# From Selection to Generation: A Survey of LLM-based Active Learning

Yu Xia[1*], Subhojyoti Mukherjee[2*], Zhouhang Xie[1], Junda Wu[1], Xintong Li[1],
Ryan Aponte[3], Hanjia Lyu[4], Joe Barrow[5], Hongjie Chen[6], Franck Dernoncourt[5],
Branislav Kveton[5], Tong Yu[5], Ruiyi Zhang[5], Jiuxiang Gu[5], Nesreen K. Ahmed[7], Yu Wang[8],
Xiang Chen[5], Hanieh Deilamsalehy[5], Sungchul Kim[5], Zhengmian Hu[5], Yue Zhao[9],
Nedim Lipka[5], Seunghyun Yoon[5], Ting-Hao 'Kenneth' Huang[10], Zichao Wang[5],
Puneet Mathur[5], Soumyabrata Pal[5], Koyel Mukherjee[5], Zhehao Zhang[11],
Namyong Park, Thien Huu Nguyen[8], Jiebo Luo[4], Ryan A. Rossi[5], Julian McAuley[1]

[1]University of California San Diego, [2]University of Wisconsin Madison,
[3]Carnegie Mellon University, [4]University of Rochester, [5]Adobe Research, [6]Dolby Labs,
[7]Cisco AI Research, [8]University of Oregon, [9]University of Southern California,
[10]Pennsylvania State University, [11]Dartmouth College

## Abstract

Active Learning (AL) has been a powerful paradigm for improving model efficiency and performance by selecting the most informative data points for labeling and training. In recent active learning frameworks, Large Language Models (LLMs) have been employed not only for selection but also for generating entirely new data instances and providing more cost-effective annotations. Motivated by the increasing importance of high-quality data and efficient model training in the era of LLMs, we present a comprehensive survey on LLM-based Active Learning. We introduce an intuitive taxonomy that categorizes these techniques and discuss the transformative roles LLMs can play in the active learning loop. We further examine the impact of AL on LLM learning paradigms and its applications across various domains. Finally, we identify open challenges and propose future research directions. This survey aims to serve as an up-to-date resource for researchers and practitioners seeking to gain an intuitive understanding of LLM-based AL techniques and deploy them to new applications.

## 1 Introduction

Active Learning (AL) has been a widely studied technique that aims to reduce data annotation efforts by actively selecting most informative data samples for labeling and subsequent model training (Cohn et al., 1994, 1996; Settles, 2009; Olsson, 2009; Fu et al., 2013; Ren et al., 2021; Zhan et al., 2022). With an effective data selection strategy, this process helps to efficiently improve model performance with fewer labeled data instances, which can be particularly valuable when data annotation
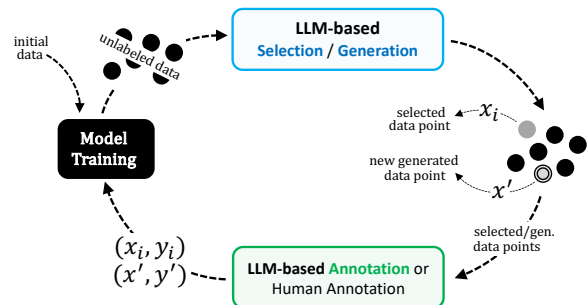
---

*Equal contributions



Figure 1: Overview of LLM-based active learning. We start with initial data, including unlabeled instances $\mathcal{U}$. There are two main steps. First, LLM-based selection and/or generation leverages an LLM $\mathcal{M}$ to select unlabeled instances $\mathbf{x}_i \in \mathcal{U}$ for annotation or *generate entirely new instances* $\mathbf{x}' \notin \mathcal{U}$. Next, given the LLM-based selected or generated instances $\mathbf{x}_i$ or $\mathbf{x}'$, LLM-based annotation uses an LLM $\mathcal{M}$ to generate labels $y_i$ and $y'$ for the instances. Note that for intuition we show how LLMs can be leveraged for both steps; however, we may also use traditional techniques for selecting unlabeled instances, use humans for annotation, or both.

is expensive or time-consuming (Aggarwal et al., 2014; Hino, 2020; Schröder and Niekler, 2020).

Despite the success of traditional active learning methods, the advent of Large Language Models (LLMs) with remarkable reasoning and generation capabilities creates a new paradigm of active learning. For example, as illustrated in Figure 1, instead of solely relying on a predefined metric to query data instances, such as uncertainty (Wang and Shang, 2014; Diao et al., 2023) or diversity (Agarwal et al., 2020; Citovsky et al., 2021), LLMs can now be used to select most informative instances after reasoning or even generate entirely new instances that are better suited for efficient model training (Bayer and Reuter, 2024; Parkar et al., 2024; Yang et al., 2024; Bhatt et al., 2024; Zhang and Nowak, 2024). Moreover, with the col-

| Class | General Mechanism | Description |
|---|---|---|
| **Querying** (Section 3) | Traditional Selection (Sec. 3.1) | This class of techniques uses traditional selection such as uncertainty sampling, disagreement, gradient-based sampling, and so on. |
| | LLM-based Selection (Sec. 3.2) | The class of LLM-based selection techniques focus on using LLMs to select the instances. |
| | LLM-based Generation (Sec. 3.3) | The class of LLM-based generation techniques focus on generating novel instances. |
| | Hybrid (Sec. 3.4) | Combines advantages of both LLM-based selection and generation |
| **Annotation** (Section 4) | Human Annotation (Sec. 4.1) | Traditional human annotation simply refers to using humans to annotate the selected or generated instances, which is costly. |
| | LLM-based Annotation (Sec. 4.2) | The class of LLM-based annotation techniques focus on leveraging LLMs for annotation and evaluation. This class of techniques are far cheaper than human annotation. |
| | Hybrid (Sec. 4.3) | This class of techniques aim to leverage the advantages of both humans and LLMs for optimal annotations while minimizing cost |

Table 1: Taxonomy of LLM-based Active Learning Techniques (Sections 3 and 4).

lected informative data instances, LLMs also enable new data annotation schemes by collaborating with a human labeler or directly simulating a human labeler (Xiao et al., 2023; Kholodna et al., 2024; Wang et al., 2024), which further reduces manual annotation efforts. LLM-based selection or generation can also help reduce training costs such as for supervised fine-tuning, in contrast to the main focus of traditional AL on reducing the labeling costs.

However, in spite of the immense potential of LLMs for active learning, particularly in high-quality data acquisition and annotation for efficient model training, existing surveys primarily focus on traditional active learning techniques, necessitating an up-to-date review of how LLMs have advanced AL in recent years. In this paper, we address this gap by presenting the first comprehensive survey of LLM-based AL techniques, which introduces a unifying taxonomy centered on the two main components of active learning: *Querying* (selecting or generating unlabeled instances) and *Annotation* (assigning labels). Table 1 and Figure 2 provide an overview of the proposed taxonomy. Table 2 further provide an intuitive comparison of existing LLM-based AL methods from the aspects of taxonomy and applications. Guided by this taxonomy, our survey organizes and systematically reviews recent works across key aspects of LLM-based active learning as follows.

- **Preliminaries** (§2): We begin by introducing and formulating LLM-based active learning.
- **Querying** (§3): We describe querying strategies, including LLM-based selection and generation.
- **Annotation** (§4): We detail various annotation

schemes, ranging from human annotation to LLM-based and hybrid approaches.
- **Stopping** (§5): We discuss how recent works consider LLM costs for stopping the AL loop.
- **Active Learning Paradigms** (§6): We examine how AL influences LLMs' learning paradigms.
- **Applications** (§7): We highlight the diverse applications of LLM-based active learning.
- **Open Problems** (§8): We discuss open problems and present future research directions.

**Survey Scope** This survey focuses mainly on recent works leveraging LLMs for AL, which creates a new paradigm driven by LLMs' reasoning and generation capabilities. While some works use traditional AL for LLMs, we may not cover them thoroughly as they use techniques reviewed by prior surveys (Zhan et al., 2022; Perlitz et al., 2023).

## 2 What is LLM-based Active Learning?

We start with basic notations and objective of traditional active learning and then introduce the LLM-based active learning loop with five main steps.

**Traditional Active Learning** Let $\mathcal{U} = \{\mathbf{x}_i\}_{i=1}^N$ be a pool of $N$ unlabeled instances, where $\mathbf{x}_i \in \mathcal{X}$ are feature vectors in the input space $\mathcal{X}$. Furthermore, let $\mathcal{L} = \{(\mathbf{x}_i, y_i)\}_{i=1}^M$ be a labeled dataset, where $y_i \in \mathcal{Y}$ are the corresponding labels from the label space $\mathcal{Y}$ and the size of the labeled set $|\mathcal{L}| = M$ grows as more data instances are labeled. We are also given an annotation budget $k$ that limits the number of instances that can be labeled by the annotator. We have a target model $f_\theta : \mathcal{X} \to \mathcal{Y}$ parameterized by $\theta$ to be iteratively trained using

the labeled data $\mathcal{L}$. Note that the target model $f_\theta$ can be any parameterized predictive model, regardless of its architecture. The objective of traditional AL is to efficiently select and label a subset of the unlabeled data instances $\mathbf{x}_i \in \mathcal{U}$ to maximize the performance of model $f_\theta$ before reaching the annotation budget $k$.

