# Learning to Rank Chain-of-Thought: An Energy-Based Approach with Outcome Supervision

Eric H. Jiang [1]    Haozheng Luo[2]    Shengyuan Pang[3]    Xiaomin Li[4]    Zhenting Qi[4]
Hengli Li[5]    Cheng-Fu Yang[1]    Zongyu Lin[1]    Xinfeng Li[6]    Hao Xu[4]
Kai-Wei Chang[1]    Ying Nian Wu [1]

[1] Univ. of California, Los Angeles    [2] Northwestern University    [3] Zhejiang University    [4] Harvard University    [5] Peking University
[6] Nanyang Technological University

**Abstract:** Mathematical reasoning presents a significant challenge for Large Language Models (LLMs), often requiring robust multi-step logical consistency. While Chain-of-Thought (CoT) prompting elicits reasoning steps, it doesn't guarantee correctness, and improving reliability via extensive sampling is computationally costly. This paper introduces the Energy Outcome Reward Model (EORM), an effective, lightweight, post-hoc verifier. EORM leverages Energy-Based Models (EBMs) to simplify the training of reward models by learning to assign a scalar energy score to CoT solutions using only outcome labels, thereby avoiding detailed annotations. It achieves this by interpreting discriminator output logits as negative energies, effectively ranking candidates where lower energy is assigned to solutions leading to correct final outcomes—implicitly favoring coherent reasoning. On mathematical benchmarks (GSM8k, MATH), EORM significantly improves final answer accuracy (e.g., with Llama 3 8B, achieving 90.7% on GSM8k and 63.7% on MATH). EORM effectively leverages a given pool of candidate solutions to match or exceed the performance of brute-force sampling, thereby enhancing LLM reasoning outcome reliability through its streamlined post-hoc verification process. The code are held at https://github.com/ericjiang18/EnergyORM/tree/main.
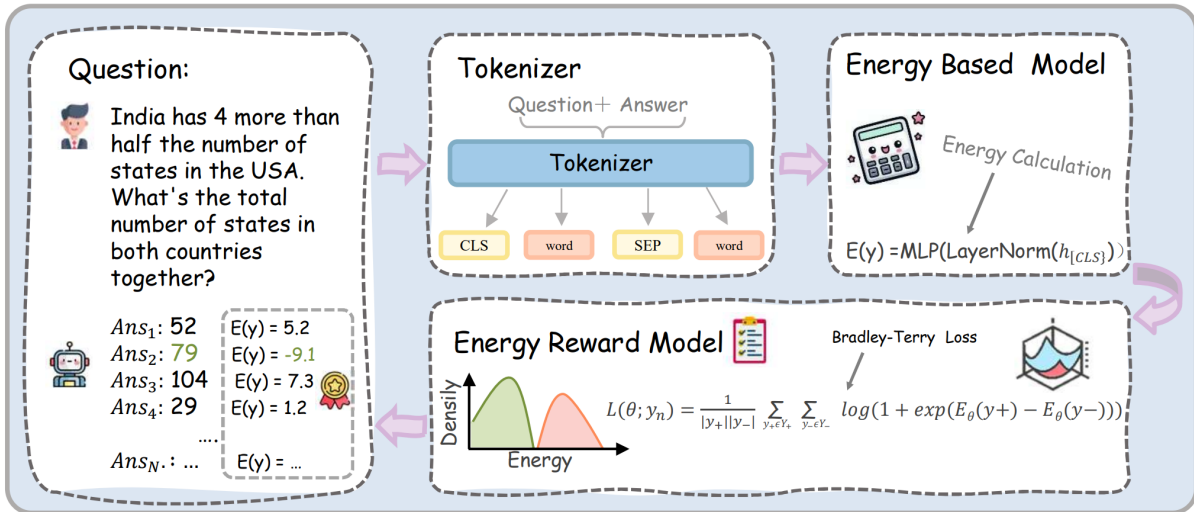
**Figure** 1: **An overview of flow chart of EORM.** In the EORM process, the model tokenizes the question-answer pair, then computes an energy score using an Energy-Based Model (EBM). The Bradley-Terry loss serves as the objective for reward-based fine-tuning. During deployment, the trained energy reward model computes energy scores for classification tasks.

| Feature | ORM | PRM | PL | EORM (ours) |
|---|---|---|---|---|
| Does not require Step-by-Step Labels? | ✓ | ✗ | ✓ | ✓ |
| Does not require Preference Pair Labels? | ✓ | ✓ | ✗ | ✓ |
| Uses Internal Indicators as Reward? | ✗ | ✗ | ✗ | ✓ |
| Low Annotation Cost? | ✓ | ✗ | ✓ | ✓ |
| Lightweight Model | ✗ | ✗ | ✗ | ✓ |

Table 1: **EORM's feature benefit compared to other reasoning baselines.** We compare EORM with ORM (Outcome Reward Model), PRM (Process Reward Model), and PL (Preference Learning) to demonstrate the feature benefit of our method over these baselines. We observe that EORM requires neither annotations nor complex labels, offering a lightweight solution for reasoning tasks. More detailed comparison analysis can be found in Appendix B

## 1. Introduction

Mathematical reasoning remains a critical and challenging domain for Large Language Models (LLMs), demanding a high degree of logical consistency and multi-step inferential accuracy that often eludes current architectures [12, 49]. While Chain-of-Thought (CoT) prompting [22, 55] has significantly advanced the ability of LLMs to articulate intermediate reasoning steps, thereby improving performance on complex tasks [14, 54, 61], it offers no intrinsic guarantee of correctness. A single erroneous step within a CoT can invalidate the entire solution [38, 49, 58], highlighting a critical need for robust verification mechanisms.

Existing approaches to bolster CoT reliability often resort to computationally intensive strategies. Self-consistency, for instance, generates numerous reasoning paths and relies on majority voting to select the final answer [54]. Other methods involve generating vast pools of candidate solutions, followed by rejection sampling or elaborate filtering protocols [33, 34, 49]. While these techniques can improve accuracy, they often entail substantial operational burdens, such as extensive complex data annotation and training regimes. This motivates the exploration of alternative verification approaches that offer a different profile of trade-offs, particularly those that simplify data requirements or streamline post-hoc application.

In this paper, we introduce the Energy Outcome Reward Model (EORM), an efficient post-hoc verifier for CoT outputs. EORM leverages principles from Energy-Based Models (EBMs) [9, 32] to effectively rerank a small set of candidate solutions. Instead of merely voting, EORM learns a rich, continuous energy landscape over the entire space of possible reasoning paths. In this landscape, logically coherent and correct solutions settle into low-energy valleys (indicating high preference), while flawed or inconsistent chains are assigned high energy. This allows EORM to capture the nuances of reasoning quality and even recognize multiple, distinct valid solutions (a multi-modal distribution), providing a more sophisticated selection mechanism than simple classification.

A cornerstone of our methodology is the insight that the output logits of standard discriminative classifiers can be directly interpreted as (negative) energies within an EBM framework [10, 32]. This crucial connection allows us to reframe the CoT reranking problem through the lens of energy-based modeling, leveraging well-understood discriminative signals related to solution correctness. By learning to score the quality of entire reasoning chains, EORM significantly simplifies the training process compared to
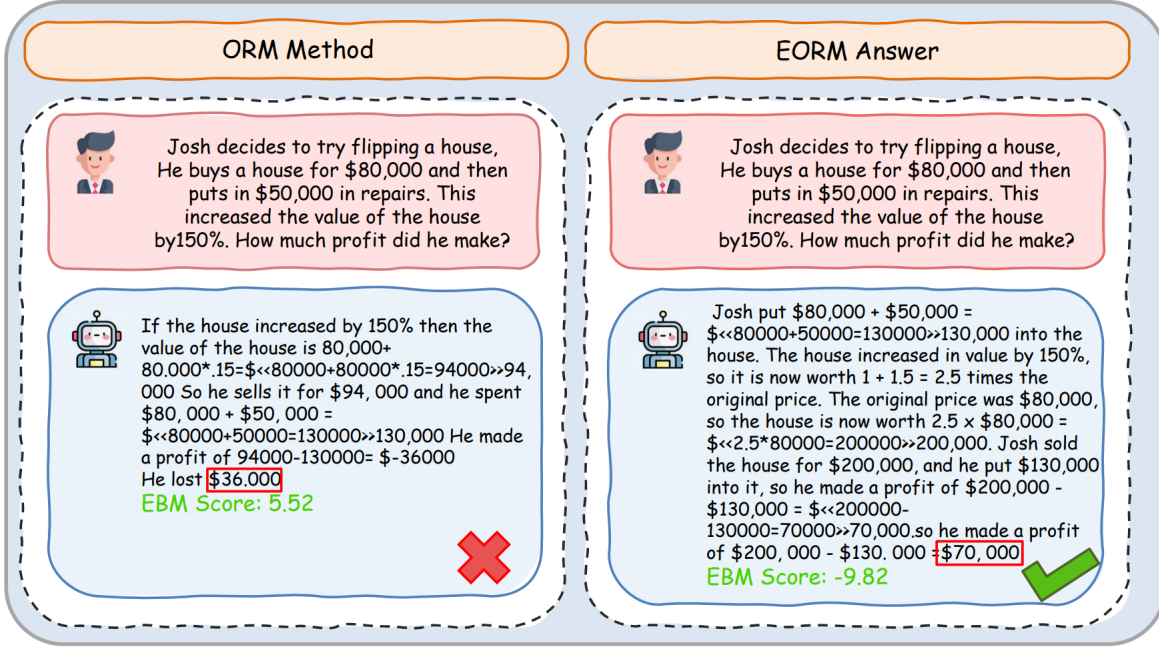
**Figure** 2: **A case study from GSM8k test set.** We present a case study comparing responses from EORM and ORM, along with the Energy-Based Model (EBM) scores for each response.

complex reinforcement learning setups or fine-grained process supervision [6, 26, 57], while remaining theoretically grounded. Consequently, EORM can substantially enhance the accuracy of the final selected answer, achieving performance comparable or superior to brute-force sampling methods but with a dramatically reduced number of candidate solutions.

Our main contributions are threefold:

- **Efficient EBM Reranker for CoT Reasoning:** We propose EORM, an energy-based verifier designed to accurately assess and rerank CoT-generated solutions. Unlike traditional methods, EORM directly assigns energy scores to generated solutions, enabling efficient selection of high-quality reasoning chains from a limited set of candidates.
- **Logits as Energy for Reward Modeling:** We demonstrate the utility of interpreting standard classifier logits directly as negative energies. This approach facilitates the straightforward training of an energy-based reward model, seamlessly bridging discriminative modeling with energy-based frameworks.
- **Significant Empirical Improvements Across Diverse Benchmarks:** We validate EORM on challenging mathematical reasoning benchmarks, including GSM8k [5] and MATH [13]. Our results show that EORM substantially improves final answer accuracy on these tasks and, importantly, also demonstrates robust generalization to various out-of-distribution (OOD) datasets, highlighting its capacity to discern valid reasoning beyond its training domains.

