

MINIPLM: Knowledge Distillation for Pre-Training Language Models

Yuxian Gu^{1,2*}, Hao Zhou², Fandong Meng², Jie Zhou², Minlie Huang¹

¹The CoAI Group, Tsinghua University

²WeChat AI, Tencent Inc., China

Abstract

Knowledge distillation (KD) is widely used to train small, high-performing student language models (LMs) using large teacher LMs. While effective in fine-tuning, KD during pre-training faces challenges in efficiency, flexibility, and effectiveness. Existing methods either incur high computational costs due to online teacher inference, require tokenization matching between teacher and student LMs, or risk losing the difficulty and diversity of the teacher-generated training data. To address these issues, we propose **MINIPLM**, a KD framework for pre-training LMs by refining the training data distribution with the teacher’s knowledge. For efficiency, MINIPLM performs offline teacher LM inference, allowing KD for multiple student LMs without adding training-time costs. For flexibility, MINIPLM operates solely on the training corpus, enabling KD across model families. For effectiveness, MINIPLM leverages the differences between large and small LMs to enhance the difficulty and diversity of the training data, helping student LMs acquire versatile and sophisticated knowledge. Extensive experiments demonstrate that MINIPLM boosts the student LMs’ performance on 9 widely used downstream tasks, improves the language modeling capabilities, and reduces pre-training computation. The benefit of MINIPLM extends to large pre-training scales, evidenced by the extrapolation of the scaling curves. Further analysis reveals that MINIPLM supports KD across model families and enhances the utilization of pre-training data. Our model, code, and data are available at <https://github.com/thu-coai/MiniPLM>.

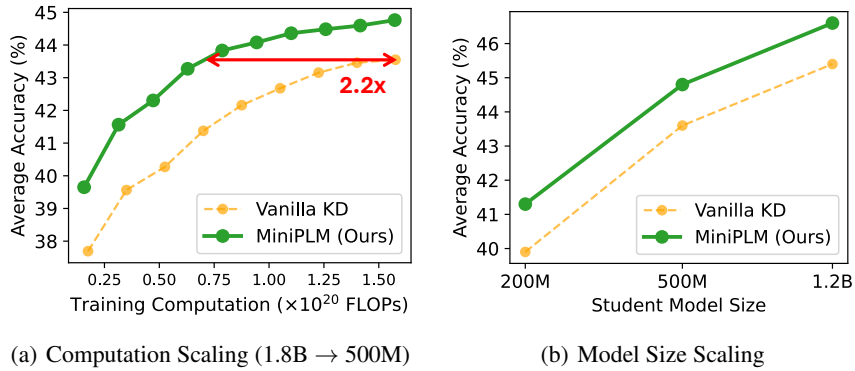


Figure 1: Computation (a) and model size (b) scaling curves of student LMs pre-trained from scratch with Vanilla KD¹ and MINIPLM. The teacher LM has 1.8B parameters. “1.8B→500M” means we use a 500M student LM. Training-time computation is kept constant for LMs of the same size in model scaling. The y-axis represents the LMs’ zero-shot performance on 9 downstream NLP tasks.

*Contribution during an internship at Tencent Inc. (guyx21@mails.tsinghua.edu.cn)

¹Vanilla KD [65, 58] minimizes the token-level forward Kullback-Leibler divergence between the output distributions of the teacher LM and student LM.

1 Introduction

Recent advances in language models (LMs; 35, 8, 60, 78) have largely been driven by scaling up model sizes, but this comes with high inference costs for deployment. At the same time, training small, deployment-friendly LMs faces training computation challenges, as small models are typically far from compute-optimal configurations according to Scaling Laws [37]. This has spurred a growing interest in exploring the limits of small LMs under the constraint of pre-training computation [51].

Knowledge Distillation (KD; 36), where a small student model learns from a large teacher model, is a promising approach for training high-performing small LMs. While KD is effective for fine-tuning LMs on specific tasks [88], its role in improving pre-training, the critical stage for small LMs to acquire versatile foundation knowledge, remains under-explored. It is non-trivial to apply KD during pre-training with the methods in fine-tuning stages, which can be categorized as **online KD** and **offline KD**. In online KD, the teacher LM has to perform inference during the student LM pre-training to provide token-level probability supervision² [65, 32], introducing additional training-time computation overhead. As shown in Figure 2, while online KD enhances performance within the same training steps, its benefits diminish if the extra computation was instead used to extend pre-training without KD. In addition, most online KD methods require the teacher and student LMs to share tokenization, limiting its flexibility across different model families. Offline KD, on the other hand, avoids extra training-time computation and allows for KD across model families, as student LMs are trained on the data offline generated by the teacher LM [41, 33]. However, ensuring sufficient difficulty and diversity in the teacher-generated data for pre-training is challenging without extensive human expertise. This often causes the student LM to overfit easy and common language patterns, hindering its generalization to versatile downstream tasks [70], as demonstrated in Figure 2.

To address these challenges, we propose **MINIPLM**, an efficient, flexible, and effective KD framework for pre-training LMs, as illustrated in Figure 3. Intuitively, MINIPLM distills the teacher LM’s knowledge into the pre-training distribution through *Difference Sampling*, which samples training instances based on the “difference” between large and small LMs. Student LMs are then pre-trained from scratch on the refined distribution. To ensure efficiency, as shown in Figure 3(a), *Difference Sampling* performs offline teacher LM inference, allowing MINIPLM to distill knowledge into multiple student LMs without incurring additional training-time costs. For flexibility, MINIPLM operates solely on the training corpus, enabling KD across model families and ensuring seamless integration with highly optimized training pipelines [15]. In terms of effectiveness, *Difference Sampling* samples training instances that the teacher LM prefers but that a small reference LM assigns low probabilities to, promoting data difficulty and diversity. As depicted in Figure 3(b), this design down-samples easy and common patterns, up-samples hard and diverse instances, and filters out noisy or harmful data points from the pre-training corpus, which encourages student LMs to acquire versatile and sophisticated knowledge, ultimately improving downstream generalization.

We apply MINIPLM to pre-train 200M, 500M, and 1.2B student LMs from scratch, using a 1.8B teacher LM. We show that MINIPLM surpasses various baselines in improving student LMs’ zero-shot performance on 9 widely used downstream tasks, enhancing language modeling capabilities, and reducing pre-training computation. By extrapolating the test loss with the Scaling Law [37], we observe that MINIPLM’s benefit remains consistent for LMs trained on $\sim 10T$ tokens. MINIPLM also facilitates KD across model families, improving Llama3.1 [19] and Mamba [27] with a teacher LM from the Qwen family [3]. Further analysis shows that MINIPLM enhances pre-training data utilization, reducing the data demand by 2.4 times, which mitigates the quick exhaustion of web-crawled corpora [80].

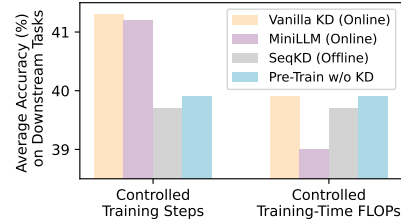


Figure 2: Results of applying KD methods in fine-tuning to pre-train a 200M student LM, using a 1.8B teacher LM. See Section 3.1 for method and evaluation details. When the training FLOPs are controlled, all KD methods perform similar or worse than Pre-Train w/o KD.

²Pre-computing this supervision is impractical, as it requires 30PB of storage for 50B tokens with a 150K vocabulary. Therefore, teacher LM’s probabilities are typically computed online during the student LM training.

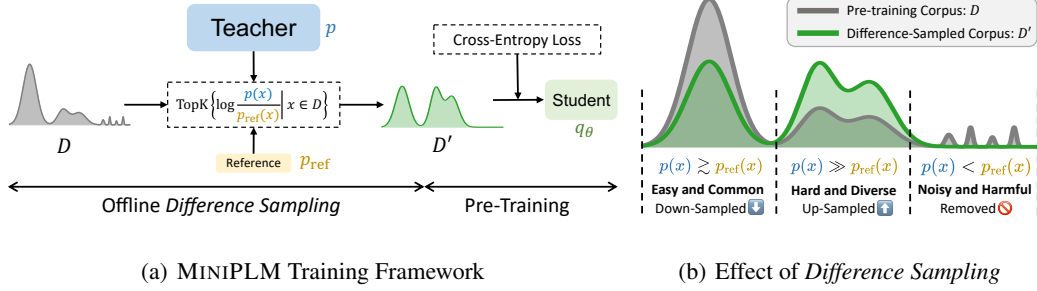


Figure 3: MINIPLM. (a): Training framework. MINIPLM distills the knowledge of the teacher LM into the student LM by adjusting the pre-training corpus of the student LM (q_θ) through **offline Difference Sampling**, based on the output probability discrepancy between the teacher LM (p) and a small reference LM (p_{ref}). (b): Illustration of the effect of *Difference Sampling*, which down-samples common easy instances, up-samples hard valuable instances, and removes noisy harmful instances.

2 MINIPLM: KD for Pre-training LMs

We consider pre-training an LM with an output distribution q_θ on a large-scale corpus \mathcal{D} consisting of N text sequences, where θ represents the model parameters. KD aids pre-training by incorporating the knowledge of a teacher LM with output distribution p , and training θ to minimize the discrepancy between p and q_θ [65]. MINIPLM formulates KD as a reward maximization problem [34], which trains the student LM to generate diverse texts receiving high preference from the teacher LM (Section 2.1). To efficiently and effectively optimize the reward, as illustrated by Figure 3, MINIPLM employs *Difference Sampling* to refine the pre-training corpus (Section 2.2). This process improves the pre-training distribution in an offline manner with the teacher LM’s knowledge and preserves the data diversity and difficulty with the help of a small reference model. The student LM is then pre-trained from scratch on the refined corpus (Section 2.3).

2.1 KD as Reward Maximization

Recent works [32, 1, 42] have shown the effectiveness of minimizing the reverse Kullback-Leibler divergence (KLD) between p and q_θ for KD of LMs in the fine-tuning stage, which avoids q_θ from over-estimating the low-probability regions of p . We reformulate this objective as a reward maximizing problem [34], which is suitable for designing efficient offline optimization methods [44]:

$$\begin{aligned} \theta &= \arg \min_{\theta} \text{KL}[q_\theta || p] = \arg \min_{\theta} \mathbb{E}_{\mathbf{x} \sim q_\theta} \log \frac{q_\theta(\mathbf{x})}{p(\mathbf{x})} \\ &= \arg \max_{\theta} \mathbb{E}_{\mathbf{x} \sim q_\theta} r(p, q_\theta, \mathbf{x}), \end{aligned} \quad (1)$$

where the reward r is defined as $r(p, q_\theta, \mathbf{x}) = \log \frac{p(\mathbf{x})}{q_\theta(\mathbf{x})}$. Intuitively, Eq. (1) trains q_θ to generate text \mathbf{x} that receives high reward values, indicating a preference from the teacher LM (i.e., high $\log p(\mathbf{x})$ values), while also ensuring high diversity (i.e., low $\mathbb{E}_{\mathbf{x} \sim q_\theta} q_\theta(\mathbf{x})$ values). To optimize Eq. (1), a simple yet effective approach is Best-of-N [71, 4, 79], where a set \mathcal{D}_{q_θ} containing M candidates is first sampled from q_θ and K instances with the highest rewards are selected from these candidates to form a new dataset \mathcal{D}'_{q_θ} :

$$\mathcal{D}'_{q_\theta} = \text{top-}K\{r(p, q_\theta, \mathbf{x}) | \mathbf{x} \in \mathcal{D}_{q_\theta}\}, \quad (2)$$

where $\mathcal{D}_{q_\theta} = \{\mathbf{x}_m | \mathbf{x}_m \sim q_\theta, 1 \leq m \leq M\}$. The student LM is then trained on \mathcal{D}'_{q_θ} to learn to generate texts with large $r(p, q_\theta, \mathbf{x})$ values. Although Best-of-N achieves performing the teacher LM’s inference prior to the student LM training similar to offline KD [41, 33], it still lacks the efficiency advantage of these methods because obtaining \mathcal{D}'_{q_θ} requires (1) **sampling data from q_θ** and (2) **computing the reward values with q_θ** , making \mathcal{D}'_{q_θ} non-transferable for pre-training other student LMs. It is also hard to ensure the diversity and difficulty of the M candidates sampled from q_θ without careful prompt engineering with human expertise. In the following, we show that these issues can be effectively and efficiently addressed by *Difference Sampling*, which is the basis of the MINIPLM training algorithm.

