

PretrainZero: Reinforcement Active Pretraining

Xingrun Xing^{1,2}, Zhiyuan Fan², Jie Lou², Guoqi Li¹, Jiajun Zhang¹, Debing Zhang²

¹ Institute of Automation, Chinese Academy of Sciences

² Xiaohongshu Inc.

loujie0822@aliyun.com, jjzhang@nlpr.ia.ac.cn

Abstract

Mimicking human behavior to actively learning from general experience and achieve artificial general intelligence has always been a human dream. Recent reinforcement learning (RL) based large-thinking models demonstrate impressive expert-level abilities, i.e., software and math, but still rely heavily on verifiable rewards in specific domains, placing a significant bottleneck to extend the performance boundary of general reasoning capabilities. In this work, we propose PretrainZero, a reinforcement active learning framework built on the pretraining corpus to extend RL from domain-specific post-training to general pretraining. PretrainZero features the following characteristics: 1) **Active pretraining**: inspired by the active learning ability of humans, PretrainZero learns a unified reasoning policy to actively identify reasonable and informative contents from pretraining corpus, and reason to predict these contents by RL. 2) **Self-supervised learning**: without any verifiable labels, pretrained reward models, or supervised fine-tuning, we directly pretrain reasoners from 3 ~ 30B base models on the general Wikipedia corpus using RL, significantly breaking the verification data-wall for general reasoning. 3) **Verification scaling**: by tackling increasingly challenging masked spans, PretrainZero substantially enhances the general reasoning abilities of pretrained base models. In reinforcement pretraining, PretrainZero improves Qwen3-4B-Base for 8.43, 5.96 and 10.60 on MMLU-Pro, SuperGPQA and math average benchmarks. In post-training, the pretrained models can also serve as reasoning foundation models for downstream RLVR tasks.

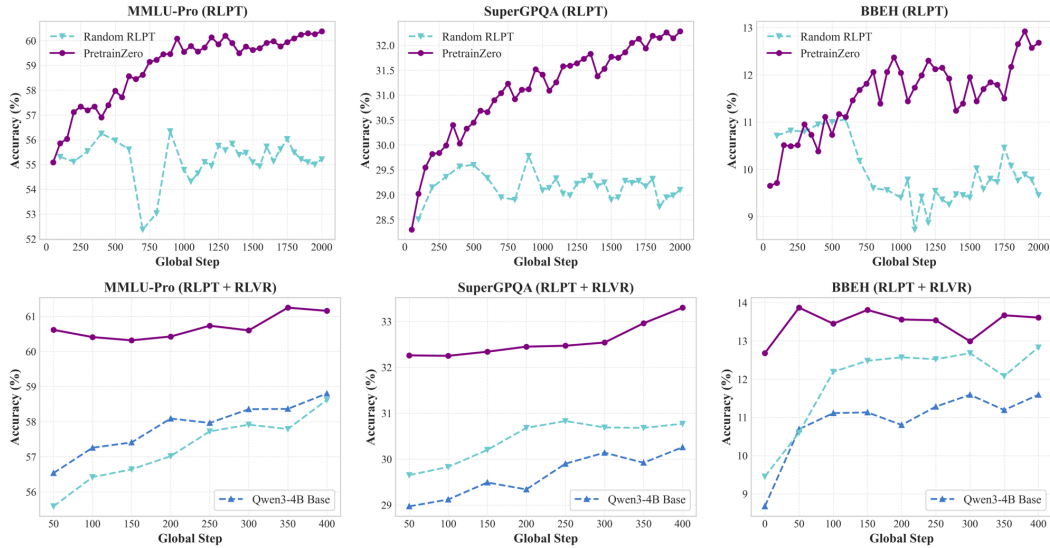


Figure 1: Reinforcement Pre-Training (RLPT) performance in pre-training and post-training stages.

1 Introduction

Recent large language models (LLMs) have achieved human-level expertise in specific domains, particularly through large-scale self-supervised learning in pretraining [KMH⁺20, AAA⁺23] and Reinforcement Learning (RL) [GYZ⁺25, YZZ⁺25, CZY⁺25] in post-training. During pretraining, self-supervised learning with a fixed next-token prediction paradigm allows models to leverage large-scale, low-cost data to improve general capabilities effectively. In contrast, the post-training RL faces a severe data-wall: Reinforcement Learning with Verifiable Rewards (RLVR) [GYZ⁺25, YCL⁺] requires domain-specific verifiers to label training samples, and Reinforcement Learning from Human Feedback (RLHF) [OWJ⁺22, BJN⁺22], relying on reward models and humans, can only train limited steps to avoid reward hacking. This motivates a natural direction—performing reinforcement learning [DDT⁺25, LLX⁺25] in a self-supervised pretraining manner [BMR⁺20], in order to use inexpensive pretraining data to extend RLVR and overcome this data-wall.

However, formulating the self-supervised pretraining as RLVR tasks is non-trivial. Towards this goal, this work first investigates stand-alone Reinforcement Learning Pre-Training (RLPT) [DDT⁺25] according to three principles: 1) Comprehensiveness: we establish both baselines including masked token prediction and next token prediction as the reasoning objective [DDT⁺25]. 2) Full self-supervision: we exclude any additional Supervised Fine-Tuning (SFT) cold-start and reward models [LLX⁺25]. 3) Generalization: we avoid Question–Answer formats or synthetic chain-of-thought (CoT) datasets, like OmniMath [GSY⁺24], and train directly on general-domain Wikipedia [DCLT19]. Experimental results demonstrate that the vanilla RLPT fails to emerge high-quality CoT: the low information density of pretraining corpus leads to inefficient learning, and the presence of noisy or incorrect tokens often causes training collapse.

This work proposes the first stand-alone RLPT method to extend RLVR on real-world pretraining corpus, termed PretrainZero. This is achieved by mimicking the human active-learning behavior [YLB⁺25, Set09]: humans can actively learn from a broad range of experiences, selectively focusing on informative elements and unfamiliar concepts. This allows effective learning even when the underlying experiences are noisy or low in information density. In contrast, current large language models—whether through supervised or reinforcement pretraining—rely on fixed prediction patterns, such as next-token or off-policy selected masked-token prediction. These rigid learning patterns limit their efficiency [BTK⁺] and prevent them from leveraging data as flexibly as humans do.

Inspired by this, this work proposes a reinforcement active learning framework in order to learn from real-world pretraining distributions. Specifically, we introduce an on-policy *mask generation task* as an auxiliary active-learning objective for the *masked-span prediction task*. The mask generation policy actively explores informative, verifiable, and not-yet-mastered content within the pretraining data and selects it as learning targets. Consequently, the masked-span prediction policy learns to produce CoT reasoning and recover these informative spans. For optimization, we formulate a min–max bilevel reinforcement learning objective, where each batch is jointly optimized using GRPO [SWZ⁺24] over both the mask generation and masked-span prediction tasks. Different from unsupervised RL methods such as self-play [HYW⁺25] and test-time scaling [SLXK24], PretrainZero provides a **verifiable RL scaling mechanism grounded in real data** in a self-supervising manner. This avoids the severe hallucination issues in self-play and test-time scaling, where majority voting from model-generated answers serves as supervision and ultimately leads to collapse in prolonged RL training [LDL⁺25].

