# History Rhymes: Accelerating LLM Reinforcement Learning with RhymeRL

Jingkai He[1], Tianjian Li[2], Erhu Feng[✉1], Dong Du[✉1], Qian Liu[2], Tao Liu[2],
Yubin Xia[1], Haibo Chen[1]

[1]Institute of Parallel and Distributed Systems, Shanghai Jiao Tong University

[2]ByteDance

## Abstract

With the rapid advancement of large language models (LLMs), reinforcement learning (RL) has emerged as a pivotal methodology for enhancing the reasoning capabilities of LLMs. Unlike traditional pre-training approaches, RL encompasses multiple stages: *rollout*, *reward*, and *training*, which necessitates collaboration among various worker types. However, current RL systems continue to grapple with substantial GPU underutilization, due to two primary factors: (1) The rollout stage dominates the overall RL process due to test-time scaling; (2) Imbalances in rollout lengths (within the same batch) result in GPU bubbles. While prior solutions like asynchronous execution and truncation offer partial relief, they may compromise training accuracy for efficiency.

Our key insight stems from a previously overlooked observation: *rollout responses exhibit remarkable similarity across adjacent training epochs*. Based on the insight, we introduce RhymeRL, an LLM RL system designed to accelerate RL training with two key innovations. First, to enhance rollout generation, we present HistoSpec, a speculative decoding inference engine that utilizes the similarity of historical rollout token sequences to obtain accurate drafts. Second, to tackle rollout bubbles, we introduce HistoPipe, a two-tier scheduling strategy that leverages the similarity of historical rollout distributions to balance workload among rollout workers. We have evaluated RhymeRL within a real production environment, demonstrating scalability from dozens to thousands of GPUs. Experimental results demonstrate that RhymeRL achieves a 2.6x performance improvement over existing methods, without compromising accuracy or modifying the RL paradigm.

## 1 Introduction

Post-training with reinforcement learning (RL) has emerged as a new paradigm for scaling and enhancing the capabilities of LLMs [1–5]. Representative post-training models, such as DeepSeek-R1 [1], have demonstrated significant improvements in areas including coding [5], mathematics [6] and many others [7–10]. A standard RL pipeline comprises three main stages: *rollout*, *reward*, and *training*. In the rollout stage, the LLM generates a large number of tokens, with extended reasoning to improve the quality of final responses (test-time scaling [11]). In the reward stage, a reward score is assigned to each response (i.e., sample). In the subsequent training stage, the model's weights are updated by computing new loss values based on the assigned rewards. As a result, RL systems are inherently complex and distributed, typically involving coordination among multiple heterogeneous workers to efficiently execute the various pipeline stages.

Current LLM RL systems face significant GPU underutilization, primarily arising from two sources. First, *overly long rollout phases*. During each RL step, the rollout phase, which entails generating a substantial number of thinking tokens, typically consumes 84% to 91% of the total time. Additionally, the autoregressive nature of LLMs prevents the rollout phase from fully utilizing GPU computational resources, resulting in substantial underutilization. Second, *bubbles caused by imbalanced rollouts*. Within a single batch, the response lengths of rollouts generated by different prompts vary dramatically, leading to a pronounced long-tail effect. As a result, completed rollout tasks must wait until the longest rollout in the batch finishes before advancing to reward and training stages. Due to these factors, we observe that SOTA RL systems, such as veRL [12], experience over 46% GPU resource idleness, significantly compromising the efficiency of RL.

To address these problems, recent research has mainly focused on scheduling optimizations. Some systems mitigate the performance degradation caused by long-tail rollouts by truncation (e.g., Kimi-K2 [13]) or reducing their batch size (e.g., StreamRL [14]). However, truncation methods inherently introduce a trade-off between accuracy and efficiency, while adjusting batch sizes provides only limited alleviation of the long-tail problem (~10% [14]). Some systems (e.g., AReaL [15]) depart from conventional RL systems by enabling full asynchronization between the rollout and training phases to reduce GPU idle time. However, this results in rollouts utilizing stale model weights, which changes the paradigm of RL. Additionally, model weight updates during rollout trigger the recomputation of all tokens in the ongoing rollouts, leading to considerable overhead. Therefore, enhancing the efficiency of rollout generation and minimizing GPU resource idleness remain critical challenges in the development of LLM RL systems.

Through a detailed analysis of the real-world RL training, we observe that the rollout process in RL is distinct from conventional LLM inference. Specifically, a complete RL train-
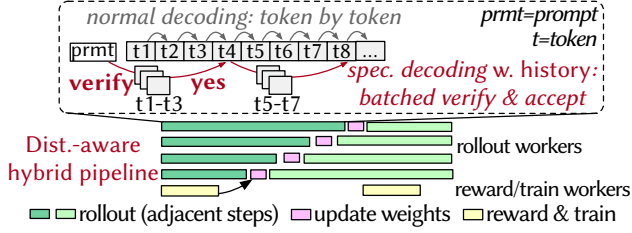
Figure 1. An overview of RhymeRL's designs.



Figure 2. Phases and time distribution in LLM RL. *We train 32B models with math/code datasets using veRL [12] and GRPO [6].*

ing consists of performing multiple inference passes on the same prompt in different RL steps (50—100 *epochs* [6, 16, 17]). While the model weights vary between these steps due to continuous updates during training, current RL algorithms (such as GRPO [6], etc. [16, 18]) apply *clipping operations* [19], which restrict the magnitude of the model's updates at each step, to maintain stable model evolution. This stability leads to a *high degree of similarity* between rollout responses across different epochs: (1) *Token sequence similarity.* The responses generated for each prompt exhibit substantial similarity to their corresponding previous rollout responses, with 75%—95% of tokens being reusable. (2) *Length distribution similarity.* Although a prompt generates responses of varying lengths in different epochs, their position within the epoch's response length ranking remains stable, with only 2%-4% of responses experiencing significant rank changes. Motivated by this observation, our central insight is to *exploit the* **similarity from historical rollouts** *to accelerate the rollout process and achieve effective load balancing across rollout workers.*

We propose a novel RL system, RhymeRL, which significantly improves LLM RL efficiency without modifying existing RL training paradigms or sacrificing accuracy, as Fig. 1 shows. First, leveraging the historical token sequence similarity, we propose a lightweight yet effective *speculative decoding* [20] mechanism, *HistoSpec*, to accelerate the rollout process and improve the computational density. HistoSpec adopts a reward-aware, tree-based management strategy to organize draft tokens, and incorporates a novel speculative strategy inspired by TCP congestion control mechanisms [21] to improve prediction accuracy. Secondly, leveraging the historical length distribution similarity, we introduce the *HistoPipe* scheduling. By complementing long and short rollouts between adjacent steps, HistoPipe effectively balances workload across rollout workers and eliminates GPU bubbles caused by the imbalanced distribution of rollout lengths. HistoPipe further tackles outliers with migration-based rebalancing, and proposes a two-tier scheduling mechanism for better complementarity.

We have implemented our system on veRL and successfully deployed it in industrial RL environments ranging from dozens to over a thousand GPUs, achieving a performance improvement of up to 2.6x. RhymeRL consistently attains SOTA performance across GPU clusters of varying scales
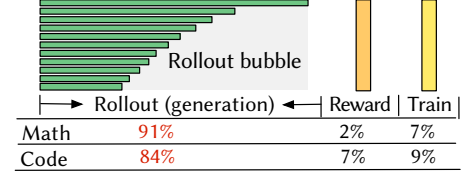
while maintaining training accuracy without compromise. RhymeRL will be open source.

## 2 Background

### 2.1 LLM Reinforcement Learning

RL-based post-training refines LLMs through interaction with environments. As the marginal benefits of pretraining diminish, RL is widely acknowledged as a pivotal approach for enhancing LLMs' reasoning capabilities [1, 2, 22, 23].

**Workflow.** A whole RL process consists of multiple repeated steps. In brief, each step comprises three stages, as Fig. 2 shows: (1) *Rollout*: the LLM generates responses to input prompts in batches (i.e., LLM inference). (2) *Reward*: The responses are scored (i.e., rewarded) using rule-based functions (e.g., running tests for code) or reward models. (3) *Train*: Loss is computed based on rewards, followed by backpropagation to optimize the model and generate new model weights.

**Algorithm.** Early algorithms like PPO [24, 25] employ not only rule-based functions/reward models to generate sample-level rewards, but also simultaneously train a critic model to produce action-level rewards. This strategy stabilized training by reducing the volatility of rewards from sample-level evaluations. Contemporary mainstream methods, exemplified by GRPO [6] and DAPO [16], have superseded critic models with *Group Relative Advantage (GRA)*. This paradigm leverages the inherent stochasticity of LLMs during rollout: for each prompt, the LLM generates *a group of responses* (usually 16 in our production). Policy updates are then guided by relative scoring within the groups. The process of generating multiple responses per prompt can be interpreted as the model exploring diverse solution pathways to a problem.

### 2.2 Speculative Decoding

During LLM decoding, each attention operation requires accessing the entire KV cache, making the decoding process memory bandwidth-bound and preventing full utilization of computational resources. Unlike conventional decoding, which generates one token per iteration (i.e., LLM forward pass), speculative decoding [20, 26–30] first predicts multiple draft tokens using a lightweight method (e.g., using a small draft model [31–34] or retrieving from corpora [35–37]). The LLM then verifies these tokens in a single forward pass by computing their logits and accepting valid tokens. Since verifying multiple tokens incurs nearly identical memory access overhead (traversing the KV cache once) as generating one

token conventionally, but with higher computational intensity, speculative decoding amortizes the memory cost across multiple tokens. This approach is theoretically proven to preserve output distribution integrity. When acceptance rates are favorable, it significantly accelerates LLM decoding [38].

RhymeRL is the first LLM RL system leveraging speculative decoding to accelerate the time-consuming rollout process.

## 3 Characterising RL Training in the Wild

Despite the widespread adoption of RL for LLM training, current RL systems still face significant performance challenges. This section presents the key implications derived from our practical experience in training state-of-the-art LLMs and analysis of real-world traces.

### 3.1 Rollout as the Major Bottleneck in RL Training

**Implication-1: Rollout dominates the RL training timeline.** During RL training, we observe that the rollout phase dominates the RL time. Specifically, as Fig. 2 shows, for LLMs trained with a maximum response length of 16K tokens, rollout accounted for 91% of the entire RL processing time for the math model and 84% for the code model. When the maximum response length was increased to 32K tokens or longer, the rollout time overhead exceeded 95%.