2.2 Difference Sampling

As shown in Figure 3(a), *Difference Sampling* refines the pre-training corpus \mathcal{D} based on the discrepancy between p and the output distribution p_{ref} from a tiny reference LM, which eliminates the dependency of Eq. (2) on q_θ , making the sampled corpus reusable in training multiple student LMs.

Top- K Sampling From \mathcal{D} , not \mathcal{D}_{q_θ} . To avoid sampling data from q_θ , we sample instances with high $r(p, q_\theta, \mathbf{x})$ values from the pre-training corpus \mathcal{D} , rather than from the \mathcal{D}_{q_θ} generated by the student LM as in Eq. (2). This way, changing the student LM does not affect the candidate set, and \mathcal{D} contains enough diverse and hard examples to be sampled for pre-training. The following proposition offers theoretical support for this approach, showing that the sampled training instances from \mathcal{D}_{q_θ} and \mathcal{D} are highly likely to be the same when the sizes of \mathcal{D}_{q_θ} and \mathcal{D} are sufficiently large:

Proposition 2.1. *Let S be the sample space of two distributions p_1 and p_2 , $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N \sim p_1$ be N i.i.d random variables, and $\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_M \sim p_2$ be M i.i.d random variables. Let $r(\cdot) : S \mapsto \mathbb{R}$ be any injective function. Assume that $\forall \mathbf{x} \in S$, $p_1(\mathbf{x}) > 0$, $p_2(\mathbf{x}) > 0$. For a fixed K satisfying $1 \leq K \leq \min\{N, M\}$, when $N \rightarrow +\infty$, $M \rightarrow +\infty$, we have*

$$P(\text{top-}K\{r(\mathbf{X}_n) | 1 \leq n \leq N\} = \text{top-}K\{r(\mathbf{Y}_m) | 1 \leq m \leq M\}) \rightarrow 1. \quad (3)$$

The proof of Proposition 2.1 is provided in Appendix A. In our context, p_1 represents the data distribution of \mathcal{D} , $p_2 = q_\theta$, which is the data sampling distribution of \mathcal{D}_{q_θ} , and $r(\mathbf{x}) = r(p, q_\theta, \mathbf{x})$. Intuitively, Proposition 2.1 reveals the fact that when $|\mathcal{D}| = N$ and $|\mathcal{D}_{q_\theta}| = M$ are sufficiently large, the top- K instances selected from both sets tend to overlap. This is well-suited to our scenario, as the pre-training corpus \mathcal{D} is typically large-scale, and it is advantageous to sample as many candidates from \mathcal{D}_{q_θ} as possible to find high-reward instances.

Decoupling the Student and the Reward-Computing LM. To avoid computing reward values with q_θ , we replace it with p_{ref} , the output distribution of a tiny reference LM, typically smaller than the student LM, for reward computation. The reference LM is pre-trained on a small subset \mathcal{D}_{ref} , uniformly sampled from \mathcal{D} , allowing p_{ref} to approximate q_θ with minimal computation. This is a reasonable approximation because q_θ evaluates the difficulty of instances, and the relative data difficulties generally remain consistent across different models [21]. As a result, the reward function in Eq. (2) is replaced with $r(p, p_{\text{ref}}, \mathbf{x}) = \log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})}$.

In summary, *Difference Sampling* constructs a pre-training corpus \mathcal{D}' from $\mathcal{D} - \mathcal{D}_{\text{ref}}$ as follows:

$$\mathcal{D}' = \text{top-}K \left\{ \log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} \mid \mathbf{x} \in \mathcal{D} - \mathcal{D}_{\text{ref}} \right\}, \quad (4)$$

which is independent of q_θ . As shown in Figure 3(b), *Difference Sampling* essentially refines the data distribution of \mathcal{D} by comparing the “difference” between p and p_{ref} , increasing difficulty and diversity while filtering out noise. A detailed discussion of these effects is provided in Section 2.5.

2.3 Pre-Training on Difference-Sampled Corpus

As illustrated in Figure 3(a), we pre-train the student LM from scratch on the difference-sampled corpus \mathcal{D}' with the cross-entropy loss for next-token prediction, which is similar to standard pre-training. The loss function $L(q_\theta, \mathcal{D}')$ is given by

$$L(q_\theta, \mathcal{D}') = -\frac{1}{|\mathcal{D}'|} \sum_{\mathbf{x} \in \mathcal{D}'} \frac{1}{|\mathbf{x}|} \sum_{t=1}^{|\mathbf{x}|} \log q_\theta(x_t | \mathbf{x}_{<t}), \quad (5)$$

where $|\mathbf{x}|$ is the length of \mathbf{x} , x_t is the t^{th} token, and $\mathbf{x}_{<t}$ denotes the prefix of \mathbf{x} with $t - 1$ tokens.

2.4 MINIPLM Training Pipeline

The general training pipeline of MINIPLM is as follows: (1) Uniformly sample a subset \mathcal{D}_{ref} from the pre-training corpus \mathcal{D} , ensuring $|\mathcal{D}_{\text{ref}}| \ll |\mathcal{D}|$. (2) Train a reference LM from scratch on \mathcal{D}_{ref} using the cross-entropy loss to obtain its output distribution p_{ref} . (3) Perform *Difference Sampling* with p_{ref} and the teacher LM’s output distribution p , generating a pre-training corpus \mathcal{D}' from $\mathcal{D} - \mathcal{D}_{\text{ref}}$ using Eq. (4). The size of \mathcal{D}' is controlled by a sampling ratio α , where $K = \alpha|\mathcal{D} - \mathcal{D}_{\text{ref}}|$. (4) Pre-train student LMs from scratch on \mathcal{D}' with the cross-entropy loss defined in Eq. (5).

2.5 Discussion

Efficiency and Flexibility of MINIPLM. As shown in Figure 3(a), MINIPLM relies only on $p(\mathbf{x})$ and $p_{\text{ref}}(\mathbf{x})$, the teacher and reference LMs’ probability over the entire sequence, which can be computed and stored offline because each instance in \mathcal{D} is associated with a single floating-point number, amounting to only 200MB storage for 50B tokens with a 1,024 sequence length. Once \mathcal{D}' is difference-sampled based on $p(\mathbf{x})$ and $p_{\text{ref}}(\mathbf{x})$, multiple student LMs can be efficiently pre-trained under the teacher LM’s guidance without extra computational cost. In addition, this process modifies only the pre-training data, imposing no restrictions on the architecture or tokenization of the student LM, making MINIPLM highly flexible for integration into optimized pre-training frameworks [15] and suitable for KD across model families. In contrast, online KD [58, 32] relies on per-token distributions, which are infeasible to store offline, as it takes $50B \times 150K \times 4\text{byte} = 30\text{PB}$ storage for an LM with 150K vocabulary trained on 50B tokens, making online teacher LM inference necessary. Aligning per-token distributions also demands matching tokenizers between the teacher and student [7], which complicates the re-implementation and optimization of pre-training workflows.

Effectiveness of MINIPLM. In essence, as illustrated in Figure 3(b), MINIPLM distills the teacher LM’s knowledge into the student LM’s pre-training distribution via *Difference Sampling*, producing three effects: (1) Down-sampling easy and common patterns that both the teacher and reference LM fit well, where $p(\mathbf{x}) \gtrsim p_{\text{ref}}(\mathbf{x})$ and $\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} \gtrsim 0$. (2) Up-sampling hard and diverse knowledge, which the larger teacher LM has mastered but the smaller reference LM struggles with, and thus $p(\mathbf{x}) \gg p_{\text{ref}}(\mathbf{x})$ and $\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} \gg 0$. (3) Discarding noisy and harmful instances that the teacher LM assigns lower probabilities to than the reference LM, where $\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} < 0$. We provide examples of these effects in Appendix C.5. Training on the distribution with these effects encourages the student model to focus more on the sophisticated world knowledge learned by the teacher without being distracted by the noise. Note that without the comparison between p and p_{ref} , the effect (1) disappears, leading to a pre-training corpus dominated by common patterns. This explains the effectiveness of MINIPLM against prior offline KD methods [41, 63], where maintaining the data difficulty and diversity, critical to pre-training [70], is challenging without extensive human efforts on prompt engineering [33].

3 Experiments

3.1 Experimental Setup

Model. We adopt the Qwen-1.5 [3] architecture in our experiments. We use the officially released 1.8B Qwen-1.5 model as the teacher LM and distill its knowledge into students with 200M, 500M, and 1.2B parameters. Detailed model configurations are provided in Appendix B.

Pre-Training. We construct pre-training corpora from the Pile [23]. To control the computation in experiments, we pre-train all LMs on a maximum of 50B tokens, where documents are merged to construct instances with sequence lengths of 1,024. For online KD methods that incur additional train-time computation, we reduce their training steps to align the total training computation with pre-training without KD or offline KD methods. See Appendix B for more pre-training details.

Baselines. We compare MINIPLM with 4 baselines:

- **Pre-Train w/o KD** pre-trains the student LM on a 50B corpus uniformly sampled from the Pile dataset, without the guidance of the teacher LM’s knowledge.
- **Vanilla KD** [58] minimizes the token-level forward KLD between p and q_{θ} , which requires **online** inference of the teacher LM to obtain the token-level output distributions.
- **SeqKD** [41] trains the student LM on the teacher-generated data. Since it is infeasible to generate all 50B tokens, we approximate Kim and Rush [41] by using the first 768 tokens of each instance from the training corpus in **Pre-Train w/o KD** as the prompts and let the teacher LM generate the remaining tokens **offline**.
- **MiniLLM** [32] minimizes the reverse KLD between p and q_{θ} with PPO [68], which requires **online** inference of the teacher LM and **online** sampling from the student LM. We treat the first 768 tokens of the instances from the training corpus of **Pre-Train w/o KD** as the prompts and sample 256 tokens from q_{θ} during the exploration of PPO.

