

# Revisiting Knowledge Distillation for Autoregressive Language Models

Anonymous ACL submission

## Abstract

Knowledge distillation (KD) is a common approach to compress a teacher model to reduce its inference cost and memory footprint, by training a smaller student model. However, in the context of autoregressive language models (LMs), we empirically find that larger teachers might dramatically result in a poorer student. In response to this problem, we conduct a series of analyses and reveal that *different tokens have different teaching modes*, neglecting which will lead to performance degradation. Motivated by this, we propose a simple yet effective **adaptive teaching** approach (ATKD) to improve the KD. The core of ATKD is to reduce rote learning and make teaching more diverse and flexible. Extensive experiments on 8 LM tasks show that, with the help of ATKD, various baseline KD methods can achieve consistent and significant performance gains (up to +3.04% average score) across all model types and sizes. More encouragingly, ATKD can improve the student model generalization effectively.

## 1 Introduction

Autoregressive language models (LMs), such as GPT-4 (OpenAI, 2023), PaLM (Chowdhery et al., 2023) and LLaMA2 (Touvron et al., 2023), have achieved great success in a numerous tasks. However, with the scaling of model size, the inference and deployment of these LMs become more computationally expensive and memory intensive, hindering the development of industrial applications. Hence, it is crucial and green to compress these LMs and accelerate the inference, while not losing much performance (Schwartz et al., 2020).

To achieve this goal, a common approach is knowledge distillation (KD), which aims to compress a large teacher model by distilling its knowledge into a small student model (Hinton et al., 2015; Kim and Rush, 2016). Recently, in the context of autoregressive LMs, various novel learning algorithms have been proposed to achieve better

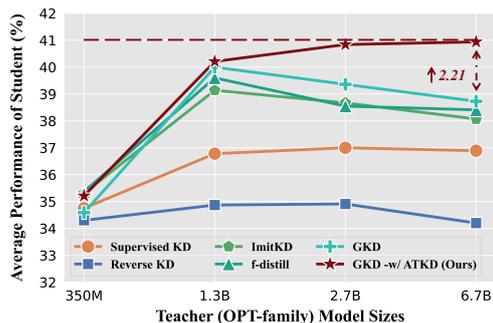


Figure 1: **Comparisons of different KD methods** for distilling the student (OPT-125M). The x-axis denotes the OPT-based teacher sizes, while the y-axis denotes the average performance of students on  $\mathcal{S}_{NLG}$  and  $\mathcal{S}_{NLU}$ . The evaluation details are in §4. Notably, ATKD can be combined with various KD methods, and we only report the results of “GKD + ATKD” for ease of illustration.

distillation performance (Wen et al., 2023; Agarwal et al., 2024). Despite their remarkable performance, we empirically find a **counter-intuitive** phenomenon, where *larger teachers might dramatically result in a poorer student, especially when the model capability gap is large*. As illustrated in Figure 1, the performance of student degrades when the teachers are too large, which is similar to the findings of Mirzadeh et al. (2020); Cho and Hariharan (2019); Zhang et al. (2023).

Although a few works aim to investigate this problem and propose to fill the gap, they are mostly studied for vision models (Mirzadeh et al., 2020; Cho and Hariharan, 2019) or discriminative language understanding models (Zhang et al., 2023), while the autoregressive KD for generative LMs is yet to be explored. In this work, we investigate this problem from the perspective of the distillation objective, which is at the core of autoregressive KD. Specifically, taking the classical token-level KD objective, *i.e.*, *forward KL-Divergence*, as an example, we first reformulate it as two parts: 1) **target-oriented knowledge distillation** (TKD), which en-

forces the student model to learn the target-related information ; 2) **diversity-oriented knowledge distillation** (DKD), which encourages the student to learn more diverse knowledge from the teacher in the non-target classes. These two parts are tied by a token-wise factor, which reflects the teacher’s uncertainty and we denote it as **uncertainty coefficient** (UNC). After reformulating the distillation objective, we conduct a series of preliminary analyses on the popular OPT-family (Zhang et al., 2022) models, and find that:

- ❶ UNC measures the learning difficulties of tokens, where the hard-to-learn ones are more important for KD.
- ❷ DKD contributes more but is greatly suppressed, especially for the larger teachers.
- ❸ TKD plays different roles in tokens with different learning difficulties.

Based on these observations, we can conclude that **different tokens have different teaching modes**, and (one of) the limitations of KD comes from the neglect of this principle. To address this limitation, we propose a simple yet effective adaptive teaching method (referred to as **ATKD**) to improve the KD. The core of ATKD is to reduce rote learning and make teaching more diverse and flexible. Specifically, ATKD skips the target-oriented teaching for the (less-informative) easy-to-learn tokens and pays more attention to the diverse learning of hard-to-learn tokens.

We evaluate ATKD on a variety of LM benchmarks, including 5 language generation tasks and 3 language understanding tasks, upon 3 types of autoregressive LMs: OPT (Zhang et al., 2022), Pythia (Biderman et al., 2023) and LLaMA (Touvron et al., 2023). Results show that ATKD can not only alleviate the problem of performance degradation in larger teachers, but also bring consistent and significant improvements (up to +3.04% average score) into various baseline KD methods among all model types and sizes. Moreover, compared to the standard KD, ATKD can effectively improve the generalization of distilled students.

**Contributions.** To summarize, our contributions are three-fold: (1) Our study reveals that *different tokens have different teaching modes*, neglecting which will cause the sub-optimal distillation performance, especially in larger teachers. (2) We

propose a simple yet effective, plug-and-play approach (ATKD) to alleviate this problem and improve the quality of teaching. (3) Extensive experiments show that ATKD outperforms the standard KD with up to +3.04% average gains and improves the student’s model generalization effectively.

## 2 Rethinking Knowledge Distillation for Autoregressive LMs

In this section, we first delve into the mechanism of classic knowledge distillation and then present the empirical analyses of this strategy in detail.

### 2.1 Recap of Knowledge Distillation

**Notations.** For autoregressive LMs, the classic KD aims to approximately minimize Kullback-Leibler (KL) divergence between the teacher and student output distribution at each token (Hinton et al., 2015). Let  $\mathbf{y} = \{y_1, \dots, y_T\}$  denote the target sequence and  $V$  denote the vocabulary, we refer to  $\mathbf{y}_{<t}$  as  $\{y_1, \dots, y_{t-1}\}$ , where  $t \in \{1, \dots, T\}$  and  $y_t \in V$ . Specifically, the loss function can be formulated as:

$$\begin{aligned} \mathcal{L}_{\text{KL}}(\mathbf{p}||\mathbf{q}) &= - \sum_{t=1}^T \text{KL}(\mathbf{p}(y_t|\mathbf{y}_{<t})||\mathbf{q}(y_t|\mathbf{y}_{<t})) \\ &= - \sum_{t=1}^T \mathbf{p}(y_t|\mathbf{y}_{<t}) \log \left( \frac{\mathbf{p}(y_t|\mathbf{y}_{<t})}{\mathbf{q}(y_t|\mathbf{y}_{<t})} \right), \end{aligned}$$

where  $\mathbf{p} = [p_1, \dots, p_C]$  and  $\mathbf{q} = [q_1, \dots, q_C]^1$  are the predicted distributions of the teacher and student, respectively;  $p_i$  is the probability of the  $i$ -th class and  $C$  is the number of vocabulary  $V$ , KL refers to the KL divergence. For simplicity, we denote  $\mathbf{p}(y_t|\mathbf{y}_{<t})$  as  $\mathbf{p}^t$ , and  $p_i^t$  as the probability of the  $i$ -th class at  $t$ -th step. Here,  $p_i^t$  is determined using a softmax function:

$$p_i^t = \frac{\exp(z_i^t)}{\sum_{j=1}^C \exp(z_j^t)}, \quad (1)$$

where  $z_i^t$  represents the logit of the  $i$ -th class in  $V$ . Let  $g_t$  denote the target token/class at  $t$ -th step, we can obtain the binary probabilities  $\mathbf{p}_{g_t}^t = [p_{g_t}^t, p_{\setminus g_t}^t]$ , where probability of the target class  $p_{g_t}^t$  and non-target classes  $p_{\setminus g_t}^t$  can be calculated as:

$$p_{g_t}^t = \frac{\exp(z_{g_t}^t)}{\sum_{j=1}^C \exp(z_j^t)}, p_{\setminus g_t}^t = \frac{\sum_{k=1, k \neq g_t}^C \exp(z_k^t)}{\sum_{j=1}^C \exp(z_j^t)}.$$

<sup>1</sup>For simplicity, we only consider the formulation of  $\mathbf{p}$  in the following context. Note that the  $\mathbf{q}$  is similar to  $\mathbf{p}$ .

