

Proxy-KD: A Proxy-Based Distillation Framework for Black-Box Large Language Models

Anonymous ACL submission

Abstract

Given the exceptional performance of proprietary large language models (LLMs) like GPT-4, recent research has increasingly focused on boosting the capabilities of smaller models through knowledge distillation (KD) from these powerful yet black-box teachers. While leveraging the high-quality outputs of these teachers is advantageous, the inaccessibility of their internal states often limits effective knowledge transfer. To overcome this limitation, we introduce Proxy-KD, a novel method that uses a proxy model to facilitate the efficient transfer of knowledge from black-box LLMs to smaller models. The white-box proxy is first aligned with the black-box teacher through supervised fine-tuning and preference optimization. Subsequently, the student model is trained using the black-box teacher’s hard labels and weighted soft logits from the aligned proxy, where the weights are based on the proxy’s alignment quality. Experimental results on multiple benchmark demonstrate that Proxy-KD significantly outperforms existing white-box and black-box knowledge distillation methods. This approach presents a compelling new avenue for distilling knowledge from advanced LLMs.

1 Introduction

Recently, proprietary large language models (LLMs) like GPT-3.5 (OpenAI, 2022) and GPT-4 (OpenAI, 2023) have demonstrated significant superiority over open-source counterparts such as the Llama series (Touvron et al., 2023a,b; MetaAI, 2024). However, their vast number of parameters leads to high inference costs, and they are only accessible via API calls, offering limited customization and transparency. To address these challenges, recent efforts like Alpaca (Taori et al., 2023), Vicuna (Chiang et al., 2023), and Orca (Mukherjee et al., 2023) have focused on transferring the capabilities of proprietary LLMs to smaller open-source

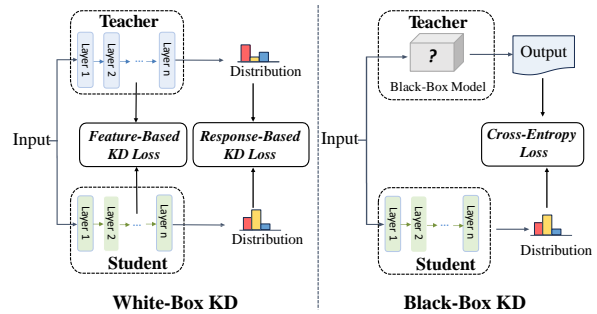


Figure 1: Comparison of white-box knowledge distillation (KD) and black-box knowledge distillation (KD).

models through knowledge distillation (Chen et al., 2023; Hsieh et al., 2023; Ho et al., 2022).

Knowledge distillation (KD) (Hinton et al., 2015) is a technique used to enhance the performance of a smaller student model by learning from a larger, more sophisticated teacher model. Depending on the level of access to the teacher model’s internals, KD methods can be categorized into two types: KD with black-box teachers and KD with white-box teachers. As illustrated in Figure 1, white-box KD allows the student model to distill more intrinsic knowledge from the teacher by mimicking the teacher model’s output distribution (Gu et al., 2023; Wen et al., 2023), hidden states (Jiao et al., 2020; Sun et al., 2019), and attention scores (Wang et al., 2021). Therefore, this method can only be applied when the teacher model’s parameters are accessible. On the other hand, black-box KD leverages the high-quality outputs from powerful proprietary LLMs to fine-tune the student model (Hsieh et al., 2023; Fu et al., 2023). Both white-box and black-box KD have their respective drawbacks. While white-box KD is hindered by the limited capacity of the teacher model, which often restricts the distillation performance of the student, black-box KD faces challenges with knowledge transfer due to the inaccessibility of the teacher model’s output distribution and internal states.

In this paper, we propose Proxy-based Knowl-

071	edge Distillation (Proxy-KD) to better transfer	122
072	knowledge from black-box teacher models. Proxy-	123
073	KD introduces a proxy model, typically a white-	124
074	box LLM, between the student and the black-box	
075	teacher. The proxy model first aligns with the capa-	125
076	bilities of the black-box teacher by leveraging the	126
077	teacher’s outputs. Moreover, preference optimiza-	127
078	tion is performed to further refine and enhance the	128
079	alignment between the proxy and teacher models.	
080	During the knowledge distillation process, the	
081	proxy model generates a dense distribution that	
082	closely approximates the black-box teacher’s out-	
083	put distribution. This enables the student model	
084	to train effectively as if it were using the black-	
085	box teacher’s guidance. To further improve the	
086	student’s learning effect, we propose incorporat-	
087	ing a sample-level weight into the distillation ob-	
088	jective. This weight reflects the quality of align-	
089	ment between the proxy and the teacher model	
090	for each sample, allowing the student to concen-	
091	trate on learning well-aligned distributions from the	
092	proxy. Moreover, the outputs from the black-box	
093	teacher serve as pseudo-labels for the supervised	
094	fine-tuning of the student model, akin to traditional	
095	white-box knowledge distillation. Introducing the	
096	proxy model also mitigates the model capacity gap	
097	issue (Cho and Hariharan, 2019), which typically	
098	occurs when there is a notable disparity in capabili-	
099	ties between the teacher and the student.	
100	To validate the effectiveness of our method, we	
101	conducted comprehensive experiments across a	
102	range of well-established benchmarks. The re-	
103	sults show that Proxy-KD consistently outperforms	
104	both black-box and white-box KD methods. We	
105	observed that the alignment between the proxy	
106	model and the black-box teacher is crucial; a poorly	
107	aligned proxy model significantly diminishes the	
108	performance of knowledge distillation. We also	
109	found that larger and more robust proxy models are	
110	generally more desirable, as they possess stronger	
111	foundational capabilities and can align more ef-	
112	fectively with the black-box teacher, enhancing	
113	the distillation process. Furthermore, we discov-	
114	ered that directly fine-tuning the proxy model with	
115	outputs from the black-box teacher is suboptimal	
116	for the alignment, requiring more effective align-	
117	ment methods. These findings highlight the impor-	
118	tance of selecting a well-aligned and capable proxy	
119	model to fully leverage the benefits of Proxy-KD.	
120	We summarize our contribution as below:	
121	• To tackle the challenge of knowledge distil-	
	lation for closed-source LLMs, we propose	122
	Proxy-KD, which introduces an aligned proxy	123
	between the teacher and student models.	124
	• We propose a DPO-based alignment strategy	125
	for the proxy to align with the teacher and	126
	demonstrate that this alignment is essential	127
	for Proxy-KD to achieve effective distillation.	128
	• We propose to include a sample-level weight	129
	in the distillation objective. This weight al-	130
	lows the student to concentrate on learning	131
	well-aligned distributions from the proxy.	132
	2 Related Work	133
	Existing knowledge distillation methods can be	134
	categorized into <i>white-box knowledge distillation</i>	135
	and <i>black-box knowledge distillation</i> .	136
	2.1 White-Box Knowledge Distillation	137
	Traditional knowledge distillation (KD) research	138
	predominantly employs white-box teachers and	139
	is typically classified into three main branches:	140
	feature-based, response-based, and relation-based	141
	methods. Feature-based methods seek to replicate	142
	the teacher’s intermediate representations, such as	143
	attention scores (Jiao et al., 2020), attribution maps	144
	(Wu et al., 2023), and hidden representations of	145
	tokens (Sun et al., 2019). Response-based methods	146
	train the student model by minimizing divergences	147
	like Kullback–Leibler (KL) divergence (Hinton	148
	et al., 2015; Sanh et al., 2019), reverse KL (Gu	149
	et al., 2023; Wen et al., 2023), Jensen–Shannon Di-	150
	vergence (JSD) (Fang et al., 2021; Yin et al., 2020),	151
	and Total Variation Distance (TVD) (Wen et al.,	152
	2023) based on the teacher’s output distribution.	153
	Relation-based methods train the student model by	154
	learning pairwise distances and triple-wise angles	155
	among token representations from the teacher (Park	156
	et al., 2021), or extracting structural relations from	157
	multi-granularity representations (Liu et al., 2022).	158
	2.2 Black-Box Knowledge Distillation	159
	Given the remarkable performance achieved by	160
	proprietary LLMs like GPT-4 (OpenAI, 2023),	161
	Claude 3 (Anthropic, 2024), and Gemini (Team	162
	et al., 2023), recent studies like Alpaca (Taori et al.,	163
	2023), Vicuna (Chiang et al., 2023), and Orca	164
	(Mukherjee et al., 2023) have focused on trans-	165
	ferring diverse capabilities from these black-box	166
	teachers into smaller open-source models. For in-	167
	stance, Li et al. (2024) and Liu et al. (2023) im-	168