This significant overhead stems from three key factors: (1) *Complex reasoning demands.* During RL, LLMs are required to generate complex reasoning chains in responses. This results in long responses, and crucially, the response length tends to increase progressively as training advances. As Fig. 3a shows, the mean response length grows from 1K to 10K in 320 steps. (2) *Memory-bound decoding.* During rollout, generating each token requires an LLM forward pass, which is constrained by memory bandwidth [39]. (3) *Sequential dependency.* Due to the dependency between rollout and the subsequent training step, the training phase cannot commence until the longest sequence within the batch completes its rollout.

**Implication-2: GPU underutilization from rollout imbalance.** Preserving the rollout-train dependency is critical for training accuracy, but it introduces significant compute bubbles and leads to GPU idleness during rollout. We observe a significant long-tail effect within a batch, as shown in Fig. 3b — rollout response lengths vary widely across sequences. Due to the imbalanced distribution of rollout lengths, some rollout workers (i.e., GPUs cooperating with tensor parallelism) finish early and become idle, yet must remain inactive until all workers complete their tasks. As Fig. 3c shows, in a step, GPU monitoring reveals that the earliest-finishing GPU remains idle for ~76% of the total rollout duration.

### 3.2 State-of-the-Art Efforts and the Limitations

We summarize existing efforts on optimizing the efficiency of LLM RL systems [12–15, 17, 40–46] in Table. 1.

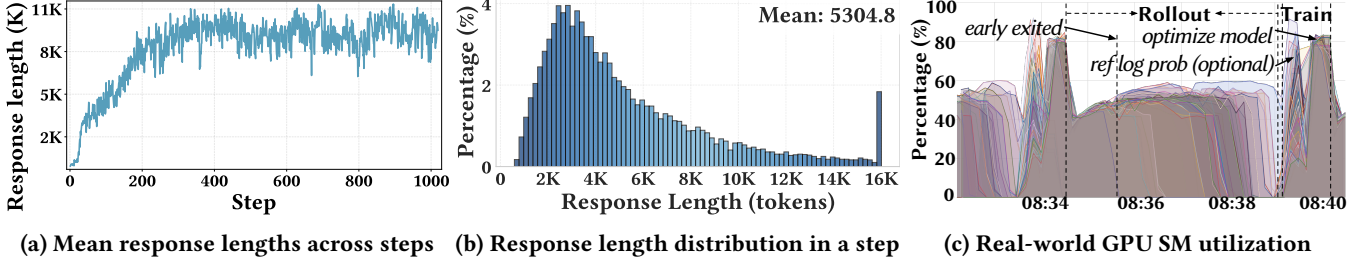| Systems | Reduce rollout time | Tackle rollout bubbles | Rollout -train pipeline | Core techniques |
|---|---|---|---|---|
| HybridFlow [12] | ✗ | ✗ | ✗ | Hybrid programming model |
| Kimi K2 [13] | ✗ | Partially | ✗ | Partial rollout |
| DeepScaleR [17] | ✗ | ✗ | ✓ | Pipelined rollout/train |
| StreamRL [14] | ✗ | Partially | ✓ | Skewness-aware scheduling |
| AReaL [15, 40] | ✗ | ✓ | ✓ | Fully async rollout |
| AsyncFlow [41] | ✗ | Partially | ✓ | Streaming pipeline |
| DistFlow [42] | ✗ | ✗ | ✓ | Distributed multi-controller |
| RhymeRL | ✓ | ✓ | ✓ | HistoSpec + HistoPipe |

**Table 1. Overview of existing systems optimizing LLM RL.**

**Pipelining rollout and training.** Early RL systems like veRL [12] and OpenRLHF [45] adopt a *colocated architecture*, where GPUs repeatedly switch between rollout and training workloads, as Fig. 4a shows. This design incurs significant overhead from context switching between workers and substantial bubbles due to strict step-wise rollout-train dependency. Consequently, DeepScalaR [17] introduces a *decoupled architecture*, relaxing the dependency by using stale weights that are one step behind for rollout (i.e., the off-policyness is $1^{1}$), which is widely proven to maintain training accuracy [14, 17, 47]. This enables coarse-grained rollout-train pipeline, and is adopted by systems including veRL [12, 47]. AsyncFlow [41] further refines the rollout-train pipeline by granularizing dependencies from batch-level to mini-batch level, enhancing pipeline efficiency (as Fig. 4b shows, the reward and train processes proceed after a mini-batch's rollout finishes). Nevertheless, imbalance persists among rollout workers, leaving significant bubbles within rollout stages.
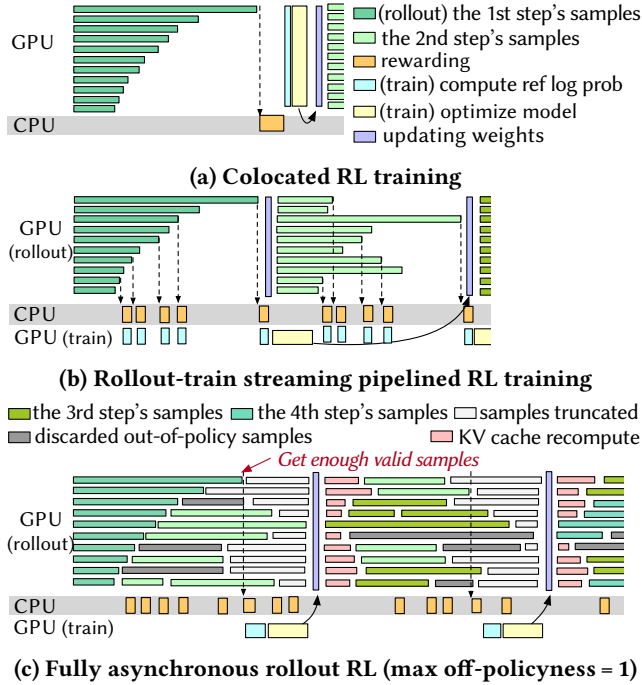
**Tackling long-tail rollout.** StreamRL [14] proposes *skewness-aware scheduling*, which allocates additional data-parallel GPUs to the prompts likely to generate long responses. By reducing batch sizes for these prompts, it alleviates the long-tail effect in rollouts. However, due to the autoregressive nature of LLM rollouts, this approach only marginally reduces the long-tail and the rollout bubbles (~10% [14]). Kimi-K2 [13] introduces *partial rollout*, which truncates excessively long responses and retains generated segments for continuation in subsequent steps. This method faces an efficiency-accuracy dilemma: Excessive truncation causes significant portions of responses to be generated by stale model weights, ultimately producing outdated reward signals; conservative truncation diminishes its effectiveness in reducing rollout bubbles.

**Fully asynchronous rollout.** AReaL [15] adopts a radical approach (Fig. 4c): (1) Fully async rollout. Rollout workers continuously generate responses while train workers select usable samples based on the maximum off-policyness threshold. (2) Aggressive dependency relaxation. Training can utilize rewards derived from rollouts generated by (potentially) stale model weights many steps prior. When new model weights are produced, ongoing rollouts are truncated to propagate updated weights. The truncated rollouts are

---

[1]The off-policyness refers to the temporal discrepancy (in steps) between the model used to generate rollouts and the current model being optimized.

(a) Mean response lengths across steps    (b) Response length distribution in a step    (c) Real-world GPU SM utilization

**Figure 3. Diving into real-world RL training of a 32B LLM.** *We train the model using math datasets (>200K samples) with 64 GPUs. The max response length is set to 16K tokens. In Fig.-c, each line represents a GPU's SM utilization. The training stage can be further divided into ref log prob computing (optional), which computes the generated tokens' logits via a forward pass of the reference LLM, and model optimizing.*



(a) Colocated RL training

(b) Rollout-train streaming pipelined RL training

(c) Fully asynchronous rollout RL (max off-policyness = 1)

**Figure 4. Existing LLM RL pipelines.** *Dashed lines represent the dependency between rollout samples and reward/train, while solid lines represent the dependency between training and weight updating.*

continued after recomputing the KV cache with new weights. This method has critical limitations: (1) Fully async rollout changes the paradigm of RL. Excessive off-policyness floods training with obsolete signals. (2) Rollout imbalance-induced GPU underutilization remains. Rollouts exceeding the off-policyness threshold are discarded, and frequent KV cache recomputation wastes GPU resources.

**Summary.** In existing systems, rollout imbalance remains a persistent bottleneck, resulting in significant GPU resource underutilization. Critically, no existing approach reduces the time required for rollout execution and improves the GPU computational underutilization of the rollout stages.

## 4 RhymeRL Overview

To resolve the persistent challenges in LLM RL: time-consuming rollouts and rollout imbalance, we design and implement

RhymeRL, which innovatively leverages historical rollout information to accelerate rollout execution and minimize imbalance, thereby achieving significant speedups in RL training.

### 4.1 Observations and Insights

**Observation.** A complete LLM RL training typically spans 50-100 *epochs* [6, 16, 17], each comprising multiple *steps* that iterate through the full dataset. In the long-term practice of LLM RL training, we observe strong ***historical similarity*** in rollout responses across epochs, characterized by:

• High *token similarity* in rollout responses generated by the same prompt across adjacent epochs, i.e., responses generated in each epoch contain a large proportion of consecutive token sequences that are identical to the sequences in the previous epoch (§5.1);

• High *response length distribution similarity* across adjacent epochs, i.e., if we rank the prompts using their response lengths, the ranking order across adjacent epochs remains similar (§6.1).

The **root cause** of historical similarity is that current RL algorithms (e.g., PPO [24], GRPO [6], DAPO [16] and GSPO [18]) apply *clipping operations* [19] to restrict the magnitude of the model's updates, which maintains stable model evolution.

**Insights.** The key insight of RhymeRL is motivated by the above novel observations. By systematically organizing and utilizing historical information, which has been overlooked in all existing RL systems, we unlock new opportunities to further enhance the overall performance of RL system. Specifically, to accelerate the rollout process, we propose a dedicated speculative decoding mechanism that leverages historical rollouts as accurate drafts. Furthermore, to balance the workload among rollout workers, we introduce a two-tier, distribution-aware pipeline scheduling strategy that exploits the distributional similarity present in historical rollouts.

