# UNDERSTANDING KNOWLEDGE DISTILLATION IN POST-TRAINING: WHEN IT HELPS AND WHEN IT FAILS

Anonymous authors
Paper under double-blind review

# **ABSTRACT**

Large language models (LLMs) achieve strong performance across many tasks, but their high computational cost limits deployment in resource-constrained environments. Knowledge Distillation (KD) offers a practical solution by transferring knowledge from a teacher model of a larger size to a smaller student model. While prior work has mainly examined task-specific or small-scale settings, the posttraining stage for building general instruction-following models has received limited attention. In this paper, we conduct a systematic study of KD in post-training using the large-scale Tulu 3 dataset. We find that KD outperforms supervised finetuning (SFT) in low-data regimes, but its advantage diminishes as more training data is added. Distilling from a stronger instruction-tuned teacher restores substantial gains even with abundant data, indicating that KD remains effective when the teacher provides knowledge that the student cannot easily acquire from the training data alone. We further study domain-specific, low-resource scenarios and propose a two-stage KD strategy that leverages synthetic teacher-labeled data followed by refinement on human annotations. This method consistently improves student performance, providing practical guidance for building compact models in data-scarce environments.

# 1 Introduction

Large Language Models (LLMs) have brought significant advancements to natural language processing, achieving state-of-the-art performance across a wide range of tasks (OpenAI, 2023; Yang et al., 2025; DeepSeek-AI et al., 2025). However, deploying these models in resource-constrained environments, such as mobile phones and edge devices, remains a considerable challenge due to their high computational and memory demands. To address this issue, model compression techniques, particularly *Knowledge Distillation* (KD), have gained substantial attention as a practical solution for improving efficiency without severely compromising performance.

KD transfers knowledge from a large, over-parameterized *teacher* model to a smaller, more efficient *student* model by encouraging the student to mimic the teacher's output distribution or internal representations (Hinton et al., 2015). This approach allows the student model to achieve competitive performance while significantly reducing resource consumption. Consequently, KD has been widely explored in the context of LLMs, yielding promising results.

Several KD methods have been proposed to enhance its effectiveness for generative language models. SeqKD (Kim & Rush, 2016) encourages the student to imitate the output sequences of the teacher directly. MiniLLM (Gu et al., 2024) replaces the commonly used forward Kullback-Leibler divergence (KLD) with reverse KLD, which is better suited to sequence generation tasks. GKD (Agarwal et al., 2024) introduces a generalized KD framework that supports a range of divergence measures, such as generalized Jensen-Shannon divergence, and reduces train-inference mismatch by incorporating on-policy samples from the student. Most recently, Direct Preference Knowledge Distillation (DPKD) (Li et al., 2024) reformulates KD as a direct preference learning problem, supplementing KL divergence with an implicit reward signal to better align the student with teacher preferences.

While these approaches demonstrate strong performance, they are typically applied in task-specific or small-scale settings. A relatively underexplored but increasingly important scenario is the *post-training* setting, where a student model is trained to acquire general instruction-following capabilities from a teacher model. This setting is particularly relevant for building smaller, more efficient models that can follow human instructions across diverse domains, yet existing work has provided limited insight into the behavior and effectiveness of KD in this context.

In this paper, we conduct a comprehensive study of KD methods applied in the post-training stage of LLM development. We focus on understanding their effectiveness across different training data scales. To this end, we utilize the large-scale instruction-following dataset Tulu 3 (Lambert et al., 2024), which contains 939k high-quality instruction-response pairs. Using this dataset, we train both teacher and student models, and apply KD using subsets of varying sizes.

Our findings reveal that KD provides clear performance benefits over supervised fine-tuning (SFT) in low-data regimes. However, as the size of the training dataset increases, the performance gap between KD and SFT narrows substantially, and KD offers little additional gain. This suggests that KD does not scale effectively to large-data settings, as the student can already recover most of the teacher's knowledge through direct supervision.

To further test this hypothesis, we replace the original teacher model with a stronger, instruction-tuned LLM (e.g., Llama3.3-70B-Instruct) trained on a much larger and more diverse corpus via reinforcement learning from human feedback (RLHF). We find that distillation from this stronger teacher significantly improves the student's performance, even in the large-data setting, highlighting that KD remains effective when the teacher possesses knowledge that the student cannot easily acquire from the data alone.

While our primary focus is on post-training KD for general instruction-following, real-world deployment often involves domain-specific applications, such as translation, summarization, or scientific QA, where high-quality labeled data is scarce. In such cases, models trained with limited supervision are prone to underfitting, and leveraging a stronger teacher becomes particularly valuable. Although KD has been applied in various such contexts, there is a lack of systematic study on its effectiveness across different low-resource domains and the optimal strategies to employ in these data-scarce scenarios.

Motivated by this gap, we further investigate the application of KD in domain adaptation settings. We propose a two-stage training paradigm that first uses teacher-annotated synthetic data to broaden the student's exposure to diverse instruction styles, and then refines the model with KD on the small set of high-quality human annotations. Our experiments show that this strategy consistently improves performance across multiple domain-specific tasks, demonstrating that carefully designed KD pipelines can substantially benefit compact student models in data-scarce environments.

In summary, our contributions are threefold:

- We present a comprehensive study of knowledge distillation in the post-training stage of LLMs, systematically evaluating its effectiveness across different data scales.
- We identify the scaling limitation of KD when the student and teacher are trained on the same dataset, and demonstrate that distilling from a stronger instruction-tuned teacher can still provide substantial benefits.
- We introduce a two-stage KD strategy leveraging synthetic data for domain-specific, low-resource scenarios, which consistently improves student performance and provides practical guidance for real-world deployment.

# 2 RELATED WORK

Knowledge Distillation and Instruction Tuning. Knowledge Distillation (KD) transfers knowledge from a large teacher model to a smaller student by matching output distributions or sequences (Hinton et al., 2015). Early methods such as SeqKD (Kim & Rush, 2016) focused on sequence-level imitation, while MiniLLM (Gu et al., 2024) introduced reverse-KL objectives better suited for generation. GKD (Agarwal et al., 2024) addressed train—inference mismatch via on-policy samples, and DPKD (Li et al., 2024) reformulated distillation as preference optimization. Despite their effectiveness, these approaches are typically applied in task-specific or small-scale settings, leaving the post-training stage—critical for building general instruction-following models—underexplored. In

contrast, recent instruction-tuning efforts such as InstructGPT (Ouyang et al., 2022), Alpaca (Taori et al., 2023), OpenAssistant (Köpf et al., 2023), and Tülu 3 (Lambert et al., 2024) highlight the importance of alignment after pretraining, yet rely mostly on supervised fine-tuning or RLHF rather than KD. Our work bridges these directions by systematically studying KD in the post-training stage, evaluating its effectiveness across data scales and highlighting scenarios where stronger teachers remain beneficial.