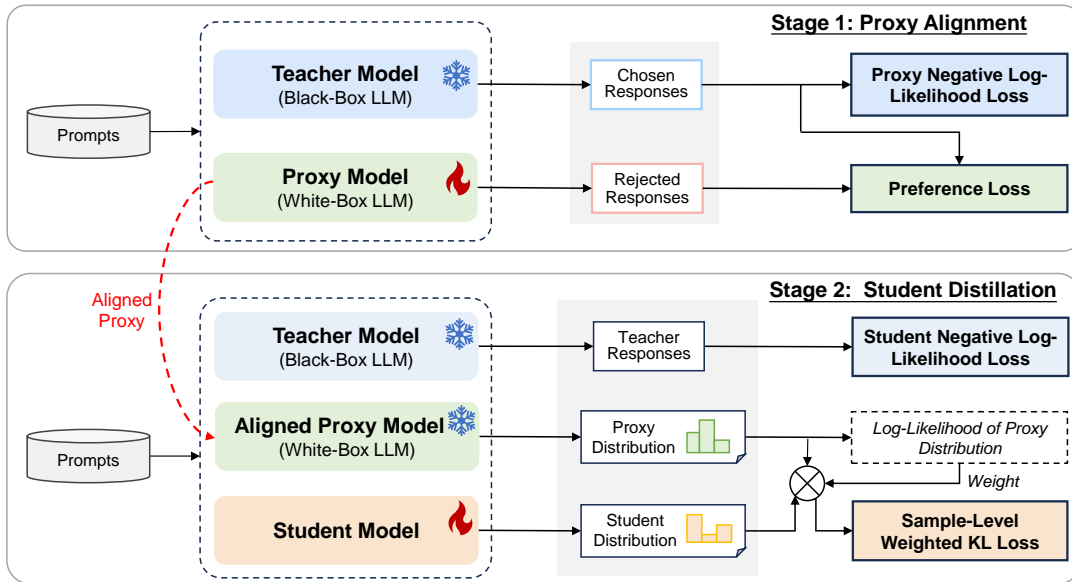


Figure 2: Overview of our proposed Proxy-based Knowledge Distillation (Proxy-KD).

proved the mathematical capability of small models by training on tailored rationale samples generated by GPT-3.5-Turbo and GPT-4. To transfer the code generation capability, [Azerbaiyev et al. \(2023\)](#) prompted Codex ([Chen et al., 2021](#)) to create natural language-code pairs and fine-tuned a smaller model on these samples. To transfer the tool usage capability, [Gou et al. \(2023\)](#) utilized GPT-4 to generate interactive tool-use trajectories as training samples for the target model. Other approaches, such as [Hsieh et al. \(2023\)](#); [Ho et al. \(2022\)](#); [Chen et al. \(2023\)](#), utilize rationales generated by black-box teachers as training data to transfer their general reasoning capabilities.

White-box knowledge distillation (KD) efficiently distills knowledge by leveraging the internal states of the teacher model. However, white-box teachers typically possess a more limited capacity compared to their black-box counterparts. In contrast, black-box KD capitalizes on the superior performance of the teacher models but is restricted to fine-tuning on teacher-generated samples. This approach captures input-output patterns without accessing the deeper, intrinsic knowledge of the teacher model. To bridge these gaps, we propose Proxy-KD, a straightforward method that combines the strengths of both white-box and black-box KD while mitigating their respective limitations.

2.3 Connection with Teacher Assistant

The proposed Proxy-KD method draws inspiration from TAKD ([Mirzadeh et al., 2020](#)), as both methods use an intermediate network to aid knowledge distillation, but they differ in three significant ways.

Firstly, the motivation behind each approach is distinct: TAKD focuses on mitigating the capacity gap between the teacher and student in white-box settings, whereas Proxy-KD addresses the challenges posed by black-box teacher models and seeks to incorporate the the benefits of white-box scenarios. Secondly, the methodologies diverge, with Proxy-KD introducing a proxy alignment phase with preference optimization to better align the proxy model with the black-box LLM and a weighting mechanism to adjust the student’s soft-label distillation loss based on proxy-teacher alignment scores. These steps reduce proxy-teacher discrepancies and prioritize samples with closer distribution alignment, thereby improving the effectiveness of the distillation process. Lastly, they operate in different domains: TAKD is applied in computer vision, while Proxy-KD is specifically designed for natural language processing, targeting the distillation of proprietary large language models (LLMs).

Some related works ([Zhou and Ai, 2024](#); [Lee et al., 2024](#)) also explore the idea of introducing an intermediate-sized teacher. [Zhou and Ai \(2024\)](#) focuses on using a teacher assistant primarily for filtering data generated by both the teacher and student models, and subsequently utilizing the filtered high-quality data for distillation. [Lee et al. \(2024\)](#) introduces an intermediate-sized teacher trained through fine-tuning, leveraging its soft labels to guide student learning during distillation. However, both of these methods overlook the importance of aligning the teacher assistant with the black-box teacher. Using an unaligned teacher assistant for black-box KD can harm student model perfor-

mance (see experiments). Proxy-KD tackles this by introducing an online preference alignment and sample-level weighting in the distillation objective. This focus on alignment is a novel contribution that has not been considered in previous related works.

3 Method

In this section, we introduce Proxy-based Knowledge Distillation (Proxy-KD), a simple yet efficient approach for knowledge distillation from black-box LLMs. As illustrated in Figure 2, Proxy-KD introduces a larger white-box LLM as the proxy aiming to capture the black-box teacher’s knowledge. The process unfolds in two main stages: (1) proxy model alignment and (2) student knowledge distillation. First, the proxy model is aligned with the teacher through supervised fine-tuning and preference optimization. Once aligned, the student model learns from both the explicit outputs (hard labels) of the black-box teacher and output distributions (soft labels) provided by the aligned proxy.

3.1 Problem Statement

To facilitate the transfer of knowledge from a black-box teacher LLM π_t to a smaller, open-source student LLM π_s , we introduce a proxy model π_p . The training dataset \mathcal{D} consists of input-output pairs (x, y) , where x represents the input prompt and y is the output sequence generated by the teacher model π_t . This dataset is strategically divided into three parts: 10% (\mathcal{D}_w) for the warm-up phase, 45% (\mathcal{D}_p) for aligning the proxy model with the teacher, and the remaining 45% (\mathcal{D}_s) for the knowledge distillation training of the student model.