	HS	LAM	Wino	OBQA	ARC-e	ARC-c	PIQA	SIQA	Story	Avg.
1.8B Teacher → 200M Student										
Pre-Train w/o KD	31.1	32.4	49.9	27.6	38.9	23.1	61.8	36.4	58.1	39.9
Vanilla KD	30.4	31.0	51.4	26.6	40.1	23.1	62.2	36.9	57.3	39.9
MiniLLM	30.2	29.4	50.0	26.6	39.0	21.3	60.5	36.6	57.6	39.0
SeqKD	30.5	31.0	51.3	27.4	39.3	22.4	61.3	36.9	57.4	39.7
MINIPLM	32.7	35.4	51.4	27.2	40.6	23.7	63.3	37.0	60.0	41.3
1.8B Teacher → 500M Student										
Pre-Train w/o KD	35.8	40.1	51.0	30.2	41.7	24.4	65.4	38.2	61.4	43.2
Vanilla KD	37.0	39.9	51.7	29.4	45.1	24.2	65.8	38.0	61.6	43.6
MiniLLM	33.0	35.4	51.2	27.5	42.1	24.2	62.3	37.3	60.2	41.5
SeqKD	34.9	37.9	50.7	28.6	42.7	23.6	65.0	38.4	58.9	42.3
MINIPLM	39.0	42.6	52.2	30.2	45.8	24.9	67.0	39.0	62.2	44.8
1.8B Teacher → 1.2B Student										
Pre-Train w/o KD	39.4	44.5	51.8	28.4	46.0	25.7	67.0	39.5	62.2	44.9
Vanilla KD	40.7	43.3	53.2	29.8	46.1	25.5	67.3	39.2	63.5	45.4
MiniLLM	36.1	42.5	51.2	28.5	44.1	25.3	65.8	37.9	61.4	43.6
SeqKD	38.5	41.4	51.9	29.2	46.5	25.1	66.3	39.0	61.0	44.3
MINIPLM	42.8	46.2	53.3	31.0	46.8	26.9	68.3	39.8	64.0	46.6

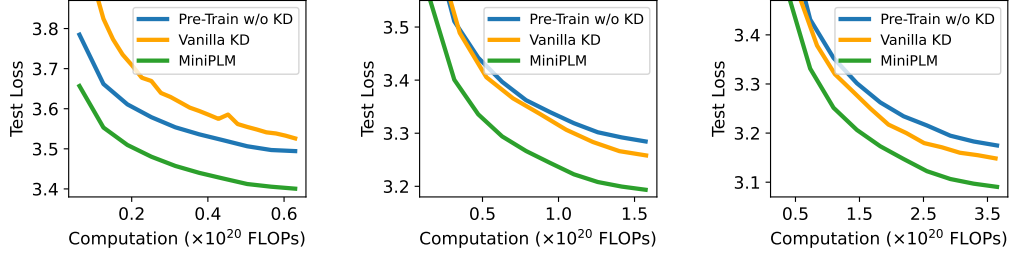
Table 1: Zero-shot accuracy scores on 9 widely-used downstream tasks and the average scores (Avg.). We use the Qwen-1.5 1.8B LM [3] as the teacher and Qwen LMs with 200M, 500M, and 1.2B parameters as the student. Student LMs with the same sizes consume the same training-time computation. The best scores of each model size are **boldfaced**.

MINIPLM. We employ a 104M reference LM trained on 5B tokens. In Section 3.2 and 3.3, we consider a setting where \mathcal{D} is sufficiently large, containing 105B tokens uniformly sampled from the Pile corpus. We reserve 5B tokens as \mathcal{D}_{ref} , and conduct *Difference Sampling* as per Eq. (4) on the other 100B tokens by setting $\alpha = 0.5$ to construct a 50B-token corpus \mathcal{D}' . In this way, the student LM is pre-trained on \mathcal{D}' for one epoch. We use the loss difference of the teacher and reference LM to sampled instances from \mathcal{D} , which is equivalent to Eq. (4) as all instances have 1,024 tokens. In Section 3.4, we evaluate MINIPLM in a data-limited setting, where \mathcal{D} is controlled to contain 50B tokens and the student LM is trained on the difference-sampled data for multiple epochs. See Appendix C.4 for the ablation studies on the reference LM and sampling ratio α .

Evaluation. We assess the zero-shot accuracy of LMs trained with different methods on 9 downstream tasks widely used in examining the foundation abilities of base models [78, 26]. We also test the language modeling capability of the LMs on a subset of DCLM [45], a high-quality corpus carefully curated with complex pipelines, to examine how well LMs capture broad and diverse knowledge. See Appendix B for more evaluation details.

3.2 Main Results

MINIPLM Improves Downstream Performance. Table 1 shows zero-shot accuracy on downstream tasks for LMs trained by different methods, leading to three key observations. *First*, among all the baselines, only Vanilla KD outperforms Pre-Train w/o KD for relatively large student LM, given a constant computation budget. This highlights the room for improvement in KD for pre-training, especially when the gap between the teacher and the student is substantial. *Second*, MINIPLM-trained model achieves the best performance across most of the tasks. Compared to Pre-Train w/o KD, MINIPLM effectively leverages the teacher LM’s knowledge to improve the student LM pre-training. Compared to online methods like Vanilla KD and MiniLLM, MINIPLM, as shown in Figure 3(a), incurs no additional training-time overhead, allowing the student LM to be optimized for more steps and leading to higher performance. Compared to offline methods like SeqKD, MINIPLM, as shown in Figure 3(b), uses a reference LM to ensure sufficient difficulty and diversity of the pre-training corpus, which is essential for the student LM to learn versatile sophisticated knowledge during pre-training and generalize across various downstream tasks. *Finally*, the improvements of MINIPLM against Vanilla KD scale well with the student LM size, which is also illustrated in Figure 1(b).



(a) 1.8B Teacher \rightarrow 200M Student (b) 1.8B Teacher \rightarrow 500M Student (c) 1.8B Teacher \rightarrow 1.2B Student

Figure 4: Language modeling loss on the DCLM [45] subset. We distill the knowledge of the 1.8B Qwen model [3] into student LMs from the Qwen family with 200M, 500M, and 1.2B parameters. We control the total training-time FLOPs of different methods to be the same.

MINIPLM Helps Language Modeling. Figure 4, compares the language modeling performance of the MINIPLM-trained LMs and baselines on the DCLM [45] subset, a diverse and high-quality dataset curated from web corpora. The results show that, given the same training-time FLOPs, LMs trained with MINIPLM achieve the lowest test losses. In Table 2, we extrapolate the test losses with the Scaling Law [37] to simulate pre-training on 1T and 10T tokens (details in Appendix C.2), showing that MINIPLM maintains its advantages at the scales of pre-training recent large LMs [78, 19]. A critical stage in DCLM’s data-cleaning pipeline involves removing easy and common patterns, thereby enhancing challenging and diverse signals. Therefore, a lower test loss on DCLM suggests that, with the teacher LM’s guidance and the reference LM, the MINIPLM-trained LMs learn the diverse and hard knowledge better due to the up-sampling of the corresponding parts in the pre-training distribution, as illustrated in Figure 3(b).

N_{stu}	Method	L_{1T}	L_{10T}
200M	Pre-Train w/o KD	3.35	3.32
	Vanilla KD	3.39	3.35
	MINIPLM	3.28	3.26
500M	Pre-Train w/o KD	3.12	3.08
	Vanilla KD	3.12	3.07
	MINIPLM	3.06	3.04
1.2B	Pre-Train w/o KD	2.98	2.94
	Vanilla KD	2.95	2.91
	MINIPLM	2.92	2.88

Table 2: Test loss predictions using the Scaling Law [37]. N_{stu} : the student LM size. L_{1T} , L_{10T} : the loss when Pre-Train w/o KD and MINIPLM process 1T and 10T tokens, with Vanilla KD consuming the same training FLOPs.

MINIPLM Reduces Training Computation. We plot the 500M student LM’s average zero-shot accuracy scores on the downstream tasks in Figure 1(a) with respect to its pre-training FLOPs. MINIPLM achieves the same performance as Vanilla KD while reducing computational costs by 2.2 times. Similar trends are observed for other student LMs (Figure 7) and in the test loss curves on DCLM corpus (Figure 4). The efficiency gains are more pronounced when comparing MINIPLM with Pre-Train w/o KD on the 500M and 1.2B models. We attribute this acceleration to *Difference Sampling*, which down-samples common patterns and filters out noisy signals from the pre-training corpus, as shown in Figure 3(b). As a result, the model trained with MINIPLM avoids wasting computation on learning the easy knowledge quickly memorized during the early training stage and is less distracted by the noisy outliers that slow the convergence down.

3.3 KD Across Model Families

A noticeable advantage of MINIPLM over Vanilla KD and MiniLLM is its flexibility to distill the knowledge of a teacher LM into student LMs with completely different tokenizers and architectures without additional strategies like Boizard et al. [7]. In Table 3, we illustrate the performance when using the LMs from the Qwen [3] family as the teacher and the reference LM to distill knowledge into a 212M Llama3.1 [19] model and a 140M Mamba [27] model. The results demonstrate the promising performance of MINIPLM in KD across model families, outperforming Pre-Train w/o KD and the existing offline KD baseline (SeqKD). This allows emerging LMs with novel architectures [47, 73] or advanced tokenization [74, 25] to inherit knowledge from existing LMs, thereby facilitating the development of more efficient and higher-performed models.

	Llama3.1		Mamba	
	Acc.	Loss	Acc.	Loss
Pre-Train w/o KD	41.0	3.52	41.6	3.24
SeqKD	40.8	3.54	41.0	3.27
MINIPLM	41.8	3.43	42.6	3.15

Table 3: Results of KD across model families. We use the teacher and reference LM from the Qwen family to distill the Llama3.1 and Mamba models. The average zero-shot accuracies on the downstream tasks and the losses on the DCLM corpus are reported. Note that Vanilla KD and MiniLLM cannot be applied when the teacher and student LMs use different tokenizations.

3.4 Data-Limited Setting

We evaluate MINIPLM in a data-limited setting where \mathcal{D} is constrained to contain 50B tokens. To this end, the student LM should be trained on the difference-sampled \mathcal{D}' over multiple epochs to ensure the total computation and trained tokens remain consistent with Pre-Train without KD. We split a $|\mathcal{D}_{\text{ref}}|$ containing 1B tokens to train the reference LM. Therefore, for a sampling ratio α , the student LM should be trained for $\frac{|\mathcal{D}|}{\alpha|\mathcal{D}-\mathcal{D}_{\text{ref}}|} \approx \frac{1}{\alpha}$ epochs, given that $|\mathcal{D}_{\text{ref}}| \ll |\mathcal{D}|$. In Figure 5, we plot the loss curves of the 200M student LMs on the DCLM corpus when using $\alpha \in [0.5, 0.25, 0.125]$ and training the LM for around 2, 4, and 8 epochs, respectively. We can see that difference-sampling 25% data (with $\alpha = 0.25$) and training the student LM for 4 epochs yields the best performance, which aligns with the observations in Muennighoff et al. [57]. The corpus sampled with a higher α does not achieve the best quality and diversity offered by *Difference Sampling*, while a lower α leads to rapid over-fitting of the student LM. These findings suggest that MINIPLM is a promising approach to enhance data utilization when high-quality web corpora become scarce [80]. By extrapolating the loss curve of Pre-Train w/o KD using the Scaling Law in a data-constrained setting [57], we estimate that it would require an additional 68B training tokens to match the performance of MINIPLM ($\alpha = 0.25$, 4 epochs), which means MINIPLM reduces the pre-training data requirement by 2.4 times. See Appendix C.2 for more details on this extrapolation.

3.5 Analysis

Impact of Teacher Model. Intuitively, larger teacher LMs will lead to better KD results, which is observed in recent works [32]. However, early works have also shown that an excessively large gap between teacher and student models can hinder effective KD [55]. In Figure 6, we plot the performance of Vanilla KD and MINIPLM when distilling teacher LMs with different sizes into a 200M student model. We observe a similar phenomenon to Mirzadeh et al. [55] on Vanilla KD and MINIPLM that larger teacher LMs are not necessarily more helpful for KD. However, we attribute this to different factors for Vanilla KD and MINIPLM. In Vanilla KD, the overhead introduced by larger LMs during training diminishes the benefits of distillation. As for MINIPLM, a 500M teacher LM proves most effective for distilling into a 200M student LM. Smaller teacher LMs (e.g., 300M) lack the capacity to identify the hard but valuable parts of the pre-training distribution, weakening the effectiveness of *Difference Sampling* as shown in Figure 3(b). On the other hand, the value scale of $\log p(\mathbf{x})$ from oversized LMs (e.g., 4B) becomes too small compared to that of the reference LM, which tends to degenerate *Difference Sampling* into sampling with $p_{\text{ref}}(\mathbf{x})$ only, losing the effect of the teacher LM. Future research could focus on optimizing teacher model size for MINIPLM or mitigating the impact of differing $\log p(\mathbf{x})$ and $\log p_{\text{ref}}(\mathbf{x})$ value scales.