Task-Specific KD and Synthetic Data. Beyond general instruction-following, KD has been widely applied in domain-specific and low-resource scenarios, such as translation, summarization, and QA (Kim & Rush, 2016). More recently, Speculative KD (Xu et al., 2025) improves efficiency and robustness by interleaving student and teacher generation. While these studies confirm the usefulness of KD under limited supervision, they primarily focus on single tasks rather than systematic analysis across domains. Another line of work explores synthetic data as a complementary supervision source: teacher-generated corpora can boost student performance (Shirgaonkar et al., 2024), and even self-training with student generations can yield competitive results (Lewis et al., 2025). However, naive mixing of synthetic and human-annotated data often introduces noise that harms performance. Our work addresses this challenge by proposing a two-stage KD strategy that leverages synthetic data as a warm-up before distillation on gold annotations, demonstrating consistent improvements in domain-specific, low-resource settings.

# 3 Post-Training Knowledge Distillation for LLMs

# 3.1 PROBLEM FORMULATION

We investigate KD in the post-training stage of LLM development, where both the teacher and student are non-intruct tuned language models. The goal is to assess whether KD can effectively transfer general instruction-following capabilities from a teacher to a student when both are fine-tuned on the same dataset.

Formally, let T denote the teacher model and S the smaller student model. Given a dataset  $\mathcal{D} = \{(x_i, y_i)\}$  of instruction-response pairs, we compare two training paradigms for the student: (1) supervised fine-tuning (SFT) directly on  $\mathcal{D}$ , and (2) KD from a teacher  $T_s$ , which is first trained on  $\mathcal{D}$  via SFT. We aim to evaluate whether the student can benefit from distillation beyond what is learned from direct supervision alone, and how this benefit varies with the size of  $\mathcal{D}$ .

### 3.2 EXPERIMENTAL SETUP

**Dataset and Evaluation Tasks** We conduct experiments using the Tulu 3 dataset (Lambert et al., 2024), which contains 939k high-quality instruction-response pairs. We split this dataset into subsets of varying sizes to evaluate the effectiveness of KD across different data scales. In specific, we create subsets of sizes ranging from 10k samples to the full training set. Note that both the teacher and student models are trained on the same subset.

For evaluation, we use five diverse benchmarks that collectively assess reasoning, scientific knowledge, and instruction-following capabilities: BBH (Srivastava et al., 2023), GPQA (Rein et al., 2023), IFEval (Zhou et al., 2023), InFoBench (Qin et al., 2024), and MMLU-Pro (Wang et al., 2024). These benchmarks cover a wide range of domains, from general reasoning to domain-specific scientific QA and fine-grained instruction following. A detailed description of each benchmark is provided in Appendix A.3.

We report the average accuracy for BBH, GPQA, and MMLU-Pro. For IFEval, we report the prompt-level loose-accuracy, which measures the percentage of prompts for which the model's response satisfies at least one of the constraints specified in the prompt. For InFoBench, we report the Decomposed Requirements Following Ratio (DRFR) (Qin et al., 2024), which measures the percentage of requirements satisfied by the model responses. We use the official evaluation scripts provided by the respective benchmarks to compute these metrics.

**Models** We use the Llama3.1-70B model (Dubey et al., 2024) as our teacher model. To investigate the impact of student model size, we use the Llama3.1-8B, Llama3.2-1B and Llama3.2-3B models as our student models.

We first fine-tune the teacher model on the Tulu 3 dataset using supervised fine-tuning (SFT) to obtain  $T_s$ . The student models are then trained either via SFT directly on the same dataset, or via knowledge distillation from  $T_s$ . The training details are listed in Appendix A.2.

# 3.3 KNOWLEDGE DISTILLATION METHOD

According to Ramesh et al. (2025), various knowledge distillation methods do not show significant differences in performance for LLMs. Therefore, we adopt the representative method GKD (Agarwal et al., 2024). GKD is a flexible framework for distilling auto-regressive language models, addressing the train-inference mismatch by incorporating *on-policy* student-generated sequences during training. Unlike standard distillation methods that rely solely on fixed datasets (e.g., ground-truth or teacher-decoded sequences), GKD enables distillation on a mixture of supervised and student-generated data, guided by token-level feedback from the teacher model.

Let x denote an input, y a target sequence, and  $\lambda \in [0,1]$  the proportion of on-policy (student-generated) data. The GKD objective is:

$$\mathcal{L}_{GKD} = (1 - \lambda) \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ KL(p_T \parallel p_S; y, x) \right] + \lambda \mathbb{E}_{x \sim \mathcal{X}} \left[ \mathbb{E}_{\hat{y} \sim p_S(\cdot \mid x)} \left[ KL(p_T \parallel p_S; \hat{y}, x) \right] \right], \quad (1)$$

where  $KL(p_T \parallel p_S; y, x)$  is the average token-level divergence:

$$KL(p_T \parallel p_S; y, x) = \frac{1}{|y|} \sum_{n=1}^{|y|} KL_{\text{token}} \left( p_T(\cdot \mid y_{< n}, x) \parallel p_S(\cdot \mid y_{< n}, x) \right), \tag{2}$$

and  $KL_{\text{token}}$  can be instantiated as forward KL, reverse KL, or generalized Jensen-Shannon divergence.

In particular, the generalized Jensen-Shannon divergence between two distributions P and Q is defined as:

$$JSD_{\beta}(P \parallel Q) = \beta \cdot KL(P \parallel M) + (1 - \beta) \cdot KL(Q \parallel M), \text{ where } M = \beta P + (1 - \beta)Q. \quad (3)$$

By adjusting the hyperparameter  $\beta \in (0,1)$ , GKD smoothly interpolates between different divergence behaviors. When  $\beta$  approaches 0,  $JSD_{\beta}$  behaves similarly to forward KL, which is modecovering: the student must assign probability mass wherever the teacher has support, leading to broader but potentially less precise coverage. Conversely, when  $\beta$  approaches 1,  $JSD_{\beta}$  resembles reverse KL, which is mode-seeking: the student focuses on the teacher's high-probability regions, yielding sharper but less diverse generations. In our experiments, we adopt  $\beta=0.5$  as the default, which balances the two effects and has been shown to perform well in practice (Agarwal et al., 2024).

