SIEVE: GENERAL PURPOSE DATA FILTERING SYSTEM MATCHING GPT-40 ACCURACY AT 1% THE COST

Anonymous authors

Paper under double-blind review

ABSTRACT

Creating specialized large language models requires vast amounts of clean, special purpose data for training and fine-tuning. With only a handful of existing large-scale, domain-specific datasets, creation of new datasets is required in most applications. This requires the development of new application-specific filtering of web-scale data. Filtering with a high-performance, general-purpose LLM such as GPT-40 can be highly effective, but this is extremely expensive at web-scale. This paper proposes SIEVE, a lightweight alternative that matches GPT-40 accuracy at a fraction of the cost. SIEVE can perform up to 500 filtering operations for the cost of one GPT-40 filtering call. The key to SIEVE is a seamless integration of GPT-40 and lightweight T5 models, using active learning to fine-tune T5 in the background with a small number of calls to GPT-40. Once trained, it performs as well as GPT-40 at a tiny fraction of the cost. We experimentally validate SIEVE on the OpenWebText dataset, using five highly customized filter tasks targeting high quality and domain-specific content. Our results demonstrate the effectiveness and efficiency of our method in curating large, high-quality datasets for language model training at a substantially lower cost (1%) than existing techniques. To further validate SIEVE, experiments show that SIEVE and GPT-40 achieve similar accuracy, with human evaluators preferring SIEVE's filtering results to those of GPT-40.

027 028 029

030

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

1 INTRODUCTION

031 Large Language Models (LLMs) have revolutionized natural language processing, demonstrating 032 remarkable capabilities across a wide range of tasks. As the field progresses, there is growing interest 033 in developing specialized LLMs tailored to specific domains or applications (Lee et al., 2023; Gupta 034 et al., 2024; Li et al., 2024). A critical component in this process is the curation of high-quality, domain-specific datasets for training and fine-tuning these models. However, the task of filtering vast 035 amounts of web-scale data to create such datasets presents significant challenges in terms of cost, time, and effectiveness. Existing approaches to data filtering and curation primarily rely on two strategies: 037 (1) sourcing data from specific, trusted sources, and (2) employing pre-existing quality and toxicity detection models to filter harmful content. For instance, medical LLMs often utilize datasets derived from PubMed to ensure domain specificity. While these data filtering methods have their merits, they 040 suffer from notable limitations in terms of both flexibility and comprehensiveness. Many domains 041 lack comprehensive, exclusive sources of high-quality text data, and pre-trained detection models 042 typically focus on general quality metrics and toxicity, limiting their applicability for domain-specific 043 queries or customized quality rubrics. Furthermore, relying solely on established data sources ignores 044 the vast amount of potentially valuable information available on the broader internet, leading to datasets that may be limited in scope and diversity. These constraints can significantly hinder the development of specialized language models that require rich, domain-specific training data. 046

Recent advancements in general-purpose language models, such as GPT-40, offer a potential solution to these challenges. These models can act as effective filtering mechanisms when provided with appropriate filtering prompts. To illustrate as an example, one could employ GPT-40 to iterate over all text snippets of the internet to determine whether each piece pertains to the 2024 presidential election. In this context, the machine learning model serves as a binary classification mechanism, evaluating one snippet at a time. However, the computational cost of applying models like GPT-40 to web-scale datasets is prohibitively expensive for most organizations. In this paper, we present a novel system SIEVE that addresses these limitations, providing versatile and high-quality data



Figure 1: **System Overview.** From user's perspective, SIEVE acts as if applying GPT-40 with the filtering prompt to all text snippets in a web-scale dataset. The output from the SIEVE system is the set of all text snippets that receive a 'pass'. To reduce the prohibitively high cost of applying GPT-40 on every snippet, SIEVE utilizes active learning to distill lightweight filtering models based on pretrained T5 encoders, effectively reducing the overall cost to less than 1%.

filtering at a fraction of the cost of using high performance, general-purpose models directly. Shown in Figure 2, our method is based on a lightweight T5 model in combination with a small number of calls to GPT-40, effectively reducing the cost per filtering to one call to T5 instead of GPT-40. This is accomplished by training and distilling a lightweight T5 model in the background that learns to mimic the filtering decisions of GPT-40 and eventually handles all/most of the filtering. This background job is optimized via a novel active learning algorithm that further reduces the need to call GPT-40.

Our active learning algorithm operates in a stream-based setting, sequentially and adaptively selecting the most informative text snippets for evaluation by GPT-40. Inspired by previous work in deep 079 active learning (Zhang et al., 2022a; Nuggehalli et al., 2023), we propose a novel stream-based 080 algorithm designed to tackle scenarios where the data distribution is imbalanced. As all of the 081 filtering decisions are highly imbalanced (see Table 3), one of our main contributions is proposing 082 the algorithm to tackle this imbalance issue. Our algorithm is designed to efficiently label a set 083 of more balanced and uncertain examples. As demonstrated in our experiments, SIEVE achieves 084 GPT-40 quality filtering performance at less than 1% of the cost. In addition, compared to random 085 sampling methods, our active learning algorithm reduces the number of queries to GPT-40 by more 086 than 5x (see Section 5). In Section 4, we provide theoretical analysis of our algorithm, establishing formal proofs of balancedness bounds for active learning in scenarios with imbalanced underlying 087 data distributions. To the best of our knowledge, this represents the first attempt to provide rigorous 088 theoretical guarantees in this context. 089

The implications of SIEVE are far-reaching, democratizing access to clean, task-specific data and
 facilitating the development of specialized language models across various domains. By dramatically
 reducing the cost and complexity of data filtering, SIEVE opens new avenues for researchers and
 organizations to create tailored LLMs for specific applications, industries, or scientific disciplines.

- ⁰⁹⁴ In summary, the core contributions of this work are:
 - A novel, cost-effective system, SIEVE, for high-quality data filtering that achieves GPT-40 level performance at less than 1% of the cost.
 - A stream-based active learning algorithm for imbalanced datasets that efficiently selects informative snippets, reducing the number of queries to GPT-40 more than 5x compared to random sampling.
 - Theoretical analysis and formal proofs of balancedness bounds for our active learning algorithm in scenarios with imbalanced data distributions.
 - Experimental validation of SIEVE on the OpenWebText dataset using five highly specific filters, demonstrating its effectiveness and versatility.

2 A GENERAL PURPOSE DATA FILTERING SYSTEM

105 106

096

098

099

100

101

102

103 104

066

067

068

069

070

2 A GENERAL FURPOSE DATA FILTERING STSTEM

107 In this section we provide a detailed description of the SIEVE system. A visualization of our system can be found in Figure 1.

A Bird's-Eye View. When viewed as a black box, SIEVE processes a web-scale dataset of N text snippets along with a filtering prompt that specifies the criteria for passing or failing each snippet. This prompt specifies the criteria for passing or failing each snippet, similar to how one would instruct any high-performance, general-purpose LLM, such as GPT-40. From this perspective, SIEVE efficiently categorizes each snippet as 'pass' or 'fail' based on the provided prompt, with filtering quality comparable to that of GPT-40.

A straightforward approach would be to apply GPT-40 with the filtering prompt to each text snippet. However, this becomes extremely costly for web-scale datasets containing billions or trillions of tokens. For example, filtering the OpenWebText dataset (Gokaslan & Cohen, 2019) used in this study, which contains 9 billion tokens, would cost approximately \$67,000 using GPT-40 directly.
For even larger datasets like the PILE (Gao et al., 2020), this cost would increase by at least 1000 times, making it prohibitively expensive. To overcome this challenge, we utilize an active distillation framework as detailed in the following section.

- 121 122
- 2.1 AN ACTIVE DISTILLATION FRAMEWORK

123 124

In this section, we describe the inner workings of SIEVE, which employs a lightweight binary

in this section, we describe the infer workings of STEVE, which employs a fightweight binary
 classification model trained on GPT-4o's filtering decisions. By leveraging active learning techniques,
 we selectively gather GPT-4o decisions on a small, informative subset of text snippets from the
 web-scale dataset. This approach significantly reduces costs while maintaining filtering quality
 comparable to GPT-4o.

