0.58
Label Flip Rate	0.98	0.86	0.86

Table 1: Generated counterfactual data quality evaluation.

Our findings indicate that our proposed pipeline maintains the quality of generated counterexamples, as measured by Pattern Keeping Rate (PKR) and Label Flip Rate (LFR). Across datasets, the PKR remains high, demonstrating the generated counterfactual examples effectively keep the original pattern rules. The LLM-based Discriminator Filtering achieves robust performance in LFR across datasets, confirming that most counterfactual examples successfully adopt the target label. However, the Soft Label Flip Rate (SLFR) varies, particularly with the MASSIVE dataset showing the highest rate and the others on the lower side. This suggests that the degree of semantic change required to remove the original label can be dataset-dependent.

4.3 Experiment 2: Generated Counterfactuals in Downstream Model Training

In the second experiment, we compare our counterfactual generation approach with five other sampling strategies in AL.

- **Random** Examples are randomly selected for each annotation iteration to train the classification model.
- **Cluster** Examples selected from a k-means clustered, pretrained Sentence Transformer model by iterating through the clusters in rotation.
- **Uncertainty** We use model confidence on the training set to choose data with the lowest confidence to be labeled. We use verbal uncertainty (Lin et al., 2022) to get model confidence in GPT-4o and model logits for the BERT model.
- **ALPS (Yuan et al., 2020)** We use ALPS a sampling strategy that addresses the cold start problem in AL.
- **Counterexamples without Variation Theory** We generate counterexamples without using the neuro-symbolic pipeline defined in Fig 1.

4.3.1 Protocol

To evaluate the generated counterfactual examples, we employ a simulated active learning task to train and evaluate a BERT model (Devlin et al., 2018) and few-shot prompting GPT-4o model for a multi-class classification task. We use the example selection conditions defined in § 4.3 to define a subset of 10, 15, 30, and progressively increasing up to 170 data points (referred to as ‘shots’), alongside their corresponding ground truths to be used as training sets. We then evaluate the classifier model using a hold-off set of the dataset.

To augment the model’s training with generated counterfactual examples, we pair each original data with its generated counterfactual examples and their assigned target label. This pairing is used to enrich the distribution and quality of the training data, hypothesizing that the inclusion of counterfactuals would enhance the model’s learning and predictive accuracy in early stages of annotation, addressing the cold start problem (Yuan et al., 2020). Similarly, the performance of the model, in this case trained with both original and counterfactual datasets, was again evaluated against the same hold-off set. This comparative analysis aimed to quantify the impact of counterfactual examples on the model’s ability to generalize and make accurate predictions on unseen data in early active learning scenarios.

4.3.2 Results

We present our findings on the efficacy of generated counterfactuals in active learning as defined in § 4.3.1. We report the macro F1-scores for the three datasets across different shots and conditions (Table 2 and Table 3) using two models - few-shot learning with GPT-4o and fine-tuning a BERT model. We use training shots ranging from 10 to 120 shots for GPT-4o to stay within OpenAI’s token limit and 10 to 170 for the BERT model.

We conducted a pair-wise t-test between the counterfactual condition and the other baseline conditions to understand the impact of the proposed approach. The results across the three datasets highlight the strong initial impact that the counterfactual condition has in addressing the cold start problem. We consistently observe a statistically significant advantage of the counterfactual condition in lower shot numbers. As the number of annotated examples increases (50 shots and above in most cases), the difference in average F1-score decreases, suggesting the advantage of the counterfactual condi-