### 4.2 System Overview

As Fig.5 shows, RhymeRL inherits the hybrid controller architecture proposed by HybridFlow [12] and the disaggregated rollout-train structure (Fig. 4b), utilizing dedicated *rollout workers* for response generation, *reward workers* for reward
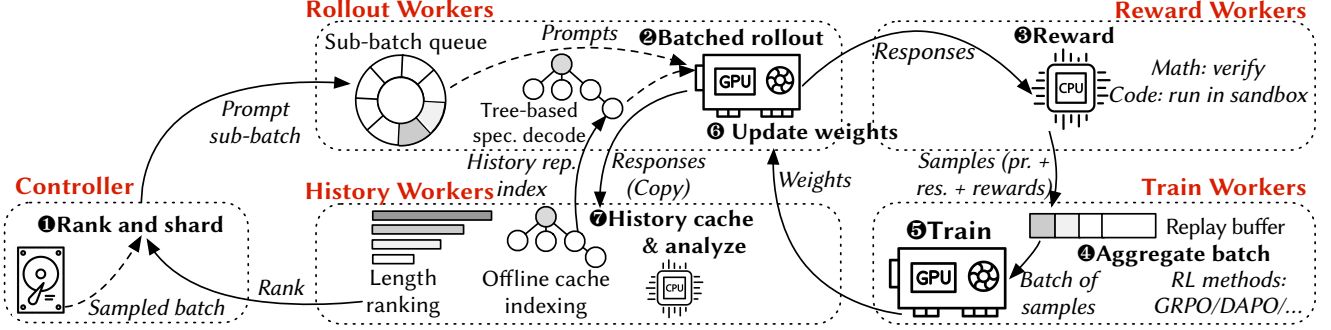
**Figure 5. RhymeRL overview.** *Solid lines indicate data flow across workers. Dashed lines indicate data dependencies within a worker.*

computation, and *train workers* for policy optimization to form a pipelined RL workflow.

To improve the RL system's overall efficiency, we adopt a *streaming pipeline* architecture. During each step, each rollout worker retrieves a sub-batch of prompts asynchronously dispatched by the controller from its *prompt sub-batch queue*, and generates responses using the inference engine (❷). Completed rollout responses (i.e., samples) proceed to reward workers for scoring (❸) before being transferred to train workers' replay buffer. Once the replay buffer accumulates sufficient samples matching the batch size (❹), the train workers execute user-defined RL algorithms to optimize the model (❺). Following full-batch policy optimization, updated model weights are propagated from train workers to the *weight buffers* (on host memory) of rollout workers. Before processing each step's sub-batch, rollout workers synchronize the weights to their GPUs (❻). This *worker-controlled weight update* strategy removes global synchronization overhead and minimizes idle waiting times.

To accelerate rollout generation via speculative decoding (HistoSpec), RhymeRL employs a *reward-aware suffix-tree-based approach* (§5.3) that efficiently generates speculation drafts from historical data with minimal overhead. We further propose an *AIMD-inspired token speculation strategy* (§5.4) to dynamically adjust the length of speculative tokens, thereby achieving higher acceptance rates while improving computational density. Moreover, to balance the workload among rollout workers, RhymeRL implements a *distribution-aware scheduling strategy* (HistoPipe), which leverages inter-step complementarity to achieve overall workload balance. It addresses anomalous outliers with *migration-based rebalancing* (§6.3), and tackles highly skewed rollout time distributions with a *two-tier scheduling* mechanism (§6.4).

To efficiently manage historical rollout information, we introduce *history workers* that operate on idle CPU resources within RL cluster. With the integration of history workers, the RL workflow incorporates two additional stages. First, the controller leverages the historical length ranking provided by history workers to enable more effective task scheduling (❶), while it also distributes relevant historical responses to the assigned rollout workers. Second, after a rollout worker

generates responses, it asynchronously sends them to history workers (❼), which update corresponding data structures.

## 5 HistoSpec: Speculative Rollout Generation

Since the rollout phase constitutes the majority of execution time in RL workflows, accelerating rollout is critical for improving overall system efficiency. Although existing RL systems utilize SOTA inference engines [48, 49] during rollout, their performance is fundamentally constrained by memory bandwidth, due to the extremely long rollout sequences. However, we observe that rollout is a specialized LLM inference that demonstrates high historical similarity. Motivated by this, we introduce a novel speculative strategy dedicated to the rollout phase, and achieve significant performance improvement in the RL training scenario.

### 5.1 Observation: Token Similarity in Rollout

RL training comprises multiple iterative epochs (typically 50−100 [6]), during which the same prompt will be repeatedly sampled across different epochs. Mainstream RL algorithms [6, 16, 18, 24] use clipping operations [19] to restrict the magnitude of the model's updates. Furthermore, they use *Group Relative Advantage (GRA) Optimization*, generating multiple responses per prompt (8−64 responses), which enables comprehensive exploration of diverse reasoning trajectories throughout each epoch. Consequently, the outputs generated for the same prompt during rollouts at adjacent epochs exhibit a high degree of similarity.

To quantitatively evaluate the similarity of tokens generated during rollouts, we analyze the response traces generated during RL training across five math/code datasets, as listed in Fig. 6c. For each response, we simulate its rollout "generation" from the beginning, and use the last three "generated" tokens as the prefix to search for exact matches in the historical responses from the previous epoch. When a match is found, we record the length of identical token sequences that follow the prefix in both historical and current responses, labeling them as "accept" tokens. Otherwise, we will forward the response to "generate" one token. The routine is repeated until the end of the response. As Fig. 6a shows, across 8 epochs, for math tasks, 93% of tokens could
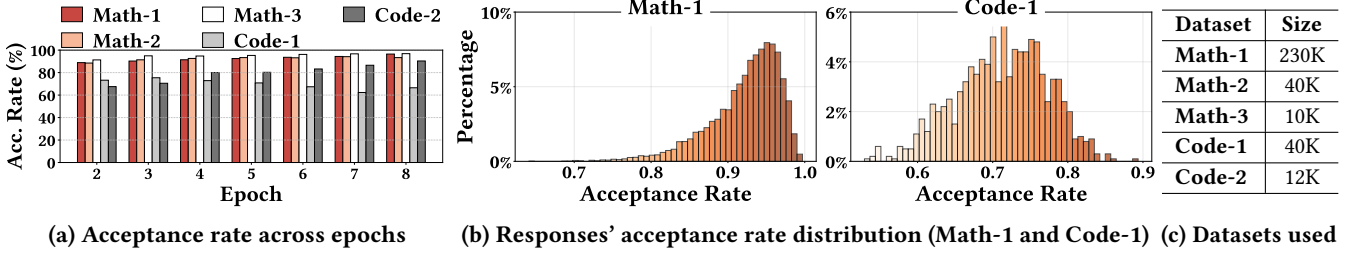
(a) Acceptance rate across epochs    (b) Responses' acceptance rate distribution (Math-1 and Code-1)   (c) Datasets used

**Figure 6. Characterizing historical token similarity.** *We use five datasets to train 14B LLMs respectively using the GRPO algorithm.*
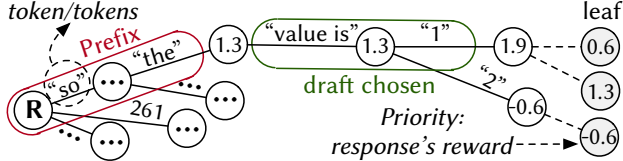


**Figure 7. Suffix-tree based draft generation.** *R = Root. For ease of understanding, we convert token ids into words.*

be successfully "accept" using this routine (on average, 75% for code). Furthermore, as training progresses (i.e., with increasing epochs), the similarity increases. Fig.6b shows the distribution of the responses' acceptance rates.

## 5.2 Speculative Rollout with History

Leveraging historical token similarity, we design HistoSpec, which utilizes speculative decoding to accelerate rollouts and uses historical responses as draft sources. In each decoding iteration, HistoSpec uses *the last few generated tokens* as the prefix to search for matches within the prompt's historical responses. Upon matching, it extracts a certain number of subsequent tokens following the prefix as drafts. The drafts are then verified via a single LLM forward pass, with some (possibly all) tokens accepted (§2.2). This history-based speculative rollout improves computational resource utilization during rollouts, and reduces the time of the rollout stage.

**Technical challenges.** However, this approach encounters two challenges: (1) Over-long rollout sequences induce significant prefix matching overhead, and the branching in historical responses complicates draft token selection; (2) Unpredictable draft acceptance length creates a trade-off between underutilizing compute capacity (when predicting too few tokens) and wasting resources on verifying rejected tokens (when over-predicting).

## 5.3 Tree-based Historical Rollout Management

Draft generation overhead is a pivotal factor determining speculative decoding's effectiveness. However, existing corpus-based speculation methods face a dilemma between retrieval costs and index building costs. E.g., n-gram [36, 50] requires time-consuming linear scans for prefix matching, and SuffixDecoding [35] limits the lengths of draft source sequences to reduce index building costs. During LLM RL, each epoch generates multiple long responses (totaling hundreds of thou-

sands of tokens per prompt). Traditional draft generation methods struggle to operate efficiently under the constraints.

**Async cache building.** The RL workloads allow us to relax the constraints on draft generation. In RL sampling strategies, the same request is generally not re-sampled until multiple steps later. Therefore, we are able to shift the computational overhead associated with retrieval to the indexing phase. RhymeRL employs history workers to index historical rollout sequences asynchronously. Upon generating new responses, the controller dispatches them to history workers for index construction. When scheduling prompts to rollout workers (asynchronously), the controller notifies the corresponding history workers to transfer the relevant indexed cache.

**Reward-aware tree-based management.** RhymeRL employs suffix trees [51] to index cached responses, which enables $O(m)$-time matching for length-$m$ prefixes. RhymeRL maintains a dedicated tree for each prompt, indexing all its historical responses generated in its last rollout. Since modern RL algorithms generate multiple responses per prompt to explore multiple solution paths, a prefix may branch to divergent suffixes, complicating draft token selection. We observe that, during training, RL algorithms optimize the model towards generating high-reward solutions with higher probability. Therefore, we add a *priority* value to each tree node, weighted by the sum of the rewards of its branch. For every suffix branch present, HistoSpec selects the one with the highest priority. Through this RL-algorithm-guided design, HistoSpec maximizes the speculation acceptance rate.

As Fig.7 shows, in the suffix tree, each node represents a subsequence formed by the path from root to that node, where each leaf node corresponds to a complete suffix of a response, and each edge represents one or more tokens extending the sequence of parent node to the sequence of the child. For each leaf node, its priority (marked on the nodes in Fig.7) is the sum of the rewards for the responses that end with this suffix (that the leaf node represents). A parent node's priority is the sum of its children's priorities.

**Resources.** Both the tree's construction time overhead and memory overhead are $O(n)$ for $n$ tokens [52]. GPU clusters have sufficient CPU resources (64–128+ cores and multi-TB host memory), which are mostly idle during previous RL flows. RhymeRL strictly limits the resources used by history workers using OS methods [53], preventing them from interfering with other workers. For RL training with a
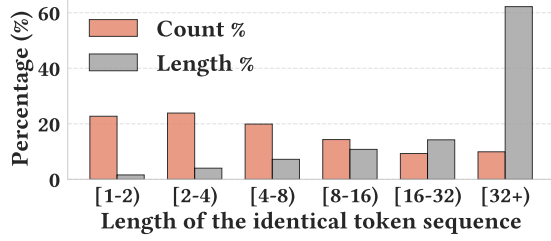
Figure 8. Distribution of speculative hit lengths.