As shown in Fig. 1, we evaluate RLPT for 2000 steps in pretraining and add general RLVR in post-training. For a Qwen3-4B-Base [YLY⁺25] model, PretrainZero consistently improves 8.43, 5.96, and 10.60 on MMLU-Pro [WMZ⁺24], SuperGPQA [DYM⁺25], and math average benchmarks during reinforcement pretraining. After general RLVR [MLJ⁺25] in the post-training stage, these improvements remain substantial, with final improvements of 2.35, 3.04, and 2.81 on MMLU-Pro, SuperGPQA, and math average respectively.

Our contributions are summarized as follows:

- We introduce PretrainZero, the first stand-alone RLPT method to operate RLVR on real-world pretraining corpus, enabling general-domain and large-scale reinforcement learning trained directly from base models using only pretraining data as grounding.
- We propose the reinforcement active pretraining mechanism inspired by human active learning. The introduced mask-generation objective enables the model to anticipate what information should be learned actively, ensuring effective training under low–information-density pretraining corpus.

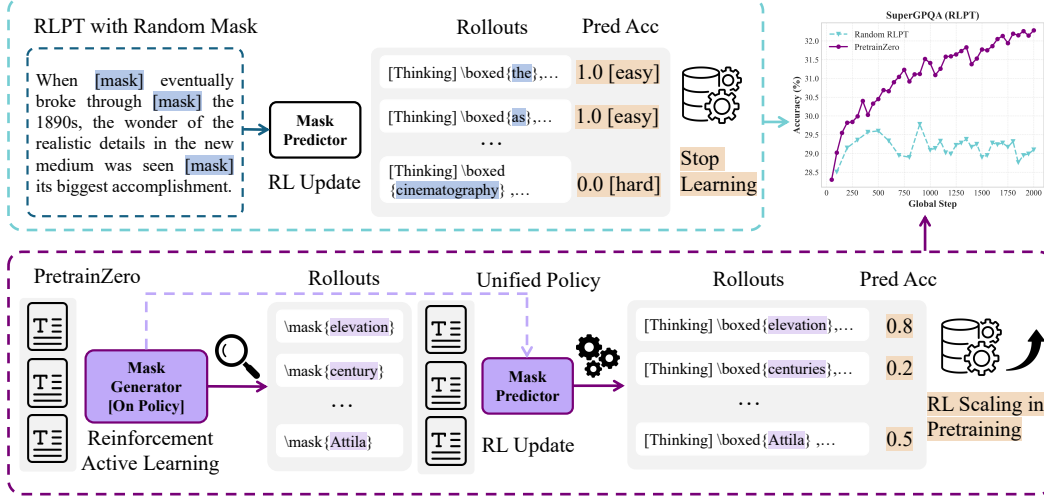


Figure 2: An overview of Reinforcement Active Pretraining. Compared with vanilla RLPT, PretrainZero actively explores and learns from the informative contexts on the pretraining corpus.

- We evaluate PretrainZero in both the pretraining and post-training stages, showing that it effectively mitigates the general reasoning data-wall with pretraining data, and finally the pretrained reasoning models can serve as the reasoning foundation models for general downstream RLVR tasks.

2 Preliminary

In order to learn from the pretraining corpus, traditional self-supervised pretraining adopts language modeling objectives to capture linguistic patterns and contextual dependencies. Recently, the emerging reinforcement pretraining constructs verifiable data through token prediction to learn the reasoning process within concepts. We briefly review different learning patterns in this section.

2.1 Self-Supervised Pretraining

Given the context, traditional self-supervised pretraining tasks include masked token prediction (MTP) [RWC⁺19] and next token prediction (NTP) [DCLT19]. As shown in Eq. 1, the NTP task predict the identity tokens x_t in each location given their preceding context $x_{<t}$ under an auto-regressive pattern:

$$\mathcal{J}_{\text{NTP}}(\theta) = \sum_{t=1}^T \log \pi_{\theta}(x_t \mid x_{<t}), \quad (1)$$

where x is the token sequence with length T and θ is the pretrained model parameters. As the counterpart, masked token prediction task jointly leverages both the preceding and succeeding contexts $x_{m<t,t>n}$ to predict the masked tokens $x_{m \leq t \leq n}$:

$$\mathcal{J}_{\text{MTP}}(\theta) = \sum_{t=m}^n \log \pi_{\theta}(x_{m \leq t \leq n} \mid x_{m<t,t>n}). \quad (2)$$

Supported by self-supervised pretraining, modern LLMs [Ope24, LFX⁺24] successfully scale up pretraining on massive Internet data. In this work, we simulate both masked token prediction and next token prediction as reinforcement reasoning tasks to explore more general RL approach [MLJ⁺25].

2.2 Reinforcement Pretraining

Recent Reinforcement Pre-Training (RPT) [DDT⁺25] extends reinforcement learning into the pre-training corpus, constructing verifiable training data directly from the pretraining corpus and thereby

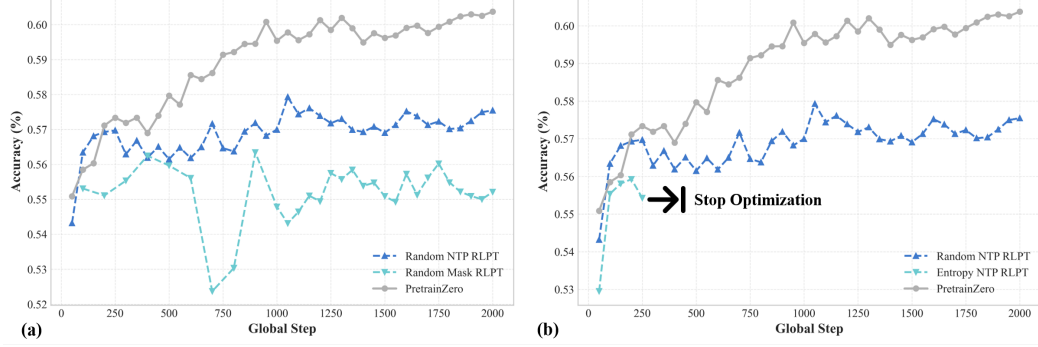


Figure 3: MMLU-Pro performance for foundational RLPT methods. (a) Reinforcement next token prediction and reinforcement masked token prediction. (b) Reinforcement next token prediction with entropy and random token selection.

alleviating the reliance on costly annotations and specific environments for verification. Specifically, RPT introduces the next-token reasoning task: given a sequence x , one token x_t is treated as ground-truth and its preceding tokens $x_{<t}$ as context for the generated output, o_t . Unlike the self-supervised NTP task, where the model directly predicts the next token, RPT [DDT⁺25] first produces a chain-of-thought reasoning process [GYZ⁺25] before generating the final predicted token. In optimization, RPT applies GRPO algorithm with group size G , and uses the exact match verifiable reward r_t^i between prediction \hat{x}_t^i and ground-truth x_t^i :

$$r_t^i = \begin{cases} 1 & \text{if } \hat{x}_t^i(x_{<t}^i) = x_t^i, \\ 0 & \text{otherwise} \end{cases}, \quad (3)$$

$$\mathcal{J}_{\text{RPT}}(\theta) = \mathbb{E}_{(x_{\leq t}) \sim \mathcal{D}, \{o_t^i\}_{i=1}^G \sim \pi_\theta(\cdot | x_{<t})} [r_t^i]. \quad (4)$$

Discussion on the weakness of vanilla reinforcement pretraining. Despite the simplicity of RPT and its potential to extend the RLVR-style method into pretraining, several significant concerns also emerge, making vanilla RPT unsuitable for practical pretraining settings:

- **Robustness on real-world corpus:** although RPT demonstrates improvements on synthetic dataset OmniMath, real-world pretraining data with more noise [DCLT19] often causes training collapse.
- **Training from base models:** vanilla RPT depends on post-training distillation; Other explorations usually rely on SFT cold-start, external reward models, significantly increasing the complexity.
- **Learning effectiveness:** due to the low information density in pretraining corpus, simple token selection methods fails to identify informative content, hindering effective optimization.
- **Training efficiency:** unlike self-supervised NTP that predicts all tokens in parallel, RPT predicts one single token in each sample, yielding limited learning information per step.