Macro F1-scores (GPT-4o)							
YELP							
Method	10	15	30	50	70	90	120
Random	.38 (±.05)***	.44 (±.06)***	.51 (±.07)***	.61 (±.05)	.65 (±.06)	.69 (±.04) ⁺	.74 (±.04)
Cluster	.41 (±.07)***	.48 (±.04)***	.57 (±.07)	.63 (±.06)	.68 (±.03)*	.69 (±.03) ⁺	.70 (±.02)
Uncertainty	.23 (±.04)***	.21 (±.05)***	.27 (±.06)***	.28 (±.05)***	.29 (±.04)***	.28 (±.06)***	.29 (±.05)
ALPS	.37 (±.04)***	.49 (±.06)*	.66 (±.05)	.68 (±.03)	.69 (±.03)	.70 (±.04)	.72 (±.03)
Counterfactuals without VT	.35 (±.10)***	.46 (±.13)*	.54 (±.05)*	.53 (±.06)*	.39 (±.08)***	.25 (±.05)***	.31 (±.05)
Counterfactuals	.55 (±.08)	.59 (±.07)	.63 (±.07)	.69 (±.07)	.59 (±.10)	.65 (±.05)	.78 (±.04)
MASSIVE							
Method	10	15	30	50	70	90	120
Random	.36 (±.06)***	.40 (±.05)*	.49 (±.12)	.51 (±.11)	.54 (±.10)*	.57 (±.09)***	.61 (±.10)
Cluster	.35 (±.06)***	.40 (±.07)*	.47 (±.08)	.49 (±.08)	.56 (±.12)*	.54 (±.12)*	.55 (±.09)
Uncertainty	.22 (±.08)***	.19 (±.10)***	.18 (±.07)***	.13 (±.06)***	.14 (±.07)***	.19 (±.09)***	.20 (±.10)
ALPS	.12 (±.03)*	.24 (±.08)	.39 (±.02)	.61 (±.03)	.65 (±.08)	.67 (±.07)	.72 (±.04)
Counterfactuals without VT	.26 (±.10)***	.37 (±.07)*	.43 (±.05)*	.40 (±.07)	.34 (±.10)	.27 (±.09)*	.37 (±.08)
Counterfactuals	.48 (±.01)	.52 (±.03)	.59 (±.3)	.63 (±.03)	.64 (±.06)	.66 (±.05)	.79 (±.03)
EMOTIONS							
Method	10	15	30	50	70	90	120
Random	.29 (±.10)	.32 (±.10)	.36 (±.07)***	.39 (±.04)***	.45 (±.04)*	.45 (±.06)	.47 (±.04)
Cluster	.32 (±.04)	.38 (±.04)	.36 (±.08)***	.39 (±.12)***	.42 (±.09)*	.42 (±.08)	.41 (±.05)
Uncertainty	.21 (±.07)***	.19 (±.05)***	.25 (±.05)***	.29 (±.04)***	.28 (±.07)***	.29 (±.06)	.33 (±.05)
ALPS	.23 (±.07)	.26 (±.03)	.34 (±.05)	.36 (±.05)	.39 (±.06)	.40 (±.05)	.44 (±.10)
Counterfactuals without VT	.28 (±.06)	.35 (±.10)	.46 (±.13)	.48 (±.13)	.49 (±.12)	.36 (±.08)	.39 (±.07)
Counterfactuals	.34 (±.08)	.43 (±.10)	.54 (±.10)	.51 (±.05)	.58 (±.10)	.47 (±.03)	.52 (±.05)

Table 2: Macro F1-scores for GPT-4o across three datasets (YELP, MASSIVE, EMOTIONS) with varying annotation shot counts. + indicates p-value<.1, * indicates p-value<.05, ** indicates p-value<.01, and *** shows p-value<.0001 between the condition and the counterfactual condition.