**Figure 9. The ranking group changes of responses across 20 epochs.** *Since RL algorithms like GRPO generate multiple responses for each prompt, we analyze the previous epoch by assigning all responses of each prompt to the same ranking group based on their* **median** *length, yielding the* predicted group. *For the current epoch, we sort all responses by their* actual *lengths to assign groups, obtaining the* real group, *and compare it with the predicted group to collect the data. In the figure, the lower parts of the bars represent the proportion of responses that remain in the same or a lower group; the upper parts represent the proportion of responses that, despite changing groups, only move up one group and have lengths in the shorter 50% of the new group (i.e., shift near the group boundary).*

230K dataset and 16K max response length, the host memory overhead of suffix trees is < 80GB per node with 8 nodes. HistoSpec also supports compression and swapping to SSD for memory saving, and checkpoint for fault tolerance.

### 5.4 AIMD-like Token Speculation

To determine the length of tokens predicted by HistoSpec in each iteration, we conduct a detailed analysis of the length distribution of identical token sequences in rollout responses. As shown in Fig. 8, short segments (1-2 tokens) dominate in quantity, while long segments account for the majority in the total length. This poses a challenge for HistoSpec: If we predict many tokens per iteration, most tokens would be rejected, wasting significant computation on verifying them. If we predict a small number of tokens, computational resources cannot be fully utilized, undermining speculative decoding's advantage in improving computational density.

We find that network congestion control encounters similar challenges. Classical TCP congestion control employs the AIMD (*Additive Increase, Multiplicative Decrease* [21]) principle: gradually increasing the window size when network is uncongested, but aggressively reducing it upon congestion. Inspired by AIMD, HistoSpec designs a dynamic speculation window for each response, initialized to two tokens. When all speculated tokens are accepted, HistoSpec additively increases the window size by 2, until it meets the upper threshold (32 by default). But when any token is rejected, it resets the size directly to 2. This approach guarantees high computational density during generating long matching sequences, and minimal wasted computation for short sequences.

As for the prefix length, HistoSpec sets it to 7 initially. If no matching suffix is found, HistoSpec progressively reduces it until it reaches 3. These hyperparameters are adjustable.

HistoSpec also considers *batch size*. At large batch sizes, increased decoding parallelism yields higher GPU computational utilization, where excessively low speculative acceptance rates degrades the overall throughput. To mitigate this, HistoSpec monitors the acceptance rates and adaptively disables speculation. Through pre-profiling, HistoSpec gets the maximum viable batch size for throughput gains at different acceptance rates. When a rollout worker's current batch size exceeds the threshold, speculation is automatically disabled to preserve system efficiency.
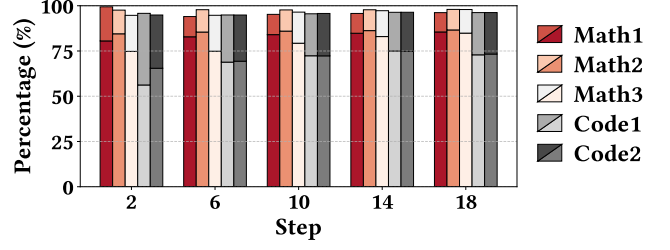
## 6 HistoPipe: Hybrid Rollout Pipeline

Although HistoSpec improves rollout efficiency, it does not address the issue of imbalanced rollout distributions, resulting in non-trivial GPU resource underutilization. We observe that leveraging the information from historical rollouts can effectively maintain load balance throughout the rollout process, and propose the HistoPipe scheduling design.

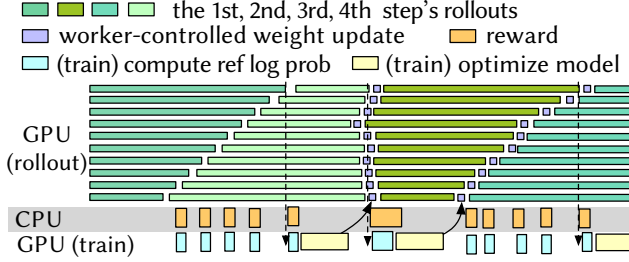### 6.1 Observation: Distribution Similarity in Rollout

Although different prompts generate varying numbers of tokens within a single rollout, we observe that the response length (i.e., token count) distribution for the same prompt across adjacent epochs remains similar. In other words, if a prompt generates a relatively large number of tokens in one rollout iteration, it is likely to produce a similarly large token count in subsequent rollouts. Therefore, *we can effectively predict the ranking of response lengths using historical lengths.*

To validate this observation, we analyze the response length rankings across multiple rollout epochs for five datasets (Fig. 6c). Specifically, we group the responses into 8 ranking groups based on their lengths, ranked from low to high. As Fig. 9 shows, across 20 epochs (300—1000 steps), for math tasks, an average of only 16% of responses change their group to higher ones (28% for code tasks). Among them, 13% of the (all) responses only shift near the group boundary without significantly altering the distribution (24% for code). These results demonstrate that not only do generated tokens exhibit similarity across rollouts, but the *distribution of rollout lengths also remains relatively consistent over epochs.*

### 6.2 Distribution-aware Hybrid Pipeline

Leveraging the similarity in rollout length distributions can enable more balanced scheduling of rollout tasks, reducing

the 1st, 2nd, 3rd, 4th step's rollouts
worker-controlled weight update    reward
(train) compute ref log prob    (train) optimize model

**Figure 10. HistoPipe design.** *Leveraging the **worker-controlled async weight updating** (§ 4.2), updated model weights are propagated asynchronously to the rollout workers' weight buffers and updated to GPUs by rollout workers upon completing a rollout step. In rare cases where the corresponding weights have not been propagated to the weight buffer, the rollout worker will await the weights.*

| Dataset name | Math-1 | | Math-2 | | Math-3 | | Code-1 | | Code-2 | |
|---|---|---|---|---|---|---|---|---|---|---|
| Rank | 8 | 16 | 8 | 16 | 8 | 16 | 8 | 16 | 8 | 16 |
| Accurate | 85.7 | 79.7 | 83.2 | 76.9 | 79.0 | 71.6 | 71.8 | 62.6 | 73.2 | 62.5 |
| Not last 10% | 12.2 | 16.8 | 12.7 | 17.5 | 16.9 | 22.8 | 23.9 | 30.4 | 22.9 | 30.1 |
| Within 1.1x | 0.61 | 1.03 | 0.55 | 1.02 | 0.86 | 1.47 | 1.44 | 2.34 | 1.46 | 2.32 |
| Migrated | **1.49** | **2.47** | **3.55** | **4.58** | **3.24** | **4.13** | **2.86** | **4.66** | **2.44** | **5.08** |

**Table 2. Accuracy of predicting the samples' ranks with historical lengths and migration rates (14B models).** *We list the average value of the metrics of the 20+ epochs.*

load imbalance during rollout. We propose a distribution-aware rollout pipeline strategy, named HistoPipe.

**Hybrid Pipeline.** Preserving the rollout-reward-train dependency is essential for algorithmic integrity. However, the imbalanced distribution of rollout lengths and the autoregressive nature of LLMs induce rollout bubbles within individual steps. Inspired by Dual-Pipe's philosophy [54], instead of seeking per-step balance, HistoPipe shifts focus to *inter-step workload complementarity*, which constructs synergistic balancing across consecutive training steps.

Historical distribution enables HistoPipe to rank rollout prompts based on their historical response lengths, obtaining several **ranking groups** by equally dividing the ranked prompts. As Fig. 10 shows, during odd-numbered rollout steps, the scheduler assigns ranking groups to workers (0 through N-1) in *ascending order* of rollout length. Conversely, during even-numbered steps, ranking groups are assigned in *descending order* of rollout length. By complementing rollout lengths in this alternating manner, we effectively fill idle times or bubbles that occur during rollout execution, thereby improving the overall efficiency of the rollout workers.

HistoPipe is not enabled during the first epoch as no historical information exists, which is acceptable because RL training typically spans 50—100 epochs, and the first epoch exhibits the shortest duration. For RL algorithms using Group Relative Advantage, HistoPipe uses the median lengths within each response group as the *historical response lengths* to rank the prompts. Although we do not preclude more complicated algorithms or model-based methods [14], the current method is sufficiently robust and effective (Fig. 9).

**Technical challenges.** HistoPipe also faces several challenges under real workloads. (1) Although the rollout length distributions exhibit high similarity, the occurrence of an anomalously long rollout can damage the complementarity and disrupt the pipelines. (2) In practice, the long-tail distribution of rollout lengths prevents consecutive steps from achieving perfect workload complementarity.

### 6.3 Migration-based Rebalancing

To mitigate the impact of anomalous rollout lengths on the hybrid pipeline, we employ two strategies. (1) *Intra-step migration.* Migrating excessively long rollouts to other groups that are still in the same step's rollout process. (2) *Inter-step migration.* Migrating the outliers to the next step.

More specifically, we set a threshold for each ranking group. When a rollout length exceeds the threshold, it will be migrated. For the first strategy, long rollouts are dynamically reassigned to other groups during execution. The generated tokens are preserved and the rollouts continue after Rhyme-RL recomputes the KV cache (i.e., prefill) of the prompt and the tokens. [2] For the second strategy, the migrated rollout is added to the next step. Similarly, the generated tokens are preserved and the rollouts continue after KV cache recomputation. Since the RL algorithms and systems inherently use oversampling, inter-step migrating a small number of rollouts and deferring their completion to the next step do not adversely affect training. In our scheduling design, anomalous rollouts in groups with short rollout lengths are preferentially intra-step migrated, whereas inter-step migration is used for anomalous rollouts in groups with generally long rollout lengths. In practice, intra-step migration efficiently handles most of the outliers.
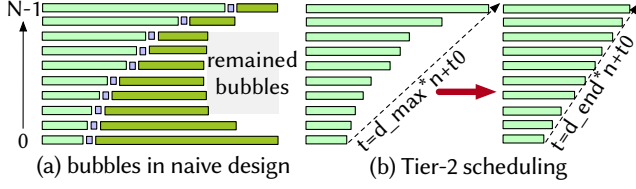
**Threshold.** RhymeRL migrates rollouts exhibiting: (1) within the last remaining $\alpha\%$ of the rollouts in the group; (2) generated response length exceeding $\beta$-times the maximum historical response length of the prompts in the same group. Currently, $\alpha$=10 and $\beta$ is determined by the 75th percentile of the growth rates in response lengths of the previous epoch ($\beta$=1.1 if the percentile < 1.1). Analysis across 5 datasets over 20 epochs reveals only 1.49%—3.55% (8 groups, on average) of responses undergo migration, demonstrating acceptable overhead and negligible impact on training accuracy.