The remainder of this paper is organized as follows. Section 2 reviews related literature on CoT prompting, verification techniques, and energy-based models. Section 3 details the architecture and theoretical underpinnings of EORM. Section 4 presents our experimental setup and results and in-depth analysis.

Finally, Section 5 concludes the paper and discusses potential avenues for future research.

## 2. Related Work

Our work intersects with advancements in Chain-of-Thought (CoT) reasoning, the verification and reranking of Large Language Model (LLM) outputs, and Energy-Based Models (EBMs). CoT prompting [22, 55] has significantly improved multi-step reasoning in LLMs [1, 12, 35, 50], with further refinements like Least-to-Most prompting [68] and Tree-of-Thoughts [62, 65]. However, the fallibility of CoT outputs [49], where a single error can invalidate a solution, necessitates robust verification. Common approaches like self-consistency [54] improve reliability but incur substantial computational costs by sampling numerous solutions [56]. To address this, various reranking and verification techniques have emerged, including training separate verifier models [6, 20, 27] or adopting learning-to-rank perspectives [7, 31]. These methods, while effective, often introduce their own complexities or computational demands. Our approach, EORM, draws inspiration from EBMs [9], which assign scalar energy scores to configurations and are well-suited for ranking [10, 31]. Specifically, we leverage the insight that classifier logits can be interpreted as negative energies [10, 32]. By training a lightweight EBM with a pairwise Bradley-Terry loss [31] on outcome labels, EORM efficiently reranks CoT candidates, offering a distinct, streamlined post-hoc verification mechanism that complements retrieval-augmented methods [21, 38, 41, 42, 61] and other specialized verifiers like Chain-of-Actions [37] or Search-in-the-Chain [59]. For a more detailed discussion of related literature, please refer to Appendix A.

## 3. Methodology

This section details the theoretical foundations and the specific architecture of our proposed model for evaluating Chain-of-Thought (CoT) reasoning. We begin by introducing the preliminary concepts of Energy-Based Models (EBMs). We then describe our proposed Energy Outcome Reward Model (EORM), including its architecture, training objective, and implementation details. Finally, we provide a theoretical analysis of EORM, focusing on the connection to classifier logits and the derivation of the training objective. The formal definitions, theorems, and their proofs related to these concepts are available in Appendix C.

### 3.1. Preliminaries

Energy-Based Models (EBMs) provide a flexible approach to modeling distributions by assigning a scalar energy $E_\theta(y)$ to each configuration $y$ from a space $\mathcal{Y}$. This energy function is parameterized by $\theta$. Lower energy typically corresponds to more desirable or probable configurations. In this work, $\mathcal{Y}$ is the space of possible Chain-of-Thought text sequences $y$, and $\theta$ represents the learnable parameters of our EORM model.

Given an energy function $E_\theta(y)$, the corresponding Boltzmann (or Gibbs) distribution $p_\theta(y)$ assigns a probability to each configuration $y \in \mathcal{Y}$ as:

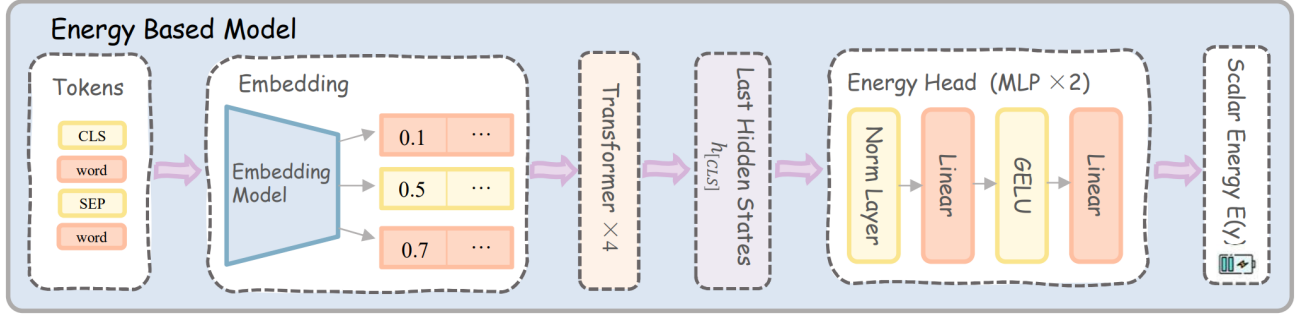$$p_\theta(y) = \frac{\exp(-E_\theta(y))}{Z_\theta}, \tag{1}$$

**Figure** 3: **An architectureof EORM.** In the EORM process, the model tokenizes the question-answer pair, then computes an energy score using an Energy-Based Model (EBM).

where $Z_\theta$ is the partition function, a normalization constant that ensures $p_\theta(y)$ sums to unity:

$$Z_\theta = \sum_{y' \in \mathcal{Y}} \exp\big(-E_\theta(y')\big) \quad \text{(for discrete } \mathcal{Y}\text{)}. \tag{2}$$

A key challenge in working with EBMs is that computing the partition function $Z_\theta$ is often computationally intractable. However, for tasks involving ranking or selecting the best candidate from a finite set $\mathcal{Y}_{\text{cand}} \subset \mathcal{Y}$, the explicit computation of $Z_\theta$ is unnecessary. Since $Z_\theta$ is constant for all $y \in \mathcal{Y}_{\text{cand}}$, comparing probabilities $p_\theta(y)$ is equivalent to comparing the unnormalized scores $\exp(-E_\theta(y))$, which in turn relates directly to comparing the energies $E_\theta(y)$. This equivalence implies that energy minimization provides a direct mechanism for identifying the most probable (preferred) solution within a candidate set, without needing to compute $Z_\theta$:

$$y^* = \arg\min_{y \in \mathcal{Y}_{\text{cand}}} E_\theta(y) = \arg\max_{y \in \mathcal{Y}_{\text{cand}}} p_\theta(y). \tag{3}$$

This makes EBMs well-suited for our task of reranking CoT solutions based on learned preferences encoded in $E_\theta(y)$. (Further details in Appendix C.1).

## 3.2. Energy Outcome Reward Model

Building upon EBM principles, we introduce EORM, an Energy Outcome Reward Model tailored for assessing CoT reasoning paths. A key aspect of EORM is its role as an energy-based verifier that directly assigns lower energy to correct solutions and higher energy to incorrect ones.

The architecture of EORM utilizes the Transformer encoder. The input CoT solution $y$ is tokenized, a special classification token (e.g., [CLS]) is prepended, and token indices are mapped to dense embedding vectors. This sequence of embeddings is processed by a stack of Transformer encoder layers, each applying multi-head self-attention and position-wise feed-forward networks. Residual connections and layer normalization are employed. After passing through the encoder, the hidden state corresponding to the prepended special token, $\mathbf{h}_{\text{CLS}}$, is treated as a summary representation of the entire input sequence $y$. A dedicated energy head, typically a small Multi-Layer Perceptron (MLP) with Layer Normalization, projects this representation $\mathbf{h}_{\text{CLS}}$ into a single scalar value:

$$E_\theta(y) = \text{MLP}\Big(\text{LayerNorm}\big(\mathbf{h}_{\text{CLS}}\big)\Big) \in \mathbb{R}. \tag{4}$$

This scalar output $E_\theta(y)$ represents the energy assigned to the sequence $y$. Lower values indicate higher assessed quality or correctness.

To train EORM effectively for the task of ranking CoT candidates, we employ a training objective based on pairwise comparisons, specifically derived from the Bradley-Terry model. The goal is to learn the parameters $\theta$ such that the model assigns consistently lower energy to correct solutions ($y_+$) compared to incorrect solutions ($y_-$) within the same context. The training data is structured into groups, where each group $\mathcal{Y}_n$ corresponds to a single problem or prompt $n$ and contains multiple candidate solutions $\{y_1, \ldots, y_k\}$. Each candidate $y_i$ is associated with a binary label $l(y_i) \in \{0, 1\}$, indicating whether it is incorrect ($l(y_i) = 0$) or correct ($l(y_i) = 1$). Within each group $\mathcal{Y}_n$, two subsets are defined:

$$\mathcal{Y}_+ = \{ y \in \mathcal{Y}_n \mid l(y) = 1\}, \quad \mathcal{Y}_- = \{ y \in \mathcal{Y}_n \mid l(y) = 0\}. \tag{5}$$

The training objective utilizes a pairwise Bradley-Terry loss (equivalent to the RankNet loss), which considers all pairs of positive and negative examples within a group. The loss aims to maximize the likelihood that positive examples have lower energy than negative examples. The loss for a single group $\mathcal{Y}_n$ is defined as:

$$\mathcal{L}(\theta; \mathcal{Y}_n) = \frac{1}{|\mathcal{Y}_+||\mathcal{Y}_-|} \sum_{y_+ \in \mathcal{Y}_+} \sum_{y_- \in \mathcal{Y}_-} \log\Big(1 + \exp\big(E_\theta(y_+) - E_\theta(y_-)\big)\Big). \tag{6}$$

This loss function penalizes pairs $(y_+, y_-)$ based on the energy difference $E_\theta(y_+) - E_\theta(y_-)$. The term $\log(1 + \exp(z))$, often referred to as the softplus function, approaches zero as $z \to -\infty$ (i.e., $E_\theta(y_+) \ll E_\theta(y_-)$) and increases as the energy of the positive candidate $E_\theta(y_+)$ increases relative to the negative candidate $E_\theta(y_-)$. Minimizing this loss encourages a separation where correct candidates consistently have lower energy scores. This pairwise formulation provides a strong learning signal, particularly beneficial for CoT ranking where subtle differences between multiple flawed reasoning paths and multiple correct paths might exist.

The training process follows a standard mini-batch gradient descent procedure, as outlined in Algorithm 1. During each epoch, the algorithm iterates through the training groups. For each non-degenerate group (where both $\mathcal{Y}_+$ and $\mathcal{Y}_-$ are non-empty), it computes the energies for positive and negative candidates, calculates the pairwise Bradley-Terry loss (Eq. 6), and updates the model parameters $\theta$ using the computed gradient.