129 Active learning minimizes annotation costs from expensive sources by selecting the most informative subset of snippets to query for labels. The goal is to train a high-performance model f with minimal 130 annotation cost. In our framework, GPT-40 serves as the expensive annotation source, while we 131 fine-tune a pretrained T5 encoder model for binary classification on the collected data and annotations 132 to achieve high performance. Active learning strategies collect annotations incrementally, retraining 133 the lightweight T5 model f after labeling every B new snippets. While most deep active learning 134 literature focuses on the pool-based setting (Ash et al., 2019; Nuggehalli et al., 2023; Fairstein et al., 135 2024; Lesci & Vlachos, 2024), these algorithms require forward inference of f on the entire dataset 136 at every iteration, incurring high computational costs for large-scale datasets (see Section 5.3 for 137 details). To mitigate this, we apply active learning in the streaming setting, an area well-studied 138 classically but with limited research for neural networks. We specifically designed our algorithm 139 to tackle the imbalance in filtering decisions generated by GPT-40 (see Table 3), aiming to query GPT-40 on a more balanced and informative set of text snippets. 140

141 Formally, we assume access to a stream of i.i.d. drawn snippets $x_1, x_2, ..., x_N$, all following the 142 same underlying data distribution \mathbb{P}_X . We let $S = x_1, ..., x_N$ denote the stream. In practice, we 143 construct the stream by randomly shuffling all snippets in the OpenWebText dataset (Gokaslan & 144 Cohen, 2019). At time i, the active learning algorithm observes $x_i \sim \mathbb{P}_X$ and decides whether to 145 query GPT-40 for its annotation based on the snippet's informativeness to model f. If queried, we 146 obtain the corresponding filtering decision from GPT-40, which we denote by $y_{GPT}(x_i) \in \{0, 1\}$. Here, the randomness comes from the nondeterministic nature of GPT-4o's response. After every 147 B new annotations, we fine-tune the T5-decoder model from its pretrained checkpoint to obtain an 148 updated model f. The distillation process terminates after a total annotation budget of T snippets. 149 The final model f fine-tuned on queried snippets is then applied to filter the entire web-scale dataset. 150

151

153

152

3 ACTIVE LEARNING ALGORITHM FOR DISTILLING GPT-40

154 In this section, we present a novel stream-based active learning algorithm for class-imbalanced 155 scenarios. We use a modified version of the uncertainty sampling strategy (Lewis & Gale, 1994; Tong 156 & Koller, 2001; Balcan et al., 2006; Settles, 2009; Kremer et al., 2014). In our setting, uncertainty 157 sampling labels snippets x_i that have predictive sigmoid score around 0.5, i.e. $f(x_i) \approx 0.5$, as 158 these are believed to be the most informative data to be labeled. However, our binary classification 159 tasks of data filtering is naturally imbalanced as shown in Table 3. Previous active learning work by Zhang et al. (2022a) and Nuggehalli et al. (2023) observed that the threshold of 0.5 is generally 160 biased towards the majority class under imbalanced scenarios, which means that most of the snippets 161 selected by uncertainty sampling will be in the majority class. This translates into poor performance

163

164

165

167

171

179



172 Figure 2: Demonstration of the TRM threshold. Snippets (shown on the bottom) are first ordered 173 based on their predictive sigmoid scores. GPT-40 class labels 0 and 1 are represented by the solid 174 or dashed borders. Queried snippets are shaded. Under imbalanced scenarios, sigmoid score of 0.5generally will not provide a good indication of where to sample, and will likely result in labeling 175 much more snippets in the majority class. The probability $\mu(s)$ denotes the likelihood of a snippet 176 with sigmoid score s belonging to class 0. The TRM threshold is defined to best separate the two 177 classes of snippets. 178

180	Algorithm 1 Stream-Based Class-Balancing Active Learning
181	Input: Data Stream S, labeling function y_{GPT} based on GPT-40 and specified filtering prompt,
182	batch size B, total budget T and confidence level δ .
183	Initialize: Query $x_1,, x_B \stackrel{iid}{\sim} \mathbb{P}_X$ to form the initial labeled set $L = \{x_i, y_i\}_{i=1}^B$.
184	for $r = 1,, \frac{T}{2}$ do
185	Fine-tune pretrained T5 encoder model on the latest labeled set L to form $f: S \to [0, 1]$.
186	Initialize confidence set $\mu, \bar{\mu} \leftarrow 0, 1$, counter $t \leftarrow 0$.
187	Let j index the head of the stream S, x_j .
188	while $ L < (s+1)B$ do // Find optimal separation threshold for f .
189	Receive the next snippet in the stream $S, x_{j+t} \sim \mathbb{P}_X$.
190	If $f(x_{j+t}) \in [\underline{\mu}, \mu]$ then
191	Query GPT-40 for label y_{j+t} by y_{GPT} , and insert to set $L \leftarrow L \cup \{x_{j+t}, y_{j+t}\}$.
192	end if if $t \in \mathbb{R}^{n+1}$ then t' Undets confidence interval
193	If $t \in 2^{n}$ then $n \in \mathbb{Z}^{n}$ then $n \in $
194	Store previously computed sigmoid scores $F_t \leftarrow \{0, f(x_j),, f(x_{j+t})\}$.
195	Compute the empirical TRM threshold as $\hat{s}_t \leftarrow \min_{s \in F_t} \mathcal{L}_t(s)$, where
196	$\mathcal{L}_t(s) := \frac{1}{t+1} \left(\sum_{i \in [j,j+t]: f(x_i) \le s} 1\{y_i \neq 0\} + \sum_{i \in [j,j+t]: f(x_i) > s} 1\{y_i \neq 1\} \right).$
197	Update $\underline{\mu}, \overline{\mu}$ be the smallest and largest thresholds $s \in F_t$ such that
198	$\hat{\boldsymbol{\alpha}}(\boldsymbol{\lambda}) = \hat{\boldsymbol{\alpha}}(\boldsymbol{\alpha}) + \hat{\boldsymbol{\alpha}} = \sqrt{(164 - 72 - 164)^2} + (162)^2 $
199	$\mathcal{L}_t(s) - \mathcal{L}_t(s_t) \le \beta_{t+1} \sqrt{(\{s' \in F_t : \min(s, s_t) \le s' \le \max(s, s_t)\} - 1)/t + \beta_{t+1}^2/2}$
200	where $\beta_{i+1} = \sqrt{2\log(2\log(t+1)^2N^2/\delta)/(t+1))}$ is chosen so that the TPM threshold
201	lies in the undated confidence threshold with probability of at least $1 - \delta/\log_2(N)$
202	end if
203	Update counter $t \leftarrow t + 1$. // Loop over snippets in S.
204	end while
205	end for
206	Return: T5 encoder model f finetuned on L . f is then used for filtering on the entire dataset.
207	
208	
209	

in detecting snippets from the target minority class. To combat this, our approach instead aims to 210 find a threshold near where the majority and minority classes are equiprobable, which we will refer 211 to as an True Risk Minimizer (TRM) threshold. Selecting and labeling snippets around the TRM 212 threshold yields a labeled dataset that is more balanced and includes uncertain snippets. In a nutshell, 213 our algorithm (detailed in Algorithm 1) alternates between 214

1. Spending *B* budget in finding and labeling close to the TRM threshold; 215

2. Fine-tuning a new T5 encoder model f on all labeled snippets thus far.

More concretely as shown in Figure 2, when ordering snippets according to their sigmoid scores, labeling snippets around the threshold of 0.5 can often result in selecting and labeling of most of snippets from a single (majority) class. Our proposed alternative, the TRM threshold, minimizes the expected number of class 1 snippets to its left and class 0 snippets to its right. Thus, the TRM threshold minimizes the expected number of misclassifications. Formally, let η_j denote the probability of snippet x_j belonging to class 0

$$\eta_j := \mathbb{P}\big(y_{\text{GPT}}(x_j) = 0\big)$$

where the possible randomness is with respect to GPT-4o's response as well as the distribution underlying the snippet $x_j \sim P_X$. The distribution P_X is unknown, so $\{\eta_j\}$ are unknown as well. For model f and a stream of i.i.d. text snippets $x_1, ..., x_N \sim \mathbb{P}_X$, the *True Risk Minimizer (TRM) threshold* is then the sigmoid score $s^* \in \{0, f(x_1), ..., f(x_N)\}$, where

222

223

 $s^{\star} := \arg\min_{s \in \{f(x_1), \dots, f(x_N)\}} \left(\sum_{j: f(x_j) \le s} (1 - \eta_j) + \sum_{j: f(x_j) > s} \eta_j \right).$ (1)

231 The optimization objective above can be viewed as the expected risk (probability of error) of a finite set of threshold classifiers with threshold locations at $\{0, f(x_1), \ldots, f(x_N)\}$. A threshold classifier 232 at s categorizes snippets into class 0 if their sigmoid score is less than or equal to s, and into class 1 233 if their score is greater than s. The TRM threshold can also be viewed as where the two classes best 234 separate from each other. Snippets around this threshold are therefore truly uncertain to the neural 235 network model. Moreover, in Section 4, we will provide a novel theoretical analysis demonstrating 236 labeling around the TRM threshold also improves the balancedness of the queried snippets, alleviating 237 the data imbalance issues when training the lightweight model f. Of course, the TRM threshold is 238 unknown (because the probabilities $\{\eta_i\}$ are unknown), so we employ an efficient active learning 239 procedure to identify a threshold close to it. 240

Our algorithm, shown in Algorithm 1, applies agnostic active learning techniques (Dasgupta et al., 241 2007; Jamieson & Jain, 2022; Katz-Samuels et al., 2021) to threshold classifiers. It focuses on 242 identifying the TRM threshold. The algorithm initializes by querying the first B snippets from the 243 stream to create a labeled set L. Each iteration begins with fine-tuning classifier f on L. To estimate 244 the TRM threshold s^* , we maintain a high-confidence interval $[\mu, \bar{\mu}]$ while annotating stream snippets. 245 Snippets with predictive sigmoid scores within this interval are queried using GPT-40. We update the 246 confidence interval on a geometric schedule, every 2^i snippets, using empirical estimates \hat{s}_t of the 247 optimal threshold. This interval, centered around the empirical estimate, ensures $s^* \in [\mu, \bar{\mu}]$ with 248 at least $1 - \delta$ probability across all $t = 2^i$ for the current iteration. For constructing the confidence 249 interval, we employ an empirical Bernstein bound, as introduced in Jamieson & Jain (2022). This 250 approach differs from the original CAL algorithm (Dasgupta et al., 2007), which uses a uniform convergence bound requiring synchronous querying. Our method allows for parallel GPT-40 queries, 251 significantly reducing computational time in practice. 252

253 254 255

256

257

258

259

4 THEORETICAL ANALYSIS

In this section, we first provide formal guarantee of the performance of our algorithm. We then proceed to analyze the balancedness of the snippets our algorithm queries, showing its improvement in collecting a more balanced set of snippets. Note that all of our analysis are conducted for any single iteration s in Algorithm 1. Define the true risk function over all snippets $1, \ldots, N$ at threshold s is denoted by

$$R(s) = \frac{1}{N} \left(\sum_{i \le N: f(x_i) \le s} (1 - \eta_i) + \sum_{i \le N: f(x_i) > s} \eta_i \right) \,.$$

Recall the goal is to find a threshold near $s^* = \arg \min_s R(s)$. This learning problem over the discrete set of threshold classifiers at thresholds $f(x_1), ..., f(x_N)$ can be exactly represented in the framework of Jamieson & Jain (2022), leading to Algorithm 1 and the following bound.