### 6.4 Two-tier Scheduling

Achieving near-bubble-free inter-step complementarity requires approximately linear execution time distributions across rollout groups. In practice, however, long-tailed rollouts (Fig. 3b) cause the time distributions to approximate exponential patterns, which leads to bubbles on some rollout workers, as Fig. 11a shows. HistoPipe proposes a two-tier

---

[2]We do not migrate the KV cache because the overall overhead of KV cache migration exceeds that of recomputation.

**Figure 11. Bubbles caused by skewed rollout time distribution, and how the two-tier scheduling reduces bubbles.**

scheduling mechanism to tackle this issue:

• **Tier-1: from prompts to ranking groups.** First, HistoPipe equally divides the ranked prompts into ranking groups.

• **Tier-2: mapping ranking groups to GPUs.** If we allocate equal GPUs per ranking group, their execution times will exhibit exponential distributions due to long-tail rollouts. Therefore, HistoPipe designs a distribution reshaping strategy: instead of evenly distributing GPUs among the groups, it allocates fewer GPUs to the short-length and medium-length groups while provisioning extra GPUs to the few groups with long responses. In this way, HistoPipe reshapes the groups' rollout time distribution from *exponential to linear*, thereby further reducing rollout bubbles.

In RhymeRL, each rollout worker manages several GPUs cooperating via model parallelism with user-specified configurations. During Tier-2 scheduling, RhymeRL allocates GPUs at rollout-worker granularity. Rollout workers operate via data parallelism, thus, increasing rollout workers reduces per-worker batch size and shortens the execution time.

**Algorithm.** HistoPipe determines the GPU allocation plan using the algorithm detailed in Fig.12. It uses the mean value of the prompts' historical response lengths within each ranking group (*representative length*) to estimate the group's execution time. Through pre-profiling of the inference engine (i.e., HistoSpec, also considering the impact of speculative decoding), we can obtain the relationship between a ranking group's execution time *(t)*, its representative rollout length *(l)*, and DP worker number *(dp)*: $t=\tau(l, dp)$. Within the RL cluster, there are $wks$ rollout workers to serve $N$ ranking groups, whose representative lengths are $lens$. We have to assign varying numbers of rollout workers to each group, and make their completion times close to linear distribution. This constitutes an integer nonlinear programming problem and we converge to the solution via binary search.

## 7 Evaluation

We evaluate RhymeRL with LLMs of sizes ranging from 8B to 32B, on real-world math and code datasets.

**Experiments.** First, we compare RhymeRL's end-to-end training throughput with SOTA LLM RL systems on training LLMs of different sizes using different datasets (§7.1). Then, we break down the improvements of RhymeRL's designs (§7.2.1). To gain a deeper understanding of HistoSpec (§7.2.2), we study HistoSpec's improvements on rollout throughput

```
►MIN_WKS/MAX_WKS: a ranking
group' min/max worker number
►Tool func: calculate workers
needed for a given d
func calWks(d, lens, wks, t0)
  wks_needed = 0
  for i in range(0, N):
    target_t = t0+i*d
    ►find min workers to meet
    the group's target time
    group_wks = -1
    for k in range(MIN_WKS,
          MAX_WKS+1):
      exec_t = τ(lens[i], k)
      if exec_t <= target_t:
        group_wks = k
        break
    if group_wks == -1:
      return (Inf, [])
    wks_needed += group_wks
    plan[i] = group_wks
  return (wks_needed, plan)

►Find best worker allocation
►The target time of group n:
t(n)=d*n+t0, 0<=n<N
func planAllocation(lens,
    wks, t_train)
  t0 = max(τ(lens[0], MAX_WKS),
      t_train)
  d_min = 0.0
  d_max = (τ(lens[N-1],
        MIN_WKS)-t0)/(N-1)
  p = 1 # precision, adjustable
  ►binary search
  while (d_max-d_min) > p:
    d_mid = (d_max+d_min)/2
    result = calWks(d_mid,
          lens, wks, t0)
    if result.wks_needed > wks:
      d_min = d_mid
    else: # feasible solution
      best_plan = result.plan
      d_max = d_mid
  return best_plan
```

**Figure 12. Tier-2 scheduling algorithm.** *It solves the problem by minimizing the execution time gradient (d) via binary search, as Fig.11b shows. We use* t_train *(time of the last finished step's training stage) to ensure the execution time of the shortest ranking group is longer than the training stage's time, avoiding the next step's rollout waiting for train workers to generate weights.* $\tau(l, dp)$ *is determined by looking up the table obtained from pre-profiling. The **cooperation of HistoPipe with HistoSpec** is also considered.* $\tau(l, dp)$ *also considers HistoSpec's current speculation information (e.g., acceptance rate), which is omitted in the algorithm for ease of understanding.*

across steps, speculation rates, and acceptance rates. To further understand HistoPipe (§7.2.3), we study its improvements on training throughput across steps, and its migration rates. Finally, we study the impact of RhymeRL on training accuracy and algorithm behavior (§7.3). We also compare HistoSpec with other speculative decoding methods (§7.4).

**Hardware settings.** We deploy RhymeRL on a GPU cluster with 16 nodes and 128 GPUs. Each node is equipped with 8 high-performance GPUs, 2 Intel Xeon Sapphire Rapids CPUs with 96 cores, 1900GB host DRAM, and 9 * 400 Gb/s NVIDIA ConnectX7 InfiniBand NIC (8 GPU-dedicated NICs and one CPU-dedicated NIC).

**Model and algorithm.** We verify RhymeRL's observations and performance improvements on Qwen-2.5 models [55], Qwen-3 models [4], DeepSeek-R1-Distill-Qwen models [1] and LLama-3 models [56]. Our evaluation utilizes Qwen3-8B-Base, Qwen3-14B-Base, and Qwen2.5-32B models, with detailed architectures listed in Table.3. We train these models using GRPO algorithm [6], which is currently the mainstream LLM RL algorithm, powering advanced LLMs such as DeepSeek-R1 [1]. We also verify RhymeRL's efficiency on recent next-generation algorithms like DAPO [16] and provide the results. Based on our practice, the response group size is set to 16, which yields peak algorithmic efficiency.

**Metrics.** As per convention [14, 43, 57], we report training throughput, measured by the average number of samples generated and processed per second. We sample 8K math-
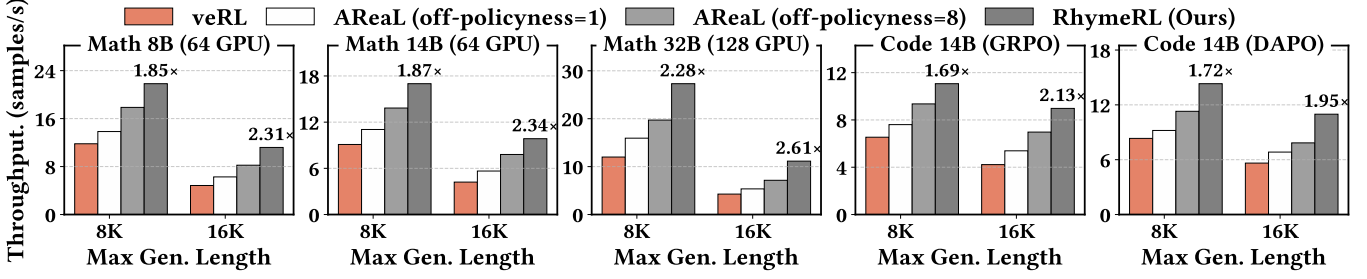
Figure 13. The training throughput of RhymeRL, veRL, and AReaL. *For each setting, we use the same configurations (e.g., number of rollout/train workers, clip ratio, etc.) for the three systems. We use 4-GPU tensor parallelism (TP) for 32B LLMs, and 2-GPU TP for 8B/14B LLMs.*

| Model arch | Size | Layers | Attn heads | K/V heads | Hidden size |
|---|---|---|---|---|---|
| Qwen-3 | 8B | 36 | 32 | 8 | 4096 |
| Qwen-3 | 14B | 40 | 40 | 8 | 5120 |
| Qwen-2.5 | 32B | 64 | 40 | 8 | 5120 |

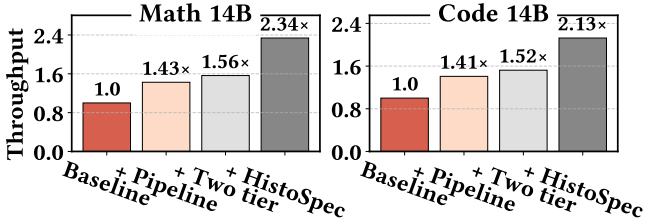Table 3. The specifications of the evaluated LLMs.



Figure 14. Breakdown of RhymeRL's improvements. *The max response length is set to 16K tokens. The throughputs are normalized.*

/code prompts from internal datasets (Fig. 6c). Under each setting, we train 80 steps for math models and 100 steps for code models, and use the first 20 steps as warm-up.

## 7.1 End-to-end Training Performance

We compare the end-to-end training performance of Rhyme-RL with the following state-of-the-art LLM RL training systems.

• **veRL [12]** (v0.4.1). Featuring the hierarchical hybrid programming model and highly-optimized 3D-HybridEngine, veRL is used by many projects and companies, and is Rhyme-RL's codebase. We use the rollout/train disaggregated architecture natively supported by veRL, which is the SOTA RL architecture and outperforms its hybrid mode [47].

• **AReaL [15]** (v0.3.0). Leveraging fully async rollout, AReaL achieves continuous rollout GPU utilization without idle time (§ 3.2). We evaluate AReaL's performance when its max off-policyness threshold is 1 (equal off-policyness with RhymeRL and veRL) and 8 (the max off-policyness recommended).

As Fig. 13 shows, RhymeRL outperforms veRL and AReaL on different model sizes (8B—32B), response lengths (8K/16K), tasks (Math/Code) and algorithms (DAPO/GRPO). Compared with veRL, RhymeRL improves the training throughput by up to 2.6x (1.9x on average for 8K max response length, and 2.3x on average for 16K max response length). This is primarily attributed to RhymeRL significantly reducing rollout
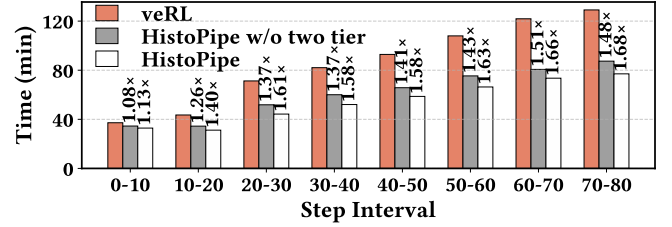


Figure 15. HistoPipe's improvements across steps. *Model: Math-14B. Max response length = 16K.*

time and effectively minimizing rollout bubbles. Compared with AReaL, when its off-policyness is 1, RhymeRL improves the training throughput by up to 2.1x (1.6x on average for 8K max response length, and 1.8x on average for 16K max response length). Although AReaL achieves continuous GPU utilization, it introduces rollout truncation and KV cache recomputation overhead. Leveraging historical similarity, RhymeRL achieves a more effective rollout pipeline. When AReaL's off-policyness is 8, RhymeRL improves the training throughput by up to 1.6x (1.3x on average for 8K max response length, and 1.4x on average for 16K max response length). RhymeRL outperforms AReaL (off-policyness = 8) as it significantly reduces rollout time through HistoSpec. Besides, RhymeRL does not change current RL paradigm.