Practical optimization employs the AdamW optimizer, which incorporates decoupled weight decay. Standard techniques like gradient norm clipping are used to prevent exploding gradients and stabilize training. A learning rate schedule, such as cosine decay with an initial warmup phase, is typically applied. For large-scale experiments, Automatic Mixed Precision (AMP) can be used to accelerate training and reduce memory consumption on compatible hardware. The combination of a relatively simple architecture (Transformer encoder + MLP head) and a well-defined, standard training procedure contributes significantly to the model's practical utility and adaptability. This reinforces the inherent model-agnostic nature and lightweight design of EORM, facilitating its easy integration with diverse pretrained language models that generate the candidate solutions. Owing to its architectural simplicity and the straightforward pairwise comparison training, EORM is expected to generalize readily across various LLM backbones and adapt to different reasoning tasks and domains without requiring extensive task-specific fine-tuning or incurring prohibitive computational overhead. This focus on relative energy differences, encouraged by the pairwise Bradley-Terry loss, potentially makes the model robust and fosters adaptability.

---

**Algorithm 2 Training procedure for the Energy Outcome Reward Model (EORM):** The model parameters $\theta$ are optimized by minimizing a pairwise Bradley-Terry loss, encouraging lower energy assignments for correct solutions $(y_+)$ compared to incorrect ones $(y_-)$ within each data group $\mathcal{Y}_n$.

---

**Require:** Labeled groups $\{\mathcal{Y}_n\}_{n=1}^N$; learning rate $\eta$; number of epochs $E$

1: Initialize EBM parameters $\theta$
2: **for** epoch $= 1$ to $E$ **do**
3:     **for** each group $\mathcal{Y}_n$ in the training set **do**
4:         $\mathcal{Y}_+ \leftarrow \{\text{positives in } \mathcal{Y}_n\}, \quad \mathcal{Y}_- \leftarrow \{\text{negatives in } \mathcal{Y}_n\}$
5:         **if** $\mathcal{Y}_+$ and $\mathcal{Y}_-$ are both non-empty **then**
6:             Compute energies $E_\theta(y)$ for all $y \in \mathcal{Y}_+ \cup \mathcal{Y}_-$
7:             $\mathcal{L}_n \leftarrow \dfrac{1}{|\mathcal{Y}_+||\mathcal{Y}_-|} \sum\limits_{y_+ \in \mathcal{Y}_+} \sum\limits_{y_- \in \mathcal{Y}_-} \log\Big(1 + \exp\big(E_\theta(y_+) - E_\theta(y_-)\big)\Big)$
8:             $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_n$
9:         **end if**
10:     **end for**
11: **end for**

---

Following training, the performance of EORM is evaluated on a held-out test set. For each group $\mathcal{Y}_n$ in the test set, the model predicts the best candidate $y^*$ by selecting the one with the minimum assigned energy: $y^* = \arg\min_{y \in \mathcal{Y}_n} E_\theta(y)$. The primary evaluation metric is Ranking Accuracy, defined as the fraction of test groups for which the selected candidate $y^*$ is a correct solution ($l(y^*) = 1$).

The theoretical grounding for our model focuses on the learned energy function and the derivation of the pairwise Bradley-Terry loss objective. EBMs learn an energy function $E_\theta(y)$ where lower energy corresponds to more desirable configurations $y$. For EORM, we train $E_\theta(y)$ to reflect the correctness of a CoT solution $y$. Learning such energy functions relates to discriminative modeling; however, instead of producing class logits, EORM learns a single scalar energy $E_\theta(y)$ via a direct head (Eq. 4). This energy acts as a discriminative score, trained via the pairwise objective (Eq. 6) to assign lower values to correct sequences, directly providing the preference score needed for ranking. The pairwise Bradley-Terry loss optimizes ranking by maximizing the likelihood that correct candidates $(y_+)$ are assigned lower energy than incorrect candidates $(y_-)$ within a group. The loss uses the smooth softplus function, enabling gradient-based optimization. The gradient of this loss, assuming $E_\theta(y)$ is differentiable with respect to $\theta$, can be expressed as:

$$\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n) = \frac{1}{|\mathcal{Y}_+||\mathcal{Y}_-|} \sum_{y_+ \in \mathcal{Y}_+} \sum_{y_- \in \mathcal{Y}_-} \sigma\Big(E_\theta(y_+) - E_\theta(y_-)\Big)\Big(\nabla_\theta E_\theta(y_+) - \nabla_\theta E_\theta(y_-)\Big), \qquad (7)$$

where $\sigma(z) = (1 + \exp(-z))^{-1}$ is the sigmoid function. Furthermore, discarding degenerate groups during training maintains an unbiased estimate of the true expected gradient. This supports using the pairwise Bradley-Terry loss as an effective and unbiased objective for training EORM. (Further details and proofs are in Appendix C.5).

## 4. Experiment

To rigorously assess our proposed Energy Outcome Reward Model (EORM), we empirically evaluated its capabilities and efficiency on mathematical reasoning tasks, specifically examining its effectiveness and generalization. For in-distribution evaluation, we utilized the GSM8k [5] and MATH [13] datasets. Out-of-distribution (OOD) generalization was assessed using AIME 2024, AMC, and AGIEval's SAT Math and Gaokao Math subsets [67]. We applied EORM to several open-source Large Language Models (LLMs);

Table 2: **Performance comparison on math reasoning benchmarks.** We evaluate EORM on GSM8K and MATH using five LLM structures (Mistral, DeepSeekMath, LLaMA3, Qwen2.5, LLaMA2), measuring accuracy. EORM achieves the highest performance across models, with best DSR values highlighted in bold.

| Model | Base | Params | GSM8k | MATH |
|---|---|---|---|---|
| Mistral-v0.1 [17] | Mistral-v0.1 | 7B | 42.9 | 12.9 |
| MathScale [47] | Mistral-v0.1 | 7B | 74.8 | 35.2 |
| MMIQC [33] | Mistral-v0.1 | 7B | 74.8 | 36.0 |
| MetaMath [63] | Mistral-v0.1 | 7B | 77.9 | 28.6 |
| KPMath-Plus [15] | Mistral-v0.1 | 7B | 82.1 | 46.8 |
| DART-Math [49] | Mistral-v0.1 | 7B | 82.6 | 43.5 |
| Math-Shepherd [53] | Mistral-v0.1 | 7B | 84.1 | 33.0 |
| MAmmoTH2-Plus [64] | Mistral-v0.1 | 7B | 84.7 | 45.0 |
| Xwin-Math [24] | Mistral-v0.1 | 7B | 89.2 | 43.7 |
| WizardMath-Mistral [34] | Mistral-v0.1 | 7B | 90.7 | 55.4 |
| **EORM (Ours)** | **Mistral-v0.1** | **7B** | **91.0** | 48.8 |
| Llama2 [50] | Llama 2 | 7B | 14.6 | 2.5 |
| MAmmoTH-CoT [64] | Llama 2 | 7B | 50.5 | 10.4 |
| MetaMath [63] | Llama 2 | 7B | 66.5 | 19.8 |
| Math-Shepherd [53] | Llama 2 | 7B | 73.2 | 21.6 |
| **EORM (Ours)** | **Llama 2** | **7B** | **75.6** | **21.8** |
| DeepSeekMath-Base [44] | DeepSeekMath | 7B | 64.2 | 36.2 |
| NuminaMath-CoT [25] | DeepseekMath | 7B | 75.4 | 55.2 |
| MMIQC [33] | DeepSeekMath | 7B | 79.0 | 45.3 |
| KPMath-Plus [15] | DeepSeekMath | 7B | 83.9 | 48.8 |
| DeepSeekMath-RL [44] | DeepSeekMath | 7B | 88.2 | 51.7 |
| **EORM (Ours)** | **DeepSeekMath** | **7B** | 84.2 | **58.7** |
| Llama 3 [48] | Llama 3 | 8B | 76.6 | 28.9 |
| MetaMath [63] | Llama 3 | 8B | 77.3 | 30.8 |
| MMIQC [33] | Llama 3 | 8B | 77.6 | 29.5 |
| DART-Math [49] | Llama 3 | 8B | 82.5 | 45.3 |
| MAmmoTH2-Plus [64] | Llama 3 | 8B | 84.1 | 42.8 |
| Llama 3.1-Instruct [48] | Llama 3 | 8B | 84.5 | 51.9 |
| JiuZhang3.0 [69] | Llama 3 | 8B | 88.6 | 51.0 |
| WizardMath-Llama [34] | Llama 3 | 8B | 90.3 | 58.8 |
| **EORM (Ours)** | **Llama 3** | **8B** | **90.7** | **63.7** |
| Qwen2.5 7B [60] | Qwen 2.5 | 7B | 89.5 | 63.4 |
| **EORM (Ours)** | **Qwen 2.5** | **7B** | **92.8** | **65.8** |

Mistral-7B-v0.1 [17], DeepSeekMath-7B [44], Llama 3 8B [11], Qwen 2.5 7B [60], and Llama 2 7B [50].

EORM was trained using a pairwise Bradley-Terry loss function [31]. The training dataset was constructed from Chain-of-Thought (CoT) solutions for problems in the in-domain GSM8k and MATH training splits. These solutions were generated using all five aforementioned LLMs. For each training problem, we generated $n = 256$ CoT candidates with a temperature of $0.7$ and a sampling probability of $0.9$, ensuring all attempts were included, regardless of their correctness. Each training instance comprised the original question, a generated solution, and a label indicating whether the solution was correct ($y^+$) or incorrect ($y^-$). EORM learned from these pairs to assign lower energy scores to preferred solutions by interpreting classifier logits as negative energies [10, 32]. The GPT-2 tokenizer [39] was employed for the energy-based tokenizer.

For evaluation, we generated Chain-of-Thought (CoT) candidates across all models using a temperature of $0.7$ and a sampling probability of $0.9$. Specifically, $n = 256$ candidates were generated for each in-distribution problem, while $n = 64$ candidates were generated for each out-of-distribution problem. Subsequently, EORM was used post-hoc to select the candidate with the lowest energy. The final answer accuracy is reported in the two subsections that follow.
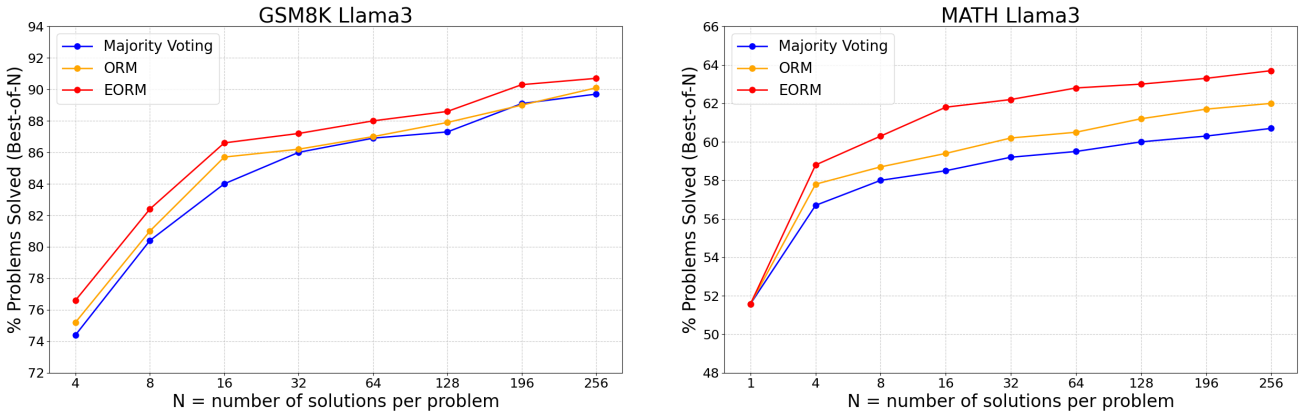


**Figure** 4: **EORM performance with varying samples per question.** We conduct experiments to show how the number of samples influences the problem-solving rate, using accuracy as the metric. The results indicate that model performance improves as the number of samples increases.