Macro F1-scores (BERT)									
YELP									
Method	10	15	30	50	70	90	120	150	170
Random	.16 (±.06)*	.18 (±.05)***	.26 (±.03)***	.33 (±.04)***	.35 (±.06)***	.45 (±.01)	.45 (±.03)	.48 (±.04)	.51 (±.02)
Cluster	.18 (±.08)***	.19 (±.06)***	.26 (±.07)***	.32 (±.06)***	.34 (±.05) ⁺	.46 (±.03)	.31 (±.08)	.42 (±.1)	.45 (±.1)
Uncertainty	.13 (±.06)	.14 (±.04)	.19 (±.07)	.33 (±.04)	.41 (±.06)	.46 (±.03)	.47 (±.04)	.53 (±.04)	.54 (±.03)
ALPS	.14 (±.05)	.16 (±.06)	.15 (±.06)	.25 (±.08)	.27 (±.08)	.27 (±.08)	.36 (±.11)	.37 (±.11)	.37 (±.10)
Counterfactuals without VT	.20 (±.06)	.16 (±.07)	.25 (±.04)	.29 (±.04)	.38 (±.08)	.45 (±.05)	.49 (±.04)	.54 (±.05)	.55 (±.04)
Counterfactuals	.38 (±.04)	.39 (±.07)	.49 (±.05)	.47 (±.04)	.51 (±.04)	.53 (±.04)	.50 (±.03)	.52 (±.02)	.53 (±.03)
MASSIVE									
Method	10	20	30	50	70	100	130	150	170
Random	.048 (±.03)***	.052 (±.03)***	.12 (±.04)***	.11 (±.05)***	.19 (±.03)***	.22 (±.02)***	.23 (±.02)***	.24 (±.02)***	1 (±.02)
Cluster	.046 (±.01)***	.058 (±.04)***	.091 (±.03)***	.13 (±.04)***	.18 (±.04)***	.20 (±.03)***	.23 (±.02)***	.24 (±.02)***	.25 (±.02)
Uncertainty	.029 (±.02)***	.035 (±.02)***	.11 (±.04)***	.14 (±.03)***	.22 (±.02)***	.23 (±.03)***	.24 (±.03)***	.25 (±.03)***	.25 (±.02)***
ALPS	.017 (±.01)***	.13 (±.01)***	.14 (±.01)***	.19 (±.01)***	.31 (±.01)	.23 (±.01)	.45 (±.02)	.45 (±.02)	.64 (±.05)
Counterfactuals without VT	.09 (±.08)***	.15 (±.07)***	.33 (±.08)***	.50 (±.07)*	.61 (±.05)⁺	.64 (±.04)	.68 (±.04)*	.68 (±.04)	.69 (±.03)⁺
Counterfactuals	.33 (±.09)	.40 (±.07)	.51 (±.08)	.58 (±.06)	.56 (±.05)	.60 (±.09)	.61 (±.06)	.66 (±.05)	.62 (±.1)
EMOTIONS									
Method	10	20	30	50	70	100	130	150	170
Random	.19 (±.04)*	.20 (±.03)***	.24 (±.08)*	.31 (±.12)	.46 (±.09)	.47 (±.09)	.53 (±.14)	.63 (±.07)	.30 (±.06)
Cluster	.18 (±.02)*	.21 (±.03)*	.23 (±.02)***	.28 (±.03)*	.41 (±.05)	.43 (±.08)	.48 (±.06)	.59 (±.05)	.52 (±.12)
Uncertainty	.23 (±.04)***	.23 (±.05)	.26 (±.08)*	.35 (±.05)	.38 (±.04) ⁺	.57 (±.07)***	.66 (±.08)***	.69 (±.07)	.70 (±.06)*
ALPS	.09 (±.04)	.15 (±.04)	.28 (±.04)	.24 (±.04)	.42 (±.04)	.44 (±.03)	.52 (±.03)	.74 (±.03)	.75 (±.03)
Counterfactuals without VT	.18 (±.05)*	.21 (±.05)*	.32 (±.09)	.36 (±.12)	.40 (±.13)	.57 (±.08)***	.62 (±.1)	.62 (±.2)	.72 (±.05)*
Counterfactuals	.27 (±.07)	.26 (±.09)	.36 (±.05)	.38 (±.12)	.49 (±.05)	.45 (±.15)	.50 (±.06)	.63 (±.06)	.56 (±.07)

Table 3: Macro F1-scores for BERT model across three datasets (YELP, MASSIVE, EMOTIONS) with varying annotation shot counts. + indicates p-value<.1, * indicates p-value<.05, ** indicates p-value<.01, and *** shows p-value<.0001 between the condition and the counterfactual condition.

tion diminishes when more data become available. Similarly, we observe significant impacts of the counterfactual condition when using a few-shot approach with the GPT-4o (Table 2). However, we did not find results that consistently indicated a substantial difference between the random, cluster, and counterfactual without variation theory conditions after 50 shots of examples have been labeled. The results demonstrated the **performance advantage** of our proposed neuro-symbolic variation theory-based counterfactual data augmentation approach in cold-start scenarios for active learning tasks.