Diversity of Difference-Sampled Data. To further verify the effectiveness of *Difference Sampling* on improving the diversity of the pre-training corpus, in Table 4, we follow Friedman and Dieng [22] to compute the semantic diversity of the difference-sampled corpus and the data used in other baseline pre-training approaches. The results show that the difference-sampled corpus has the highest diversity,

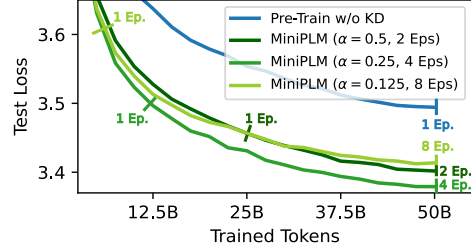


Figure 5: MINIPLM in the data-constrained setting. We fix \mathcal{D} to contain 50B tokens and alter the sampling ratio α to obtain \mathcal{D}' with *Difference Sampling*, which will be trained on for multiple epochs to achieve the constant total trained tokens. The y-axis represents the test loss on the DCLM corpus.

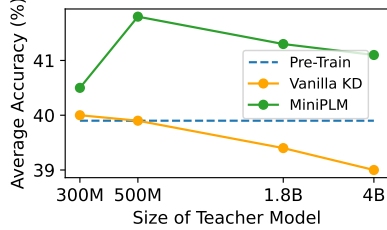


Figure 6: Impact of the teacher LM’s sizes on Vanilla KD and MINIPLM, with the pre-training FLOPs aligned. The y-axis represents the average zero-shot accuracy on the downstream tasks.

Pre-Training Corpus	Usage	Diversity
Original	Pre-Train w/o KD & Vanilla KD	32.25
Teacher-Generated	SeqKD	30.16
Difference-Sampled	MINIPLM	36.70

Table 4: Semantic diversity of the original pre-training corpus \mathcal{D} used in Pre-Train w/o KD and Vanilla KD, the teacher-generated corpus used in SeqKD, and the difference-sampled corpus \mathcal{D}' in MINIPLM. *Difference Sampling* increases the diversity of the refined pre-training distribution, which helps LM pre-training.

despite that *Difference Sampling* is derived from minimizing the reverse KLD, which exhibits the mode-seeking behavior [54]. We suspect the reason is that *Difference Sampling* down-samples the easy parts of the corpus containing repeated contents while up-sampling the hard parts consisting of diverse texts. These two components, with large $p(x)$ values as seen in Figure 3(b), constitute the major modes of the teacher LM that q_θ seeks during optimizing Eq. (1). Therefore, the mode-seeking behavior helps remove the noisy parts of the pre-training distribution, and the loss of diversity due to noise reduction is compensated by the up-sampling of the hard and diverse data points. We provide a case study in Appendix C.5 to further explain our argument.

Combining Vanilla KD and MINIPLM. As shown in Table 1 and Figure 4, Vanilla KD improves the 500M and 1.2B student LM performance compared to Pre-Train w/o KD, suggesting the potential of further improving MINIPLM by combining it with Vanilla KD. This combined approach, termed “MINIPLM + Vanilla KD”, applies Vanilla KD to the difference-sampled corpus used in MINIPLM. In Table 5, we compare “MINIPLM + Vanilla KD” with the individual use of Vanilla KD and MINIPLM, using a 1.8B teacher LM. The results show that when pre-training student LMs with 500B and 1.2B parameters, “MINIPLM + Vanilla KD” further improves the performance, given the same training-time FLOPs. This demonstrates that MINIPLM and Vanilla KD complement each other: MINIPLM distills the coarse-grained sequence-level knowledge of the teacher LM into the student LM via the pre-training data, while Vanilla KD directly aligns the token-level probability distribution between p and q_θ , providing fine-grained token-level signals.

N_{stu}	Method	Acc.
200M	Vanilla KD	39.9
	MINIPLM	41.3
	MINIPLM + Vanilla KD	40.7
500M	Vanilla KD	43.6
	MINIPLM	44.8
	MINIPLM + Vanilla KD	44.9
1.2B	Vanilla KD	45.4
	MINIPLM	46.6
	MINIPLM + Vanilla KD	48.1

Table 5: Average accuracy on downstream tasks when combining MINIPLM and Vanilla KD. “MINIPLM + Vanilla KD”: applying Vanilla KD to pre-train student LMs on the difference-sampled corpus in MINIPLM. N_{stu} : the size of student LMs.

4 Related Work

Language Model Pre-Training. Pre-training is the critical phase for language models (LMs; 10, 60, 77, 15, 78, 79, 19) to obtain their foundation abilities for various downstream tasks. To improve pre-training, some works focus on data curation, such as adjusting domain mixing [86, 89], selecting valuable data points relevant to desired tasks [9, 20, 31], or transforming the instances based on downstream requirements [13, 28, 29]. Another line of work improves the optimization during pre-training by solving better data reweighting strategies [30], designing more effective optimizers [48, 69], or discovering better training recipes [39, 37]. Variations in model architectures [87] and training objectives [75] are also explored to boost pre-training stability and final LM performance. In this work, we utilize the knowledge from existing LMs to enhance pre-training.

Small Language Models. Given the high computational demands of large LMs during inference, there has been growing interest in pre-training small LMs [67]. However, achieving high performance

with limited parameter sizes remains challenging because the training computation that small LMs need to match the capabilities of large LMs often exceeds Chinchilla-optimal [37] and scales as a power law with respect to the model size gap [40]. Despite these challenges, recent efforts have made promising progress by data quality improvement [52, 5, 91] or model pruning [58, 85]. We explore knowledge distillation as a complementary approach.

Knowledge Distillation. Knowledge distillation (KD; [36]) uses a large teacher model to improve the performance of a small student model, which is widely used to build efficient neural network systems [62, 17, 64]. In NLP, early works primarily apply KD to encoder-only models [18, 49] for text classification by aligning token-level distribution [65], hidden states [72], and attention matrices [81, 82]. For generative LMs, a straightforward KD approach is training small LMs on the texts generated by large LMs [14, 63, 38]. Other works [32, 1, 83, 46, 84] explore better optimization objectives. However, these works focus on KD for fine-tuning LMs, while pre-training is critical for establishing core LM capabilities [2]. Therefore, we investigate KD in the pre-training stage to develop strong small base LMs.

5 Conclusion

Summary. In this work, we find it non-trivial to adapt existing KD approaches for fine-tuning LMs to the pre-training stage because of the high overhead brought by the teacher LM inference in online KD and the tendency of losing data difficulty and diversity in offline KD. Therefore, we propose MINIPLM to address these issues through *Difference Sampling*, which refines the training distribution by down-sampling easy patterns, up-sampling hard instances, and filtering out noisy data points, with the knowledge of the difference between the large teacher LM and a small reference LM. The offline nature of MINIPLM makes it both efficient and flexible to distill student LMs with diverse configurations. The use of the large-small-model differences ensures the difficulty and diversity of the refined pre-training distribution. Using a 1.8B LM as the teacher to guide the pre-training of 200M, 500M, and 1.2B LMs, we demonstrate that MINIPLM improves the student LMs’ performance on 9 downstream tasks, enhances their language modeling capability, and reduces pre-training computation. Additionally, MINIPLM improves the utilization of limited pre-training data and can distill teacher LM’s knowledge into student LMs from completely different families.

Limitations. One limitation of MINIPLM is that it requires the large LMs’ probabilities on texts from the pre-training corpus, making black-box KD for close-source LMs [60, 76] with MINIPLM challenging. For some APIs [59], this issue can be solved by specifying user-provided bias in the softmax operation, which allows obtaining the probabilities of given tokens one by one [12], at the expense of a large number of API calls.

Future Work. A promising future direction to explore is applying the difference-sampled corpus to pre-train LMs larger than the teacher LM, enabling weak-to-strong generalization [11]. Since data properties are critical for pre-training LMs across various sizes, the improvement of data diversity and difficulty is likely to be beneficial for LMs larger than the difference-sampling models.

References

- [1] R. Agarwal, N. Vieillard, Y. Zhou, P. Stanczyk, S. R. Garea, M. Geist, and O. Bachem. On-policy distillation of language models: Learning from self-generated mistakes. In *Proceedings of ICLR*, 2024. URL <https://openreview.net/forum?id=3zKtaqxLhW>.
- [2] Z. Allen-Zhu and Y. Li. Physics of language models: Part 3.1, knowledge storage and extraction. In *Proceedings of ICML*, 2024. URL <https://openreview.net/forum?id=5x788rqbcbj>.
- [3] J. Bai, S. Bai, Y. Chu, Z. Cui, K. Dang, X. Deng, Y. Fan, W. Ge, Y. Han, F. Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023. URL <https://arxiv.org/pdf/2309.16609>.
- [4] Y. Bai, S. Kadavath, S. Kundu, A. Askell, J. Kernion, A. Jones, A. Chen, A. Goldie, A. Mirhoseini, C. McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022. URL <https://arxiv.org/abs/2212.08073>.

- [5] M. Bellagente, J. Tow, D. Mahan, D. Phung, M. Zhuravinskyi, R. Adithyan, J. Baicoianu, B. Brooks, N. Cooper, A. Datta, et al. Stable LM 2 1.6B technical report. *arXiv preprint arXiv:2402.17834*, 2024. URL <https://arxiv.org/abs/2402.17834>.
- [6] Y. Bisk, R. Zellers, J. Gao, Y. Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of AAAI*, 2020. URL <https://ojs.aaai.org/index.php/AAAI/article/view/6239/6095>.
- [7] N. Boizard, K. El-Haddad, C. Hudelot, and P. Colombo. Towards cross-tokenizer distillation: the universal logit distillation loss for llms. *arXiv preprint arXiv:2402.12030*, 2024. URL <https://arxiv.org/abs/2402.12030>.
- [8] R. Bommasani, D. A. Hudson, E. Adeli, R. Altman, S. Arora, S. von Arx, M. S. Bernstein, J. Bohg, A. Bosselut, E. Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021. URL <https://arxiv.org/abs/2108.07258>.
- [9] D. Brandfonbrener, H. Zhang, A. Kirsch, J. R. Schwarz, and S. Kakade. CoLoR-Filter: Conditional loss reduction filtering for targeted language model pre-training. *arXiv preprint arXiv:2406.10670*, 2024. URL <https://arxiv.org/abs/2406.10670>.
- [10] T. Brown, B. Mann, N. Ryder, M. Subbiah, et al. Language models are few-shot learners. In *Proceedings of NeurIPS*, 2020. URL <https://papers.nips.cc/paper/2020/hash/1457c0d6bfbcb4967418bfb8ac142f64a-Abstract.html>.
- [11] C. Burns, P. Izmailov, J. H. Kirchner, B. Baker, L. Gao, L. Aschenbrenner, Y. Chen, A. Ecoffet, M. Joglekar, J. Leike, et al. Weak-to-strong generalization: Eliciting strong capabilities with weak supervision. *arXiv preprint arXiv:2312.09390*, 2023. URL <https://openai.com/index/weak-to-strong-generalization/>.
- [12] N. Carlini, D. Paleka, K. D. Dvijotham, T. Steinke, J. Hayase, A. F. Cooper, K. Lee, M. Jagielski, M. Nasr, A. Conmy, et al. Stealing part of a production language model. In *Proceedings of ICML*, 2024. URL <https://arxiv.org/abs/2403.06634>.
- [13] D. Cheng, Y. Gu, S. Huang, J. Bi, M. Huang, and F. Wei. Instruction pre-training: Language models are supervised multitask learners. In *Proceedings of EMNLP*, 2024. URL <https://arxiv.org/abs/2406.14491>.
- [14] W.-L. Chiang, Z. Li, Z. Lin, Y. Sheng, Z. Wu, H. Zhang, L. Zheng, S. Zhuang, Y. Zhuang, J. E. Gonzalez, I. Stoica, and E. P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- [15] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022. URL <https://arxiv.org/abs/2204.02311>.
- [16] P. Clark, I. Cowhey, O. Etzioni, T. Khot, A. Sabharwal, C. Schoenick, and O. Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv:1803.05457v1*, 2018. URL <https://arxiv.org/abs/1803.05457>.
- [17] W. M. Czarnecki, R. Pascanu, S. Osindero, S. Jayakumar, G. Swirszcz, and M. Jaderberg. Distilling policy distillation. In *Proceedings of AISTATS*, 2019. URL <http://proceedings.mlr.press/v89/czarnecki19a/czarnecki19a.pdf>.
- [18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, 2019. URL <https://aclanthology.org/N19-1423.pdf>.
- [19] A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Yang, A. Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. URL <https://arxiv.org/abs/2407.21783>.
- [20] L. Engstrom, A. Feldmann, and A. Madry. Dsdm: Model-aware dataset selection with datamodels. *arXiv preprint arXiv:2401.12926*, 2024. URL <https://arxiv.org/abs/2401.12926>.