### 3.4 KD Training Paradigms

While the distillation objective specifies *how* knowledge is transferred from the teacher to the student, the initialization of the student model determines *what* prior capabilities it possesses before distillation begins. This initialization choice can substantially influence learning dynamics and final performance, especially under different data regimes. To investigate this factor, we compare two KD training paradigms that differ in whether the student starts from a raw pre-trained checkpoint or from an SFT-adapted model:

- Base model as student (Base-S): The student is initialized with the pre-trained weights of the base model (e.g., Llama3.1-8B) and trained via KD from the teacher model  $T_s$ .
- SFTed model as student (SFT-S): The student is first fine-tuned on the training dataset via SFT to obtain  $S_s$ , and then further trained via KD from the teacher model  $T_s$ .

We evaluate both paradigms on Llama3.1-8B as the student model, varying the training set size from 10k to the full 939k samples. Results in Figure 1 show that the SFT-initialized student consistently outperforms the base-initialized student across most data sizes, though the gap narrows as more data is used. Even with the full dataset, the SFT-initialized student matches or slightly exceeds the base variant, indicating that SFT provides a stronger starting point for KD. This suggests that prior adaptation to instruction-following enables the student to learn more effectively from the teacher.

Based on these results, we adopt the SFT-initialized student as the default configuration for KD in subsequent experiments, as it offers a stronger prior for learning from the teacher.

# 3.5 RESULTS AND ANALYSIS

We report the average performance of student models trained via SFT and KD on five evaluation benchmarks across varying data sizes in Figure 2. In low-data regimes (fewer than 80k samples), KD

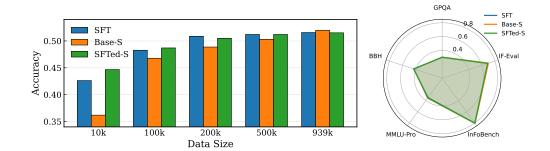


Figure 1: **Left:** Performance of the student model (Llama3.1-8B) trained via SFT and KD under different initialization paradigms across varying training set sizes. **Right:** Task-level performance on the full training set. Across most data scales, the SFT-initialized student outperforms the base-initialized student, with the gap narrowing as data increases. When trained on the full dataset, their performances converge to nearly the same level, indicating that sufficient supervision largely closes the initialization gap. This suggests that prior SFT adaptation offers a stronger starting point for KD.

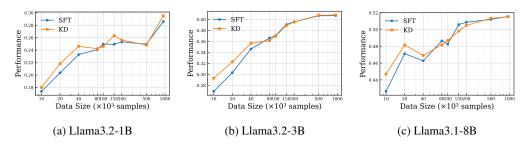


Figure 2: Performance of three student models (Llama3.2-1B, Llama3.2-3B, Llama3.1-8B) trained with SFT and KD across different training set sizes (logarithmic scale on the x-axis). Results are averaged over five benchmarks: BBH, GPQA, IFEval, InFoBench, and MMLU-Pro. KD provides clear gains in low-data regimes, while the advantage diminishes as more training data is used. Complete numerical results for all models and data sizes are provided in Appendix A.3.

consistently outperforms SFT, with the largest gain of up to 5% absolute observed at 10k samples. This suggests that KD can effectively transfer knowledge from the teacher, enabling more efficient learning when training data is scarce. As the training set grows, the performance gap narrows, and in some cases SFT even surpasses KD (e.g., Llama3.1-8B at 150k samples), indicating that with sufficient data, the student can acquire most of the teacher's knowledge through direct supervision, reducing the benefits of KD.

To examine this further, we replace GKD with SeqKD, a more traditional approach that omits onpolicy samples, and train Llama3.1-8B on the full Tulu 3 dataset. For each instance, we sample five outputs from the teacher and train the student to minimize the cross-entropy loss against these sequences. As shown in Table 1, SeqKD achieves similar performance to GKD, confirming that the KD method choice has limited impact when ample training data is available, consistent with prior findings (Ramesh et al., 2025). This supports our hypothesis that, in large-data regimes, little additional information remains to be distilled.

We further hypothesize that this limitation arises because the teacher and student are trained on the same dataset, leaving the student with few opportunities to acquire novel knowledge. In such cases, the supervision provided by KD largely duplicates the ground-truth labels, limiting the marginal utility of distillation. Consistent with this, we observed that when distilling from a same-dataset teacher, the KD loss of the student began at a low value and fluctuated without further reduction, suggesting that the student had little additional signal to absorb. To test this hypothesis, we replace the teacher with a stronger instruction-tuned model, Llama3.3-70B-Instruct (denoted as GKD-IT), which has broader instruction-following capabilities and exposure to more diverse tasks. As shown in Figure 3, GKD-IT substantially outperforms the original GKD teacher, achieving an average performance improvement of around 4% across the five benchmarks. This result supports our intuition:

| - | 1 | U |  |
|---|---|---|--|
| 2 | 7 | 1 |  |
| 2 | 7 | 2 |  |
| 2 | 7 | 3 |  |
| 2 | 7 | 4 |  |

| Model | BBH   | GPQA  | IF-Eval | InfoBench | MMLU-pro | Average |
|-------|-------|-------|---------|-----------|----------|---------|
| SFT   | 43.42 | 30.04 | 69.50   | 80.25     | 34.65    | 51.57   |
| GKD   | 42.93 | 30.22 | 68.95   | 80.33     | 35.16    | 51.52   |
| SeqKD | 42.22 | 30.59 | 63.59   | 81.63     | 34.64    | 50.53   |

2 2 275

277

278 279

281

283

284

285 286

287

288

289 290

291

292

293

295

296

297

298 299

300

301

302

303

304

305 306

307

308

309

310

311

312

313

314 315

316

317

318 319

320

321

322

323

Table 1: Performance of different methods trained on full Tulu3 training set across five evaluation benchmarks.

Average

0.55

0.54

0.5

0.5

0.51

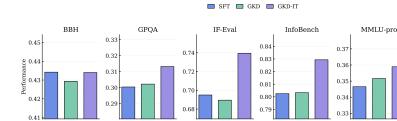


Figure 3: Performance of Llama3.1-8B with GKD and GKD-IT on the full Tulu 3 dataset. Using a stronger instruction-tuned teacher (GKD-IT) yields notable improvements, indicating that distillation from a more capable teacher can still provide substantial benefits.

distillation remains beneficial when the teacher brings in knowledge beyond the training data, enabling the student to learn patterns and reasoning strategies that it would not acquire through direct supervision alone.