Our approach introduces useful data to address the lack of label distribution and representation in cold start scenarios. Compared to the *counterfactuals without Variation Theory* condition, the counterexamples generated through Variation Theory have a significantly higher F1-score, showing the impact of the pipeline in generating useful data in early AL. Moreover, the ablation study in Appendix B evaluating the impact of the filtering components in the pipeline shows there is a statistically significant difference in the downstream performance of a model trained on filtered data compared to data that does not have the complete filtering pipeline.

As we get more annotated data, we observe either minimal improvement or a decline in the model’s performance. We believe that this occurs because after a certain point, the generated counterfactuals begin to replicate previously observed patterns, and there is a limit to the amount of information that can be extracted from these patterns. We also see similar patterns of model decline in the *counterfactuals without VT condition*. This ultimately may cause the model to overly rely on itself, resulting in the performance not scaling. To address this, it is important to heuristically understand the amount of data distribution that can be captured by generated data and switch gears back to using real data when needed.

5 Conclusion

Li et al. (2023b) find that the performance of synthetic data is highly dependent on the distribution of the generated data, suggesting that enhancing data diversity could significantly improve the utility of synthetic data. Our approach achieves this by generating counterfactual examples along dynamic neuro-symbolic boundaries to allow the synthetic data to represent underlying concepts for better

generalizability. In our evaluation, we find that models trained on counterfactual examples have a statistically significant advantage in the early stage of active learning, where there is a limited number of annotated data. When there is only a small amount of annotated data available, the distribution of the ‘*real data*’ may not sufficiently cover the latent space.

Notably, the performance benefit of the counterfactual condition begins to decline when more than 70 labeled data points are used in model training. This reduction in advantage could potentially be attributed to model collapse. This happens when the model fails to capture the full diversity of the data on which it is trained (Wang et al., 2023; Su et al., 2023). With the introduced distribution shift, after the 70-shot threshold, the model might overfit to the specific characteristics of the synthetic examples it has seen, rather than generalizing to the broader real data distribution.

In terms of cost, we spent approximately 400USD in total on API calls to OpenAI for running Experiments 1 and 2. Since then, the cost of GPT-4o has decreased by 79%, and we anticipate that the cost of using LLMs will continue to decline as the technology advances. Recent models are already demonstrating state-of-the-art performance at significantly lower costs. To explore the feasibility of more scalable alternatives, we also conducted an experiment using an open-weight model and found the results to be comparable (see Appendix C).

6 Limitations

Our neuro-symbolic pipeline enables the automatic, real-time creation of counterfactual data using a pattern-based program synthesis approach. This method defines the concept space varied during counterfactual generation. Although the current pattern building blocks are designed for general domains, they rely on predefined rules, which may need augmentation with domain-specific lexical rules for specialized applications. Additionally, our use of a GPT-based discriminator to assign target labels for each counterfactual introduces potential biases or limitations inherent to the discriminator model itself. Future work could focus on understanding how human annotators understand and label the generated counterfactual examples.