Moreover, for independently analyzing the probabilities among non-target classes, we declare  $\hat{\mathbf{p}}^t = [\hat{p}_1^t, \dots, \hat{p}_{g_t-1}^t, \hat{p}_{g_t+1}^t, \dots, \hat{p}_C^t]$ , where  $\hat{p}_i^t$  is:

$$\hat{p}_i^t = \frac{\exp(z_i^t)}{\sum_{j=1, j \neq g_t}^C \exp(z_j^t)}. \quad (2)$$

**Reformulation of  $\mathcal{L}_{\text{KL}}$ .** Here, we are inspired by Zhao et al. (2022)<sup>2</sup>, and attempt to reformulate  $\mathcal{L}_{\text{KL}}$  with the binary probabilities  $\mathbf{p}_b^t$  and the probabilities among non-target classes  $\hat{\mathbf{p}}^t$ , which can be reformulated as:

$$\mathcal{L}_{\text{KL}} = - \sum_{t=1}^T \left( p_{g_t}^t \log \left( \frac{p_{g_t}^t}{q_{g_t}^t} \right) + \sum_{j=1, j \neq g_t}^C p_j^t \log \left( \frac{p_j^t}{q_j^t} \right) \right). \quad (3)$$

According to Eq. 1 and 2, we have  $p_i^t = \hat{p}_i^t * p_{g_t}^t$ , and can further rewrite Eq. 3 as:

$$\begin{aligned} \mathcal{L}_{\text{KL}} &= - \sum_{t=1}^T \left( p_{g_t}^t \log \left( \frac{p_{g_t}^t}{q_{g_t}^t} \right) \right. \\ &\quad \left. + p_{g_t}^t \sum_{j=1, j \neq g_t}^C \hat{p}_i^t \left( \log \left( \frac{\hat{p}_j^t}{\hat{q}_j^t} \right) + \log \left( \frac{p_{g_t}^t}{q_{g_t}^t} \right) \right) \right) \\ &= - \sum_{t=1}^T \left( p_{g_t}^t \log \left( \frac{p_{g_t}^t}{q_{g_t}^t} \right) + p_{g_t}^t \log \left( \frac{p_{g_t}^t}{q_{g_t}^t} \right) \right. \\ &\quad \left. + p_{g_t}^t \sum_{j=1, j \neq g_t}^C \hat{p}_i^t \log \left( \frac{\hat{p}_j^t}{\hat{q}_j^t} \right) \right) \\ &= - \sum_{t=1}^T \left( \text{KL}(\mathbf{p}_b^t || \mathbf{q}_b^t) + p_{g_t}^t \text{KL}(\hat{\mathbf{p}}^t || \hat{\mathbf{q}}^t) \right). \quad (4) \end{aligned}$$

As seen, we can reformulate the classic KD objective as a combination of binary classification loss on the target class, and KL loss on the non-target classes. The former forces the student to learn the target-related information, and we thus denote it as **target-oriented knowledge distillation** (TKD). Conversely, the latter encourages the student to distill the diverse knowledge among non-target classes, and we denote it as **diversity-oriented knowledge distillation** (DKD). Moreover, we find that TKD and DKD are tied by a token-wise factor  $p_{g_t}^t$ , which could reflect the

<sup>2</sup>Although the reformulation of  $\mathcal{L}_{\text{KL}}$  is inspired by the previous work (Zhao et al., 2022), we take a further step by exploring the potential mechanism of autoregressive KD from the perspective of teaching modes among different tokens, which are our main contributions.

teacher’s uncertainty on the tokens, *i.e.*, the larger  $p_{g_t}^t$  denotes the more uncertainty<sup>3</sup> in the teacher output distribution. Hence, we refer to  $p_{g_t}^t$  as **uncertainty coefficient** (UNC).

## 2.2 Empirical Analyses

**Setting.** We conduct experiments by first fine-tuning larger LMs on the instruction-response dataset  $\mathcal{D}$  as teachers. Then, we use different KD methods to distill a smaller student on  $\mathcal{D}$  with the teacher’s guidance. Here, we use the original OPT-125M as the student and use the other OPT-family models (*i.e.*, OPT-350M/-1.3B/-2.7B/-6.7B) as teachers. Alpaca-GPT4 (Peng et al., 2023) is used as training data, and the models are evaluated on three instruction-following datasets, *i.e.*, DollyEval (Gu et al., 2023), VicunaEval (Chiang et al., 2023) and SelfInst (Wang et al., 2022). We follow (Gu et al., 2023) and use the LLM-based metric, *i.e.*, **LLM-as-a-Judge**, to quantify the model responses. Specifically, we ask GPT-3.5-Turbo-1106<sup>4</sup> to compare model responses with the ground-truth answers and raise 1-10 scores for both responses and report the ratio of the total score of model responses and ground-truth answers.

**Findings.** To reveal the drawbacks of  $\mathcal{L}_{\text{KL}}$  and explore the reasons for performance degradation in large teachers, we conduct systematic analyses to investigate the different effects of UNC, TKD and DKD, respectively. Through the extensive analyses, we empirically observe that:

❶ **UNC measures the learning difficulties of tokens, where the hard-to-learn ones are more important for KD.** Motivated by the token imbalance nature (Piantadosi, 2014) and the truth that different tokens in a sequence contribute differently to the sentence meaning (Church and Hanks, 1990; Chen et al., 2020), we conjecture that different tokens play different roles in autoregressive KD.

Intuitively, the tokens with less uncertainty have simple learning patterns and *easy-to-learn*, while the more uncertain tokens are more informative and are *hard-to-learn*. To verify our conjecture, we rank the training tokens according to the UNC for each mini-batch and evenly split them into two subsets. For clarity, one subset (denoted as “hard-to-learn”) includes samples with top-50% uncertainty, while the remaining samples are in the other subset

<sup>3</sup>For example, the token with  $p_{g_t} = 0.7$  is more uncertain than the one with  $p_{g_t} = 0.1$ .

<sup>4</sup>The analysis of this evaluator is shown in Appendix A.4.

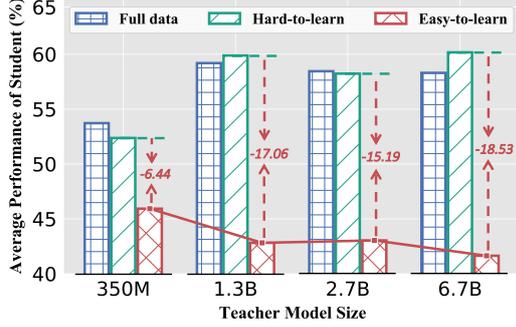


Figure 2: **Comparisons of different training tokens.** The y-axis denotes the average performance of students (OPT-125M) on the evaluated tasks, while the x-axis denotes the sizes of OPT-based teachers.

(denoted as “easy-to-learn”). We train the student model with vanilla  $\mathcal{L}_{KL}$  on different training sets, and illustrate the results in Figure 2.