The process begins with a warm-up phase where the proxy model π_p is trained on \mathcal{D}_w . This phase helps π_p develop a basic capability to generate responses to input prompts. Following this, the proxy model undergoes alignment with the teacher model π_t using the next dataset, \mathcal{D}_p . This alignment is achieved through two methods: hard-label knowledge distillation (KD) and preference learning. These methods enable π_p to approximate the behavior and outputs of the teacher model. Once aligned, π_p acts as an intermediary, facilitating the transfer of knowledge to the student π_s on \mathcal{D}_s .

3.2 Preliminary

Hard-Label Knowledge Distillation. In this approach, the student model is trained using the outputs generated by the teacher model by minimizing

the negative log-likelihood (NLL) function:

$$\mathcal{L}_{\text{NLL}} = \mathbb{E}_{(x,y) \sim \mathcal{D}} [-\log \pi_s(y|x)], \quad (1)$$

where $\pi_s(y|x)$ is the probability of π_s generating y given x . This approach is essentially a form of supervised fine-tuning and typically employed when the teacher is a black-box model.

Soft-Label Knowledge Distillation. In this approach, the student is trained to imitate the token-level probabilities of the teacher, by minimizing the Kullback-Leibler (KL) divergence:

$$\mathcal{L}_{\text{KL}} = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\mathbb{D}_{\text{KL}}(\pi_t(y|x) || \pi_s(y|x))]. \quad (2)$$

This knowledge distillation approach is typically employed when the teacher is a white-box model.

While the KL divergence objective provides richer supervision signals by using the token-level output distributions of the teacher model, it cannot be applied to black-box teachers due to the inaccessibility of these distributions. Consequently, current methods (Chiang et al., 2023; Mukherjee et al., 2023) rely on supervised fine-tuning using the outputs generated by black-box models to transfer their knowledge. Proxy-KD addresses this limitation by using a proxy model to incorporate the KL objective. The proxy mimics the black-box teacher, allowing access to its output distributions and enabling a more effective knowledge transfer.

3.3 Proxy Model Alignment

The proxy model π_p is typically a larger white-box LLM than the student model π_s . For effective knowledge transfer, it’s crucial to first align the output distribution of the proxy model with that of the black-box teacher model π_t . This alignment ensures that the proxy accurately captures the teacher’s behavior.

The proxy model π_p first undergoes supervised fine-tuning on a warm-up dataset \mathcal{D}_w . Following this, the proxy is further trained on the \mathcal{D}_p dataset by minimizing the NLL loss:

$$\mathcal{L}_{\text{Proxy-NLL}} = \mathbb{E}_{(x,y) \sim \mathcal{D}_p} [-\log \pi_p(y|x)]. \quad (3)$$

To enhance the alignment of the proxy model with the teacher, we further introduce a preference learning-based alignment objective, with the hypothesis that the teacher model’s responses are of higher quality compared to those from the unaligned proxy model. The objective is to iteratively

adjust the proxy model so that it increasingly favors responses similar to those of the teacher while reducing its preference for its own initial outputs. To implement this, we employ the Direct Preference Optimization (DPO) algorithm (Rafailov et al., 2024), which refines the proxy model by systematically preferring the teacher’s responses.

Specifically, for a given input x , we iteratively sample a response y from the teacher and \hat{y} from the proxy. These responses form a preference pair (x, y, \hat{y}) . To train the proxy model to prefer y over \hat{y} , we define the following preference loss function:

$$\mathcal{L}_{\text{DPO}}^{(i)}(x, y, \hat{y}) = \log \sigma \left[\beta \log \frac{\pi_p^{(i)}(y|x)}{\pi_p^{(i-1)}(y|x)} - \beta \log \frac{\pi_p^{(i)}(\hat{y}|x)}{\pi_p^{(i-1)}(\hat{y}|x)} \right], \quad (4)$$

where $\pi_p^{(i-1)}$ is the proxy model from the previous training iteration. The overall preference loss over all the preference samples is defined as:

$$\mathcal{L}_{\text{Pref}}^{(i)} = \mathbb{E}_{(x,y) \sim \mathcal{D}_p, \hat{y} \sim \pi_p^{(i)}(x)} \mathcal{L}_{\text{DPO}}^{(i)}(x, y, \hat{y}). \quad (5)$$

At each iteration i , the proxy model is updated based on the combined objective that includes both the NLL loss and the preference loss:

$$\mathcal{L}_{\text{Proxy}}^{(i)} = \mathcal{L}_{\text{Proxy-NLL}}^{(i)} + \mathcal{L}_{\text{Pref}}^{(i)}. \quad (6)$$

This iterative process continues for a fixed number of iterations k or until the proxy model converges. Through this method, the proxy model π_p is aligned to emulate the distribution of the black-box teacher π_t , becoming an effective intermediary for transferring knowledge to the student model.

3.4 Knowledge Distillation

To transfer knowledge from the black-box teacher to the student model π_s , we define the first training objective using teacher-generated sequences and the hard-label knowledge distillation objective:

$$\mathcal{L}_{\text{Student-NLL}} = \mathbb{E}_{(x,y) \sim \mathcal{D}_s} [-\log \pi_s(y|x)]. \quad (7)$$

Based on the proxy model aligned with the black-box teacher, which delivers accessible output distributions, we define another training objective for the student through soft-label knowledge distillation:

$$\mathcal{L}_{\text{Student-KL}} = \mathbb{E}_{(x,y) \sim \mathcal{D}_s} [\mathbb{D}_{\text{KL}}(\pi_p(y|x) || \pi_s(y|x))]. \quad (8)$$

In this process, the proxy model functions as an intermediary for the black-box teacher, facilitating the transfer of knowledge to the student

model. However, as illustrated in Figure 5 in Appendix, discrepancies between the teacher’s and the proxy’s output distributions persist even after aligning the proxy model, potentially degrading the effectiveness of knowledge distillation. To address these discrepancies, we propose a weighted approach to the soft-label knowledge distillation objective. By introducing weights, we dynamically adjust the influence of each sample based on the alignment quality between the proxy and the black-box teacher. This approach ensures that the student model prioritizes samples where the proxy’s distribution closely matches the teacher’s distribution and reduces focus on samples where it does not. The weights are calculated based on the log-likelihood of the teacher’s output generated by the proxy, normalized by the mean and variance of these log-likelihoods:

$$\begin{aligned} w(x, y) &= \sigma \left[\frac{\log \pi_p(y|x) - \mu}{\gamma} \right], \\ \mu &= \mathbb{E}_{(x,y) \sim \mathcal{D}_s} [\log \pi_p(y|x)], \\ \gamma^2 &= \mathbb{V}\text{ar}_{(x,y) \sim \mathcal{D}_s} [\log \pi_p(y|x)], \end{aligned} \quad (9)$$

where $w(x, y)$ is a weight reflecting the quality of the proxy’s prediction for the sample (x, y) , $\mathbb{V}\text{ar}(\cdot)$ is the variance operation, γ is the standard deviation, σ is the sigmoid function. Based on Equation (8), we derive the sample-level weighted version of $\mathcal{L}_{\text{Student-KL}}$ as follow:

$$\begin{aligned} \mathcal{L}_{\text{Weight-KL}} &= \\ &\mathbb{E}_{(x,y) \sim \mathcal{D}_s} [w(x, y) \mathbb{D}_{\text{KL}}(\pi_p(y|x) || \pi_s(y|x))]. \end{aligned} \quad (10)$$

Therefore, the overall objective for student knowledge distillation can be derived as:

$$\mathcal{L}_{\text{Student}} = \mathcal{L}_{\text{Student-NLL}} + \alpha \mathcal{L}_{\text{Weight-KL}}, \quad (11)$$

where α is a hyperparameter utilized to adjust the strength of the weighted KL loss.