3 Reinforcement Active Pretraining

To solve these questions, this work first establishes a Reinforcement Learning Pre-Training (RLPT) baseline on the widely used general domain Wikipedia dataset building upon the Qwen3-4B-Base model. Based on the empirical observations, we then propose a unified and active pretraining task to confirm the general and practical reinforcement pretraining.

3.1 Reinforcement Pretraining Baselines

We establish an RLPT baseline with three masking strategies for training corpus. The model is required to predict masked tokens x_t through CoT reasoning, receiving binary rewards from exact match with ground truth:

$$\mathcal{J}_{\text{RLPT}}(\theta) = \mathbb{E}_{x \sim \mathcal{D}, t, \{o_t^i\}_{i=1}^G \sim \pi_\theta(\cdot | x_{\setminus t})} [\mathbb{I}[\hat{x}_t^i = x_t^i]]. \quad (5)$$

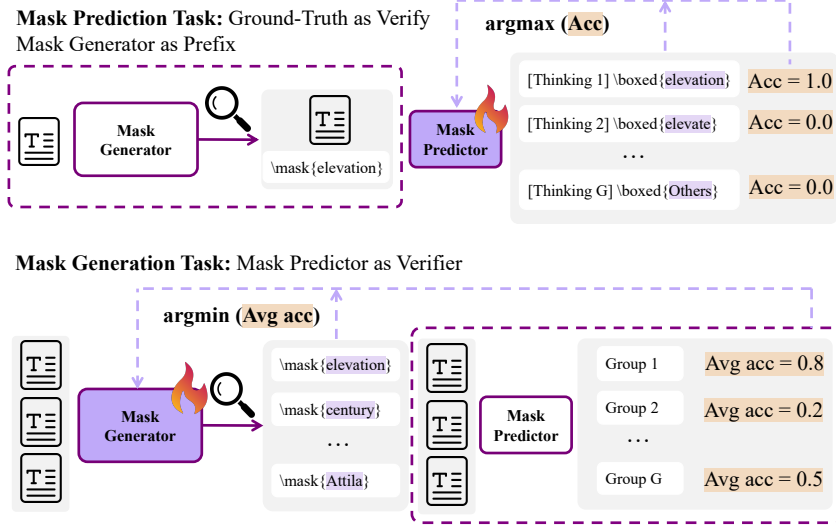


Figure 4: Pretraining Mask Prediction and Mask Generation tasks with GRPO.

Three mask prediction strategies for RLPT are investigated:

- **Random Next Token Reasoning.** The sequence is randomly truncated and the last token before truncation is masked for prediction. For each sample, the model first generates a CoT and only one selected token is predicted according to the generated CoT.
- **Random Masked Span Reasoning.** A word span containing several tokens in the sequence is randomly selected and masked [JCL⁺20], allowing the CoT to predict more than one tokens.
- **Entropy-based Next Token Reasoning.** The token with the top 20% entropy in the sequence is randomly selected and masked, with all subsequent tokens truncated, which consists with RPT.

Empirical Observation. Preliminary experiments are conducted to evaluate these masking strategies on the Wikipedia corpus, with performance measured by MMLU-Pro. As shown in Figure 3 (left), Random NPT RLPT outperforms Random Mask RLPT with more stable training dynamics:

Findings 1. *Although Mask RLPT increases the predicted tokens, the vanilla random word-span selection strategy cannot effectively capture richer semantics in pretraining.*

To further investigate the effect of token selection, Random NPT is compared with Entropy NPT, where the token with higher entropy is selected for masking. As shown in Figure 3 (right), Entropy NPT leads to training collapse and rapid reward degradation. At the position marked *stop optimization*, the reward signal becomes degenerate—all samples within a group yield either 0 or 1 accuracy. The reason is the data quality discrepancy between the synthetic and raw corpus. While entropy-based selection performs well on OmniMath (in RPT), a high-quality synthetic dataset where high-entropy tokens consistently represent challenging but learnable patterns, the same strategy fails on Wikipedia. Raw Wikipedia data contains noise and inconsistencies, causing high-entropy tokens to be either genuinely difficult or simply noisy and unpredictable, which creates unstable learning signals:

Findings 2. *In real-world pretraining data distributions, selecting high-entropy tokens is no more effective than a random selection strategy, and learning actively from noisy data is necessary.*

3.2 Active Pretraining Tasks

The limited performance of these passive masking strategies motivates a more informative and effective learning approach. Consider how humans learn: students focus on informative and valuable content in their experience that maximizes their improvement, rather than randomly selecting practice materials. Inspired by the active learning behavior, we propose an active masking strategy where the model learns to identify beneficial masking positions during training. Rather than relying on fixed heuristics like random sampling or entropy thresholding, the model discovers which tokens provide

Prompt Example: Mask Generation / Prediction

Generate a mask to mask important words in the following paragraph, satisfying the requirements below:

- 1) The mask should mask one or more entities in the passage. The masked words should be continuous.
- 2) The masked words should exactly match words in the original passage.
- 3) The masked words could be predicted according to the context. The difficulty to predict should be moderately challenging for you, so the answer would be short and as unique as possible.

Paragraph: <paragraph>

The final generated masked words must be placed inside \mask{ }.

There is a passage with masked words by [mask]:

<paragraph with mask>

Please reason step by step, and put the predicted masked words within \boxed{ }.

Figure 5: Prompt for Mask Generation and Prediction.

the strongest learning signals. The training process consists of two tasks for the shared LLM, as shown in Fig. 4:

Mask Generation: given a text sequence s from pretraining data \mathcal{D} , the pretraining LLM π_ω first generates a thinking process and then generates a word span, $m \sim \pi_\omega(\cdot | s)$, to mask in this sequence. As shown in Fig. 5, we initially prompt the policy $\pi_\omega(\cdot | s)$ to generate a span mask with one or several words verifiable for reasoning. During pretraining, the policy π_ω continuously learns to explore and capture semantic contents from the noisy pretraining corpus by RL. During the early training stage, the mask prediction policy is relatively weak and requires explicit clues, while in later stages, it needs to focus on the harder and unsolved words and domains.

Mask Prediction: We introduce the masked span prediction as a verifiable reinforcement learning task. Different from next token reasoning, a single CoT process predicts multiple masked tokens in a continuous span, $s_{[p:q]}$. Given the generated mask $m \sim \pi_\omega(\cdot | s)$, we replace the word span $s_{[p:q]}$ with the mark [mask] in the sequence, and then recover the masked content through CoT reasoning. As shown in Fig. 5, we prompt the policy $\psi_\omega(\cdot | m, s)$ to directly generate a CoT at the initial stage before the final mask prediction $\hat{x} \sim \psi_\omega(\cdot | m, s)$. During optimization, the CoT reasons verifiable and semantic targets from the prefixed mask generation task.