General Takeaway: Knowledge distillation provides the greatest benefits in low-data regimes, making it particularly valuable for scenarios where only small or domain-specific datasets are available. Moreover, when distilling from a stronger instruction-tuned teacher, substantial gains can still be observed even in large-data settings, indicating that KD remains effective whenever the teacher contributes knowledge beyond the training set.

# TASK-SPECIFIC KNOWLEDGE DISTILLATION

Motivated by the findings in Section 3.5, we examine the use of KD in domain-specific, low-resource settings. Such scenarios are common in real-world applications, where high-quality labeled data for a specialized domain is often scarce. We evaluate the performance of the student model on several domain-specific tasks and compare multiple KD strategies to assess their effectiveness under these constraints.

### 4.1 EXPERIMENTAL SETUP

We use the same student models as in Section 3.5, but now focus on domain-specific tasks. Specifically, we consider the following tasks:

Low-resource Translation We adopt the Assamese-English subset of the Flores-200 dataset (Costa-jussà et al., 2022) in the low-resource setting, using the processed data splits provided by (Xu et al., 2025). Specifically, 997 instances from the development set are used for training, while the original test set (1012 instances) is split into 500-instance development and 512-instance test sets. Translation quality is evaluated with the COMET metric (Rei et al., 2022).

**Dialogue Summarization** We use the DialogSum dataset (Chen et al., 2021) following the preprocessed splits in (Xu et al., 2025). The training set contains 1k instances, and evaluation is conducted on the official development set (500 instances) and a 1500-instance test set. Summarization quality is measured using ROUGE-L (Lin, 2004).

**ARC-Challenge** We use the ARC-Challenge dataset (Clark et al., 2018), which consists of multiple-choice grade-school science questions that are deliberately constructed to be difficult for surface-level methods such as retrieval or word co-occurrence. The Challenge Set contains questions that require deeper knowledge and reasoning. Following the processed splits provided by (Ramesh et al., 2025), we use 1.1k instances for training, 500 for development, and 672 for testing. Performance is evaluated using accuracy.

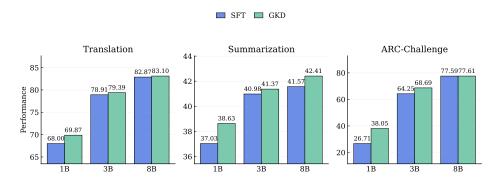


Figure 4: Performance of student models on Translation, Summarization, and ARC-Challenge. GKD improves over SFT across all tasks, but the gains diminish as the student model size increases.

Before doing task-specific training, we first fine-tune the student models on the Tulu 3 dataset via SFT to obtain student models. We then apply knowledge distillation from the teacher model  $T_s$  to the student models, using the GKD. The training details are listed in Appendix A.2.

# 4.2 KD RESULTS AND ANALYSIS IN DOMAIN-SPECIFIC TASKS

We report the results of student models trained with SFT and GKD on the three domain-specific tasks in Figure 4. Overall, GKD consistently improves performance compared to SFT, confirming the effectiveness of distillation in these settings. This observation aligns with our earlier findings in Section 3.5, where KD provided the largest benefits in low-data regimes: when task-specific training data is scarce, the student can better leverage the additional knowledge provided by the teacher.

We also observe that the magnitude of improvement varies with student model size. The performance gain from GKD is most pronounced for the 1B student, moderate for the 3B student, and becomes marginal for the 8B student. This trend suggests that smaller models depend more on distillation to acquire knowledge that cannot be fully captured from limited supervision, whereas larger models are capable of learning much of the teacher's knowledge directly from data, leaving less room for additional gains.

In practice, these findings imply that knowledge distillation is particularly valuable for building compact student models that need to operate in domain-specific, data-scarce environments. Such scenarios are common in real-world applications (e.g., specialized translation systems or domain-specific assistants), where the ability to improve small models with limited data is often more crucial than optimizing already strong larger models.

# 4.3 STRENGTHENING KNOWLEDGE DISTILLATION WITH SYNTHETIC DATA

The preceding results suggest that knowledge distillation is most beneficial in low-data regimes, yet its advantages diminish as the amount of available training data increases. One key limitation lies in the heavy reliance on human-annotated data: when the student and teacher are exposed to the same dataset, the student can already recover much of the teacher's knowledge through direct supervision, leaving limited room for further gains. In practice, labeled data in specialized domains is scarce, and scaling up high-quality annotation is often impractical. To address this challenge, we propose to augment KD with synthetic data.

**Synthetic Data Generation.** Concretely, we employ the Llama 3.1-70B model fine-tuned on Tulu 3 as the generator  $T_{\rm gen}$ . Given a small set of demonstrations  $\{(x_j,y_j)\}_{j=1}^k$  sampled from the training data  $\mathcal{D}$  and an instruction prompt  $I^1$ , we construct the in-context prompt

$$Prompt = I \oplus [In: x_1 \text{ Out: } y_1 \text{ ... In: } x_k \text{ Out: } y_k \text{ In:}], \tag{4}$$

where  $\oplus$  denotes concatenation. Conditioned on this prompt, the generator samples an unlabeled input

$$x^s \sim T_{\rm gen}(\cdot \mid {\tt Prompt}).$$
 (5)

<sup>&</sup>lt;sup>1</sup>Instruction for unlabeled data generation is provided in Appendix A.4

| Model             | Translation |       |       | Summarization |       |       | ARC-Challenge |         |       |       |       |         |
|-------------------|-------------|-------|-------|---------------|-------|-------|---------------|---------|-------|-------|-------|---------|
| 1710461           | SFT         | GKD   | Mix   | 2-Stage       | SFT   | GKD   | Mix           | 2-Stage | SFT   | GKD   | Mix   | 2-Stage |
| Teacher (70B SFT) | 86.31       | -     | -     | -             | 44.99 | -     | -             | -       | 91.38 | -     | -     | -       |
| 1B                | 68.00       | 69.87 | 78.13 | 78.65         | 37.03 | 38.63 | 40.48         | 42.84   | 26.71 | 38.05 | 70.22 | 70.31   |
| 3B                | 78.91       | 79.39 | 82.41 | 82.67         | 40.98 | 41.37 | 42.17         | 44.02   | 64.25 | 68.69 | 80.22 | 80.46   |
| 8B                | 82.87       | 83.10 | 83.15 | 83.63         | 41.57 | 42.41 | 43.40         | 44.29   | 77.99 | 77.05 | 85.35 | 85.67   |

Table 2: Performance of different models on Translation, Summarization, and ARC-Challenge. We compare SFT, KD, Mix (GKD with synthetic data mixing), and 2-Stage (two-stage GKD). Best results for each model are highlighted in **bold**.