References

- Nabiha Asghar. 2016. Yelp dataset challenge: Review rating prediction. *arXiv preprint arXiv:1605.05362*.
- Samuel Budd, Emma C Robinson, and Bernhard Kainz. 2021. A survey on active learning and human-in-the-loop deep learning for medical image analysis. *Medical Image Analysis*, 71:102062.
- Zeming Chen, Qiyue Gao, Antoine Bosselut, Ashish Sabharwal, and Kyle Richardson. 2023. Disco: Distilling counterfactuals with large language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5514–5528.
- Wai Lun Eddie Cheng. 2016. [Learning through the variation theory: A case study](#). *The International Journal of Teaching and Learning in Higher Education*, 28:283–292.
- Emily Denton, Ben Hutchinson, Margaret Mitchell, Timnit Gebru, and Andrew Zaldivar. 2020. [Image counterfactual sensitivity analysis for detecting unintended bias](#). *Preprint*, arXiv:1906.06439.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Nidula Elgiriye withana. 2024. [Emotions Dataset](#).
- Jack FitzGerald, Christopher Hench, Charith Peris, Scott Mackie, Kay Rottmann, Ana Sanchez, Aaron Nash, Liam Urbach, Vishesh Kakarala, Richa Singh, Swetha Ranganath, Laurie Crist, Misha Britan, Wouter Leeuwis, Gokhan Tur, and Prem Natarajan. 2022. [Massive: A 1m-example multilingual natural language understanding dataset with 51 typologically-diverse languages](#). *Preprint*, arXiv:2204.08582.
- Yifan Fu, Xingquan Zhu, and Bin Li. 2013. A survey on instance selection for active learning. *Knowledge and information systems*, 35:249–283.
- Simret Araya Gebreegziabher, Zheng Zhang, Xiaohang Tang, Yihao Meng, Elena L. Glassman, and Toby Jia-Jun Li. 2023. [Patat: Human-ai collaborative qualitative coding with explainable interactive rule synthesis](#). In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*, CHI '23, New York, NY, USA. Association for Computing Machinery.
- Abbas Ghaddar, Philippe Langlais, Ahmad Rashid, and Mehdi Rezagholizadeh. 2021. [Context-aware Adversarial Training for Name Regularity Bias in Named Entity Recognition](#). *Transactions of the Association for Computational Linguistics*, 9:586–604.
- Tamara Gog and Nikol Rummel. 2010. [Example-based learning: Integrating cognitive and social-cognitive research perspectives](#). *Educational Psychology Review*, 22:155–174.
- Sumit Gulwani. 2011. Automating string processing in spreadsheets using input-output examples. *ACM Sigplan Notices*, 46(1):317–330.
- Pascal Hitzler and Md Kamruzzaman Sarker. 2022. Neuro-symbolic artificial intelligence: The state of the art.
- Nitish Joshi and He He. 2022. [An investigation of the \(in\)effectiveness of counterfactually augmented data](#). *Preprint*, arXiv:2107.00753.
- Divyansh Kaushik, Eduard Hovy, and Zachary Lipton. 2020. [Learning the difference that makes a difference with counterfactually-augmented data](#). In *International Conference on Learning Representations*.
- Divyansh Kaushik, Amrith Setlur, Eduard Hovy, and Zachary C. Lipton. 2021. [Explaining the efficacy of counterfactually augmented data](#). *Preprint*, arXiv:2010.02114.
- Yongqi Li, Mayi Xu, Xin Miao, Shen Zhou, and Tiejun Qian. 2023a. [Large language models as counterfactual generator: Strengths and weaknesses](#). *Preprint*, arXiv:2305.14791.
- Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, and Ming Yin. 2023b. Synthetic data generation with large language models for text classification: Potential and limitations. *arXiv preprint arXiv:2310.07849*.
- Henry Lieberman. 2001. *Your wish is my command: Programming by example*. Morgan Kaufmann.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. Teaching models to express their uncertainty in words. *arXiv preprint arXiv:2205.14334*.
- Alisa Liu, Swabha Swayamdipta, Noah A. Smith, and Yejin Choi. 2022. [WANLI: Worker and AI collaboration for natural language inference dataset creation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6826–6847, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Ference Marton. 2014. *Necessary conditions of learning*. Routledge.
- Sewon Min, Xinxin Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. [Rethinking the role of demonstrations: What makes in-context learning work?](#) In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Abbavaram Gowtham Reddy, Saketh Bachu, Saloni Dash, Charchit Sharma, Amit Sharma, and Vineeth N Balasubramanian. 2023. [Rethinking counterfactual data augmentation under confounding](#). *Preprint*, arXiv:2305.18183.