- [21] K. Ethayarajh, Y. Choi, and S. Swayamdipta. Understanding dataset difficulty with v-usable information. In *Proceedings of ICML*, 2022. URL <https://proceedings.mlr.press/v162/ethayarajh22a/ethayarajh22a.pdf>.
- [22] D. Friedman and A. B. Dieng. The vendi score: A diversity evaluation metric for machine learning. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=g970HbQyk1>.
- [23] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima, et al. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020. URL <https://arxiv.org/abs/2101.00027>.
- [24] L. Gao, J. Tow, B. Abbasi, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, et al. A framework for few-shot language model evaluation, 07 2024. URL <https://zenodo.org/records/12608602>.
- [25] N. Godey, R. Castagné, É. de la Clergerie, and B. Sagot. MANTa: Efficient gradient-based tokenization for end-to-end robust language modeling. In *Findings of EMNLP*, 2022. URL <https://aclanthology.org/2022.findings-emnlp.207>.
- [26] D. Groeneveld, I. Beltagy, P. Walsh, A. Bhagia, R. Kinney, O. Tafjord, A. H. Jha, H. Ivison, I. Magnusson, Y. Wang, et al. OLMo: Accelerating the science of language models. *arXiv preprint arXiv:2402.00838*, 2024. URL <https://arxiv.org/abs/2402.00838>.
- [27] A. Gu and T. Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. URL <https://arxiv.org/abs/2312.00752>.
- [28] Y. Gu, X. Han, Z. Liu, and M. Huang. PPT: Pre-trained prompt tuning for few-shot learning. In *Proceedings of ACL*, 2022. URL <https://arxiv.org/abs/2109.04332>.
- [29] Y. Gu, L. Dong, F. Wei, and M. Huang. Pre-training to learn in context. In *Proceedings of ACL*, 2023. URL <https://aclanthology.org/2023.acl-long.267/>.
- [30] Y. Gu, L. Dong, Y. Hao, Q. Dong, M. Huang, and F. Wei. Towards optimal learning of language models. *arXiv preprint arXiv:2402.17759*, 2024. URL <https://arxiv.org/abs/2402.17759>.
- [31] Y. Gu, L. Dong, H. Wang, Y. Hao, Q. Dong, F. Wei, and M. Huang. Data selection via optimal control for language models. *arXiv preprint arXiv:2410.07064*, 2024. URL <https://arxiv.org/abs/2410.07064>.
- [32] Y. Gu, L. Dong, F. Wei, and M. Huang. MiniLLM: Knowledge distillation of large language models. In *Proceedings of ICLR*, 2024. URL <https://openreview.net/forum?id=5h0qf7IBZZ>.
- [33] S. Gunasekar, Y. Zhang, J. Aneja, C. C. T. Mendes, A. Del Giorno, S. Gopi, M. Javaheripi, P. Kauffmann, G. de Rosa, O. Saarikivi, et al. Textbooks are all you need. *arXiv preprint arXiv:2306.11644*, 2023.
- [34] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of ICML*, 2017. URL <http://proceedings.mlr.press/v70/haarnoja17a.html?ref=https://githubhelp.com>.
- [35] X. Han, Z. Zhang, N. Ding, Y. Gu, et al. Pre-trained models: Past, present and future. *AI Open*, 2021. URL <https://www.sciencedirect.com/science/article/pii/S26666510210002319>.
- [36] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. URL <https://arxiv.org/pdf/1503.02531.pdf>.
- [37] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. de Las Casas, L. A. Hendricks, et al. Training compute-optimal large language models. In *Proceedings of NeurIPS*, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/c1e2faff6f588870935f114ebe04a3e5-Paper-Conference.pdf.

- [38] C.-Y. Hsieh, C.-L. Li, C.-k. Yeh, H. Nakhost, Y. Fujii, A. Ratner, R. Krishna, C.-Y. Lee, and T. Pfister. Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes. In *Findings of the ACL*, 2023. URL <https://aclanthology.org/2023.findings-acl.507/>.
- [39] S. Hu, Y. Tu, X. Han, C. He, G. Cui, X. Long, Z. Zheng, Y. Fang, Y. Huang, W. Zhao, et al. MiniCPM: Unveiling the potential of small language models with scalable training strategies. In *Proceedings of COLM*, 2024. URL <https://openreview.net/forum?id=3X2L2TfR0f¬eId=QvwPc5chyd>.
- [40] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020. URL <https://arxiv.org/abs/2001.08361>.
- [41] Y. Kim and A. M. Rush. Sequence-level knowledge distillation. In *Proceedings of EMNLP*, 2016. URL <https://aclanthology.org/D16-1139.pdf>.
- [42] J. Ko, S. Kim, T. Chen, and S.-Y. Yun. Distillm: Towards streamlined distillation for large language models. In *Proceedings of ICML*, 2024. URL <https://openreview.net/forum?id=lsHZNNOC7r>.
- [43] H. Levesque, E. Davis, and L. Morgenstern. The winograd schema challenge. In *Proceedings of KR*, 2012. URL <https://dl.acm.org/doi/10.5555/3031843.3031909>.
- [44] S. Levine, A. Kumar, G. Tucker, and J. Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020. URL <https://arxiv.org/pdf/2005.01643.pdf>.
- [45] J. Li, A. Fang, G. Smyrnis, M. Ivgi, M. Jordan, S. Gadre, H. Bansal, E. Guha, S. Keh, K. Arora, et al. DataComp-LM: In search of the next generation of training sets for language models. *arXiv preprint arXiv:2406.11794*, 2024. URL <https://arxiv.org/abs/2406.11794>.
- [46] Y. Li, Y. Gu, L. Dong, D. Wang, Y. Cheng, and F. Wei. Direct preference knowledge distillation for large language models. *arXiv preprint arXiv:2406.19774*, 2024. URL <https://arxiv.org/abs/2406.19774>.
- [47] O. Lieber, B. Lenz, H. Bata, G. Cohen, J. Osin, I. Dalmedigos, E. Safahi, S. Meirom, Y. Belinkov, S. Shalev-Shwartz, et al. Jamba: A hybrid transformer-mamba language model. *arXiv preprint arXiv:2403.19887*, 2024. URL <https://arxiv.org/abs/2403.19887>.
- [48] H. Liu, Z. Li, D. L. W. Hall, P. Liang, and T. Ma. Sophia: A scalable stochastic second-order optimizer for language model pre-training. In *Proceedings of ICLR*, 2024. URL <https://openreview.net/forum?id=3xHDeA8Noi>.
- [49] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. RoBERTa: A robustly optimized BERT pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019. URL <https://arxiv.org/abs/1907.11692>.
- [50] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. In *Proceedings of ICLR*, 2019. URL <https://openreview.net/forum?id=Bkg6RiCqY7>.
- [51] Z. Lu, X. Li, D. Cai, R. Yi, F. Liu, X. Zhang, N. D. Lane, and M. Xu. Small language models: Survey, measurements, and insights. *arXiv preprint arXiv:2409.15790*, 2024. URL <https://arxiv.org/abs/2409.15790>.
- [52] S. Mehta, M. H. Sekhvat, Q. Cao, M. Horton, Y. Jin, C. Sun, I. Mirzadeh, M. Najibi, D. Belenko, P. Zatloukal, et al. OpenELM: An efficient language model family with open-source training and inference framework. *arXiv preprint arXiv:2404.14619*, 2024. URL <https://arxiv.org/abs/2404.14619>.
- [53] T. Mihaylov, P. Clark, T. Khot, and A. Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *Proceedings of EMNLP*, 2018. URL <https://api.semanticscholar.org/CorpusID:52183757>.

- [54] T. Minka et al. Divergence measures and message passing. Technical report, Citeseer, 2005. URL <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tr-2005-173.pdf>.
- [55] S. I. Mirzadeh, M. Farajtabar, A. Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh. Improved knowledge distillation via teacher assistant. In *Proceedings of AAAI*, 2020. URL <https://ojs.aaai.org/index.php/AAAI/article/view/5963/5819>.
- [56] N. Mostafazadeh, N. Chambers, X. He, D. Parikh, D. Batra, L. Vanderwende, P. Kohli, and J. Allen. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of NAACL-HLT*, 2016. URL <https://aclanthology.org/N16-1098/>.
- [57] N. Muennighoff, A. M. Rush, B. Barak, T. L. Scao, N. Tazi, A. Piktus, S. Pyysalo, T. Wolf, and C. Raffel. Scaling data-constrained language models. In *Proceedings of NeurIPS*, 2023. URL <https://openreview.net/forum?id=j5BuTrEj35>.
- [58] S. Muralidharan, S. T. Sreenivas, R. Joshi, M. Chochowski, M. Patwary, M. Shoenybi, B. Catanzaro, J. Kautz, and P. Molchanov. Compact language models via pruning and knowledge distillation. *arXiv preprint arXiv:2407.14679*, 2024. URL <https://www.arxiv.org/abs/2407.14679>.
- [59] OpenAI. OpenAI: Introducing ChatGPT, 2022. URL <https://openai.com/blog/chatgpt>.
- [60] OpenAI. GPT-4 technical report, 2023. URL <https://cdn.openai.com/papers/gpt-4.pdf>.
- [61] D. Paperno, G. Kruszewski, A. Lazaridou, N.-Q. Pham, R. Bernardi, S. Pezzelle, M. Baroni, G. Boleda, and R. Fernández. The lambda dataset: Word prediction requiring a broad discourse context. In *Proceedings of ACL*, 2016. URL <https://aclanthology.org/P16-1144.pdf>.
- [62] W. Park, D. Kim, Y. Lu, and M. Cho. Relational knowledge distillation. In *Proceedings of CVPR*, 2019. URL https://openaccess.thecvf.com/content_CVPR_2019/papers/Park_Relational_Knowledge_Distillation_CVPR_2019_paper.pdf.
- [63] B. Peng, C. Li, P. He, M. Galley, and J. Gao. Instruction tuning with GPT-4. *arXiv preprint arXiv:2304.03277*, 2023. URL <https://arxiv.org/abs/2304.03277>.
- [64] T. Salimans and J. Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022. URL <https://openreview.net/pdf/3b30857a628099896b6123e85d6cf04c59abe77b.pdf>.
- [65] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. DistilBERT, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019. URL <https://arxiv.org/pdf/1910.01108.pdf>.
- [66] M. Sap, H. Rashkin, D. Chen, R. Le Bras, and Y. Choi. Social IQa: Commonsense reasoning about social interactions. In *Proceedings of EMNLP*, 2019. URL <https://aclanthology.org/D19-1454>.
- [67] N. Sardana, J. Portes, S. Doubov, and J. Frankle. Beyond chinchilla-optimal: Accounting for inference in language model scaling laws. In *Proceedings of ICML*, 2024. URL <https://openreview.net/forum?id=0bmXrtTDUu>.
- [68] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017. URL <https://arxiv.org/pdf/1707.06347.pdf>.
- [69] N. Shazeer and M. Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In *Proceedings of ICML*, 2018. URL <https://proceedings.mlr.press/v80/shazeer18a/shazeer18a.pdf>.
- [70] I. Shumailov, Z. Shumaylov, Y. Zhao, N. Papernot, R. Anderson, and Y. Gal. Ai models collapse when trained on recursively generated data. *Nature*, 631(8022):755–759, 2024. URL <https://www.nature.com/articles/s41586-024-07566-y>.