## 4.1. In Distribution Learning

Our primary evaluation focused on EORM's performance within the domains it was trained on, namely the GSM8k and MATH test sets. The quantitative outcomes, detailed in Table 2, consistently show that employing EORM as a post-hoc reranker leads to substantial improvements in final answer accuracy compared to the baseline performance of the underlying LLMs generating the candidate solutions. By effectively identifying and selecting the most plausible reasoning path from the generated set, EORM significantly enhances problem-solving capabilities. For instance, when integrated with Llama 3 8B and reranking $n = 256$ candidate solutions, EORM achieved high accuracy levels of 90.7% on GSM8k and 63.7% on MATH, demonstrating its ability to leverage multiple reasoning attempts effectively. The relationship between the number of candidate samples ($n$) and the resulting accuracy improvement is
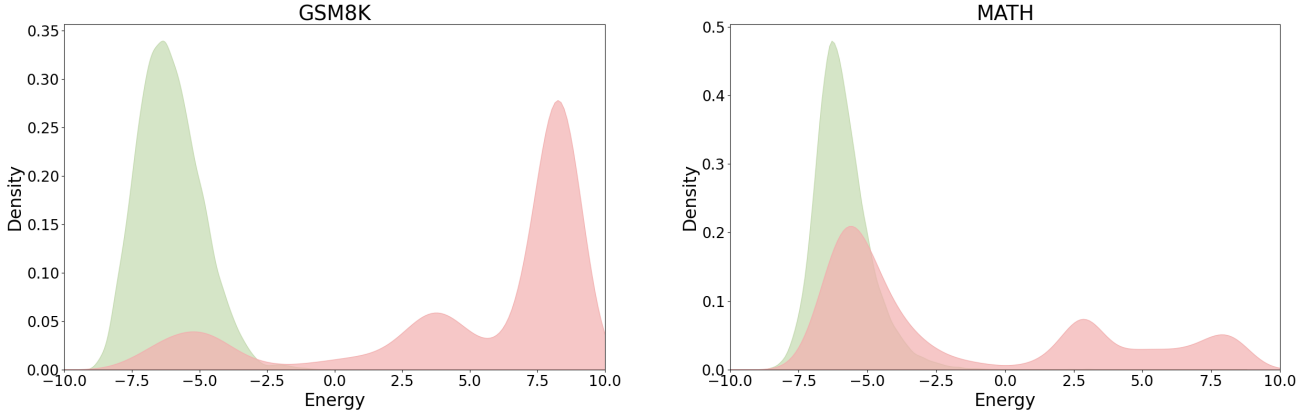
**Figure** 5: **Sample of Energy distribution of EORM under LLaMA-Sample energy distribution of EORM under LLaMA-3.** We conduct experiments on GSM8K (left) and Math (right) using EORM with LLaMA-3, evaluating 256 samples per question. The energy distribution of EORM shows correct responses in green and incorrect ones in red.

visualized in Figure 4, illustrating how performance generally scales with more samples but strong gains are often achievable with moderate sampling budgets. Beyond quantitative accuracy, qualitative analyses provide deeper insights into EORM's function. The learned energy distributions, examples of which are depicted via Kernel Density Estimation plots in Figure 6, typically reveal a clear separation between the energy scores assigned to correct versus incorrect CoT solutions. This gap indicates that EORM successfully learns to discriminate between high-quality and flawed reasoning based on the learned energy function. Additionally, a case study presented in Figure 2 provides a concrete example of EORM identifying specific logical errors within intermediate reasoning steps – a capability crucial for robust verification that goes beyond merely checking the final numerical answer.

## 4.2. Out of Distribution Learning

A critical aspect of evaluating any reasoning model is its ability to generalize to new, unseen problems and potentially different reasoning styles. To rigorously test this, we assessed the performance of EORM (trained exclusively on GSM8k and MATH data) on a suite of OOD benchmarks known for their difficulty: AIME 2024, AMC, AGIE SAT Math, and AGIE Gaokao Math [67]. These tasks often require more complex mathematical insights or different problem-solving strategies than those predominant in the training datasets. For this phase, evaluations primarily utilized Llama-3 8B and Qwen-2.5 7B as the base models, generating $n = 64$ candidate solutions for each OOD problem instance. The results, documented in Table 4, demonstrate EORM's strong generalization capabilities. When applied to the Llama-3 base model, EORM generally yielded superior performance compared to alternative baseline reranking methods like TTRL [71] and MathWizard [34] across the majority of the evaluated OOD datasets. Impressively, EORM achieved significant scores on highly challenging benchmarks such as AIME 2024 (reaching 10.0%) and the AGIE Gaokao Math problems (achieving 70.3%). Additionally the Qwen-2.5 base model also find similar finding, and provides further evidence that EORM is not merely overfitting to patterns specific to GSM8k and MATH. Instead, it appears to learn more fundamental, transferable principles related to the structure and validity of mathematical reasoning, allowing it to effectively discern higher-quality

solutions even when faced with novel problem types and domains.

Table 4: **Comparison of EORM with other reasoning methods.** We evaluate EORM against two reasoning baselines, TTRL and MathWizard, using accuracy as the metric across four mathematical datasets. Each method uses 64 samples per question for a fair comparison. EORM consistently achieves the highest performance, with the best DSR values highlighted in bold.

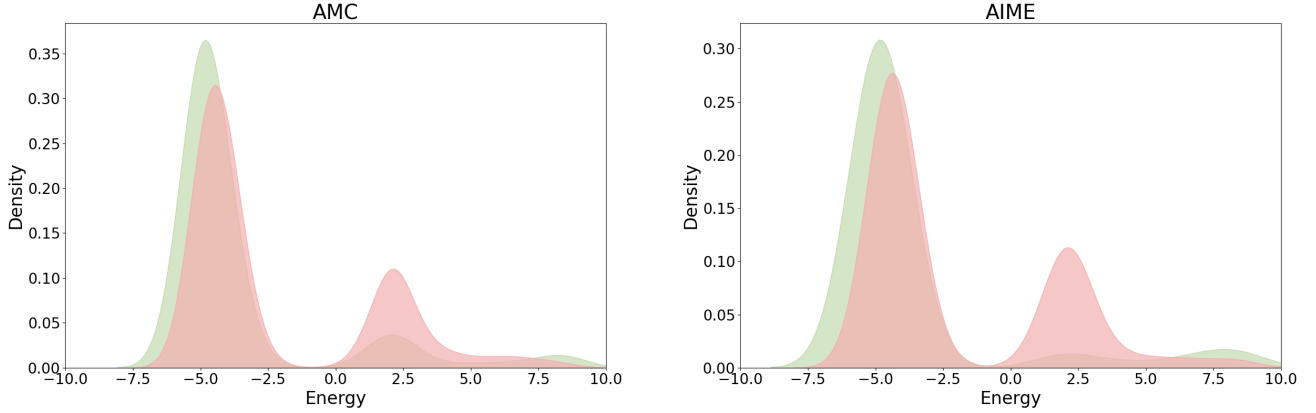| Name | AIME 2024 | AMC | SAT Math | Gaokao Math | Avg |
|---|---|---|---|---|---|
| Llama-3 7B | 3.3 | 19.3 | 77.3 | 48.7 | 37.2 |
| + TTRL [71] | 3.3 | 32.5 | 89.1 | 61.0 | 46.5 |
| + DART-MATH (Prop2Diff) [49] | 6.7 | 26.5 | **90.5** | 67.6 | 47.8 |
| + EORM Ours | **10.0** | **28.9** | **90.5** | **70.3** | **49.9** |
| Qwen-2.5 7B | 16.7 | 53.0 | 91.4 | 83.3 | 61.1 |
| + TTRL [71] | **43.3** | 67.5 | 95.0 | 88.2 | 73.5 |
| + EORM Ours | **43.3** | **68.7** | **96.4** | **88.9** | **74.3** |



**Figure** 6: **Sample of Energy distribution of EORM under LLaMA-Sample energy distribution of EORM under LLaMA-3.** We conduct experiments on AMC (left) and AIME (right) using EORM with LLaMA-3, evaluating 64 samples per question. The energy distribution of EORM shows correct responses in green and incorrect ones in red.

## 5. Conclusion

In conclusion, our work introduces EORM, an efficient and effective post-hoc verifier for Chain-of-Thought reasoning that leverages an energy-based framework to rerank candidate solutions, achieving significant accuracy improvements on both in-distribution and out-of-distribution mathematical reasoning benchmarks with reduced computational overhead compared to traditional sampling-based methods, thus presenting a promising approach for enhancing the reliability of large language models in complex reasoning tasks. This efficiency is achieved by training EORM to learn a discriminative energy function directly from simple outcome labels, thereby circumventing the need for more complex process-level supervision or the extensive computational costs associated with large-scale Monte Carlo sampling

approaches common in reasoning verification. As such, EORM represents a significant step towards developing more dependable and practically deployable LLMs capable of robust multi-step inference, particularly in settings where computational resources and detailed annotations are constrained.

# References

[1] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.

[2] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheeraz El-Showk, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022.

[3] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Chris Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Stanislav Fort, Zac Hatfield-Dodds, Tom Henighan, Nicholas Joseph, Nelson Elhage, Neel Nanda, Scott Johnston, Shauna Kravec, Liane Lovitt, Tristan Hume, Kamal Ndousse, Sam McCandlish, Tom B Brown, Dario Amodei, Jack Clark, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.

[4] Paul F Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in Neural Information Processing Systems*, 30, 2017.