## 7.2 Extended Studies

### 7.2.1 Ablation Study.
We break down the improvements brought by RhymeRL's designs. As Fig. 14 shows, for the Math-14B LLM, HistoPipe (§ 6.2) achieves 1.43x (1.41x for Code-14B) throughput boosts, with Two-tier Scheduling (§ 6.4) further improving the throughput of the naive hybrid pipeline by 1.10x (1.08x for Code-14B). HistoSpec (§ 5) further improves the training throughput by 1.50x (1.40x for Code-14B).

### 7.2.2 Detailed Analysis of HistoSpec.
**Improvement across steps.** We present HistoSpec's rollout throughput gains in each step at 8K max response length in Fig. 16. HistoSpec delivers up to 1.86x per-step rollout throughput gains, with progressive enhancement observed throughout training. This acceleration stems from: (1) increasing proportion of tokens generated by speculation and
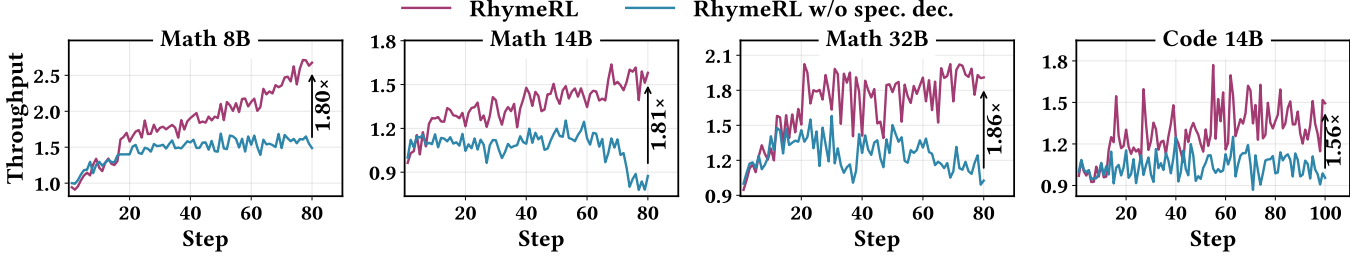
**Figure 16. HistoSpec's rollout throughput improvements.** *Max response length = 8K. The data are normalized to the baseline's first step.*
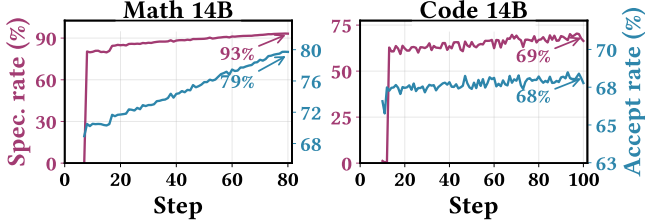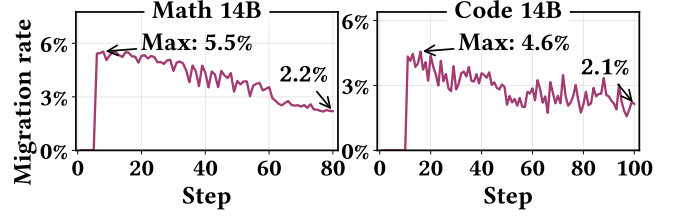


**Figure 17. Speculation rate and acceptance rate.**



**Figure 18. The proportion of samples migrated.**



**Figure 19. The reward scores of Math and Code 14B LLMs.**

acceptance rate, and (2) growing decoding dominance due to extending response lengths, which exacerbates memory-bound limitations.

**Speculation rate and acceptance rate.** As shown in Fig. 17, the proportion of tokens generated from speculative decoding (speculation rate) increases as training continues. The acceptance rate, i.e., the proportion of accepted tokens among the draft, remains 65%–79% and also increases as training continues. Utilizing the reward-aware tree-based history management (§ 5.3) and AIMD-like token speculation (§ 5.4), HistoSpec achieves high speculation rate and acceptance rate without wasting much computational resources.

### 7.2.3 Detailed Analysis of HistoPipe.

**Improvement across steps.** We evaluate the improvements of HistoPipe with/without Two-tier Scheduling across steps. As Fig. 15 shows, HistoPipe shortens the training time per 10 steps by up to 1.68x, and the Two-tier Scheduling (§ 6.4) accelerates the naive hybrid pipe by up to 1.14x. As the training progresses, the response length distribution becomes more stable, and the improvement becomes more significant.

**Migration Rate.** We quantify the proportion of rollout samples migrated (intra-step or inter-step, § 6.3) as outliers across steps. We use 16 ranking groups for math models and 8 ranking groups for code models. As Fig. 18 shows, for math tasks, 2.2%—5.5% of samples undergo migration (1.6%—4.6% for code tasks), with the outlier proportion progressively decreasing during training. Such minor migration maintains the efficiency of HistoPipe, while preserving algorithmic integrity (§ 7.3) and introducing ignoreable overhead.

### 7.3 Accuracy and Algorithmic Integrity

We present the reward scores during RL training. As Fig. 19 shows, for both GRPO and DAPO training, the curve of RhymeRL closely overlaps with that of veRL. RhymeRL main-
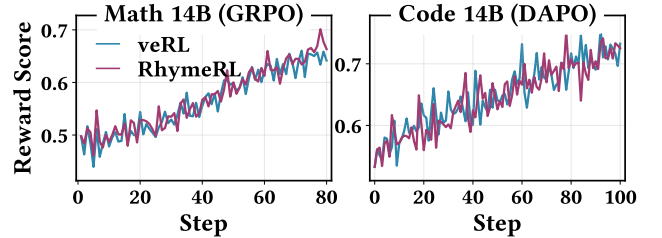
tains the overall training accuracy because: (1) Speculative decoding is theoretically proven to guarantee output equivalence. (2) The off-policyness of RhymeRL is strictly limited to one step, which is proven to maintain training accuracy [14, 17, 47]. (3) Leveraging historical distribution, RhymeRL only inter-step migrated a small portion of the rollout samples and continues their rollouts in the next step.

### 7.4 HistoSpec vs. Model-based Speculation

Although model-based speculation methods can also accelerate LLM inference, they face significant challenges in LLM RL rollout, making HistoSpec outperform them.

**Draft model adaptability.** Small LLMs cannot have peer reasoning capabilities as large LLMs [58]; therefore, SOTA methods utilize the target LLM's hidden states for draft generation [31–34]. However, they still face fundamental challenges in RL due to continuous model evolution. During RL, the LLMs are updated iteratively for thousands of steps. However, SOTA methods typically acquire draft models through distillation from the static model (Eagle1-3 [32–34]) or incremental fine-tuning of frozen base models (Medusa [31]), limiting their adaptability to the dynamically evolving LLMs. While concurrently training the draft model is possible, we have observed that, in practice, due to the effects of *algorithmic stochasticity and uncertainty*, such concurrently trained

draft models often fail to deliver high-accuracy drafts consistently throughout the entire training process.

**Performance.** HistoSpec outperforms SOTA model-based methods, stemming from three key advantages: (1) ultra-low draft generation overhead (hundreds of CPU cycles vs. millisecond-level costs in model-based speculation), (2) zero GPU computation or HBM footprint for drafting, and (3) consistently high acceptance rates. SOTA methods like Eagle3 [34] suffer throughput degradation beyond a batch size of 64 despite using a single-layer transformer (also reported by Eagle3 [34] and vLLM [59]). In contrast, HistoSpec sustains significant acceleration even at extreme batch sizes. At step 80, automatic oversampling leads to batch sizes of 4,928 rollouts/TP for Math-32B and 2,176 rollouts/TP for Math-8B, where HistoSpec still accelerates rollout by 1.80x—1.86x.

# 8 Conclusion

We present RhymeRL, an LLM RL system leveraging historical similarity for efficiency optimization. As the first RL system leveraging speculative decoding to shorten rollout time, RhymeRL utilizes historical rollout sequences as draft sources. Leveraging historical distribution, it proposes distribution-aware scheduling to reduce rollout bubbles. RhymeRL improves RL efficiency without compromising accuracy.

# References

[1] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren,

Zhangli Sha, Zhe Fu, Zhean Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. 2025. DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning. arXiv:2501.12948 [cs.CL] https://arxiv.org/abs/2501.12948

[2] Google. 2025. Gemini 2.5: Pushing the Frontier with Advanced Reasoning, Multimodality, Long Context, and Next Generation Agentic Capabilities. arXiv:2507.06261 [cs.CL] https://arxiv.org/abs/2507.06261

[3] 2025. The Llama 4 herd: The beginning of a new era of natively multimodal AI innovation. https://ai.meta.com/blog/llama-4-multimodal-intelligence/.

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

[5] 2025. Introducing Claude 4. https://www.anthropic.com/news/claude-4.

