Adaptive Curriculum Reinforcement Learning for Long-form Writing

Anonymous ACL submission

Abstract

Recent advances in Large Language Models (LLMs) have enabled strong performance in long-form writing, yet existing supervised finetuning (SFT) approaches suffer from limitations such as data saturation and restricted learning capacity bounded by teacher signals. In this work, we present an Adaptive Curriculum Reinforcement Learning (ACRL) framework to advance long-form writing capabilities beyond SFT. The framework consists of three key components: Margin-aware Data Se*lection* strategy that prioritizes samples with high learning potential, Pairwise Comparison Reward mechanism that enhances reward discriminability, and Dynamic Reference Scheduling approach, which plays a particularly critical role by adaptively adjusting task difficulty based on evolving model performance. Experiments on 7B-scale writer models show that our RL framework largely improves long-form writing performance over strong SFT baselines. Furthermore, we observe that models trained with long-output RL generalize surprisingly well to long-input reasoning tasks, potentially offering a promising perspective for rethinking long-context training.

1 Introduction

002

012

016

017

021

028

034

042

Recent years have witnessed the remarkable advance of Large Language Models (LLMs) (OpenAI, 2023; DeepSeek-AI et al., 2025; Zhao et al., 2023) to follow instructions and provide helpful responses. Among their impressive capabilities, long-form writing, which aims to generate long and high-quality articles, has drawn increasing attention (Wu et al., 2025a; Bai et al., 2024b; Wu et al., 2025b) due to its broad practical applications.

However, generating articles of both fulfilled long length and satisfactory quality is non-trivial for current LLMs. Previous research has identified several challenges to employ LLMs for long-form generation, including inherently limited output ceiling (Bai et al., 2024b; Tu et al., 2025) and performance degradation as output length grows (Wu et al., 2025b; Tu et al., 2025). To address these issues, recent efforts perform targeted Supervised Fine-Tuning (SFT) on LLMs to extend their output lengths, with long-generation datasets constructed by iterative agent pipelines (Bai et al., 2024b; Quan et al., 2024; Wu et al., 2025b) or instruction backtranslation (Pham et al., 2024; Wang et al., 2024). Though effective, these approaches introduce heavy burdens of dataset construction due to the broad coverage of writing tasks and potential copyright issues (Maini et al., 2024) when incorporating human-written texts. Furthermore, training LLMs to imitate the collected long-generation responses inherently imposes a capability upper bound determined by teacher models or human experts, which may cause data saturation and sample inefficiency. 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

077

078

079

081

Meanwhile, the recent success of outcome-based Reinforcement Learning (RL) method (DeepSeek-AI et al., 2025; Team et al., 2025; Yuan et al., 2025) in reasoning-intensive areas reveals a promising direction to advance model capabilities beyond supervised fine-tuning. Despite its potential, the practice of online outcome-based RL on long-form writing is relatively underexplored and therefore poses the following challenges:

- **Data Selection**: Data quality and difficulty play a critical role in eliciting model potential. However, the optimal approach for selecting data for RL in long-form writing tasks remains unclear.
- **Reward Design**: Rule-based outcome rewards (DeepSeek-AI et al., 2025) cannot be directly applied to generative writing tasks. Without ground-truth labels, constructing an effective reward mechanism for long-form writing poses a significant challenge.
- **Curriculum Scheduling**: Curriculum Learning (Bengio et al., 2009) is widely used to progressively improve model performance, but



Figure 1: Overall framework of Adaptive Curriculum Reinforcement Learning (ACRL). 1) *Margin-aware Data Selection*: prioritizes samples with high learning potential; 2) *Pairwise Comparison Reward*: provides more discriminative reward signals; 3) *Dynamic Reference Scheduling*: adaptively incentivizes the model to surpass progressively stronger references.

static scheduling fails to adapt to the model's evolving competence, thereby reducing training effectiveness.

To tackle these challenges, our work proposes an **Adaptive Curriculum Reinforcement Learning** (ACRL) framework tailored for long-form writing, as illustrated in Figure 1. Our framework begins with Margin-aware Data Selection strategy which leverages the quality differential between the policy model response and the highest-quality reference as a measure of *learning potential*, diverging from the conventional difficulty-prioritized selection approach. Considering the limited discriminative capacity of pointwise scoring, we construct a Pairwise Comparison Reward mechanism which challenges the policy model to generate responses of better quality than provided references to earn positive rewards. To facilitate progressive model enhancement, we propose a Dynamic Reference Scheduling approach that assigns each query a set of references with progressively increasing quality. The scheduling approach dynamically updates the references per sample when the evolving policy model surpasses the current reference during training. In this way, the dynamic curriculum adjusts sample-level task difficulty based on the current model performance, encouraging the model to consistently outperform a marginally superior reference. The motivation behind is also aligned with the insights from R1-like RL practices (Shi et al., 2025; Bae et al., 2025) that samples neither

097

100

102

103

106

107

109

110

111

112

113

too easy nor too difficult help to achieve the best learning efficiency.

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

To evaluate our ACRL framework, we conduct continuous reinforcement training on top of supervised fine-tuned writer models. The results indicate that our RL framework effectively boosts the long-form writing capability, advancing the SOTA performances of 7B-level writer models. Besides the improvement in long-form generation, we also observe an interesting generalization phenomenon: our RL-trained writer model (average input length < 1k) shows a surprising improvement in long-text reasoning tasks (*input length: 8k–2M*), in contrast to the performance degradation of the SFT-trained model. The results may suggest a novel perspective on long-context training that training on longoutput tasks may also enhance their reasoning abilities on long inputs, thereby offering training insights into the relationship between long-context understanding and generation.

In summary, the contributions of our work are:

- We propose an Adaptive Curriculum Reinforcement Learning framework for long-form writing, which integrates three key components: *Margin-aware Data Selection, Pairwise Comparison Reward*, and *Dynamic Reference Scheduling*.
- Particularly, we propose **Dynamic Reference Scheduling**, which adaptively adjusts samplelevel task difficulty based on the model's evolving performance. This dynamic curriculum en-

- courages the model to continually outperformprogressively stronger references.
 - Our RL-trained 7B-scale writer model achieves state-of-the-art performance, demonstrating the effectiveness of our framework. Furthermore, we observe inspiring **generalization** from *longoutput* generation to *long-input* reasoning, revealing a novel benefit of long-form RL training for long-context understanding.