3.3 Reinforcement Active Learning

Active learning objective. We cast mask generation and mask prediction as a coupled adversarial process, implemented with a shared LLM parameterized by ω . The generator $\pi_{\omega'}(\cdot | s)$ proposes masking patterns, while the predictor $\psi_\omega(\cdot | m, s)$ seeks to recover the masked content. Based on the final mask prediction rewards $R(s, m, \hat{x})$, this interaction is governed by the objective:

$$J(\omega) := \mathbb{E}_{s \sim \mathcal{D}, m \sim \pi_{\omega'}(\cdot | s)} \left[\mathbb{E}_{\hat{x} \sim \psi_\omega(\cdot | m, s)} [R(s, m, \hat{x})] \right], \quad (6)$$

which evaluates the predictor’s performance under the generator’s masking strategy. To encourage increasingly informative and challenging masks, we define the generator’s objective as $V(\omega) = \min_{\omega'} J(\omega')$, while the predictor optimizes in the opposite direction, i.e., $\arg \max_{\omega} J(\omega)$, thereby forming a coupled min–max formulation:

$$\omega^* = \arg \max_{\omega} V(\omega) = \arg \min_{\omega' \in \Omega} \max_{\omega \in \Omega} \mathbb{E}_{s \sim \mathcal{D}} \left[\mathbb{E}_{m \sim \pi_{\omega'}(\cdot | s), \hat{x} \sim \psi_\omega(\cdot | m, s)} [R(s, m, \hat{x})] \right]. \quad (7)$$

This adversarial min–max structure naturally mirrors the principle of active learning, where the generator actively selects reasonable and informative masks to probe the model’s weaknesses, thereby driving the predictor toward improved robustness and generalization.

Reinforcement optimization. To optimize the min-max active-learning objective in Eq. (7), we implement both the mask prediction and mask generation as RL problems. For the prediction policy $\psi_\omega(\hat{x} \mid m, s)$, the reward is simply defined as an exact match between the predicted token span \hat{x}^i and ground-truth x^i :

$$r_{\text{pred}}^i = R(s, m, \hat{x}) = \mathbb{I}[\hat{x}^i = x^i], \quad (8)$$

which directly optimizes the inner maximization $\mathbb{E}_{\hat{x} \sim \psi_\omega(\cdot \mid m, s)}[R(s, m, \hat{x})]$ in Eq. (6). For the generation policy $\pi_{\omega'}(m \mid s)$, the reward is defined as the negative prediction accuracy under its own masks:

$$r_{\text{gen}}^j = 1 - \mathbb{E}_{x \sim \psi_\omega(\cdot \mid m, s)}[R(s, m, x)] = 1 - \frac{1}{G} \sum_{i=1}^G \mathbb{I}[\hat{x}^{i,j} = x^{i,j}], \quad (9)$$

which aligns with the outer minimization $\mathbb{E}_{s \sim \mathcal{D}, m \sim \pi_{\omega'}(\cdot \mid s)} \left[\mathbb{E}_{\hat{x} \sim \psi_\omega(\cdot \mid m, s)} [R(s, m, \hat{x})] \right]$ in Eq. (6).

The mask generator is rewarded when its masks lead to lower prediction accuracy, indicating that the induced masks contain higher information content for the model. In addition, when the mask prediction accuracy is zero, we further define the generator’s reward to be $r_{\text{gen}}^i = 0$, in order to avoid rewarding the noisy masks that are not predictable for $\psi_\omega(\cdot \mid m, s)$.

Given the reward definitions, r_{pred}^i and r_{gen}^j , we optimize Eq. (7) using GRPO. By substituting r_{gen}^j into the generator’s advantage, we obtain $A_{\text{gen}}^j = -\mathbb{E}[A_{\text{pred}}^{i,j}]$, which proves that the GRPO update is fully consistent with the min-max objective in Eq. (7):

$$\hat{A}_{\text{gen}}^j = \frac{r_{\text{gen}}^j - \text{mean}(r^1, \dots, r^G)}{\text{std}(r^1, \dots, r^G)} = -\mathbb{E}[\hat{A}_{\text{pred}}^{i,j}]. \quad (10)$$

We directly concatenate and uniformly optimize the mask generation and prediction batches in each step:

$$\mathcal{L}_{\text{GRPO}}(\omega) = -\frac{1}{G} \sum_{i=1}^G \min\left(\frac{\pi_\omega(x_i)}{\pi_{\omega_{\text{old}}}(x_i)} \hat{A}_i, \text{clip}\left(\frac{\pi_\omega(x_i)}{\pi_{\omega_{\text{old}}}(x_i)}, 1 - \epsilon, 1 + \epsilon\right) \hat{A}_i\right). \quad (11)$$

4 Experimental Results

4.1 Implementation Details

Model. To evaluate stand-alone reinforcement pretraining, we directly continue pretraining the base models using reinforcement learning without introducing any intermediate supervised finetuning (SFT) cold start. Specifically, we pretrain base models in 3 ~ 30 billion parameters, including the Qwen3-4B-Base, Qwen3-8B-Base, Qwen3-30B-A3B-MoE-Base, and SmoLLM3-3B-Base.

Dataset. To evaluate on real-world distributed pretraining corpus, we use only the most general Wikipedia dataset. Notice that existing RLPT often includes explicit Question-Answer pairs or synthetic datasets such as OmniMath that contain strong reasoning CoTs; this risks allowing the RL objective to copy these reasoning traces directly, implicitly degrading to supervised learning.

Training. For RLPT, we train 2000 steps using GRPO without KL regularization [SWD⁺17]. Following DAPO [YZZ⁺25], we filter samples whose accuracies are exactly 0.0 or 1.0, and we adopt the clip-higher strategy for stability. For Qwen-Base models, we directly perform the PretrainZero strategy; for SmoLLM3-3B-Base, we first use random RLPT for 100 steps as RL cold-start, and then perform PretrainZero for the remaining 1900 steps. During reasoning, the max length of the prompt and response is limited to 1536 and 4096 tokens, respectively. We adopt the 5×10^{-7} learning rate and the cosine scheme. In the mask-generation task, each batch contains 32 pretraining paragraphs, with 8 rollouts for each to produce masks. In the mask-prediction task, we evaluate 256 masks from the prefixed mask generation task (32×8), and each mask is also paired with 8 rollouts for prediction. Consequently, the overall prompt batch size becomes 288 for RL ($32 + 32 \times 8$).

Evaluation. We evaluate on both general-domain and math-domain reasoning benchmarks. For general domain reasoning, we evaluate on the MMLU-Pro [WMZ⁺24], SuperGPQA [DYM⁺25], and BBEH [KFB⁺25]; for math domain reasoning, we evaluate on 6 widely used benchmarks, including Math 500 [HBK⁺21], Olympiad [HLB⁺24], Minerva [LAD⁺22], GSM8K [CKB⁺21], AMC23, and AIME24.