- Raphael Schumann and Ines Rehbein. 2019. Active learning via membership query synthesis for semi-supervised sentence classification. In *Proceedings of the 23rd conference on computational natural language learning (CoNLL)*, pages 472–481.
- Burr Settles. 2009. Active learning literature survey.
- Yi Su, Yixin Ji, Juntao Li, Hai Ye, and Min Zhang. 2023. Beware of model collapse! fast and stable test-time adaptation for robust question answering. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.
- Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, and 1 others. 2023. On the robustness of chatgpt: An adversarial and out-of-distribution perspective. *arXiv preprint arXiv:2302.12095*.
- Zhao Wang and Aron Culotta. 2020. [Robustness to spurious correlations in text classification via automatically generated counterfactuals](#). *Preprint*, arXiv:2012.10040.
- Sherry Wu, Hua Shen, Daniel S Weld, Jeffrey Heer, and Marco Tulio Ribeiro. 2023. Scattershot: Interactive in-context example curation for text transformation. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*, pages 353–367.
- Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel S Weld. 2021. Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models. *arXiv preprint arXiv:2101.00288*.
- Linyi Yang, Jiazheng Li, Pádraig Cunningham, Yue Zhang, Barry Smyth, and Ruihai Dong. 2022a. [Exploring the efficacy of automatically generated counterfactuals for sentiment analysis](#). *Preprint*, arXiv:2106.15231.
- Suorong Yang, Weikang Xiao, Mengcheng Zhang, Suhan Guo, Jian Zhao, and Furao Shen. 2022b. [Image data augmentation for deep learning: A survey](#). *Preprint*, arXiv:2204.08610.
- Michelle Yuan, Hsuan-Tien Lin, and Jordan Boyd-Graber. 2020. Cold-start active learning through self-supervised language modeling. *arXiv preprint arXiv:2010.09535*.

A Appendix

A.1 Generation Pipeline

In this section, we provide the details of all the prompts and models we use to construct the whole counterfactual generation pipeline.

A.1.1 GPT-4o Multi-label Separator

As shown in Fig. 2 Step-1, we utilize zero-shot GPT-4 to preprocess the raw data in order to separate the given multi-labeled sentences into several single-labeled parts. We call GPT-4 through the API provided by OpenAI, set the temperature parameter to 0, and restrict the maximum token number to 512, which ensures the reliability of the generated results. The prompt used is shown below:

- {"role": "system", "content": "The assistant will separate the given multi-labeled sentences into different parts, each corresponds to a label and a pattern (if the pattern is viable)"}
- {"role": "user", "content": "The assistant will generate outputs based on the following example. New content should be in the format: 'text' + 'pattern' + 'label'; 'text' + 'pattern' + 'label'. All the text, patterns and labels are already given as input, if there is no corresponding pattern, just use '' to indicate empty."}
- {"role": "user", "content": "Each separated text must only have a single label, but may contain several patterns. Each label or pattern must appear at least once in the completion. The patterns can be composed with AND (+) or OR (|) operators."}

- {"role": "user", "content": "Conversation: Great customer service, reasonable prices, and a chill atmosphere. Pattern: ['(customer)+*+[service]', '(pay)|(sale)', '(environment)'] Label: price, service, environment"}
- {"role": "assistant", "content": "'Great customer service, ' + '(customer)+*+[service]' + 'service'; 'reasonable prices, ' + '(pay)|(sale)' + 'price'; 'and a chill atmosphere.' + '(environment)' + 'environment' "}
- {"role": "user", "content": "Conversation: {text} Pattern: {pattern} Label: {label}"}

A.1.2 GPT-4o Candidate Phrases Generator

As we are generating counterfactuals that keeps neuro-symbolic patterns, the first step of this task is to generate candidate phrases that keep the pattern but variate semantically, which make up crucial branches of generated counterfactual variations. For this part, we call GPT-4o through the API provided by OpenAI, set the temperature parameter to 0 and restrict the maximum token number to 256. The prompt used is shown below:

- {"role": "system", "content": "The assistant will create a list of phrases that match the given domain specific language based on the given definition."}
- {"role": "user", "content": "For the following text and pattern, generate as many diverse example phrases that match the given pattern and can be part of the given target label. Try to not use the word {label} or {target_label} in the phrases you generate. Separated your answer by a comma"}
- {"role": "user", "content": "text: {matched_phrase}, pattern: {pattern}, current label: {label} target label: {target_label}"}
- {"role": "user", "content": "The word '{match}' is a soft match, you can only use {soft-match_words} as its synonyms to replace it. You can not use other words for {match}"}