2 Related Work

147

148

149

150

151

152

153

154

155

157

158

159

160

161

162

163

164

166

167

Training Methods for Long-form Writing. Recent efforts to advance long-form writing capabilities (Bai et al., 2024b; Wu et al., 2025b) mainly focuses on constructing long-generation post-training datasets for fine-tuning. Main approaches include teacher model distillation (Wu et al., 2025b), iterative agent pipelines for extended output (Bai et al., 2024b; Tu et al., 2025; Quan et al., 2024) and instruction back-translation (Pham et al., 2024; Wang et al., 2024). However, the application of online reinforcement learning methods (Schulman et al., 2017; Shao et al., 2024) are relatively underexplored, hindering further improvement.

Long-form Writing Evaluation. Long-form writ-168 ing (Wu et al., 2025a) requires LLMs to write long-169 form articles, posing challenges for evaluation due 170 to the lack of ground-truths. Researchers estab-171 lish writing benchmarks (Wu et al., 2025b; Que 172 et al., 2024), with proprietary models (Bai et al., 173 2024b; Paech, 2023; Liu et al., 2024) or fine-tuned 174 LLMs (Wu et al., 2025b; Ke et al., 2024) to serve as 175 judges. However, there exists several bias of includ-176 177 ing position bias and self-enhancement bias (Zheng et al., 2023), challenging the reliability of LLM-as-178 Judge evaluation methods. 179

Curriculum Learning. Reinforcement Learning 180 methods (Schulman et al., 2017; Shao et al., 2024; 181 DeepSeek-AI et al., 2025) have become a critical 182 step to elicit LLM capabilities. To boost efficiency, 183 Curriculum Learning (Bengio et al., 2009) has been widely adopted in RL practices (Team et al., 2025; 185 Xie et al., 2025; Wen et al., 2025), including static difficulty-based scheduling (Luo et al., 2025; Song 187 et al., 2025) and dynamic data selection (Bae et al., 189 2025; Shi et al., 2025). However, these methods use rule-based correctness as a measure for difficulty 190 and perform sample selection, which increases roll-191 outs and may cause imbalanced learning across samples. 193

3 Adaptive Curriculum RL

In this work, we propose **ACRL** (Adaptive Curriculum Reinforcement Learning), an adaptive reinforcement learning framework aimed at further improving long-form writing capabilities after instruction fine-tuning. The framework comprises three key components: *Margin-aware Data Selection* strategy, *Pairwise Comparison Reward* mechanism and *Dynamic Reference Scheduling* approach. By integrating outcome-based RL into long-form writing tasks, our approach improves model writing capabilities through more effective sample selection, reward design, and learning scheduling. We will describe the components in detail respectively.

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

3.1 Margin-aware Data Selection

Previous data selection approaches typically take question difficulty as a key criteria, measured by the accuracy of the policy model (Shi et al., 2025; Bae et al., 2025), simplistic indicators (Cheng et al., 2021; Yang et al., 2025) like solution step counts or simple heuristics grounded in human intuition (Hendrycks et al., 2021). However, we argue that questions favored by difficulty-prioritized data selection algorithm may not be the most suitable for effective reinforcement learning.

To validate this assumption, we propose *Marginaware Data Selection*, which uses the performance gap between the policy model and the highest-quality reference as a measure of *learning potential*. Our intuition is simple: a question suitable for learning is a question with sufficient room for performance improvement. Specifically, the procedure is detailed as follows.

Generation with Multiple LLMs. Instead of relying on a single model as the difficulty estimator (Shi et al., 2025; Bae et al., 2025), we leverage a set of competitive LLMs $C = \{\pi, M_1, M_2, ...\}$, including the policy model, to generate diverse candidate responses for each writing instruction.

Multi-dimensional Grading. Each generated response r_j from model $M_j \in C$ is graded using a multi-dimensional pointwise LLM-as-a-Judge approach (Liu et al., 2024; Wu et al., 2025b), with averaged quality score denoted as s_j per response. **Data Selection on Learning Potential.** To prioritize samples from which the policy model can benefit most, we define the *model-grounded learning potential* p as the quality gap between the best competitor and the policy model:

$$p = \max_{j \in \mathcal{C}, \ j \neq \pi} (s_j - s_\pi)$$
 24

where s_{π} is the score of the policy model's response. A higher *p* indicates greater headroom for improvement. To filter out noisy instructions, we first discard samples where all the competitors produce under-performing responses, as such instructions are often overly difficult or suffer from quality issues themselves. After filtering, we rank the remaining samples by their learning potential *p*, and retain the top-*k* examples to construct the training set.

245

247

248

251

253

257

261

262

263

265

266

269

270

271

273

274

277

278

283

287

291

3.2 Pairwise Comparison Reward Mechanism

Reward function is a critical component to guide policy optimization in RL practice. While rulebased outcome reward (DeepSeek-AI et al., 2025; Team et al., 2025) has been proven to be remarkably effective in eliciting long-CoT (Wei et al., 2022) reasoning in reasoning-intensive tasks, it can not be directly applied to long-form writing tasks due to the lack of ground-truths and its subjective nature, posing challenges to reward design.

Recent efforts utilize LLM-as-a-Judge (Zheng et al., 2023; Wu et al., 2025b) to measure the quality of model-generated responses, achieving high agreement with human judges. There exists two evaluation approaches including pointwise grading and pairwise comparison. Though widely adopted in writing evaluation due to its simplicity, pointwise grading exhibits limited discriminative capabilities and relatively high variance. On the contrary, pairwise comparison compares the response with a high-quality reference, capturing the subtle differences and potential direction of improvement. By providing more discriminative reward signals, pairwise grading incentivizes the policy model to generate better response and defeat high-quality references for positive rewards. Therefore, our reward design is as follows:

$$r_{\text{quality}}(\mathbf{x}) = \begin{cases} 1 & \text{if } \text{Judge}(\mathbf{ref}, \mathbf{x}) = \mathbf{x} \succ \mathbf{ref} \\ 0.5 & \text{if } \text{Judge}(\mathbf{ref}, \mathbf{x}) = \mathbf{x} \equiv \mathbf{ref} \\ 0 & \text{if } \text{Judge}(\mathbf{ref}, \mathbf{x}) = \mathbf{x} \prec \mathbf{ref} \end{cases}$$

where $r_{\text{quality}}(\mathbf{x})$ denotes the reward for a generated response \mathbf{x} ; ref represents the high-quality reference response; and $\text{Judge}(\text{ref}, \mathbf{x})$ is the evaluation function performed by the LLM-based judge to compare \mathbf{x} with ref.