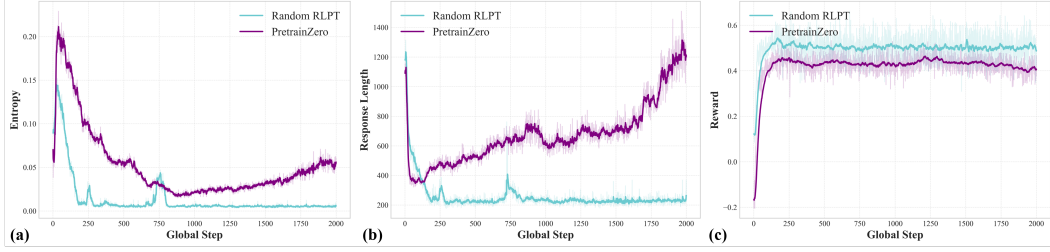


Figure 6: Training dynamic comparisons between PretrainZero and Random RLPT on Qwen3-4B-Base: (a) entropy of model outputs; (b) response length of overall samples; (c) the overall reward.

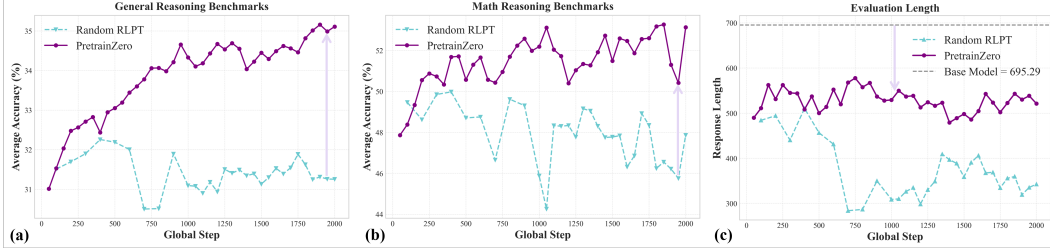


Figure 7: Evaluation comparisons between PretrainZero and Random RLPT on Qwen3-4B-Base: (a) the average accuracy on 3 general reasoning benchmarks; (b) the average accuracy on 6 math reasoning benchmarks; (c) response length on a fixed subset from MMLU-Pro.

For AMC23 and AIME24, we evaluate 32 times and report the mean@32 accuracy. We use the Qwen-Math-eval [YZH⁺24] as the math verifier.

4.2 Pretraining Results

Baselines. To compare Reinforcement Learning Pre-Training (RLPT) with conventional training patterns, we primarily establish the following baselines: 1) the base model as the initial baseline. 2) Continue Pre-Training: We continuously pretrain with the self-supervised next token prediction on the same Wikipedia data. 3) Supervised Fine-Tuning: we formulate the masked token prediction task as question-answer pairs as Fig. 5, and remove the CoT. 4) Random RLPT: We use the introduced random masked span prediction introduced in Sec. 3.1 as the strong RL baseline.

Comparison with Supervised Pretraining. We summarize the overall and detailed math performance in Table 1 and 2 respectively. Compared with the base model, Continued Pre-Training and Supervised Fine-Tuning lead to performance drops of 16.47 and 8.09 on Qwen3-4B, and 24.75 and 10.27 on Qwen3-8B, respectively. This occurs because, for highly optimized models, supervised learning on low-quality Wikipedia passages offers limited meaningful information and can even disrupt the model’s pretrained distribution. For RLPT, Random RLPT and PretrainZero improve 3.05 and 7.25 on Qwen3-4B, and 3.89 and 5.71 on Qwen3-8B. These results suggest that reinforcement learning is capable of extracting meaningful supervision from relatively low-quality data.

Comparison with Reinforcement Pretraining. As mentioned in Sec. 3.1, previous RPT training on high-entropy tokens quickly stops optimization when applied to real-world pretraining corpus. We compare the training dynamics between Random RLPT and PretrainZero in Fig. 6 and 7. As training steps increase, we observe that both the reasoning length of PretrainZero and its performance on general- and math-reasoning benchmarks improve consistently. This indicates that PretrainZero’s reasoning ability is gradually strengthened, similar to RLVR in DeepSeek-R1. Compared with Random RLPT, the active-learning strategy arouses longer CoT trajectories and noticeably stronger reasoning performance. As shown in Table 1, PretrainZero consistently outperforms Random RLPT by 4.20, 1.82, 3.17, and 3.16 points on the Qwen3-4B, Qwen3-8B, Qwen3-30B-A3B-MoE, and SmolLM3-3B base models, respectively.

Reasoning Efficiency. As shown in Fig. 6 (b), although the growth of CoT length in training, we need not worry about the inference efficiency of the reasoning process. The growth mainly comes from improvements in the mask-prediction capability. To verify this, we sample 10% of the MMLU-Pro

Table 1: Results on general-domain reasoning benchmarks. We compare the Base Model, Continue Pre-Training, Supervised Fine-Tuning, our Random RLPT baseline and PretrainZero. We highlight the best performance in **bold** and the second performance in underline.

Model Name	Overall AVG	MATH AVG	SuperGPQA	BBEH	MMLU-Pro
<i>Qwen3-4B-Base</i>					
Base Model	32.36	42.53	26.32	8.67	51.94
Continue PT	15.89	24.65	9.67	0.04	29.21
Supervised FT	24.27	15.55	26.38	<u>12.28</u>	42.88
Random RLPT	<u>35.41</u>	<u>47.87</u>	<u>29.10</u>	9.45	<u>55.21</u>
PretrainZero	39.61	53.13	32.28	12.68	60.37
<i>Qwen3-8B-Base</i>					
Base Model	37.07	47.48	31.12	10.49	59.19
Continue PT	12.32	27.78	9.94	0.04	11.51
Supervised FT	26.80	19.23	29.02	11.17	47.78
Random RLPT	<u>40.96</u>	<u>55.08</u>	<u>34.19</u>	<u>12.96</u>	<u>61.59</u>
PretrainZero	42.78	57.72	34.46	14.67	64.28
<i>SmolLM3-3B-Base</i>					
Base Model	16.23	32.31	12.62	3.32	16.66
Random RLPT	<u>20.25</u>	<u>35.95</u>	<u>14.48</u>	7.85	<u>22.74</u>
PretrainZero	23.41	38.03	19.44	<u>3.78</u>	32.41
<i>Qwen3-30B-A3B-MoE-Base</i>					
Base Model	38.88	52.49	33.73	10.51	58.79
Random RLPT	<u>40.38</u>	<u>52.62</u>	<u>36.33</u>	<u>12.99</u>	<u>59.57</u>
PretrainZero	43.55	58.12	36.58	14.91	64.59

Table 2: Results on math-domain reasoning benchmarks. We highlight the best performance in **bold** and the second performance in underline.