- [71] N. Stiennon, L. Ouyang, J. Wu, D. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and P. F. Christiano. Learning to summarize with human feedback. In *Proceedings of NeurIPS*, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/1f89885d556929e98d3ef9b86448f951-Paper.pdf>.
- [72] S. Sun, Y. Cheng, Z. Gan, and J. Liu. Patient knowledge distillation for BERT model compression. In *Proceedings EMNLP*, 2019. URL <https://aclanthology.org/D19-1441>.
- [73] Y. Sun, L. Dong, S. Huang, S. Ma, Y. Xia, J. Xue, J. Wang, and F. Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023. URL <https://arxiv.org/abs/2307.08621>.
- [74] Y. Tay, V. Q. Tran, S. Ruder, J. Gupta, H. W. Chung, D. Bahri, Z. Qin, S. Baumgartner, C. Yu, and D. Metzler. Charformer: Fast character transformers via gradient-based subword tokenization. In *Proceedings of ICLR*, 2022. URL <https://openreview.net/forum?id=JtBRnr10EFN>.
- [75] Y. Tay, M. Dehghani, V. Q. Tran, X. Garcia, J. Wei, X. Wang, H. W. Chung, D. Bahri, T. Schuster, S. Zheng, et al. U12: Unifying language learning paradigms. In *Proceedings of ICLR*, 2023. URL <https://openreview.net/pdf?id=6ruVLB727MC>.
- [76] G. Team, R. Anil, S. Borgeaud, Y. Wu, J.-B. Alayrac, J. Yu, R. Soricut, J. Schalkwyk, A. M. Dai, A. Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. URL <https://arxiv.org/abs/2312.11805>.
- [77] G. Team, T. Mesnard, C. Hardin, R. Dadashi, S. Bhupatiraju, S. Pathak, L. Sifre, M. Rivière, M. S. Kale, J. Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024. URL <https://arxiv.org/abs/2403.08295>.
- [78] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample. LLaMA: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. URL <https://arxiv.org/pdf/2302.13971.pdf>.
- [79] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. URL <https://arxiv.org/abs/2307.09288>.
- [80] P. Villalobos, J. Sevilla, L. Heim, T. Besiroglu, M. Hobbhahn, and A. Ho. Will we run out of data? an analysis of the limits of scaling datasets in machine learning. *arXiv preprint arXiv:2211.04325*, 2022. URL <https://arxiv.org/abs/2211.04325>.
- [81] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou. MiniLM: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In *Proceedings of NeurIPS*, 2020. URL <https://proceedings.neurips.cc/paper/2020/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.
- [82] W. Wang, H. Bao, S. Huang, L. Dong, and F. Wei. MiniLMv2: Multi-head self-attention relation distillation for compressing pretrained transformers. In *Findings of ACL*, 2021. URL <https://aclanthology.org/2021.findings-acl.188>.
- [83] Y. Wen, Z. Li, W. Du, and L. Mou. f-divergence minimization for sequence-level knowledge distillation. In *Proceedings of ACL*, 2023. URL <https://aclanthology.org/2023.acl-long.605.pdf>.
- [84] T. Wu, C. Tao, J. Wang, Z. Zhao, and N. Wong. Rethinking kullback-leibler divergence in knowledge distillation for large language models. *arXiv preprint arXiv:2404.02657*, 2024. URL <https://arxiv.org/abs/2404.02657>.
- [85] M. Xia, T. Gao, Z. Zeng, and D. Chen. Sheared llama: Accelerating language model pre-training via structured pruning. In *Proceedings of ICLR*, 2024. URL <https://openreview.net/pdf?id=6s77hjBNfS>.

- [86] S. M. Xie, H. Pham, X. Dong, N. Du, H. Liu, Y. Lu, P. S. Liang, Q. V. Le, T. Ma, and A. W. Yu. DoReMi: Optimizing data mixtures speeds up language model pretraining. In *Proceedings of NeurIPS*, 2024. URL <https://openreview.net/forum?id=lXuByUeHhd>.
- [87] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T. Liu. On layer normalization in the transformer architecture. In *Proceedings of ICML*, 2020. URL <https://openreview.net/forum?id=B1x8anVFPr>.
- [88] X. Xu, M. Li, C. Tao, T. Shen, R. Cheng, J. Li, C. Xu, D. Tao, and T. Zhou. A survey on knowledge distillation of large language models. *arXiv preprint arXiv:2402.13116*, 2024. URL <https://arxiv.org/abs/2402.13116>.
- [89] J. Ye, P. Liu, T. Sun, Y. Zhou, J. Zhan, and X. Qiu. Data mixing laws: Optimizing data mixtures by predicting language modeling performance. *arXiv preprint arXiv:2403.16952*, 2024. URL <https://arxiv.org/abs/2403.16952>.
- [90] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of ACL*, 2019. URL <https://aclanthology.org/P19-1472/>.
- [91] P. Zhang, G. Zeng, T. Wang, and W. Lu. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*, 2024. URL <https://arxiv.org/abs/2401.02385>.

A Proof of Proposition 2.1

To prove Proposition 2.1, We start from the $K = 1$ case, corresponding to selecting an instance with the maximal reward and Eq. (3) becomes

$$P\left(\arg \max_{1 \leq n \leq N} r(\mathbf{X}_n) = \arg \max_{1 \leq m \leq M} r(\mathbf{Y}_m)\right) \rightarrow 1, \quad (6)$$

which is equivalent to

$$\sum_{\mathbf{x} \in S} P\left(\arg \max_{1 \leq n \leq N} r(\mathbf{X}_n) = \mathbf{x} \text{ and } \arg \max_{1 \leq m \leq M} r(\mathbf{Y}_m) = \mathbf{x}\right) \rightarrow 1. \quad (7)$$

Since \mathbf{X}_n and \mathbf{Y}_m are independent random variables for $1 \leq n \leq N$ and $1 \leq m \leq M$, Eq. (7) can be further written as

$$\sum_{\mathbf{x} \in S} P\left(\arg \max_{1 \leq n \leq N} r(\mathbf{X}_n) = \mathbf{x}\right) P\left(\arg \max_{1 \leq m \leq M} r(\mathbf{Y}_m) = \mathbf{x}\right) \rightarrow 1. \quad (8)$$

We focus on the term $P\left(\arg \max_{1 \leq n \leq N} r(\mathbf{X}_n) = \mathbf{x}\right)$, which can be expanded as follows based on the fact that $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ are i.i.d random variables:

$$\begin{aligned} & P\left(\arg \max_{1 \leq n \leq N} r(\mathbf{X}_n) = \mathbf{x}\right) \\ &= P(r(\mathbf{X}_n) \leq r(\mathbf{x}), \text{ for } 1 \leq n \leq N \text{ and at least one variable in } \mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N \text{ equals } \mathbf{x}) \\ &= P(r(\mathbf{X}_n) \leq r(\mathbf{x}), \text{ for } 1 \leq n \leq N) - P(r(\mathbf{X}_n) \leq r(\mathbf{x}) \text{ and } \mathbf{X}_n \neq \mathbf{x}, \text{ for } 1 \leq n \leq N) \\ &= \prod_{n=1}^N P(r(\mathbf{X}_n) \leq r(\mathbf{x})) - \prod_{n=1}^N [P(r(\mathbf{X}_n) \leq r(\mathbf{x})) - P(\mathbf{X}_n = \mathbf{x})] \\ &= P(r(\mathbf{X}_1) \leq r(\mathbf{x}))^N - [P(r(\mathbf{X}_1) \leq r(\mathbf{x})) - p_1(\mathbf{x})]^N \\ &= P_{\mathbf{X}}(\mathbf{x})^N \left[1 - \left[1 - \frac{p_1(\mathbf{x})}{P_{\mathbf{X}}(\mathbf{x})}\right]^N\right], \end{aligned} \quad (9)$$

where $P_{\mathbf{X}}(\mathbf{x}) = P(r(\mathbf{X}_1) \leq r(\mathbf{x}))$. Note that $0 < p_1(\mathbf{x}) = P(\mathbf{X}_1 = \mathbf{x}) \leq P(r(\mathbf{X}_1) \leq r(\mathbf{x})) = P_{\mathbf{X}}(\mathbf{x})$, we have $0 < \frac{p_1(\mathbf{x})}{P_{\mathbf{X}}(\mathbf{x})} \leq 1$, and thus $\lim_{N \rightarrow +\infty} \left[1 - \frac{p_1(\mathbf{x})}{P_{\mathbf{X}}(\mathbf{x})}\right]^N = 0$. Similarly, by setting $P_{\mathbf{Y}}(\mathbf{x}) = P(r(\mathbf{Y}_1) \leq r(\mathbf{x}))$, we have

$$P\left(\arg \max_{1 \leq m \leq M} r(\mathbf{Y}_m) = \mathbf{x}\right) = P_{\mathbf{Y}}(\mathbf{x})^M \left[1 - \left[1 - \frac{p_2(\mathbf{x})}{P_{\mathbf{Y}}(\mathbf{x})}\right]^M\right], \quad (10)$$

and $\lim_{M \rightarrow +\infty} \left[1 - \frac{p_2(\mathbf{x})}{P_{\mathbf{Y}}(\mathbf{x})}\right]^M = 0$. Let $\mathbf{x}^* = \arg \max_{\mathbf{x} \in S} r(\mathbf{x})$. Therefore, $P_{\mathbf{X}}(\mathbf{x}^*) = P_{\mathbf{Y}}(\mathbf{x}^*) = 1$ and $P_{\mathbf{X}}(\mathbf{x}) < 1, P_{\mathbf{Y}}(\mathbf{x}) < 1$ for $\mathbf{x} \neq \mathbf{x}^*$. When $N, M \rightarrow +\infty$, we have $P_{\mathbf{X}}(\mathbf{x}^*)^N = P_{\mathbf{Y}}(\mathbf{x}^*)^M = 1$, and $P_{\mathbf{X}}(\mathbf{x})^N, P_{\mathbf{Y}}(\mathbf{x})^M \rightarrow 0$ for $\mathbf{x} \neq \mathbf{x}^*$. Therefore, we have

$$\begin{aligned} & \lim_{N, M \rightarrow +\infty} \sum_{\mathbf{x} \in S} P\left(\arg \max_{1 \leq n \leq N} r(\mathbf{X}_n) = \mathbf{x}\right) P\left(\arg \max_{1 \leq m \leq M} r(\mathbf{Y}_m) = \mathbf{x}\right) \\ &= \sum_{\mathbf{x} \in S} \lim_{N, M \rightarrow +\infty} P_{\mathbf{X}}(\mathbf{x})^N P_{\mathbf{Y}}(\mathbf{x})^M \lim_{N, M \rightarrow +\infty} \left[1 - \left[1 - \frac{p_1(\mathbf{x})}{P_{\mathbf{X}}(\mathbf{x})}\right]^N\right] \left[1 - \left[1 - \frac{p_2(\mathbf{x})}{P_{\mathbf{Y}}(\mathbf{x})}\right]^M\right] \\ &= \lim_{N, M \rightarrow +\infty} P_{\mathbf{X}}(\mathbf{x}^*)^N P_{\mathbf{Y}}(\mathbf{x}^*)^M + \sum_{\mathbf{x} \in S, \mathbf{x} \neq \mathbf{x}^*} \lim_{N, M \rightarrow +\infty} P_{\mathbf{X}}(\mathbf{x})^N P_{\mathbf{Y}}(\mathbf{x})^M \\ &= 1, \end{aligned} \quad (11)$$

Model Size	d_{model}	d_{FFN}	n_{layers}	n_{head}	d_{head}	learning rate
104M	512	1,408	8	8	64	6×10^{-4}
200M	768	2,112	12	12	64	6×10^{-4}
300M	768	2,112	18	12	64	6×10^{-4}
500M	1,024	2,816	24	16	64	3×10^{-4}
1.2B	1,536	4,224	24	16	96	2.5×10^{-4}

Table 6: Model configurations and corresponding learning rates.