Furthermore, LLM judges are known to exhibit position bias (Zheng et al., 2023) in pairwise comparisons, systematically favoring the first response. To impose additional learning pressure, we deliberately place the model-generated response in the



Figure 2: Sample-wise asynchronous learning schedule during training enabled by ACRL. Each line represents a sample, where an upward step indicates LLM surpassing its current reference and advancing to a better one.

second position, thereby introducing *positional disadvantage* in training. This avoids the need for position-swapped comparisons and halves the evaluation cost, while encouraging the model to generate stronger outputs from a less favorable position. 292

293

294

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

3.3 Dynamic Reference Scheduling

Curriculum Learning (Bengio et al., 2009) schedules progressive task difficulty for better learning efficiency. Previous efforts utilize offlinecalculated difficulty for scheduling (Shi et al., 2025; Song et al., 2025) or introducing additional rollouts during training for adaptive sample selection (Bae et al., 2025; Yu et al., 2025). Though effective in reasoning-centered RL, these methods suffer from either non-adaptive difficulty estimates or increased inference overhead.

Faced with the disadvantages of insufficient adaptivity of current curriculum scheduling, we propose a *Dynamic Reference Scheduling* approach that encourages the policy model to sequentially outperform references of ascending quality. With the algorithm detailed in Algorithm 1, our framework introduces a more competitive reference as the policy model beats the current one, enabling asynchronous per-sample difficulty updates and dynamic adaptivity with the evolving model capability.

Pre-training: Data Preparation. Given a set of writing instructions W, we first apply the Marginaware Data Selection strategy as elaborated in Sec 3.1, obtaining multiple competitive references $\mathcal{R} = \{r_{\pi}, r_1, r_2, ...\}$ and their corresponding LLM-judged quality scores $\mathcal{S} = \{s_{\pi}, s_1, s_2, ...\}$ for each instruction. The references are then sorted in ascending order of quality to produce a stage-

Algorithm 1 Dynamic Reference Scheduling for Long-form Writing

1: Pre-processing: For each instruction $w \in W$, apply Margin-aware Data Selection (Section 3.1) to obtain a stage-wise				
reference list $\mathcal{R}^{(w)} = \{r_{\pi}^{(w)}, r_1^{(w)}, r_2^{(w)}, \dots\}$ ordered by ascending quality.				
2: Input: Instruction set W; reference lists $\{\mathcal{R}^{(w)}\}_{w \in W}$; policy model π_{θ} ; RL update	ater \mathcal{A} (e.g., PPO); batch size B .			
3: Initialize reference pointer $t_w \leftarrow 1$ for all $w \in W$	▷ current reference index			
4: while training not finished do				
5: Sample batch $\mathcal{B} = \{w_k\}_{k=1}^B$ from W				
6: for all $w_k \in \mathcal{B}$ do				
7: $r_k \leftarrow \mathcal{R}^{(w_k)}[t_{w_k}]$	⊳ current reference			
8: Generate response $g_k \leftarrow \pi_{\theta}(w_k)$				
9: Compute reward $R_k \leftarrow \text{Judge}(r_k, g_k)$	ightarrow 1 (win), 0.5 (tie), 0 (loss)			
10: end for				
11: Update policy $\pi_{\theta} \leftarrow \mathcal{A}(\pi_{\theta}, \{(w_k, g_k, R_k)\}_{k=1}^B)$				
12: for all $w_k \in \mathcal{B}$ such that $R_k = 1$ do	▷ reference surpassed			
13: if $t_{w_k} < \mathcal{R}^{(w_k)} $ then				
14: $t_{w_k} \leftarrow t_{w_k} + 1$	▷ promote to next stronger reference			
15: end if				
16: end for				
17: end while				

351

353

327

wise reference list $\mathcal{R}_s = \{r_{q1}, r_{q2}, \dots\}$. To maintain sufficient positive feedback early in training, we deliberately include the response from the initial policy model π in the reference set, as the other reference-generation LLMs are generally larger in size and more competent.

In-training: Dynamic Scheduling. At the start of training, each instruction is initialized with the lowest-quality reference r_{q1} , which is comparable to the initial policy model's response. As the model evolves during training, the model gradually generates higher-quality responses during rollouts and receives positive rewards in some of the LLM-judged pairwise comparisons. Subsequently, the defeated references r_t are replaced with marginally stronger ones r_{t+1} while the undefeated references are retained, progressively increasing the challenge without overwhelming the model, in alignment with the model's evolving capability. This dynamic and adaptive reference update mechanism establishes an asynchronous learning schedule for each writing instruction and effectively incentivize the model to consistently perform better. As shown in Figure 2, our approach enables sample-wise asynchronous scheduling to dynamically adapt task difficulty to model capability.

4 Experiments

354To demonstrate the effectiveness of ACRL, we con-355duct experiments on writing-oriented fine-tuned356LLMs to see whether ACRL can further advance357long-form writing capabilities beyond supervised358fine-tuning.

4.1 Datasets

We use two carefully-constructed generative writing datasets primarily designed for supervised finetuning, including LongWriter training set (Bai et al., 2024b) and WritingBench training set (Wu et al., 2025b). As detailed in Section 3.1, we perform the Margin-aware Data Selection procedure on these two datasets respectively. Specifically, we first generate references for each writing instruction with the initial policy model and four competent larger-size LLMs to construct competitive references, including Qwen-Plus (Yang et al., 2024), GPT-40 (Hurst et al., 2024), Claude-3.7 (Anthropic Team, 2025) and Deepseek R1 (DeepSeek-AI et al., 2025). Then, we utilize a fine-tuned judge model (Wu et al., 2025b), which is optimized for evaluating long-form writing responses and reaches high agreement with human judges, to grade the responses in multiple dimensions. Finally, after the selection process, we obtain 1.5k chosen samples each dataset for further reinforcement learning. Each sample contains a writing instruction and references ordered by ascending quality.

359

360

361

362

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

379

380

381

382

384

386

387

388

389

390

391

392

4.2 Training Setup

To fully realize the effectiveness of reinforcement learning, we use two writing-expert LLMs as the base models for RL, which are primarily fine-tuned with the full WritingBench training set, denoted as *Qwen2.5-7B-WritingBench-SFT* and *Llama3.1-8B-WritingBench-SFT* respectively.