[6] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. arXiv:2402.03300 [cs.CL] https://arxiv.org/abs/2402.03300

[7] Yuxiang Zheng, Dayuan Fu, Xiangkun Hu, Xiaojie Cai, Lyumanshan Ye, Pengrui Lu, and Pengfei Liu. 2025. DeepResearcher: Scaling Deep Research via Reinforcement Learning in Real-world Environments. arXiv:2504.03160 [cs.AI] https://arxiv.org/abs/2504.03160

[8] Junde Wu, Jiayuan Zhu, Yuyuan Liu, Min Xu, and Yueming Jin. 2025. Agentic Reasoning: A Streamlined Framework for Enhancing LLM Reasoning with Agentic Tools. arXiv:2502.04644 [cs.AI] https://arxiv.org/abs/2502.04644

[9] Vignesh Prabhakar, Md Amirul Islam, Adam Atanas, Yao-Ting Wang, Joah Han, Aastha Jhunjhunwala, Rucha Apte, Robert Clark, Kang Xu, Zihan Wang, and Kai Liu. 2025. OmniScience: A Domain-Specialized LLM for Scientific Reasoning and Discovery. arXiv:2503.17604 [cs.AI] https://arxiv.org/abs/2503.17604

[10] Yifei Zhou, Song Jiang, Yuandong Tian, Jason Weston, Sergey Levine, Sainbayar Sukhbaatar, and Xian Li. 2025. SWEET-RL: Training Multi-Turn LLM Agents on Collaborative Reasoning Tasks. arXiv:2503.15478 [cs.LG] https://arxiv.org/abs/2503.15478

[11] Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. 2025. s1: Simple test-time scaling. arXiv:2501.19393 [cs.CL] https://arxiv.org/abs/2501.19393

[12] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2025. HybridFlow: A Flexible and Efficient RLHF Framework. In *Proceedings of the Twentieth European Conference on Computer Systems* (Rotterdam, Netherlands) *(EuroSys '25)*. Association for Computing Machinery, New York, NY, USA, 1279–1297. https://doi.org/10.1145/3689031.3696075

[13] Kimi Team, Yifan Bai, Yiping Bao, Guanduo Chen, Jiahao Chen, Ningxin Chen, Ruijue Chen, Yanru Chen, Yuankun Chen, Yutian Chen, Zhuofu Chen, Jialei Cui, Hao Ding, Mengnan Dong, Angang Du, Chenzhuang Du, Dikang Du, Yulun Du, Yu Fan, Yichen Feng, Kelin Fu, Bofei

Gao, Hongcheng Gao, Peizhong Gao, Tong Gao, Xinran Gu, Longyu Guan, Haiqing Guo, Jianhang Guo, Hao Hu, Xiaoru Hao, Tianhong He, Weiran He, Wenyang He, Chao Hong, Yangyang Hu, Zhenxing Hu, Weixiao Huang, Zhiqi Huang, Zihao Huang, Tao Jiang, Zhejun Jiang, Xinyi Jin, Yongsheng Kang, Guokun Lai, Cheng Li, Fang Li, Haoyang Li, Ming Li, Wentao Li, Yanhao Li, Yiwei Li, Zhaowei Li, Zheming Li, Hongzhan Lin, Xiaohan Lin, Zongyu Lin, Chengyin Liu, Chenyu Liu, Hongzhang Liu, Jingyuan Liu, Junqi Liu, Liang Liu, Shaowei Liu, T. Y. Liu, Tianwei Liu, Weizhou Liu, Yangyang Liu, Yibo Liu, Yiping Liu, Yue Liu, Zhengying Liu, Enzhe Lu, Lijun Lu, Shengling Ma, Xinyu Ma, Yingwei Ma, Shaoguang Mao, Jie Mei, Xin Men, Yibo Miao, Siyuan Pan, Yebo Peng, Ruoyu Qin, Bowen Qu, Zeyu Shang, Lidong Shi, Shengyuan Shi, Feifan Song, Jianlin Su, Zhengyuan Su, Xinjie Sun, Flood Sung, Heyi Tang, Jiawen Tao, Qifeng Teng, Chensi Wang, Dinglu Wang, Feng Wang, Haiming Wang, Jianzhou Wang, Jiaxing Wang, Jinhong Wang, Shengjie Wang, Shuyi Wang, Yao Wang, Yejie Wang, Yiqin Wang, Yuxin Wang, Yuzhi Wang, Zhaoji Wang, Zhengtao Wang, Zhexu Wang, Chu Wei, Qianqian Wei, Wenhao Wu, Xingzhe Wu, Yuxin Wu, Chenjun Xiao, Xiaotong Xie, Weimin Xiong, Boyu Xu, Jing Xu, Jinjing Xu, L. H. Xu, Lin Xu, Suting Xu, Weixin Xu, Xinran Xu, Yangchuan Xu, Ziyao Xu, Junjie Yan, Yuzi Yan, Xiaofei Yang, Ying Yang, Zhen Yang, Zhilin Yang, Zonghan Yang, Haotian Yao, Xingcheng Yao, Wenjie Ye, Zhuorui Ye, Bohong Yin, Longhui Yu, Enming Yuan, Hongbang Yuan, Mengjie Yuan, Haobing Zhan, Dehao Zhang, Hao Zhang, Wanlu Zhang, Xiaobin Zhang, Yangkun Zhang, Yizhi Zhang, Yongting Zhang, Yu Zhang, Yutao Zhang, Yutong Zhang, Zheng Zhang, Haotian Zhao, Yikai Zhao, Huabin Zheng, Shaojie Zheng, Jianren Zhou, Xinyu Zhou, Zaida Zhou, Zhen Zhu, Weiyu Zhuang, and Xinxing Zu. 2025. Kimi K2: Open Agentic Intelligence. arXiv:2507.20534 [cs.LG] https://arxiv.org/abs/2507.20534

[14] Yinmin Zhong, Zili Zhang, Xiaoniu Song, Hanpeng Hu, Chao Jin, Bingyang Wu, Nuo Chen, Yukun Chen, Yu Zhou, Changyi Wan, Hongyu Zhou, Yimin Jiang, Yibo Zhu, and Daxin Jiang. 2025. StreamRL: Scalable, Heterogeneous, and Elastic RL for LLMs with Disaggregated Stream Generation. arXiv:2504.15930 [cs.LG] https://arxiv.org/abs/2504.15930

[15] Wei Fu, Jiaxuan Gao, Xujie Shen, Chen Zhu, Zhiyu Mei, Chuyi He, Shusheng Xu, Guo Wei, Jun Mei, Jiashu Wang, Tongkai Yang, Binhang Yuan, and Yi Wu. 2025. AReaL: A Large-Scale Asynchronous Reinforcement Learning System for Language Reasoning. arXiv:2505.24298 [cs.LG] https://arxiv.org/abs/2505.24298

[16] Qiying Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Weinan Dai, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. 2025. DAPO: An Open-Source LLM Reinforcement Learning System at Scale. arXiv:2503.14476 [cs.LG] https://arxiv.org/abs/2503.14476

[17] Michael Luo, Sijun Tan, Justin Wong, Xiaoxiang Shi, William Y. Tang, Manan Roongta, Colin Cai, Jeffrey Luo, Li Erran Li, Raluca Ada Popa, and Ion Stoica. 2025. DeepScaleR: Surpassing O1-Preview with a 1.5B Model by Scaling RL. https://pretty-radio-b75.notion.site/DeepScaleR-Surpassing-O1-Preview-with-a-1-5B-Model-by-Scaling-RL-19681902c1468005bed8ca303013a4e2. Notion Blog.

[18] Chujie Zheng, Shixuan Liu, Mingze Li, Xiong-Hui Chen, Bowen Yu, Chang Gao, Kai Dang, Yuqiong Liu, Rui Men, An Yang, Jingren Zhou, and Junyang Lin. 2025. Group Sequence Policy Optimization. arXiv:2507.18071 [cs.LG] https://arxiv.org/abs/2507.18071

[19] Nai-Chieh Huang, Ping-Chun Hsieh, Kuo-Hao Ho, and I-Chen Wu. 2024. PPO-Clip Attains Global Optimality: Towards Deeper Understandings of Clipping. arXiv:2312.12065 [cs.LG] https://arxiv.org/abs/2312.12065

[20] Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *Proceedings of the 40th International Conference on Machine Learning* (Honolulu, Hawaii, USA) *(ICML'23)*. JMLR.org, Article 795, 13 pages.

[21] Y.R. Yang and S.S. Lam. 2000. General AIMD congestion control. In *Proceedings 2000 International Conference on Network Protocols*. 187–198. https://doi.org/10.1109/ICNP.2000.896303

[22] OpenAI, :, Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, Ilge Akkaya, Ilya Kostrikov, Ilya Sutskever, Irina Kofman, Jakub Pachocki, James Lennon, Jason Wei, Jean Harb, Jerry Twore, Jiacheng Feng, Jiahui Yu, Jiayi Weng, Jie Tang, Jieqi Yu, Joaquin Quiñonero Candela, Joe Palermo, Joel Parish, Johannes Heidecke, John Hallman, John Rizzo, Jonathan Gordon, Jonathan Uesato, Jonathan Ward, Joost Huizinga, Julie Wang, Kai Chen, Kai Xiao, Karan Singhal, Karina Nguyen, Karl Cobbe, Katy Shi, Kayla Wood, Kendra Rimbach, Keren Gu-Lemberg, Kevin Liu, Kevin Lu, Kevin Stone, Kevin Yu, Lama Ahmad, Lauren Yang, Leo Liu, Leon Maksin, Leyton Ho, Liam Fedus, Lilian Weng, Linden Li, Lindsay McCallum, Lindsey Held, Lorenz Kuhn, Lukas Kondraciuk, Lukasz Kaiser, Luke Metz, Madelaine Boyd, Maja Trebacz, Manas Joglekar, Mark Chen, Marko Tintor, Mason Meyer, Matt Jones, Matt Kaufer, Max Schwarzer, Meghan Shah, Mehmet Yatbaz, Melody Y. Guan, Mengyuan Xu, Mengyuan Yan, Mia Glaese, Mianna Chen, Michael Lampe, Michael Malek, Michele Wang, Michelle Fradin, Mike McClay, Mikhail Pavlov, Miles Wang, Mingxuan Wang, Mira Murati, Mo Bavarian, Mostafa Rohaninejad, Nat McAleese, Neil Chowdhury, Neil Chowdhury, Nick Ryder, Nikolas Tezak, Noam Brown, Ofir Nachum, Oleg Boiko, Oleg Murk, Olivia Watkins, Patrick Chao, Paul Ashbourne, Pavel Izmailov, Peter Zhokhov, Rachel Dias, Rahul Arora, Randall Lin, Rapha Gontijo Lopes, Raz Gaon, Reah Miyara, Reimar Leike, Renny Hwang, Rhythm Garg, Robin Brown, Roshan James, Rui Shu, Ryan Cheu, Ryan Greene, Saachi Jain, Sam Altman, Sam Toizer, Sam Toyer, Samuel Miserendino, Sandhini Agarwal, Santiago Hernandez, Sasha Baker, Scott McKinney, Scottie Yan, Shengjia Zhao, Shengli Hu, Shibani Santurkar, Shraman Ray Chaudhuri, Shuyuan Zhang, Siyuan Fu, Spencer Papay, Steph Lin, Suchir Balaji, Suvansh Sanjeev, Szymon Sidor, Tal Broda, Aidan Clark, Tao Wang, Taylor Gordon, Ted Sanders, Tejal Patwardhan, Thibault Sottiaux, Thomas Degry, Thomas Dimson, Tianhao Zheng, Timur Garipov, Tom Stasi, Trapit Bansal, Trevor Creech, Troy Peterson, Tyna Eloundou, Valerie Qi, Vineet Kosaraju, Vinnie Monaco, Vitchyr Pong, Vlad Fomenko, Weiyi Zheng, Wenda Zhou, Wes McCabe, Wojciech Zaremba, Yann Dubois, Yinghai Lu, Yining Chen, Young Cha, Yu Bai, Yuchen He, Yuchen Zhang,

Yunyun Wang, Zheng Shao, and Zhuohan Li. 2024. OpenAI o1 System Card. arXiv:2412.16720 [cs.AI] https://arxiv.org/abs/2412.16720

[23] 2025. Introducing OpenAI o3 and o4-mini. https://openai.com/index/introducing-o3-and-o4-mini/.