which proves Eq. (6). For $K > 1$, the equality of two top- K subsets requires the elements ranked from 1 to K to be equal, respectively. Therefore, the left hand of Eq. (3) can be decomposed using Bayes’s Law. Let $\mathcal{X} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$, $\mathcal{Y} = \{\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_M\}$, $\mathbf{X}^* = \arg \max_{1 \leq n \leq N} r(\mathbf{X}_n)$, and $\mathbf{Y}^* = \arg \max_{1 \leq m \leq M} r(\mathbf{Y}_m)$, we have $P(\mathbf{X}^* = \mathbf{Y}^*) \rightarrow 1$ (Eq. (6)) and

$$P(\text{top-2}\{r(\mathbf{X}_n) | 1 \leq n \leq N\} = \text{top-2}\{r(\mathbf{Y}_m) | 1 \leq m \leq M\}) \\ = P\left(\arg \max_{\mathbf{X} \in \mathcal{X} - \{\mathbf{X}^*\}} r(\mathbf{X}) = \arg \max_{\mathbf{Y} \in \mathcal{Y} - \{\mathbf{Y}^*\}} r(\mathbf{Y}) \mid \mathbf{X}^* = \mathbf{Y}^*\right) P(\mathbf{X}^* = \mathbf{Y}^*). \quad (12)$$

Since $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$ and $\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_M$ are i.i.d. random variables, we can still decompose the term $P\left(\arg \max_{\mathbf{X} \in \mathcal{X} - \{\mathbf{X}^*\}} r(\mathbf{X}) = \mathbf{x} \mid \mathbf{X}^* = \mathbf{Y}^*\right)$ as Eq. (9), which means when $N, M \rightarrow +\infty$:

$$P\left(\arg \max_{\mathbf{X} \in \mathcal{X} - \{\mathbf{X}^*\}} r(\mathbf{X}) = \arg \max_{\mathbf{Y} \in \mathcal{Y} - \{\mathbf{Y}^*\}} r(\mathbf{Y}) \mid \mathbf{X}^* = \mathbf{Y}^*\right) \rightarrow 1. \quad (13)$$

Eq. (13) means the elements with the secondary large $r(\mathbf{x})$ values are highly likely to be the same. The decomposition in Eq. (12) can be conducted K times for the elements ranked from 1 to K , and each decomposed term approaches 1 similar to Eq. (13). So far, we have proved that the probability of \mathcal{X} and \mathcal{Y} having the same top- K subsets measured by $r(\mathbf{x})$ approaches to 1 when $N, M \rightarrow +\infty$, which is formally written as

$$P(\text{top-}K\{r(\mathbf{X}_n) | 1 \leq n \leq N\} = \text{top-}K\{r(\mathbf{Y}_m) | 1 \leq m \leq M\}) \rightarrow 1. \quad (14)$$

This completes the proof of Proposition 2.1.

B More Experimental Details

Model and Training Configurations. We mostly follow Brown et al. [10] to set the model and learning rate configurations, as summarized in Table 6. The 500M, 1.8B, and 4B teacher models are the officially released Qwen-1.5 checkpoints³ and the 300M teacher model used in Section 3.5 to analysis the effect of teacher LM sizes. It is pre-trained on 200B tokens from our pre-training corpus. We train all the LMs with the AdamW [50] optimizer, with $\beta_1 = 0.9$, $\beta_2 = 0.98$, and a 0.1 weight decay. We set the batch size to 512 and the max sequence length to 1,024, corresponding to 100K total training steps for roughly 50B tokens in Pre-Train w/o KD, SeqKD, and MiniPLM. For MiniLLM and Vanilla KD, we limit the training steps to align their training computation with Pre-Train w/o KD. Specifically, we assume the computation of a forward and backward pass are $2ND$ and $4ND$, respectively, where N is the model size and D is the number of trained tokens. The training steps of MiniLLM and Vanilla KD are listed in Table 7. We linearly warm up the learning rate for 2K steps and apply cosine learning rate decay until 1/10 of the max values. All experiments are conducted on NVIDIA 40G A100 and NVIDIA 32G V100 GPUs.

Evaluation Details. Our downstream datasets for evaluation include Hellaswag (HS; 90), LAMBADA (LAM; 61), Winograde (Wino; 43), OpenbookQA (OBQA; 53), ARC-Easy/Challenge (ARC-e/c; 16), PIQA [6], SIQA [66], and StoryCloze (Story; 56). We apply the LM-Eval-Harness [24]⁴ framework to conduct zero-shot evaluation. We sample 10K documents from the DCLM [45] corpus to construct our test set for language modeling evaluation.

³<https://huggingface.co/Qwen>

⁴

	Vanilla KD			MiniLLM		
Formula	$\frac{3N_{\text{stu}}}{3N_{\text{stu}}+N_{\text{tch}}}T$			$\frac{3N_{\text{stu}}}{4N_{\text{stu}}+2N_{\text{tch}}}T$		
Student Model Size N_{stu}	200M	500M	1.2B	200M	500M	1.2B
Training Steps	25K	45K	65K	15K	30K	40K

Table 7: Training steps in Vanilla KD and MiniLLM, which is set to ensure training-time computation to be the same as Pre-Train w/o KD. N_{stu} and N_{tch} are model sizes of student and teacher LMs respectively. $N_{\text{tch}} = 1.8B$ in our experiments. $T = 100K$ is the training steps in Pre-Train w/o KD.

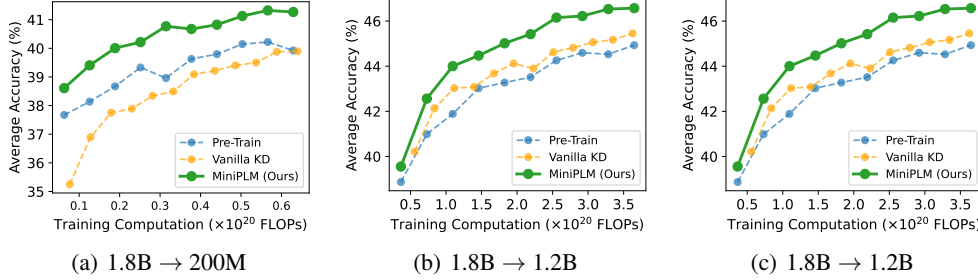


Figure 7: Computation scaling curves of student LMs when trained with Pre-Train w/o KD, Vanilla KD, and MINIPLM. The y-axis represents the average zero-shot accuracy on downstream tasks.

C More Results

C.1 Computation Scaling Curves of More Sizes

In Figure 7, we plot the scaling curves of average accuracy on the downstream tasks with respect to the pre-training FLOPs for student LMs with 200M and 1.2B parameters. We can see that MINIPLM saves pre-training computation for both student LM sizes and constantly outperforms Vanilla KD.

C.2 Test Losses Extrapolation with Scaling Laws

Data-Unlimited Setting. In Table 2, we follow Hoffmann et al. [37] to fit the scaling law curves with the test losses on the DCLM corpus. Then, we used the fitted constants to predict the test losses for 1T and 10T training data in Pre-Train w/o KD and MINIPLM, and that for Vanilla KD when it consumes the same training FLOPs. Specifically, we consider the following power law suggested by Hoffmann et al. [37]:

$$L(N_m, N_d) = L_{\text{irr}} + \frac{A_m}{N_m^{\alpha_m}} + \frac{A_d}{N_d^{\alpha_d}} \quad (15)$$

where N_m is the model sizes, N_d is the number of trained tokens, and L_{irr} is the irreducible test loss. Since the required computation C of training on a token of Vanilla KD and other methods are different, we re-write Eq. (15) using the fact that $C \propto N_m N_d$ for a fixed model size:

$$L(C) = L_{\infty} + \frac{A_c}{C^{\alpha_c}}, \quad (16)$$

where L_{∞} and A_c depends on the model size N_m , and $\alpha_c = \alpha_d$. In this way, we can fit the loss curves in Figure 4 with Eq. (16). Table 8 includes the fitted values of L_{∞} , A_c , and α_c . We also include the total FLOPs of Pre-Train w/o KD and MINIPLM on 1T (C_{1T}) and 10T (C_{10T}) tokens, which Vanilla KD aligns with. The numbers in Table 8 can be computed with Eq. (16) and the constants in Table 8.

Data-Limited Setting. In Section 3.4, we extrapolate the loss curve of Pre-Train w/o KD with the Data-Constrained Scaling Law [57], which is almost the same as Eq. (15), except that N_d represents the total number of tokens in the pre-training corpus repeated 4 times, as suggested by Muennighoff

N_{stu}	Method	A_c	α_c	L_∞	C_{1T} (FLOPs)	C_{10T} (FLOPs)
200M	Pre-Train w/o KD	2.19×10^7	0.41	3.30		
	Vanilla KD	9.77×10^7	0.44	3.34	1.26×10^{21}	1.26×10^{22}
	MINIPLM	8.56×10^{10}	0.59	3.25		
500M	Pre-Train w/o KD	2.73×10^8	0.45	3.06		
	Vanilla KD	3.14×10^8	0.45	3.05	3.14×10^{21}	3.14×10^{21}
	MINIPLM	6.64×10^9	0.52	3.03		
1.2B	Pre-Train w/o KD	1.88×10^8	0.43	2.91		
	Vanilla KD	1.10×10^{10}	0.52	2.90	7.30×10^{21}	7.30×10^{21}
	MINIPLM	4.29×10^8	0.45	2.86		

Table 8: Scaling Law constants in Eq. (16) fitted using the loss curves in Figure 4. N_{stu} means the student LM size. C_{1T} and C_{10T} are the compute spent on processing 1T and 10T tokens in Pre-Train w/o KD and MINIPLM, which Vanilla KD aligns with.

et al. [57]. Therefore, after N_d is solved by letting the loss of Pre-Train w/o KD equal that of MINIPLM ($\alpha = 0.25$, 4 Eps.), we divide the value by 4, resulting in a training tokens number of 118B. This means 68B extra training tokens are needed for Pre-Train w/o KD to achieve the performance of MINIPLM using 50B training tokens.

C.3 Difference Sampling with a Proxy Model

Motivation. In this section, we show that the offline computational overhead of *Difference Sampling* can be further reduced by conducting teacher and reference LM inference on a small proxy dataset and then transferring the value $r(p, p_{\text{ref}}, \mathbf{x})$ to the entire corpus \mathcal{D} with a proxy model. Specifically, we first uniformly sample a proxy subset \mathcal{D}_{prx} from \mathcal{D} , satisfying $|\mathcal{D}_{\text{prx}}| \ll |\mathcal{D}|$. Then, we compute the $r(p, p_{\text{ref}}, \mathbf{x})$ value for each instance \mathbf{x} in \mathcal{D}_{prx} . Note that since $|\mathcal{D}_{\text{prx}}| \ll |\mathcal{D}|$, the computational overhead of teacher LM inference is significantly reduced. After that, we fine-tune a small proxy model on \mathcal{D}_{prx} to fit the $r(p, p_{\text{ref}}, \mathbf{x})$ values, which is used to infer the values for instances from \mathcal{D} . Finally, the Top- K operation in Eq. (4) is based on the inferred values.