With the proposed ACRL, we use the PPO algorithm (Schulman et al., 2017) to optimize the two selected based models for long-form writing. During the training process, we adopt Qwen-Plus

Model	Writing-Oriented Training		Long-form Writing Evaluation			
	SFT	RL	WritingBench	EQ-Bench	LongBench-Write	Average
Qwen-Plus	-	_	77.62	76.78	95.42	83.27
GPT-4o	-	-	83.42	80.45	92.92	85.60
Suri-7B	1	×	49.70	18.44	33.44	33.86
Longwriter-9B	1	DPO	79.10	44.15	80.83	68.03
Qwen2.5-7B-Instruct	X	×	73.26	49.59	85.03	69.29
Qwen2.5-7B-WritingBench-SFT (12k)	1	X	83.71	70.02	92.22	81.98
Qwen2.5-7B-WritingBench-SFT (24k)	1	X	83.71	69.55	92.57	81.94
Qwen2.5-7B-Writing-RL (Ours)	1	PPO	87.23	73.19	93.06	84.49
Llama3.1-8B-Instruct	X	×	66.40	48.40	73.89	62.89
Llama3.1-8B-WritingBench-SFT	1	X	83.98	78.11	90.66	84.25
Llama3.1-8B-Writing-RL (Ours)	1	PPO	87.10	82.73	92.36	87.40

Table 1: Evaluation results of the models trained with ACRL, with the highest score in each model family **bold**. Notably, ACRL-trained models perform the best within their model family, on par with the proprietary models.

to serve as pairwise-comparison judge, providing rewards for policy optimization. The resulting models are denoted as *Qwen2.5-7B-Writing-RL* and *Llama3.1-8B-Writing-RL* respectively. More implementation details and training parameters can be found in Appendix A.

4.3 Benchmarks and Baselines

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

To comprehensively evaluate long-form writing capabilities of LLMs, we use three established benchmarks including WritingBench (Wu et al., 2025b), LongBench-Write (Bai et al., 2024b), and EQ-Bench creative writing split (Paech, 2023). The benchmarks are of broad coverage and use strong judge LLMs to evaluate the quality of generated responses. Note that the judge LLMs adopted for evaluation are diverse and different from the rewarding judge LLM used in training, mitigating the risk of overfitting particular judge preferences to ensure a fair evaluation.

Our selected baselines include strong proprietary models (Yang et al., 2024; Hurst et al., 2024), instruction fine-tuned LLMs (Yang et al., 2024; Dubey et al., 2024), and writing-oriented fine-tuned LLMs (Wu et al., 2025b; Bai et al., 2024b; Pham et al., 2024). More evaluation details can be found in Appendix B.

4.4 Results

As detailed in Table 1, the evaluation results 420 demonstrate that models trained with ACRL out-421 perform other models across all the three bench-422 423 marks. Specifically, Llama3.1-8B-Writing-RL (Ours) achieves the highest average score of 87.14, 424 with Qwen2.5-7B-Writing-RL (Ours) follows with 425 an average of 84.49, both showing strong perfor-426 mance in 10B-level. Notably, our trained models 427

exhibit long-form writing capabilities that match or even surpass those of proprietary models, positioning them as strong open-source alternatives for long-form generation tasks. 428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

Meanwhile, we observe distinct performance trends when applying RL and SFT to relatively strong models. Despite utilizing identically constructed datasets from the same expert model and agent pipeline, the fine-tuned model on 24k samples exhibits performance equivalent to, or slightly below, that of the variant trained with 12k samples. This observation potentially underscores the phenomenon of data saturation, where beyond a certain capability threshold, simply increasing data volume fails to enhance model performance. In contrast, models continuously trained by reinforcement learning, such as Llama3.1-8B-Writing-RL (Ours) compared to Llama3.1-8B-WritingBench-SFT within the same model family, demonstrate consistent performance improvements and thereby indicates the promising potential of RL to further advance model capabilities where SFT encounters limitations.

5 Generalization from Output to Input

To understand the influence on long-context capabilities of long-output RL, we adopt the challenging long-context reasoning benchmark LongBench v2 (Bai et al., 2024a) to evaluate long-input reasoning. Notably, as shown in Figure 3, the input lengths in LongBench v2 are substantially longer than those in our training set, mostly exceeding not only the input lengths but also the total input–output lengths.

As detailed in Table 2, our findings are inspiring. Beyond improved performance in long-form

Model	Writing	g-Oriented Training		Evaluation				
	SFT	RL	Easy	Hard	Short	Medium	Long	Overall
Qwen2.5-7B-Instruct	X	×	31.8	28.3	38.9	26.0	21.3	29.6
Qwen2.5-7B-WritingBench-SFT	1	X	27.6	27.7	35.0	25.1	20.4	27.6
Qwen2.5-7B-Writing-RL (Ours)	✓	PPO	35.8	29.3	42.1	25.7	26.5	31.8
Llama3.1-8B-Instruct	X	×	32.3	28.9	35.6	27.4	26.9	30.2
Llama3.1-8B-WritingBench-SFT	1	X	29.7	27.7	36.7	23.7	24.1	28.4
Llama3.1-8B-Writing-RL (Ours)	1	PPO	31.2	33.8	42.2	29.3	24.1	32.8

Table 2: Evaluation results of the models trained with ACRL on LongBench v2, demonstrating the generalization potential from long-output generation to long-input reasoning.



Figure 3: Length distribution of our *long-output* RL training dataset and *long-input* evaluation dataset.

generation, the writer models fine-tuned with our RL recipe also exhibit surprising generalization to long-context reasoning tasks with substantially longer inputs, while the SFT-trained counterparts show slight performance degradation in this regime. To further understand and utilize this interesting phenomenon, we give an intuitive explanation to the following research questions.

Why does long-output training generalize to longinput reasoning? Generating high-quality longform text inherently requires a deep and holistic understanding of the preceding context. Therefore, long-generation RL encourages LLMs to develop long-input understanding capabilities as a prerequisite for producing coherent long-form outputs.

Why does long-output RL generalize better than
SFT? SFT forces the model to imitate and memorize the behaviors of the training samples, while RL aligns model behavior with outcome-based objectives via reward signals. Therefore, by empowering the model to enhance its underlying capabilities, RL generalizes better. This observation is also consistent with recent findings in other domains (Chu et al., 2025; Shen et al., 2025).