Obviously, training on the “hard-to-learn” tokens achieves much better performance than on the “easy-to-learn” tokens, and even outperforms the full-data training. This indicates that *tokens with more uncertainty contain more “dark knowledge” and are more important for KD.* Conversely, due to the shallow patterns of easy-to-learn tokens, forcing the student to learn from them might suffer from over-fitting, leading to poorer performance. More interestingly, *this phenomenon seems to be more significant in larger teachers.*

❷ **DKD contributes more (than TKD) but is greatly suppressed, especially for the larger teachers.** Here, we delve into the individual effect of TKD and DKD by comparing the performance of (1) “TKD-only”, (2) “DKD-only” and (3) “TKD+DKD” (where both are decoupled and simply added, *i.e.*, ignoring the effect of UNC). The contrastive results among different training sets (as mentioned in ❶) are listed in Table 1. As seen, “DKD-only” outperforms the “TKD-only” among all model sizes and training sets by a large margin, indicating that the diversity-oriented knowledge is of vital importance to autoregressive KD. However, in Eq. 4, we can find that the effect of DKD is suppressed by the UNC (ranging from 0 to 1), which might lead to the sub-optimal performance. To verify it, we further analyze the distributions of UNC across different model sizes. In practice, we randomly sample 100 instances from the training dataset and illustrate the distributions of UNC in Figure 3. It can be seen that UNC is generally smaller (tends to be 0) in large

Method	350M	1.3B	2.7B	6.7B
1) Full data are used.				
TKD-only	49.19	48.01	47.21	48.29
<b>DKD-only</b>	<b>54.00</b>	<b>57.78</b>	<b>59.43</b>	<b>60.42</b>
<b>TKD+DKD</b>	52.97	57.01	58.66	58.70
2) Easy-to-learn tokens are used.				
TKD-only	39.21	43.82	42.37	41.43
<b>DKD-only</b>	<b>48.68</b>	<b>54.43</b>	<b>58.26</b>	<b>60.02</b>
<b>TKD+DKD</b>	45.59	44.97	45.09	44.66
3) Hard-to-learn tokens are used.				
TKD-only	47.40	45.15	44.63	48.32
DKD-only	51.42	58.51	55.47	59.88
<b>TKD+DKD</b>	<b>53.26</b>	<b>60.49</b>	<b>60.60</b>	<b>61.47</b>

Table 1: **Comparisons of different teaching objectives.** The best results within the same training set are in **bold**.

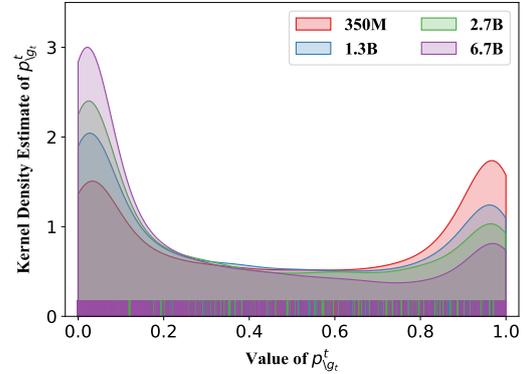


Figure 3: **Illustration of distributions of UNC ( $p_{gt}^t$ )** among different OPT-based teachers on 100 training samples (about 10K tokens). In particular, we use the kernel density estimate for visualizing, where the larger density refers to more tokens.

models than in small models, *i.e.*, the larger models, the more suppressed the effect of DKD. This is also indicated by the results of “TKD+DKD”, as removing the UNC seems to alleviate the performance degradation problem in the large models (except training on easy-to-learn tokens, where the further analyses are shown in ❸). In general, these analyses prove that *DKD is more important but is greatly suppressed by the UNC in the larger models*, which could be the main reason why a larger teacher leads to a poorer student.

❸ **TKD plays different roles in tokens with different learning difficulties.** We can observe an interesting phenomenon in Table 1, where adding TKD upon DKD (“TKD+DKD”) seems to dramatically result in performance degrades when training on the easy-to-learn set, compared to the singly DKD (*e.g.*, decreasing from 60.02% to 44.66%).

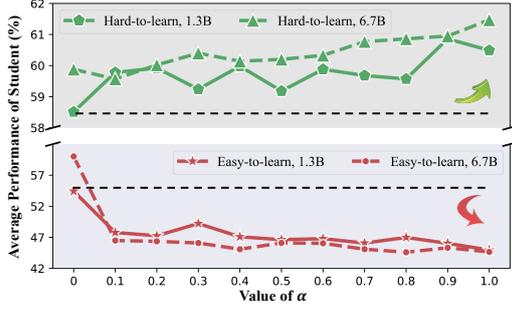


Figure 4: **Effect of TKD in different training tokens.** Here, we report the performance of students distilled with “ $\alpha \times \text{TKD} + \text{DKD}$ ”, where  $\alpha$  is varied from 0 to 1. For ease of illustration, we only illustrate the results of using OPT-1.3B and OPT-6.7B as teachers.

Conversely, in the case of hard-to-learn tokens, adding TKD brings remarkable performance gains. These results motivate us to investigate the special effect of TKD on different tokens, by comparing the performance of different combinations of TKD and DKD in the setting of “ $\alpha \times \text{TKD} + \text{DKD}$ ”. The contrastive performance of varied  $\alpha$  is illustrated in Figure 4. It can be seen that TKD indeed behaves differently in different training sets. TKD hurts the knowledge transfer of easy-to-learn tokens, but is beneficial to the learning of hard-to-learn tokens. We attribute it to the different learning difficulties of tokens, as *the target-oriented learning on easy-to-learn tokens might damage the diversity of students (Tan et al., 2008)*. On the other hand, *adding target-related supervision signals could reduce the learning difficulties on the hard-to-learn tokens*, thus leading to better performance.

### 3 Improving Knowledge Distillation with Adaptive Teaching Modes

Based on the observations in §2, we recognize that *different tokens have different teaching modes*, and the side effect (*i.e.*, problem degrades in larger teachers) of KD mainly comes from the neglect of this principle. To this end, we propose to improve the autoregressive KD with adaptive teaching modes (ATKD). In this section, we introduce the ATKD approach in detail.

**Motivation and Overview of ATKD.** In addition to the empirical findings in §2, our ATKD is also inspired by a famous education initiative (Tan et al., 2008), “Teach Less, Learn More”, which highlights that *reducing rote learning and making education more diverse and flexible can improve the quality of teaching and enhance student learn-*

*ing*. Motivated by this, our ATKD aims to encourage the students to learn from different perspectives for different tokens. In short, ATKD skips the target-oriented teaching for the easy-to-learn tokens, and pays more attention to the learning of diverse knowledge in the hard-to-learn tokens.

To achieve this goal, we should first obtain the easy-/hard-to-learn tokens. As mentioned in ❶ of §2.2, UNC can effectively measure the learning difficulties of tokens, and we thus use it as a metric to select the easy-/hard-to-learn tokens. Specifically, for each mini-batch, we rank the training tokens according to UNC and select the top- $k^5$  tokens as hard-to-learn tokens, while the others are easy-to-learn. Then, ATKD performs the KD processes with adaptive teaching modes as follows:

**Adaptive Teaching Modes of ATKD.** As aforementioned, TKD and DKD contribute differently in easy-/hard-to-learn tokens. Thus, instead of using a unified teaching mode for all tokens, we use adaptive teaching modes for easy-to-learn and hard-to-learn tokens, respectively. Specifically, we decouple the TKD and DKD (*i.e.*, DKD will not be suppressed by the UNC) to enhance the diverse learning of students. Moreover, for the easy-to-learn tokens, considering that the student can easily learn the target-class information, we skip the target-oriented teaching, *i.e.*, removing TKD. On the other hand, both TKD and DKD are used for hard-to-learn tokens, as we empirically found that target-oriented teaching is essential to the learning of hard-to-learn tokens. The learning objectives of different tokens can be formulated as:

$$\begin{aligned} \mathcal{L}_{\text{KL}}^e &= - \sum_{t \in \mathcal{D}_e} \text{KL}(\hat{\mathbf{p}}^t || \hat{\mathbf{q}}^t), \\ \mathcal{L}_{\text{KL}}^h &= - \sum_{t \in \mathcal{D}_h} \text{KL}(\mathbf{p}_b^t || \mathbf{q}_b^t) + \text{KL}(\hat{\mathbf{p}}^t || \hat{\mathbf{q}}^t), \end{aligned}$$

where  $\mathcal{D}_e$  and  $\mathcal{D}_h$  denote the sets of easy-to-learn and hard-to-learn tokens, respectively.

Additionally, since the hard-to-learn tokens contain more informative knowledge and are more important, we adaptively combine the easy-to-learn  $\mathcal{L}_{\text{KL}}^e$  and hard-to-learn  $\mathcal{L}_{\text{KL}}^h$  objectives and formulate the overall learning objective of ATKD as:

$$\mathcal{L}_{\text{KL}}^{\text{all}} = \lambda * \mathcal{L}_{\text{KL}}^e + (1 - \lambda) * \mathcal{L}_{\text{KL}}^h, \quad (5)$$

where  $\lambda$  is a weight factor to balance the different objectives, which is empirically set as 0.2.