**LLM-based Active Learning**  LLM-based AL shares a similar goal. However, it is not bounded by the unlabeled set $\mathcal{U}$, but can also use an LLM $\mathcal{M}$ to generate entirely new data instances denoted as $\mathbf{x}' \notin \mathcal{U}$ as well as generating corresponding labels $y'$ by simulating a human annotator. We define LLM $\mathcal{M}$ here as a decoder, encoder, or encoder-decoder architecture trained on a corpus of hundreds of millions to trillions of tokens following Gallegos et al. (2024). We formulate in details the LLM-based AL loop as follows.

- **Initialize**: For a good starting point of the active learning loop, an LLM $\mathcal{M}$ can be used to annotate an initial set of labels or generate an initial dataset $\mathcal{L}_{\text{init}}$ to warm up the target model $f_\theta$. This approach overcomes the cold start problem that traditional AL methods face when there is no labeled data instance available and the initial model $f_\theta$ does not offer sufficient information for selecting informative data instances, especially when $f_\theta$ is not a pre-trained model.

- **Query**: With an initialized model $f_\theta$, the **Querying** (§3) module is implemented to acquire the most informative data instances. Extending traditional AL methods that only select from the unlabeled set $\mathcal{U}$ with certain uncertainty or diversity metrics, the LLM $\mathcal{M}$ can now be used to select instances $\mathbf{x}_i \in \mathcal{U}$ either by scoring or directly choosing, augment existing examples with generated paraphrases or contrast examples, or even generate entirely new instances $\mathbf{x}' \notin \mathcal{U}$.

- **Annotate**: The acquired data instances $\mathbf{x}_i$ or $\mathbf{x}'$ are then passed to the **Annotation** (§4) module to obtain corresponding labels $y_i$ or $y'$. The LLM $\mathcal{M}$ can be used either in conjunction with, to augment, or even to entirely replace the human labeler. Data instances $\mathbf{x}_i \in \mathcal{U}$ that are selected from the unlabeled set are then excluded from $\mathcal{U}$. All labeled instances $(\mathbf{x}_i, y_i)$ or $(\mathbf{x}', y')$ are then added to the labeled set $\mathcal{L}$.

- **Train**: With newly annotated data instances added to the labeled dataset $\mathcal{L}$, the target model

---

**Algorithm 1** LLM-based Active Learning
**Input:** Unlabeled dataset $\mathcal{U}$, LLM $\mathcal{M}$, Annotation budget $k$
**Output:** Trained model $f_\theta$, Labeled dataset $\mathcal{L}$

1: $\mathcal{L}_{\text{init}}, \mathcal{U} \leftarrow$ Initialize($\mathcal{U}, \mathcal{M}$)
2: $f_\theta \leftarrow$ Train($\mathcal{L}_{\text{init}}$)
3: **while not** Stop($k, f_\theta, \mathcal{M}$) **do**      ▷ **Stopping** (§5)
4:     $\mathbf{x} \leftarrow$ Query($f_\theta, \mathcal{U}, \mathcal{M}$)      ▷ **Querying** (§3)
5:     $(\mathbf{x}, y) \leftarrow$ Annotate($\mathbf{x}, \mathcal{M}$)  ▷ **Annotation** (§4)
6:     **if** $\mathbf{x} \in \mathcal{U}$ **then** $\mathcal{U} \leftarrow \mathcal{U} \setminus \{\mathbf{x}\}$
7:     $\mathcal{L} \leftarrow \mathcal{L} \cup \{(\mathbf{x}, y)\}$
8:     $f_\theta \leftarrow$ Train($\mathcal{L}$)
9: **return** $f_{\theta^*}, \mathcal{L}$

---

$f_\theta$ is trained with a step update on its parameters $\theta$. The updated model $f_\theta$ is then used to provide information for the querying module in the next iteration before a stopping criterion is met.

- **Stop**: The active learning loop reaches **Stopping** (§5) when a fixed annotation budget $k$ is reached, or when some property of the model $f_\theta$ being optimized, such as convergence, is satisfied. With LLMs, the budget can be not only based on the cost of human annotators, but also the cost of prompting the LLM or a combination of both.

We also summarize the LLM-based active learning loop in Algorithm 1. The goal of LLM-based AL in the general setting is to iteratively select or generate with an LLM most informative instances $\mathbf{x}_i$ or $\mathbf{x}'$ for human or LLM labeling and subsequent model training. We discuss other variants of LLM-based AL goals in Appendix A.

## 3 Query: From Selection to Generation

Active learning fundamentally seeks to maximize model performance with minimal annotation cost by carefully selecting most informative examples. Traditionally, this process has relied on uncertainty-based and diversity-based metrics (Settles, 2011; Wang and Shang, 2014; Geifman and El-Yaniv, 2017; Citovsky et al., 2021). With the advent of LLMs, however, the paradigm is evolving from merely selecting examples from a fixed unlabeled pool to also generating new, high-value queries on demand, which can be especially helpful in reducing the training cost for supervised fine-tuning.

### 3.1 Traditional Selection

In this section, we briefly survey some of the key existing active learning query selection strategies. To effectively select informative examples, traditional methods capture how uncertain the model
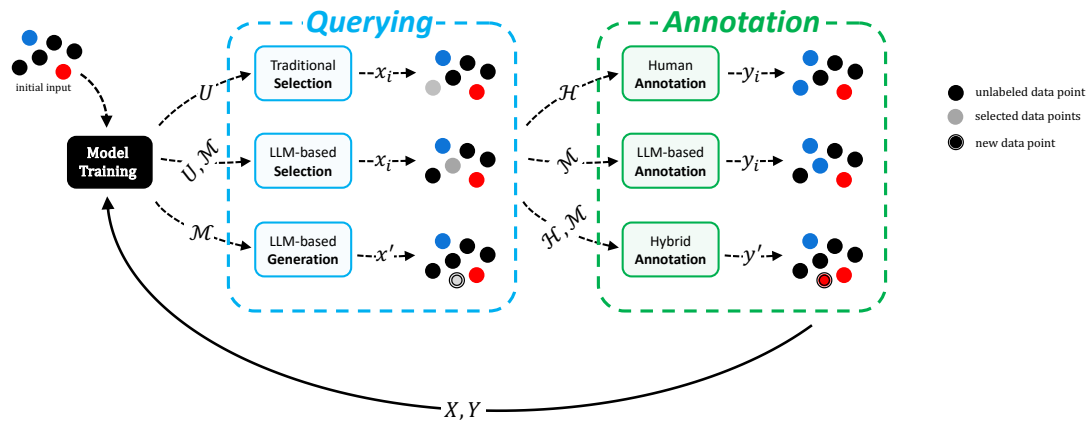
Figure 2: Our proposed taxonomy classifies LLM-based active learning (AL) methods by their querying and annotation processes, the two key components of AL. Beyond selection, LLMs enable the querying module to also generate unlabeled instances, while the annotation module assigns labels using LLMs, human annotators, or both.

is in predicting the label of the example, i.e., uncertainty, or how different the chosen example is from the already selected examples in the labeled pool, i.e., diversity. Uncertainty-based methods, such as Least Confidence (Settles, 2009, 2011; Wang and Shang, 2014), Margin Sampling (Tong and Koller, 2001; Balcan et al., 2009; Settles, 2009), and Max-Entropy (Wang and Shang, 2014; Kremer et al., 2014; Diao et al., 2023), quantify how unsure a model is about its predictions. Complementary to these, diversity-based strategies—like CoreSet (Sener and Savarese, 2017; Geifman and El-Yaniv, 2017; Citovsky et al., 2021) and CDAL (Agarwal et al., 2020)—ensure that selected examples cover varied regions of the input space. Optimal designs (Pukelsheim, 2006) are also a form of diversity based on information gain. Hybrid approaches, including BADGE (Ash et al., 2019, 2021) and BALD (Kirsch et al., 2019; Gal et al., 2017), strike a balance between these two aspects. For a comprehensive discussion of these traditional strategies, see Appendix B and recent surveys in Ren et al. (2021); Li et al. (2024a).

## 3.2 LLM-based Selection

With the emergence of LLMs, active learning strategies are being reimagined to exploit their powerful in-context reasoning and few-shot capabilities. In LLM-based selection, the model itself plays a dual role—both as a predictor and as a selector of informative queries. For instance, ActiveLLM (Bayer and Reuter, 2024) leverages an LLM to assess uncertainty and diversity in a completely unsupervised manner, making it particularly suitable in few-shot and model mismatch settings.

Similarly, ActivePrune (Azeemi et al., 2024) applies an LLM-driven approach to prune large unlabeled pools, reducing the computational burden of traditional acquisition functions for tasks such as translation, sentiment analysis, topic classification, and summarization. In another line of work, SelectLLM (Parkar et al., 2024) prompts LLMs directly to evaluate and rank the usefulness of unlabeled examples; the ranked instances are then refined using k-NN clustering to form effective few-shot demonstrations. Ask-LLM (Sachdeva et al., 2024) also directly prompts an instruction-tuned LLM to assess the quality of a training example. Recent work by Jeong et al. (2024) further demonstrates that LLMs can generate meaningful rankings of examples, which in turn can inform fine-tuning for downstream tasks.