This knowledge distillation strategy effectively blends the advantages of both black-box and white-box knowledge distillation methods, employing the proxy model to bridge the gap between black-box LLMs and open-source student LLMs.

4 Experimental Setup

In this section, we introduce the experimental settings of models, datasets, and method baselines.

4.1 Models and Datasets

Teacher/Proxy/Student Models. In Proxy-KD, we choose GPT-4 (OpenAI, 2023) as the teacher, which is a powerful proprietary large language model. We select Llama-2-70b (Touvron et al., 2023b) and Llama-2-13b (MetaAI, 2024) as the proxy, respectively. Our student models come from two model types: Llama-1-7B (Touvron et al., 2023a) and Llama-2-7B (Touvron et al., 2023b).

Training Corpus. We combine the OpenOrca (Lian et al., 2023) and Nectar (Zhu et al., 2023) datasets as our training corpus, containing a total of 1M output sequences generated by the black-box teacher GPT-4. The OpenOrca dataset consists of instruction-following tasks, where GPT-4 is prompted to generate responses based on diverse input instructions. Nectar is a 7-wise comparison dataset, we filter and select those responses derived from GPT-4. Following Li et al. (2024), we also incorporate synthetic data generated by GPT-4, based on existing benchmark training sets. We split the original training corpus \mathcal{D} into three parts: 10% as \mathcal{D}_w with 100K samples, 45% as \mathcal{D}_p with 450K samples, and 45% as \mathcal{D}_s with 450K samples.

Evaluation Benchmarks. Evaluation benchmarks include complex reasoning dataset BBH (Suzgun et al., 2022), knowledge-based datasets AGIEval (Zhong et al., 2023), ARC-challenge (Clark et al., 2018), and MMLU (Zeng, 2023), commonsense reasoning dataset CSQA (Talmor et al., 2019), and mathematical reasoning dataset GSM8K (Cobbe et al., 2021). All evaluated models apply a zero-shot greedy decoding strategy.

4.2 Training Configurations

All experiments are conducted on 8xA100 Nvidia GPUs with 80GB memory. All proxy and student models are trained for only one epoch. We use a constant learning rate of $1e-5$ and the Adam optimizer, with a max sequence length of 1024. We set hyperparameter $\alpha = 100$ in Equation (11), and $k = 16$ for the number of proxy alignment iterations. All models are trained using LoRA (Hu et al., 2021) with mixed-precision: frozen parameters in bfloat16 and LoRA-trained parameters in float32.

4.3 Baselines

We compare Proxy-KD with different white-box KD and black-box KD methods.

White-Box KD. For knowledge distillation with

white-box teachers, we compare forward KL methods (Hinton et al., 2015; Agarwal et al., 2024) and reverse KL methods including MiniLLM (Gu et al., 2023) and GKD (Agarwal et al., 2023) (with the same hyperparameters set in the paper). The chat version of Llama-2-70b is utilized as the white-box teacher. We also compare with using the aligned proxy as white-box teacher to perform distillation.

Black-Box KD. For knowledge distillation with black-box teachers, we compare the black-box KD method (i.e., data distillation) (Mukherjee et al., 2023; Mitra et al., 2023; Xu et al., 2023), which directly fine-tunes the student on the data generated by the black-box teacher. We also compare Proxy-KD with the TAKD (Mirzadeh et al., 2020) and Mentor-KD (Lee et al., 2024) method.

For baselines implemented by us, we start from the same student checkpoint as Proxy-KD and use the same input prompts. In white-box KD, output sequences are generated by the white-box teacher, while in black-box KD, output sequences are generated by the black-box teacher.

5 Result and Analysis

In this section, we present the main results and additional experiments of Proxy-KD.

5.1 Overall Results

We show the comparison of Proxy-KD against baselines in Table 1, the proxy models in Proxy-KD are based on Llama-2-70B backbone. Overall, the performance of black-box KD methods outperforms that of white-box KD methods, demonstrating the efficacy of distilling knowledge from powerful black-box models. Additionally, we conduct extended experiments on Proxy-KD using Qwen models in Table 5 in Appendix.

Proxy-KD outperforms white-box KD and black-box KD methods. Notably, Proxy-KD further enhances the performance, consistently achieving higher scores across most evaluated benchmarks compared to the white-box KD methods (e.g. MiniLLM and GKD) and the black-box KD methods. Improvement is particularly pronounced in the challenging datasets like ARC, BBH, and GSM8K, where Proxy-KD obtains accuracy of 71.09%, 53.40%, and 53.07%, respectively.

Proxy-KD outperforms TAKD consistently. TAKD performs even worse than Black-Box KD. When using Llama-1-7B as the student, Black-Box KD achieves an average of 49.11%, while TAKD

Table 1: Overall results on evaluated benchmarks. We report accuracy (%) for all tasks. Best performances are shown in **bold**, while suboptimal ones underlined. All models utilize a zero-shot greedy decoding strategy for evaluation. Llama-2-70B-Proxy indicates that we use the aligned proxy as the white-box teacher for distillation.

Method	Student	Teacher	AGIEval	ARC	BBH	CSQA	GSM8K	MMLU	Avg
<i>Black-Box Teacher</i>									
GPT-4	-	-	56.40	93.26	88.0	-	92.0	86.4	-
<i>White-Box KD</i>									
Forward KL	Llama-1-7B	Llama-2-70B-Chat	25.16	62.18	37.27	74.20	37.39	45.43	46.94
Forward KL	Llama-2-7B	Llama-2-70B-Chat	35.16	66.87	35.68	74.40	44.12	51.42	51.27
Forward KL	Llama-2-7B	Llama-2-70B-Proxy	35.56	<u>69.34</u>	45.72	74.97	46.34	51.13	53.84
MiniLLM (Gu et al., 2023)	Llama-2-7B	Llama-2-70B-Chat	<u>35.77</u>	63.25	<u>53.11</u>	75.15	44.64	51.32	<u>53.87</u>
GKD (Agarwal et al., 2023)	Llama-2-7B	Llama-2-70B-Chat	34.22	62.28	52.58	<u>75.16</u>	42.79	50.64	52.95
<i>Black-Box KD</i>									
Black-Box KD	Llama-1-7B	GPT-4	28.01	63.17	41.98	74.43	41.83	45.21	49.11
Black-Box KD	Llama-2-7B	GPT-4	34.71	66.85	46.68	74.43	49.51	49.82	53.66
Mentor-KD (Lee et al., 2024)	Llama-2-7B	GPT-4	35.49	70.14	51.33	<u>75.16</u>	<u>52.11</u>	50.13	55.73
TAKD (Mirzadeh et al., 2020)	Llama-1-7B	GPT-4	26.74	64.31	39.55	72.01	40.49	38.78	46.98
TAKD (Mirzadeh et al., 2020)	Llama-2-7B	GPT-4	35.46	68.31	49.41	74.52	49.49	48.12	53.22
Proxy-KD (ours)	Llama-1-7B	GPT-4	35.47	67.48	43.74	74.08	44.89	41.88	52.09
Proxy-KD (ours)	Llama-2-7B	GPT-4	36.59	71.09	53.40	75.18	53.07	<u>51.35</u>	56.78

only reaches 46.98%. Similarly, with Llama-2-7B as the student, Black-Box KD attains 53.66% compared to TAKD’s average of 53.22%. This decline in performance is likely due to TAKD’s lack of proxy alignment, a critical step in closed-source KD. Introducing an unaligned proxy fails to enhance and even degrades the student model’s performance. Additionally, the results show Proxy-KD outperforms Mentor-KD across benchmarks, demonstrating the effectiveness of Proxy-KD.