Method. To test this approach, we uniformly sample a \mathcal{D}_{prx} containing 0.1B tokens from the 50B-token \mathcal{D} . Computing $\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})}$ values on \mathcal{D}_{prx} takes only 0.2% computation of that on \mathcal{D} . Then, we employ the reference LM as the proxy model and fine-tune it with the following regression loss:

$$\mathbf{w}^*, b^*, \boldsymbol{\theta}_{\text{ref}}^* = \arg \min_{\mathbf{w}, b, \boldsymbol{\theta}_{\text{ref}}} \frac{1}{|\mathcal{D}_{\text{prx}}|} \sum_{\mathbf{x} \in \mathcal{D}_{\text{prx}}} [\mathbf{w}^\top \bar{\mathbf{h}}(\mathbf{x}, \boldsymbol{\theta}_{\text{ref}}) + b - r(p, p_{\text{ref}}, \mathbf{x})]^2, \quad (17)$$

where $\bar{\mathbf{h}}(\mathbf{x}, \boldsymbol{\theta}_{\text{ref}}) \in \mathbb{R}^d$ is the average output hidden states of the reference LM, with d representing the model’s hidden size and $\boldsymbol{\theta}_{\text{ref}}$ representing the parameters of the reference LM. $\mathbf{w} \in \mathbb{R}^d$, $b \in \mathbb{R}$ are the parameters of a linear head outputting a scalar as the predicted $r(p, p_{\text{ref}}, \mathbf{x})$ values, given the average hidden states. The inferred values on \mathcal{D} are given by $\hat{r}(\mathbf{x}) = \mathbf{w}^{*\top} \bar{\mathbf{h}}(\mathbf{x}, \boldsymbol{\theta}_{\text{ref}}^*) + b^*$. Since this inference process is based on the reference LM, it still saves computation compared to inference with the teacher LM. The difference-sampled corpus is then obtained by selecting $K = \alpha|\mathcal{D}|$ instances from \mathcal{D} with the highest $\hat{r}(\mathbf{x})$ values.

Method	FLOPs	Acc.
Vanilla KD	Online	39.9
MINIPLM	2×10^{20}	41.3
MINIPLM _{prx}	9×10^{18}	40.9

Table 9: Offline FLOPs and average accuracy (Acc.) on downstream tasks of MINIPLM using a proxy model, compared with that of standard MINIPLM and Vanilla KD.

Results. We term this method as MINIPLM_{prx} and compare its performance with Vanilla KD and MINIPLM in Table 9. We can see that MINIPLM_{prx} requires much less offline inference computation compared to standard MINIPLM while still maintaining substantial improvement over Vanilla KD.

C.4 Ablation Study

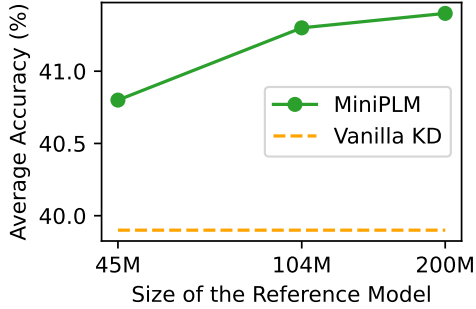


Figure 8: Impact of the reference model size. We use the 1.8B LM as the teacher and the 200M LM as the student. We report the average zero-shot accuracy on the downstream tasks of the LMs trained with MINIPLM and compare it with Vanilla KD.

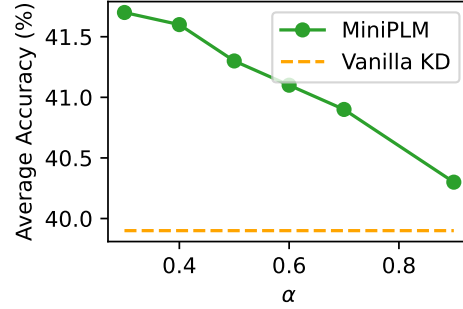


Figure 9: Impact of the difference sampling ratio α . We report the average zero-shot accuracy on the downstream tasks of the LMs trained with MINIPLM, using $\alpha \in [0.3, 0.4, 0.5, 0.6, 0.7, 0.9]$ and compare it with Vanilla KD.

Impact of the Reference Model. In Figure 8, we show the impact of the reference model size on the performance of LMs trained with MINIPLM. We can see that larger reference models lead to better performance of MINIPLM, but the improvement saturates as the reference LM grows larger.

Difference Sampling Ratio. In Figure 9, we plot the impact of the difference sampling ratio on the student LM’s performance in the data-unlimited setting. We can see that small sampling ratios result in better zero-shot accuracy on downstream tasks. However, to ensure that the difference-sampled data contains 50B tokens, we need a larger original corpus \mathcal{D} . To control the computation overhead, we mainly use $\alpha = 0.5$ in our experiments.

C.5 Case Study

We present a case study in Table 10, 11, and 12 to show the instances corresponding to the three parts of data distribution shown in Figure 3(b).

$p(\mathbf{x}) \gtrsim p_{\text{ref}}(\mathbf{x})$: Easy and common instances				
Instance #1	$-\log p(\mathbf{x}) = 1.24$	$-\log p_{\text{ref}}(\mathbf{x}) = 1.28$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = 0.04$	Discarded
	<pre> <node id='-659' action='modify' visible='true' lat= '0.05467069248579575' lon='0.014892969670168166' /> <node id='-657' action='modify' visible='true' lat= '0.05243834625645345' lon='0.023237696125627347' /> <node id='-655' action='modify' visible='true' lat= '0.04940873338995851' lon='0.03152927145717611' /> <node id='-653' action='modify' visible='true' lat= '0.04786735135170292' lon='0.040405509151806455' /> <node id='-651' action='modify' visible='true' lat= '0.04733584029593642' lon='0.04513595918070425' /> </pre>			
Instance #2	$-\log p(\mathbf{x}) = 0.44$	$-\log p_{\text{ref}}(\mathbf{x}) = 0.51$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = 0.07$	Discarded
	<pre> {\sum\limits_{j = 1}^{\hat{n} - 1}\operatornamename{size}\left({\mathcal{X},j} \right) \cdot \operatornamename{size} \left({\mathcal{Z},n - j} \right)\quad} \& {\mathcal{I} = \mathcal{M}_{\{1\}},} \& \& {\sum\limits_{j = 1}^{\hat{n} - 1}\operatornamename{size}\left({\mathcal{X},j} \right) \cdot \operatornamename{size} \left({\mathcal{X},n - j} \right)\quad} \& {\mathcal{I} = \mathcal{M}_{\{2\}},} \& \& {\sum\limits_{j = 1}^{\hat{n} - 1}\operatornamename{size}\left({\mathcal{X},j} \right) \cdot \operatornamename{size} \left({\mathcal{A},n - j} \right)\quad} \& {\mathcal{I} = \mathcal{M}_{\{3\}},} \& \& {\sum\limits_{j = 1}^{\hat{n} - 1}\operatornamename{size}\left({\mathcal{A},j} \right) \cdot \operatornamename{size} \left({\mathcal{Z},n - j} \right)\quad} \& {\mathcal{I} = \mathcal{M}_{\{4\}},} \& \& </pre>			
Instance #3	$-\log p(\mathbf{x}) = 2.83$	$-\log p_{\text{ref}}(\mathbf{x}) = 3.86$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = 1.02$	Selected
	<p>"I felt bad for Coach Rod to have to deal with it," said senior tight end Mike Massey, a St. Ignatius grad. "But for the players it was a non-issue. We didn't even talk about it."</p> <p>But some other people in college football did.</p> <p>"[Boren] was a legacy guy and when he makes comments like [that], that's like getting kicked square in the shorts when you're Rich Rodriguez," ESPN analyst and former OSU quarterback Kirk Herbstreit said Thursday. "With all the other things that are happening, when you get that comment from a legacy guy, a guy whose dad started three or four years for Bo, all of a sudden everybody's tentacles go up a little bit."</p>			

Table 10: Easy and common instances down-sampled by *Difference Sampling*. Instance #1 and #2: HTML and LaTeX code data that contain repeated patterns, which is easy to fit by both the reference and teacher LM. Instance #3: Dialogues from a story, which is relative easy for LMs but are still selected by *Difference Sampling* to help student LMs learn basic language skills.

$p(\mathbf{x}) \gg p_{\text{ref}}(\mathbf{x})$: Hard and valuable instances				
Instance #1	$-\log p(\mathbf{x}) = 1.26$	$-\log p_{\text{ref}}(\mathbf{x}) = 4.20$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = 2.94$	Selected
	Legal along with Environmental Responsibility! Dumpster rentals in the user side may seem as fundamental as placing a phone, having a dumpster sent and hurling all your disposals inside to be carted away. Nonetheless, there are legal issues attached to appropriate disposal connected with certain products which tie up into environmental issues. The 10 Yard Dumpster For Rent in Pocahontas customer or perhaps demolition purchaser should be informed about these issues by means of careful screening so as to reduce a firm's liability which inturn keeps a firm's overhead all the way down and makes for prompt fall off, pick up along with disposal of the dumpster and it's articles.			
Instance #2	$-\log p(\mathbf{x}) = 2.36$	$-\log p_{\text{ref}}(\mathbf{x}) = 5.59$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = 3.23$	Selected
	有利 you3li4 yǒulì advantageous; beneficial 谨慎 jin3shen4 jǐnshèn cautious; prudent 甲 jia3 jiǎ one; armor (1st Heavenly Stem) 犹豫 you2yu4 yóuyù hesitate; hesitant; undecided 从此 cong2ci3 cóngcǐ from now on; since then 企业 qi3ye4 qǐyè company; business; firm 下载 xia4zai3 xiàzǎi to download 狮子 shi1zi5 shīzi lion 青少年 qing1shao4nian2 qīngshàonián teenager			
Instance #3	$-\log p(\mathbf{x}) = 0.16$	$-\log p_{\text{ref}}(\mathbf{x}) = 2.73$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = 2.56$	Selected
	<pre>function WritableState(options, stream) { var Duplex = require('./_stream_duplex'); options = options {}; // the point at which write() starts returning false // Note: 0 is a valid value, means that we always return false if // the entire buffer is not flushed immediately on write() var hwm = options.highWaterMark; var defaultHwm = options.objectMode?16:16*1024; this.highWaterMark = (hwm hwm === 0) ? hwm : defaultHwm; // object stream flag to indicate whether or not this stream // contains buffers or objects. this.objectMode = !!options.objectMode; ... }</pre>			

Table 11: Hard and valuable instances up-sampled by *Difference Sampling*. Instance #1: High-quality long documents containing versatile world knowledge. Instance #2: instance that contains translation tasks, presented in an in-context learning form. Instance #3: High-quality code data with detailed comments.

$p(\mathbf{x}) \lesssim p_{\text{ref}}(\mathbf{x})$: Noisy and harmful instances				
Instance #1	$-\log p(\mathbf{x}) = 9.50$	$-\log p_{\text{ref}}(\mathbf{x}) = 6.60$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = -2.90$	Discarded
] {} *****			
Instance #2	$-\log p(\mathbf{x}) = 1.01$	$-\log p_{\text{ref}}(\mathbf{x}) = 0.90$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = -0.11$	Discarded
	(91.00,60.00) (1.00,35.00) (2.00,35.00) [(1,0) [13.00] {}] {} (16.00,35.00) (16.70,34.30) [(1,-1) [8.60] {}] {} (26.00,25.00) (26.00,11.00) (26.00,12.00) [(0,1) [12.00] {}] {} (36.00,35.00) (26.00,45.00) (16.70,35.70) (35.30,34.30) [(-1,-1) [8.60] {}] {} (56.00,35.00) (66.00,25.00) (76.00,35.00) (66.00,45.00) (66.00,46.00)			
Instance #3	$-\log p(\mathbf{x}) = 2.53$	$-\log p_{\text{ref}}(\mathbf{x}) = 0.26$	$\log \frac{p(\mathbf{x})}{p_{\text{ref}}(\mathbf{x})} = -2.26$	Discarded
	Z = 8 ----- Data collection #tablewrapdatacollectionlong ===== ----- Bruker APEXII area-detector diffractometer Radiation source: fine-focus sealed tube graphite p and w scans Absorption correction: multi-scan (*SADABS*; Sheldrick, 2004) *T* _{min} = 0.970, *T* _{max} = 0.981 20408 measured reflections -----			

Table 12: Noisy and harmful instances discarded by *Difference Sampling*. Instance #1: irregular symbols. Both the reference LM and the teacher LM are hard to fit. Instance #2: a string primarily made of meaningless numbers. Both the reference LM and the teacher LM are easy to fit, by predicting numbers. Instance #3: data collection in scientific paper, but lack of contexts. The reference LM fits the patterns, which the teacher LM finds useless.