## 3.3 LLM-based Generation

Beyond selection from a fixed unlabeled set, LLMs also enable the generation of entirely new examples and labels, thereby extending the AL paradigm to an effectively infinite search space. Such extension is fundamentally different from an earlier concept (Angluin, 1988) that only queries memberships for hypothetical instances. In the following, we distinguish between generation strategies that remain within the confines of an existing unlabeled pool and those that extend beyond it.

**Generation for Selection within Unlabeled Set**
Several recent works integrate traditional selection metrics with LLM capabilities to improve query selection for few-shot learning. For example, Margatina et al. (2023) employ a combination of k-NN

and perplexity-based strategies, demonstrating that uncertainty sampling generally underperforms in few-shot in-context learning settings with some evidence that this may change with larger models. In a similar vein, Mukherjee et al. (2024a) highlight the effectiveness of experimental design techniques, such as G-Optimal design, for selecting high-impact unlabeled examples. Additionally, Diao et al. (2023) harness LLMs to generate multiple answers to a given question, using the variability among these answers as a proxy for uncertainty—albeit without venturing outside the initial dataset. Meanwhile, EAGLE (Bansal and Sharma, 2023) first samples examples based on a conditional informativeness criterion and then leverages LLMs to generate in-context labels, streamlining the annotation process.

**Generation for Selection outside Unlabeled Set** A more radical departure from traditional AL involves generating new examples that are not present in the original unlabeled pool. The APE framework (Qian et al., 2024) uses a Query-by-Committee strategy combined with chain-of-thought prompting to synthesize new prompts that are then sent to human annotators. Other works, such as those by Yang et al. (2024) and Mukherjee et al. (2024b), generate both new examples and their labels using a trained LLM. These generated examples are then subjected to a rejection sampling process, ensuring that only those meeting predefined accuracy thresholds are retained. Similarly, Yao et al. (2023) use an explanation-generation model to produce human-guided rationales, which not only enhance label quality but also inform a novel diversity-based selection strategy akin to coreset sampling.

### 3.4 Hybrid

Recognizing that neither pure selection nor generation can fully address all challenges in AL, recent work has begun to explore hybrid strategies that combine both. For instance, NoiseAL (Yuan et al., 2024) employs a two-stage process: small LLMs first identify promising unlabeled examples, which are then passed to an annotator LLM for labeling. Similarly, the Causal-guided Active Learning (CAL) framework (Du et al., 2024) integrates density-based clustering with LLM-driven query selection (e.g., via GPT-4) to autonomously identify and correct bias patterns in unlabeled data. Such hybrid methods seek to leverage the complementary strengths of LLMs and traditional human-in-the-loop strategies, ultimately pushing the boundaries of active learning toward more efficient and robust systems.

## 4 Annotation: From Human to LLMs

Annotation has traditionally relied on human experts for high-quality labels. Recently, leveraging LLMs as annotators has further reduced annotation expenses, though challenges such as bias and label inconsistency. A hybrid approach that integrates human expertise with LLM-based annotation offers a promising solution, balancing efficiency and accuracy through dynamic task routing, verification mechanisms, and prompt engineering strategies.

### 4.1 Human Annotation

Traditional human annotation involves sending selected or generated instances to annotators for labeling, which remains the most accurate approach but is often costly. Several recent works have explored active learning strategies to optimize the annotation process. ActivePrune (Azeemi et al., 2024) and CAL (Du et al., 2024) reduce annotation costs by actively selecting the most valuable instances for labeling. For instance, in the case of imbalanced data, enhancing the model's performance on the minority class can not only improve overall accuracy but also mitigate biases. One effective approach to achieving this is by increasing the size of the minority class. Providing human annotators with data that includes more minority samples has been shown to be effective (Lyu et al., 2022). Active-Prompt (Diao et al., 2023) and AL-Principle (Margatina et al., 2023) focus on guiding LLMs by selecting instances where annotators verify final answers or prompt examples before model predictions. Furthermore, Beyond-Labels (Yao et al., 2023) extends traditional annotation by collecting short rationales or natural language explanations alongside labels, improving both model interpretability and performance. Additionally, APL (Muldrew et al., 2024) and BAL-PM (Melo et al., 2024) incorporate human preference learning by asking annotators to compare or rank multiple model outputs. However, there are still several challenges in human annotation. Human annotator variability, as differences in expertise, cognitive biases, and annotation consistency can lead to label noise and disagreements, ultimately affecting model performance. Moreover, bias and fairness considerations in active learning

| | Querying (Section 3) | | | | Annotation (Section 4) | | | Applications (Section 7) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Traditional Selection (§3.1) | LLM-based Selection (§3.2) | LLM-based Generation (§3.3) | Hybrid (§3.4) | Human Annotation (§4.1) | LLM-based Annotation (§4.2) | Hybrid (§4.3) | Text Classification | Text Summarization | Classification | Question Answering | Entity Matching | Debiasing | Translation | Sentiment Analysis | Other |
| **APE** (Qian et al., 2024) | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| **LDCAL** (Li et al., 2024b) | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **ActiveLLM** (Bayer and Reuter, 2024) | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **ActivePrune** (Azeemi et al., 2024) | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ |
| **AutoLabel** (Ming et al., 2024) | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| **LLMaAA** (Zhang et al., 2023) | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| **Active-Prompt** (Diao et al., 2023) | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **HybridAL** (Rouzegar and Makrehchi, 2024) | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **NoiseAL** (Yuan et al., 2024) | ✗ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| **CAL** (Du et al., 2024) | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| **APL** (Muldrew et al., 2024) | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| **AL-Loop** (Kholodna et al., 2024) | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| **BAL-PM** (Melo et al., 2024) | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **FreeAL** (Xiao et al., 2023) | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ |
| **AL-Principle** (Margatina et al., 2023) | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ |
| **Beyond-Labels** (Yao et al., 2023) | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |

Table 2: Overview of the proposed taxonomy for LLM-based active learning techniques and their applications. Using this taxonomy, we provide a qualitative and quantitative comparison of LLM-based active learning methods.

are also an active area of research, as human annotations can inadvertently reflect societal biases, leading to skewed model predictions. There are also works on active learning in RLHF and DPO, which we refer to Section C for more details.

## 4.2 LLM-based Annotation

To address the challenges in human annotation, ongoing research explores leveraging LLMs as annotators in active learning, primarily to reduce annotation costs. FreeAL (Xiao et al., 2023), which distills task-specific knowledge with the help of a downsteam small language model, demonstrates improved performance without any human supervision. Similarly, by employing GPT-4-Turbo for annotating low-resource languages, Kholodna et al. (2024) reported substantial reductions in estimated annotation costs compared to human annotation. However, a key challenge in LLM-based annotation is ensuring high-quality labels. To mitigate quality issues, LLMaAA (Zhang et al., 2023) incorporates

in-context examples, demonstrating improved annotation reliability. Despite such advancements, LLMs, like human annotators, are susceptible to biases inherited from their training data. Research has shown that LLMs' responses to cognitive psychological tasks often resemble those of individuals from western, educated, industrialized, rich, and democratic societies (Atari et al., 2023). Biases in LLM-based annotations also persist in certain domains, such as political science (Zhang et al., 2024b). For example, Qi et al. (2024) identified three dimensions of bias in LLM-generated political samples: societal and cultural contexts, demographic groups, and political institutions. A particular concerning case arises when LLMs are used to annotate content that has also been generated by LLMs. In such scenarios, there is a risk of *self-reinforcement bias*, where the model's inherent tendencies are amplified rather than corrected. This could create a feedback loop where the model's performance may appear artificially

inflated. LLM-based annotations can also be sensitive to input variations (Mizrahi et al., 2024). Slight changes in prompt phrasing, context, or model sampling parameters can lead to inconsistent annotations (Zhu et al., 2024). Ensuring consistency in LLM-generated annotations remains an ongoing challenge, requiring further research into prompt engineering, calibration techniques, and hybrid human-LLM validation strategies.

### 4.3 Hybrid

While LLM-based annotation has significantly reduced costs, it remains prone to errors, particularly in complex or domain-specific tasks (Lu et al., 2023). To address this issue, researchers have developed methods to evaluate annotation quality and dynamically route data to either LLMs or human annotators, balancing efficiency and accuracy. For example, Wang et al. (2024) proposed a multi-step human-LLM collaborative approach where LLMs first generate labels and explanations. A verifier then assesses the quality of these labels, and human annotators re-annotate a subset of low-quality cases. Similarly, Rouzegar and Makrehchi (2024) investigates using LLMs with human annotations for text classification to achieve lower costs while maintaining accuracy based on confidence thresholds. Another approach to combining human expertise with LLMs involves having humans curate a set of annotated examples, which are then incorporated into LLM prompts to enable annotation in a few-shot learning manner. While this method is both intuitive and effective, selecting optimal examples remains challenging. Qiu et al. (2024) show that examples included in prompts can sometimes overly constrain LLM decision-making, leading models to favor labels that align closely with provided examples rather than considering a broader range of possibilities. Refining strategies for selecting and structuring prompt examples is therefore an important direction for improving LLM-assisted annotation.