[5] Karl Cobbe, Jacob Hilton, and John Schulman. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

[6] Karl Cobbe, Gregory Druck, Anand Kirubarajan, , et al. Improving mathematical reasoning with process supervision. 2023.

[7] Yunfu Deng, Qian Liu, Xiaodong He, Zihan Chen, and Hongchao Dai. Residual energy-based models for text generation. In *International Conference on Learning Representations (ICLR)*, 2020.

[8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423/.

[9] Yilun Du and Igor Mordatch. Implicit generation and modeling with energy-based models. *Advances in Neural Information Processing Systems*, 32, 2019.

[10] Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. *arXiv preprint arXiv:1912.03263*, 2019.

[11] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[12] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

[13] Dan Hendrycks, Colin Burns, Spencer Chen, , et al. Measuring mathematical problem solving with the MATH dataset. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 34, pages 3512–3524, 2021.

[14] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International Conference on Machine Learning*, pages 9118–9147. PMLR, 2022.

[15] Yiming Huang, Xiao Liu, Yeyun Gong, Zhibin Gou, Yelong Shen, Nan Duan, and Weizhu Chen. Key-point-driven data synthesis with its enhancement on mathematical reasoning, 2024.

[16] Alon Jacovi, Swaroop Mishra, Giovanni Napolitano, , et al. Reveal: A benchmark for step-level reasoning verification. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 10042–10063, Bangkok, Thailand, 2024.

[17] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023.

[18] Eric Hanchen Jiang, Zhi Zhang, Dinghuai Zhang, Andrew Lizarraga, Chenheng Xu, Yasi Zhang, Siyan Zhao, Zhengjie Xu, Peiyu Yu, Yuer Tang, Deqian Kong, and Ying Nian Wu. Dodt: Enhanced online decision transformer learning through dreamer's actor-critic trajectory forecasting, 2024. URL https://arxiv.org/abs/2410.11359.

[19] Ryo Karakida, Masato Okada, and Shun-ichi Amari. Dynamical analysis of contrastive divergence learning: Restricted boltzmann machines with gaussian visible units. *Neural Networks*, 79:78–87, 2016.

[20] Nadir Khalifa, Abdullah Alnasrallah, Joke Demuynck, , et al. GRACE: Discriminator-guided chain-of-thought reasoning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 11710–11725, Singapore, 2023.

[21] Omar Khattab, Keshav Santhanam, Xiang Lisa Li, David Hall, Percy Liang, Christopher Potts, and Matei Zaharia. Demonstrate-search-predict: Composing retrieval and language models for knowledge-intensive NLP. *arXiv preprint arXiv:2212.14024*, 2022.

[22] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35: 22199–22213, 2022.

[23] Yann LeCun, Sumit Chopra, and Raia Hadsell. A tutorial on energy-based learning. In *Predicting Structured Data*, pages 1–59. MIT Press, 2006.

[24] Chen Li, Weiqi Wang, Jingcheng Hu, Yixuan Wei, Nanning Zheng, Han Hu, Zheng Zhang, and Houwen Peng. Common 7b language models already possess strong math capabilities, 2024.

[25] Jia Li, Edward Beeching, Lewis Tunstall, Ben Lipkin, Roman Soletskyi, Shengyi Costa Huang, Kashif Rasul, Longhui Yu, Albert Jiang, Ziju Shen, Zihan Qin, Bin Dong, Li Zhou, Yann Fleureau, Guillaume Lample, and Stanislas Polu. Numinamath. [https://github.com/project-numina/aimo-progress-prize](https://github.com/project-numina/aimo-progress-prize/blob/main/report/numina_dataset.pdf), 2024.

[26] Miaoran Li, Baolin Peng, and Zhu Zhang. Self-checker: Plug-and-play modules for fact-checking with large language models. *arXiv preprint arXiv:2305.14623*, 2023.

[27] Shuohang Li, Liangming Pan, Jinpeng Zhang, , et al. Diverse: Step-aware verifier for chain-of-thought reasoning. *arXiv preprint arXiv:2306.04637*, 2023.

[28] Xiaomin Li, Xupeng Chen, Jingxuan Fan, Eric Hanchen Jiang, and Mingye Gao. Multi-head reward aggregation guided by entropy, 2025. URL https://arxiv.org/abs/2503.20995.

[29] Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. *arXiv preprint arXiv:2305.20050*, 2023.

[30] Zongyu Lin, Yao Tang, Xingcheng Yao, Da Yin, Ziniu Hu, Yizhou Sun, and Kai-Wei Chang. Qlass: Boosting language agent inference via q-guided stepwise search. *arXiv preprint arXiv:2502.02584*, 2025.

[31] Tie-Yan Liu. Learning to rank for information retrieval. In *Foundations and Trends® in Information Retrieval*, volume 3, pages 225–331. Now Publishers Inc., 2009.

[32] Weitang Liu, Xiaoyun Wang, John D Owens, and Yixuan Li. Energy-based out-of-distribution detection. *Advances in Neural Information Processing Systems*, 33:21464–21475, 2020.

[33] Xiang Liu, Baolin Peng, Yichong Zhang, , et al. MMIQC: Multi-modal in-context learning for math reasoning. *arXiv preprint arXiv:2407.13690*, 2023.

[34] Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jian-Guang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, Yansong Tang, and Dongmei Zhang. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. In *The Thirteenth International Conference on Learning Representations*, 2025.

[35] OpenAI. Gpt-4 technical report, 2023.

[36] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

[37] Zhenyu Pan, Haozheng Luo, Manling Li, and Han Liu. Chain-of-action: Faithful and multimodal question answering through large language models. In *The Thirteenth International Conference on Learning Representations*, 2025.

[38] Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt, Noah A Smith, and Mike Lewis. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*, 2022.

[39] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

[40] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *arXiv preprint arXiv:2305.18290*, 2023.

[41] Salman Rakin, Md AR Shibly, Zahin M Hossain, Zeeshan Khan, and Md Mostofa Akbar. Leveraging the domain adaptation of retrieval augmented generation models for question answering and reducing hallucination. *arXiv preprint arXiv:2410.17783*, 2024.

[42] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.

[43] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[44] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.

[45] Lin Song, Kuan Zhao, Yilun Du, and Jun Qi. Trainable energy-based models for solvable lattice models. *arXiv preprint arXiv:2101.03288*, 2021.

[46] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. In *Advances in Neural Information Processing Systems*, volume 33, pages 3008–3021, 2020.

[47] Zhengyang Tang, Xingxing Zhang, Benyou Wang, and Furu Wei. Mathscale: Scaling instruction tuning for mathematical reasoning, 2024.

[48] Llama Team. The llama 3 herd of models, 2024. URL https://arxiv.org/abs/2407.21783.

[49] Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu, and Junxian He. DART-math: Difficulty-aware rejection tuning for mathematical problem-solving. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=zLU21oQjD5.

[50] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[51] Jonathan Uesato, Nate Kushman, Ramana Kumar, H Francis Song, Noah Siegel, Lisa Wang, Antonia Creswell, Geoffrey Irving, and Irina Higgins. Solving math word problems with process-and outcome-based feedback. *arXiv preprint arXiv:2211.14275*, 2022.

[52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.

[53] Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui. Math-shepherd: Verify and reinforce LLMs step-by-step without human annotations. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9426–9439, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.510.

[54] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H. Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *The Eleventh International Conference on Learning Representations*, 2023.

[55] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi, Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large language models. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.

[56] Haoran Wu, Avishkar Balika Singh, McKay Andrus, Bohyung Yang, Abdelrahman Mourad, , et al. Mitigating misleading chain-of-thought reasoning with selective filtering. *arXiv preprint arXiv:2403.19167*, 2024.

[57] Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. A brief overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal of Automatica Sinica*, 10(5):1122–1136, 2023.

[58] Siheng Xiong, Ali Payani, Ramana Kompella, and Faramarz Fekri. Large language models can learn temporal reasoning. *arXiv preprint arXiv:2401.06853*, 2024.

[59] Shicheng Xu, Liang Pang, Huawei Shen, Xueqi Cheng, and Tat-seng Chua. Search-in-the-chain: Towards the accurate, credible and traceable content generation for complex knowledge-intensive tasks. *arXiv preprint arXiv:2304.14732*, 2023.

[60] An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, et al. Qwen2. 5 technical report. *arXiv preprint arXiv:2412.15115*, 2024.

[61] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.

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

[63] Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for large language models. In *The Twelfth International Conference on Learning Representations*, 2024.

[64] Xiang Yue, Tianyu Zheng, Ge Zhang, and Wenhu Chen. MAmmoTH2: Scaling instructions from the web. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

[65] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. In *The Eleventh International Conference on Learning Representations (ICLR 2023)*, 2023.

[66] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. URL https://openreview.net/forum?id=uccHPGDlao.

[67] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv preprint arXiv:2304.06364*, 2023.

[68] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.

[69] Kun Zhou, Beichen Zhang, jiapeng wang, Zhipeng Chen, Xin Zhao, Jing Sha, Zhichao Sheng, Shijin Wang, and Ji-Rong Wen. Jiuzhang3.0: Efficiently improving mathematical reasoning by training small data synthesis models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.

[70] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.

[71] Yuxin Zuo, Kaiyan Zhang, Shang Qu, Li Sheng, Xuekai Zhu, Biqing Qi, Youbang Sun, Ganqu Cui, Ning Ding, and Bowen Zhou. Ttrl: Test-time reinforcement learning, 2025. URL https://arxiv.org/abs/2504.16084.

# Appendix

# A. Extended Related Work

This section provides a more detailed discussion of research areas related to our work on the Energy Outcome Reward Model (EORM), covering Chain-of-Thought reasoning, techniques for verifying and reranking LLM outputs, and Energy-Based Models.

## A.1. Chain-of-Thought and Multi-Step Reasoning

Large Language Models (LLMs) [1, 12, 35, 50] have demonstrated remarkable capabilities, particularly when prompted to generate step-by-step solutions. Chain-of-Thought (CoT) prompting [22, 55] has become a cornerstone technique for eliciting complex multi-step reasoning from LLMs. This approach encourages models to "think aloud," breaking down problems into manageable intermediate steps. Various refinements to CoT have been proposed, such as Least-to-Most prompting [68], which guides the model through progressively more complex steps, and Tree-of-Thoughts (ToT) [62, 65], which explores multiple reasoning paths in a structured manner. Retrieval-augmented CoT methods [21, 37, 38, 41, 42, 61] further enhance reasoning by allowing models to consult external knowledge sources or tools. Despite these advancements, CoT outputs are not infallible; a single incorrect step can derail the entire reasoning process [49]. This has led to the development of post-hoc processing techniques. Among the most prominent is self-consistency [54], which samples multiple CoT outputs and selects the final answer via majority vote. While effective, self-consistency and similar ensemble methods often require generating a large number of candidate solutions, leading to substantial computational overhead [56]. Our work seeks to mitigate this cost by providing a more efficient mechanism for identifying high-quality CoT solutions.