Theorem 4.1 (Jamieson & Jain (2022)). During iteration r of Algorithm 1, given the classifier model f, with probability at least $1 - \delta$, both $R(\mu) - R(s^*)$ and $R(\bar{\mu}) - R(s^*)$ are upper bounded by

$$c_0 R(s^*) + c_1 \beta_t \sqrt{R(s^*)} + c_2 \beta_t^2.$$
 (2)

270 for all the confidences intervals $[\mu, \bar{\mu}]$ updated at time $t \in 2^{N^+}$. Here, c_0, c_1 and c_2 are some 271 universal constants. 272

273 The theorem suggests that the gap in the true risk of the confidence intervals shrinks roughly on the scale of $\frac{1}{\sqrt{t}}$ over time since $\beta_t = \widetilde{O}(\frac{1}{\sqrt{t}})$. For large t, this gap goes to 0. 274

275 **Definition 4.2** (Score Re-Ordering). Let π denote the permutation of $\{1, \ldots, N\}$ such that 276 $f(x_{\pi(1)}) \leq \cdots \leq f(x_{\pi(N)}).$ 277

Definition 4.3 (Discrete Smoothness). Let $L = \max_{j \in [N-1]} |\eta_{\pi(j)} - \eta_{\pi(j+1)}|$ denote the maximum 278 change in probabilities $\eta_{(i)}$. This mirrors the Lipschitz smoothness in a discrete fashion. Note we 279 always have L < 1. 280

281 Without loss of generality, assume class 0 to be the minority class. The expected imbalance ratio is 282 then $\frac{\sum_{j=1}^{N} \eta_j}{\sum_{j=1}^{N} 1 - \eta_j} < 1$. When $s^* = 0$, it means even the lowest sigmoid score snippets are more likely 283 284 to be in the majority class. In this case, a reasonable strategy is simply query the lowest sigmoid 285 score snippets, which is exactly what our algorithm does.

286 **Theorem 4.4** (Balancedness of Labeled Snippets). Assume class 0 is the minority class and $s^* \neq 0$. 287 Consider an interval of scores $[\mu, \bar{\mu}]$ with $s^* \in [\mu, \bar{\mu}]$. Let the corresponding gaps in risk denoted 288 by $\gamma_0 := R(\mu) - R(s^*) > 0$ and $\gamma_1 := R(\bar{\mu}) - R(s^*) > 0$. When labeling snippets indexed within 289 this interval uniformly at random, the imbalance ratio $\lambda(\mu, \bar{\mu})$ between the minority class and the 290 majority class must satisfy 291

$$\lambda(\underline{\mu}, \bar{\mu}) = \frac{\sum_{j:f(x_j)\in[\underline{\mu}, \bar{\mu}]} \eta_j}{\sum_{j:f(x_j)\in[\underline{\mu}, \bar{\mu}]} 1 - \eta_j} \ge 1 - \min(\frac{N\bar{\gamma} + LN}{1.5 - 2L}, \sqrt{L} \cdot \frac{N\bar{\gamma} + LN + 1}{(1 - L)\sqrt{N\underline{\gamma}}})$$

where $\gamma := \min(\gamma_0, \gamma_1)$ and $\bar{\gamma} := \max(\gamma_0, \gamma_1)$, with $L, \underline{\gamma}, \bar{\gamma} < 1$. This implies,

(a) if $\bar{\gamma} \to 0$ and $L \to \lambda(\mu, \bar{\mu}) \ge 1 - \frac{N\bar{\gamma} + LN}{1.5 - 2L} \to 1$ (perfect balance);

(b) if $\underline{\gamma} \geq \frac{c}{N}$ for some constants c > 0, then $\lambda(\underline{\mu}, \overline{\mu}) \geq 1 - \sqrt{L} \cdot \frac{N \overline{\gamma} + LN + 1}{c(1-L)}$. When $L \to 0$, we again recovers $\lambda(\mu, \bar{\mu}) \rightarrow 1$ (perfect balance). 300

301 *Proof sketch.* Let $A = \sum_{j:f(x_j)\in[\underline{\mu},\overline{\mu}]} \eta_j$ and $B = \sum_{j:f(x_j)\in[\underline{\mu},\overline{\mu}]} 1 - 2\eta_j$, we can rewrite the imbalance ratio into $\lambda(\underline{\mu},\overline{\mu}) = \frac{A}{A+B}$. Since $N(\gamma_1 - \gamma_0) = \sum_{j:f(x_j)\in(\underline{\mu},\overline{\mu}]} 1 - 2\eta_j \ge B - LN$, 302 303 we can lower bound the balancedness by $\frac{A}{A+N(\gamma_1-\gamma_0)+LN}$. We therefore need to prove that A is 304 305 sufficiently large as compared to $N(\gamma_1 - \gamma_0) + LN$. 306

To prove this, we note that the η values around the optimal separation threshold s^* are close to .5. 307 Therefore, by the smoothness condition, when L is small, $\eta_j \approx .5$ for all $f(x_j) \in [\mu, \bar{\mu}]$, so A roughly 308 scales linearly in the number of elements within $[\underline{\mu}, \overline{\mu}]$. We can prove that there are at least $O(\sqrt{\frac{N\bar{\gamma}}{L}})$ 309 310 elements in this confidence interval. A therefore follows a similar scale, which is much greater than 311 $N(\gamma_1 - \gamma_0) + LN$ when $L \to 0$. Under such case, we recover a balancedness bound of 1, i.e. the annotated snippets are perfectly balanced. Our detailed proof is provided in Appendix B. 312

313 314

315

292 293

295 296

297

298 299

5 **EXPERIMENTS**

316 In this section, we present our experiments conducted for 317 the five highly customized data filters in Table 1: poli-318 tics, climate, AI, mainstream knowledge and text qual-319 ity filters. These filters are applied to the OpenWebText 320 dataset (Gokaslan & Cohen, 2019), which is divided into 321 around 13.5M snippets of 1024 tokens. The first three filters identify text snippets related to a particular topic, with 322 highly detailed specifications of filtering prompts to GPT-323 40 about the subdomain of topics that should be included.

Filter	Imbalance Ratio λ
Politics	0.153
Climate	0.043
AI	0.026
Mainstream	0.208
Quality	0.457

Figure 3: Imbalance ratio of minority vs majority decisions, calculated based on 5000 randomly sampled snippets. $\lambda =$ #Minority/#Majority.

Filter	Method	Bal. Accuracy (GPT-40 as GT)	Human Preference Over GPT-40	#Queries to GPT-40	Lightweight Model Cost	Total Cost
Politics	GPT-40 SIEVE (Ours)	95.6%* 95.6%	_	13.5M 60K	\$0 \$270	\$67,000 \$570
Climate	GPT-40 SIEVE (Ours)	96.6%* 96.7%	_	13.5M 7.5K	\$0 \$180	\$67,000 \$220
AI	GPT-40 SIEVE (Ours)	95.5%* 95.7%	_	13.5M 6K	\$0 \$180	\$67,000 \$210
Mainstream	GPT-40 SIEVE (Ours)	92.4%* 91.0%	50% 54%	13.5M 100K	\$0 \$400	\$67,000 \$900
Quality	GPT-40 SIEVE (Ours)	88.2%* 86.3%	50% 53%	13.5M 60K	\$0 \$270	\$67,000 \$570
	SIEVE (Ours)	80.3%	53%	60K	\$270	2 0

Table 1: Performance results of applying SIEVE on five highly specialized filters (see Appendix A). SIEVE can match or exceed GPT-4o's quality in terms of balanced accuracy and human preference. On the other hand, SIEVE saves more than 99% of the cost compared to using GPT-4o. See Section 5 for experiment details. *We assess GPT-4o's accuracy by measuring output consistency between identical API calls using greedy decoding (temperature set to 0). Some inconsistency persists, possibly due to hardware non-determinism. This may be amplified when using our CoT prompts.