## 5 Stopping: From Criterion to LLMs

Stopping criterion in active learning loop is crucial for balancing model performance improvements with annotation costs. In LLM-based AL, stopping criteria must consider not only traditional factors such as annotation budget $k$, model performance gains or uncertainty reduction, but also the variable costs associated with LLMs.

### 5.1 Traditional Approaches

In traditional active learning, one widely adopted strategy is to stop querying when performance improvements on a validation set fall below a predefined threshold (Settles, 2009). Other approaches monitor the stability of the model's predictions or the uncertainty estimates, terminating annotation once these metrics indicate that the model has sufficiently converged (Tong, 2001). In addition, several theoretical frameworks provide sample complexity bounds as termination criteria, ensuring that marginal returns are decreasing (Zhu et al., 2003; Hanneke, 2007; Dasgupta et al., 2009). However, such criteria assume that the cost per annotation is uniform and homogeneous, which simplifies the stopping rule.

### 5.2 Cost-Aware Termination

In LLM-based AL, estimating the annotation cost can be challenging and complex. While traditional AL relies on a discrete budget $k$ to represent the number of human annotations, LLM-based AL may combine both human and LLM annotations. In such cases, the cost of an annotation depends not only on the source (human vs. LLM) but also on variable factors like input and output token counts. Even in the case that LLM-based annotations are used without any human annotations, the cost cannot be easily approximated as the discrete budget that often represents the amount of examples that can be labeled. It is straightforward to see that since the cost depends on the input and more so the output tokens, then the budget may be better represented as a real-valued amount that pertains directly to a monetary cost. Recent exploration in hybrid stopping criteria (Akins et al., 2024; Pullar-Strecker et al., 2024) including cost-aware termination criterion that integrates token-level cost analysis and performance plateau detection. Other works have suggested a combined cost-performance metric that balances the fixed costs of human annotations with the variable costs of LLM-based annotations (Xia et al., 2024a; Zhang et al., 2024a).

## 6 AL Paradigms with LLMs

With the rise of LLMs, AL has evolved to address new challenges and opportunities across various learning paradigms. In this section, we briefly outline four LLM-based AL paradigms, where more details can be found at Appendix C. We also take a step further and briefly discuss in Appendix D

14558

| Task | Description |
|---|---|
| **Text Classification** (Sec.E.1) | Selects uncertain or ambiguous texts for classification. |
| **Text Summarization** (Sec.E.2) | Chooses diverse or complex document types to improve summarization. |
| **Non-Text Classification** (Sec.E.3) | The pairing of non-text samples, such as images, with labels. |
| **Question Answering** (Sec.E.4) | Chooses ambiguous or difficult questions for QA systems to refine. |
| **Entity Modeling** (Sec.E.5) | A binary classification task to pair entities with one of two labels. |
| **Debiasing** (Sec.E.6) | The process of reducing measured bias in machine-generated output. |
| **Translation** (Sec.E.7) | Selects sentences or phrases with high uncertainty in translation. |
| **Sentiment Analysis** (Sec.E.8) | Focuses on sentences where the model is uncertain about sentiment. |

Table 3: Applications of LLM-based Active Learning.

a unifying view of the LLM-based AL problem through the lens of bandits, online learning, RLHF, pandoras box, among others.

**Active In-Context Learning**   Recent studies frame few-shot demonstration selection as an active learning problem, leveraging semantic coverage and ambiguity-driven sampling to optimize LLM performance (Margatina et al., 2023; Mavromatis et al., 2024; Qian et al., 2024).

**Active Supervised Fine-Tuning**   To reduce labeling costs, active learning has been integrated into supervised fine-tuning via uncertainty-based querying, self-training on low-uncertainty data, and strategic sample selection (Yu et al., 2022; Xia et al., 2024b; Bayer and Reuter, 2024).

**Active Preference Alignment**   Efficient label selection is critical in reinforcement learning from human feedback (RLHF). Recent approaches employ targeted preference queries to accelerate alignment and improve data efficiency (Ji et al., 2024; Muldrew et al., 2024; Chen et al., 2024).

**Active Knowledge Distillation**   Selective knowledge transfer from LLMs to smaller models reduces computational costs. Methods using uncertainty-based sample selection and iterative student feedback improve distillation efficiency while maintaining performance (Zhang et al., 2024c; Liu et al., 2024; Palo et al., 2024).

## 7   Applications

LLM-based active learning (AL) has been applied across diverse tasks, reducing annotation costs and improving performance in data-scarce settings. We summarize in Table 3 key applications and specific use cases. Table 2 also bridges these applications with our taxonomy of techniques, which provides an intuitive comparison of the state-of-the-art LLM-based AL methods. These applications include

text classification (Rouzegar and Makrehchi, 2024), text summarization (Li et al., 2024b), non-text classification (Margatina et al., 2023), question answering (Diao et al., 2023), entity matching (Qian et al., 2024), debiasing (Du et al., 2024), translation (Kholodna et al., 2024), and sentiment analysis (Xiao et al., 2023). Beyond these, AL has also been used for optimizing system design (Taneja and Goel, 2024) and question generation (Piriyakulkij et al., 2023). For a more detailed discussion of the AL applications on each task, please refer to Appendix E, where we also include a pairing of active learning applications and datasets in Table 5. In Appendix F, we also discuss in further detail the benefits of applying LLM-based AL.

## 8   Open Problems & Challenges

In this section, we discuss open problems and challenges of LLM-based AL for future works.

**Heterogeneous Annotation Costs**   Traditional AL assumes a fixed annotation budget, but LLM-based AL introduces complex cost structures, including human labeling, LLM query costs, and annotation expenses. Future work should develop algorithms that optimize selection strategies while accounting for these heterogeneous costs.

**Multi-LLM AL Algorithms**   Different LLMs present trade-offs in performance and cost, suggesting opportunities for hybrid approaches that combine weak-cheap and strong-expensive models. Inspired by retrieval-augmented models (Huang et al., 2024) and multi-oracle clustering (Silwal et al., 2023), optimizing multi-LLM frameworks for AL remains an open challenge.

**LLM Agents and Active Learning**   Integrating AL into LLM agents presents new possibilities, such as improving retrieval-augmented generation (RAG) (Xu et al., 2024) and in-context example selection (Mukherjee et al., 2024a). Conversely,

exploring LLM agents as annotation tools could enhance existing AL pipelines by reducing human labeling costs.

**Multimodal Active Learning** Most LLM-based AL applications focus on text, leaving open questions on extending these methods to images, audio, and behavioral data. Future work should investigate how well LLMs generalize to non-text domains and whether they can match human-level annotation quality in these settings.

**Complex Feedback and Multi-Aspect Labels** Traditional AL usually relies on single-label supervision, but reward models in the era of LLMs are able to return structured feedback that scores multiple facets of a response (Bai et al., 2022), such as factual accuracy, completeness, and fluency. Developing LLM-based AL strategies that can leverage these rich signals without overfitting to any one dimension may be a promising solution.

**Unstable Performance in LLM-based Annotation** While LLMs are increasingly used to simulate human annotations (Zheng et al., 2023b; Lambert et al., 2024), their reliability varies across domains (Tan et al., 2024). Future research should explore adaptive routing mechanisms that dynamically allocate annotation tasks between LLMs and human annotators based on model competency.

## 9 Conclusion

In this survey, we present an intuitive taxonomy of LLM-based Active Learning, detailing how LLMs can act as sample selectors, data generators, and annotators within the AL loop. We show how these techniques are reshaping traditional AL paradigms and enabling more efficient data acquisition and model training across various applications. By reviewing existing methods, highlighting emerging trends and discussing open challenges, we aim to offer a useful foundation for researchers and practitioners looking to incorporate LLM-based AL techniques into their applications.

## Limitations

LLM-based Active Learning (AL) gives rise to many important applications with critical advantages over traditional AL techniques. Despite the fundamental importance of LLM-based AL, there remain some challenges. The reliance on high-quality initial labeled data may introduce biases, while the computational overhead, cost, and robustness of iterative query generation and selection may limit its application in practice. Query generation and selection techniques via LLMs remain sensitive to model uncertainty, often lacking theoretical guarantees and may lead to inconsistent performance. Furthermore, such techniques may also suffer from reproducibility and robustness issues. Ethical risks such as bias amplification through generation of examples and labels outside known set remains important to handle when deployed in practice.

## References

Sharat Agarwal, Himanshu Arora, Saket Anand, and Chetan Arora. 2020. Contextual diversity for active learning. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVI 16*, pages 137–153. Springer.

Charu C Aggarwal, Xiangnan Kong, Quanquan Gu, Jiawei Han, and S Yu Philip. 2014. Active learning: A survey. In *Data Classification*, pages 599–634. Chapman and Hall/CRC.

Sapphira Akins, Hans Mertens, and Frances Zhu. 2024. Cost-aware query policies in active learning for efficient autonomous robotic exploration. *arXiv preprint arXiv:2411.00137*.

Dana Angluin. 1988. Queries and concept learning. *Machine learning*, 2:319–342.

David Arthur and Sergei Vassilvitskii. 2006. k-means++: The advantages of careful seeding. Technical report, Stanford.

Jordan Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham Kakade. 2021. Gone fishing: Neural active learning with fisher embeddings. *Advances in Neural Information Processing Systems*, 34.

Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. 2019. Deep batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint arXiv:1906.03671*.

Nicolás Astorga, Tennison Liu, Nabeel Seedat, and Mihaela van der Schaar. 2024. Partially observable cost-aware active-learning with large language models. In *The Thirty-Eighth Annual Conference on Neural Information Processing Systems*.

Mohammad Atari, Mona J Xue, Peter S Park, Damián Blasi, and Joseph Henrich. 2023. Which humans?

Abdul Hameed Azeemi, Ihsan Ayyub Qazi, and Agha Ali Raza. 2024. Language model-driven data pruning enables efficient active learning. *arXiv preprint arXiv:2410.04275*.

Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. 2022. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*.

Maria-Florina Balcan, Alina Beygelzimer, and John Langford. 2009. Agnostic active learning. *Journal of Computer and System Sciences*, 75(1):78–89.

Parikshit Bansal and Amit Sharma. 2023. Large language models as annotators: Enhancing generalization of nlp models at minimal cost. *arXiv preprint arXiv:2306.15766*.

Markus Bayer and Christian Reuter. 2024. Activellm: Large language model-based active learning for textual few-shot scenarios. *arXiv preprint arXiv:2405.10808*.

Gantavya Bhatt, Yifang Chen, Arnav M Das, Jifan Zhang, Sang T Truong, Stephen Mussmann, Yinglun Zhu, Jeffrey Bilmes, Simon S Du, Kevin Jamieson, et al. 2024. An experimental design framework for label-efficient supervised finetuning of large language models. *arXiv preprint arXiv:2401.06692*.

Yifang Chen, Shuohang Wang, Ziyi Yang, Hiteshi Sharma, Nikos Karampatziakis, Donghan Yu, Kevin Jamieson, Simon Shaolei Du, and Yelong Shen. 2024. Cost-effective proxy reward model construction with on-policy and active learning. *arXiv preprint arXiv:2407.02119*.

Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin Rostamizadeh, and Sanjiv Kumar. 2021. Batch active learning at scale. *Advances in Neural Information Processing Systems*, 34.

David Cohn, Les Atlas, and Richard Ladner. 1994. Improving generalization with active learning. *Machine learning*, 15(2):201–221.

David A Cohn, Zoubin Ghahramani, and Michael I Jordan. 1996. Active learning with statistical models. *Journal of artificial intelligence research*, 4:129–145.

Sanjoy Dasgupta, Adam Tauman Kalai, and Adam Tauman. 2009. Analysis of perceptron-based active learning. *Journal of Machine Learning Research*, 10(2).

Shizhe Diao, Pengcheng Wang, Yong Lin, Rui Pan, Xiang Liu, and Tong Zhang. 2023. Active prompting with chain-of-thought for large language models. *arXiv preprint arXiv:2302.12246*.

Li Du, Zhouhao Sun, Xiao Ding, Yixuan Ma, Yang Zhao, Kaitao Qiu, Ting Liu, and Bing Qin. 2024. Causal-guided active learning for debiasing large language models. *arXiv preprint arXiv:2408.12942*.

Mostafa Fathy. 2025. Dblp scholar dataset.

Yifan Fu, Xingquan Zhu, and Bin Li. 2013. A survey on instance selection for active learning. *Knowledge and information systems*, 35(2):249–283.

Yarin Gal, Riashat Islam, and Zoubin Ghahramani. 2017. Deep Bayesian active learning with image data. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1183–1192. PMLR.

Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. 2024. Bias and fairness in large language models: A survey. *Computational Linguistics*, 50(3):1097–1179.

Yonatan Geifman and Ran El-Yaniv. 2017. Deep active learning over the long tail. *arXiv preprint arXiv:1711.00941*.

Steve Hanneke. 2007. A bound on the label complexity of agnostic active learning. In *Proceedings of the 24th international conference on Machine learning*, pages 353–360.

Hideitsu Hino. 2020. Active learning: Problem settings and recent developments. *arXiv preprint arXiv:2012.04225*.

Junjie Huang, Jizheng Chen, Jianghao Lin, Jiarui Qin, Ziming Feng, Weinan Zhang, and Yong Yu. 2024. A comprehensive survey on retrieval methods in recommender systems.

Daniel P Jeong, Zachary C Lipton, and Pradeep Ravikumar. 2024. Llm-select: Feature selection with large language models. *arXiv preprint arXiv:2407.02694*.

Kaixuan Ji, Jiafan He, and Quanquan Gu. 2024. Reinforcement learning from human feedback with active queries. *arXiv preprint arXiv:2402.09401*.

Nataliia Kholodna, Sahib Julka, Mohammad Khodadadi, Muhammed Nurullah Gumus, and Michael Granitzer. 2024. Llms in the loop: leveraging large language model annotations for active learning in low-resource languages. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 397–412. Springer.

Andreas Kirsch, Joost Van Amersfoort, and Yarin Gal. 2019. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *Advances in neural information processing systems*, 32.

Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.

Jan Kremer, Kim Steenstrup Pedersen, and Christian Igel. 2014. Active learning with support vector machines. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 4(4):313–326.

Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2024. Rewardbench: Evaluating reward models for language modeling.

Dongyuan Li, Zhen Wang, Yankai Chen, Renhe Jiang, Weiping Ding, and Manabu Okumura. 2024a. A survey on deep active learning: Recent advances and new frontiers. *IEEE Transactions on Neural Networks and Learning Systems*.

Dongyuan Li, Ying Zhang, Zhen Wang, Shiyin Tan, Satoshi Kosugi, and Manabu Okumura. 2024b. Active learning for abstractive text summarization via llm-determined curriculum and certainty gain maximization. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 8959–8971.

Tao Li, Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Vivek Srikumar. 2020. UNQOVERing stereotyping biases via underspecified questions. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 3475–3489, Online. Association for Computational Linguistics.

Chengyuan Liu, Fubang Zhao, Kun Kuang, Yangyang Kang, Zhuoren Jiang, Changlong Sun, and Fei Wu. 2024. Evolving knowledge distillation with large language models and active learning. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 6717–6731, Torino, Italia. ELRA and ICCL.

Yuxuan Lu, Bingsheng Yao, Shao Zhang, Yun Wang, Peng Zhang, Tun Lu, Toby Jia-Jun Li, and Dakuo Wang. 2023. Human still wins over llm: An empirical study of active learning on domain-specific annotation tasks. *arXiv preprint arXiv:2311.09825*.

Hanjia Lyu, Junda Wang, Wei Wu, Viet Duong, Xiyang Zhang, Timothy D Dye, and Jiebo Luo. 2022. Social media study of public opinions on potential covid-19 vaccines: informing dissent, disparities, and dissemination. *Intelligent medicine*, 2(01):1–12.

Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. 2011. Learning word vectors for sentiment analysis. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 142–150, Portland, Oregon, USA. Association for Computational Linguistics.

Katerina Margatina, Timo Schick, Nikolaos Aletras, and Jane Dwivedi-Yu. 2023. Active learning principles for in-context learning with large language models. *arXiv preprint arXiv:2305.14264*.

Costas Mavromatis, Balasubramaniam Srinivasan, Zhengyuan Shen, Jiani Zhang, Huzefa Rangwala, Christos Faloutsos, and George Karypis. 2024. Covericl: Selective annotation for in-context learning via active graph coverage. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21268–21286.

Luckeciano C. Melo, Panagiotis Tigas, Alessandro Abate, and Yarin Gal. 2024. Deep bayesian active learning for preference modeling in large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Xuran Ming, Shoubin Li, Mingyang Li, Lvlong He, and Qing Wang. 2024. Autolabel: Automated textual data annotation method based on active learning and large language model. In *International Conference on Knowledge Science, Engineering and Management*, pages 400–411. Springer.

Moran Mizrahi, Guy Kaplan, Dan Malkin, Rotem Dror, Dafna Shahaf, and Gabriel Stanovsky. 2024. State of what art? a call for multi-prompt LLM evaluation. *Transactions of the Association for Computational Linguistics*, 12:933–949.

Subhojyoti Mukherjee, Anusha Lalitha, Aniket Deshmukh, Ge Liu, Yifei Ma, and Branislav Kveton. 2024a. Experimental design for active transductive inference in large language models. *arXiv preprint arXiv:2404.08846*.

Subhojyoti Mukherjee, Anusha Lalitha, Sailik Sengupta, Aniket Deshmukh, and Branislav Kveton. 2024b. Multi-objective alignment of large language models through hypervolume maximization. *arXiv preprint arXiv:2412.05469*.

William Muldrew, Peter Hayes, Mingtian Zhang, and David Barber. 2024. Active preference learning for large language models. *arXiv preprint arXiv:2402.08114*.

Fredrik Olsson. 2009. A literature survey of active machine learning in the context of natural language processing.

Flavio Di Palo, Prateek Singhi, and Bilal H Fadlallah. 2024. Performance-guided LLM knowledge distillation for efficient text classification at scale. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3675–3687, Miami, Florida, USA. Association for Computational Linguistics.

Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales.

Ritik Sachin Parkar, Jaehyung Kim, Jong Inn Park, and Dongyeop Kang. 2024. Selectllm: Can llms select important instructions to annotate? *arXiv preprint arXiv:2401.16553*.

Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel R. Bowman. 2022. Bbq: A hand-built bias benchmark for question answering.

Yotam Perlitz, Ariel Gera, Michal Shmueli-Scheuer, Dafna Sheinwald, Noam Slonim, and Liat Ein Dor. 2023. Active learning for natural language generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9862–9877.

Wasu Top Piriyakulkij, Volodymyr Kuleshov, and Kevin Ellis. 2023. Active preference inference using language models and probabilistic reasoning. *arXiv preprint arXiv:2312.12009*.

Friedrich Pukelsheim. 2006. *Optimal Design of Experiments*. Society for Industrial and Applied Mathematics.

Zac Pullar-Strecker, Katharina Dost, Eibe Frank, and Jörg Wicker. 2024. Hitting the target: stopping active learning at the cost-based optimum. *Machine Learning*, 113(4):1529–1547.

Weihong Qi, Hanjia Lyu, and Jiebo Luo. 2024. Representation bias in political sample simulations with large language models. *arXiv preprint arXiv:2407.11409*.

Kun Qian, Yisi Sang, Farima Fatahi Bayat, Anton Belyi, Xianqi Chu, Yash Govind, Samira Khorshidi, Rahul Khot, Katherine Luna, Azadeh Nikfarjam, et al. 2024. Ape: Active learning-based tooling for finding informative few-shot examples for llm-based entity matching. *arXiv preprint arXiv:2408.04637*.

Zhongyi Qiu, Kangyi Qiu, Hanjia Lyu, Wei Xiong, and Jiebo Luo. 2024. Semantics preserving emoji recommendation with large language models. In *2024 IEEE International Conference on Big Data (BigData)*, pages 7131–7140. IEEE.

Utkarsh Rajput. 2024. Fake news dataset.

Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B Gupta, Xiaojiang Chen, and Xin Wang. 2021. A survey of deep active learning. *ACM Computing Surveys (CSUR)*, 54(9):1–40.

Hamidreza Rouzegar and Masoud Makrehchi. 2024. Enhancing text classification through llm-driven active learning and human annotation. *arXiv preprint arXiv:2406.12114*.

Noveen Sachdeva, Benjamin Coleman, Wang-Cheng Kang, Jianmo Ni, Lichan Hong, Ed H Chi, James Caverlee, Julian McAuley, and Derek Zhiyuan Cheng. 2024. How to train data-efficient llms. *arXiv preprint arXiv:2402.09668*.

Christopher Schröder and Andreas Niekler. 2020. A survey of active learning for text classification using deep neural networks. *arXiv preprint arXiv:2008.07267*.

Ozan Sener and Silvio Savarese. 2017. Active learning for convolutional neural networks: A core-set approach. *arXiv preprint arXiv:1708.00489*.

Burr Settles. 2009. Active learning literature survey.

Burr Settles. 2011. From theories to queries: Active learning in practice. In *Active learning and experimental design workshop in conjunction with AISTATS 2010*, pages 1–18. JMLR Workshop and Conference Proceedings.

Sandeep Silwal, Sara Ahmadian, Andrew Nystrom, Andrew McCallum, Deepak Ramachandran, and Seyed Mehran Kazemi. 2023. Kwikbucks: Correlation clustering with cheap-weak and expensive-strong signals. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.

Zhen Tan, Dawei Li, Song Wang, Alimohammad Beigi, Bohan Jiang, Amrita Bhattacharjee, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. 2024. Large language models for data annotation and synthesis: A survey.

Karan Taneja and Ashok Goel. 2024. Can active label correction improve llm-based modular ai systems? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 9019–9031.

Simon Tong. 2001. *Active learning: theory and applications*. Stanford University.

Simon Tong and Daphne Koller. 2001. Support vector machine active learning with applications to text classification. *Journal of machine learning research*, 2(Nov):45–66.

Dan Wang and Yi Shang. 2014. A new active labeling method for deep learning. In *2014 International joint conference on neural networks (IJCNN)*, pages 112–119. IEEE.

Xinru Wang, Hannah Kim, Sajjadur Rahman, Kushan Mitra, and Zhengjie Miao. 2024. Human-llm collaborative annotation through effective verification of llm labels. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–21.

Yu Xia, Fang Kong, Tong Yu, Liya Guo, Ryan A Rossi, Sungchul Kim, and Shuai Li. 2024a. Which llm to play? convergence-aware online model selection with time-increasing bandits. In *Proceedings of the ACM on Web Conference 2024*, pages 4059–4070.

Yu Xia, Xu Liu, Tong Yu, Sungchul Kim, Ryan Rossi, Anup Rao, Tung Mai, and Shuai Li. 2024b. Hallucination diversity-aware active learning for text summarization. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8665–8677, Mexico City, Mexico. Association for Computational Linguistics.

Ruixuan Xiao, Yiwen Dong, Junbo Zhao, Runze Wu, Minmin Lin, Gang Chen, and Haobo Wang. 2023. FreeAL: Towards human-free active learning in the

era of large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14520–14535, Singapore. Association for Computational Linguistics.

Zhipeng Xu, Zhenghao Liu, Yukun Yan, Shuo Wang, Shi Yu, Zheni Zeng, Chaojun Xiao, Zhiyuan Liu, Ge Yu, and Chenyan Xiong. 2024. Activerag: Autonomously knowledge assimilation and accommodation through retrieval-augmented agents.

Rui Yang, Xiaoman Pan, Feng Luo, Shuang Qiu, Han Zhong, Dong Yu, and Jianshu Chen. 2024. Rewards-in-context: Multi-objective alignment of foundation models with dynamic preference adjustment. *arXiv preprint arXiv:2402.10207*.

Bingsheng Yao, Ishan Jindal, Lucian Popa, Yannis Katsis, Sayan Ghosh, Lihong He, Yuxuan Lu, Shashank Srivastava, Yunyao Li, James Hendler, et al. 2023. Beyond labels: Empowering human annotators with natural language explanations through a novel active-learning architecture. *arXiv preprint arXiv:2305.12710*.

Qinyuan Ye, Bill Yuchen Lin, and Xiang Ren. 2021. Crossfit: A few-shot learning challenge for cross-task generalization in nlp.

Yue Yu, Lingkai Kong, Jieyu Zhang, Rongzhi Zhang, and Chao Zhang. 2022. AcTune: Uncertainty-based active self-training for active fine-tuning of pretrained language models. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1422–1436, Seattle, United States. Association for Computational Linguistics.

Bo Yuan, Yulin Chen, Yin Zhang, and Wei Jiang. 2024. Hide and seek in noise labels: Noise-robust collaborative active learning with llms-powered assistance. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10977–11011.

Xueying Zhan, Qingzhong Wang, Kuan-hao Huang, Haoyi Xiong, Dejing Dou, and Antoni B Chan. 2022. A comparative survey of deep active learning. *arXiv preprint arXiv:2203.13450*.

Jiaxin Zhang, Zhuohang Li, Kamalika Das, and Sricharan Kumar. 2024a. Interactive multi-fidelity learning for cost-effective adaptation of language model with sparse human supervision. *Advances in Neural Information Processing Systems*, 36.

Jifan Zhang and Robert Nowak. 2024. Sieve: General purpose data filtering system matching gpt-4o accuracy at 1% the cost. *arXiv preprint arXiv:2410.02755*.

Rui Zhang and Joel Tetreault. 2019. This email could save your life: Introducing the task of email subject line generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 446–456, Florence, Italy. Association for Computational Linguistics.

Ruoyu Zhang, Yanzeng Li, Yongliang Ma, Ming Zhou, and Lei Zou. 2023. LLMaAA: Making large language models as active annotators. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 13088–13103, Singapore. Association for Computational Linguistics.

Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.

Xinnong Zhang, Jiayu Lin, Libo Sun, Weihong Qi, Yihang Yang, Yue Chen, Hanjia Lyu, Xinyi Mou, Siming Chen, Jiebo Luo, et al. 2024b. Electionsim: Massive population election simulation powered by large language model driven agents. *arXiv preprint arXiv:2410.20746*.

Yifei Zhang, Bo Pan, Chen Ling, Yuntong Hu, and Liang Zhao. 2024c. Elad: Explanation-guided large language models active distillation. *arXiv preprint arXiv:2402.13098*.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E Gonzalez, and Ion Stoica. 2023a. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems 36 (NeurIPS)*, volume 36, pages 46595–46623.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023b. Judging llm-as-a-judge with mt-bench and chatbot arena.

Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Yue Zhang, Neil Gong, and Xing Xie. 2024. Promptrobust: Towards evaluating the robustness of large language models on adversarial prompts. In *Proceedings of the 1st ACM Workshop on Large AI Systems and Models with Privacy and Safety Analysis*, LAMPS '24, page 57–68, New York, NY, USA. Association for Computing Machinery.

Xiaojin Zhu, John Lafferty, and Zoubin Ghahramani. 2003. Combining active learning and semi-supervised learning using gaussian fields and harmonic functions. In *ICML 2003 workshop on the continuum from labeled to unlabeled data in machine learning and data mining*, volume 3, pages 58–65.