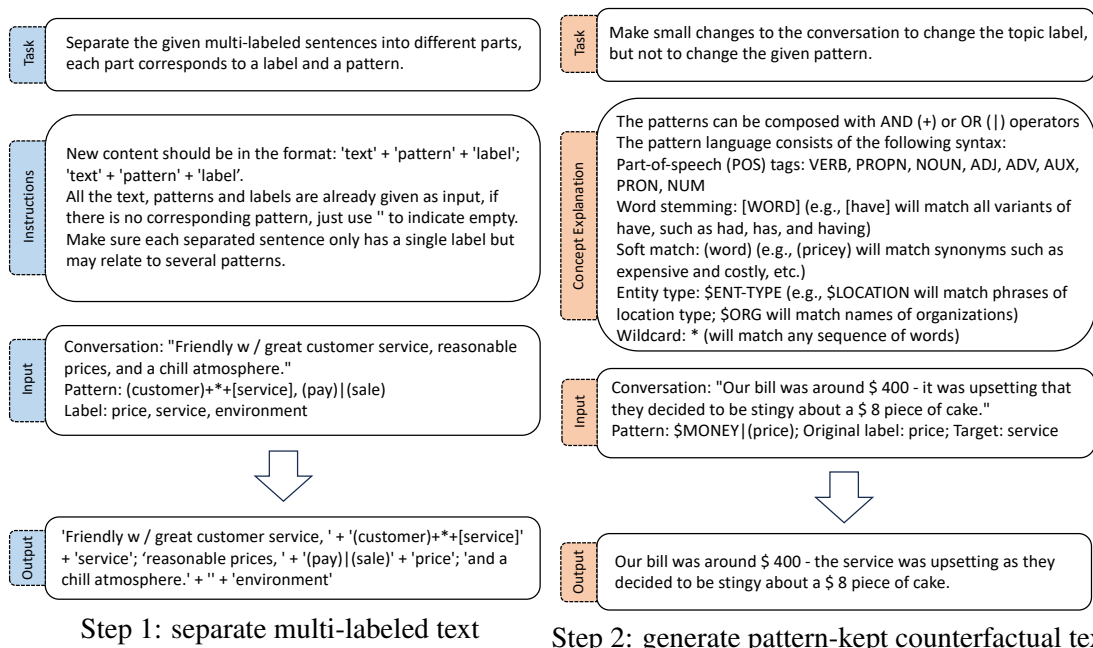


Figure 2: Illustration of LLM prompts used for preparing training datapoints and generating counterfactual datapoints

No of Shots	10	15	30	50	70	90	120
No Filters	0.10	0.12	0.15	0.23	0.23	0.21	0.21
SD	0.03	0.04	0.05	0.04	0.04	0.03	0.03
Herustic Filter	0.15	0.17	0.19	0.28	0.27	0.28	0.28
SD	0.08	0.1	0.1	0.07	0.09	0.1	0.1
Herustic + Symbolic Filters	0.12	0.13	0.13	0.17	0.16	0.18	0.20
SD	0.04	0.03	0.01	0.02	0.03	0.02	0.01
Herustic + LLM Discriminator	0.17	0.21	0.23	0.34	0.42	0.45	0.49
SD	0.08	0.04	0.09	0.07	0.02	0.02	0.05
Herustic + Symbolic + LLM Discriminator	0.38	0.39	0.49	0.47	0.51	0.53	0.50
SD	0.04	0.08	0.06	0.04	0.05	0.05	0.04

Table 4: Average F1-score and SD from an ablation study with the YELP dataset on BERT model

A.1.3 GPT-4o Counterfactual Generator

The GPT-4o generator will finish the second step of counterfactual generation, making use of candidate phrases generated in the first step and combining these semantic pieces into reasonable sentences. We set the temperature parameter to 0 and restrict the maximum token number to 256. The prompt used is shown below:

- {"role": "system", "content": "The assistant will generate a counterfactual example close to the original sentence that contains one of the given phrases."}
- {"role": "user", "content": "Your task is to change the given sentence from the current label to the target."}

For example: 'Find me a train ticket next monday to new york city' with original label "transport" would be turned to 'Play me a song called New York City by Taylor Swift' with a label "audio".