The mainstream knowledge prompt aims to exclude any obscure and niche content determined by
 GPT-40. Lastly, the quality filter aims to identify text snippets that are considered high quality by
 GPT-40. Detailed prompts can be found in Appendix A.

Throughout our experiments, we use the encoder part of pretrained T5-large (Raffel et al., 2020) as
the lightweight model. A linear layer is attached to the encoder model for binary classification. The
classification model has less than 770M parameters, orders of magnitude smaller than GPT-40. To
mitigate the imbalanced nature of the snippets shown in Table 3, we utilize the focal loss (Lin, 2017).
We refer the readers to Appendix C for more training details.

352 353

354

337

338

339

340

341

342 343

5.1 RESULTS: GPT-40-BASED AND HUMAN-BASED EVALUATION

355 **GPT-4o-Based Evaluation.** Table 1 demonstrates that SIEVE achieves comparable balanced 356 accuracy to GPT-40 on politics, climate, and AI filters, while closely matching its performance on 357 mainstream knowledge and text quality filters. We evaluated these filters using a set of test snippets 358 randomly sampled from OpenWebText. Ground truth labels were established through individual GPT-359 40 API calls for each snippet. To assess GPT-40's performance against this ground truth, we conducted additional API calls using identical prompts for each test snippet, effectively measuring the model's 360 self-consistency rate. Surprisingly, we observed inherent noise in GPT-40's decisions, even when 361 using greedy decoding (temperature set to 0). This variability likely stems from non-deterministic 362 factors in the hardware infrastructure. While we employed chain-of-thought (CoT) reasoning in our 363 filtering prompts (detailed in Appendix A) to enhance decision quality, the compounding noise from 364 each generated token appears to have contributed to the noticeable inconsistencies. We also note that CoT plays an increasing role in inference-time scaling as demonstrated by OpenAI's o1 model, 366 which may inevitably cause inconsistencies in the models' decisions.

367

368 **Human Evaluation.** In the above, we compared our distilled lightweight model's accuracy to GPT-369 40 for both mainstream and quality filters. When evaluated by GPT-40 itself, our model's performance 370 appeared lower, raising the question: Is our lightweight model actually worse, or is GPT-40 biased 371 when judging its own decisions? To investigate, we conducted a human evaluation. For each filter, 372 we randomly selected 100 text snippets where GPT-40 and our lightweight model disagreed. We 373 then recruited two groups of 13 annotators per filter to manually assess these challenging cases. The 374 ground truth for each snippet was determined by the majority vote of the annotators, allowing us 375 to compare both models' performance against human judgment. As shown in the fourth column of Table 1, we see even a slight edge of our lightweight model over GPT-40 when judged by human. 376 This suggests our lightweight model is at least comparable to the decisions made by GPT-40, while 377 incurring much less computational cost when filtering web-scale datasets.



Figure 4: Active vs Random Distillation: Performance of the distilled lightweight model across different number of queries made to GPT-40. With active learning, we can save the number of queries to GPT-40 by more than 5x for the politics filter and more than 3x for the climate filter.

Filter	Method	Bal. Accuracy (GPT-40 as GT)	#Queries to GPT-40	GPT-40 Cost	Lightweight Model Training	Lightweight Model Inference	Total Cost
Politics	GPT-40	95.6%*	13.5M	\$67,000	\$0	\$0	\$67,000
	SIEVE (Ours)	95.6%	60K	\$300	\$120	\$150	\$570
	SIEVE (Ours)	95%	25K	\$125	\$30	\$150	\$305
	Random	95%	100K	\$500	\$20	\$150	\$670
Climate	GPT-40	96.6%*	13.5M	\$67,000	\$0	\$0	\$67,000
	SIEVE (Ours)	96.7%	7.5K	\$40	\$30	\$150	\$220
	Random	96.6%	25K	\$125	\$20	\$150	\$315

Table 2: Cost breakdown of politics and climate filters. Total cost is consisted of GPT-40 querying cost, lightweight model training cost and T5 inference cost.

405 5.2 ACTIVE VS RANDOM DISTILLATION

406 In this section, we compare the effectiveness of active learning versus random querying for distilling 407 our lightweight model. The random distillation strategy involves querying GPT-4 on a predetermined 408 number of randomly selected snippets from the OpenWebText dataset, followed by fine-tuning 409 the T5 encoder model on this queried data. As illustrated in Figure 4, our comparison of the 410 lightweight model's performance when trained with active versus random sampling for both politics 411 and climate filters reveals significant advantages for active distillation. Specifically, active distillation 412 demonstrates remarkable efficiency, requiring over 5 times fewer queries for the politics filter and 413 more than 3 times fewer for the climate filter compared to random distillation. Figure 4c also 414 demonstrates our active learning algorithm's effectiveness in querying a much more balanced set of snippets. It's noteworthy that even with 100,000 queries, random distillation fails to match GPT-40's 415 performance. Due to budgetary constraints, we did not extend random distillation experiments beyond 416 this query count. Importantly, we observe that as accuracy increases, active learning typically yields 417 greater query budget savings compared to random sampling (Citovsky et al., 2021; Ash et al., 2021; 418 Zhang et al., 2024a), suggesting that for the politics filter, in particular, the efficiency gains from 419 active distillation could potentially exceed the observed 5-fold reduction in query costs. As inference 420 cost of state-of-art LLMs increases, active distillation could play an increasingly important role in 421 SIEVE.

422 423 424

389

390

391 392 393

402

403 404

5.3 COMPUTATIONAL COST COMPARISONS

Cost Breakdown and Comparison. Table 2 presents a comprehensive breakdown of SIEVE's computational costs, consisting of GPT-40 query costs, lightweight model training costs, and inference costs for filtering the entire OpenWebText dataset. The sum of GPT-40 query and training costs represents the total model distillation expense for SIEVE. Our experiments indicate that querying GPT-40 for 1000 snippets costs approximately \$5. The lightweight model inference cost, a conservative estimate for processing 13.5M snippets using a 770M parameter model, is expected to be lower in practice by parallelization across multiple CPUs or cheap inference GPUs. The lightweight model training costs includes costs for model fine-tuning at each iteration of Algorithm 1, inference costs

for computing sigmoid scores of data in stream, and negligible uncertainty update expenses. We
 calculated these costs based on the hourly rate for 8xA100 80GB GPUs, multiplied by the actual
 time spent on each experiment. Together, these components provide a detailed overview of the
 computational resources required for SIEVE's implementation and application to the entire dataset.

As shown in Table 2, active distillation offers a more cost-effective approach compared to random distillation in achieving equivalent model performance. This cost advantage becomes even more pronounced when considering more expensive teacher models, such as the recently released o1, which significantly increases query costs and potentially dominates the total expenses. Given that active distillation substantially reduces the number of required queries, the cost disparity between active and random distillation methods is expected to widen further, especially when utilizing more advanced and costly teacher models.

443

444 Stream-Based vs Pool-Based Active Distillation. The choice of a stream-based active learning approach for SIEVE is justified by its significant cost advantages over pool-based methods. While 445 stream-based algorithms process snippets only once, pool-based active learning requires repeated 446 forward inference on the entire dataset of 13.5M snippets for each batch of queried snippets. This 447 difference translates to substantial additional costs: approximately \$1800 for the politics filter and 448 \$750 for the climate filter. These extra expenses would dramatically increase the overall cost of 449 active distillation. Our decision to develop a stream-based active learning algorithm for SIEVE, 450 particularly effective in imbalanced scenarios, is thus strongly supported by these cost considerations, 451 ensuring a more economically viable solution for our system.

452 453 454

6 RELATED WORK

455 **Data Filtering for Large Language Models.** Data curation is fundamental to the development 456 of LLMs (Longpre et al., 2023; Zhou et al., 2023). As interest in domain-specific LLMs grows, the 457 need for extensive, relevant data collection becomes increasingly important. For a comprehensive 458 overview of existing datasets, we direct readers to Raffel et al. (2020); Gao et al. (2020); Liu et al. 459 (2024). Current methods for acquiring domain-specific data predominantly rely on a few established 460 large-scale databases, including textbooks (Gunasekar et al., 2023), code repositories (Muennighoff 461 et al., 2023; Gao et al., 2020), medical literature, and other specialized sources (Gao et al., 2020; 462 Cheng et al., 2023). However, for rapidly evolving topics like the 2024 presidential election, climate 463 change and artificial intelligence, relevant information is often dispersed across the internet, making comprehensive data collection challenging. Our approach involves training lightweight, task-specific 464 data filtering models distilled from GPT-4. These models are then applied to web-scale datasets to 465 identify pertinent information across various domains. Existing data filtering techniques span a wide 466 range, from basic rule-based methods utilizing sentence-level statistical features (Rae et al., 2021; 467 Yang, 2019; Laurençon et al., 2022; Zhang et al., 2022b) to sophisticated filters leveraging pretrained 468 neural networks for text quality (Brown, 2020; Du et al., 2022; Chowdhery et al., 2023; Touvron 469 et al., 2023; Enomoto et al., 2024; Qian et al., 2024) and toxicity (Lees et al., 2022; Friedl, 2023) 470 filtering. To our knowledge, this work represents the first attempt to develop domain-specific data 471 filtering models adaptable to a diverse array of filtering requirements and specialized domains.

