

IMPROVE MATHEMATICAL REASONING IN LANGUAGE MODELS WITH AUTOMATED PROCESS SUPERVISION

Anonymous authors

Paper under double-blind review

ABSTRACT

Complex multi-step reasoning tasks, such as solving mathematical problems or generating code, remain a significant hurdle for even the most advanced large language models (LLMs). Verifying LLM outputs with an Outcome Reward Model (ORM) is a standard inference-time technique aimed at enhancing the reasoning performance of LLMs. However, this still proves insufficient for reasoning tasks with a lengthy or multi-hop reasoning chain, where the intermediate outcomes are neither properly rewarded nor penalized. Process supervision addresses this limitation by assigning intermediate rewards during the reasoning process. To date, the methods used to collect process supervision data have relied on either human annotation or per-step Monte Carlo estimation, both prohibitively expensive to scale, thus hindering the broad application of this technique. In response to this challenge, we propose a novel divide-and-conquer style Monte Carlo Tree Search (MCTS) algorithm named *OmegaPRM* for the efficient collection of high-quality process supervision data. This algorithm swiftly identifies the first error in the Chain of Thought (CoT) with binary search and balances the positive and negative examples, thereby ensuring both efficiency and quality. As a result, we are able to collect over 1.5 million process supervision annotations to train Process Reward Models (PRMs). This fully automated process supervision alongside the weighted self-consistency algorithm is able to enhance LLMs' math reasoning performances. We improved the success rates of the instruction-tuned Gemini Pro model from 51% to 69.4% on MATH500 and from 86.4% to 93.6% on GSM8K. Similarly, we boosted the success rates of Gemma2 27B from 42.3% to 58.2% on MATH500 and from 74.0% to 92.2% on GSM8K. The entire process operates without any human intervention or supervision, making our method both financially and computationally cost-effective compared to existing methods.

1 INTRODUCTION

Despite the impressive advancements achieved by scaling Large Language Models (LLMs) on established benchmarks (Wei et al., 2022a), cultivating more sophisticated reasoning capabilities, particularly in domains like mathematical problem-solving and code generation, remains an active research area. Chain-of-thought (CoT) generation is crucial for these reasoning tasks, as it decomposes complex problems into intermediate steps, mirroring human reasoning processes. Prompting LLMs with CoT examples (Wei et al., 2022b) and fine-tuning them on question-CoT solution pairs (Perez et al., 2021; Ouyang et al., 2022) have proven effective, with the latter demonstrating superior performance. Furthermore, the advent of Reinforcement Learning with Human Feedback (RLHF; Ouyang et al., 2022) has enabled the alignment of LLM behaviors with human preferences through reward models, significantly enhancing model capabilities.

Beyond prompting and further training, developing effective decoding strategies is another crucial avenue for improvement. Self-consistency decoding (Wang et al., 2023) leverages multiple reasoning paths to arrive at a voted answer. Incorporating a verifier, such as an off-the-shelf LLM (Huang et al., 2022; Luo et al., 2023), can further guide LLMs in reasoning tasks by providing a feedback loop to verify final answers, identify errors, and suggest corrections. However, the gain of such approaches remains limited for complex multi-step reasoning problems. Reward models offer a promising alternative to verifiers, enabling the reranking of candidate outcomes based on reward signals to ensure higher accuracy. Two primary types of reward models have emerged: Outcome

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Reward Models (ORMs; Yu et al., 2024; Cobbe et al., 2021), which provide feedback only at the end of the problem-solving process, and Process Reward Models (PRMs; Li et al., 2023; Uesato et al., 2022; Lightman et al., 2023), which offer granular feedback at each reasoning step. PRMs have demonstrated superior effectiveness for complex reasoning tasks by providing such fine-grained supervision.

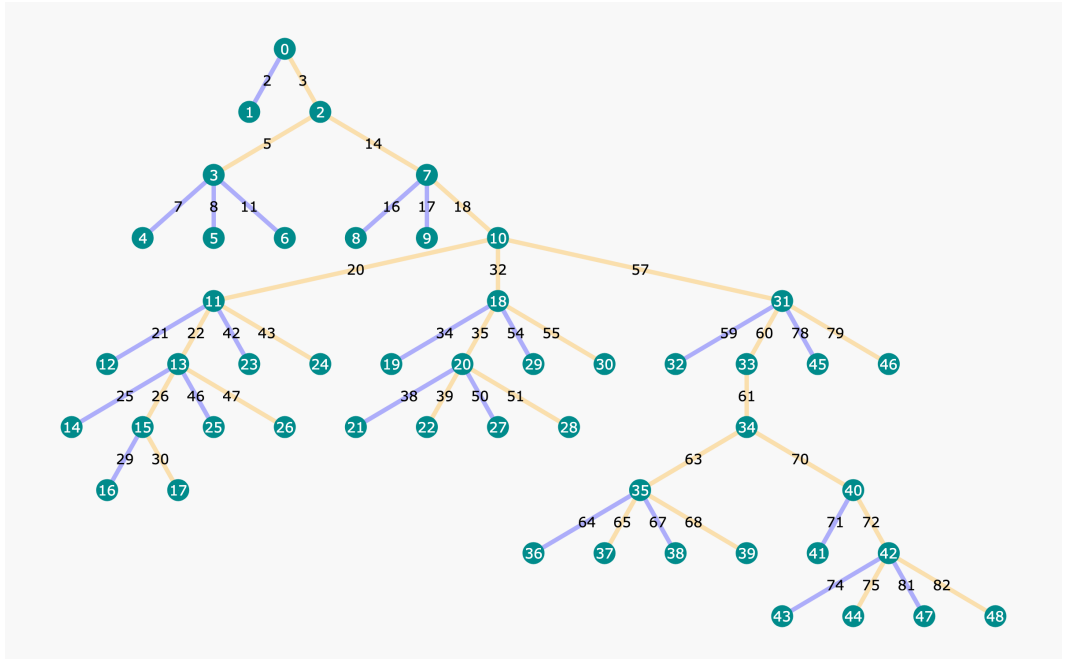


Figure 1: Example tree structure built with our proposed OmegaPRM algorithm. Each node in the tree indicates a state of partial chain-of-thought solution, with information including accuracy of rollouts and other statistics. Each edge indicates an action, *i.e.*, a reasoning step, from the last state. Yellow edges are correct steps and blue edges are wrong.

The primary bottleneck in developing PRMs lies in obtaining process supervision signals, which require supervised labels for each reasoning step. Current approaches rely heavily on costly and labor-intensive human annotation (Uesato et al., 2022; Lightman et al., 2023). Automating this process is crucial for scalability and efficiency. While recent efforts using per-step Monte Carlo estimation have shown promise (Wang et al., 2024a;b), their efficiency remains limited due to the vast search space. To address this challenge, we introduce OmegaPRM, a novel divide-and-conquer Monte Carlo Tree Search (MCTS) algorithm inspired by AlphaGo Zero (Silver et al., 2017) for automated process supervision data collection. For each question, we build a Monte Carlo Tree, as shown in Fig. 1, with the details explained in §3.3. This algorithm enables efficient collection of over 1.5 million high-quality process annotations without human intervention. Our PRM, trained on this dataset and combined with weighted self-consistency decoding, significantly improves the performance of instruction-tuned Gemini Pro from 51% to 69.4% on MATH500 (Lightman et al., 2023) and from 86.4% to 93.6% on GSM8K (Cobbe et al., 2021). We also boosted the success rates of Gemma2 27B from 42.3% to 58.2% on MATH500 and from 74.0% to 92.2% on GSM8K.

Our main contributions are as follows:

- We propose a novel divide-and-conquer style Monte Carlo Tree Search algorithm for automated process supervision data generation.
- The algorithm enables the efficient generation of over 1.5 million process supervision annotations, representing the largest and highest quality dataset of its kind to date. Additionally, the entire process operates without any human annotation, making our method both financially and computationally cost-effective.
- We combine our verifier with weighted self-consistency to further boost the performance of LLM reasoning. We significantly improves the success rates from 51% to 69.4% on

MATH500 and from 86.4% to 93.6% on GSM8K for instruction-tuned Gemini Pro. For Gemma2 27B, we also improved the success rates of from 42.3% to 58.2% on MATH500 and from 74.0% to 92.2% on GSM8K.