Repeating this process yields a synthetic unlabeled set  $\mathcal{X}^s = \{x_1^s, \dots, x_N^s\}$ . These inputs are subsequently annotated by the teacher model T, which is the Llama 3.1-70B fine-tuned on the corresponding training subset  $\mathcal{D}$ , to obtain pseudo-labeled pairs

$$\mathcal{D}^{s} = \{ (x_{i}^{s}, y_{i}^{s}) \}_{i=1}^{N}, \quad y_{i}^{s} \sim T(\cdot \mid x_{i}^{s}).$$
 (6)

**Synthetic Data Integration.** A straightforward approach is to simply combine the synthetic dataset  $\mathcal{D}^s$  with the original human-annotated dataset  $\mathcal{D}$  and train the student under the KD objective:

$$\mathcal{L}_{\text{mix}} = \mathbb{E}_{(x,y) \sim \mathcal{D} \cup \mathcal{D}^s} \left[ \text{KL}(p_T(\cdot \mid x) \parallel p_S(\cdot \mid x)) \right]. \tag{7}$$

As we will show in Section 4.4, this simple mixing strategy indeed provides improvements over SFT, indicating that synthetic data can be beneficial despite its noise However, its effect remains limited, as the noisy synthetic samples  $\mathcal{D}^s$  dilute the high-quality supervision from  $\mathcal{D}$ , preventing the student from fully exploiting the additional data. This motivates the development of a more structured integration scheme.

To this end, we adopt a two-stage training paradigm. In the first stage, the student S is exposed only to the synthetic dataset, performing an initial training step:

$$S^{(0)} = \operatorname{Train}(S; \mathcal{D}^s), \tag{8}$$

where the objective is the standard cross-entropy loss with pseudo-labels  $y_i^s$ . Although  $\mathcal{D}^s$  is noisy, this stage helps the student adapt to a broader distribution of instruction formats and response structures, providing a better starting point.

In the second stage, we initialize the student with  $S^{(0)}$  and then optimize the GKD objective on the human-annotated dataset  $\mathcal{D}$  (see Sec. 3.3 for details):

$$S^* = \arg\min_{S} \mathcal{L}_{GKD}(S; T, \mathcal{D})$$
 with initialization  $S \leftarrow S^{(0)}$ . (9)

This staged design leverages synthetic data as a "warm-up" that broadens the student's exposure, while preserving the high-quality guidance of the teacher on  $\mathcal{D}$  in the final distillation step. Empirically, we find that this method significantly improves performance compared to both direct mixing and vanilla KD.

# 4.4 RESULTS OF SYNTHETIC DATA AUGMENTATION

We empirically set the number of synthetic samples N for each task to 100k. Besides SFT and GKD, we report the results of the following two strategies:

- **Mixing**: Directly combining the synthetic dataset  $\mathcal{D}^s$  with the original dataset  $\mathcal{D}$  and training the student via GKD on the mixed dataset.
- Two-stage: The proposed two-stage training paradigm, where the student is first trained on the synthetic dataset  $\mathcal{D}^s$ , and then fine-tuned via GKD on the original dataset  $\mathcal{D}$ .

We report the results in Table 2, from which we draw the following findings:

**Mixing synthetic data yields clear improvements.** Across all tasks and model sizes, the mixing strategy consistently outperforms SFT, showing that synthetic data can effectively enhance KD despite its noisiness. For example, the 3B model on *Translation* improves from 78.91 (SFT) to 82.41 with mixing, and the 8B model on *Summarization* improves from 41.57 (SFT) to 43.40 with mixing. These results indicate that even a direct combination of human and synthetic data provides noticeable gains.

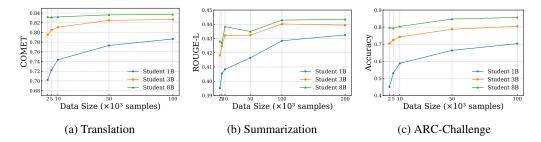


Figure 5: Impact of synthetic data size on knowledge distillation. Performance improves as more synthetic data is added, with smaller models benefiting more, but the gains quickly saturate as data size grows.

**Two-stage training maximizes the benefit of synthetic data.** The two-stage paradigm consistently achieves the best results across all settings. For example, the 1B model on *ARC-Challenge* improves to **44.97**, compared to 26.71 with SFT and 28.66 with mixing. This demonstrates that, while synthetic data is inherently noisy, it becomes substantially more beneficial when used as a warm-up stage before distillation on human-annotated data.

**General takeaway.** Overall, our results show that synthetic data is indeed useful for improving KD in low-resource, domain-specific tasks. Even the simple mixing strategy brings consistent gains over SFT. At the same time, the two-stage paradigm further amplifies these benefits by structuring how synthetic data is leveraged. This demonstrates that synthetic data, whether used directly or in a staged manner, can substantially enhance student models in low-resource settings, with two-stage training providing the most effective integration.

# 4.5 IMPACT OF SYNTHETIC DATA SIZE

To better understand the role of synthetic data, we investigate how the size of the synthetic dataset influences distillation performance. Specifically, we vary the number of synthetic samples N from 5k to 100k for all 1B, 3B, and 8B student models, and report the results in Figure 5. We find that adding synthetic data consistently improves performance across different tasks, with the 1B student benefiting the most. However, the improvement is most significant when moving from very limited to moderate amounts of synthetic data, after which the marginal gains diminish. This indicates that while synthetic data is an effective complement to KD in low-resource settings, its utility does not scale linearly with quantity. Instead, the main advantage comes from a relatively small but diverse synthetic set that broadens the student's exposure beyond what human annotations alone can provide.

# 5 Conclusion

In this work, we conducted a systematic study of KD in the post-training stage of LLMs. Through extensive experiments on the large-scale Tulu 3 dataset, we found that KD consistently outperforms SFT in low-data regimes, but its benefits diminish as training data grows. Nevertheless, distilling from a stronger instruction-tuned teacher restores substantial gains even in high-data settings, high-lighting that KD remains effective when the teacher possesses knowledge beyond the training set.

We further explored domain-specific, low-resource scenarios and demonstrated that KD is particularly valuable for smaller student models. To address the limitations of scarce human annotations, we introduced a two-stage KD paradigm that first leverages synthetic teacher-labeled data before refining on human annotations. This method consistently improved student performance, surpassing both direct mixing and standard KD, thereby offering a practical recipe for building compact and capable models under resource constraints.