472

473 Active Learning Active learning is a strategy aimed at reducing data annotation costs by selectively 474 choosing which examples to label. Traditional approaches iteratively update machine learning models 475 based on newly labeled data, using various informativeness metrics to guide the selection process. 476 These metrics typically include uncertainty (Lewis & Gale, 1994; Tong & Koller, 2001; Settles, 2009; Balcan et al., 2006; Kremer et al., 2014; Gal et al., 2017; Ducoffe & Precioso, 2018; Beluch et al., 477 2018), diversity (Sener & Savarese, 2017; Geifman & El-Yaniv, 2017; Citovsky et al., 2021), and 478 expected model change (Ash et al., 2019; 2021; Wang et al., 2021; Elenter et al., 2022; Mohamadi 479 et al., 2022). However, most existing methods are designed for pool-based settings with balanced 480 data distributions, which may not be suitable for all real-world scenarios. 481

In contrast, our work focuses on stream-based algorithms capable of handling data imbalance, an
area that has received limited attention in deep learning contexts. While pool-based algorithms have
shown success in addressing class imbalance (Aggarwal et al., 2020; Kothawade et al., 2021; Emam
et al., 2021; Zhang et al., 2022a; Coleman et al., 2022; Jin et al., 2022; Cai, 2022; Nuggehalli et al.,
2023; Zhang et al., 2024b; Lesci & Vlachos, 2024; Fairstein et al., 2024), stream-based approaches

486 for deep neural networks remain understudied. The recent work by Saran et al. (2023) introduces an 487 online volume sampling technique for annotating diverse sets of examples in the representation space. 488 Representation diversity, however, has been shown to struggle on class-imbalanced data distribution 489 under the pool based setting (Zhang et al., 2022a; Nuggehalli et al., 2023; Lesci & Vlachos, 2024; 490 Fairstein et al., 2024). Specifically these studies suggest that sampling diversely in the representation space does not necessarily improve the class-balancedness of the annotated examples. To address this 491 gap, we propose what we believe to be the first stream-based algorithm specifically designed for class 492 imbalance scenarios in active learning. 493

494 As part of our algorithm, we use an agnostic active learning algorithm for identifying the TRM 495 threshold. Agnostic active learning has been widely studied in the classical PAC learning setups, 496 where the labels of any particular example is inherently noisy. Our procedure is a direct application of Jamieson & Jain (2022), which was inspired by Dasgupta et al. (2007). The algorithm is proven to 497 be near minimax optimal in these literature. In addition, our algorithm can also be seen as an instance 498 of the algorithm proposed by Katz-Samuels et al. (2021) for threshold classifiers, where they also 499 prove such algorithm is near instance-optimal. In this paper, we also prove the first bound towards 500 balancedness of labeled examples in agnostic active learning, focusing on the class of threshold 501 classifiers. 502

Knowledge Distillation Knowledge distillation is a technique where a smaller "student" model 504 learns to emulate a larger, more sophisticated "teacher" model. With the increasing capabilities of 505 large language models (LLMs) across diverse tasks, recent research has explored using LLMs as 506 annotators to train domain-specific student models. For a comprehensive review, see Tan et al. (2024). 507 While most research in LLM knowledge distillation focuses on knowledge extraction methods, few 508 studies address the high computational cost of using LLMs for large-scale annotation. Recent work 509 by Zhang et al. (2023) and Rouzegar & Makrehchi (2024) has begun to tackle this issue by using active learning. Our paper applies knowledge distillation to the specific problem of data filtering. We 510 employ a straightforward approach, using LLM annotations as binary classification labels to actively 511 fine-tune an encoder-based model. Future research could explore using GPT-4's chain-of-thought 512 outputs to distill a decoder-based student model within the multi-task learning framework proposed 513 by Hsieh et al. (2023). It's worth noting that classic knowledge distillation work, such as Hinton 514 (2015), trains student classifiers to match the teacher models' output logits. However, in our case, 515 using chain-of-thought filtering prompts makes it impractical to obtain probabilities for the binary 516 decision.

- 517 518
- 519

7 CONCLUSION, LIMITATIONS AND FUTURE WORK

In this paper, we introduced SIEVE, demonstrating the feasibility of achieving GPT-40 quality data filtering across a diverse range of user-specified filtering prompts. Our comprehensive study showcases the effectiveness of SIEVE in curating large-scale, high-quality datasets for language model training at a fraction of the cost of existing techniques. The experimental results, validated on the OpenWebText using five highly customized filter tasks, provide strong evidence of SIEVE's capability to match GPT-40's accuracy while significantly reducing computational expenses.

While our initial study presents promising results, we acknowledge that there are several avenues for 527 future enhancement and exploration. For future work, we are particular excited in scaling SIEVE to 528 even larger datasets like the PILE, and more data modalities beyond text. The modular nature of 529 SIEVE also allows for the integration of more advanced active learning algorithms. Additionally, 530 SIEVE serves as an excellent testbed for these algorithms, offering immediate real-world impact. 531 Future work could also investigate the incorporation of semi-supervised learning techniques to further 532 reduce annotation costs following the framework proposed by Zhang et al. (2024a). Moreover, 533 while our current implementation focuses on T5 architectures, future research could examine the 534 efficacy of SIEVE with a broader range of pretrained model architectures for transfer learning. Lastly, 535 exploring the use of more powerful models beyond GPT-40, such as 01, for handling complex filtering prompts could extend the capabilities of SIEVE to even more challenging scenarios. In such cases, 536 the importance of active learning becomes even more pronounced, as the increased querying costs 537 associated with these advanced models necessitate highly efficient sampling strategies. 538

540 REFERENCES

546

547

548

552

553

554

560

570

- 542 Umang Aggarwal, Adrian Popescu, and Céline Hudelot. Active learning for imbalanced datasets.
 543 In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1428–1437, 2020.
 - Jordan Ash, Surbhi Goel, Akshay Krishnamurthy, and Sham Kakade. Gone fishing: Neural active learning with fisher embeddings. *Advances in Neural Information Processing Systems*, 34:8927–8939, 2021.
- Jordan T Ash, Chicheng Zhang, Akshay Krishnamurthy, John Langford, and Alekh Agarwal. Deep
 batch active learning by diverse, uncertain gradient lower bounds. *arXiv preprint arXiv:1906.03671*, 2019.
 - Maria-Florina Balcan, Alina Beygelzimer, and John Langford. Agnostic active learning. In *Proceedings of the 23rd international conference on Machine learning*, pp. 65–72, 2006.
- William H Beluch, Tim Genewein, Andreas Nürnberger, and Jan M Köhler. The power of ensembles for active learning in image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 9368–9377, 2018.
- ⁵⁵⁹ Tom B Brown. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165*, 2020.
- Xinmeng Cai. Active learning for imbalanced data: The difficulty and proportions of class matter.
 Wireless Communications and Mobile Computing, 2022, 2022.
- Daixuan Cheng, Shaohan Huang, and Furu Wei. Adapting large language models via reading comprehension. In *The Twelfth International Conference on Learning Representations*, 2023.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.
- Gui Citovsky, Giulia DeSalvo, Claudio Gentile, Lazaros Karydas, Anand Rajagopalan, Afshin
 Rostamizadeh, and Sanjiv Kumar. Batch active learning at scale. *Advances in Neural Information Processing Systems*, 34:11933–11944, 2021.
- 574
 575
 576
 576
 576
 577
 578
 Cody Coleman, Edward Chou, Julian Katz-Samuels, Sean Culatana, Peter Bailis, Alexander C Berg, Robert Nowak, Roshan Sumbaly, Matei Zaharia, and I Zeki Yalniz. Similarity search for efficient active learning and search of rare concepts. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6402–6410, 2022.
- Sanjoy Dasgupta, Daniel J Hsu, and Claire Monteleoni. A general agnostic active learning algorithm.
 Advances in neural information processing systems, 20, 2007.
- Nan Du, Yanping Huang, Andrew M Dai, Simon Tong, Dmitry Lepikhin, Yuanzhong Xu, Maxim
 Krikun, Yanqi Zhou, Adams Wei Yu, Orhan Firat, et al. Glam: Efficient scaling of language
 models with mixture-of-experts. In *International Conference on Machine Learning*, pp. 5547–5569.
 PMLR, 2022.
- Melanie Ducoffe and Frederic Precioso. Adversarial active learning for deep networks: a margin based approach. *arXiv preprint arXiv:1802.09841*, 2018.
- Juan Elenter, Navid NaderiAlizadeh, and Alejandro Ribeiro. A lagrangian duality approach to active
 learning. *arXiv preprint arXiv:2202.04108*, 2022.
- 591

581

 Zeyad Ali Sami Emam, Hong-Min Chu, Ping-Yeh Chiang, Wojciech Czaja, Richard Leapman,
 Micah Goldblum, and Tom Goldstein. Active learning at the imagenet scale. *arXiv preprint arXiv:2111.12880*, 2021.