Model Name	AVG	MATH-500	Olympiad	Minerva	GSM8K	AMC	AIME24
<i>Qwen3-4B-Base</i>							
Base Model	42.53	73.30	37.30	<u>22.10</u>	86.30	36.17	0.00
Continue PT	24.65	38.00	13.60	11.00	67.00	15.00	3.30
Supervised FT	15.55	28.50	8.10	14.30	27.40	15.00	0.00
Random RLPT	<u>47.87</u>	<u>74.80</u>	<u>38.50</u>	<u>22.10</u>	<u>87.50</u>	<u>54.30</u>	<u>10.00</u>
PretrainZero	53.13	79.10	42.70	33.80	92.90	56.95	13.30
<i>Qwen3-8B-Base</i>							
Base Model	47.48	70.10	35.30	25.40	91.50	52.58	10.00
Continue PT	27.78	42.70	16.90	11.80	55.30	33.28	6.70
Supervised FT	19.23	30.50	11.70	15.40	32.80	25.00	0.00
Random RLPT	<u>55.08</u>	<u>79.20</u>	42.70	39.00	<u>93.80</u>	<u>62.50</u>	<u>13.30</u>
PretrainZero	57.72	81.90	<u>42.50</u>	43.40	93.50	65.00	20.00
<i>SmolLM3-3B-Base</i>							
Base Model	32.31	53.80	20.40	14.00	81.20	22.81	<u>1.65</u>
Random RLPT	<u>35.95</u>	<u>59.00</u>	<u>21.50</u>	<u>20.20</u>	<u>82.50</u>	32.50	0.00
PretrainZero	38.03	62.60	25.60	22.10	83.70	<u>27.50</u>	6.70
<i>Qwen3-30B-A3B-MoE-Base</i>							
Base Model	52.49	74.70	<u>43.00</u>	22.80	<u>91.10</u>	<u>66.95</u>	<u>16.36</u>
Random RLPT	<u>52.62</u>	<u>79.20</u>	<u>41.20</u>	<u>38.60</u>	<u>82.40</u>	<u>59.77</u>	<u>14.58</u>
PretrainZero	58.12	81.70	43.40	41.20	94.40	70.62	17.40

prompts and evaluate the reasoning length. As shown in Fig. 7 (c), for the same questions, the reasoning length remains similar during RLPT. Moreover, compared with the base model, RLPT actually improves the efficiency of CoT reasoning.

Reasoning Pattern. As shown in Fig. 8, we compare the reasoning patterns between Random RLPT and PretrainZero from the Qwen3-8B-Base. Given the same masked target, Random RLPT directly outputs the answer without any explicit reasoning. In contrast, PretrainZero first explores multiple

Reasoning Patterns: Random RLPT / PretrainZero

Paragraph: In addition to his 1983 Triple Crown wins, Ralph Hanover won seventeen additional [mask] events, including the very important Adios and Meadowlands Pace ...

Random RLPT: Here is the passage with the masked words filled in:
In addition to his 1983 Triple Crown wins, Ralph Hanover won seventeen additional stakes ...
\boxed{stakes} ...

PretrainZero: To determine the masked word in the sentence:
"In addition to his 1983 Triple Crown wins, Ralph Hanover won seventeen ..."

1. **Analyze the sentence structure**: The sentence is: "Ralph Hanover won seventeen ..."
2. **Identify the missing word**: The sentence structure suggests that the missing word should complete the phrase "seventeen additional [mask] events." **Common phrases include** "seventeen additional stakes events" or "seventeen additional championship events."
3. **Consider the context**: The sentence mentions specific races like ... These are all racing events, particularly in harness racing. The mention of world records and specific track distances further indicates that the missing word is related to racing events.
4. **Determine the most appropriate word**: The most common completion ... in this context is ... However, the most relevant term in the context of racing achievements is "seventeen additional stakes events."
5. **Conclusion**: The missing word is "stakes."

Therefore, the predicted masked word is:

In addition to his 1983 Triple Crown wins, Ralph Hanover won seventeen additional \boxed{stakes} events ...

Figure 8: Reasoning pattern comparison between Random RLPT and PretrainZero.

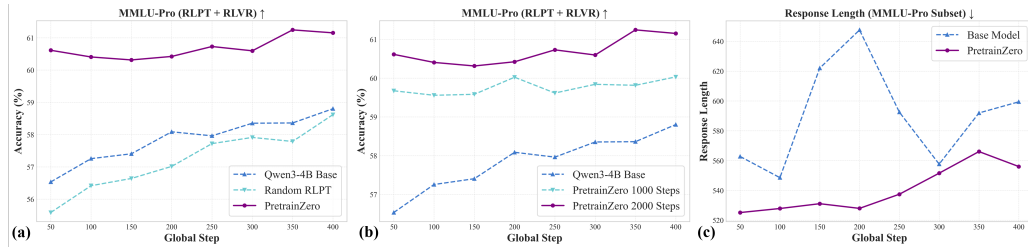


Figure 9: RLPT performance after the same RLVR post-training. (a) Comparison of Qwen3-4B-base, Random RLPT, and PretrainZero. (b) Comparison of Qwen3-4B-base, PretrainZero with 1000 and 2000 steps RLPT. (c) Response length comparison in the same MMLU-Pro subset.

possibilities, analyzes and verifies them step by step, and finally summarizes to reach a conclusion. Since the mask-prediction objective does not appear in downstream RL tasks, the emergence of such reasoning behavior during pretraining provides a stronger reasoning ability for generalization.

4.3 Post-Training Results

To investigate whether PretrainZero can improve the general reasoning capabilities of the foundation model for efficient RL finetuning, we apply RLVR as a post-training stage on PretrainZero. For the general RLVR task, we follow the General Reasoner recipe [MLJ⁺25]. Specifically, we apply the Web-Instruct dataset [MLJ⁺25] in a Question-Answer format, and the same pretrained reward model

Table 3: Results on general-domain reasoning benchmarks after the RLVR post-training. We perform the general RLVR post-training [MLJ⁺25] from the Qwen3-4B-Base model, Random RLPT, and PretrainZero with 1000 / 2000 step RLPT. **RLPT / RLVR** indicates RL steps in RLPT and RLVR stages respectively. We highlight the best performance in **bold**.

Model Name	RLPT / RLVR	Overall AVG	MATH AVG	SuperGPQA	BBEH	MMLU-Pro
Base Model	– / 400	37.90	50.96	30.26	11.59	58.80
Random RLPT	2000 / 400	38.43	51.49	30.77	12.83	58.62
PretrainZero	1000 / 400	39.15	51.84	32.32	12.39	60.03
PretrainZero	2000 / 400	40.46	53.77	33.30	13.61	61.15

Table 4: Results on math-domain reasoning benchmarks after the RLVR post-training. **RLPT / RLVR** indicates RL steps in RLPT and RLVR stages respectively.

Model Name	RLPT / RLVR	AVG	MATH-500	Olympiad	Minerva	GSM8K	AMC	AIME24
Base Model	– / 400	50.96	75.70	41.80	31.60	91.30	52.03	13.30
Random RLPT	2000 / 400	47.87	74.80	38.50	22.10	87.50	54.30	10.00
PretrainZero	1000 / 400	51.84	77.00	43.00	32.40	92.50	55.00	11.13
PretrainZero	2000 / 400	53.77	78.80	43.00	39.70	93.00	54.84	13.30

as the verifier. For efficient RL finetuning, we train 400 steps on Qwen3-4B series models with a single node with 8× H800 GPUs, which supports at most the 128 batchsize, 1/8 compared with General Reasoner.

We evaluate the training process in Fig. 9, and report the final general-domain and math-domain performance in Table 3 and Table 4, respectively. As shown in Fig. 9 (b), the performance consistently improves as the training starts progressing from the base model to PretrainZero at 1000 RLPT steps and further to PretrainZero at 2000 RLPT steps on MMLU-Pro. As shown in Fig. 9 (c), PretrainZero has more stable and efficient CoT in downstream RLVR. Compared with the base model, PretrainZero significantly improves the math average and overall accuracy by 2.18 and 2.56 points end-to-end.