<sup>5</sup> $k$  ranges from 0% to 100%, and is set as 50% by default. The analysis of  $k$  can be found in §4.3.

Method	OPT-350M			OPT-1.3B			OPT-2.7B			OPT-6.7B		
	$\mathcal{S}_{NLG}$	$\mathcal{S}_{NLU}$	Avg.									
Teacher	58.33	20.36	<u>39.35</u>	68.90	22.60	<u>45.75</u>	74.21	22.28	<u>48.25</u>	78.71	23.43	<u>51.07</u>
Supervised KD	50.62	18.88	<u>34.75</u>	55.57	17.99	<u>36.78</u>	55.30	18.69	<u>37.00</u>	55.45	18.33	<u>36.89</u>
+ATKD	52.16	19.58	<u>35.87</u>	56.76	19.73	<u>38.25</u>	57.26	19.48	<u>38.37</u>	57.56	19.31	<u>38.43</u>
$\Delta$ ( $\uparrow$ )	+1.54	+0.69	<b>+1.12</b>	+1.20	+1.74	<b>+1.47</b>	+1.96	+0.78	<b>+1.37</b>	+2.11	+0.98	<b>+1.54</b>
Reverse KD	50.54	18.05	<u>34.30</u>	51.60	18.15	<u>34.87</u>	51.26	18.56	<u>34.91</u>	50.08	18.33	<u>34.20</u>
+ATKD	50.86	19.13	<u>34.99</u>	54.40	19.40	<u>36.90</u>	54.34	19.27	<u>36.80</u>	54.37	19.16	<u>36.76</u>
$\Delta$ ( $\uparrow$ )	+0.32	+1.08	<b>+0.70</b>	+2.80	+1.25	<b>+2.03</b>	+3.08	+0.70	<b>+1.89</b>	+4.29	+0.83	<b>+2.56</b>
ImitKD	52.27	18.35	<u>35.31</u>	59.87	18.41	<u>39.14</u>	59.88	17.46	<u>38.67</u>	58.86	17.28	<u>38.07</u>
+ATKD	52.36	18.66	<u>35.51</u>	60.76	19.29	<u>40.02</u>	60.77	19.18	<u>39.97</u>	62.66	19.56	<u>41.11</u>
$\Delta$ ( $\uparrow$ )	+0.09	+0.31	<b>+0.20</b>	+0.89	+0.88	<b>+0.88</b>	+0.89	+1.71	<b>+1.30</b>	+3.80	+2.28	<b>+3.04</b>
f-distill	52.18	18.57	<u>35.37</u>	59.74	19.46	<u>39.60</u>	60.01	17.08	<u>38.55</u>	59.02	17.80	<u>38.41</u>
+ATKD	52.69	18.80	<u>35.75</u>	61.30	19.54	<u>40.42</u>	60.70	19.02	<u>39.86</u>	61.25	19.18	<u>40.22</u>
$\Delta$ ( $\uparrow$ )	+0.51	+0.23	<b>+0.37</b>	+1.55	+0.08	<b>+0.82</b>	+0.68	+1.94	<b>+1.31</b>	+2.23	+1.38	<b>+1.80</b>
GKD	51.87	17.32	<u>34.59</u>	61.23	18.77	<u>40.00</u>	61.24	17.48	<u>39.36</u>	60.59	16.87	<u>38.73</u>
+ATKD	51.90	18.52	<u>35.21</u>	61.36	19.07	<u>40.21</u>	62.46	19.21	<u>40.84</u>	62.62	19.26	<u>40.94</u>
$\Delta$ ( $\uparrow$ )	+0.04	+1.20	<b>+0.62</b>	+0.13	+0.30	<b>+0.21</b>	+1.22	+1.73	<b>+1.48</b>	+2.03	+2.39	<b>+2.21</b>

Table 2: **Results (%) of students (OPT-125M) distilling with different teachers and KD methods.** “Avg.” means the average performance of  $\mathcal{S}_{NLG}$  and  $\mathcal{S}_{NLU}$ . “ $\Delta$  ( $\uparrow$ )” denotes the performance gains of ATKD against the baselines.

## 4 Evaluation

### 4.1 Setup

**Tasks and Datasets.** We conduct extensive experiments on various LM benchmarks, covering a diversity of language generation tasks (denoted as  $\mathcal{S}_{NLG}$ ) and language understanding tasks (denoted as  $\mathcal{S}_{NLU}$ ). Specifically,  $\mathcal{S}_{NLG}$  consists of 5 widely-used generation tasks, *i.e.*, DollyEval (Gu et al., 2023), VicunaEval (Chiang et al., 2023), SelfInst (Wang et al., 2022), Koala (Geng et al., 2023), and WizardLM (Xu et al., 2023) benchmarks.  $\mathcal{S}_{NLU}$  includes 3 popular classification tasks, *i.e.*, MMLU (Hendrycks et al., 2020), Drop (Dua et al., 2019) and BBH (Suzgun et al., 2022). For evaluation, we use the LLM-based metric<sup>6</sup> to quantify the model response for  $\mathcal{S}_{NLG}$  and report the performance with Exact-Match metric for  $\mathcal{S}_{NLU}$ . We report the averaged results over 5 random seeds to avoid stochasticity. The details of all tasks are shown in Appendix A.1.

**Models.** We evaluate ATKD on three types of LMs with various sizes: OPT (Zhang et al., 2022) (*student*: 125M, *teachers*: 350M, 1.3B, 2.7B, 6.7B), Pythia (Biderman et al., 2023) (*student*: 410M, *teachers*: 1.4B, 2.8B), and LLaMA (*student*: 68M (Miao et al., 2023), *teachers*: 1.1B (Zhang et al., 2024), 7B (Touvron et al., 2023)). Alpaca-GPT4 (Peng et al., 2023) consisting of 52K GPT4-

<sup>6</sup>The details of this metric can be found in §2.2.

generated instruction-response pairs is used as training data. The details of training and evaluation can be found in Appendix A.2 and A.3.

**Baselines.** We consider 5 cutting-edge KD baselines in our main experiment: Supervised KD (Hinton et al., 2015), Reverse KD (Gu et al., 2023), ImitKD (Lin et al., 2020), f-distill (Wen et al., 2023) and GKD (Agarwal et al., 2024). For reference, we also report the performance of teachers as the upper bound. We use the codebase of Liu et al. (2023) to implement these baselines and distill students.

### 4.2 Compared Results

Results of distilled models are shown in Table 2 and 3. For ease of illustration, we only report the overall performance of  $\mathcal{S}_{NLG}$  and  $\mathcal{S}_{NLU}$ , respectively, where the detailed results are listed in Table 6 and 9. From these results, we can find that:

**ATKD effectively alleviates the problem of performance degrades in larger teachers.** As seen, various baseline KD methods suffer from this problem, *e.g.*, distilling OPT using GKD (*1.3B*: 40.00% *v.s.* *6.7B*: 38.73%). However, with the help of our ATKD, the students can generally achieve better performance in larger teachers among various baseline KD methods, *i.e.*, alleviating the problem. These results can prove the effectiveness of ATKD in improving the quality of teaching.

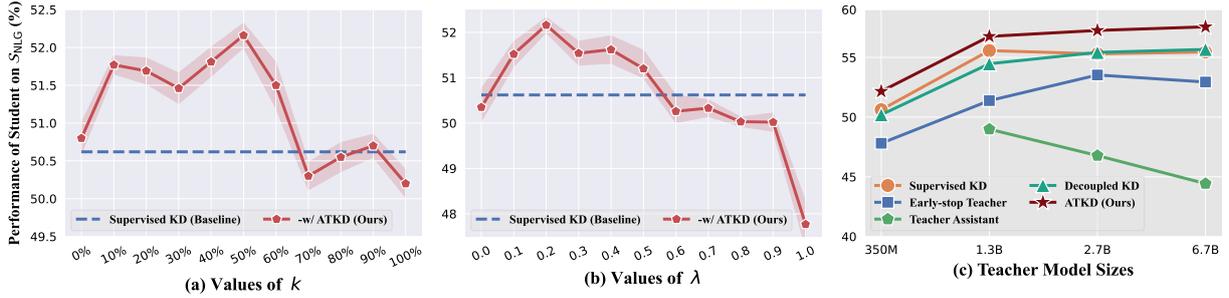


Figure 5: (a) Effect of different ratios (top- $k$ ) for selecting hard-to-learn tokens, (b) Parameter analysis of  $\alpha$  in Eq. 5, and (c) Comparison of different KD methods that aim to alleviate the problem of performance degrades in larger teachers. We use the Supervised KD as the baseline and report the performance of OPT-125M on  $\mathcal{S}_{NLG}$ .