2 RELATED WORK

Improving mathematical reasoning ability of LLMs. Mathematical reasoning poses significant challenges for LLMs, and it is one of the key tasks for evaluating the reasoning ability of LLMs. With a huge amount of math problems in pretraining datasets, the pretrained LLMs (OpenAI, 2023; Gemini Team et al., 2024; Touvron et al., 2023) are able to solve simple problems, yet struggle with more complicated reasoning. To overcome that, the chain-of-thought (Wei et al., 2022b; Fu et al., 2023) type prompting algorithms were proposed. These techniques were effective in improving the performance of LLMs on reasoning tasks without modifying the model parameters. The performance was further improved by supervised fine-tuning (SFT; Cobbe et al., 2021; Liu et al., 2024; Yu et al., 2023) with high quality question-response pairs with full CoT reasoning steps.

Application of reward models in mathematical reasoning of LLMs. To further improve the LLM’s math reasoning performance, verifiers can help to rank and select the best answer when multiple rollouts are available. Several works (Huang et al., 2022; Luo et al., 2023) have shown that using LLM as verifier is not suitable for math reasoning. For trained verifiers, two types of reward models are commonly used: Outcome Reward Model (ORM) and Process Reward Model (PRM). Both have shown performance boost on math reasoning over self-consistency (Cobbe et al., 2021; Uesato et al., 2022; Lightman et al., 2023), yet evidence has shown that PRM outperforms ORM (Lightman et al., 2023; Wang et al., 2024a). Generating high quality process supervision data is the key for training PRM, besides expensive human annotation (Lightman et al., 2023), Math-Shepherd (Wang et al., 2024a) and MiPS (Wang et al., 2024b) explored Monte Carlo estimation to automate the data collection process with human involvement, and both observed large performance gain. Our work shared the essence with MiPS and Math-Shepherd, but we explore further in collecting the process data using MCTS.

Monte Carlo Tree Search (MCTS). MCTS (Świechowski et al., 2021) has been widely adopted in reinforcement learning (RL). AlphaGo (Silver et al., 2016) and AlphaGo Zero (Silver et al., 2017) were able to achieve great performance with MCTS and deep reinforcement learning. For LLMs, there are planning algorithms that fall in the category of tree search, such as Tree-of-Thought (Yao et al., 2023) and Reasoning-via-Planing (Hao et al., 2023). Recently, utilizing tree-like decoding to find the best output during the inference-time has become a hot topic to explore as well, multiple works (Feng et al., 2023; Ma et al., 2023; Zhang et al., 2024; Tian et al., 2024; Feng et al., 2024; Kang et al., 2024) have observed improvements in reasoning tasks.

3 METHODS

3.1 PROCESS SUPERVISION

Process supervision is a concept proposed to differentiate from outcome supervision. The reward models trained with these objectives are termed Process Reward Models (PRMs) and Outcome Reward Models (ORMs), respectively. In the ORM framework, given a query q (e.g., a mathematical problem) and its corresponding response x (e.g., a model-generated solution), an ORM is trained to predict the correctness of the final answer within the response. Formally, an ORM takes q and x and outputs the probability $p = \text{ORM}(q, x)$ that the final answer in the response is correct. With a training set of question-answer pairs available, an ORM can be trained by sampling outputs from a policy model (e.g., a pretrained or fine-tuned LLM) using the questions and obtaining the correctness labels by comparing these outputs with the golden answers.

In contrast, a PRM is trained to predict the correctness of each intermediate step x_t in the solution. Formally, $p_t = \text{PRM}([q, x_{1:t-1}], x_t)$, where $x_{1:i} = [x_1, \dots, x_i]$ represents the first i steps in the solution. This provides more precise and fine-grained feedback than ORMs, as it identifies the exact location of errors. Process supervision has also been shown to mitigate incorrect reasoning in the domain of mathematical problem solving. Despite these advantages, obtaining the intermediate signal

for each step’s correctness to train such a PRM is non-trivial. Previous work (Lightman et al., 2023) has relied on hiring domain experts to manually annotate the labels, which is and difficult to scale.

3.2 PROCESS ANNOTATION WITH MONTE CARLO METHOD

In two closely related works, Math-Shepherd (Wang et al., 2024a) and MiPS (Wang et al., 2024b), the authors propose an automatic annotation approach to obtain process supervision signals using the Monte Carlo method. Specifically, a “completer” policy is established that can take a question q and a prefix solution comprising the first t steps $x_{1:t}$ and output the completion — often referred to as a “rollout” in reinforcement learning — of the subsequent steps until the final answer is reached. As shown in Fig. 2(a), for any step of a solution, the completer policy can be used to randomly sample k rollouts from that step. The final answers of these rollouts are compared to the golden answer, providing k labels of answer correctness corresponding to the k rollouts. Subsequently, the ratio of correct rollouts to total rollouts from the t -th step, as represented in Eq. (1), estimates the “correctness level” of the prefix steps up to t . Regardless of false positives, $x_{1:t}$ should be considered correct as long as any of the rollouts is correct in the logical reasoning scenario.

$$c_t = \text{MonteCarlo}(q, x_{1:t}) = \frac{\text{num}(\text{correct rollouts from } t\text{-th step})}{\text{num}(\text{total rollouts from } t\text{-th step})} \quad (1)$$

Taking a step forward, a straightforward strategy to annotate the correctness of intermediate steps in a solution is to perform rollouts for every step from the beginning to the end, as done in both Math-Shepherd and MiPS. However, this brute-force approach requires a large number of policy calls. To optimize annotation efficiency, we propose a binary-search-based Monte Carlo estimation.

Monte Carlo estimation using binary search. As suggested by Lightman et al. (2023), supervising up to the first incorrect step in a solution is sufficient to train a PRM. Therefore, our objective is locating the first error in an efficient way. We achieve this by repeatedly dividing the solution and performing rollouts. Assuming no false positives or negatives, we start with a solution with potential errors and split it at the midpoint m . We then perform rollouts for $s_{1:m}$ with two possible outcomes: (1) $c_m > 0$, indicating that the first half of the solution is correct, as at least one correct answer can be rolled out from m -th step, and thus the error is in the second half; (2) $c_m = 0$, indicating the error is very likely in the first half, as none of the rollouts from m -th step is correct. This process narrows down the error location to either the first or second half of the solution. As shown in Fig. 2(b), by repeating this process on the erroneous half iteratively until the partial solution is sufficiently small (*i.e.*, short enough to be considered as a single step), we can locate the first error with a time complexity of $O(k \log M)$ rather than $O(kM)$ in the brute-force setting, where M is the total number of steps in the original solution.

3.3 MONTE CARLO TREE SEARCH

Although binary search improves the efficiency of locating the first error in a solution, we are still not fully utilizing policy calls as rollouts are simply discarded after stepwise Monte Carlo estimation. In practice, it is necessary to collect multiple PRM training examples (*a.k.a.*, triplets of question, partial solution and correctness label) for a question (Lightman et al., 2023; Wang et al., 2024a). Instead of starting from scratch each time, we can store all rollouts during the process and conduct binary searches from any of these rollouts whenever we need to collect a new example. This approach allows for triplets with the same solution prefix but different completions and error locations. Such reasoning structures can be represented as a tree, as described in previous work like Tree of Thought (Yao et al., 2023).

Formally, consider a *state-action tree* representing detailed reasoning paths for a question, where a state s contains the question and all preceding reasoning steps, and an action a is a potential subsequent step from a specific state. The root state is the question without any reasoning steps: $r_{\text{root}} = q$. The policy can be directly modeled by a language model as $\pi(a|s) = \text{LM}(a|s)$, and the state transition function is simply the concatenation of the preceding steps and the action step, *i.e.*, $s' = \text{Concatenate}(s, a)$.