How might these findings inform long-context training? The generalization from long-output generation to long-input reasoning may suggest a mutu-

Data Selection Strategy	Sample Num	Policy Model Initial Score	Learning Potential	WritingBench Score
Baseline (w/o RL)	-	-	-	83.71
Full (w/o Selection)	5k	84.20	3.64	85.64
Difficulty-prioritized	1.5k	77.61	8.18	86.40
Margin-aware (Ours)	1.5k	78.84	9.16	87.02

 Table 3: Comparison of different data selection strategies.

ally beneficial relationship between long-input and long-output training. Integrating both perspectives may lead to more effective long-context training strategies, and we leave the systematic exploration of this promising approach to future work. 490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

6 Discussion

6.1 Analysis on Data Selection Strategy

Our Margin-aware Data Selection strategy aims to prioritize training samples with greater room for improvement. Unlike prior work that employs single-model difficulty estimates (Shi et al., 2025; Bae et al., 2025), our method measures the *learning potential* of each sample using the performance gap between the policy model and other competent LLMs, thereby amplifying sample-wise *learning potential*.

To validate this approach, we conduct data selection experiments on WritingBench (Wu et al., 2025b) Hard training dataset, training *Qwen2.5-7B-WritingBench-SFT* model with high-quality references generated by *Qwen-plus* (Yang et al., 2024). We adopt WritingBench (Wu et al., 2025b) to benchmark writing capabilities due to its broad coverage and evaluation efficiency. As shown in Table 3, the results indicate that our strategy can boost learning efficiency by choosing samples with higher learning potential. Compared to difficulty-prioritized approaches, our selected samples are slightly less difficult—as reflected by higher initial score measured with the policy model—highlighting the effectiveness of using

7

488

489

521

525

531

533

534

535

536

537

538

539

540

541

542

545

549

550

551

552

558

562

learning potential rather than absolute difficulty for data selection.

6.2 Analysis on Reward Design 523

To provide effective rewards, we construct a reward mechanism based on pairwise comparison with high-quality references. To validate our reward design, we compare our reward mechanism with the widely-adopted pointwise grading method (Zheng et al., 2023; Liu et al., 2025), which utilizes Judge LLM to provide a scalar rating representing response quality. We follow the experiment setting in Section 6.1. The results shown in Table 4 demonstrate the superiority of our approach to provide more discriminative rewards, incentivizing the model to further advance writing capabilities.

6.3 Analysis on Reference Quality

Under the Pairwise Comparison Reward Mechanism, the quality of references directly influences the difficulty for the policy model to obtain positive rewards, thereby impacting training stability and final performance. To examine the effect of reference quality, we conduct training experiments using multiple static reference sets, each generated by a different LLM, as well as a combined set consisting of the highest-quality references selected from all candidates. Specifically, we also include a reference set generated by the initial policy model itself to serve as a baseline, denoted as Self-Generated.

Reward Strategy Score		Reference Quality	Score
Baseline (w/o RL) Pointwise Pairwise (Ours)	83.71 84.59 87.02	Self-Generated Qwen-Plus Deepseek R1 Best Reference	86.80 87.02 86.15 82.51

Table 4: Comparison of Table 5: Comparison of different reward designs.

different reference quality.

As shown in Table 5, the results demonstrate that reference quality plays a critical role in effective training. Specifically, when statically using relatively low-quality references (e.g., Self-Generated), the policy model initially receives sufficient positive rewards to improve but quickly saturates, achieving near-perfect win rates without further progress. In contrast, overly high-quality references (e.g., Best Reference) suffer from the sparsity of positive rewards early in training, thereby reducing learning efficiency and destabilizing optimization. These observations highlight a key limitation of static reference scheduling: it requires careful

Curriculum	WritingBench	EQ-Bench	LongBench-Write	Average
Baseline (w/o RL)	83.71	70.02	92.22	81.98
None	86.82	71.78	90.83	83.15
Static	87.32	72.73	91.56	83.87
Dynamic (Ours)	87.23	73.19	93.06	84.49

Table 6: Comparison of different curriculum scheduling approaches.

reference selection and fails to adapt to the evolving capability of the policy model during training.

563

564

565

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

586

587

588

589

590

591

592

593

594

595

597

598

599

600

601

602

603

604

605

6.4 Ablation on Curriculum Scheduling

Considering the importance of reference quality and the disadvantages of fixed references as discussed in Section 6.3, we propose Dynamic Reference Scheduling which encourages the model to surpass increasingly higher-quality references as the model evolves. To demonstrate the effectiveness of this scheduling approach, we ablate the scheduling methods in RL training, including mixed training without scheduling, static scheduling and our proposed dynamic scheduling. As shown in Table 3.3, the results demonstrate the effectiveness of our approach.

Given the importance of reference quality and the limitations of fixed references discussed in Section 6.3, we propose Dynamic Reference Scheduling, which encourages the model to progressively surpass higher-quality references as it evolves. To evaluate the effectiveness of this scheduling strategy, we conduct an ablation study comparing three RL training setups: mixed training without scheduling (*None*), static scheduling which partitions the training set into two subsets with references of different quality, and our proposed dynamic scheduling. As shown in Table 3.3, the results confirm the superiority of our approach. Furthermore, both static and dynamic scheduling outperform the nocurriculum baseline, demonstrating the effectiveness of incorporating curriculum into the RL training process.

7 Conclusion

In this work, we propose an Adaptive Curriculum Reinforcement Learning (ACRL) framework, which consists of Margin-aware Data Selection, Pairwise Comparison Reward and Dynamic Reference Scheduling. Our experiments demonstrate its effectiveness on enhancing long-form writing capabilities and the performance gain successfully generalizes from long-output generation to long-input reasoning, indicating a promising perspective for long-context training.

Limitations

606

631

637

638

643

647

648

653

- Here we discuss several limitations of this work. To scale up model size. While the performance gain by training 7B-scale writer models with ACRL 609 is relatively large, there remains considerable room 610 for exploration at larger model scales. Prior re-611 search has shown that the underlying capability of 612 the base model plays a crucial role in the effectiveness of RL (Gandhi et al., 2025). Therefore, 614 applying ACRL to stronger models may lead to 615 even greater performance improvements, as well 616 as more pronounced generalization effects from 617 long-output generation to long-input reasoning. 618
- 619To explore the zero phenomenon of RL. This620work demonstrates that reinforcement learning,621when applied to long-form generation, can elicit622strong performance gains and even induce general-623ization to long-input reasoning. While an intriguing624research direction is to investigate this phenomenon625from a more fundamental perspective by directly626applying RL to base models without prior super-627vised fine-tuning. Such a setup may offer clearer628insight into whether RL alone is sufficient to induce629strong long-form generation capabilities.