### ATKD brings consistent and significant performance gains among all model sizes and types.

From Table 2, we can see that, compared with the baseline methods, our ATKD consistently achieves better performance (up to +3.04% average gains) across various model sizes. Moreover, as seen in Table 3, in addition to OPT, ATKD also works well in Pythia-family and LLaMA-family models. These results demonstrate the universality of our ATKD and indicate that ATKD has great potential to expand to more LMs.

**ATKD is beneficial to various baseline KD methods.** In the preliminary analyses, we only conducted experiments on the typical Supervised KD. Here, we additionally investigate the combinability of ATKD and other baseline KD methods. As observed in Table 2, ATKD can bring consistent performance gains among all baseline KD methods. For example, with the help of ATKD, Reverse KD and ImitKD achieve +1.80% and +1.36% average performance gains, respectively.

### 4.3 Ablation Study

**Impact of ratio  $k$ .** The ratio  $k$  that is used to select the hard-to-learn tokens, is an important hyperparameter in ATKD. In this study, we analyze its influence by evaluating the performance with different  $k$  spanning from 0% to 100% at 10% intervals on  $\mathcal{S}_{NLG}$  tasks. Figure 5 (a) illustrates the average results, in which we can find that: 1) Too large  $k$  values (e.g., 70%) lead to performance degradation, as many of the selected tokens are “false” hard-to-learn and might distort the adaptive teaching. 2) The model’s performance stably increases between 10% and 50%, and ATKD performs best with  $k = 50%$ , thus leaving as our default settings.

**Impact of coefficient  $\lambda$ .** The factor  $\lambda$  in Eq. 5, which is used to balance different objectives, is also

Method	Pythia-410M		LLaMA-68M	
	1.4B	2.8B	1.1B	7B
Teacher	67.86	73.50	75.23	84.17
Supervised KD	60.66	59.91	30.06	27.94
+ATKD	61.81	61.22	31.19	30.19
$\Delta$ ( $\uparrow$ )	+1.15	+1.31	+1.13	+2.25
Reverse KD	55.92	54.67	26.15	25.94
+ATKD	57.05	57.94	26.73	26.99
$\Delta$ ( $\uparrow$ )	+1.14	+3.27	+0.58	+1.05

Table 3: Results (%) of students (Pythia-410M and LLaMA-68M). Due to the space limitation, we only report the results upon two typical KD baselines.

needed to be investigated. Figure 5 (b) illustrates the results of varied  $\lambda$  ranging from 0 to 1. As seen, compared to the single learning of hard-to-learn tokens, incorporating some supervision signals from easy-to-learn tokens results in better performance. However, too large  $\lambda$  values (e.g., 0.9) would be harmful to the effectiveness of ATKD, as paying much attention to the learning of easy-to-learn tokens might lead to overfitting. More specifically, the case of  $\lambda = 0.2$  performs best, and we thereby use this setting in our experiments.

### 4.4 Discussion

Here, we conduct further analyses to discuss: 1) whether ATKD outperforms the other counterparts, and 2) whether it gains better model generalization.

**Comparison with other counterparts.** To the best of our knowledge, there are no existing KD methods that involve solving the problem of performance degradation for autoregressive LLMs. Thus, we compare ATKD with the related methods in the vision community: “Early-stop Teacher” (Cho and Hariharan, 2019), “Teacher Assistant”<sup>7</sup> (Mirzadeh

<sup>7</sup>We use the OPT-350M as the assistant model and only report the results distilling from teachers larger than 350M.

et al., 2020) and “Decoupled KD” (Zhao et al., 2022). The contrastive results are illustrated in Figure 5 (c), from which we can find that: 1) Suppressing the teacher’s performance via early stopping or leveraging a smaller assistant might not be effective and even lead to worse performance, 2) Although “Decoupled KD” could alleviate this problem, it achieves sub-optimal performance, as it equally adopts the same teaching modes for all tokens. Takeaway: *among all methods, our ATKD can not only alleviate this problem but also bring further performance gains in a simple manner, proving its superiority.*

**Model Generalization.** Enforcing the student to learn more diverse knowledge could improve its generalization. To verify this conjecture, we visualize the loss landscapes of different distilled OPT-125M models on the VicunaEval task. In practice, we follow He et al. (2021); Zhong et al. (2022) to plot the 1D loss curve by linear interpolation between the model weights before (denoted as  $\theta_0$ ) and after (denoted as  $\theta_1$ ) distilling, *i.e.*, “ $\theta_1 + \beta \cdot (\theta_1 - \theta_0)$ ”, where  $\beta$  is a scalar parameter that is ranged from -1 to 1. The 1D visualization results are illustrated in Figure 6, and we find that “-w/ ATKD (Ours)” shows a flatter and optimal property against the baseline Supervised KD. Takeaway: *These results prove that ATKD can smooth the loss landscape and improve the model generalization effectively.*

## 5 Related Works

Recently, autoregressive LMs (OpenAI, 2023; Chowdhery et al., 2023; Touvron et al., 2023) have shown their superior performance by solving various NLP tasks in a generative manner. Despite their success, they usually suffer from unbearable inference latency (Leviathan et al., 2023). To this end, several model compression approaches are proposed to reduce the model size and accelerate the inference (Hinton et al., 2015; Jaszczur et al., 2021; Zhu et al., 2023). Among these efforts, KD strategy (Hinton et al., 2015), which aims at training a smaller student model with the guidance of a teacher model, has attracted great attention recently (Wen et al., 2023; Gu et al., 2023; Agarwal et al., 2024). Although these KD methods realize promising performance when distilling (relatively) smaller LMs, they might fall short in distilling larger LMs (*e.g.*, OPT-6.7B) especially when the student is of a small scale. In fact, this phe-

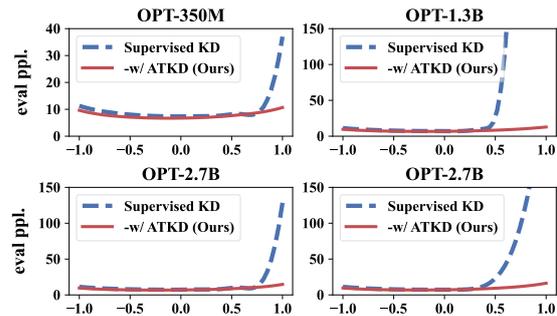


Figure 6: 1D visualization of loss landscapes of OPT-125M distilled by different methods and teachers. The y-axis denotes the model perplexity on VicunaEval.

nomenon has been observed in the vision community (Mirzadeh et al., 2020; Cho and Hariharan, 2019) and language understanding models (Zhang et al., 2023). To alleviate this problem, a few studies including teacher assistant-based (Mirzadeh et al., 2020) and student-friendly (Cho and Hariharan, 2019; Zhao et al., 2022; Zhang et al., 2023) distillation have been recently explored.

The above efforts are generally used for vision models or discriminative LMs, while the autoregressive KD for generative LMs is yet to be explored. To the best of our knowledge, we are the (nearly) first to alleviate the problem of performance degradation in larger autoregressive teacher LMs. Different from the previous methods that aim to directly bridge the performance gap between teacher and student, we attempt to improve the quality of teaching by exploring and addressing the limitations of existing KD objectives.

## 6 Conclusion

In this paper, we reveal and address the limitations of KD in compressing the larger autoregressive teachers. Based on a series of preliminary analyses, we find that equally adopting the same teaching modes for all tokens is sub-optimal, as learning more target-oriented knowledge of the easy-to-learn tokens might lead to overfitting and result in poor performance. To address these limitations, we improve KD with a novel adaptive teaching algorithm. It skips the target-oriented teaching for easy-to-learn tokens and pays more attention to the diverse learning of hard-to-learn tokens. Experiments show that our approach consistently and significantly improves distillation performance across all model architectures. In-depth analyses prove that our approach indeed alleviates the problem, and further improves the model generalization.