Proxy-KD outperforms white-box KD with an aligned proxy as the teacher. Relying solely on an aligned proxy for white-box KD offers limited knowledge to the student, attaining average accuracy of 53.84%, compared to Proxy-KD’s average of 56.78%. This suggests that the capabilities of closed-source teachers are more beneficial than those of open-source teachers, even after alignment, underscoring the superiority of distilling from closed-source LLMs.

5.2 Ablation Studies

In this section, we examine the impact of different components within Proxy-KD. Llama-2-7B and Llama-2-70B are utilized as the backbones of the student and the proxy models, respectively.

Effect of the Proxy Model. The proxy model π_p is crucial for the effectiveness of Proxy-KD. Removing it forces the distillation process to revert to hard-label knowledge distillation, leading to significant performance drops across multiple benchmarks: a decrease of 4.24 on ARC, 6.72 on BBH, and 3.56 on GSM8K, as shown in Table 2. These declines underscore the proxy model’s

essential role in capturing and transferring the distributional knowledge from the black-box teacher, which is particularly vital for complex reasoning and mathematical tasks. Without the proxy, the student model lacks detailed distributional guidance, resulting in markedly lower performance.

Effect of Proxy Model Alignment. The proxy model alignment, facilitated by the loss $\mathcal{L}_{\text{Proxy}}$, is vital for effective knowledge transfer. Table 2 shows that when the proxy is initialized from the Llama-2-70B checkpoint without alignment, the performance drops notably on BBH (-10.40), GSM8K (-5.53), and MMLU (-3.26). This decline illustrates the adverse effect of an unaligned proxy, which fails to approximate the teacher’s distribution and underperforms against models directly fine-tuned on teacher data. The slight increase on CSQA (+0.86) without alignment might be due to the simplicity of the task, indicating potential overfitting to teacher outputs without proxy guidance. This reinforces the necessity of the alignment process to ensure the proxy effectively bridges the knowledge transfer from the teacher to the student model across diverse and complex tasks.

Effect of Preference Optimization. Table 2 illustrates the significant role of preference optimization in enhancing the performance of both the proxy and student models. When the proxy preference loss $\mathcal{L}_{\text{Pref}}$ is removed, reducing the proxy alignment loss to $\mathcal{L}_{\text{Proxy-NLL}}$, we observe notable performance drops across various benchmarks. Specifically, the alignment of the proxy model with the black-box teacher deteriorates, as evidenced by decreases in scores on benchmarks like BBH and

Table 2: Ablation studies of Proxy-KD. We examine the impact of the proxy model π_p , proxy model alignment loss $\mathcal{L}_{\text{Proxy}}$, proxy preference loss $\mathcal{L}_{\text{Pref}}$, and weighted KL loss $\mathcal{L}_{\text{Weight-KL}}$ on the performance of the student model training, as well as the impact of the proxy preference loss $\mathcal{L}_{\text{Pref}}$ on the performance of the proxy model alignment.

Method	AGIEval	ARC	BBH	CSQA	GSM8K	MMLU
<i>Student Model Distillation</i>						
$\mathcal{L}_{\text{Student}}$	36.59	71.09	53.40	75.18	53.07	51.35
w/o π_p	34.71 (-1.88)	66.85 (-4.24)	46.68 (-6.72)	74.43 (-0.75)	49.51 (-3.56)	49.82 (-1.53)
w/o $\mathcal{L}_{\text{Proxy}}$	35.05 (-1.54)	67.18 (-3.91)	43.0 (-10.40)	76.04 (+0.86)	47.54 (-5.53)	48.09 (-3.26)
w/o $\mathcal{L}_{\text{Pref}}$	35.38 (-1.21)	66.11 (-4.98)	52.51 (-0.89)	75.51 (+0.33)	52.49 (-0.58)	48.79 (-2.56)
w/o $\mathcal{L}_{\text{Weight-KL}}$	33.99 (-2.60)	71.81 (+0.72)	51.50 (-1.90)	75.11 (-0.07)	52.91 (-0.16)	49.47 (-1.88)
<i>Proxy Model Alignment</i>						
$\mathcal{L}_{\text{Proxy}}$	49.12	87.67	66.04	82.18	78.24	68.62
w/o $\mathcal{L}_{\text{Pref}}$	48.79 (-0.33)	87.11 (-0.56)	64.87 (-1.17)	81.33 (-0.85)	78.56 (+0.32)	66.94 (-1.68)
w/o $\mathcal{L}_{\text{Pref}}$ and $\mathcal{L}_{\text{Proxy-NLL}}$	48.31 (-0.81)	86.93 (-0.74)	62.16 (-3.88)	80.95 (-1.23)	79.15 (+0.91)	66.38 (-2.24)

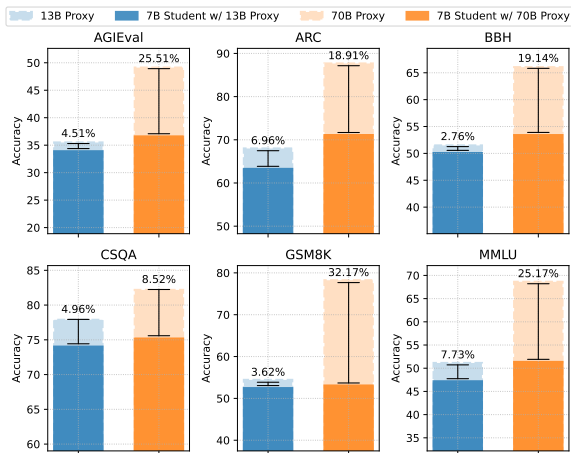


Figure 3: Performance of student models under different proxy models. We also show the ratio of performance gap between the proxy models and the student models.

MMLU, which subsequently impacts the student model. The overall trend confirms that preference optimization is crucial for refining the proxy model’s ability to emulate the teacher effectively.

Effect of Weighted KL. When $\mathcal{L}_{\text{Weight-KL}}$ is replaced with the standard KL loss $\mathcal{L}_{\text{Student-KL}}$, we also observe declines in performance across most benchmarks, indicating that the effectiveness of the distillation process diminishes. The results shown in Table 2 highlight that focusing on high log-likelihood distributions from the proxy, as facilitated by the weighted KL loss, significantly enhances the quality of knowledge transfer. The overall declines underscore that this weighting mechanism significantly improves the quality of knowledge distillation, enhancing the student’s ability to learn from a well-aligned proxy.

5.3 Impact of Proxy Model’s Capability

How well the proxy aligned with the teacher can directly affect the performance of the student. Proxy alignment effectiveness depends on two factors: the design of the alignment algorithm and the in-

herent alignment capability of the proxy backbone model itself. This section investigate the impact of the latter. We hypothesize that the size of the proxy model’s parameters is crucial for its capacity to align with the black-box teacher’s capability, especially when the teacher’s parameter size is significantly larger than the proxy’s. Experiments are conducted with Llama-2-70B and Llama-2-13B as the proxy backbone models. We show the performance of these aligned proxy models. As depicted in Figure 3, the proxy model based on Llama-2-70B performs better than the one based on Llama-2-13B, the latter has fewer parameters. We also examine the impact of proxy models with different capacities on student performance. We observe that the stronger Llama-2-70B proxy yields better student performance than the weaker proxy based on Llama-2-13B. Furthermore, when using a proxy based on a backbone model with a larger capacity, the student demonstrates a greater potential for achieving higher performance.