References

- Anthropic Team. 2025. Claude 3.7 sonnet system card. https://assets.anthropic. com/m/785e231869ea8b3b/original/ claude-3-7-sonnet-system-card.pdf.
 - Sanghwan Bae, Jiwoo Hong, Min Young Lee, Hanbyul Kim, JeongYeon Nam, and Donghyun Kwak.
 2025. Online difficulty filtering for reasoning oriented reinforcement learning. arXiv preprint arXiv:2504.03380.
 - Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024a.
 LongBench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. *CoRR*, abs/2412.15204.
 - Yushi Bai, Jiajie Zhang, Xin Lv, Linzhi Zheng, Siqi Zhu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024b. LongWriter: Unleashing 10,000+ word generation from long context LLMs. *CoRR*, abs/2408.07055.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009, volume 382 of ACM International Conference Proceeding Series, pages 41–48. ACM.

Yi Cheng, Siyao Li, Bang Liu, Ruihui Zhao, Sujian Li, Chenghua Lin, and Yefeng Zheng. 2021. Guiding the growth: Difficulty-controllable question generation through step-by-step rewriting. In *Proceedings* of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021, pages 5968–5978. Association for Computational Linguistics. 658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

- Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V. Le, Sergey Levine, and Yi Ma. 2025. SFT memorizes, RL generalizes: A comparative study of foundation model post-training. *CoRR*, abs/2501.17161.
- 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, and 81 others. 2025. DeepSeek-R1: Incentivizing reasoning capability in Ilms via reinforcement learning. *CoRR*, abs/2501.12948.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, and 82 others. 2024. The Llama 3 herd of models. *CoRR*, abs/2407.21783.
- Kanishk Gandhi, Ayush Chakravarthy, Anikait Singh, Nathan Lile, and Noah D. Goodman. 2025. Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective stars. *CoRR*, abs/2503.01307.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the MATH dataset. In Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual.
- Aaron Hurst, Adam Lerer, Adam P. Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, Aleksander Madry, Alex Baker-Whitcomb, Alex Beutel, Alex Borzunov, Alex Carney, Alex Chow, Alex Kirillov, Alex Nichol, Alex Paino, and 79 others. 2024. GPT-40 system card. *CoRR*, abs/2410.21276.
- Pei Ke, Bosi Wen, Andrew Feng, Xiao Liu, Xuanyu Lei, Jiale Cheng, Shengyuan Wang, Aohan Zeng, Yuxiao Dong, Hongning Wang, Jie Tang, and Minlie Huang. 2024. CritiqueLLM: Towards an informative critique generation model for evaluation of large language model generation. In *Proceedings of the* 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL

- 774
- 775 776 777 778 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824

825

826

827

828

772

2024, Bangkok, Thailand, August 11-16, 2024, pages 13034-13054. Association for Computational Linguistics.

716

718

719

720

721

722

724

725

731

735

740

741

742

743

745

746

747

748

751

752

754

755

761

767

- Xiao Liu, Xuanyu Lei, Shengyuan Wang, Yue Huang, Andrew Feng, Bosi Wen, Jiale Cheng, Pei Ke, Yifan Xu, Weng Lam Tam, Xiaohan Zhang, Lichao Sun, Xiaotao Gu, Hongning Wang, Jing Zhang, Minlie Huang, Yuxiao Dong, and Jie Tang. 2024. AlignBench: Benchmarking chinese alignment of large language models. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 11621-11640. Association for Computational Linguistics.
- Zijun Liu, Peiyi Wang, Runxin Xu, Shirong Ma, Chong Ruan, Peng Li, Yang Liu, and Yu Wu. 2025. Inference-time scaling for generalist reward modeling. Preprint, arXiv:2504.02495.
- 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. Notion Blog.
- Pratyush Maini, Skyler Seto, Richard He Bai, David Grangier, Yizhe Zhang, and Navdeep Jaitly. 2024. Rephrasing the web: A recipe for compute and dataefficient language modeling. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024, pages 14044-14072. Association for Computational Linguistics.
- OpenAI. 2023. GPT-4 technical report. CoRR. abs/2303.08774.
- Samuel J. Paech. 2023. EQ-Bench: An emotional intelligence benchmark for large language models. CoRR, abs/2312.06281.
- Chau Pham, Simeng Sun, and Mohit Iyyer. 2024. Suri: Multi-constraint instruction following in long-form text generation. In Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024, pages 1722-1753. Association for Computational Linguistics.
- Shanghaoran Quan, Tianyi Tang, Bowen Yu, An Yang, Dayiheng Liu, Bofei Gao, Jianhong Tu, Yichang Zhang, Jingren Zhou, and Junyang Lin. 2024. Language models can self-lengthen to generate long texts. CoRR, abs/2410.23933.
- Haoran Que, Feiyu Duan, Liqun He, Yutao Mou, Wangchunshu Zhou, Jiaheng Liu, Wenge Rong, Zekun Moore Wang, Jian Yang, Ge Zhang, Junran Peng, Zhaoxiang Zhang, Songyang Zhang, and Kai Chen. 2024. HelloBench: Evaluating long text generation capabilities of large language models. CoRR, abs/2409.16191.

- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. CoRR, abs/1707.06347.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y. K. Li, Y. Wu, and Daya Guo. 2024. DeepSeekMath: Pushing the limits of mathematical reasoning in open language models. CoRR, abs/2402.03300.
- Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qiangian Zhang, and 1 others. 2025. Vlm-r1: A stable and generalizable r1style large vision-language model. arXiv preprint arXiv:2504.07615.
- Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. 2024. HybridFlow: A flexible and efficient RLHF framework. arXiv preprint arXiv: 2409.19256.
- Taiwei Shi, Yiyang Wu, Linxin Song, Tianyi Zhou, and Jieyu Zhao. 2025. Efficient reinforcement finetuning via adaptive curriculum learning. arXiv preprint arXiv:2504.05520.
- Mingyang Song, Mao Zheng, Zheng Li, Wenjie Yang, Xuan Luo, Yue Pan, and Feng Zhang. 2025. FastCuRL: Curriculum reinforcement learning with progressive context extension for efficient training R1-like reasoning models. *CoRR*, abs/2503.17287.
- Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, and 75 others. 2025. Kimi k1.5: Scaling reinforcement learning with llms. CoRR, abs/2501.12599.
- Shangqing Tu, Yucheng Wang, Daniel Zhang-Li, Yushi Bai, Jifan Yu, Yuhao Wu, Lei Hou, Huigin Liu, Zhiyuan Liu, Bin Xu, and Juanzi Li. 2025. LongWriter-V: Enabling ultra-long and high-fidelity generation in vision-language models. CoRR. abs/2502.14834.
- Tiannan Wang, Jiamin Chen, Qingrui Jia, Shuai Wang, Ruoyu Fang, Huilin Wang, Zhaowei Gao, Chunzhao Xie, Chuou Xu, Jihong Dai, Yibin Liu, Jialong Wu, Shengwei Ding, Long Li, Zhiwei Huang, Xinle Deng, Teng Yu, Gangan Ma, Han Xiao, and 27 others. 2024. Weaver: Foundation models for creative writing. CoRR, abs/2401.17268.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022.

914

886

Liang Wen, Yunke Cai, Fenrui Xiao, Xin He, Qi An, Zhenyu Duan, Yimin Du, Junchen Liu, Lifu Tang, Xiaowei Lv, Haosheng Zou, Yongchao Deng, Shousheng Jia, and Xiangzheng Zhang. 2025. Light-R1: Curriculum SFT, DPO and RL for long COT from scratch and beyond. *CoRR*, abs/2503.10460.

829

830

837

838

839

840

841

844

845

851

852

853

854

856

857

864

868

870

871

873

874

875

876

878

879

881

- Yuhao Wu, Yushi Bai, Zhiqing Hu, Shanqqing Tu, Ming Shan Hee, Juanzi Li, and Roy Ka-Wei Lee. 2025a. Shifting long-context LLMs research from input to output. *CoRR*, abs/2503.04723.
- Yuning Wu, Jiahao Mei, Ming Yan, Chenliang Li, Shaopeng Lai, Yuran Ren, Zijia Wang, Ji Zhang, Mengyue Wu, Qin Jin, and Fei Huang. 2025b. WritingBench: A comprehensive benchmark for generative writing. *CoRR*, abs/2503.05244.
- Tian Xie, Zitian Gao, Qingnan Ren, Haoming Luo, Yuqian Hong, Bryan Dai, Joey Zhou, Kai Qiu, Zhirong Wu, and Chong Luo. 2025. Logic-RL: Unleashing LLM reasoning with rule-based reinforcement learning. *CoRR*, abs/2502.14768.
- 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, and 22 others. 2024. Qwen2.5 technical report. *CoRR*, abs/2412.15115.
- Shuxun Yang, Cunxiang Wang, Yidong Wang, Xiaotao Gu, Minlie Huang, and Jie Tang. 2025. Step-MathAgent: A step-wise agent for evaluating mathematical processes through tree-of-error. *CoRR*, abs/2503.10105.
- Qiying 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, Hang Zhu, and 16 others. 2025. DAPO: an open-source LLM reinforcement learning system at scale. *CoRR*, abs/2503.14476.
- Yufeng Yuan, Qiying Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaze Chen, Chengyi Wang, TianTian Fan, Zhengyin Du, Xiangpeng Wei, and 1 others. 2025. VAPO: Efficient and reliable reinforcement learning for advanced reasoning tasks. *arXiv preprint arXiv:2504.05118*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, and 3 others. 2023. A survey of large language models. *CoRR*, abs/2303.18223.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-Judge with MT-Bench and chatbot arena.

In Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023.

A Implementation and Training Settings

A.1 Implementation Details

In this section, we introduce the implementation details of our proposed RL framework.

Margin-aware Data Selection. We use several close-sourced LLMs to generate high-quality references for further training, including Qwenplus (Yang et al., 2024), GPT-40 (Hurst et al., 2024), Claude 3.7 (Anthropic Team, 2025) and Deepseek R1 (DeepSeek-AI et al., 2025). We set the inference temperature to 0.1 for balanced diversity and quality, and we remain other parameters to the default setting.

In our pointwise grading process, we utilize the state-of-the-art evaluation procedure proposed by WritingBench (Wu et al., 2025b), which includes generating sample-dependent evaluation criteria, then uses a fine-tuned LLM to grade the answers from multiple dimensions, finally averages the dimensional scores to give a scalar rating. We use Qwen-Plus (Yang et al., 2024) to generate the evaluation dimensions and we use the same evaluation prompt as WritingBench (Wu et al., 2025b) for the Judge Model.

Evaluation Prompt Template

Evaluate the Response based on the Query and criteria provided. ** Criteria **

"{criteria}"

** Query ** "`{query}"'

** Response **
'``{response}'''

Provide your evaluation based on the criteria:

"{criteria}"

Provide reasons for each score, indicating where and why any strengths or deficiencies occur within the Response. Reference specific passages or elements from the text to support your justification. Ensure that each reason is concrete, with explicit references to the text that aligns with the criteria requirements.

Scoring Range: Assign an integer score between 1 to 10

** Output format **

Return the results in the following JSON format, Only output this JSON format and nothing else:

"'json { {

"score": an integer score between 1 to 10, "reason": "Specific and detailed justification for the score using text elements." }} "

Pairwise Comparison Reward Mechanism.

We use the Qwen-Plus (Yang et al., 2024) model to judge the quality of the generated responses. The pairwise comparison prompts used in our experiment are adapted from (Zheng et al., 2023) and (Wu et al., 2025b).

For the training samples in LongWriter (Bai et al., 2024b) dataset, we use the original evaluation dimensions and the prompt is as follows.

Default Pairwise Comparison Prompt

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user's instructions and answers the user's question better. Your evaluation should consider factors such as the helpfulness, relevance, accuracy, depth, creativity, and level of detail of their responses. Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output

your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie. NOTE: If the response contains severe repetition or redundancy, it should be viewed as low quality score, losing the comparison.

User Question {question}

The Start of Assistant A's Answer {answer_a} The End of Assistant A's Answer

The Start of Assistant B's Answer {answer_b} The End of Assistant B's Answer

For the training samples in WritingBench (Wu et al., 2025b) training dataset, we use the generated criteria as the original paper recommends and the prompt is as follows.