## 557 Limitations

558 Our work has several potential limitations. First,  
559 given the limited computational budget, we only  
560 validate our ATKD on up to 7B autoregressive LMs.  
561 It will be more convincing if scaling up to the larger  
562 model size (e.g., 70B) and applying ATKD to more  
563 cutting-edge model architectures. On the other  
564 hand, besides the distillation performance, we be-  
565 lieve that there are still other properties, e.g., train-  
566 ing efficiency and model robustness, of LMs that  
567 can be improved by our ATKD approach, which  
568 are not fully explored in this work.

## 569 Ethics and Reproducibility Statements

570 **Ethics** We take ethical considerations very seri-  
571 ously and strictly adhere to the ACL Ethics Policy.  
572 This paper proposes an adaptive teaching algorithm  
573 to improve existing KD strategies. It aims to com-  
574 press the existing larger LMs into smaller students,  
575 instead of encouraging them to learn privacy knowl-  
576 edge that may cause the ethical problem. Moreover,  
577 all training and evaluation datasets used in this pa-  
578 per are publicly available and have been widely  
579 adopted by researchers. Thus, we believe that this  
580 research will not pose ethical issues.

581 **Reproducibility** In this paper, we discuss the de-  
582 tailed experimental setup, such as hyper-parameters  
583 and statistic descriptions. More importantly, *we*  
584 *have provided our code in the supplementary ma-*  
585 *terials* to help reproduce our experimental results.

## 586 References

587 Rishabh Agarwal, Nino Vieillard, Piotr Stanczyk,  
588 Sabela Ramos, Matthieu Geist, and Olivier Bachem.  
589 2024. [On-policy distillation of language models:  
590 Learning from self-generated mistakes](#). In *ICLR*.

591 Stella Biderman, Hailey Schoelkopf, Quentin Gregory  
592 Anthony, Herbie Bradley, Kyle O’Brien, Eric Hal-  
593 lahan, Mohammad Aflah Khan, Shivanshu Purohit,  
594 USVSN Sai Prashanth, Edward Raff, et al. 2023.  
595 [Pythia: A suite for analyzing large language models  
596 across training and scaling](#). In *ICML*.

597 Kehai Chen, Rui Wang, Masao Utiyama, and Eiichiro  
598 Sumita. 2020. [Content word aware neural machine  
599 translation](#). In *ACL*.

600 Lichang Chen, Shiyang Li, Jun Yan, Hai Wang, Kalpa  
601 Gunaratna, Vikas Yadav, Zheng Tang, Vijay Srini-  
602 vasan, Tianyi Zhou, Heng Huang, et al. 2023. [Al-  
603 pagasus: Training a better alpaca with fewer data](#).  
604 *arXiv preprint*.

605 Yew Ken Chia, Pengfei Hong, Lidong Bing, and Sou-  
606 janya Poria. 2023. [Instructeval: Towards holistic  
607 evaluation of instruction-tuned large language mod-  
608 els](#). *arXiv preprint*.

609 Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng,  
610 Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan  
611 Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al.  
612 2023. [Vicuna: An open-source chatbot impressing  
613 gpt-4 with 90%\\* chatgpt quality](#).

614 Jang Hyun Cho and Bharath Hariharan. 2019. [On the  
615 efficacy of knowledge distillation](#). In *ICCV*.

616 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,  
617 Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul  
618 Barham, Hyung Won Chung, Charles Sutton, Sebas-  
619 tian Gehrman, et al. 2023. [Palm: Scaling language  
620 modeling with pathways](#). *Journal of Machine Learn-  
621 ing Research*.

622 Kenneth Church and Patrick Hanks. 1990. [Word associ-  
623 ation norms, mutual information, and lexicography](#).  
624 *Computational linguistics*.

625 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel  
626 Stanovsky, Sameer Singh, and Matt Gardner. 2019.  
627 [Drop: A reading comprehension benchmark requir-  
628 ing discrete reasoning over paragraphs](#). In *NAACL-  
629 HLT*.

630 Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang,  
631 Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy  
632 Liang, and Tatsunori B Hashimoto. 2023. [Alpaca-  
633 farm: A simulation framework for methods that learn  
634 from human feedback](#). *arXiv preprint*.

635 Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wal-  
636 lace, Pieter Abbeel, Sergey Levine, and Dawn Song.  
637 2023. [Koala: A dialogue model for academic re-  
638 search](#). *Blog post, April*.

639 Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang.  
640 2023. [Knowledge distillation of large language mod-  
641 els](#). *arXiv preprint*.

642 Ruidan He, Linlin Liu, Hai Ye, Qingyu Tan, Bosheng  
643 Ding, Liying Cheng, Jiawei Low, Lidong Bing, and  
644 Luo Si. 2021. [On the effectiveness of adapter-based  
645 tuning for pretrained language model adaptation](#). In  
646 *ACL*.

647 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,  
648 Mantas Mazeika, Dawn Song, and Jacob Steinhardt.  
649 2020. [Measuring massive multitask language under-  
650 standing](#). In *ICLR*.

651 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015.  
652 [Distilling the knowledge in a neural network](#). *arXiv  
653 preprint*.

654 Sebastian Jaszczur, Aakanksha Chowdhery, Afroz Mo-  
655 hiuddin, Lukasz Kaiser, Wojciech Gajewski, Henryk  
656 Michalewski, and Jonni Kanerva. 2021. [Sparse is  
657 enough in scaling transformers](#). *NeurIPS*.

658	Yoon Kim and Alexander M Rush. 2016. <a href="#">Sequence-level knowledge distillation</a> . In <i>EMNLP</i> .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. <a href="#">Llama 2: Open foundation and fine-tuned chat models</a> . <i>arXiv preprint</i> .	710
659			711
660	Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. <a href="#">Fast inference from transformers via speculative decoding</a> . In <i>ICML</i> .		712
661			713
662			714
663	Alexander Lin, Jeremy Wohlwend, Howard Chen, and Tao Lei. 2020. <a href="#">Autoregressive knowledge distillation through imitation learning</a> . In <i>EMNLP</i> .	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananeh Hajishirzi. 2022. <a href="#">Self-instruct: Aligning language model with self generated instructions</a> . <i>arXiv preprint</i> .	715
664			716
665			717
666	Xiaoxuan Liu, Lanxiang Hu, Peter Bailis, Ion Stoica, Zhijie Deng, Alvin Cheung, and Hao Zhang. 2023. <a href="#">Online speculative decoding</a> . <i>arXiv preprint</i> .		718
667			719
668		Yuqiao Wen, Zichao Li, Wenyu Du, and Lili Mou. 2023. <a href="#">f-divergence minimization for sequence-level knowledge distillation</a> . In <i>ACL</i> .	720
669	Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Rae Ying Yee Wong, Zhuoming Chen, Daiyaan Arfeen, Reyna Abhyankar, and Zhihao Jia. 2023. <a href="#">Specinfer: Accelerating generative llm serving with speculative inference and token tree verification</a> . <i>arXiv preprint</i> .		721
670			722
671		Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. <a href="#">Wizardlm: Empowering large language models to follow complex instructions</a> . <i>arXiv preprint</i> .	723
672			724
673			725
674			726
675	Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. 2020. <a href="#">Improved knowledge distillation via teacher assistant</a> . In <i>AAAI</i> .	Chen Zhang, Yang Yang, Jiahao Liu, Jingang Wang, Yunsen Xian, Benyou Wang, and Dawei Song. 2023. <a href="#">Lifting the curse of capacity gap in distilling language models</a> . In <i>ACL</i> .	728
676			729
677			730
678			731
679	OpenAI. 2023. <a href="#">Gpt-4 technical report</a> .	Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. <a href="#">Tinyllama: An open-source small language model</a> . <i>arXiv preprint</i> .	732
680	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. <a href="#">Instruction tuning with gpt-4</a> . <i>arXiv preprint</i> .		733
681			734
682		Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. 2022. <a href="#">Opt: Open pre-trained transformer language models</a> . <i>arXiv preprint</i> .	735
683	Steven T Piantadosi. 2014. <a href="#">Zipf’s word frequency law in natural language: A critical review and future directions</a> . <i>Psychonomic bulletin &amp; review</i> .		736
684			737
685			738
686	Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. 2020. <a href="#">Green ai</a> . <i>Communications of the ACM</i> .	Borui Zhao, Quan Cui, Renjie Song, Yiyu Qiu, and Jiajun Liang. 2022. <a href="#">Decoupled knowledge distillation</a> . In <i>CVPR</i> .	740
687			741
688			742
689	Andrea Sottana, Bin Liang, Kai Zou, and Zheng Yuan. 2023. <a href="#">Evaluation metrics in the era of GPT-4: reliably evaluating large language models on sequence to sequence tasks</a> . <i>arXiv preprint</i> .	Jiaxu Zhao, Meng Fang, Shirui Pan, Wenpeng Yin, and Mykola Pechenizkiy. 2023. <a href="#">GPTBIAS: A comprehensive framework for evaluating bias in large language models</a> . <i>arXiv preprint</i> .	743
690			744
691			745
692			746
693	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2023. <a href="#">Beyond the imitation game: Quantifying and extrapolating the capabilities of language models</a> . <i>Transactions on Machine Learning Research</i> .	Qihuang Zhong, Liang Ding, Li Shen, Peng Mi, Juhua Liu, Bo Du, and Dacheng Tao. 2022. <a href="#">Improving sharpness-aware minimization with fisher mask for better generalization on language models</a> . In <i>Findings of EMNLP</i> .	747
694			748
695			749
696			750
697			751
698		Miaoxi Zhu, Qihuang Zhong, Li Shen, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2023. <a href="#">Zero-shot sharpness-aware quantization for pre-trained language models</a> . In <i>EMNLP</i> .	752
699			753
700	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. <a href="#">Challenging big-bench tasks and whether chain-of-thought can solve them</a> . <i>arXiv preprint</i> .		754
701			755
702		<b>A Appendix</b>	756
703		<b>A.1 Details of Tasks and Datasets</b>	757
704		In this work, we conduct extensive experiments on several language generation and understanding tasks. Here, we introduce the descriptions of these tasks and datasets in detail. Firstly, we present the	758
705			759
706	Kelvin HK Tan, Charlene Tan, and Jude SM Chua. 2008. <a href="#">Innovation in education: The "teach less, learn more" initiative in singapore schools</a> . <i>Innovation in education</i> .		760
707			761
708			
709			