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
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

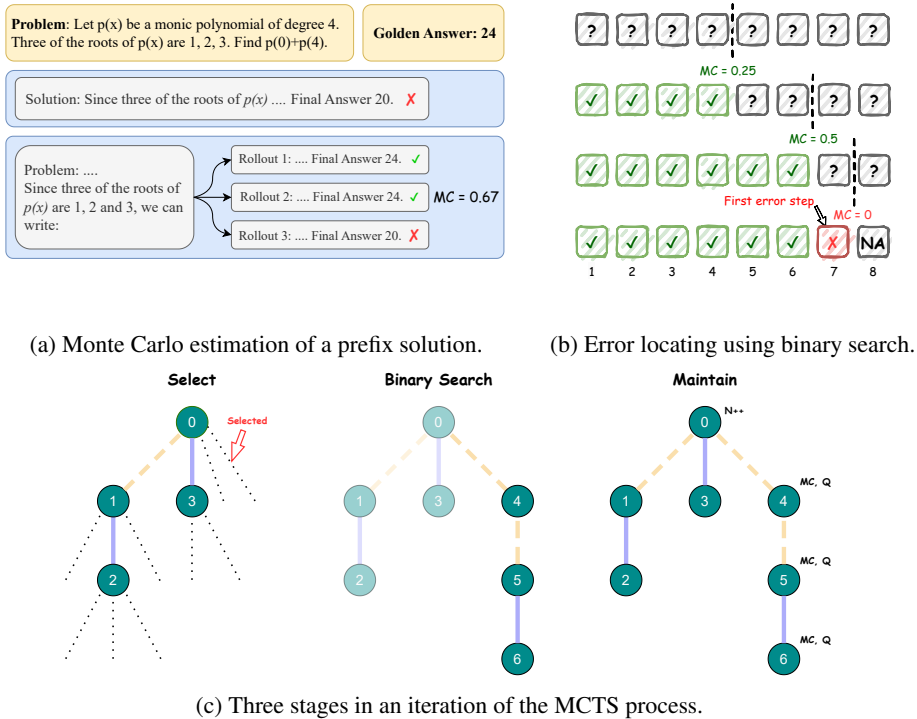


Figure 2: Illustration of the process supervision rollouts, Monte Carlo estimation using binary search and the MCTS process. (a) An example of Monte Carlo estimation of a prefix solution. Two out of the three rollouts are correct, producing the Monte Carlo estimation $MC(q, x_{1:t}) = 2/3 \approx 0.67$. (b) An example of error locating using binary search. The first error step is located at the 7th step after three divide-and-rollouts, where the rollout positions are indicated by the vertical dashed lines. (c) The MCTS process. The dotted lines in Select stage represent the available rollouts for binary search. The bold colored edges represent steps with correctness estimations. The yellow color indicates a correct step, *i.e.*, with a preceding state s that $MC(s) > 0$ and the blue color indicates an incorrect step, *i.e.*, with $MC(s) = 0$. The number of dashes in each colored edge indicates the number of steps.

Collecting PRM training examples for a question can now be formulated as constructing such a state-action tree. This reminds us the classic Monte Carlo Tree Search (MCTS) algorithm, which has been successful in many deep reinforcement learning applications (Silver et al., 2016; 2017). However, there are some key differences when using a language model as the policy. First, MCTS typically handles an environment with a finite action space, such as the game of Go, which has fewer than 361 possible actions per state (Silver et al., 2017). In contrast, an LM policy has an infinite action space, as it can generate an unlimited number of distinct actions (sequences of tokens) given a prompt. In practice, we use temperature sampling to generate a fix number of k completions for a prompt, treating the group of k actions as an approximate action space. Second, an LM policy can sample a full rollout until the termination state (*i.e.*, reaching the final answer) without too much overhead than generating a single step, enabling the possibility of binary search. Consequently, we propose an adaptation of the MCTS algorithm named **OmegaPRM**, primarily based on the one introduced in AlphaGo (Silver et al., 2016), but with modifications to better accommodate the scenario of PRM training data collection. We describe the algorithm details as below.

Tree Structure. Each node s in the tree contains the question q and prefix solution $x_{1:t}$, together with all previous rollouts $\{(s, r_i)\}_{i=1}^k$ from the state. Each edge (s, a) is either a single step or a sequence of consecutive steps from the node s . The nodes also store a set of statistics,

$$\{N(s), MC(s), Q(s, r)\},$$

where $N(s)$ denotes the visit count of a state, $\text{MC}(s)$ represents the Monte Carlo estimation of a state as specified in Eq. (1), and $Q(s, r)$ is a state-rollout value function that is correlated to the chance of selecting a rollout during the selection phase of tree traversal. Specifically,

$$Q(s, r) = \alpha^{1-\text{MC}(s)} \cdot \beta^{\frac{\text{len}(r)}{L}}, \quad (2)$$

where $\alpha, \beta \in (0, 1]$ and $L > 0$ are constant hyperparameters; while $\text{len}(r)$ denotes the length of a rollout in terms of number of tokens. Q is supposed to indicate how likely a rollout will be chosen for each iteration and our goal is to define a heuristic that selects the most valuable rollout to search with. The most straightforward strategy is uniformly choosing rollout candidates generated by the policy in previous rounds; however, this is obviously not an effective way. Lightman et al. (2023) suggests surfacing the *convincing wrong-answer* solutions for annotators during labeling. Inspired by this, we propose to prioritize *supposed-to-be-correct wrong-answer* rollouts during selection. We use the term *supposed-to-be-correct* to refer to the state with a Monte Carlo estimation $\text{MC}(s)$ closed to 1; and use *wrong-answer* to refer that the specific rollout r has a wrong final answer. The rollout contains mistakes made by the policy that should have been avoided given its high $\text{MC}(s)$. We expect a PRM that learns to detect errors in such rollouts will be more useful in correcting the mistakes made by the policy. The first component in Eq. (2), $\alpha^{1-\text{MC}(s)}$, has a larger value as $\text{MC}(s)$ is closer to 1. Additionally, we incorporate a length penalty factor $\beta^{\frac{\text{len}(r)}{L}}$, to penalize excessively long rollouts.

Select. The selection phase in our algorithm is simpler than that of AlphaGo (Silver et al., 2016), which involves selecting a sequence of actions from the root to a leaf node, forming a trajectory with multiple states and actions. In contrast, we maintain a pool of all rollouts $\{(s_i, r_j^i)\}$ from previous searches that satisfy $0 < \text{MC}(s_i) < 1$. During each selection, a rollout is popped and selected according to tree statistics, $(s, r) = \arg \max_{(s,r)} [Q(s, r) + U(s)]$, using a variant of the PUCT (Rosin, 2011) algorithm,

$$U(s) = c_{\text{puct}} \frac{\sqrt{\sum_i N(s_i)}}{1 + N(s)}, \quad (3)$$

where c_{puct} is a constant determining the level of exploration. This strategy initially favors rollouts with low visit counts but gradually shifts preference towards those with high rollout values.

Binary Search. We perform a binary search to identify the first error location in the selected rollout, as detailed in §3.2. The rollouts with $0 < \text{MC}(s) < 1$ during the process are added to the selection candidate pool. All divide-and-rollout positions before the first error become new states. For the example in Fig. 2(b), the trajectory $s[q] \rightarrow s[q, x_{1:4}] \rightarrow s[q, x_{1:6}] \rightarrow s[q, x_{1:7}]$ is added to the tree after the binary search. The edges $s[q] \rightarrow s[q, x_{1:4}]$ and $s[q, x_{1:4}] \rightarrow s[q, x_{1:6}]$ are correct, with MC values of 0.25 and 0.5, respectively; while the edge $s[q, x_{1:6}] \rightarrow s[q, x_{1:7}]$ is incorrect with MC value of 0.

Maintain. After the binary search, the tree statistics $N(s)$, $\text{MC}(s)$, and $Q(s, r)$ are updated. Specifically, $N(s)$ is incremented by 1 for the selected (s, r) . Both $\text{MC}(s)$ and $Q(s, r)$ are updated for the new rollouts sampled from the binary search. This phase resembles the *backup* phase in AlphaGo but is simpler, as it does not require recursive updates from the leaf to the root.

Tree Construction. By repeating the aboved process, we can construct a state-action tree as the example illustrated in Fig. 1. The construction ends either when the search count reaches a predetermined limit or when no additional rollout candidates are available in the pool.