Overall, our findings provide a clearer understanding of when and why KD is most effective, establish its limitations when teacher and student share the same supervision, and present a simple yet effective two-stage framework for integrating synthetic data. We hope these insights will guide future research on efficient LLM post-training and inform the development of deployable models in real-world low-resource applications.

# REFERENCES

Rishabh Agarwal, Nino Vieillard, Yongchao Zhou, Piotr Stanczyk, Sabela Ramos Garea, Matthieu Geist, and Olivier Bachem. On-policy distillation of language models: Learning from self-generated mistakes. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024.* OpenReview.net, 2024. URL https://openreview.net/forum?id=3zKtaqxLhW.

Yulong Chen, Yang Liu, Liang Chen, and Yue Zhang. Dialogsum: A real-life scenario dialogue summarization dataset. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), Findings of the Association for Computational Linguistics: ACL/IJCNLP 2021, Online Event, August 1-6, 2021, volume ACL/IJCNLP 2021 of Findings of ACL, pp. 5062–5074. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.FINDINGS-ACL.449. URL https://doi.org/10.18653/v1/2021.findings-acl.449.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the AI2 reasoning challenge. *CoRR*, abs/1803.05457, 2018. URL http://arxiv.org/abs/1803.05457.

Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Y. Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. No language left behind: Scaling human-centered machine translation. *CoRR*, abs/2207.04672, 2022. doi: 10.48550/ARXIV.2207.04672. URL https://doi.org/10.48550/arXiv.2207.04672.

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. CoRR, abs/2501.12948, 2025. doi: 10. 48550/ARXIV.2501.12948. URL https://doi.org/10.48550/arXiv.2501.12948.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak,

Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The llama 3 herd of models. *CoRR*, abs/2407.21783, 2024. doi: 10.48550/ARXIV.2407.21783. URL https://doi.org/10.48550/arXiv.2407.21783.

- Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models. In *The Twelfth International Conference on Learning Representations, ICLR* 2024, *Vienna, Austria, May* 7-11, 2024. OpenReview.net, 2024. URL https://openreview.net/forum?id=5h0qf7IBZZ.
- Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015. URL http://arxiv.org/abs/1503.02531.
- Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation. In Jian Su, Xavier Carreras, and Kevin Duh (eds.), *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pp. 1317–1327. The Association for Computational Linguistics, 2016. doi: 10.18653/V1/D16-1139. URL https://doi.org/10.18653/v1/d16-1139.
- Andreas Köpf, Yannic Kilcher, Dimitri von Rütte, Sotiris Anagnostidis, Zhi Rui Tam, Keith Stevens, Abdullah Barhoum, Duc Nguyen, Oliver Stanley, Richárd Nagyfi, Shahul ES, Sameer Suri, David Glushkov, Arnav Dantuluri, Andrew Maguire, Christoph Schuhmann, Huu Nguyen, and Alexander Mattick. Openassistant conversations democratizing large language model alignment. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023. URL http://papers.nips.cc/paper\_files/paper/2023/hash/949f0f8f32267d297c2d4e3ee10a2e7e-Abstract-Datasets\_and\_Benchmarks.html.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tülu 3: Pushing frontiers in open language model post-training. *CoRR*, abs/2411.15124, 2024. doi: 10. 48550/ARXIV.2411.15124. URL https://doi.org/10.48550/arXiv.2411.15124.
- Ashley Lewis, Michael White, Jing Liu, Toshiaki Koike-Akino, Kieran Parsons, and Ye Wang. Winning big with small models: Knowledge distillation vs. self-training for reducing hallucination in QA agents. *CoRR*, abs/2502.19545, 2025. doi: 10.48550/ARXIV.2502.19545. URL https://doi.org/10.48550/arXiv.2502.19545.
- Yixing Li, Yuxian Gu, Li Dong, Dequan Wang, Yu Cheng, and Furu Wei. Direct preference knowledge distillation for large language models. *CoRR*, abs/2406.19774, 2024. doi: 10.48550/ARXIV. 2406.19774. URL https://doi.org/10.48550/arXiv.2406.19774.
- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pp. 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL https://aclanthology.org/W04-1013/.
- OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774. URL https://doi.org/10.48550/arXiv.2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022. URL http://papers.nips.cc/paper\_files/paper/2022/hash/blefde53be364a73914f58805a001731-Abstract-Conference.html.

Yiwei Qin, Kaiqiang Song, Yebowen Hu, Wenlin Yao, Sangwoo Cho, Xiaoyang Wang, Xuansheng Wu, Fei Liu, Pengfei Liu, and Dong Yu. Infobench: Evaluating instruction following ability in large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pp. 13025–13048. Association for Computational Linguistics, 2024. doi: 10.18653/V1/2024.FINDINGS-ACL.772. URL https://doi.org/10.18653/v1/2024.findings-acl.772.

Suhas Kamasetty Ramesh, Ayan Sengupta, and Tanmoy Chakraborty. On the generalization vs fidelity paradox in knowledge distillation. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Findings of the Association for Computational Linguistics, ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pp. 17930–17951. Association for Computational Linguistics, 2025. URL https://aclanthology.org/2025.findings-acl.923/.

Ricardo Rei, José G. C. de Souza, Duarte M. Alves, Chrysoula Zerva, Ana C. Farinha, Taisiya Glushkova, Alon Lavie, Luísa Coheur, and André F. T. Martins. COMET-22: unbabel-ist 2022 submission for the metrics shared task. In Philipp Koehn, Loïc Barrault, Ondrej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno-Yepes, Tom Kocmi, André F. T. Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo Negri, Aurélie Névéol, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri (eds.), *Proceedings of the Seventh Conference on Machine Translation, WMT 2022, Abu Dhabi, United Arab Emirates (Hybrid), December 7-8, 2022*, pp. 578–585. Association for Computational Linguistics, 2022. URL https://aclanthology.org/2022.wmt-1.52.

David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. GPQA: A graduate-level google-proof q&a benchmark. *CoRR*, abs/2311.12022, 2023. doi: 10.48550/ARXIV.2311.12022. URL https://doi.org/10.48550/arXiv.2311.12022.

Anup Shirgaonkar, Nikhil Pandey, Nazmiye Ceren Abay, Tolga Aktas, and Vijay Aski. Knowledge distillation using frontier open-source llms: Generalizability and the role of synthetic data. *CoRR*, abs/2410.18588, 2024. doi: 10.48550/ARXIV.2410.18588. URL https://doi.org/10.48550/arXiv.2410.18588.

Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartlomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, Cèsar Ferri Ramírez, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodolà, Emma Lam, Eric Chu, Eric Tang, Erkut

650

651

652

653

654

655

656

657

658

659

660

661

662

666

667

668

669

670

671

672

673

674

675

676

677

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

696

699

700

Erdem, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan J. Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, François Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kocon, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse H. Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, José Hernández-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory W. Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, Luke Metz, Lütfi Kerem Senel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, María José Ramírez-Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Baitemirova, Melody Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael I. Ivanitskiy, Michael Starritt, Michael Strube, Michal Swedrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T., Nanyun Peng, Nathan A. Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Milkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima (Shammie) Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantadosi, Stuart M. Shieber, Summer Misherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay V. Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. Trans. Mach. Learn. Res., 2023, 2023. URL https://openreview.net/forum?id=uyTL5Bvosj.

Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/stanford\_alpaca, 2023.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024, 2024. URL http://papers.nips.cc/paper\_files/paper/2024/hash/ad236edc564f3e3156e1b2feafb99a24-Abstract-Datasets\_and\_Benchmarks\_Track.html.

Wenda Xu, Rujun Han, Zifeng Wang, Long T. Le, Dhruv Madeka, Lei Li, William Yang Wang, Rishabh Agarwal, Chen-Yu Lee, and Tomas Pfister. Speculative knowledge distillation: Bridging the teacher-student gap through interleaved sampling. In *The Thirteenth International Conference on Learning Representations, ICLR 2025, Singapore, April 24-28, 2025.* OpenReview.net, 2025. URL https://openreview.net/forum?id=EgJhwYR2tB.

An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jian Yang, Jiaxi Yang, Jingren Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *CoRR*, abs/2505.09388, 2025. doi: 10.48550/ARXIV.2505.09388. URL https://doi.org/10.48550/arxiv.2505.09388.

Jeffrey Zhou, Tianjian Lu, Swaroop Mishra, Siddhartha Brahma, Sujoy Basu, Yi Luan, Denny Zhou, and Le Hou. Instruction-following evaluation for large language models, 2023. URL https://arxiv.org/abs/2311.07911.

# A APPENDIX

# A.1 DETAILS OF EVALUATION TASKS

- **BIG-bench(BBH)** (Srivastava et al., 2023): A large-scale benchmark comprising 204 diverse tasks spanning linguistics, reasoning, math, science, social bias, and more. Designed to probe capabilities believed to be beyond current language models, it evaluates both quantitative performance and qualitative behaviors across a wide range of domains.
- **GPQA** (Rein et al., 2023): A challenging multiple-choice benchmark of 448 questions in biology, physics, and chemistry, authored by domain experts. The questions are designed to be "Google-proof" and extremely difficult, with expert-level accuracy around 65% and GPT-4 baselines achieving only 39%, making it suitable for evaluating advanced reasoning in specialized scientific domains.
- **IFEval** (Zhou et al., 2023): An instruction-following benchmark of around 500 prompts containing verifiable constraints, such as word-count limits or required keywords. It provides a reproducible and objective way to assess LLMs' ability to follow natural language instructions without relying on costly human evaluation.
- InFoBench (Qin et al., 2024): A benchmark of 500 diverse instructions decomposed into 2,250 fine-grained requirements, designed to evaluate LLMs' instruction-following ability under the Decomposed Requirements Following Ratio (DRFR) metric. It enables detailed assessment of compliance with multiple constraint categories and supports evaluation using human or LLM-based annotators.

• MMLU-Pro (Wang et al., 2024): A benchmark built upon the Massive Multitask Language Understanding (MMLU) benchmark, which tests language understanding and reasoning across a wide range of subjects. While MMLU mainly contains knowledge-driven multiple-choice questions, MMLU-Pro increases difficulty by replacing trivial or noisy items with reasoning-focused questions and expanding the choice set from four to ten options. This design better discriminates between advanced LLMs and reduces score sensitivity to prompt variations.

# A.2 IMPLEMENTATION DETAILS

We adopt the GKD trainer from the TRL library for all KD experiments. We use the AdamW optimizer with a learning rate of  $5 \times 10^{-6}$  and a batch size of 128. The models are trained for 2 epochs. We use a linear learning rate scheduler with a warm-up phase of 3% of the total training steps. The maximum sequence length is set to 4096 tokens. All experiments are conducted on NVIDIA H100 GPUs. To aviod the instability of gradient accumulation  $^2$ , we use sum instead of average to compute the loss over multiple batches.

We set the hyperparameter  $\lambda$  in Eqn. 2 to 0.5, balancing the contributions of the training data and the on-policy samples. The hyperparameter  $\beta$  in the generalized Jensen-Shannon divergence is set to 0.5, which balances the mode-covering and mode-seeking behaviors.

For SFT training, we use the same hyperparameters as in KD training.

For generating synthetic data, we randomly sample 10 examples from the training set as in-context demonstrations. We use nucleus sampling with p=1.0 and a temperature of 0.8 to generate synthetic inputs. The maximum generation length is set to 4096 tokens. We set the temperature to 0.6 when the teacher annotates the synthetic inputs.

# A.3 DETAILED EXPERIMENTAL RESULTS

Table 3 reports the detailed performance of all student models across different training set sizes, complementing the averaged trends shown in Figure 2. Each entry corresponds to accuracy (or task-specific metric) on the five evaluation benchmarks: BBH, GPQA, IFEval, InFoBench, and MMLU-Pro.

The results confirm that knowledge distillation (GKD) yields clear gains over supervised fine-tuning (SFT) in low-data regimes, particularly for smaller students such as Llama3.2-1B and Llama3.2-3B. However, as the training set size increases, the advantage of GKD diminishes, and in some cases SFT matches or slightly surpasses GKD (e.g., Llama3.1-8B at 939k samples). These detailed results provide quantitative evidence for the scaling limitations of KD and support our conclusion that its primary benefits lie in low-resource settings.

### A.4 SYNTHETIC DATA GENERATION PROMPT

The instruction prompt used for generating synthetic data is as follows:

```
You are a data generation assistant. You will be given 10 demonstrations
Task: Based on these examples, produce exactly one new input message that matches the same task, language, and style.
Strict requirements: Output only the message content itself. Do NOT include any explanations, quotes, labels, or the answer.

Here are 10 demonstrations:
<10 examples from the training set>
Now, generate exactly one brand-new input message that follows the same task and formatting.
```

<sup>&</sup>lt;sup>2</sup>https://unsloth.ai/blog/gradient

Output only the input message content itself. Do not output any answer or extra text.