4.4 Ablation Studies

Specific Domain. To explore the impact of data domain on RLPT, we compare RLPT performance on the Wikipedia corpus in the general-domain versus the MathPile [WLXL24] dataset in the math-domain. As shown in Fig. 10 (a), directly using general-domain Wikipedia data yields better performance. Since curating high-quality mathematical data requires substantial expert effort, we recommend using general-domain data for a much lower cost of data acquisition.

Training Robustness. To confirm RLPT over 2000 steps, we evaluate different mask regularization strategies: 1) PretrainZero: For the generated masks, we retain only those whose underlying spans appear fewer than eight times within the paragraph. 2) PretrainZero-OneMask: Based on PretrainZero, if a generated mask appears multiple times in the paragraph, we randomly replace only one occurrence with [mask] and make prediction. 3) PretrainZero-Words: Since PretrainZero may produce masks that cover incomplete words—reducing interpretability—we filter masks that keep only complete word spans. As shown in Fig. 10 (b), three recipes can be trained stably, and PretrainZero consistently achieves better performance.

5 Related Works and Discussion

Self-Supervised Pretraining for LLMs. Scalable self-supervised pretraining [KMH⁺20] formulates the foundation of advanced large language models. Under the simple and fixed learning pattern, the next token prediction, autoregressive LLMs [VSP⁺17, TBB⁺25] can be trained on massive corpus at the Internet-scale, establishing strong general-purpose capabilities. Beyond this pattern, token-masked prediction objectives [DCLT19, JCL⁺19] continue to play an important role in the pretraining of language models, such as in BERT-style embedding models [CXZ⁺24], diffusion language models [NZY⁺25], and code-focused pretraining [HYC⁺24]. The reliability, scalability, and broad applicability of self-supervised learning offer key insights for reinforcement pretraining and highlight its potential as a fundamental training strategy.

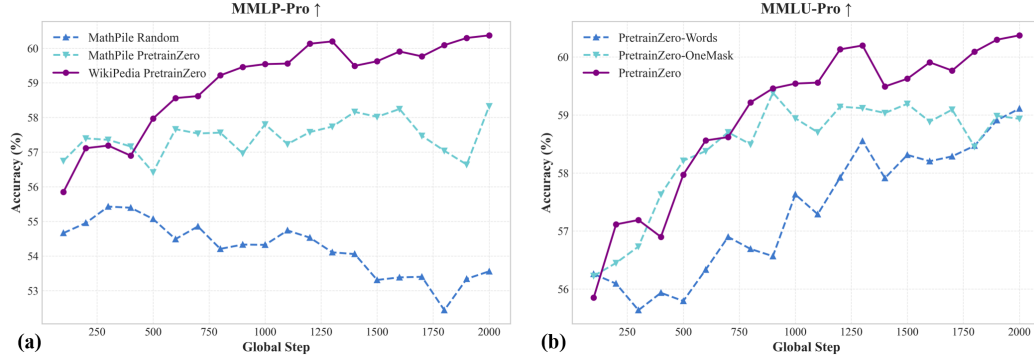


Figure 10: Comparisons for data domain and mask regularization. (a) MMLU-Pro performance on MathPile and Wikipedia. (b) MMLU-Pro performance with different mask regularization strategies.

Reinforcement Learning for LLMs. Recent large reasoning models are largely driven by post-training reinforcement learning, enabling human-expert performance in specialized domains such as web agents [TLZ⁺25], tool use [PMY⁺], software development [JYW⁺23], and mathematics [GYZ⁺25]. Despite this progress, existing RLHF [OWJ⁺22, BJN⁺22] and RLVR [LMP⁺25] approaches rely heavily on human annotation and domain-specific verification environments, leading to a severe data bottleneck in general domains [MLJ⁺25, ZLS⁺25]. For RLHF, reward models must be continuously updated with human-labeled data to avoid reward hacking. For RLVR, training data must come from domains with verifiable ground-truth answers, and the construction of verifiable environments fundamentally limits its scalability for general reasoning tasks [MLJ⁺25].

Reinforcement Pretraining. To overcome the substantial verification data-wall, Reinforcement Learning Pre-Training (RLPT) has recently emerged as a promising direction, which constructs general-purpose RLVR directly on pretraining corpus using self-supervised objectives. Early works including Quiet-STaR [ZHS⁺24] and Fast Quiet-STaR [HXY⁺25] focus on token-level reasoning. Reinforcement Pre-Training (RPT) [DDT⁺25] is the first to apply the next-token-prediction as the RLVR objective, demonstrating the feasibility of general-purpose RL. However, RPT remains limitations, such as relying on synthetic OmniMath data with CoT annotations rather than real pretraining distributions, and training from a post-trained model instead of a base model, which prevents RPT from being practical and prolonged [LDL⁺25] RLPT.

Recently, PRT [HAP⁺25] and RLPT¹ [LLX⁺25] are proposed around the similar period as this work. PRT incorporates reinforcement learning as an auxiliary objective to the self-supervised pretraining and does not exclude some QA-style training data. RLPT¹ employs an additional reward model as a verifier for the sentence-level prediction objective and further introduces a high-quality SFT cold-start. Despite these advantages, the foundational questions in RLPT remain unexplored: under fully self-supervised conditions—removing reward models, SFT cold start, and supervised cross-entropy losses—can stand-alone RLPT be effectively trained on noisy, real-world pretraining corpus? And how to improve learning efficiency in low-information-density pretraining data? Addressing these fundamental questions becomes the primary focus of this work.

6 Conclusion

This work introduces the stand-alone reinforcement pretraining method in a real-world pretraining corpus, named PretrainZero. Coupled with PretrainZero, a new reinforcement active pretraining framework is proposed to explore informative, verifiable, and not-yet-mastered content in noisy pretraining data. Thanks to active learning ability, PretrainZero significantly surpasses previous fixed learning patterns, such as continued pretraining, supervised fine-tuning, and random or entropy-based reinforcement pretraining. We reveal that even Wikipedia, which has already been trained during base model pretraining, can successfully improve end-task performance with reinforcement and active learning methods. We believe that there would be great potential to explore more efficient learning patterns to discover latent information from the pretraining corpus in the future.