3.4 PRM TRAINING

Each edge (s, a) with a single-step action in the constructed state-action tree can serve as a training example for the PRM. It can be trained using the standard classification loss

$$\mathcal{L}_{\text{pointwise}} = \sum_{i=1}^N \hat{y}_i \log y_i + (1 - \hat{y}_i) \log(1 - y_i), \quad (4)$$

where \hat{y}_i represents the correctness label and $y_i = \text{PRM}(s, a)$ is the prediction score of the PRM. Wang et al. (2024b) have used the Monte Carlo estimation as the correctness label, denoted as

324 $\hat{y} = \text{MC}(s)$. Alternatively, Wang et al. (2024a) have employed a binary labeling approach, where
 325 $\hat{y} = \mathbf{1}[\text{MC}(s) > 0]$, assigning $\hat{y} = 1$ for any positive Monte Carlo estimation and $\hat{y} = 0$ otherwise.
 326 We refer the former option as *pointwise soft* label and the latter as *pointwise hard* label. In addition,
 327 considering there are many cases where a common solution prefix has multiple single-step actions, we
 328 can also minimize the cross-entropy loss between the PRM predictions and the normalized pairwise
 329 preferences following the Bradley-Terry model (Christiano et al., 2017). We refer this training method
 330 as *pairwise* approach, and the detailed pairwise loss formula can be found in Section Appendix B.

331 We use the pointwise soft label when evaluating the main results in §4.1, and a comparison of the
 332 three objectives are discussed in §4.3.
 333
 334

335 4 EXPERIMENTS

336
 337 **Data Generation.** We conduct our experiments on the challenging MATH dataset (Hendrycks
 338 et al., 2021). We use the same training and testing split as described in Lightman et al. (2023),
 339 which consists of 12K training examples and a subset with 500 holdout representative problems from
 340 the original 5K testing examples introduced in Hendrycks et al. (2021). We observe similar policy
 341 performance on the full test set and the subset. For creating the process annotation data, we use the
 342 questions from the training split and set the search limit to 100 per question, resulting 1.5M per-step
 343 process supervision annotations. To reduce the false positive and false negative noise, we filtered out
 344 questions that are either too hard or too easy for the model. Please refer to Appendix A for details.
 345 We use $\alpha = 0.5$, $\beta = 0.9$ and $L = 500$ for calculating $Q(s, r)$ in Eq. (2); and $c_{\text{puct}} = 0.125$ in Eq. (3).
 346 We sample $k = 8$ rollouts for each Monte Carlo estimation.
 347

348 **Models.** In previous studies (Lightman et al., 2023; Wang et al., 2024a;b), both proprietary models
 349 such as GPT-4 (OpenAI, 2023) and open-source models such as Llama2 (Touvron et al., 2023) were
 350 explored. In our study, we perform experiments with both proprietary Gemini Pro (Gemini Team
 351 et al., 2024) and open-source Gemma2 (Gemma Team et al., 2024) models. For Gemini Pro, we
 352 follow Lightman et al. (2023); Wang et al. (2024a) to initially fine-tune it on math instruction data,
 353 achieving an accuracy of approximately 51% on the MATH test set. The instruction-tuned model is
 354 then used for solution sampling. For open-source models, to maximize reproducibility, we directly
 355 use the pretrained Gemma2 27B checkpoint with the 4-shot prompt introduced in Gemini Team et al.
 356 (2024). The reward models are all trained from the pretrained checkpoints.
 357

358 **Metrics and baselines.** We evaluate the PRM-based majority voting results on GSM8K (Cobbe
 359 et al., 2021) and MATH500 (Lightman et al., 2023) using PRMs trained on different process super-
 360 vision data. We choose the product of scores across all steps as the final solution score following
 361 Lightman et al. (2023), where the performance difference between product and minimum of scores
 362 was compared and the study showed the difference is minor. Baseline process supervision data
 363 include PRM800K (Lightman et al., 2023) and Math-Shepherd (Wang et al., 2024a), both publicly
 364 available. Additionally, we generate a process annotation dataset with our Gemini policy model using
 365 the brute-force approach described in Wang et al. (2024a;b), referred to as Math-Shepherd (our impl)
 366 in subsequent sections.
 367

368 4.1 MAIN RESULTS

369 Table 1 and Fig. 3 presents the performance comparison of PRMs trained on various process an-
 370 notation datasets. OmegaPRM consistently outperforms the other process supervision datasets.
 371 Specifically, the fine-tuned Gemini Pro achieves 69.4% and 93.6% accuracy on MATH500 and
 372 GSM8K, respectively, using OmegaPRM-weighted majority voting. For the pretrained Gemma2 27B,
 373 it also performs the best with 58.2% and 92.2% accuracy on MATH500 and GSM8K, respectively.
 374 It shows superior performance comparing to both human annotated PRM800K but also automatic
 375 annotated Math-Shepherd. More specifically, when the number of samples is small, almost all the
 376 PRM models outperforme the majority vote. However, as the number of samples increases, the
 377 performance of other PRMs gradually converges to the same level of the majority vote. In contrast,
 our PRM model continues to demonstrate a clear margin of accuracy.

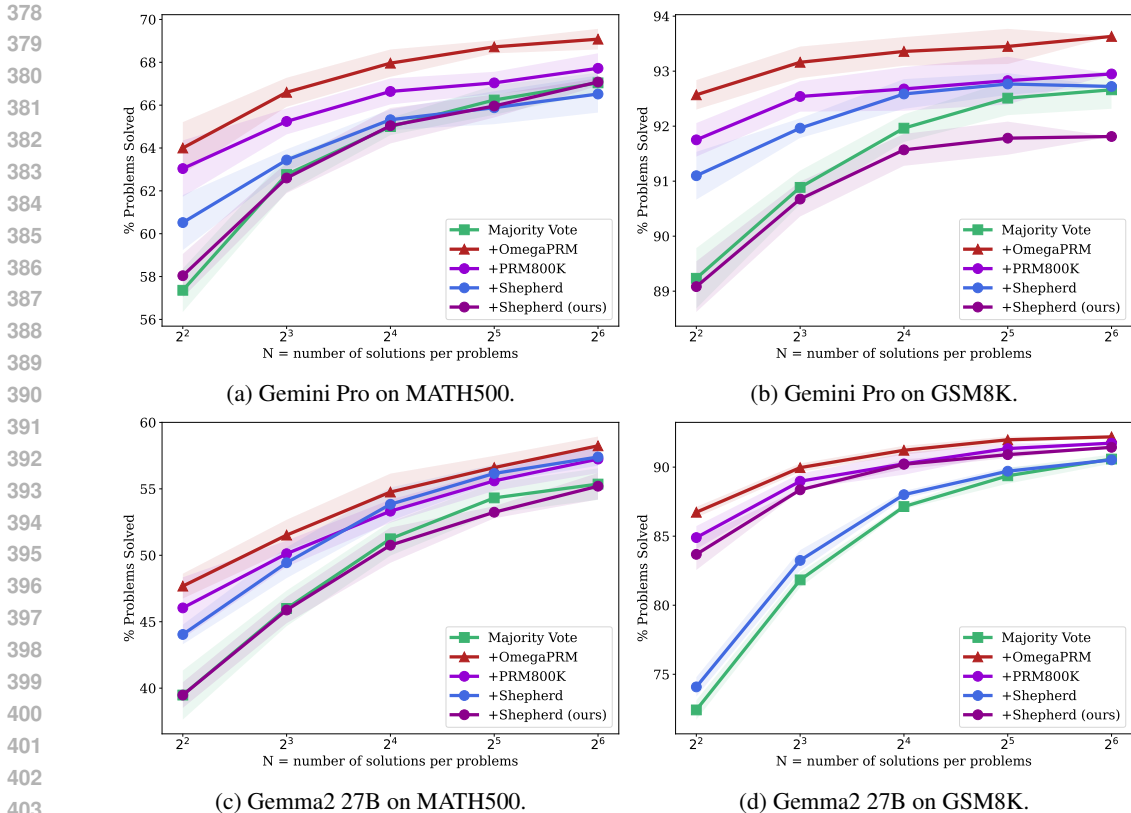


Figure 3: A comparison of PRMs trained with different process supervision datasets, evaluated by their ability to search over many test solutions using a PRM-weighted majority voting. We visualize the variance across many sub-samples of the 128 solutions we generated in total per problem.

Table 1: The performance comparison of PRMs trained with different process supervision datasets. The numbers represent the percentage of problems solved using PRM-weighted majority voting with $k = 64$.