14564

## A LLM-based AL Goals

In this section, we highlight the various ways LLM-based AL can be used by providing an intuitive categorization based on the LLM-based AL goals. In particular, LLM-based AL techniques can focus on selecting or generating specific data points $\mathbf{X}$, learning an optimal prompt $P$, selecting contexts $C$ to include in the prompt, deciding on the LLM model $\mathcal{M}_i \in M$, features to leverage $x_j$, hyperparameters $\theta$, model architectures $\mathcal{A}$, data for annotation $\mathbf{X}_S$, and evaluation data $\mathbf{X}_{\text{eval}}$.

- **Data Points** ($X$): LLM-based active learning can be used to select specific data points $x \in X$ from a larger pool for different purposes, such as training, budget-constrained training, and domain adaptation. This also includes selecting data for annotation, $x_{\text{annot}}$, which maximizes learning gains, and choosing evaluation subsets, $x_{\text{eval}}$, that ensure comprehensive model assessment. Furthermore, we can also leverage LLMs to generate entirely novel data points that lie outside the set of unlabeled data. This

- **Prompts** ($P$): The technique can be applied to select the most effective prompts $p \in P$ and prompt variations that optimize LLM outputs for specific tasks.

- **Contexts** ($C$): LLM-based active learning can be utilized to select the most relevant contexts $c \in C$ or contextual inputs that enhance model performance on context-dependent tasks.

- **LLM Model Variants** ($\mathcal{M}$): The method can be used to select the best LLM model variant $\mathcal{M}_i \in \mathcal{M}$ from a set of available models for a given input. Alternatively, we can also use LLM-based AL to identify models that contribute effectively to ensembles.

- **Features** ($x_j$): LLM-based active learning can be used to select important features $x_j \in \mathbf{x}$ or identify new features that can be integrated to improve model accuracy and explainability. Furthermore, it can also be leveraged to estimate missing features, e.g., if there are several missing values in the feature vector of a specific instance.

- **Hyperparameters** ($\theta$): The technique is useful for selecting optimal hyperparameter values $\theta_i \in \theta$, such as learning rate or batch size, and adjusting them dynamically during training to optimize performance.

- **Model Architectures** ($\mathcal{A}$): LLM-based active learning can assist in selecting the most appropriate model architecture $\mathcal{A}_i \in \mathcal{A}$ or choosing between different versions of a model for specific use cases.

- **Data for Annotation** ($\mathbf{X}_S$): The technique can be used to select specific data subsets $\mathbf{X}_S \subset \mathbf{X}$ that should be labeled by human annotators, focusing on those that provide the greatest potential improvement in model learning.

- **Evaluation Data Subsets** ($\mathbf{X}_{\text{eval}}$): LLM-based active learning is beneficial for selecting evaluation data points $\mathbf{X}_{\text{eval}} \subset X$ that maximize the effectiveness of model evaluation, ensuring coverage of edge cases and comprehensive testing.

## B Traditional Selection Strategies

We provide additional discussions on traditional selection strategies from Section 3.1. The CoreSet is a pure diversity-based strategy where unlabeled examples are selected using a greedy furthest-first traversal conditioned on all labeled examples (Sener and Savarese, 2017; Geifman and El-Yaniv, 2017; Citovsky et al., 2021). Similarly Agarwal et al. (2020) uses a coreset-based strategy on features to select unlabeled examples using CNN. The Least confidence is an uncertainty-based active learning algorithm where the uncertainty score of an unlabeled example is its predicted class probability and the algorithm then samples unlabeled examples with the smallest uncertainty score (Settles, 2009, 2011; Wang and Shang, 2014). The Margin-based selection strategy is also an uncertainty-based strategy (Tong and Koller, 2001; Balcan et al., 2009; Settles, 2009). Margin first sorts unlabeled examples according to their multiclass margin score and then selects examples that are the hardest to discriminate and can be thought of as examples closest to their class margin. The Max-entropy strategy (Wang and Shang, 2014; Kremer et al., 2014; Diao et al., 2023) is an uncertainty-based strategy that selects unlabeled examples according to the entropy of the example's predictive class probability distribution. The BADGE algorithm combines both uncertainty and diversity sampling (Ash et al., 2019, 2021). Badge chooses a batch of unlabeled examples by applying $k$-Means++ (Arthur and Vassilvitskii, 2006) on the gradient embeddings computed from the penultimate layer of the model. The value of the gradient

embedding captures the uncertainty score of the examples. Finally, BALD (Bayesian Active Learning by Disagreements) (Kirsch et al., 2019; Gal et al., 2017) chooses unlabeled examples that are expected to maximize the information gained from the model parameters i.e. the mutual information between predictions and model posterior.

## C  AL Paradigms with LLMs

With the rise of LLMs, the active learning has evolved to address new challenges and opportunities across various learning paradigms.

### C.1  Active In-Context Learning

Recent advances have recast demonstration selection for in-context learning as an active learning problem, where the goal is to identify and annotate the most informative examples under stringent labeling budgets. For example, Margatina et al. (2023) demonstrate that similarity-based sampling can consistently outperform traditional uncertainty-based methods when selecting single-round demonstrations for LLMs. Building on these insights, Mavromatis et al. (2024) propose CoverICL—a graph-based algorithm that integrates uncertainty sampling with semantic coverage to enhance both performance and budget efficiency across diverse tasks and LLM architectures. In a complementary approach, Qian et al. (2024) introduce APE, a human-in-the-loop tool that iteratively pinpoints ambiguous examples for few-shot prompts, thus progressively refining LLM performance through active learning principles.

### C.2  Active Supervised Finetuning

Active learning strategies have also been adapted to reduce the labeling cost associated with supervised finetuning. For instance, Yu et al. (2022) present AcTune, which actively queries high-uncertainty instances while simultaneously leveraging self-training on low-uncertainty unlabeled data. This dual strategy, further refined by region-aware sampling, effectively mitigates redundancy in the training data. Similarly, Xia et al. (2024b) propose an active learning framework tailored for text summarization that systematically identifies and annotates diverse instances exhibiting various types of hallucinations, thereby improving both data efficiency and factual correctness. Extending these ideas to overcome cold-start challenges, Bayer and Reuter (2024) develop ActiveLLM, which employs state-of-the-art LLMs (e.g., GPT-4, Llama 3) to select informative instances, significantly boosting fine-tuning performance in both few-shot and iterative settings.

### C.3  Active Preference Alignment

Active preference alignment targets label-efficient methods for refining LLMs via human or AI feedback. Ji et al. (2024) frame reinforcement learning from human feedback (RLHF) as a contextual dueling bandit problem and develop an active-query-based algorithm (ADPO) that drastically cuts down the number of preference queries needed for LLM alignment. Building on simpler, more stable methods, Muldrew et al. (2024) propose an active learning extension to Direct Preference Optimization (DPO), showing notable gains in both convergence speed and final quality through selective preference labeling. Complementarily, Chen et al. (2024) introduce a cost-effective approach for constructing reward models, combining on-policy querying and active data selection to maximize the impact of limited human feedback, and achieving strong performance improvements in DPO with minimal expert annotation.

### C.4  Active Knowledge Distillation

Active knowledge distillation has emerged to reduce the computational costs of LLMs by selectively transferring their knowledge into smaller models. Zhang et al. (2024c) propose ELAD, which leverages reasoning-step uncertainties to guide sample selection and employs teacher-driven explanation revisions to optimize the distillation process. Liu et al. (2024) introduce EvoKD, an iterative strategy that identifies student model weaknesses and dynamically generates labeled data, continuously refining the student's capabilities through LLM feedback. Meanwhile, Palo et al. (2024) present PGKD, an active feedback loop that uses performance signals such as hard-negative mining to inform new data creation, yielding substantial efficiency gains and reducing inference costs for text classification at scale.

## D  From AL to Bandits and Beyond

In this section, we briefly discuss a unifying view of the LLM-based AL problem through the lens of bandits, online learning, RLHF, pandoras box, among others. We provide an intuitive summary of such connections in Table 4 across a variety of

| | Sequential Learning | Active Querying | Human Feedback | Exploration-Exploitation | Statistical Foundations | Dynamic Environments |
|---|---|---|---|---|---|---|
| **Active Learning** | ✓ | ✓ | | ✓ | ✓ | ✓ |
| **RLHF** | ✓ | ✓ | ✓ | ✓ | | ✓ |
| **Pandora's Box** | ✓ | ✓ | | ✓ | ✓ | ✓ |
| **Bandits** | ✓ | | | ✓ | | ✓ |
| **Optimal Design** | | | | | ✓ | |
| **Online Learning** | ✓ | | | ✓ | ✓ | ✓ |

Table 4: Comparison of problem settings across a range of important properties. *Sequential Learning*: Learning incrementally over time, rather than all at once. *Active Querying*: Actively selecting or querying data points to improve the model. *Human Feedback*: Utilizing human input or feedback to guide and improve learning. *Exploration-Exploitation*: Balancing the need to explore new options versus exploiting known ones. *Cost-Awareness*: Considering the costs of acquiring data or feedback during the learning process. *Statistical Foundations*: Grounding the technique in statistical principles or experimental design. *Dynamic Environments*: Adapting to environments that change over time or where new information continuously emerges.

problem settings. Intuitively, they overlap in goals, that is, many of these techniques aim to optimize resource use (e.g., labels, feedback, compute) while improving learning. They nearly all operate sequentially, but they differ in what drives decisions (e.g., reward, uncertainty, cost). However, they may have different constraints, for instance, some settings assume explicit costs (e.g., pandora's box), while others lack such constraints (e.g., online learning). These insights allow practitioners to choose the right framework depending on the data, task, and constraints of their problem.