594 595 596 597 598	Rintaro Enomoto, Arseny Tolmachev, Takuro Niitsuma, Shuhei Kurita, and Daisuke Kawahara. Inves- tigating web corpus filtering methods for language model development in japanese. In <i>Proceedings</i> of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 4: Student Research Workshop), pp. 154–160, 2024.
599 600 601 602	Yaron Fairstein, Oren Kalinsky, Zohar Karnin, Guy Kushilevitz, Alexander Libov, and Sofia Tolmach. Class balancing for efficient active learning in imbalanced datasets. In <i>Proceedings of The 18th</i> <i>Linguistic Annotation Workshop (LAW-XVIII)</i> , pp. 77–86, 2024.
603 604	Paul Friedl. Dis/similarities in the design and development of legal and algorithmic normative systems: the case of perspective api. <i>Law, Innovation and Technology</i> , 15(1):25–59, 2023.
605 606	Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In <i>International Conference on Machine Learning</i> , pp. 1183–1192. PMLR, 2017.
608 609 610	Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. <i>arXiv preprint arXiv:2101.00027</i> , 2020.
611 612	Yonatan Geifman and Ran El-Yaniv. Deep active learning over the long tail. <i>arXiv preprint arXiv:1711.00941</i> , 2017.
613 614 615	Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007.github.io/ OpenWebTextCorpus, 2019.
616 617 618	Suriya Gunasekar, Yi Zhang, Jyoti Aneja, Caio César Teodoro Mendes, Allie Del Giorno, Sivakanth Gopi, Mojan Javaheripi, Piero Kauffmann, Gustavo de Rosa, Olli Saarikivi, et al. Textbooks are all you need. <i>arXiv preprint arXiv:2306.11644</i> , 2023.
619 620 621 622 623	Aman Gupta, Anup Shirgaonkar, Angels de Luis Balaguer, Bruno Silva, Daniel Holstein, Dawei Li, Jennifer Marsman, Leonardo O Nunes, Mahsa Rouzbahman, Morris Sharp, et al. Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture. <i>arXiv preprint arXiv:2401.08406</i> , 2024.
624 625	Geoffrey Hinton. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> , 2015.
626 627 628 629	Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh, Hootan Nakhost, Yasuhisa Fujii, Alexander Ratner, Ranjay Krishna, Chen-Yu Lee, and Tomas Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. <i>arXiv preprint arXiv:2305.02301</i> , 2023.
630 631	Kevin Jamieson and Lalit Jain. Interactive machine learning. 2022.
632 633	Qiuye Jin, Mingzhi Yuan, Haoran Wang, Manning Wang, and Zhijian Song. Deep active learning models for imbalanced image classification. <i>Knowledge-Based Systems</i> , 257:109817, 2022.
635 636 637	Julian Katz-Samuels, Jifan Zhang, Lalit Jain, and Kevin Jamieson. Improved algorithms for agnostic pool-based active classification. In <i>International Conference on Machine Learning</i> , pp. 5334–5344. PMLR, 2021.
638 639 640	Suraj Kothawade, Nathan Beck, Krishnateja Killamsetty, and Rishabh Iyer. Similar: Submodular information measures based active learning in realistic scenarios. <i>Advances in Neural Information Processing Systems</i> , 34:18685–18697, 2021.
641 642 643 644	Jan Kremer, Kim Steenstrup Pedersen, and Christian Igel. Active learning with support vector machines. <i>Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery</i> , 4(4):313–326, 2014.
645 646 647	Hugo Laurençon, Lucile Saulnier, Thomas Wang, Christopher Akiki, Albert Villanova del Moral, Teven Le Scao, Leandro Von Werra, Chenghao Mou, Eduardo González Ponferrada, Huu Nguyen, et al. The bigscience roots corpus: A 1.6 tb composite multilingual dataset. <i>Advances in Neural</i> <i>Information Processing Systems</i> , 35:31809–31826, 2022.

648 649	Peter Lee, Sebastien Bubeck, and Joseph Petro. Benefits, limits, and risks of gpt-4 as an ai chatbot for medicine. <i>New England Journal of Medicine</i> , 388(13):1233–1239, 2023.
650 651 652 653 654	Alyssa Lees, Vinh Q Tran, Yi Tay, Jeffrey Sorensen, Jai Gupta, Donald Metzler, and Lucy Vasserman. A new generation of perspective api: Efficient multilingual character-level transformers. In <i>Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining</i> , pp. 3197–3207, 2022.
655 656	Pietro Lesci and Andreas Vlachos. Anchoral: Computationally efficient active learning for large and imbalanced datasets. <i>arXiv preprint arXiv:2404.05623</i> , 2024.
657 658 659 660	DD Lewis and WA Gale. A sequential algorithmfor training text classifiers. In SIGIR'94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University, pp. 3–12, 1994.
661 662	Beibin Li, Yi Zhang, Sébastien Bubeck, Jeevan Pathuri, and Ishai Menache. Small language models for application interactions: A case study. <i>arXiv preprint arXiv:2405.20347</i> , 2024.
663 664	T Lin. Focal loss for dense object detection. arXiv preprint arXiv:1708.02002, 2017.
665 666	Yang Liu, Jiahuan Cao, Chongyu Liu, Kai Ding, and Lianwen Jin. Datasets for large language models: A comprehensive survey. <i>arXiv preprint arXiv:2402.18041</i> , 2024.
668 669 670 671	Shayne Longpre, Gregory Yauney, Emily Reif, Katherine Lee, Adam Roberts, Barret Zoph, Denny Zhou, Jason Wei, Kevin Robinson, David Mimno, et al. A pretrainer's guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity. <i>arXiv preprint arXiv:2305.13169</i> , 2023.
672 673	Mohamad Amin Mohamadi, Wonho Bae, and Danica J Sutherland. Making look-ahead active learning strategies feasible with neural tangent kernels. <i>arXiv preprint arXiv:2206.12569</i> , 2022.
674 675 676 677	Niklas Muennighoff, Qian Liu, Armel Zebaze, Qinkai Zheng, Binyuan Hui, Terry Yue Zhuo, Swayam Singh, Xiangru Tang, Leandro Von Werra, and Shayne Longpre. Octopack: Instruction tuning code large language models. <i>arXiv preprint arXiv:2308.07124</i> , 2023.
678 679	Shyam Nuggehalli, Jifan Zhang, Lalit Jain, and Robert Nowak. Direct: Deep active learning under imbalance and label noise. <i>arXiv preprint arXiv:2312.09196</i> , 2023.
680 681 682	Crystal Qian, Emily Reif, and Minsuk Kahng. Understanding the dataset practitioners behind large language model development. <i>arXiv preprint arXiv:2402.16611</i> , 2024.
683 684 685	Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. <i>arXiv preprint arXiv:2112.11446</i> , 2021.
686 687 688 689	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of machine learning research</i> , 21(140):1–67, 2020.
690 691	Hamidreza Rouzegar and Masoud Makrehchi. Enhancing text classification through llm-driven active learning and human annotation. <i>arXiv preprint arXiv:2406.12114</i> , 2024.
692 693 694 695	Akanksha Saran, Safoora Yousefi, Akshay Krishnamurthy, John Langford, and Jordan T Ash. Stream- ing active learning with deep neural networks. In <i>International Conference on Machine Learning</i> , pp. 30005–30021. PMLR, 2023.
696 697	Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. <i>arXiv preprint arXiv:1708.00489</i> , 2017.
698 699	Burr Settles. Active learning literature survey. 2009.
700	Zhen Tan, Alimohammad Beigi, Song Wang, Ruocheng Guo, Amrita Bhattacharjee, Bohan Jiang, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. Large language models for data

Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. Large language models for data annotation: A survey. *arXiv preprint arXiv:2402.13446*, 2024.

702 703 704	Simon Tong and Daphne Koller. Support vector machine active learning with applications to text classification. <i>Journal of machine learning research</i> , 2(Nov):45–66, 2001.
705 706 707	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
708 709 710	Haonan Wang, Wei Huang, Andrew Margenot, Hanghang Tong, and Jingrui He. Deep active learning by leveraging training dynamics. <i>arXiv preprint arXiv:2110.08611</i> , 2021.
710 711 712	Zhilin Yang. Xlnet: Generalized autoregressive pretraining for language understanding. <i>arXiv</i> preprint arXiv:1906.08237, 2019.
713 714 715	Jifan Zhang, Julian Katz-Samuels, and Robert Nowak. Galaxy: Graph-based active learning at the extreme. <i>arXiv preprint arXiv:2202.01402</i> , 2022a.
716 717 718 719	Jifan Zhang, Yifang Chen, Gregory Canal, Arnav Mohanty Das, Gantavya Bhatt, Stephen Mussmann, Yinglun Zhu, Jeff Bilmes, Simon Shaolei Du, Kevin Jamieson, et al. Labelbench: A comprehensive framework for benchmarking adaptive label-efficient learning. <i>Journal of Data-centric Machine Learning Research</i> , 2024a.
720 721	Jifan Zhang, Shuai Shao, Saurabh Verma, and Robert Nowak. Algorithm selection for deep active learning with imbalanced datasets. <i>Advances in Neural Information Processing Systems</i> , 36, 2024b.
722 723 724	Ruoyu Zhang, Yanzeng Li, Yongliang Ma, Ming Zhou, and Lei Zou. Llmaaa: Making large language models as active annotators. <i>arXiv preprint arXiv:2310.19596</i> , 2023.
725 726 727	Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. <i>arXiv preprint arXiv:2205.01068</i> , 2022b.
728 729 730 731	Tong Zhou, Yubo Chen, Pengfei Cao, Kang Liu, Jun Zhao, and Shengping Liu. Oasis: Data curation and assessment system for pretraining of large language models. <i>arXiv preprint arXiv:2311.12537</i> , 2023.
732 733	
734 735 736	
737 738	
739 740	
741 742 743	
744 745	
746 747	
748 749 750	
751 752	
753 754	
())	

756 FILTERING PROMPTS А

A.1 POLITICS

Please analyze the following text snippet and determine if it is relevant to aspects of a presidential election. The snippet may be arbitrarily cut off from a longer article, so please ignore any oddities caused by the truncation and focus on the overall relevance to presidential elections.