Test set	Task	# Types	# Samples
$\mathcal{S}_{\text{NLG}}$	DollyEval	Generation	500
	VicunaEval	Generation	80
	SelfInst	Generation	242
	Koala	Generation	180
	WizardLM	Generation	218
$\mathcal{S}_{\text{NLU}}$	MMLU	Classification	14,079
	Drop	Classification	9,540
	BBH	Classification	6,511

Table 4: **Statistics of all test sets** used in this paper.

statistics of all evaluated datasets in Table 4. Then, each task is described as:

**DollyEval.** DollyEval (Gu et al., 2023) is a 500-sample test set that is splitted from the databricks-dolly-15k<sup>8</sup> dataset.

**VicunaEval.** VicunaEval (Chiang et al., 2023) contains 80 challenging questions used in the Vicuna evaluation.

**SelfInst.** SelfInst (Wang et al., 2022) is a user-oriented instruction-following test set with 252 samples.

**Koala.** This test set consists of 180 queries that Geng et al. (2023) source from publicly available user-written language model prompts.

**WizardLM.** WizardLM (Xu et al., 2023) consists of 218 instances, each of which is an instruction for a specific skill, such as Math, Reasoning, Complex Formats, and so on.

**MMLU.** Massive Multitask Language Understanding (MMLU) (Hendrycks et al., 2020) is a popular benchmark designed to measure the multi-task accuracy of LLMs, covering 57 tasks.

**Drop.** Discrete Reasoning Over Paragraphs (DROP) (Dua et al., 2019) is a math-based reading comprehension task that requires a system to perform discrete reasoning over passages extracted from Wikipedia articles.

**BBH.** BIG-Bench Hard (BBH) (Suzgun et al., 2022) is a subset of 23 challenging tasks from the BIG-Bench benchmark (Srivastava et al., 2023), which focuses on tasks believed to be beyond the capabilities of current language models.

## A.2 Training Hyper-parameters.

For teachers, we train each model with a batch size of 128 and a peak learning rate of  $2e-5$ . For distilling students, the learning rate is selected in

<sup>8</sup><https://github.com/databrickslabs/dolly/tree/master>

Evaluator	Method	350M	1.3B	2.7B	6.7B
ChatGPT	Supervised KD	46.93	51.92	53.02	53.78
	+ATKD	<b>52.75</b>	<b>52.99</b>	<b>53.69</b>	<b>54.74</b>
GPT-4	Supervised KD	30.09	32.28	32.45	33.28
	+ATKD	<b>32.48</b>	<b>33.21</b>	<b>33.53</b>	<b>34.49</b>

Table 5: **Comparison between ChatGPT-based and GPT-4-based automatic evaluators.** Here, we report the evaluation results of students (OPT-125M) on the Koala benchmark, and we can see that ChatGPT makes similar judgments to GPT-4.

{ $2e-4$ ,  $2e-5$ } depending on model sizes, while the batch size is 256 and the maximum tokenizer length is 512. All models are trained for 3 epochs, and all experiments are conducted on 8 NVIDIA A800 (80GB) GPUs.

## A.3 Evaluation Details.

For  $\mathcal{S}_{\text{NLG}}$ , we report the zero-shot performance by directly evaluating the instruction-following responses using the **LLM-as-judge** metric<sup>9</sup>. We use the same evaluation prompt in Gu et al. (2023) to instruct the ChatGPT to judge the usefulness of model responses. Notably, for each query in  $\mathcal{S}_{\text{NLG}}$ , we set the maximum number of output tokens as 256. As for  $\mathcal{S}_{\text{NLU}}$ , we follow Chen et al. (2023) and use the code provided by Chia et al. (2023) to conduct benchmark evaluation. Specifically, we use 5-shot direct prompting and measure the exact-match score for MMLU (Hendrycks et al., 2020). Regarding the Drop (Dua et al., 2019) and BBH (Suzgun et al., 2022), 3-shot direct prompting is used and exact-match scores are reported.

## A.4 ChatGPT v.s. GPT-4

Although the GPT-4 is more commonly used as the automatic evaluator for the “LLM-as-Judge” metric (Chen et al., 2023; Chiang et al., 2023), it requires a much higher cost, especially for our extensive experiments. As an alternative, we use the cheaper ChatGPT as the automatic evaluator to evaluate the model responses. Here, to verify whether ChatGPT is enough to reflect the behavior of LMs, we conduct a comparative study on ChatGPT and GPT-4. Specifically, taking the responses of OPT-125M on Koala as an example, we

<sup>9</sup>Although some studies show that LLM-as-Judge may exhibit a certain degree of bias (Zhao et al., 2023; Sottana et al., 2023), strong proprietary LLMs, e.g., ChatGPT and GPT-4, are capable of making preference determinations that are highly consistent with those of human annotators (Dubois et al., 2023).