You can use the following phrases to help you generate the counterfactuals. Please make the sentence about {**target_label**}. Make sure that the new sentence is

not about {**label**}. You must use one of the following phrases without rewording it in the new sentence: {**generated_phrases**}"}

- {"role": "user", "content": "You must follow three criteria:
criteria 1: the phrase should change the label from {**label**} to {**target_label**} to the highest degree.
criteria 2: the modified sentence can not also be about {**label**} and make sure the word {**target_label**} is not part of the modified sentence.
criteria 3: the modified sentence should be grammatically correct."}
- {"role": "user", "content": "If you find that you cannot generate new sentence that fulfill all the requirements above, just response 'cannot generate counterfactual' and don't feel bad about this"}
- {"role": "user", "content": "original text:{**text**}, original label:{**label**}, modified label:{**target_label**}, generated phrases:{**generated_phrases**}, modified text: "}

Counterfactual Generated using Llama3.3							
Method	10	15	30	50	70	90	120
YELP	.31 (\pm .05)	.32 (\pm .06)	.34 (\pm .08)	.44 (\pm .05)	.51 (\pm .08)	.53 (\pm .04)	.64 (\pm .04)
MASSIVE	.28 (\pm .08)	.36 (\pm .06)	.40 (\pm .02)	.54 (\pm .07)	.58 (\pm .06)	.60 (\pm .03)	.66 (\pm .03)
EMOTIONS	.21 (\pm .09)	.24 (\pm .05)	.32 (\pm .1)	.34 (\pm .07)	.39 (\pm .1)	.47 (\pm .06)	.51 (\pm .08)

Table 5: Macro F1-scores for the counterfactual condition using Llama3.3 as the counterfactual generator model for BERT, evaluated across three datasets (YELP, MASSIVE, EMOTIONS) at varying annotation shot counts.

B Ablation Study on Counterfactual Filtering Methods

We performed an ablation study to investigate the impact of the different components in our filtering pipeline. We follow the same approach as § 4.3.1 where each condition is run with different seeds 8 times. For each condition, we report an average F1 score and the standard deviation (SD) in Table 4. Our approach involves generating counterexamples with a fine-tuned GPT-4o model and applying all three filters defined in § 3.2 before using the data for active learning.

In this study, we investigate the impact of different configurations by varying the filtering mechanisms used with the generator model.

The ablation study is conducted using the YELP dataset with a BERT model for the downstream active learning tasks. The configurations tested include:

- No Filters: Counterexamples generated without any filters applied
- Heuristic Filter: Applying only the heuristic filter
- Heuristic + Symbolic Filters: Applying both heuristic and symbolic filters
- All Filters: Applying all three filters defined in § 3.2

The results indicate that the use of all filters significantly improves the performance of the trained model (See Table 4). The average F1-score with all filters applied reaches 0.51 for 70 shots and peaks at 0.53 for 90 shots, demonstrating a 2X improvement over the baseline with no filters (F1-score of 0.23 for 70 shots). Using a pairwise t-test, we find that this is statistically significant ($p < 0.0001$), underscoring the value of carefully filtering LLM-generated counterfactuals to produce usable data for model training.

Surprisingly, we found that incorporating the symbolic filter without the LLM discriminator decreases the performance of downstream training. Further analysis of the included examples revealed that some generated sentences included the original sentence with additional parts that corresponded to the target label. While the LLM discriminator would filter these out, without its use in the pipeline, these generated counterfactuals are mistakenly treated as negative examples, when technically they are just multi-labeled positive examples. However, we observe a substantial improvement in performance when the symbolic filter is used in conjunction with the LLM discriminator, as opposed to using the LLM discriminator alone. This demonstrates the effectiveness of combining both methods to enhance the quality and accuracy of the generated counterfactuals.

The ablation study highlights the crucial role of the filtering pipeline. By systematically evaluating the impact of each component, we demonstrate that the integration of heuristic, symbolic filters, and the LLM discriminator leads to significant improvements in the downstream active learning task. This validates our hypothesis that filtering LLM-generated data is essential in determining usable and useful data for achieving higher performance and reliability in model training.

C Experiment using Open-weight Counterfactual Generator

We evaluate the effectiveness of our counterfactual generation approach using an open-weight model Llama3.3 model as the generator with BERT as the classifier model across three datasets (YELP, MASSIVE, EMOTIONS), under increasing annotation shot counts. Our findings as seen in Table 5 show that performance using the Llama3.3 model is comparable across all datasets, showing the viability of our method beyond proprietary models.