763 764 765

769

770

771

772

773

774

775

776

777

779

781

758 759

760 761

762

1. Read the text snippet carefully.

- 766 2. Identify any key terms, phrases, or concepts related to presidential election. These include by are 767 not limited to Candidate information (biographies, backgrounds, policy positions) 768
 - Economic policies and their potential impacts
 - Social issues and proposed solutions
 - Foreign policy stances and international relations
 - Campaign events, debates, and public appearances
 - · Polling data, electoral projections, and voter demographics
 - Media coverage, endorsements, and fact-checking
 - Campaign finance and fundraising efforts
 - Party dynamics and internal politics
 - · Electoral processes, including voting systems and potential reforms
 - Controversies or scandals involving candidates or their campaigns
 - Vice presidential candidates and potential cabinet members
- 778 Analysis of key battleground states or regions
 - Digital campaigning strategies and social media presence
 - · Grassroots organizing and volunteer efforts
 - External events or crises that may influence the election
 - 3. Ignore any abrupt beginning or ending of the snippet and focus on the main content.
- 782 4. Assign either "PASS" for content relevant to presidential election and "FAIL" for those that are 783 irrelevant. 784
- 785 Based on your analysis, determine if the snippet is relevant or irrelevant to the general knowledge of 786 presidential elections. Think step by step and provide your final answer as either "PASS" or "FAIL" 787 at the end of your response and nothing else. 788
- Text snippet: <Insert Text Snippet> 789

790

791 792

793 794

796

797

799

800

801

A.2 CLIMATE

Please analyze the following text snippet and determine if it is relevant to the general knowledge of climate change. The snippet may be arbitrarily cut off from a longer article, so please ignore any oddities caused by the truncation and focus on the overall relevance to climate change.

- 1. Read the text snippet carefully.
- 798 2. Identify any key terms, phrases, or concepts related to climate change, such as global warming, carbon emission, energy, technology innovation, agriculture, natural resources, pollution, landfill, chemistry, rising sea levels, extreme weather events, climate policies, environmental justice and sustainability.
- 802 3. Assess whether the snippet discusses causes, effects, or solutions to climate change, or provides 803 information that contributes to the understanding of climate change.
- 804 4. Ignore any abrupt beginning or ending of the snippet and focus on the main content.
- 805 5. Assign either "PASS" for content relevant to climate change and "FAIL" for those that are irrelevant.

806 Based on your analysis, determine if the snippet is relevant or irrelevant to the general knowledge of 807 climate change. Think step by step and provide your final answer as either "PASS" or "FAIL" at the 808 end of your response and nothing else. 809

Text snippet: <Insert Text Snippet>

810 A.3 AI

811 812

Please analyze the following text snippet and determine if it is relevant to aspects of AI. The snippet 813 may be arbitrarily cut off from a longer article, so please ignore any oddities caused by the truncation 814 and focus on the overall relevance to artificial intelligence.

- 815 1. Read the text snippet carefully.
- 816 2. Find content related to artificial intelligence that discusses computer systems or software designed 817 to perform tasks typically requiring human intelligence, such as visual perception, speech recogni-818 tion, decision-making, and language translation. Look for explanations of technologies enabling 819 machines to learn from experience, adjust to new inputs, and perform human-like tasks without explicit programming. Include descriptions of systems that can interpret visual information from 820 the world, as well as content about software capable of processing, analyzing, generating, or 821 understanding human language. Seek information on machines or programs that can improve 822 their performance through experience or data. Include discussions of AI applications in various 823 fields such as healthcare, finance, transportation, education, or entertainment. Consider content 824 addressing ethical considerations surrounding AI, including bias, privacy, job displacement, or 825 long-term implications of advanced AI systems. Look for historical accounts of AI development, major milestones, breakthroughs, or setbacks in the field. Include explanations of AI algorithms, 827 their workings, strengths, limitations, and potential applications. Capture debates or discussions 828 about the future of AI, including topics like artificial general intelligence, superintelligence, or 829 potential societal impacts of widespread AI adoption. Include reports on current research, new 830 methodologies, experimental results, or theoretical advancements in AI. Consider content about 831 prominent figures, organizations, or companies significantly contributing to AI research and development. Look for discussions of AI policy, regulation, or governance at organizational, national, 832 or international levels. Include explanations of the relationship between AI and other fields such as 833 robotics, Internet of Things, big data, or quantum computing. Finally, capture content addressing 834 challenges in AI development, such as data quality, computational requirements, or the need for 835 explainable AI systems. 836
 - 3. Ignore any abrupt beginning or ending of the snippet and focus on the main content.
- 837 4. Assign either "PASS" for content relevant to AI and "FAIL" for those that are irrelevant. 838

Based on your analysis, determine if the snippet is relevant or irrelevant to aspects of artificial 839 intelligence, its development, applications, or implications. Think step by step and provide your final 840 answer as either "PASS" or "FAIL" at the end of your response and nothing else. 841

- 842 Text snippet: <Insert Text Snippet>
- 843 844 845

851

MAINSTREAM KNOWLEDGE A.4

846 We asked for obscure knowledge instead, so any snippet that "Failed" would be considered 847 mainstream knowledge.

848 Please analyze the following text snippet and determine if it contains obscure or niche knowledge that 849 less than 10000 people know and understand. The snippet may be arbitrarily cut off from a longer 850 article, so please ignore any oddities caused by the truncation and focus on the overall relevance to artificial intelligence. 852

- 1. Read the text snippet carefully.
- 853 2. Filter the dataset for content related to obscure, specialized, or highly niche knowledge that is 854 not commonly known or easily accessible to the general public. Include information on rare 855 historical events, obscure scientific theories, uncommon philosophical concepts, extinct languages, 856 highly specialized mathematics, niche literary works, rare medical conditions, uncommon species, esoteric subcultures, mystical practices, experimental technologies, obscure laws, rare geological 858 formations, lesser-known art movements, specialized crafting techniques, rare musical instruments, 859 uncommon culinary practices, obscure sports, theoretical cosmological concepts, and highly specialized areas of archaeology or anthropology. Exclude any content that is commonly known, part of standard education, regularly discussed in popular media, or widely understood by the 861 general public. The ideal content should require specialized knowledge, extensive research, or 862 access to uncommon sources of information, rather than being something an average person would encounter in daily life or through casual exposure to media and education.

- 3. Ignore any abrupt beginning or ending of the snippet and focus on the main content.
 - 4. Assign either "PASS" for obscure and niche knowledge and "FAIL" for those that are common knowledge.

Think step by step. Then, you must provide your final answer as either "PASS" or "FAIL" at the end of your response and nothing else.

Text snippet: <Insert Text Snippet>

A.5 QUALITY

Please analyze the following text snippet and determine if it is high quality or low quality training data for a large language model. The snippet may be cut off abruptly from a longer piece of text, but focus your analysis on the quality factors present in the provided text rather than the awkward truncation. Quality factors to consider include:

- 1. Evaluate the spelling, grammar, and overall writing quality of the snippet. Note any errors or inconsistencies that could negatively impact the model's learning.
- 2. Assess the factual accuracy and reliability of the information presented in the snippet. Consider whether the content appears trustworthy and well-researched.
- 3. Analyze the clarity, coherence, and logical flow of ideas in the snippet. Determine if the text is easy to understand and follow.
- 4. Gauge the breadth and depth of knowledge conveyed in the snippet. Consider whether the content provides valuable information or insights on the topic at hand.
 - 5. Examine the neutrality and objectivity of the tone and perspective presented in the snippet. Consider if the text appears biased or presents a balanced viewpoint.
 - 6. Based on the above factors, determine if the snippet is: PASS: High quality training data FAIL: Low quality training data
- Think step by step and answer with either PASS or FAIL as your final decision in the end and nothing else.

- Text snippet: <Insert Text Snippet>

ANALYSIS PROOF В

D (a (

Proof. Let $A = \sum_{j:f(x_j) \in [\underline{\mu}, \overline{\mu}]} \eta_j$ and $B = \sum_{j:f(x_j) \in [\underline{\mu}, \overline{\mu}]} 1 - 2\eta_j$, we can rewrite the imbalance ratio into $\lambda(\underline{\mu}, \overline{\mu}) = \frac{A}{A+B}$. Since $N(\gamma_1 - \gamma_0) = \sum_{j:f(x_j) \in (\mu, \overline{\mu}]} 1 - 2\eta_j \ge B - LN$, we can lower bound the balancedness by $\frac{A}{A+N(\gamma_1-\gamma_0)+LN}$.