831 use the ChatGPT and GPT-4 to measure the score,  
832 respectively. As listed in Table 5, GPT-4 seems to  
833 be more strict in evaluating the model responses, as  
834 the evaluated scores of GPT-4 are generally lower  
835 than those of ChatGPT. Nevertheless, both auto-  
836 matic evaluators make similar judgments, *i.e.*, our  
837 ATKD performs better than baselines among all  
838 model sizes. Thus, we believe that *ChatGPT is*  
839 *enough to reflect whether the model generates a*  
840 *useful response, and it is credible to use ChatGPT*  
841 *as the automatic evaluator in this study.*

Method	$S_{NLG}$					$S_{NLU}$			Average	
	DollyEval	VicunaEval	SelfInst	Koala	WizardLM	MMLU	Drop	BBH	$S_{NLG}$	$S_{NLU}$
SFT -w/o KD	55.05	38.45	52.52	45.27	42.35	21.66	3.8	27.5	49.75	17.65
<b>Teacher-OPT-350M</b>	64.96	51.09	61.98	52.13	46.84	26.03	6.98	28.08	58.33	20.36
Supervised KD	54.90	45.93	53.86	46.93	41.98	24.44	4.85	27.36	50.62	18.88
+ATKD	56.14	43.35	52.98	52.75	44.88	24.52	6.94	27.27	<b>52.16</b>	<b>19.58</b>
Reverse KD	55.53	44.30	50.25	50.14	42.04	22.54	4.88	26.73	50.54	18.05
+ATKD	54.81	43.68	52.66	51.05	42.26	23.74	6.95	26.70	<b>50.86</b>	<b>19.13</b>
ImitKD	55.68	41.30	54.73	52.51	45.55	24.06	4.25	26.73	52.27	18.35
+ATKD	55.10	43.64	55.37	52.49	45.84	25.07	4.39	26.51	<b>52.36</b>	<b>18.66</b>
f-distill	56.31	43.52	52.67	52.69	44.93	24.71	4.50	26.49	52.18	18.57
+ATKD	54.86	42.61	57.22	53.00	46.15	24.60	5.04	26.76	<b>52.69</b>	<b>18.80</b>
GKD	53.76	44.41	53.82	54.62	45.83	23.93	1.42	26.61	51.87	17.32
+ATKD	54.43	44.35	53.88	54.79	44.31	25.40	2.29	27.88	<b>51.90</b>	<b>18.52</b>
<b>Teacher-OPT-1.3B</b>	72.29	68.86	74.35	65.02	58.30	24.78	14.00	29.01	68.90	22.60
Supervised KD	60.89	52.35	57.95	51.92	44.92	22.27	4.57	27.13	55.57	17.99
+ATKD	62.35	51.52	59.59	52.99	45.86	25.08	6.43	27.67	<b>56.76</b>	<b>19.73</b>
Reverse KD	57.16	46.36	50.75	50.10	42.94	23.02	4.22	27.21	51.60	18.15
+ATKD	59.08	48.41	57.17	52.04	44.71	26.06	5.44	26.71	<b>54.40</b>	<b>19.40</b>
ImitKD	64.55	50.74	61.99	59.15	50.73	23.45	4.31	27.47	59.87	18.41
+ATKD	65.27	53.70	63.41	60.00	50.70	25.76	4.90	27.20	<b>60.76</b>	<b>19.29</b>
f-distill	64.80	51.45	61.57	59.00	49.78	26.59	4.71	27.08	59.74	19.46
+ATKD	65.72	51.56	62.96	60.72	53.35	26.58	4.84	27.21	<b>61.30</b>	<b>19.54</b>
GKD	63.48	56.08	64.73	61.54	53.83	25.99	4.42	25.89	61.23	18.77
+ATKD	64.84	56.75	64.43	60.66	52.25	25.69	4.69	26.82	<b>61.36</b>	<b>19.07</b>
<b>Teacher-OPT-2.7B</b>	75.64	74.43	80.99	74.12	63.39	24.74	12.86	29.25	74.21	22.28
Supervised KD	59.16	52.89	58.31	53.02	45.88	22.89	5.63	27.56	55.30	18.69
+ATKD	62.47	54.47	60.22	53.69	46.01	23.83	6.48	28.12	<b>57.26</b>	<b>19.48</b>
Reverse KD	56.09	48.58	49.46	51.07	43.34	24.08	4.23	27.38	51.26	18.56
+ATKD	59.79	50.96	55.73	50.70	44.54	24.65	5.76	27.39	<b>54.34</b>	<b>19.27</b>
ImitKD	63.30	57.55	62.98	59.23	50.01	22.82	4.50	25.07	59.88	17.46
+ATKD	65.04	57.27	63.11	59.93	50.37	25.11	6.17	26.25	<b>60.77</b>	<b>19.18</b>
f-distill	63.78	58.58	62.79	58.57	50.00	22.21	4.40	24.63	60.01	17.08
+ATKD	64.45	57.00	63.03	59.92	51.49	24.57	5.33	27.17	<b>60.70</b>	<b>19.02</b>
GKD	64.13	57.42	64.41	63.59	50.56	22.78	3.42	26.24	61.24	17.48
+ATKD	66.84	60.73	63.23	63.02	51.72	25.42	4.57	27.65	<b>62.46</b>	<b>19.21</b>
<b>Teacher-OPT-6.7B</b>	81.03	77.38	84.92	78.65	67.01	24.67	15.16	30.45	78.71	23.43
Supervised KD	60.01	49.41	58.22	53.78	45.51	23.46	5.43	26.10	55.45	18.33
+ATKD	63.08	53.75	60.05	54.74	45.84	24.23	5.95	27.74	<b>57.56</b>	<b>19.31</b>
Reverse KD	53.73	47.33	49.70	49.50	43.61	23.95	4.30	26.73	50.08	18.33
+ATKD	59.13	52.24	57.63	52.21	42.38	25.62	4.80	27.05	<b>54.37</b>	<b>19.16</b>
ImitKD	62.32	57.64	63.02	57.08	48.24	22.59	4.02	25.23	58.86	17.28
+ATKD	65.07	58.07	65.93	63.76	54.29	25.89	6.68	26.11	<b>62.66</b>	<b>19.56</b>
f-distill	63.25	55.97	62.06	57.23	48.56	24.25	4.03	25.12	59.02	17.80
+ATKD	64.51	59.48	64.04	62.28	50.48	25.15	5.57	26.82	<b>61.25</b>	<b>19.18</b>
GKD	64.37	58.47	61.63	62.19	50.23	22.03	3.53	25.04	60.59	16.87
+ATKD	66.68	60.87	65.29	63.19	50.51	25.84	4.36	27.58	<b>62.62</b>	<b>19.26</b>

Table 6: Full results of Table 2, *i.e.*, performance of student (OPT-125M) on  $S_{NLG}$  and  $S_{NLU}$  across different teachers and KD methods. “Average” denotes the average results of  $S_{NLG}$  and  $S_{NLU}$ , and “SFT -w/o KD” refers to the results of the vanilla student that is tuned on the ground-truth data. Better results among baseline KD methods and ours are in **bold**.

Method	$S_{\text{NLG}}$ , Pythia-410M				
	Dolly	Vicuna	SelfInst	Koala	WizardLM
SFT-w/o KD	61.81	57.38	60.62	50.67	50.06
<b>Teacher-1.4B/-1.1B</b>	69.51	73.13	69.59	65.17	62.44
Supervised KD	62.62	64.61	63.47	56.36	55.15
+ATKD (Ours)	63.87	62.70	65.64	59.13	54.72
Reverse KD	58.82	56.17	58.87	53.17	48.15
+ATKD (Ours)	61.14	57.80	57.30	53.85	49.77
<b>Teacher-2.8B/-7B</b>	75.84	76.63	72.99	70.95	69.63
Supervised KD	61.10	60.38	63.51	56.99	55.43
+ATKD (Ours)	63.37	64.31	63.06	59.22	54.79
Reverse KD	58.80	54.99	53.74	52.61	47.79
+ATKD (Ours)	61.21	63.23	56.38	57.23	50.81

Table 7: Results of Pythia-410M.

Method	$S_{\text{NLG}}$ , LLaMA-68M				
	Dolly	Vicuna	SelfInst	Koala	WizardLM
SFT-w/o KD	26.37	26.67	28.37	27.27	23.72
<b>Teacher-1.1B/7B</b>	78.82	77.50	75.02	72.01	69.08
Supervised KD	29.63	28.74	31.51	31.84	28.45
+ATKD (Ours)	30.29	29.50	34.65	33.25	28.34
Reverse KD	25.69	25.74	27.98	29.23	22.78
+ATKD (Ours)	26.02	25.31	28.95	29.23	24.35
<b>Teacher-2.8B/7B</b>	86.50	83.25	83.70	84.18	79.68
Supervised KD	27.27	28.31	29.26	30.02	26.15
+ATKD (Ours)	30.08	28.95	30.97	32.09	28.47
Reverse KD	25.65	25.11	27.85	28.03	23.05
+ATKD (Ours)	26.70	27.38	28.35	29.53	23.91

Table 8: Results of LLaMA-68M.

Table 9: Full results of Table 3, *i.e.*, performance of students (Pythia-410M, Table 7 and LLaMA-68M, Table 8) on  $S_{\text{NLG}}$ . Notably, for Pythia-410M, we use the Pythia-1.4B/2.8B as teachers, while LLaMA-1.1B/7B are used as teachers for LLaMA-68M.