	MATH500		GSM8K	
	Gemini Pro	Gemma 2 27B	Gemini Pro	Gemma 2 27B
MajorityVote@64	67.2	54.7	92.7	90.6
+ Math-Shepherd	67.2	57.4	92.7	90.5
+ Math-Shepherd (our impl)	67.2	55.2	91.8	91.4
+ PRM800K	67.6	57.2	92.9	91.7
+ OmegaPRM	69.4	58.2	93.6	92.2

4.2 STEP DISTRIBUTION

An important factor in process supervision is the number of steps in a solution and the length of each step. Previous works (Lightman et al., 2023; Wang et al., 2024a;b) use rule-based strategies to split a solution into steps, *e.g.*, using newline as delimiters. In contrast, we propose a more flexible method for step division, treating any sequence of consecutive tokens in a solution as a valid step. We observe that many step divisions in Math-Shepherd lack semantic coherence to some extent. Therefore, we hypothesize that semantically explicit cutting is not necessary for training a PRM.

In practice, we first examine the distribution of the number of steps per solution in PRM800K and Math-Shepherd, as shown in Fig. 4, noting that most solutions have less than 20 steps. During binary search, we aim to divide a full solution into 16 pieces. To calculate the expected step length, we

432
433
434
435
436
437
438
439
440
441
442

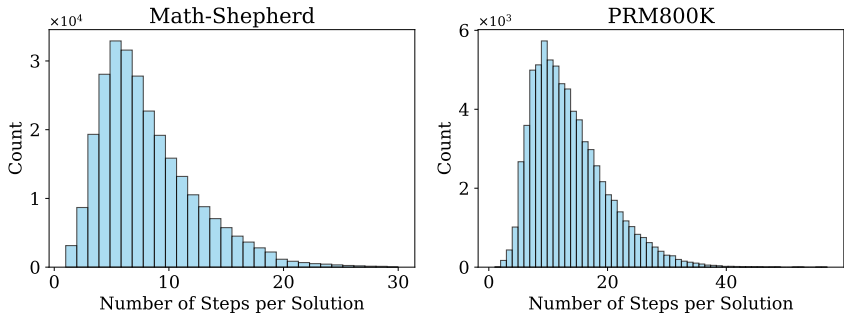


Figure 4: Number of steps distribution.

443
444
445
446
447
448
449

divide the average solution length by 16. The binary search terminates when a step is shorter than this value. The resulting distributions of step lengths for OmegaPRM and the other two datasets are shown in Fig. 5. This flexible splitting strategy produces a step length distribution similar to that of the rule-based strategy.

450
451
452
453
454
455
456
457
458

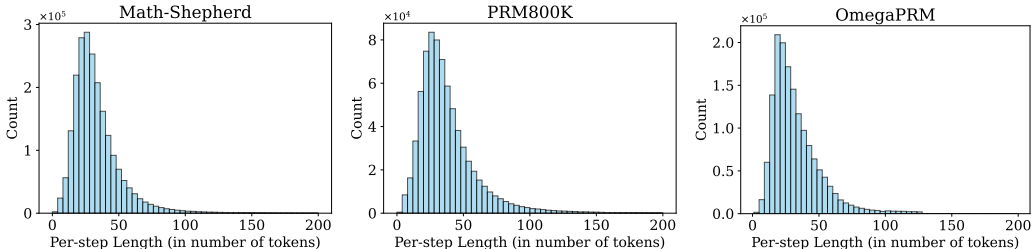


Figure 5: Step length distribution in terms of number of tokens.

459
460
461

4.3 PRM TRAINING OBJECTIVES

462
463
464
465
466
467
468
469
470

Table 2: Comparison of different training objectives for PRMs.

	Soft Label	Hard Label	Pairwise
PRM Accuracy (%)	70.1	63.3	64.2

471
472
473
474
475
476

As outlined in §3.4, PRMs can be trained using multiple objectives. We construct a small process supervision test set using the problems from the MATH test split. We train PRMs using pointwise soft label, pointwise hard label and pairwise loss respectively, and evaluate how accurately they can classify the per-step correctness. Table 2 presents the comparison of different objectives, and the pointwise soft label is the best among them with 70.1% accuracy.

477
478

4.4 ALGORITHM EFFICIENCY

479
480
481
482
483
484
485

As described in Section §3.2 and §3.3, we utilize binary search and Monte Carlo Tree Search to improve the efficiency of OmegaPRM process supervision data collection by effectively identifying the first incorrect step and reusing rollouts in Monte Carlo estimation. To quantitatively measure the efficiency of OmegaPRM, we collected process supervision data using both brute-force-style method (Wang et al., 2024a;b) and OmegaPRM with the same computational budget. As a result, we were able to generate 200K data points using the brute-force algorithm compared to 15 million data points with OmegaPRM, demonstrating a 75-times efficiency improvement. In practice, we randomly down-sampled OmegaPRM data to 1.5 million for PRM training.

5 LIMITATIONS

There are some limitations with our paper, which we reserve for future work:

Automatic process annotation is noisy. Our method for automatic process supervision annotation introduces noise in the form of false positives and negatives, but experiments indicate that it can still effectively train a PRM. The PRM trained on our dataset performs better than one trained on the human-annotated PRM800K dataset. The precise impact of noise on PRM performance remains uncertain. For future research, a comprehensive comparison of human and automated annotations should be conducted. One other idea is to integrate human and automated annotations, which could result in more robust and efficient process supervision.

Human supervision is still necessary. Unlike the work presented in AlphaGo Zero (Silver et al., 2017), our method requires the question and golden answer pair. The question is necessary for LLM to start the MCTS and the golden answer is inevitable for the LLM to compare its rollouts with and determine the correctness of the current step. This will limit the method to the tasks with such question and golden answer pairs. Therefore, we need to adapt the current method further to make it suitable for open-ended tasks.

6 CONCLUSION

In conclusion, we introduce OmegaPRM, a divide-and-conquer Monte Carlo Tree Search algorithm, designed to automate the process supervision data collection for LLMs. By efficiently pinpointing the first error in the Chain-of-Thought and balancing data quality, OmegaPRM addresses the shortcomings of existing methods. Our automated approach enables the collection of over 1.5 million process supervision annotations, which are used to train a PRM. Leveraging this automated process supervision with the weighted self-consistency algorithm, we improve LLM mathematical reasoning performance, achieving a 69.4% success rate on the MATH benchmark — a 18.4% absolute increase over the base model which amounts to a relative improvement of 36%. Additionally, our method significantly reduces data collection costs compared to human annotation and brute force Monte-Carlo sampling. These findings highlight OmegaPRM’s potential to enhance LLM capabilities in complex multi-step reasoning tasks.

REFERENCES

- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *arXiv preprint arXiv:1706.03741*, 2017.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. Alphazero-like tree-search can guide large language model decoding and training. *arXiv preprint arXiv:2309.17179*, 2023.
- Xidong Feng, Ziyu Wan, Muning Wen, Stephen Marcus McAleer, Ying Wen, Weinan Zhang, and Jun Wang. Alphazero-like tree-search can guide large language model decoding and training. *arXiv preprint arXiv:2309.17179*, 2024.
- Yao Fu, Hao Peng, Ashish Sabharwal, Peter Clark, and Tushar Khot. Complexity-based prompting for multi-step reasoning. In *Proceedings of the 11th International Conference on Learning Representations (ICLR)*, May 2023.
- Gemini Team, Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, Ioannis Antonoglou, Rohan Anil, Sebastian Borgeaud, Andrew Dai, Katie Millican, Ethan Dyer, Mia Glaese, Thibault Sottiaux, Benjamin Lee, Fabio Viola, Malcolm Reynolds, Yuanzhong Xu, James