As the lower bound $\frac{A}{A+N(\gamma_1-\gamma_0)+LN}$ increases as A increases, we would now like to prove a lower bound of $A = \sum_{j:f(x_i) \in [\mu, \overline{\mu}]} \eta_j$.

Recall $\pi(1), ..., \pi(N)$ is the ordering of examples based on sigmoid score. We let $\pi^{-1}(\cdot)$ denote the inverse mapping of π , so that $\pi^{-1}(\pi(i)) = i$. We let $\underline{r} = \min(\{j : f(x_{\pi(j)}) \in [\underline{\mu}, s^*]\}),$ $\bar{r} = \max\{j : f(x_{\pi(j)}) \in (s^*, \bar{\mu}]\}$ and r^* denote the index where $f(x_{\pi(r^*)}) = s^*$.

First note since class 0 is the minority class, we must have $r^* \neq N$. Since $r^* \neq 0$ and $r^* \neq N$, we must have $\eta_{\pi(r^*)} \ge 0.5$ and $\eta_{\pi(r^*+1)} \le 0.5$. Otherwise, $r^* + 1$ or $r^* - 1$ will have lower risk than $R(\eta_{\pi(r^*)})$. By the smoothness definition above, we further have $\forall j \leq r^*, 0.5 - (r^* - j)L \leq \eta_{\pi(j)} \leq r^*$ $0.5 + (r^* - j + 1)L$, and $\forall j \ge r^*, \eta_{\pi(j)} \ge 0.5 - (j - r^* + 1)L$.

 $A = \sum_{j:f(x_j) \in [\mu,\bar{\mu}]} \eta_j$ can then be rewritten in the ranked format as $A = \sum_{j \in [r,\bar{r}]} \eta_j$.

First, we divide the sampling range into $[\underline{r}, r^*]$ and $[r^* + 1, \overline{r}]$. When sampling in $[\underline{r}, r^*]$, we have

$$R(f(x_{\pi(\underline{r})})) - R(s^{\star}) = \gamma_0 > 0 \implies$$

$$N\gamma_0 = \sum_{j \in (\underline{r}, r^{\star}]} \eta_{(j)} - \sum_{j \in (\underline{r}, r^{\star}]} (1 - \eta_{(j)}) \ge -1 + \sum_{j \in [\underline{r}, r^{\star}]} \eta_{(j)} - \sum_{j \in [\underline{r}, r^{\star}]} (1 - \eta_{(j)}). \tag{3}$$

When sampling in $[r^* + 1, \bar{r}]$, since $R(\bar{r}) - R(r^*) = \gamma_1$, we must have

$$N\gamma_1 = N \cdot (R(\bar{r}) - R(r^*))) = \sum_{j \in (r^*, \bar{r}]} 1 - 2\eta_{(j)} = \sum_{j \in [r^* + 1, \bar{r}]} (1 - \eta_{(j)}) - \sum_{j \in [r^* + 1, \bar{r}]} \eta_{(j)}.$$
 (4)

To obtain the lower bound of $\sum_{j \in [r, \bar{r}]} \eta_{(j)}$, we start by bounding \bar{r} and \underline{r} . Specifically, by equation 4, we have

$$N\gamma_{1} = \sum_{j \in [r^{\star}+1,\bar{r}]} (1 - 2\eta_{(j)}) \le \sum_{j \in [r^{\star}+1,\bar{r}]} (1 - 2(0.5 - (j - r^{\star})L))$$
$$= \sum_{j \in [r^{\star}+1,\bar{r}]} 2(j - r^{\star})L = \sum_{j=1}^{\bar{r}-r^{\star}} jL = (\bar{r} - r^{\star})(\bar{r} - r^{\star} + 1)L \le L(\bar{r} - r^{\star} + 1)^{2}$$

As a result $\bar{r} \ge \sqrt{\frac{N\gamma_1}{L}} + r^{\star} - 1$. Let $\alpha_1 = \sqrt{\frac{N\gamma_1}{L}} - 1$, we then have $\bar{r} \ge r^{\star} + \alpha_1$. Similarly, we have

$$N\gamma_0 \le \sum_{j \in [\underline{r}, r^\star]} (2\eta_{(j)} - 1) \le \sum_{j \in [\underline{r}, r^\star]} (2 \cdot (0.5 + (r^\star - j + 1)L) - 1)$$

m^{*} m ⊨ 1

$$= \sum_{j \in [\underline{r}, r^{\star}]} 2(r^{\star} - j + 1)L = \sum_{j=1}^{r} 2jL \le (r^{\star} - \underline{r} + 2)^{2}L,$$

so
$$\underline{r} \leq r^{\star} - \sqrt{\frac{N\gamma_0}{L}} + 2$$
. Let $\alpha_0 := \sqrt{\frac{N\gamma_0}{L}} - 2$, we then have $\underline{r} \leq r^{\star} - \alpha_0$

972 Now, we can bound $\sum_{j \in [r^*+1,\bar{r}]} \eta_{(j)}$ by the following

$$\sum_{j \in [r^{\star}+1,\bar{r}]} \eta_{(j)} = \sum_{j=r^{\star}+1}^{r^{\star}+\alpha_1} \eta_{(j)} \ge \sum_{j=r^{\star}+1}^{r^{\star}+\alpha_1} 0.5 - (j-r^{\star})L$$

$$= \sum_{j=1}^{n-1} 0.5 - jL = \frac{\alpha_1}{2} - \frac{\alpha_1(\alpha_1 + 1)L}{2}$$

979
980
981
$$= \frac{j=1}{(1-L)\alpha_1 - \alpha_1^2 L}{2}.$$

Similarly, we can bound $\sum_{j \in [r,r^*]} \eta_{(j)}$ by the following

$$\sum_{j \in [\underline{r}, r^{\star}]} \eta_{(j)} = \sum_{j=r^{\star}-\alpha_0}^{r^{\star}} \eta_{(j)} \ge \sum_{j=r^{\star}-\alpha_0}^{r^{\star}} 0.5 - (r^{\star}-j)L$$
$$= \sum_{j=0}^{\alpha_0} 0.5 - jL = \frac{\alpha_0 + 1}{2} - \frac{\alpha_0(\alpha_0 + 1)L}{2}$$
$$= \frac{(1-L)\alpha_0 - \alpha_0^2L + 1}{2}$$

Together, we have

$$\sum_{j \in [r,\bar{r}]} \eta_{(j)} \ge \frac{(1-L)(\alpha_0 + \alpha_1) - (\alpha_0^2 + \alpha_1^2)L + 1}{2}.$$
$$\ge \frac{1}{2}((1-L)(\alpha_0 + \alpha_1) + 1) \ge \frac{1}{2}(2(1-L)\sqrt{\frac{N\gamma}{L}} - 2) = (1-L)\sqrt{\frac{N\gamma}{L}} - 1$$

For the edge case, since the interval must have more than three snippets, we can bound $\sum_{j \in [\underline{r}, \overline{r}]} \eta_{(j)} \geq 1.5 - 2L$. Therefore, we have $\sum_{j \in [\underline{r}, \overline{r}]} \eta_{(j)} \geq \max(1.5 - 2L, (1 - L)\sqrt{\frac{N\gamma}{L}} - 1)$. Finally, we can bound the balancedness by

$$\frac{\sum_{j\in[\underline{r},\bar{r}]}\eta_{(j)}}{\sum_{j\in[\underline{r},\bar{r}]}(1-\eta_{(j)})} \ge \frac{\sum_{j\in[\underline{r},\bar{r}]}\eta_{(j)}}{\sum_{j\in[\underline{r},\bar{r}]}\eta_{(j)} + N(\gamma_1 - \gamma_0) + LN}$$

$$\ge 1 - \frac{N(\gamma_1 - \gamma_0) + LN}{\sum_{j\in[\underline{r},\bar{r}]}\eta_{(j)} + N(\gamma_1 - \gamma_0) + LN}$$

$$\ge 1 - \frac{N\bar{\gamma} + LN}{\sum_{j\in[\underline{r},\bar{r}]}\eta_{(j)} + N\bar{\gamma} + LN}$$

$$\ge 1 - \min(\frac{N\bar{\gamma} + LN}{1.5 + N\bar{\gamma} + LN - 2L}, \frac{N\bar{\gamma} + LN}{(1-L)\sqrt{\frac{N\gamma}{L}} + N\bar{\gamma} + LN - 1}$$

$$N\bar{\gamma} + LN = C, N\bar{\gamma} + LN + 1$$

$$\geq 1 - \min(\frac{N\bar{\gamma} + LN}{1.5 - 2L}, \sqrt{L} \cdot \frac{N\bar{\gamma} + LN + 1}{(1 - L)\sqrt{N\underline{\gamma}}})$$

1021 C TRAINING DETAILS

1023 Our model is fine-tuned using the AdamW optimizer with a cosine learning rate schedule. For every 1024 training of f in Algorithm 1, we train up to 5 epochs. When measuring performance, we use a 1025 separate validation set to find the highest performance checkpoint. For focal loss, we use $\gamma = 5$, and α is set to the imbalance ratio estimated from Table 3 for the minority class.