## A.2. Reranking and Verification of LLM Outputs

Given the variability in the quality of CoT paths, a subsequent verification or reranking stage is often crucial for improving the final answer accuracy [16, 27, 30, 56]. Researchers have explored various strategies for this purpose. Some approaches focus on training separate verifier models to score the correctness of generated solutions or even individual reasoning steps [6, 18, 20]. For example, DI-VeRSe [27] trains a model to perform step-aware verification. Other methods adopt a learning-to-rank perspective [7, 31], training models to compare and order candidate solutions based on their perceived quality using pairwise or listwise objectives. While these verification and reranking techniques can be effective, many existing methods either rely on large, specialized verifier models that add significant computational load or employ complex, computationally expensive decoding strategies. The goal of EORM is to provide a lightweight yet powerful verifier that can efficiently identify correct CoT solutions without imposing such heavy overhead. Techniques like Chain-of-Actions [37] and Search-in-the-Chain [59] also introduce verification but often focus on different aspects or involve more intricate integration with the generation process. EORM distinguishes itself by its post-hoc applicability and its foundation in energy-based principles for efficient scoring.

## A.3. Energy-Based Models (EBMs)

Energy-Based Models (EBMs) provide a flexible and powerful framework for modeling complex data distributions by assigning a scalar "energy" value to each input configuration [9, 23]. Lower energy values typically correspond to more plausible or desirable configurations. EBMs have found successful applications in diverse domains, including computer vision [10], generative modeling [9, 45], out-of-

distribution detection [32], and ranking tasks [10, 31]. A key characteristic of EBMs is that they do not necessarily require explicit normalization of the probability distribution (i.e., computation of the partition function), which is often intractable for high-dimensional or complex output spaces [19? ]. This property makes EBMs particularly well-suited for ranking scenarios, where the primary goal is to compare the relative "goodness" of different candidates rather than to compute their absolute probabilities. This is directly applicable to CoT reranking, as we aim to identify the best solution among several alternatives. A pivotal insight, leveraged by our work, is that the logits produced by standard discriminative classifiers can be interpreted as negative unnormalized log-probabilities, or effectively, as negative energies [10]. By adopting this perspective, an EBM can be trained using objectives similar to those used for classifiers, without the need for complex sampling procedures often associated with training generative EBMs. EORM builds on this by using a pairwise Bradley-Terry loss [31], a common technique in learning-to-rank, to train the energy function to distinguish between correct and incorrect CoT solutions. Our work thus extends the application of EBM principles to the domain of multi-step reasoning in LLMs, offering an efficient and theoretically grounded alternative to existing verification methods.

## B. Comparative Analysis: EORM versus Reward Models for Guiding CoT Generation

This section provides a focused comparison between our Energy Outcome Reward Model (EORM) and the paradigm of Reward Models (RMs) specifically architected to guide Chain-of-Thought (CoT) generation, often through Reinforcement Learning (RL). Conceptually, EORM can be understood as a specialized model, built on a Transformer architecture similar to BERT [8], designed to evaluate and assign a scalar quality score, or "points"), to complex textual CoT responses. Its core evaluative goal, discerning the quality of language-based reasoning, shares similarities with how a classification model might discern patterns within its specific domain, though EORM is ultimately optimized for ranking via its energy-based framework. Functionally, this makes EORM operate very much like an automated and efficient "LLM as a judge"[66] for CoT solutions. While Table 1 in the main paper offers a high-level feature comparison showcasing EORM's advantages, here we delve into more specific distinctions regarding operational mechanisms, data requirements, training procedures, and inference characteristics when contrasted with RM-guided CoT systems, emphasizing aspects of simplicity and efficiency.

### B.1. Overview of RM-guided CoT Systems

Reward Models designed to actively guide CoT generation typically involve training a dedicated model to score the quality of partial or complete CoT paths. This learned scalar reward signal is subsequently used to improve a base Large Language Model (LLM). A common strategy is to fine-tune the LLM using RL algorithms, such as Proximal Policy Optimization (PPO) [43], to maximize the rewards obtained from the RM. Alternatively, the RM can steer the LLM's generation process more directly at inference time, for instance, by influencing sophisticated decoding strategies like best-of-N sampling or rewarded beam search. The training of these RMs frequently relies on human preference data, where annotators rank or select preferred CoT solutions from multiple options generated in response to the same prompt. This methodology is central to Reinforcement Learning from Human Feedback (RLHF) or preference learning (PL) techniques [4, 28, 36, 46, 70]. AI-generated feedback, sometimes based on predefined principles as in Constitutional AI [3], also serves as a data source for training such RMs.

## B.2. Data Requirements and Annotation Simplicity

A primary distinction favoring Eorm lies in its simpler data requirements and consequently lower annotation costs. RM-guided CoT approaches that leverage PL through RLHF often necessitate extensive and costly collection of human preference data (e.g., pairwise comparisons of CoTs [2]). If these RMs aim for fine-grained, step-by-step guidance, their data needs can escalate further, resembling those of Process Reward Models (PRM), which require detailed annotations for each reasoning step [29, 51]. In stark contrast, Eorm, which can be considered an advanced type of Outcome Reward Model (ORM), is trained solely on binary outcome labels indicating the correctness of the final answer for a complete CoT. As highlighted in Table 1, this reliance on simpler outcome labels aligns with Eorm's features of "Low Annotation Cost" and not requiring step-by-step supervision, making data acquisition significantly more straightforward and scalable.

## B.3. Training Complexity and Stability

The training paradigms also differ substantially in complexity and stability. The RL fine-tuning phase, integral to many RM-guided CoT systems (a common application of PL), is well-known for its inherent complexities. It often involves intricate actor-critic architectures, careful reward shaping, management of exploration-exploitation trade-offs, and can be prone to instability, demanding meticulous hyperparameter tuning and substantial computational resources [46]. Eorm, conversely, undergoes a standard supervised learning process. As detailed in our methodology (Section 3), it is trained by minimizing a pairwise Bradley-Terry loss. This approach is generally more stable, computationally efficient, and simpler to implement, effectively sidestepping the challenges associated with RL loops. Furthermore, Eorm's EBM framework, by interpreting logits as energies, avoids the complex training procedures of traditional generative EBMs (e.g., MCMC sampling).

## B.4. Model Characteristics and Deployment

Eorm is intentionally designed as a "lightweight" model (see Table 1), which can be significantly smaller than the base LLM generating the CoTs. While the RMs used in PL can vary in size, they are sometimes comparable to the policy model they supervise. A key advantage of Eorm is its post-hoc nature: it does not require any modification or fine-tuning of the base LLM. It simply reranks a set of pre-generated CoT outputs. This simplifies deployment significantly and allows Eorm to be versatile, readily applicable to outputs from any CoT-generating LLM without retraining or integrating into the generation model itself. Direct Preference Optimization (DPO) [40] simplifies the PL pipeline by avoiding explicit RM training and RL, but it still aims to fine-tune the policy LLM based on preference pairs, differing from Eorm's post-hoc verification role.

## B.5. Methodological Distinctions and Inference Overhead

Methodologically, many RMs used for guiding CoT (especially those trained on preferences for PL) are regression-based, learning to predict a scalar score intended to reflect human preference or solution quality. Eorm's energy-based formulation, trained with a pairwise loss, directly optimizes for the task of ranking by learning energy differences. This focus on relative quality can be more direct and potentially more robust for verification and reranking than learning an absolute reward scale that might require careful calibration. The work by Cobbe et al. [6] also developed verifiers (a form of ORM) for mathematical

problems; EORM distinguishes itself through its explicit energy-based framework and training on diverse outputs from multiple LLMs, aiming for broader generalization.

Regarding inference overhead, if an RM is used to actively guide the generation process at each step or token, it can introduce significant latency. EORM's post-hoc application to a set of $N$ completed candidates is designed for efficiency. While it adds a reranking step, the cost is often offset by the ability to achieve high accuracy from a smaller $N$ compared to methods like extensive self-consistency, thereby offering a favorable balance between accuracy improvement and computational cost, as supported by our experimental results (Section 4).

Therefore, while RM-guided CoT generation via PL and RL can produce powerful, specialized LLMs, EORM provides a simpler, more data-efficient, and computationally more tractable pathway to enhancing CoT reasoning reliability. It achieves this by focusing on post-hoc verification using readily available outcome supervision and a straightforward, EBM-inspired training regime.

## C. Theoretical Foundations and Details

This appendix furnishes comprehensive definitions, rigorous proofs, and supplementary remarks that underpin the theoretical concepts and analyses presented throughout the main text of the paper.

### C.1. Preliminaries: Foundations of Energy-Based Models

We begin by establishing the fundamental concepts of Energy-Based Models (EBMs) that form the basis of our proposed EORM.

**Definition C.1** (Energy Function). *An* energy function $E_\theta : \mathcal{Y} \to \mathbb{R}$ *is a function, parameterized by $\theta$, that assigns a scalar value $E_\theta(y)$ (its energy) to each possible configuration or input $y$ from a space $\mathcal{Y}$. By convention, lower energy values typically correspond to more desirable or probable configurations. In the EORM framework, $\mathcal{Y}$ denotes the space of possible Chain-of-Thought (CoT) text sequences, which are characterized by their variable length and complex linguistic and logical structures.*

**Definition C.2** (Boltzmann Distribution). *Given an energy function $E_\theta(y)$, the corresponding* Boltzmann (or Gibbs) distribution $p_\theta(y)$ *assigns a probability (or probability density for continuous $\mathcal{Y}$) to each configuration $y \in \mathcal{Y}$ as:*

$$p_\theta(y) \;=\; \frac{\exp\big(-E_\theta(y)\big)}{Z_\theta}, \tag{8}$$

*where $Z_\theta$ is the partition function.*

**Definition C.3** (Partition Function). *The* partition function $Z_\theta$ *serves as a normalization constant, ensuring that the Boltzmann distribution $p_\theta(y)$ integrates (or sums, for discrete $\mathcal{Y}$) to unity over the entire space $\mathcal{Y}$:*

$$Z_\theta \;=\; \int_{y' \in \mathcal{Y}} \exp\big(-E_\theta(y')\big)\, dy' \quad \Big(\text{or } \sum_{y' \in \mathcal{Y}} \exp\big(-E_\theta(y')\big) \text{ for discrete } \mathcal{Y}\Big). \tag{9}$$

*The computational intractability of $Z_\theta$ for complex, high-dimensional spaces, such as those involving text sequences, is a principal motivation for developing methods like EORM that can operate effectively without its explicit calculation.*

A key property of EBMs relevant to ranking tasks is the equivalence between energy minimization and probability maximization.