Criteria Pairwise Comparison Prompt

Please act as an impartial judge and evaluate the quality of the responses provided by two AI assistants to the user question displayed below. You should choose the assistant that follows the user's instructions and answers the user's question better. Your evaluation should consider the following dimensions. criteria

Begin your evaluation by comparing the two responses and provide a short explanation. Avoid any position biases and ensure that the order in which the responses were presented does not influence your decision. Do not allow the length of the responses to influence your evaluation. Do not favor certain names of the assistants. Be as objective as possible. After providing your explanation, output your final verdict by strictly following this format: "[[A]]" if assistant A is better, "[[B]]" if assistant B is better, and "[[C]]" for a tie. NOTE: If the response contains severe repetition or redundancy, it should be viewed as low quality score, losing the comparison.

921

923

924

915

926

927

928

929

User Question {question}

The Start of Assistant A's Answer {answer_a} The End of Assistant A's Answer

The Start of Assistant B's Answer {answer_b} The End of Assistant B's Answer

A.2 Training Parameters

932

933

934

935

936

937

939

941

951

953

955

956

957

960

961

963

964

965

966

967

969

We display the key training parameters used in our training experiments. We adopt the effective reinforcement training framework VeRL (Sheng et al., 2024) to train our models. In our experiment, we use the proximal policy optimization (PPO) (Schulman et al., 2017) algorithm with generalized advantage estimation (GAE) as the advantage estimator. The training process is conducted using a batch size of 32 for training, with a maximum prompt length of 4096 tokens and response length capped at 10,000 tokens to accommodate long-form generation tasks. We enable the parameter/optimizer offloading via Fully Sharded Data Parallel (FSDP) to support efficient multi-GPU training and the training is conducted on 8x A100 GPUs. we use dynamic batch sizing and a low learning rate (1e-6) with a warm-up ratio of 0.4 to train the actor model, while the critic adopts a higher learning rate (1e-5) with a warm-up ratio of 0.05. We utilize a rollout strategy based on the vLLM engine with a tensor model parallel size of 2. The KL divergence penalty is set to a modest coefficient of 0.001. We train each model for about 400 steps and evaluate the checkpoints on the validation set each 50 steps.

B Benchmarks and Evaluation Methods

In this section, we introduce the benchmarks and evaluation prompt templates used in our experiments.

LongBench-Write LongBench-Write (Bai et al., 2024b) is designed to evaluate the LLM long-form generation abilities, which focuses on generating coherent outputs exceeding 10000 words, addressing challenges in maintaining consistency and quality over extended text. Key evaluation metrics include coherence, fluency and topic relevance. The evaluation prompt template used is as follows:

Evaluation Prompt Template

You are an expert in evaluating text quality. Please evaluate the quality of an AI assistant's response to a user's writing request. Be as strict as possible.

You need to evaluate across the following six dimensions, with scores ranging from 1 to 5. The scoring criteria from 5 to 1 for each dimension are as follows:

1. Relevance: From content highly relevant and fully applicable to the user's request to completely irrelevant or inapplicable.

2. Accuracy: From content completely accurate with no factual errors or misleading information to content with numerous errors and highly misleading.

3. Coherence: From clear structure with smooth logical connections to disorganized structure with no coherence.

4. Clarity: From clear language, rich in detail, and easy to understand to confusing expression with minimal details.

5. Breadth and Depth: From both broad and deep content with a lot of information to seriously lacking breadth and depth with minimal information.

6. Reading Experience: From excellent reading experience, engaging and easy to understand content to very poor reading experience, boring and hard to understand content.

Please evaluate the quality of the following response to a user's request according to the above requirements.

<User Request>

\$INST\$

</User Request>

<Response>

\$RESPONSE\$

</Response>

Please evaluate the quality of the response. You must first provide a brief analysis of its quality, then give a comprehensive analysis with scores for each dimension. The output must strictly follow the JSON format: "Analysis": ..., "Relevance": ..., "Accuracy": ..., "Coherence": ..., "Clarity": ..., "Breadth and Depth": ..., "Reading Experience": You do not need to consider whether the response meets the user's length

requirements in your evaluation. Ensure that only one integer between 1 and 5 is output for each dimension score.

WritingBench WritingBench (Wu et al., 2025b) is designed to evaluate the LLM long-form generation capabilities across six domains: creative, persuasive, informative, technical, business, and legal writing. It includes over 1200 tasks, further divided into 100 subdomains, with each task evaluated based on five dynamic criteria, such as coherence and relevance. The evaluation prompt template used is the same as the one used for pointwise grading in data selection process, detailed in Appendix A.1.

EQ-Bench EQ-Bench (Paech, 2023) is designed to evaluate the emotional intelligence (EQ) of LLMs, which focuses on assessing models' ability to predict and quantify emotions such as anger, surprise, empathy and confidence in dialogue-based cotexts. In this work, we use the creative writing split of EQ-Bench. The evaluation prompt template used is as follows:

Evaluation Prompt Template

You are an expert in assessing creative writing. Your task is to score the quality of the test model's response below, by several metrics, on a 0-20 scale. [PROMPT START] {writing_prompt} [PROMPT END] [TEST MODEL RESPONSE] {test_model_response} [TEST MODEL RESPONSE END] [Task] You are an expert in assessing creative writing. Your task is to score the quality of

the test model's response below, by several metrics, on a 0-20 scale.

Scoring notes:

- A scores of 20 represents a masterpiece.
- If the criteria is not relevant to the piece (e.g. criteria relating to dialogue but the piece has none), exclude it from your assessment.
- Do not add any commentary or explanation to the scores section.

- Do not penalize NSFW output; that is the default for some models. Just assess it on its merits.

- Everything within the "TEST MODEL RE-SPONSE" section was written by the test model. Sometimes models like to write comments on the piece after the piece is concluded; if this happens you should ignore their comments.

- In the output, write the metric names exactly as below so they can be parsed.

- Do not use markdown in your response. Use the designated output format exactly.

- You are to write a comprehensive analysis of the piece, then give your scores.

- For these criteria, lower is better: {lower_is_better_criteria}

- You are a critic, and your job is to be critical, especially of any failings or amateurish elements.

Output format is: [Analysis]
Write your detailed analysis. [Scores]
Metric 1 name: [Score 0-20]
Metric 2 name: ...

Now, rate the supplied model output on the following criteria: {creative_writing_criteria}

971

972

973

974

975

976

979

981

982

983