6 Conclusion

This paper aims to tackle the challenge of knowledge distillation for black-box large language models (LLMs), where we can only access the outputs generated by the teacher model. Given the inaccessibility of the internal states of these black-box models, we introduce Proxy-KD, a novel approach that leverages a proxy model to enhance the distillation process. The proxy model is first aligned with the black-box teacher, closely mimicking its behavior. Then, the student model is trained using the combined knowledge from both the black-box teacher and the proxy model. Extensive experiments and analyses across a variety of well-established benchmarks demonstrate that Proxy-KD significantly outperforms existing black-box and white-box knowledge distillation methods.

633 Limitations

634 The limitations of this work include the training
635 time overhead associated with proxy model align-
636 ment, particularly when the proxy model has a large
637 number of parameters. Additionally, the proposed
638 preference optimization requires online sampling
639 from the proxy model, further increasing the train-
640 ing time overhead. Another limitation is the type
641 of experimental backbone models used. Due to
642 resource constraints, this work only conducts ex-
643 periments with the Llama model series, without
644 including other model backbones such as Mistral
645 (Jiang et al., 2023).

646 References

647 Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Pi-
648 otr Stanczyk, Sabela Ramos Garea, Matthieu Geist,
649 and Olivier Bachem. 2024. On-policy distillation
650 of language models: Learning from self-generated
651 mistakes. In *The Twelfth International Conference
652 on Learning Representations*.

653 Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr
654 Stanczyk, Sabela Ramos, Matthieu Geist, and Olivier
655 Bachem. 2023. [Generalized knowledge distillation
656 for auto-regressive language models](#).

657 Anthropic. 2024. [Claude 3 family](#). Accessed: 2024-06-
658 04.

659 Zhangir Azerbayev, Ansong Ni, Hailey Schoelkopf, and
660 Dragomir Radev. 2023. [Explicit knowledge transfer
661 for weakly-supervised code generation](#).

662 Hongzhan Chen, Siyue Wu, Xiaojun Quan, Rui Wang,
663 Ming Yan, and Ji Zhang. 2023. [MCC-KD: Multi-
664 CoT consistent knowledge distillation](#). In *Findings
665 of the Association for Computational Linguistics:
666 EMNLP 2023*, pages 6805–6820, Singapore. Associ-
667 ation for Computational Linguistics.

668 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
669 Henrique Ponde, Jared Kaplan, Harrison Edwards,
670 Yura Burda, Nicholas Joseph, et al. 2021. [Evaluat-
671 ing large language models trained on code](#). *ArXiv*,
672 abs/2107.03374.

673 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,
674 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan
675 Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion
676 Stoica, and Eric P. Xing. 2023. [Vicuna: An open-
677 source chatbot impressing gpt-4 with 90%* chatgpt
678 quality](#).

679 Jang Hyun Cho and Bharath Hariharan. 2019. On the
680 efficacy of knowledge distillation. In *Proceedings of
681 the IEEE/CVF international conference on computer
682 vision*, pages 4794–4802.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot,
Ashish Sabharwal, Carissa Schoenick, and Oyvind
Tafjord. 2018. [Think you have solved question an-
swering? try arc, the ai2 reasoning challenge](#). *ArXiv*,
abs/1803.05457. 683
684
685
686
687

Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,
Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias
Plappert, Jerry Tworek, Jacob Hilton, Reiichiro
Nakano, et al. 2021. Training verifiers to solve math
word problems. *arXiv preprint arXiv:2110.14168*. 688
689
690
691
692

Gongfan Fang, Yifan Bao, Jie Song, Xinchao Wang,
Donglin Xie, Chengchao Shen, and Mingli Song.
2021. Mosaicking to distill: Knowledge distillation
from out-of-domain data. *Advances in Neural Infor-
mation Processing Systems*, 34:11920–11932. 693
694
695
696
697

Yao Fu, Hao Peng, Litu Ou, Ashish Sabharwal, and
Tushar Khot. 2023. Specializing smaller language
models towards multi-step reasoning. *arXiv preprint
arXiv:2301.12726*. 698
699
700
701

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen,
Yujia Yang, Minlie Huang, Nan Duan, and Weizhu
Chen. 2023. [Tora: A tool-integrated reasoning
agent for mathematical problem solving](#). *ArXiv*,
abs/2309.17452. 702
703
704
705
706

Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. 2023.
Minillm: Knowledge distillation of large language
models. In *The Twelfth International Conference on
Learning Representations*. 707
708
709
710

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.
Distilling the knowledge in a neural network. *arXiv
preprint arXiv:1503.02531*. 711
712
713

Namgyu Ho, Laura Schmid, and Se-Young Yun. 2022.
Large language models are reasoning teachers. *arXiv
preprint arXiv:2212.10071*. 714
715
716

Cheng-Yu Hsieh, Chun-Liang Li, Chih-kuan Yeh,
Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay
Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. [Dis-
tilling step-by-step! outperforming larger language
models with less training data and smaller model
sizes](#). In *Findings of the Association for Compu-
tational Linguistics: ACL 2023*, pages 8003–8017,
Toronto, Canada. Association for Computational Lin-
guistics. 717
718
719
720
721
722
723
724
725

Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu,
Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
et al. 2021. Lora: Low-rank adaptation of large lan-
guage models. In *International Conference on Learn-
ing Representations*. 726
727
728
729
730

Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-
sch, Chris Bamford, Devendra Singh Chaplot, Diego
de las Casas, Florian Bressand, Gianna Lengyel, Guil-
laume Lample, Lucile Saulnier, et al. 2023. Mistral
7b. *arXiv preprint arXiv:2310.06825*. 731
732
733
734
735

736	Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. TinyBERT: Distilling BERT for natural language understanding . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4163–4174, Online. Association for Computational Linguistics.	792
737		793
738		794
739		
740		795
741		796
742		797
743		798
744		799
745		800
746		
747		801
748		802
749		803
		804
		805
750		
751		806
752		807
753		808
754		809
755		810
756		811
757		812
758		813
759		814
		815
760		816
761		817
762		
763		818
764		819
		820
765		821
766		822
767		823
768		824
769		
770		825
771		826
		827
772		828
773		829
		830
774		831
775		832
776		833
777		
778		834
779		835
		836
780		837
781		838
782		
783		839
784		840
785		841
		842
786		843
787		844
788		
789		845
790		846
		847
791		

848 Azhar, et al. 2023a. Llama: Open and effi-
849 cient foundation language models. *arXiv preprint*
850 *arXiv:2302.13971*.

851 Hugo Touvron, Louis Martin, Kevin R. Stone, Peter
852 Albert, et al. 2023b. *Llama 2: Open foundation and*
853 *fine-tuned chat models*. *ArXiv*, abs/2307.09288.

854 Wenhui Wang, Hangbo Bao, Shaohan Huang, Li Dong,
855 and Furu Wei. 2021. *MiniLMv2: Multi-head self-*
856 *attention relation distillation for compressing pre-*
857 *trained transformers*. In *Findings of the Association*
858 *for Computational Linguistics: ACL-IJCNLP 2021*,
859 pages 2140–2151, Online. Association for Computa-
860 tional Linguistics.

861 Yuqiao Wen, Zichao Li, Wenyu Du, and Lili Mou. 2023.
862 *f-divergence minimization for sequence-level knowl-*
863 *edge distillation*. In *Proceedings of the 61st An-*
864 *ual Meeting of the Association for Computational*
865 *Linguistics (Volume 1: Long Papers)*, pages 10817–
866 10834, Toronto, Canada. Association for Computa-
867 tional Linguistics.