**Theorem C.1** (Energy–Probability Equivalence for Ranking). *For any finite candidate set $\mathcal{Y}_{\text{cand}} \subset \mathcal{Y}$, the configuration $y^*$ that minimizes the energy function $E_\theta(y)$ over this set is also the configuration that maximizes the Boltzmann probability $p_\theta(y)$:*

$$y^* \;=\; \arg \min_{y \in \mathcal{Y}_{\text{cand}}} E_\theta(y) \;=\; \arg \max_{y \in \mathcal{Y}_{\text{cand}}} p_\theta(y). \tag{10}$$

*Proof of Theorem C.1 (Energy–Probability Equivalence for Ranking).* Consider a finite candidate set $\mathcal{Y}_{\text{cand}}$ and the Boltzmann probability $p_\theta(y) = \frac{\exp(-E_\theta(y))}{Z_\theta}$ for $y \in \mathcal{Y}_{\text{cand}}$. The partition function $Z_\theta$, whether defined globally over $\mathcal{Y}$ or locally over $\mathcal{Y}_{\text{cand}}$ (if considering probabilities conditional on this set), is a positive constant with respect to the specific choice of $y$ from $\mathcal{Y}_{\text{cand}}$. Consequently, maximizing $p_\theta(y)$ over $y \in \mathcal{Y}_{\text{cand}}$ is equivalent to maximizing its numerator, $\exp(-E_\theta(y))$:

$$\arg \max_{y \in \mathcal{Y}_{\text{cand}}} p_\theta(y) = \arg \max_{y \in \mathcal{Y}_{\text{cand}}} \frac{\exp(-E_\theta(y))}{Z_\theta} = \arg \max_{y \in \mathcal{Y}_{\text{cand}}} \exp(-E_\theta(y)).$$

The function $f(x) = \exp(-x)$ is strictly monotonically decreasing, as its derivative, $f'(x) = -\exp(-x)$, is always negative for $x \in \mathbb{R}$. A property of strictly monotonically decreasing functions is that maximizing $f(x)$ is equivalent to minimizing $x$. Applying this, maximizing $\exp(-E_\theta(y))$ is equivalent to minimizing its argument, $E_\theta(y)$:

$$\arg \max_{y \in \mathcal{Y}_{\text{cand}}} \exp(-E_\theta(y)) = \arg \min_{y \in \mathcal{Y}_{\text{cand}}} E_\theta(y).$$

Combining these steps, we conclude that $\arg\max_{y \in \mathcal{Y}_{\text{cand}}} p_\theta(y) = \arg\min_{y \in \mathcal{Y}_{\text{cand}}} E_\theta(y)$. This equivalence is fundamental to using energy functions for ranking and selection tasks. $\qquad\square$

## C.2. Details for EORM Architecture and Training Objective

The EORM model architecture employs a Transformer encoder [52] to process input CoT sequences. Each input sequence $y$, typically a concatenation of a question and a candidate CoT solution, is first tokenized. A special classification token (e.g., [CLS]) is prepended to the tokenized sequence; its final hidden state representation is conventionally used in Transformer models to capture a summary of the entire sequence. Let $\mathbf{h}_{\text{CLS}}$ be the final hidden state vector from the Transformer encoder corresponding to this [CLS] token. This vector is then passed through a small Multi-Layer Perceptron (MLP) head, which includes Layer Normalization [**?** ] for improved training stability. The MLP projects $\mathbf{h}_{\text{CLS}}$ to a single scalar value, which EORM interprets as the energy $E_\theta(y)$ of the input sequence:

$$E_\theta(y) \;=\; \text{MLP}\Big(\text{LayerNorm}\big(\mathbf{h}_{\text{CLS}}\big)\Big) \in \mathbb{R}. \tag{11}$$

The model is trained such that lower energy values $E_\theta(y)$ correspond to higher-quality (i.e., correct) CoT solutions.

For training EORM, the data is organized into groups. Each group $\mathcal{Y}_n$ pertains to a single problem instance $n$ and comprises multiple candidate CoT solutions $\{y_1, \ldots, y_k\}$ generated for that problem. Each candidate $y_i$ is accompanied by a binary label $l(y_i) \in \{0, 1\}$, where $l(y_i) = 1$ indicates a correct

solution and $l(y_i) = 0$ indicates an incorrect one. Within each group $\mathcal{Y}_n$, we delineate two subsets: $\mathcal{Y}_+ = \{y \in \mathcal{Y}_n \mid l(y) = 1\}$ (correct solutions) and $\mathcal{Y}_- = \{y \in \mathcal{Y}_n \mid l(y) = 0\}$ (incorrect solutions). The optimization of EORM's parameters $\theta$ is driven by the pairwise Bradley-Terry loss [? ]. For a single group $\mathcal{Y}_n$ (assuming it is non-degenerate, i.e., $|\mathcal{Y}_+| > 0$ and $|\mathcal{Y}_-| > 0$), the loss is defined as:

$$\mathcal{L}(\theta; \mathcal{Y}_n) = \frac{1}{|\mathcal{Y}_+||\mathcal{Y}_-|} \sum_{y_+ \in \mathcal{Y}_+} \sum_{y_- \in \mathcal{Y}_-} \log\Big(1 + \exp\big(E_\theta(y_+) - E_\theta(y_-)\big)\Big). \tag{12}$$

This loss function penalizes instances where a correct solution $y_+$ is assigned a higher (or not sufficiently lower) energy than an incorrect solution $y_-$. The training procedure, outlined in Algorithm 1 of the main paper (and duplicated as **??** below for convenience), aims to minimize this loss across all training groups.

### C.3. Interpreting Classifier Logits as Energies: Connection to EORM

A pivotal insight that connects discriminative machine learning models to EBMs is the interpretation of classifier outputs (logits) as quantities related to energy [10, 32]. This subsection elaborates on this connection and situates EORM's methodology within this context.

**Definition C.4** (Classifier Logits and Softmax Probability). *For a $K$-class discriminative classifier, let $f_\theta(y) = [f_1(y), \ldots, f_K(y)] \in \mathbb{R}^K$ denote the vector of output scores (logits) for an input $y$, where $\theta$ represents the classifier's parameters. The conditional probability $P(k|y; \theta)$ of $y$ belonging to class $k$ is commonly computed using the softmax function:*

$$P(k|y; \theta) = \frac{\exp(f_k(y))}{\sum_{j=1}^K \exp(f_j(y))}. \tag{13}$$

**Proposition C.1** (Implicit Energy Functions from Classifier Logits). *A $K$-class classifier, as described by its logits $f_\theta(y)$ and softmax probabilities $P(k|y; \theta)$ (Definition C.4), can be understood as implicitly defining $K$ class-conditional energy functions $E_k(y; \theta) = -f_k(y)$ for each class $k$.*

1. *Using these energy functions, the conditional probability $P(k|y; \theta)$ can be expressed in an explicitly energy-based form:*

$$P(k|y; \theta) = \frac{\exp(-E_k(y; \theta))}{\sum_{j=1}^K \exp(-E_j(y; \theta))}. \tag{14}$$

2. *Furthermore, a joint probability distribution $P(y, k; \theta)$ over inputs and classes can be consistently defined within an EBM framework. If we posit a joint energy $E(y, k; \theta) = E_k(y; \theta) + E_0(y; \theta)$, where $E_0(y; \theta)$ is an energy function associated with the input $y$ itself (representing the marginal energy of $y$), then:*

$$P(y, k; \theta) = \frac{\exp\big(-(E_k(y; \theta) + E_0(y; \theta))\big)}{Z'_\theta}, \tag{15}$$

*where $Z'_\theta = \sum_{y' \in \mathcal{Y}} \sum_{l=1}^K \exp\big(-(E_l(y'; \theta) + E_0(y'; \theta))\big)$ is the global partition function. This joint distribution correctly recovers the classifier's original conditional probabilities $P(k|y; \theta)$.*

*Proof of Proposition C.1.* :
Part 1: Derivation of the energy-based form for $P(k|y; \theta)$.

We begin with the standard softmax definition for $P(k|y;\theta)$:

$$P(k|y;\theta) = \frac{\exp(f_k(y))}{\sum_{j=1}^{K} \exp(f_j(y))}.$$

By defining the class-conditional energy $E_k(y;\theta) = -f_k(y)$, it follows that the logit $f_k(y) = -E_k(y;\theta)$. Substituting $f_k(y) = -E_k(y;\theta)$ into the numerator yields $\exp(-E_k(y;\theta))$. Similarly, substituting into each term of the sum in the denominator yields $\sum_{j=1}^{K} \exp(-E_j(y;\theta))$. Thus, the conditional probability $P(k|y;\theta)$ becomes:

$$P(k|y;\theta) = \frac{\exp(-E_k(y;\theta))}{\sum_{j=1}^{K} \exp(-E_j(y;\theta))},$$

which is Eq. (14). The denominator, $\sum_{j=1}^{K} \exp(-E_j(y;\theta))$, serves as a local (input-dependent) partition function, often denoted $Z(y;\theta)$, for the conditional distribution $P(\cdot|y;\theta)$.

Part 2: Consistent definition of the joint probability $P(y,k;\theta)$.

Let us define a joint energy function for the pair $(y,k)$ as $E(y,k;\theta) = E_k(y;\theta) + E_0(y;\theta)$. Here, $E_k(y;\theta) = -f_k(y)$ are the class-conditional energies derived from classifier logits, and $E_0(y;\theta)$ represents an energy associated with the input $y$, effectively capturing the energy of $y$ under its marginal distribution. The joint probability $P(y,k;\theta)$ is then given by the Boltzmann distribution corresponding to this joint energy:

$$P(y,k;\theta) = \frac{\exp(-E(y,k;\theta))}{Z_\theta'} = \frac{\exp\big(-(E_k(y;\theta) + E_0(y;\theta))\big)}{Z_\theta'},$$

where $Z_\theta' = \sum_{y' \in \mathcal{Y}} \sum_{l=1}^{K} \exp\big(-(E_l(y';\theta) + E_0(y';\theta))\big)$ is the global partition function, summing over all possible inputs $y'$ and all classes $l$.