[24] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. arXiv:1707.06347 [cs.LG] https://arxiv.org/abs/1707.06347

[25] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. 2015. Trust Region Policy Optimization. In *Proceedings of the 32nd International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 37)*, Francis Bach and David Blei (Eds.). PMLR, Lille, France, 1889–1897. https://proceedings.mlr.press/v37/schulman15.html

[26] Xupeng Miao, Gabriele Oliaro, Zhihao Zhang, Xinhao Cheng, Zeyu Wang, Zhengxin Zhang, Rae Ying Yee Wong, Alan Zhu, Lijie Yang, Xiaoxiang Shi, Chunan Shi, Zhuoming Chen, Daiyaan Arfeen, Reyna Abhyankar, and Zhihao Jia. 2024. SpecInfer: Accelerating Large Language Model Serving with Tree-based Speculative Inference and Verification. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3* (La Jolla, CA, USA) *(ASPLOS '24)*. Association for Computing Machinery, New York, NY, USA, 932–949. https://doi.org/10.1145/3620666.3651335

[27] Zhihao Zhang, Alan Zhu, Lijie Yang, Yihua Xu, Lanting Li, Phitchaya Mangpo Phothilimthana, and Zhihao Jia. 2024. Accelerating Iterative Retrieval-augmented Language Model Serving with Speculation. In *Forty-first International Conference on Machine Learning.* https://openreview.net/forum?id=CDnv4vg02f

[28] Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023. Accelerating Large Language Model Decoding with Speculative Sampling. arXiv:2302.01318 [cs.CL] https://arxiv.org/abs/2302.01318

[29] Baohao Liao, Yuhui Xu, Hanze Dong, Junnan Li, Christof Monz, Silvio Savarese, Doyen Sahoo, and Caiming Xiong. 2025. Reward-Guided Speculative Decoding for Efficient LLM Reasoning. arXiv:2501.19324 [cs.CL] https://arxiv.org/abs/2501.19324

[30] Yingpeng Du, Tianjun Wei, Zhu Sun, and Jie Zhang. 2025. Reinforcement Speculative Decoding for Fast Ranking. arXiv:2505.20316 [cs.AI] https://arxiv.org/abs/2505.20316

[31] Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri Dao. 2024. Medusa: Simple LLM Inference Acceleration Framework with Multiple Decoding Heads. arXiv:2401.10774 [cs.LG] https://arxiv.org/abs/2401.10774

[32] Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2025. EAGLE: Speculative Sampling Requires Rethinking Feature Uncertainty. arXiv:2401.15077 [cs.LG] https://arxiv.org/abs/2401.15077

[33] Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2024. EAGLE-2: Faster Inference of Language Models with Dynamic Draft Trees. arXiv:2406.16858 [cs.CL] https://arxiv.org/abs/2406.16858

[34] Yuhui Li, Fangyun Wei, Chao Zhang, and Hongyang Zhang. 2025. EAGLE-3: Scaling up Inference Acceleration of Large Language Models via Training-Time Test. arXiv:2503.01840 [cs.CL] https://arxiv.org/abs/2503.01840

[35] Gabriele Oliaro, Zhihao Jia, Daniel Campos, and Aurick Qiao. 2025. SuffixDecoding: Extreme Speculative Decoding for Emerging AI Applications. arXiv:2411.04975 [cs.CL] https://arxiv.org/abs/2411.04975

[36] Apoorv Saxena. 2023. Prompt Lookup Decoding. https://github.com/apoorvumang/prompt-lookup-decoding/

[37] Nan Yang, Tao Ge, Liang Wang, Binxing Jiao, Daxin Jiang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. Inference with Reference: Lossless Acceleration of Large Language Models. arXiv:2304.04487 [cs.CL] https://arxiv.org/abs/2304.04487

[38] 2025. (vLLM official blog) How Speculative Decoding Boosts vLLM Performance by up to 2.8x. https://blog.vllm.ai/2024/10/17/spec-decode.html.

[39] Sheng-Chun Kao, Suvinay Subramanian, Gaurav Agrawal, Amir Yazdanbakhsh, and Tushar Krishna. 2023. FLAT: An Optimized Dataflow for Mitigating Attention Bottlenecks. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2* (Vancouver, BC, Canada) *(ASPLOS 2023)*. Association for Computing Machinery, New York, NY, USA, 295–310. https://doi.org/10.1145/3575693.3575747

[40] Zhiyu Mei, Wei Fu, Kaiwei Li, Guangju Wang, Huanchen Zhang, and Yi Wu. 2025. ReaL: Efficient RLHF Training of Large Language Models with Parameter Reallocation. In *Eighth Conference on Machine Learning and Systems.* https://openreview.net/forum?id=yLU1zRf95d

[41] Zhenyu Han, Ansheng You, Haibo Wang, Kui Luo, Guang Yang, Wenqi Shi, Menglong Chen, Sicheng Zhang, Zeshun Lan, Chunshi Deng, Huazhong Ji, Wenjie Liu, Yu Huang, Yixiang Zhang, Chenyi Pan, Jing Wang, Xin Huang, Chunsheng Li, and Jianping Wu. 2025. AsyncFlow: An Asynchronous Streaming RL Framework for Efficient LLM Post-Training. arXiv:2507.01663 [cs.LG] https://arxiv.org/abs/2507.01663

[42] Zhixin Wang, Tianyi Zhou, Liming Liu, Ao Li, Jiarui Hu, Dian Yang, Jinlong Hou, Siyuan Feng, Yuan Cheng, and Yuan Qi. 2025. DistFlow: A Fully Distributed RL Framework for Scalable and Efficient LLM Post-Training. arXiv:2507.13833 [cs.DC] https://arxiv.org/abs/2507.13833

[43] Yinmin Zhong, Zili Zhang, Bingyang Wu, Shengyu Liu, Yukun Chen, Changyi Wan, Hanpeng Hu, Lei Xia, Ranchen Ming, Yibo Zhu, and Xin Jin. 2025. Optimizing RLHF Training for Large Language Models with Stage Fusion. In *22nd USENIX Symposium on Networked Systems Design and Implementation (NSDI 25)*. USENIX Association, Philadelphia, PA, 489–503. https://www.usenix.org/conference/nsdi25/presentation/zhong

[44] 2025. NeMo: A scalable generative AI framework built for researchers and developers working on Large Language Models, Multimodal, and Speech AI. https://github.com/NVIDIA/NeMo.

[45] Jian Hu, Jason Klein Liu, Haotian Xu, and Wei Shen. 2025. REINFORCE++: An Efficient RLHF Algorithm with Robustness to Both Prompt and Reward Models. arXiv:2501.03262 [cs.CL] https://arxiv.org/abs/2501.03262

[46] Weixun Wang, Shaopan Xiong, Gengru Chen, Wei Gao, Sheng Guo, Yancheng He, Ju Huang, Jiaheng Liu, Zhendong Li, Xiaoyang Li, Zichen Liu, Haizhou Zhao, Dakai An, Lunxi Cao, Qiyang Cao, Wanxi Deng, Feilei Du, Yiliang Gu, Jiahe Li, Xiang Li, Mingjie Liu, Yijia Luo, Zihe Liu, Yadao Wang, Pei Wang, Tianyuan Wu, Yanan Wu, Yuheng Zhao, Shuaibing Zhao, Jin Yang, Siran Yang, Yingshui Tan, Huimin Yi, Yuchi Xu, Yujin Yuan, Xingyao Zhang, Lin Qu, Wenbo Su, Wei Wang, Jiamang Wang, and Bo Zheng. 2025. Reinforcement Learning Optimization for Large-Scale Learning: An Efficient and User-Friendly Scaling Library. arXiv:2506.06122 [cs.LG] https://arxiv.org/abs/2506.06122

[47] 2025. (veRL) One Step Off Policy Async Trainer. https://verl.readthedocs.io/en/latest/advance/one_step_off.html.

[48] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the 29th Symposium on Operating Systems Principles* (Koblenz, Germany) *(SOSP '23)*. Association for Computing Machinery, New York, NY, USA, 611–626. https://doi.org/10.1145/3600006.3613165

[49] Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E. Gonzalez, Clark Barrett, and Ying Sheng. 2024. SGLang: Efficient Execution of Structured Language Model Programs. In *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (Eds.), Vol. 37. Curran Associates, Inc., 62557–62583. https://proceedings.neurips.cc/paper_files/paper/2024/file/724be4472168f31ba1c9ac630f15dec8-Paper-Conference.pdf

[50] 2025. (vLLM documentation) vLLm Speculative Decoding. https://docs.vllm.ai/en/latest/features/spec_decode.html.

[51] Peter Weiner. 1973. Linear pattern matching algorithms. In *Proceedings of the 14th Annual Symposium on Switching and Automata Theory (Swat 1973) (SWAT '73)*. IEEE Computer Society, USA, 1–11. https://doi.org/10.1109/SWAT.1973.13

[52] E. Ukkonen. 1995. On-line construction of suffix trees. *Algorithmica* 14, 3 (Sept. 1995), 249–260. https://doi.org/10.1007/BF01206331

[53] 2024. cgroups - Linux control groups. http://man7.org/linux/man-pages/man7/cgroups.7.html.

[54] DeepSeek-AI. 2025. DeepSeek-V3 Technical Report. arXiv:2412.19437 [cs.CL] https://arxiv.org/abs/2412.19437

[55] Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. 2025. Qwen2.5 Technical Report. arXiv:2412.15115 [cs.CL] https://arxiv.org/abs/2412.15115

[56] AI @ Meta Llama Team. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] https://arxiv.org/abs/2407.21783

[57] Chang Hyun Park, Taekyung Heo, Jungi Jeong, and Jaehyuk Huh. 2017. Hybrid tlb coalescing: Improving tlb translation coverage under diverse fragmented memory allocations. In *Computer Architecture (ISCA), 2017 ACM/IEEE 44th Annual International Symposium on*. IEEE, 444–456.

[58] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling Laws for Neural Language Models. arXiv:2001.08361 [cs.LG] https://arxiv.org/abs/2001.08361

[59] 2025. (vLLM issue) vLLM Eagle performance is worse than expected. https://github.com/vllm-project/vllm/issues/9565.