540 Molloy, Jilin Chen, Michael Isard, Paul Barham, Tom Hennigan, Ross McIlroy, Melvin Johnson,
 541 Johan Schalkwyk, Eli Collins, Eliza Rutherford, Erica Moreira, Kareem Ayoub, Megha Goel,
 542 Clemens Meyer, Gregory Thornton, Zhen Yang, Henryk Michalewski, Zaheer Abbas, Nathan
 543 Schucher, Ankesh Anand, Richard Ives, James Keeling, Karel Lenc, Salem Haykal, Siamak
 544 Shakeri, Pranav Shyam, Aakanksha Chowdhery, Roman Ring, Stephen Spencer, Eren Sezener,
 545 Luke Vilnis, Oscar Chang, Nobuyuki Morioka, George Tucker, Ce Zheng, Oliver Woodman, Nithya
 546 Attaluri, Tomas Kocisky, Evgenii Eltyshev, Xi Chen, Timothy Chung, Vittorio Selo, Siddhartha
 547 Brahma, Petko Georgiev, Ambrose Slone, Zhenkai Zhu, James Lottes, Siyuan Qiao, Ben Caine,
 548 Sebastian Riedel, Alex Tomala, Martin Chadwick, Juliette Love, Peter Choy, Sid Mittal, Neil
 549 Hounsby, Yunhao Tang, Matthew Lamm, Libin Bai, Qiao Zhang, Luheng He, Yong Cheng, Peter
 550 Humphreys, Yujia Li, Sergey Brin, Albin Cassirer, Yingjie Miao, Lukas Zilka, Taylor Tobin,
 551 Kelvin Xu, Lev Proleev, Daniel Sohn, Alberto Magni, Lisa Anne Hendricks, Isabel Gao, Santiago
 552 Ontanon, Oskar Bunyan, Nathan Byrd, Abhanshu Sharma, Biao Zhang, Mario Pinto, Rishika
 553 Sinha, Harsh Mehta, Dawei Jia, Sergi Caelles, Albert Webson, Alex Morris, Becca Roelofs, Yifan
 554 Ding, Robin Strudel, Xuehan Xiong, Marvin Ritter, Mostafa Dehghani, Rahma Chaabouni, Abhijit
 555 Karmarkar, Guangda Lai, Fabian Mentzer, Bibo Xu, YaGuang Li, Yujing Zhang, Tom Le Paine,
 556 Alex Goldin, Behnam Neyshabur, Kate Baumli, Anselm Levskaya, Michael Laskin, Wenhao
 557 Jia, Jack W. Rae, Kefan Xiao, Antoine He, Skye Giordano, Lakshman Yagati, Jean-Baptiste
 558 Lespiau, Paul Natsev, Sanjay Ganapathy, Fangyu Liu, Danilo Martins, Nanxin Chen, Yunhan
 559 Xu, Megan Barnes, Rhys May, Arpi Vezer, Junhyuk Oh, Ken Franko, Sophie Bridgers, Ruizhe
 560 Zhao, Boxi Wu, Basil Mustafa, Sean Sechrist, Emilio Parisotto, Thanumalayan Sankaranarayanan
 561 Pillai, Chris Larkin, Chenjie Gu, Christina Sorokin, Maxim Krikun, Alexey Guseynov, Jessica
 562 Landon, Romina Datta, Alexander Pritzel, Phoebe Thacker, Fan Yang, Kevin Hui, Anja Hauth,
 563 Chih-Kuan Yeh, David Barker, Justin Mao-Jones, Sophia Austin, Hannah Sheahan, Parker Schuh,
 564 James Svensson, Rohan Jain, Vinay Ramasesh, Anton Briukhov, Da-Woon Chung, Tamara von
 565 Glehn, Christina Butterfield, Priya Jhakra, Matthew Wiethoff, Justin Frye, Jordan Grimstad, Beer
 566 Changpinyo, Charline Le Lan, Anna Bortsova, Yonghui Wu, Paul Voigtlaender, Tara Sainath, Shane
 567 Gu, Charlotte Smith, Will Hawkins, Kris Cao, James Besley, Srivatsan Srinivasan, Mark Omernick,
 568 Colin Gaffney, Gabriela Surita, Ryan Burnell, Bogdan Damoc, Junwhan Ahn, Andrew Brock,
 569 Mantas Pajarskas, Anastasia Petrushkina, Seb Noury, Lorenzo Blanco, Kevin Swersky, Arun Ahuja,
 570 Thi Avrahami, Vedant Misra, Raoul de Liedekerke, Mariko Inuma, Alex Polozov, Sarah York,
 571 George van den Driessche, Paul Michel, Justin Chiu, Rory Blevins, Zach Gleicher, Adria Recasens,
 572 Alban Rrustemi, Elena Gribovskaya, Aurko Roy, Wiktor Gworek, Sébastien M. R. Arnold, Lisa
 573 Lee, James Lee-Thorp, Marcello Maggioni, Enrique Piqueras, Kartikeya Badola, Sharad Vikram,
 574 Lucas Gonzalez, Anirudh Baddepudi, Evan Senter, Jacob Devlin, James Qin, Michael Azzam, Maja
 575 Trebacz, Martin Polacek, Kashyap Krishnakumar, Shuo yiin Chang, Matthew Tung, Ivo Penchev,
 576 Rishabh Joshi, Kate Olszewska, Carrie Muir, Mateo Wirth, Ale Jakse Hartman, Josh Newlan,
 577 Sheleem Kashem, Vijay Bolina, Elahe Dabir, Joost van Amersfoort, Zafarali Ahmed, James
 578 Cobon-Kerr, Aishwarya Kamath, Arnar Mar Hrafnkelsson, Le Hou, Ian Mackinnon, Alexandre
 579 Frechette, Eric Noland, Xiance Si, Emanuel Taropa, Dong Li, Phil Crone, Anmol Gulati, Sébastien
 580 Cevey, Jonas Adler, Ada Ma, David Silver, Simon Tokumine, Richard Powell, Stephan Lee, Kiran
 581 Vodrahalli, Samer Hassan, Diana Mincu, Antoine Yang, Nir Levine, Jenny Brennan, Mingqiu
 582 Wang, Sarah Hodkinson, Jeffrey Zhao, Josh Lipschultz, Aedan Pope, Michael B. Chang, Cheng Li,
 583 Laurent El Shafey, Michela Paganini, Sholto Douglas, Bernd Bohnet, Fabio Pardo, Seth Odoom,
 584 Mihaela Rosca, Cicero Nogueira dos Santos, Kedar Soparkar, Arthur Guez, Tom Hudson, Steven
 585 Hansen, Chulayuth Asawaroengchai, Ravi Addanki, Tianhe Yu, Wojciech Stokowiec, Mina Khan,
 586 Justin Gilmer, Jaehoon Lee, Carrie Grimes Bostock, Keran Rong, Jonathan Caton, Pedram Pejman,
 587 Filip Pavetic, Geoff Brown, Vivek Sharma, Mario Lučić, Rajkumar Samuel, Josip Djolonga,
 588 Amol Mandhane, Lars Lowe Sjöstrand, Elena Buchatskaya, Elspeth White, Natalie Clay, Jiepu
 589 Jiang, Hyeontaek Lim, Ross Hemsley, Zeynep Cankara, Jane Labanowski, Nicola De Cao, David
 590 Steiner, Sayed Hadi Hashemi, Jacob Austin, Anita Gergely, Tim Blyth, Joe Stanton, Kaushik
 591 Shivakumar, Aditya Siddhant, Anders Andreassen, Carlos Araya, Nikhil Sethi, Rakesh Shivanna,
 592 Steven Hand, Ankur Bapna, Ali Khodaei, Antoine Miech, Garrett Tanzer, Andy Swing, Shantanu
 593 Thakoor, Lora Aroyo, Zhufeng Pan, Zachary Nado, Jakub Sygnowski, Stephanie Winkler, Dian
 Yu, Mohammad Saleh, Loren Maggiore, Yamini Bansal, Xavier Garcia, Mehran Kazemi, Piyush
 Patil, Ishita Dasgupta, Iain Barr, Minh Giang, Thais Kagohara, Ivo Danihelka, Amit Marathe,
 Vladimir Feinberg, Mohamed Elhawaty, Nimesh Ghelani, Dan Horgan, Helen Miller, Lexi Walker,
 Richard Tanburn, Mukarram Tariq, Disha Shrivastava, Fei Xia, Qingze Wang, Chung-Cheng
 Chiu, Zoe Ashwood, Khuslen Baatarsukh, Sina Samangooei, Raphaël Lopez Kaufman, Fred