References

- [AAA⁺23] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat,
et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [BJN⁺22] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma,
Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and
harmless assistant with reinforcement learning from human feedback. *arXiv preprint
arXiv:2204.05862*, 2022.
- [BMR⁺20] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla
Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini
Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, et al. Language
models are few-shot learners. In *Advances in Neural Information Processing Systems*,
volume 33, pages 1877–1901. Curran Associates, Inc., 2020.
- [BTK⁺] L Berglund, M Tong, M Kaufmann, M Balesni, AC Stickland, T Korbak, and O Evans.
The reversal curse: Llms trained on “a is b” fail to learn “b is a”. arxiv 2023. *arXiv
preprint arXiv:2309.12288*.
- [CKB⁺21] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz
Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training
verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- [CXZ⁺24] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge
m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings
through self-knowledge distillation. *arXiv preprint arXiv:2402.03216*, 2024.
- [CZY⁺25] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schu-
urmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A
comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*,
2025.
- [DCLT19] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-
training of deep bidirectional transformers for language understanding. In *Proceedings
of the 2019 Conference of the North American Chapter of the Association for Computa-
tional Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*,
pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational
Linguistics.
- [DDT⁺25] Qingxiu Dong, Li Dong, Yao Tang, Tianzhu Ye, Yutao Sun, Zhifang Sui, and Furu Wei.
Reinforcement pre-training. *arXiv preprint arXiv:2506.08007*, 2025.
- [DYM⁺25] Xinrun Du, Yifan Yao, Kaijing Ma, Bingli Wang, Tianyu Zheng, King Zhu, Minghao
Liu, Yiming Liang, Xiaolong Jin, Zhenlin Wei, et al. SuperGPQA: Scaling LLM
evaluation across 285 graduate disciplines. *arXiv preprint arXiv:2502.14739*, 2025.
- [GSY⁺24] Bofei Gao, Feifan Song, Zhe Yang, Zefan Cai, Yibo Miao, Qingxiu Dong, Lei Li,
Chenghao Ma, Liang Chen, Runxin Xu, Zhengyang Tang, Benyou Wang, Daoguang Zan,
Shanghaoran Quan, Ge Zhang, Lei Sha, Yichang Zhang, Xuancheng Ren, Tianyu Liu,
and Baobao Chang. Omni-MATH: A universal Olympiad level mathematic benchmark
for large language models. *ArXiv*, abs/2410.07985, 2024.
- [GYZ⁺25] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao
Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning
capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [HAP⁺25] Ali Hatamizadeh, Syeda Nahida Akter, Shrimai Prabhumoye, Jan Kautz, Mostofa
Patwary, Mohammad Shoeybi, Bryan Catanzaro, and Yejin Choi. Rlp: Reinforcement
as a pretraining objective. *arXiv preprint arXiv:2510.01265*, 2025.

- [HBK⁺21] Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Xiaodong Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *ArXiv*, abs/2103.03874, 2021.
- [HLB⁺24] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3828–3850, 2024.
- [HXY⁺25] Wei Huang, Yizhe Xiong, Xin Ye, Zhijie Deng, Hui Chen, Zijia Lin, and Guiguang Ding. Fast quiet-star: Thinking without thought tokens. *arXiv preprint arXiv:2505.17746*, 2025.
- [HYC⁺24] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
- [HYW⁺25] Chengsong Huang, Wenhao Yu, Xiaoyang Wang, Hongming Zhang, Zongxia Li, Ruosen Li, Jiaxin Huang, Haitao Mi, and Dong Yu. R-zero: Self-evolving reasoning llm from zero data. *arXiv preprint arXiv:2508.05004*, 2025.
- [JCL⁺19] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. SpanBERT: Improving pre-training by representing and predicting spans. *arXiv preprint arXiv:1907.10529*, 2019.
- [JCL⁺20] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the association for computational linguistics*, 8:64–77, 2020.
- [JYW⁺23] Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- [KFB⁺25] Mehran Kazemi, Bahare Fatemi, Hritik Bansal, John Palowitch, Chrysovalantis Anastasiou, Sanket Vaibhav Mehta, Lalit K Jain, Virginia Aglietti, Disha Jindal, Yuanzhu Peter Chen, et al. Big-bench extra hard. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 26473–26501, 2025.
- [KMH⁺20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *CoRR*, abs/2001.08361, 2020.
- [LAD⁺22] Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Venkatesh Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. Solving quantitative reasoning problems with language models. *ArXiv*, abs/2206.14858, 2022.
- [LDL⁺25] Mingjie Liu, Shizhe Diao, Ximing Lu, Jian Hu, Xin Dong, Yejin Choi, Jan Kautz, and Yi Dong. Prorl: Prolonged reinforcement learning expands reasoning boundaries in large language models. *arXiv preprint arXiv:2505.24864*, 2025.
- [LFX⁺24] Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
- [LLX⁺25] Siheng Li, Kejiao Li, Zenan Xu, Guanhua Huang, Evander Yang, Kun Li, Haoyuan Wu, Jiajia Wu, Zihao Zheng, Chenchen Zhang, et al. Reinforcement learning on pre-training data. *arXiv preprint arXiv:2509.19249*, 2025.

- [LMP⁺25] 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. Tulu 3: Pushing frontiers in open language model post-training, 2025.
- [MLJ⁺25] Xueguang Ma, Qian Liu, Dongfu Jiang, Ge Zhang, Zejun Ma, and Wenhua Chen. General-reasoner: Advancing llm reasoning across all domains. *arXiv preprint arXiv:2505.14652*, 2025.
- [NZY⁺25] Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, and Chongxuan Li. Large language diffusion models. *arXiv preprint arXiv:2502.09992*, 2025.
- [Ope24] OpenAI. Hello, gpt-4o. <https://openai.com/index/hello-gpt-4o>, 2024.
- [OWJ⁺22] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [PMY⁺] Shishir G Patil, Huanzhi Mao, Fanjia Yan, Charlie Cheng-Jie Ji, Vishnu Suresh, Ion Stoica, and Joseph E Gonzalez. The berkeley function calling leaderboard (bfc1): From tool use to agentic evaluation of large language models. In *Forty-second International Conference on Machine Learning*.
- [RWC⁺19] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- [Set09] Burr Settles. Active learning literature survey. 2009.
- [SLXK24] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- [SWD⁺17] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [SWZ⁺24] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- [TBB⁺25] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijie Chen, Yanru Chen, Yuankun Chen, Yutian Chen, et al. Kimi k2: Open agentic intelligence. *arXiv preprint arXiv:2507.20534*, 2025.
- [TLZ⁺25] Tongyi DeepResearch Team, Baixuan Li, Bo Zhang, Dingchu Zhang, Fei Huang, Guangyu Li, Guoxin Chen, Huifeng Yin, Jialong Wu, Jingren Zhou, et al. Tongyi deepresearch technical report. *arXiv preprint arXiv:2510.24701*, 2025.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 6000–6010, 2017.
- [WLXL24] Zengzhi Wang, Xuefeng Li, Rui Xia, and Pengfei Liu. Mathpile: A billion-token-scale pretraining corpus for math. *Advances in Neural Information Processing Systems*, 37:25426–25468, 2024.

- [WMZ⁺24] Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, et al. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- [YCL⁺] Yang Yue, Zhiqi Chen, Rui Lu, Andrew Zhao, Zhaokai Wang, Yang Yue, Shiji Song, and Gao Huang. Does reinforcement learning really incentivize reasoning capacity in llms beyond the base model?, 2025. URL <https://arxiv.org/abs/2504.13837>.
- [YLB⁺25] Jason Yang, Ravi G Lal, James C Bowden, Raul Astudillo, Mikhail A Hameedi, Sukhvinder Kaur, Matthew Hill, Yisong Yue, and Frances H Arnold. Active learning-assisted directed evolution. *Nature Communications*, 16(1):714, 2025.
- [YLY⁺25] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- [YZH⁺24] An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. Qwen2. 5-math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122*, 2024.
- [YZZ⁺25] Qiyang Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, et al. DAPO: An open-source LLM reinforcement learning system at scale. *ArXiv*, abs/2503.14476, 2025.
- [ZHS⁺24] Eric Zelikman, Georges Harik, Yijia Shao, Varuna Jayasiri, Nick Haber, and Noah D Goodman. Quiet-star: Language models can teach themselves to think before speaking. *arXiv preprint arXiv:2403.09629*, 2024.
- [ZLS⁺25] Xiangxin Zhou, Zichen Liu, Anya Sims, Haonan Wang, Tianyu Pang, Chongxuan Li, Liang Wang, Min Lin, and Chao Du. Reinforcing general reasoning without verifiers. *arXiv preprint arXiv:2505.21493*, 2025.