# B LLM USAGE

In preparing this paper, we used GPT5 and Gemini solely as a writing assistant to polish the writing. Specifically, LLMs were employed to improve the fluency, clarity, and grammar of the text, while the core research ideas, experimental design, implementation, analysis, and conclusions were entirely developed by the authors.

| Model       | Data Size | Method     | BBH            | GPQA           | IF-Eval        | InfoBench      | MMLU-Pro       |          |
|-------------|-----------|------------|----------------|----------------|----------------|----------------|----------------|----------|
|             | 10k       | SFT<br>GKD | 2.95<br>3.84   | 25.09<br>24.18 | 12.38<br>15.71 | 34.25<br>34.15 | 12.17<br>12.13 | 17<br>18 |
|             | 20k       | SFT<br>GKD | 11.06<br>14.88 | 24.73<br>25.27 | 15.90<br>17.93 | 37.84<br>38.68 | 12.25<br>12.37 | 20<br>21 |
|             | 40k       | SFT<br>GKD | 14.59<br>15.02 | 23.81<br>25.09 | 21.44<br>24.77 | 44.06<br>45.70 | 12.53<br>12.48 | 23<br>24 |
| llama3.2-1b | 80k       | SFT<br>GKD | 13.64<br>10.81 | 25.46<br>25.09 | 23.84<br>24.58 | 45.10<br>48.00 | 12.36<br>12.67 | 24<br>24 |
|             | 100k      | SFT<br>GKD | 13.75<br>11.96 | 26.19<br>23.44 | 25.69<br>27.54 | 46.65<br>47.41 | 12.58<br>12.61 | 24<br>24 |
|             | 150k      | SFT<br>GKD | 13.79<br>13.84 | 26.01<br>26.92 | 25.69<br>27.73 | 46.83<br>50.64 | 12.36<br>12.52 | 24<br>26 |
|             | 200k      | SFT<br>GKD | 13.21<br>11.37 | 26.01<br>26.37 | 26.99<br>27.36 | 48.09<br>50.40 | 12.33<br>12.61 | 25<br>25 |
|             | 500k      | SFT<br>GKD | 9.32<br>4.72   | 23.81<br>24.54 | 30.13<br>30.31 | 49.05<br>51.39 | 12.48<br>13.10 | 24<br>24 |
|             | 939k      | SFT<br>GKD | 13.85<br>14.02 | 27.84<br>25.64 | 36.60<br>39.93 | 52.03<br>55.11 | 12.66<br>12.92 | 28<br>29 |
|             | 10k       | SFT<br>GKD | 0.35<br>3.96   | 29.30<br>28.39 | 27.73<br>29.21 | 51.73<br>59.05 | 25.48<br>25.96 | 26       |
|             | 20k       | SFT<br>GKD | 12.70<br>16.43 | 26.74<br>27.11 | 33.09<br>35.12 | 53.44<br>56.55 | 25.61<br>26.28 | 30       |
|             | 40k       | SFT<br>GKD | 23.22<br>22.68 | 26.92<br>27.84 | 37.71<br>39.37 | 60.36<br>62.52 | 24.92<br>25.94 | 34<br>35 |
| llama3.2-3b | 80k       | SFT<br>GKD | 24.57<br>20.93 | 28.02<br>26.92 | 39.74<br>41.04 | 65.36<br>65.70 | 25.41<br>26.13 | 36<br>36 |
|             | 100k      | SFT<br>GKD | 23.67<br>21.01 | 26.74<br>26.56 | 43.44<br>44.55 | 66.92<br>67.93 | 24.46<br>24.68 | 37<br>36 |
|             | 150k      | SFT<br>GKD | 25.03<br>21.76 | 28.39<br>29.30 | 47.69<br>48.61 | 69.17<br>69.39 | 25.07<br>25.36 | 39       |
|             | 200k      | SFT<br>GKD | 27.12<br>21.55 | 28.21<br>29.30 | 48.98<br>51.76 | 68.50<br>69.23 | 24.97<br>25.79 | 39       |
|             | 500k      | SFT<br>GKD | 25.93<br>21.07 | 25.46<br>25.64 | 55.08<br>57.30 | 73.27<br>75.24 | 23.58<br>24.51 | 40       |
|             | 939k      | SFT<br>GKD | 21.30<br>18.14 | 27.11<br>27.66 | 57.30<br>58.60 | 73.79<br>74.86 | 23.96<br>24.71 | 40       |
|             | 10k       | SFT<br>GKD | 21.98<br>25.50 | 29.85<br>31.14 | 48.98<br>51.76 | 73.90<br>76.27 | 38.41<br>38.71 | 42<br>44 |
|             | 20k       | SFT<br>GKD | 38.00<br>34.74 | 31.32<br>32.05 | 53.05<br>57.49 | 75.76<br>77.47 | 37.52<br>39.09 | 47<br>48 |
| llama3.1-8b | 40k       | SFT<br>GKD | 31.72<br>28.23 | 30.95<br>32.05 | 55.08<br>58.96 | 76.01<br>77.13 | 37.61<br>38.23 | 46       |
|             | 80k       | SFT<br>GKD | 33.94<br>28.29 | 33.15<br>31.68 | 59.52<br>63.03 | 79.29<br>79.56 | 37.38<br>38.42 | 48<br>48 |
|             | 100k      | SFT<br>GKD | 30.56<br>28.54 | 31.68<br>30.77 | 64.51<br>68.76 | 78.47<br>78.86 | 36.14<br>36.56 | 48<br>48 |
|             | 150k      | SFT<br>GKD | 40.15<br>34.43 | 31.50<br>29.12 | 65.62<br>68.39 | 79.10<br>79.45 | 36.65<br>37.70 | 50<br>49 |
|             | 200k      | SFT<br>GKD | 41.47<br>39.32 | 29.85<br>29.12 | 68.58<br>68.02 | 78.58<br>78.90 | 35.90<br>37.13 | 50<br>50 |
|             | 500k      | SFT<br>GKD | 40.98<br>38.77 | 30.95<br>30.04 | 69.13<br>69.50 | 79.24<br>82.64 | 35.74<br>35.80 | 51<br>51 |
|             | 939k      | SFT<br>GKD | 43.42<br>42.93 | 30.04<br>30.22 | 69.50<br>68.95 | 80.25<br>80.33 | 34.65<br>35.16 | 51<br>51 |

Table 3: Detailed performance of student models (Llama3.2-1B, Llama3.2-3B, and Llama3.1-8B) trained with SFT and GKD across different training set sizes. Results are reported on five benchmarks (BBH, GPQA, IFEval, InFoBench, and MMLU-Pro).