868 Siyue Wu, Hongzhan Chen, Xiaojun Quan, Qifan Wang,
869 and Rui Wang. 2023. *AD-KD: Attribution-driven*
870 *knowledge distillation for language model compres-*
871 *sion*. In *Proceedings of the 61st Annual Meeting of*
872 *the Association for Computational Linguistics (Vol-*
873 *ume 1: Long Papers)*, pages 8449–8465, Toronto,
874 Canada. Association for Computational Linguistics.

875 Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng,
876 Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin
877 Jiang. 2023. *Wizardlm: Empowering large language*
878 *models to follow complex instructions*.

879 Hongxu Yin, Pavlo Molchanov, Jose M Alvarez,
880 Zhizhong Li, Arun Mallya, Derek Hoiem, Niraj K
881 Jha, and Jan Kautz. 2020. Dreaming to distill: Data-
882 free knowledge transfer via deepinversion. In *Pro-*
883 *ceedings of the IEEE/CVF Conference on Computer*
884 *Vision and Pattern Recognition*, pages 8715–8724.

885 Hui Zeng. 2023. *Measuring massive multitask chinese*
886 *understanding*. *ArXiv*, abs/2304.12986.

887 Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang,
888 Shuai Lu, Yanlin Wang, Amin Saied Sanosi Saied,
889 Weizhu Chen, and Nan Duan. 2023. *Agieval: A*
890 *human-centric benchmark for evaluating foundation*
891 *models*. *ArXiv*, abs/2304.06364.

892 Yuhang Zhou and Wei Ai. 2024. *Teaching-assistant-*
893 *in-the-loop: Improving knowledge distillation from*
894 *imperfect teacher models in low-budget scenarios*.
895 In *Findings of the Association for Computational*
896 *Linguistics: ACL 2024*, pages 265–282, Bangkok,
897 Thailand. Association for Computational Linguistics.

898 Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu,
899 and Jiantao Jiao. 2023. *Starling-7b: Improving llm*
900 *helpfulness and harmlessness with rlaiif*.

Table 3: Training time overhead. We show the training hours per round for different methods. SFT is the supervised fine-tuning method, Distill is the knowledge distillation method, Pref is the preference optimization method. For GKD (Agarwal et al., 2023), student model is based on 7B, teacher model is based on 70B. Each round contains 40K training samples.

Models	#GPUs	Hours/Round
Llama-7B-SFT	4	1.0
Llama-7B-Distill	4	2.0
Llama-7B-GKD	8	10.0
Llama-13B-SFT	8	1.8
Llama-13B-Pref	8	9.0
Llama-70B-SFT	8	5.5
Llama-70B-Pref	8	28.0

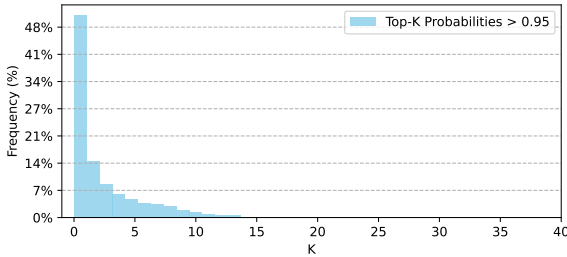


Figure 4: The statistics of the cumulative probability within the Top K exceeding 0.95. The x-axis represents different values of K, while the y-axis shows the percentage of instances meeting this threshold.

A Experimental Analysis

A.1 Analysis of Training Efficiency

We show the training time overhead for different methods in Table 3. We show the training hours per round for supervised fine-tuning, knowledge distillation, and preference optimization methods across various model sizes. Each round contains 40K training samples. We note that preference optimization is the main time overhead due to online sampling from the proxy model. In Proxy-KD, we obtain the proxy model’s output distribution offline during student distillation. As Figure 4 shows, most probability mass is concentrated on a few tokens. To save memory, only the top 10 token indices and their logits are retained.

A.2 Output Token Agreement

To serve as a stand-in for the teacher model’s output distribution, it’s important for the proxy model’s output to align with the teacher model’s output distribution, which is achieved through proxy model alignment. We measure the change in agreement

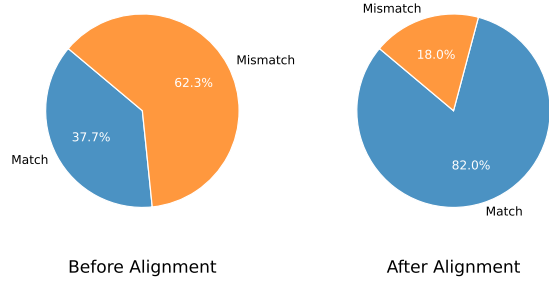


Figure 5: The match ratio between the proxy and teacher’s output tokens before and after alignment. If the top-1 token given by the proxy equals the token given by the teacher in a current step, it is considered a match; otherwise, it is considered a mismatch..

between the top-1 token given by the proxy and the token provided by teacher in current step, before and after alignment. To visualize this alignment, at each step, consider the top-1 token given by the proxy’s output distribution and the token given by the teacher. If the top-1 token given by the proxy matches the token given by the teacher at the current step, it is considered a match; otherwise, it is considered a mismatch. As shown in Figure 5, We find that after the proxy model alignment, the matched portions show a significant upward trend, indicating a trend towards alignment.

A.3 Evaluation of Alignment Performance

We provide detailed analysis on alignment in this section.

Limitation of NLL Loss. The NLL loss (also referred to as SFT loss) provides only positive reinforcement signals. After simple alignment through SFT, the proxy model may still exhibit undesirable distributional characteristics. In contrast, the DPO loss incorporates both positive reinforcement and penalty signals, enabling the proxy distribution to better align with the teacher’s distribution.

Experimental Results. We evaluate the performance of alignment through three aspects.

(1) Student Performance: This study focuses on aligning the proxy model with the black-box teacher to effectively transfer knowledge and facilitate the student’s knowledge distillation. Therefore, we primarily assess the effectiveness of proxy alignment based on the performance of the student model. The results, as shown in Table 2, demonstrate that the proxy model trained with preference loss outperforms the one trained with simple additional training, underscoring the importance of the alignment method.

Table 4: KL Divergence between the distributions of the proxy and the teacher. The student models are based on Llama-2-7B and the proxy model are based on Llama-2-70B.

Method	KL Divergence
Proxy	1.53
Proxy + $\mathcal{L}_{\text{Proxy-NLL}}$	0.87
Proxy + $\mathcal{L}_{\text{Proxy-NLL}}$ + $\mathcal{L}_{\text{Pref}}$	0.56

(2) Proxy Performance: Additionally, we recognize that the current proxy model still exhibits performance gaps compared to the black-box teacher across various datasets. One of the objectives of alignment is to reduce this capability gap between the proxy model and the teacher model. Thus, we also evaluate the effectiveness of proxy alignment by measuring the improvement in the proxy model’s performance across different datasets. As shown in Table 2, aligning using only $\mathcal{L}_{\text{Proxy-NLL}}$ underperforms compared to alignment using $\mathcal{L}_{\text{Proxy-NLL}}$ and $\mathcal{L}_{\text{Pref}}$ on most benchmarks, indicating that incorporating $\mathcal{L}_{\text{Pref}}$ helps reduce the performance gap between the proxy and the black-box teacher on the benchmarks.