594 Alcober, Axel Stjerngren, Paul Komarek, Katerina Tsihlias, Anudhyan Boral, Ramona Comanescu,
595 Jeremy Chen, Ruibo Liu, Chris Welty, Dawn Bloxwich, Charlie Chen, Yanhua Sun, Fangxiaoyu
596 Feng, Matthew Mauger, Xerxes Dotiwalla, Vincent Hellendoorn, Michael Sharman, Ivy Zheng,
597 Krishna Haridasan, Gabe Barth-Maroon, Craig Swanson, Dominika Rogozińska, Alek Andreev,
598 Paul Kishan Rubenstein, Ruoxin Sang, Dan Hurt, Gamaleldin Elsayed, Renshen Wang, Dave
599 Lacey, Anastasija Ilić, Yao Zhao, Adam Iwanicki, Alejandro Lince, Alexander Chen, Christina Lyu,
600 Carl Lebsack, Jordan Griffith, Meenu Gaba, Paramjit Sandhu, Phil Chen, Anna Koop, Ravi Rajwar,
601 Soheil Hassas Yeganeh, Solomon Chang, Rui Zhu, Soroush Radpour, Elnaz Davoodi, Ving Ian
602 Lei, Yang Xu, Daniel Toyama, Constant Segal, Martin Wicke, Hanzhao Lin, Anna Bulanova,
603 Adrià Puigdomènech Badia, Nemanja Rakićević, Pablo Sprechmann, Angelos Filos, Shaobo Hou,
604 Víctor Campos, Nora Kassner, Devendra Sachan, Meire Fortunato, Chimezie Iwuanyanwu, Vitaly
605 Nikolaev, Balaji Lakshminarayanan, Sadegh Jazayeri, Mani Varadarajan, Chetan Tekur, Doug
606 Fritz, Misha Khalman, David Reitter, Kingshuk Dasgupta, Shourya Sarcar, Tina Ornduff, Javier
607 Snaider, Fantine Huot, Johnson Jia, Rupert Kemp, Nejc Trdin, Anitha Vijayakumar, Lucy Kim,
608 Christof Angermueller, Li Lao, Tianqi Liu, Haibin Zhang, David Engel, Somer Greene, Anaïs
609 White, Jessica Austin, Lilly Taylor, Shereen Ashraf, Dangyi Liu, Maria Georgaki, Irene Cai,
610 Yana Kulizhskaya, Sonam Goenka, Brennan Saeta, Ying Xu, Christian Frank, Dario de Cesare,
611 Brona Robenek, Harry Richardson, Mahmoud Alnahlawi, Christopher Yew, Priya Ponnappalli,
612 Marco Tagliasacchi, Alex Korchemniy, Yelin Kim, Dinghua Li, Bill Rosgen, Kyle Levin, Jeremy
613 Wiesner, Praseem Banzal, Praveen Srinivasan, Hongkun Yu, Çağlar Ünlü, David Reid, Zora Tung,
614 Daniel Finchelstein, Ravin Kumar, Andre Elisseeff, Jin Huang, Ming Zhang, Ricardo Aguilar,
615 Mai Giménez, Jiawei Xia, Olivier Dousse, Willi Gierke, Damion Yates, Komal Jalan, Lu Li,
616 Eri Latorre-Chimoto, Duc Dung Nguyen, Ken Durden, Praveen Kallakuri, Yaxin Liu, Matthew
617 Johnson, Tomy Tsai, Alice Talbert, Jasmine Liu, Alexander Neitz, Chen Elkind, Marco Selvi,
618 Mimi Jasarevic, Livio Baldini Soares, Albert Cui, Pidong Wang, Alek Wenjiao Wang, Xinyu
619 Ye, Krystal Kallarackal, Lucia Loher, Hoi Lam, Josef Broder, Dan Holtmann-Rice, Nina Martin,
620 Bramandia Ramadhana, Mrinal Shukla, Sujoy Basu, Abhi Mohan, Nick Fernando, Noah Fiedel,
621 Kim Paterson, Hui Li, Ankush Garg, Jane Park, DongHyun Choi, Diane Wu, Sankalp Singh,
622 Zhishuai Zhang, Amir Globerson, Lily Yu, John Carpenter, Félix de Chaumont Quitry, Carey
623 Radebaugh, Chu-Cheng Lin, Alex Tudor, Prakash Shroff, Drew Garmon, Dayou Du, Neera Vats,
624 Han Lu, Shariq Iqbal, Alex Yakubovich, Nilesh Tripuraneni, James Manyika, Haroon Qureshi, Nan
625 Hua, Christel Ngani, Maria Abi Raad, Hannah Forbes, Jeff Stanway, Mukund Sundararajan, Victor
626 Ungureanu, Colton Bishop, Yunjie Li, Balaji Venkatraman, Bo Li, Chloe Thornton, Salvatore
627 Scellato, Nishesh Gupta, Yicheng Wang, Ian Tenney, Xihui Wu, Ashish Shenoy, Gabriel Carvajal,
628 Diana Gage Wright, Ben Bariach, Zhuyun Xiao, Peter Hawkins, Sid Dalmia, Clement Farabet,
629 Pedro Valenzuela, Quan Yuan, Ananth Agarwal, Mia Chen, Wooyeol Kim, Brice Hulse, Nandita
630 Dukkipati, Adam Paszke, Andrew Bolt, Kiam Choo, Jennifer Beattie, Jennifer Prendki, Harsha
631 Vashisht, Rebeca Santamaria-Fernandez, Luis C. Cobo, Jarek Wilkiewicz, David Madras, Ali
632 Elqursh, Grant Uy, Kevin Ramirez, Matt Harvey, Tyler Liechty, Heiga Zen, Jeff Seibert, Clara Huiyi
633 Hu, Andrey Khorlin, Maigo Le, Asaf Aharoni, Megan Li, Lily Wang, Sandeep Kumar, Norman
634 Casagrande, Jay Hoover, Dalia El Badawy, David Soergel, Denis Vnukov, Matt Miecnikowski, Jiri
635 Simsa, Praveen Kumar, Thibault Sellam, Daniel Vlasic, Samira Daruki, Nir Shabat, John Zhang,
636 Guolong Su, Jiageng Zhang, Jeremiah Liu, Yi Sun, Evan Palmer, Alireza Ghaffarkhah, Xi Xiong,
637 Victor Cotruta, Michael Fink, Lucas Dixon, Ashwin Sreevatsa, Adrian Goedeckemeyer, Alek
638 Dimitriev, Mohsen Jafari, Remi Crocker, Nicholas FitzGerald, Aviral Kumar, Sanjay Ghemawat,
639 Ivan Philips, Frederick Liu, Yannie Liang, Rachel Sterneck, Alena Repina, Marcus Wu, Laura
640 Knight, Marin Georgiev, Hyo Lee, Harry Askham, Abhishek Chakladar, Annie Louis, Carl Crous,
641 Hardie Cate, Dessie Petrova, Michael Quinn, Denese Owusu-Afriyie, Achintya Singhal, Nan
642 Wei, Solomon Kim, Damien Vincent, Milad Nasr, Christopher A. Choquette-Choo, Reiko Tojo,
643 Shawn Lu, Diego de Las Casas, Yuchung Cheng, Tolga Bolukbasi, Katherine Lee, Saaber Fatehi,
644 Rajagopal Ananthanarayanan, Miteyan Patel, Charbel Kaed, Jing Li, Shreyas Rammohan Belle,
645 Zhe Chen, Jaclyn Konzelmann, Siim Pöder, Roopal Garg, Vinod Koverkathu, Adam Brown, Chris
646 Dyer, Rosanne Liu, Azade Nova, Jun Xu, Alanna Walton, Alicia Parrish, Mark Epstein, Sara
647 McCarthy, Slav Petrov, Demis Hassabis, Koray Kavukcuoglu, Jeffrey Dean, and Oriol Vinyals.
Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv
preprint arXiv:2403.05530*, 2024.

Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya
Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan

- 648 Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar,
649 Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stanczyk, Sertan Girgin,
650 Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur,
651 Olivier Bachem, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison,
652 Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia
653 Paterson, Ben Bastian, Bilal Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris
654 Perry, Chris Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger,
655 Dimple Vijaykumar, Dominika Rogozińska, Dustin Herbison, Elisa Bandy, Emma Wang, Eric
656 Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary
657 Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Plucińska, Harleen Batra,
658 Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha
659 Chan, Jin Peng Zhou, Joana Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fernandez, Joost
660 van Amersfoort, Josh Gordon, Josh Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mohamed,
661 Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia,
662 Kish Greene, Lars Lowe Sjoesund, Lauren Usui, Laurent Sifre, Lena Heuermann, Leticia Lago,
663 Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel
664 Reid, Manvinder Singh, Mark Iverson, Martin Görner, Mat Velloso, Mateo Wirth, Matt Davidow,
665 Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan,
666 Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad
667 Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda,
668 Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep
669 Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Ardeshtir Rokni, Rishabh
670 Agarwal, Ryan Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah Perrin, Sébastien M. R. Arnold,
671 Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy
672 Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomas Kocisky, Tulsee Doshi, Vihan Jain, Vikas
673 Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun
674 Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe
675 Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Warkentin, Eli Collins, Joelle Barral, Zoubin
676 Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals,
677 Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Sebastian
678 Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev.
Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*,
2024.
- 679 Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu.
680 Reasoning with language model is planning with world model. *arXiv preprint arXiv:2305.14992*,
681 2023.
- 682 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,
683 and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. *NeurIPS*,
684 2021.
- 685 Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han.
686 Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.
- 687 Jikun Kang, Xin Zhe Li, Xi Chen, Amirreza Kazemi, Qianyi Sun, Boxing Chen, Dong Li, Xu He,
688 Quan He, Feng Wen, Jianye Hao, and Jun Yao. Mindstar: Enhancing math reasoning in pre-trained
689 llms at inference time. *arXiv preprint arXiv:2405.16265*, 2024.
- 690 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making
691 large language models better reasoners with step-aware verifier. *arXiv preprint arXiv:2206.02336*,
692 2023.
- 693 Hunter Lightman, Vineet Kosaraju, Yura Burda, Harri Edwards, Bowen Baker, Teddy Lee, Jan
694 Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let’s verify step by step. *arXiv preprint
695 arXiv:2305.20050*, 2023.
- 696 Haoxiong Liu, Yifan Zhang, Yifan Luo, and Andrew Chi-Chih Yao. Augmenting math word problems
697 via iterative question composing. *arXiv preprint arXiv:2401.09003*, 2024.
- 698 Liangchen Luo, Zi Lin, Yinxiao Liu, Lei Shu, Yun Zhu, Jingbo Shang, and Lei Meng. Critique ability
699 of large language models. *arXiv preprint arXiv:2310.04815*, 2023.
- 700
701

- 702 Qianli Ma, Haotian Zhou, Tingkai Liu, Jianbo Yuan, Pengfei Liu, Yang You, and Hongxia Yang.
703 Let’s reward step by step: Step-level reward model as the navigators for reasoning. *arXiv preprint*
704 *arXiv:2310.10080*, 2023.
705
- 706 OpenAI. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
707
- 708 Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong
709 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton,
710 Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and
711 Ryan Lowe. Training language models to follow instructions with human feedback. *arXiv preprint*
712 *arXiv:2203.02155*, 2022.
- 713 Ethan Perez, Douwe Kiela, and Kyunghyun Cho. True few-shot learning with language models.
714 *arXiv preprint arXiv:2105.11447*, 2021.
715
- 716 Christopher D. Rosin. Multi-armed bandits with episode context. *Annals of Mathematics and*
717 *Artificial Intelligence*, 61(3):203–230, 2011.
- 718 David Silver, Aja Huang, Christopher J. Maddison, Arthur Guez, Laurent Sifre, George van den
719 Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander
720 Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap,
721 Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game
722 of go with deep neural networks and tree search. *Nature*, 2016. URL [https://doi.org/10.](https://doi.org/10.1038/nature16961)
723 [1038/nature16961](https://doi.org/10.1038/nature16961).
- 724 David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez,
725 Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, Yutian Chen, Timothy Lillicrap, Fan
726 Hui, Laurent Sifre, George van den Driessche, Thore Graepel, and Demis Hassabis. Mastering the
727 game of go without human knowledge. *nature*, 550(7676):354–359, 2017.
728
- 729 Ye Tian, Baolin Peng, Linfeng Song, Lifeng Jin, Dian Yu, Haitao Mi, and Dong Yu. Toward self-
730 improvement of llms via imagination, searching, and criticizing. *arXiv preprint arXiv:2404.12253*,
731 2024.
732
- 733 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
734 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cris-
735 tian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu,
736 Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn,
737 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel
738 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee,
739 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra,
740 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi,
741 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh
742 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen
743 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic,
744 Sergey Edunov, and Thomas Scialom. LLaMA 2: Open foundation and fine-tuned chat models.
745 *arXiv preprint arXiv:2307.09288*, 2023.
- 746 Jonathan Uesato, Nate Kushman, Ramana Kumar, Francis Song, Noah Siegel, Lisa Wang, Antonia
747 Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and
748 outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.
- 749 Peiyi Wang, Lei Li, Zhihong Shao, R. X. Xu, Damai Dai, Yifei Li, Deli Chen, Y. Wu, and Zhifang
750 Sui. Math-Shepherd: Verify and reinforce LLMs step-by-step without human annotations. *arXiv*
751 *preprint arXiv:2312.08935*, 2024a.
752
- 753 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-
754 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models.
755 In *Proceedings of the 11th International Conference on Learning Representations (ICLR)*, May
2023.

756 Zihan Wang, Yunxuan Li, Yuexin Wu, Liangchen Luo, Le Hou, Hongkun Yu, and Jingbo Shang.
757 Multi-step problem solving through a verifier: An empirical analysis on model-induced process
758 supervision. *arXiv preprint arXiv:2402.02658*, 2024b.

759 Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama,
760 Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals,
761 Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models, 2022a.
762

763 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, and Denny
764 Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *Proceedings of
765 the 36th Conference on Neural Information Processing Systems (NeurIPS)*, Dec 2022b.
766

767 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik
768 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *arXiv
769 preprint arXiv:2305.10601*, 2023.

770 Fei Yu, Anningzhe Gao, and Benyou Wang. Ovm, outcome-supervised value models for planning in
771 mathematical reasoning. *arXiv preprint arXiv:2311.09724*, 2024.

772 Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T. Kwok,
773 Zhenguo Li, Adrian Weller, and Weiyang Liu. MetaMath: Bootstrap your own mathematical
774 questions for large language models. *arXiv preprint arXiv:2309.12284*, 2023.
775

776 Dan Zhang, Sining Zhou, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. Rest-mcts*: Llm
777 self-training via process reward guided tree search. *arXiv preprint arXiv:2406.03816*, 2024.

778 Maciej Świechowski, Konrad Godlewski, Bartosz Sawicki, and Jacek Mańdziuk. Monte carlo tree
779 search: A review of recent modifications and applications. *arXiv preprint arXiv:2103.04931*, 2021.
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
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

APPENDIX

A QUESTION FILTERING

During the evaluation of partial solution correctness using MC estimation, false negative noise may be introduced when a question is too hard for the model, thus no correct rollout can be found even with correct partial solution. Or false positive noise may be introduced when a question is too easy, that model can conclude in correct answer given partial solution with wrong step. It is not possible to exclude such noise completely, but we can reduce the chance by filtering out questions that are either too hard or too easy for the model. Specifically, we ran a $k = 32$ rollouts for each question in the 12K training data, and filter out the questions that with no correct answer (too hard) or no wrong answer (too easy) in the 32 rollouts.

B PAIRWISE LOSS FORMULA

When training with pairwise labels, the Bradley-Terry model (people typically use this objective to train reward models in RLHF) generally accepts two probability scalars summing up to 1. When we select the two actions as a pair, there are two cases. The first case is that one sample with a zero MC value, and the other sample with a positive MC value. The second case is that both samples are with positive MC values. The first case is straight-forward, and a normalization step is required for the second case.

Assume the two MC values are p and q , and they follow the Bernoulli distribution: $P(X = 1) = p$ and $P(Y = 1) = q$. We need to calculate the probability that action X is preferred over action Y and vice versa.

$$\begin{aligned} P(X > Y) &= P(X = 1, Y = 0) = p(1 - q), \\ P(X < Y) &= P(X = 0, Y = 1) = (1 - p)q, \\ P(X = Y) &= P(X = 0, Y = 0) + P(X = 1, Y = 1) = (1 - p)(1 - q) + pq. \end{aligned} \tag{5}$$

For the tied situation, each action has half the chance of being preferred. Thus,

$$\begin{aligned} P(\text{action X is preferred}) &= P(X > Y) + 1/2 * P(X = Y) = 1/2 * (1 + p - q), \\ P(\text{action Y is preferred}) &= P(X < Y) + 1/2 * P(X = Y) = 1/2 * (1 + q - p). \end{aligned} \tag{6}$$

Now the MC values are normalized and we can train with the pairwise loss.