# E   Applications

LLM-based active learning offers many promising applications; from LLM-to-SLM knowledge distillation to low-resource language translation, to entity matching. Active learning reduces the cost of data supervision, broadening the tasks machine learning can be applied to. Active learning, combined with causal learning, may also be useful for reducing measured bias (Sec. E.6). Table 1

offers a taxonomy of active learning techniques. Overall, active learning is likely to continue to offer benefits in data-constrained environments. In particular, active learning may increase the utility of LLMs in domains requiring expert knowledge, such as medicine, law, and engineering; a historical weakness (Yao et al., 2023). Finally, Table 5, we include a pairing of active learning applications and datasets used.

## E.1   Text Classification

Rouzegar and Makrehchi (2024) apply active learning to identify the most relevant samples for labeling text. The framework is applied to, among other datasets, Fake News for document authenticity (Rajput, 2024). The method reduces the cost of obtaining supervision and also enables a trade-off between cost efficiency and performance. The authors use GPT-3.5 and require some human supervision, unlike some more recent work (Du et al., 2024). ActiveLLM (Bayer and Reuter, 2024) uses LLMs for selecting instances for few-shot text classification with model mismatch, in which the selection, or query, model is different from the model used for the final task. ActiveLLM use an LLM to estimate data point uncertainty and diversity without external supervision, improving few-shot learning.

## E.2   Text Summarization

Li et al. (2024b) proposes LLM-Determined Curriculum Active Learning (LDCAL) that improves the stability of the active learner by selecting instances from easy to hard by using LLMs to determine the difficulty of a document. LDCAL also leverages a new AL technique termed Certainty Gain Maximization that captures how well the unselected instances are represented by the selected ones, which is then used to select instances that maximize the certainty gain for the unselected ones. More formally, the certainty gain (CG) measures the gain in representation certainty for an unlabeled instance $\mathbf{x}_u$ when a candidate instance $\mathbf{x}_s$ is selected for annotation, that is,

$$\text{CG}(\mathbf{x}_s, \mathbf{x}_u) = \max\left(f(\mathbf{x}_s, \mathbf{x}_u) - \max_{\mathbf{x}_i \in \mathcal{D}_\ell} f(\mathbf{x}_u, \mathbf{x}_i), 0\right)$$

It is derived as the maximum of the difference between $f(\mathbf{x}_s, \mathbf{x}_u)$, which measures the similarity between $\mathbf{x}_s$ and $\mathbf{x}_u$, and the current maximum similarity $\max_{\mathbf{x}_i \in \mathcal{D}_\ell} f(\mathbf{x}_u, \mathbf{x}_i)$, where $\mathbf{x}_i$ represents instances in the labeled set $\mathcal{D}_\ell$. To ensure non-negative values, any negative difference is clipped

to 0. This formulation ensures that the selection of $\mathbf{x}_s$ positively contributes to the representational coverage of the unlabeled instances, balancing the sampling process across high-density and low-density regions in the data distribution. Finally, the Average Certainty Gain (ACG) for a candidate instance $\mathbf{x}_s$ is

$$\text{ACG}(\mathbf{x}_s) = \frac{1}{L} \sum_{\mathbf{x}_u \in \mathcal{D}_u} \text{CG}(\mathbf{x}_s, \mathbf{x}_u)$$

where $L$ is the total number of unlabeled instances and $\mathcal{D}_u$ is the set of all unlabeled instances. A key advantage of LDCAL over uncertainty-based acquisition strategies is that the acquisition model does not need to be trained after each AL iteration.

### E.3 Non-Textual Classification

Margatina et al. (2023) apply active learning to several areas, including multiple-choice question answering and test on 15 different models from the GPT and OPT families. The selection of similar in-context samples for multiple-choice questions was the most effective sampling criterion. For classification, diversity was more effective than similarity as selection criterion. The authors also found larger models to have higher performance.

### E.4 Question Answering

Diao et al. (2023), in recognition of the utility of chain-of-thought prompting, ActivePrompt uses active learning to design more effective prompts for LLMs based on human-designed chains of thought. With only eight exemplars made by humans, the method achieves higher performance on complex reasoning tasks. Uncertainty is estimated by querying an LLM with the same prompt repeatedly to and response disagreement is measured.

### E.5 Entity Matching

An approach called APE (Qian et al., 2024) focuses on an active prompt engineering approach for entity matching where at each iteration, a set of prompts are derived, and then evaluated by a committee of models, where the best is selected, and then the approach repeats. APE selects the most informative samples of a dataset reduces the cost to humans of identifying samples in most need of human feedback. Ming et al. (2024) proposed AutoLabel that starts by selecting the most representative seed data using traditional techniques such as density clustering and sampling. Then uses an LLM with chain-of-thought prompting to obtain labels, and then

| Task | Datasets |
| --- | --- |
| **Text Classification** | IMDB (Maas et al., 2011) |
| **Text Summarization** | AESLC (Zhang and Tetreault, 2019) |
| **Non-Text Classification** | Crossfit (Ye et al., 2021) |
| **Question Answering** | Crossfit (Ye et al., 2021) |
| **Entity Modeling** | DBLP-Scholar (Fathy, 2025) |
| **Debiasing** | BBQ (Parrish et al., 2022), MT-Bench (Zheng et al., 2023a), UNQOVER (Li et al., 2020) |
| **Translation** | IT domain (Koehn and Knowles, 2017) |
| **Sentiment Analysis** | IMDB (Maas et al., 2011), AESLC (Zhang and Tetreault, 2019), AG-News (Zhang et al., 2015) |

Table 5: LLM-based Active Learning Datasets

leverages human feedback to rectify the labeled results for entity recognition. Zhang et al. (2023) proposed LLMaAA that leverages LLMs for annotation in an active learning loop. LLMaAA is used for both named entity recognition and relation extraction.

### E.6 Active Debiasing of LLMs

Du et al. (2024) proposed a causal-guided AL approach for debiasing LLMs by leveraging the LLMs to select data samples that contribute bias to the dataset. The method works by applying active learning to identify the most important samples within the dataset and the model looks for causally invariant relationships; this approach is less compute-intensive than fine-tuning a model on a debiasing dataset.

### E.7 Translation

Kholodna et al. (2024) apply active learning for annotations in low-resource languages and find near-SOTA performance, with a reduction in annotation cost (relative to human annotators) of 42.45 times. Like other active learning methods, samples with the highest prediction uncertainty are selected. The authors test on 20 low-resource languages spoken in Sub-Saharan Africa. The authors found larger LLMs like GPT-4-Turbot and Claude 3 Opus had more consistent performance than Llama 2-70B and Mistral 7B. FreeAL use active learning to collect data for task-specific knowledge, such as

translation (Xiao et al., 2023), without requiring human annotation. Tested on eight benchmarks, FreeAL achieves near-human-supervised performance, without requiring any human annotation. With additional feedback rounds, FreeAL was able to improve performance. In FreeAL, an LLM and SLM are paired and the LLM provides annotation, while the SLM is a weak learner. The authors propose the use of limited human supervision to further improve performance.

### E.8 Sentiment Analysis

Xiao et al. (2023) use active learning for movie sentiment analysis with the Movie Review dataset Seeing Stars (Pang and Lee, 2005) in the method ActivePrune, which uses ordinal (movie stars, 1-5) labels. ActivePrune outperforms other pruning methods and with its increased compute efficiency, they reduce end-to-end active learning time by 74%. ActivePrune works by reducing dataset size and increasing representation of underrepresented data, based on perplexity.

### E.9 Other

In this section, we describe additional noteworthy applications of active learning.

#### E.9.1 Question Generation

Piriyakulkij et al. (2023) proposed an active preference inference approach that infers the preferences of individual users by minimizing the amount of questions to ask the user to obtain their preferences. The approach uses more informative questions to improve the user experience of such systems.

#### E.9.2 System Design

Taneja and Goel (2024) leverages an approach for actively correcting labels to enhance LLM-based systems that have multiple components. Astorga et al. (2024) introduced a partially observable cost-aware AL method focused on the setting where features and/or labels may be partially observed.

## F LLM-based AL Advantages

By iteratively selecting and generating instances (to label) for training, LLM-based active learning can have the following advantages:

- **Better Accuracy:** LLM-based active learning can achieve better accuracy with fewer instances by selecting and generating the most informative instances.

- **Reduced Annotation Costs:** LLM-based active learning techniques can reduce labeling and other annotation costs by selecting or generating the most informative samples to use for model training and fine-tuning.

- **Faster Convergence:** Using LLM-based active learning often leads to faster convergence as the model can be learned more quickly by iteratively selecting and generating the most informative examples.

- **Improved Generalization:** By leveraging LLM-based active learning techniques that iteratively select and generate the most informative and diverse examples, the model can generalize better to new data.

- **Robustness:** Iteratively selecting or generating the best examples for training can often improve the robustness to noise, as the model is learned from a set of high-quality representative examples.