(3) KL Divergence: We have supplemented our evaluation by measuring the average KL Divergence between the distributions of the proxy and the teacher on to assess alignment effectiveness. We compare three methods: the proxy without alignment, the proxy aligned using $\mathcal{L}_{\text{Proxy-NLL}}$, and the proxy aligned using $\mathcal{L}_{\text{Proxy-NLL}}$ and $\mathcal{L}_{\text{Pref}}$. The results, as shown in Table 4, indicate that the use of $\mathcal{L}_{\text{Proxy-NLL}}$ and $\mathcal{L}_{\text{Pref}}$ achieves lower KL divergence, demonstrating the effectiveness of preference optimization in improving alignment.

Explore Alignment Method. In this work, we propose the concept of proxy alignment in black-box KD and address this challenge by introducing preference optimization methods (e.g., DPO); however, exploring the effectiveness of alternative preference optimization methods (e.g., SimPO, ORPO) will be a direction for future research.

A.4 Additional Results

We present the performance changes of student models during the distillation in Figure 6 and 7. The student models are based on Llama-2-7B and Llama-1-7B backbone, and the proxy models are based on Llama-2-70B backbone. We plot their accuracy curves on benchmark test sets every 40K training steps in Figure 6 and every 20K train-

ing steps in Figure 7. We compare Proxy-KD with black-box KD method and white-box KD method (Forward KL with Llama-2-70b-chat as white-box teacher). The results show Proxy-KD stands out with the most significant enhancements across benchmarks. Its steeper and more consistent improvement curves in complex tasks like BBH and GSM8K underscore its robust approach for effectively leveraging proxy models in knowledge distillation.

A.5 Extended Evaluation

We incorporate models from different families, such as Qwen, which exhibit distinct architectural and training data characteristics compared to LLaMA. Within the Proxy-KD framework, we use Qwen-2.5-7B as the student model and Qwen-2.5-32B as the proxy model. The results in table 5 indicate that Proxy-KD maintains its performance advantages across different model families. This supports our hypothesis that the proposed proxy alignment and sample-weighted distillation are model-agnostic mechanisms that can generalize beyond LLaMA-based architectures.

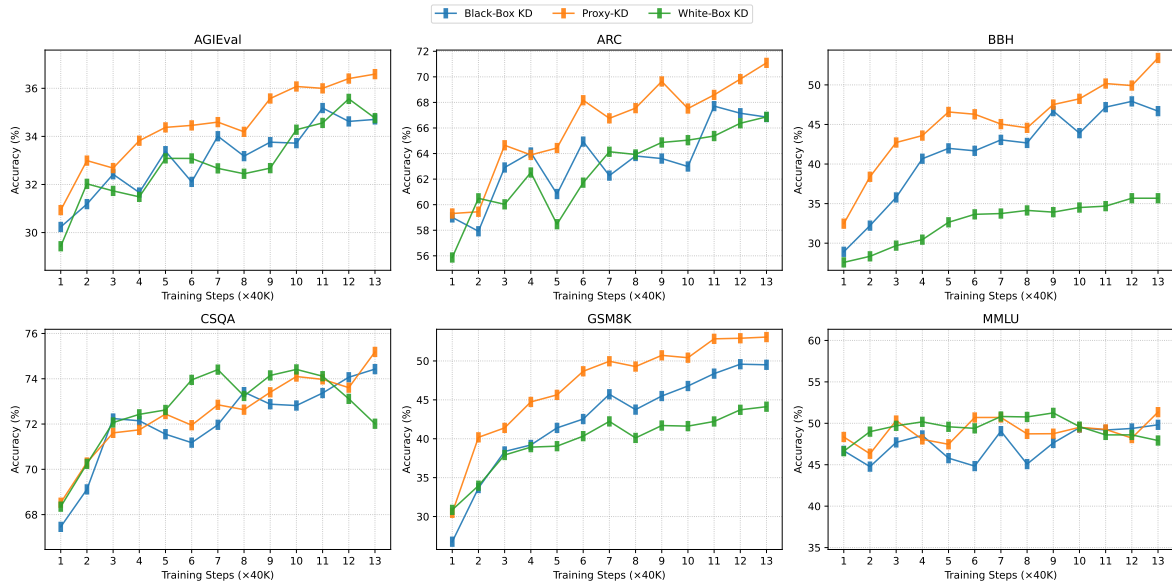


Figure 6: Accuracy curves for student during distillation process. The y-axis is the accuracy on the benchmark test sets, and the x-axis is the number of training steps. We compare Proxy-KD with black-box KD and white-box KD (forward KL) baselines. Notably, Proxy-KD did not show sign of saturation on some benchmarks, such as AGIEval, ARC, and BBH benchmarks.

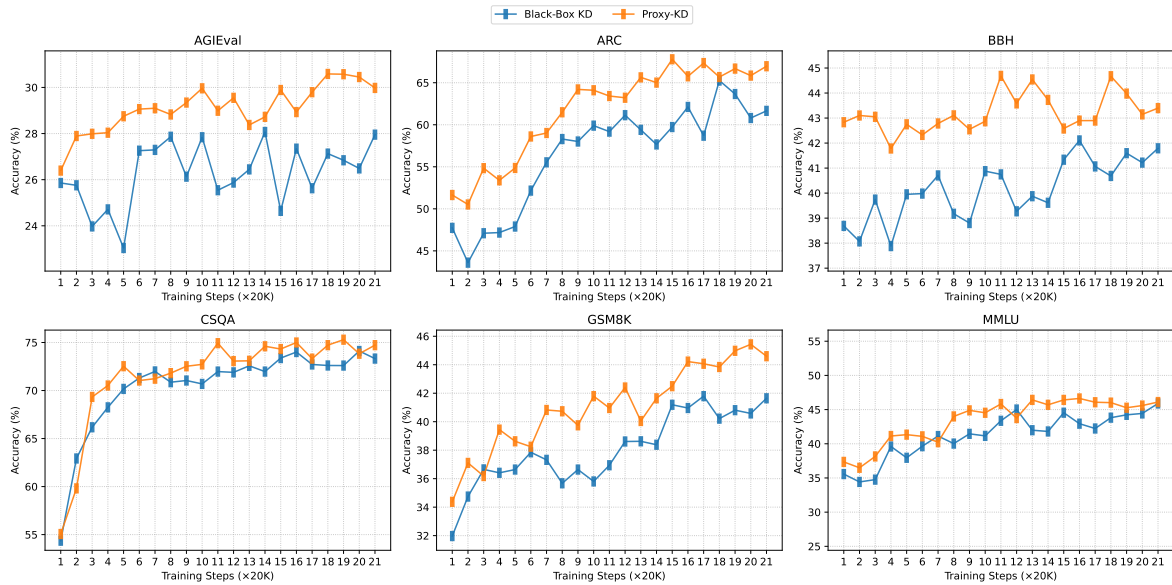


Figure 7: Accuracy curves for student models during knowledge distillation process. The y-axis is the accuracy of students on the benchmark test sets, and x-axis is the number of training steps. We compare Proxy-KD with black-box KD. The students are based on Llama-1-7B, and the proxy is based on Llama-2-70B.

Table 5: Extended Evaluation of Qwen Models in Proxy-KD. We report accuracy (%) for all tasks. All models utilize a zero-shot greedy decoding strategy for evaluation.

Method	Student	Proxy	Teacher	AGIEval	ARC	BBH	CSQA	GSM8K	MMLU	Avg
Forward KL	Qwen-2.5-7B	-	Qwen-2.5-32B	51.34	85.33	58.94	75.12	70.11	68.34	68.20
Black-Box KD	Qwen-2.5-7B	-	GPT-4	52.88	86.53	66.35	74.59	73.47	70.35	70.70
Proxy-KD	Qwen-2.5-7B	Qwen-2.5-32B	GPT-4	53.11	89.64	69.37	75.21	74.56	71.28	72.20