To confirm consistency, we derive the conditional probability $P(k|y;\theta)$ from this joint distribution using the fundamental relation $P(k|y;\theta) = \frac{P(y,k;\theta)}{P(y;\theta)}$. First, the marginal probability of $y$, $P(y;\theta)$, is obtained by summing (or marginalizing) the joint probability $P(y,l;\theta)$ over all classes $l$:

$$\begin{aligned}
P(y;\theta) &= \sum_{l=1}^{K} P(y,l;\theta) \\
&= \sum_{l=1}^{K} \frac{\exp\big(-(E_l(y;\theta) + E_0(y;\theta))\big)}{Z_\theta'} \\
&= \frac{\exp(-E_0(y;\theta))}{Z_\theta'} \sum_{l=1}^{K} \exp(-E_l(y;\theta)).
\end{aligned}$$

Now, we compute the conditional probability $P(k|y;\theta)$:

$$
\begin{aligned}
P(k|y;\theta) &= \frac{P(y,k;\theta)}{P(y;\theta)} \\
&= \frac{\frac{\exp\left(-(E_k(y;\theta)+E_0(y;\theta))\right)}{Z'_\theta}}{\frac{\exp(-E_0(y;\theta))}{Z'_\theta}\sum_{l=1}^{K}\exp(-E_l(y;\theta))} \\
&= \frac{\exp\left(-(E_k(y;\theta)+E_0(y;\theta))\right)}{\exp(-E_0(y;\theta))\sum_{l=1}^{K}\exp(-E_l(y;\theta))} \\
&= \frac{\exp(-E_k(y;\theta))\exp(-E_0(y;\theta))}{\exp(-E_0(y;\theta))\sum_{l=1}^{K}\exp(-E_l(y;\theta))} \\
&= \frac{\exp(-E_k(y;\theta))}{\sum_{l=1}^{K}\exp(-E_l(y;\theta))}.
\end{aligned}
$$

This resulting expression for $P(k|y;\theta)$ is identical to the energy-based form derived in Part 1 (Eq. (14)) and, consequently, consistent with the original softmax formulation of the classifier (Definition C.4), given the definition $E_j(y;\theta) = -f_j(y)$. This demonstrates that a joint energy definition of the form $E(y,k;\theta) = E_k(y;\theta) + E_0(y;\theta)$ allows a standard classifier to be seamlessly integrated into a joint EBM framework. The specific choice of $E_0(y;\theta)$ influences the modeled marginal $P(y;\theta)$, but the conditional $P(k|y;\theta)$ remains determined by the classifier's logits. For example, if one aims to have $P(y,k) = P(k|y)P(y)$, then $E(y,k)$ can be set to $-\log P(k|y) - \log P(y)$, which implies a specific form for $E_0(y;\theta)$ related to $-\log P(y)$ and the local partition function $Z(y;\theta)$. However, the proposition's assertion of consistency holds for a general $E_0(y;\theta)$ representing the energy contribution of $y$. $\qquad\square$

**Specialization to Binary Classification.** In the context of binary classification (e.g., a CoT solution $y$ being correct, $c = 1$, or incorrect, $c = 0$), let $f_1(y)$ and $f_0(y)$ be the logits for the "correct" and "incorrect" classes, respectively. The corresponding energies are $E_1(y) = -f_1(y)$ and $E_0(y) = -f_0(y)$. The probability of the solution being correct is $P(c = 1|y) = \frac{\exp(f_1(y))}{\exp(f_1(y))+\exp(f_0(y))}$. This can be conveniently expressed using the sigmoid function $\sigma(z) = 1/(1 + e^{-z})$ as $P(c = 1|y) = \sigma(f_1(y) - f_0(y))$. The log-odds ratio, $\log\frac{P(c=1|y)}{P(c=0|y)} = f_1(y) - f_0(y)$, is directly equal to the difference in the implicit energies: $E_0(y) - E_1(y)$. If a binary classifier is structured to output a single logit $s(y)$ (e.g., representing the evidence for the "correct" class), this is often equivalent to setting $f_1(y) = s(y)$ and $f_0(y) = 0$ (or vice-versa, potentially with a sign change). In such a scenario, the energy for the "correct" class would be $E_1(y) = -s(y)$, while the energy for the "incorrect" class would be $E_0(y) = 0$.

**EORM's Direct Energy Prediction Framework.** EORM adopts a distinct approach compared to standard classifiers that model $P(c|y)$. Instead of outputting class logits that are subsequently converted to probabilities, EORM is architected to directly output a single scalar value $E_\theta(y)$ for each input CoT sequence $y$, as specified in Eq. (11). This scalar $E_\theta(y)$ is, by design, the energy assigned to the solution $y$. EORM does not explicitly compute or model $P(c|y)$. The learning process is governed by the pairwise Bradley-Terry loss (Eq. (12)), which directly shapes the energy function $E_\theta(y)$. Specifically, the loss $\mathcal{L}(\theta; \mathcal{Y}_n)$ encourages the energy of a correct solution $y_+$ to be substantially lower than that of an incorrect solution $y_-$, i.e., $E_\theta(y_+) \ll E_\theta(y_-)$. This training regime shapes an energy landscape where desirable properties (such as correctness) correspond to low-energy states.

**Remark C.1** (EORM as a Specialized Energy-Based Verifier). *While the EORM architecture incorporates components common to discriminative classifiers (e.g., a Transformer encoder followed by an MLP head), its direct energy output and pairwise ranking-based training paradigm establish it as a specialized "energy-based verifier." Key distinctions include:*

1. ***Direct Energy Output:*** *The scalar output $E_\theta(y)$ is interpreted directly as the energy (or a measure of "undesirability"/"cost") of solution $y$, not as a logit for a specific class in the conventional sense.*
2. ***Ranking-Oriented Training:*** *Optimization occurs via a loss function that enforces relative energy orderings between pairs of solutions, rather than through a typical classification loss like cross-entropy applied to $P(c|y)$. Consequently, EORM is not required to produce calibrated class probabilities.*
3. ***Simplified EBM Application:*** *This framework avoids the complexities often associated with training generative EBMs (such as MCMC-based sampling for gradient estimation or explicit calculation of partition functions). By leveraging discriminative signals (outcome labels) within a pairwise comparison structure, EORM learns an energy function optimized for ranking tasks without needing to model the complete data distribution $p(y)$.*

*In essence, EORM capitalizes on the principle of interpreting discriminative information through an energy-based perspective but tailors this for efficient and direct learning of an energy function suited for ranking and verification.*

## C.4. Conceptual Analysis of the Learned Energy Landscape

The energy function $E_\theta(y)$ learned by EORM effectively defines an "energy landscape" across the high-dimensional, discrete space $\mathcal{Y}$ of CoT solutions. The characteristics of this landscape are paramount to EORM's ability to discriminate effectively between high-quality, correct reasoning paths and flawed or incorrect ones. While EORM does not employ $E_\theta(y)$ for generative sampling (e.g., to create novel CoTs by navigating this landscape), an understanding of the intended structure of this landscape, as sculpted by the training objective, offers valuable insights into its operational mechanism.

**Definition C.5** (Optimal Energy Separation (Idealized Goal)). *Consider a specific problem or question $q$. Let $\mathcal{Y}_{C,q} \subseteq \mathcal{Y}$ denote the set of all correct CoT solutions for $q$, and $\mathcal{Y}_{I,q} \subseteq \mathcal{Y}$ denote the set of all incorrect CoT solutions for $q$. An ideally structured energy landscape, from the perspective of discriminating solutions for question $q$, would satisfy the condition:*

$$\sup_{y_c \in \mathcal{Y}_{C,q}} E_\theta(y_c) < \inf_{y_i \in \mathcal{Y}_{I,q}} E_\theta(y_i). \tag{16}$$

*This inequality implies the existence of an energy threshold $\tau_q$ such that all correct solutions for question $q$ possess an energy $E_\theta(y_c) < \tau_q$, while all incorrect solutions have an energy $E_\theta(y_i) > \tau_q$. Achieving such perfect global separation for all possible CoTs is a highly stringent condition and unlikely to be fully realized in practice. However, the training process of EORM is designed to approximate this separation, particularly for the types of candidate solutions encountered in the training data.*

**Role of the Pairwise Bradley-Terry Loss in Shaping the Landscape.**   The pairwise Bradley-Terry loss, as defined in Eq. (12), is instrumental in sculpting the desired energy landscape. For each pair of solutions $(y_+, y_-)$ from the same problem context, where $y_+$ is correct and $y_-$ is incorrect, the loss term

is $\mathcal{L}_{pair}(y_+, y_-) = \log(1 + \exp(E_\theta(y_+) - E_\theta(y_-)))$. To understand its effect, let us define the "energy margin" for a correctly ordered pair as $\delta(y_+, y_-) = E_\theta(y_-) - E_\theta(y_+)$. The loss term can then be expressed as $\log(1 + \exp(-\delta(y_+, y_-)))$. Minimizing this loss is equivalent to maximizing the margin $\delta(y_+, y_-)$, thereby actively pushing $E_\theta(y_+)$ to be lower than $E_\theta(y_-)$. The gradient of this loss with respect to the energies (see components in Eq. (19)) illustrates this dynamic: The term $\sigma(E_\theta(y_+) - E_\theta(y_-))$ acts as a dynamic weighting factor.

- If $E_\theta(y_+) \geq E_\theta(y_-)$ (i.e., the pair is misordered or has no margin, $\delta(y_+, y_-) \leq 0$), the sigmoid term $\sigma(E_\theta(y_+) - E_\theta(y_-))$ is $\geq 0.5$. This results in a relatively strong gradient signal that pushes to decrease $E_\theta(y_+)$ and increase $E_\theta(y_-)$, effectively trying to correct the order and increase the margin.
- If $E_\theta(y_+) \ll E_\theta(y_-)$ (i.e., the pair is well-ordered with a large positive margin $\delta(y_+, y_-) \gg 0$), the sigmoid term approaches $0$. The gradient signal becomes weak, as the desired ordering is already satisfied.

This adaptive mechanism concentrates the learning effort on problematic pairs, carving out low-energy regions for correct solutions relative to incorrect ones. The loss enforces relative energy orderings rather than absolute energy targets, which can make the training more robust, for example, to imbalances in the number of correct versus incorrect examples per problem.

**Implications for Reranking and Implicit Feature Learning.** A well-structured energy landscape, characterized by such relative energy differences, enables effective reranking. During inference, EORM selects the candidate $y^* = \arg\min_{y \in \mathcal{Y}_{cand}} E_\theta(y)$ from a given pool $\mathcal{Y}_{cand}$. If the landscape reliably assigns lower energies to better solutions, this selection rule is more likely to yield a correct or high-quality CoT. The empirical performance improvements detailed in Section 4 (main paper) provide evidence for the success of this strategy. The EORM model, $E_\theta(y) = g(\phi(y; \theta_{enc}))$, learns to extract salient features $\phi(y; \theta_{enc})$ from the CoT sequence $y$ via its Transformer encoder. The pairwise training objective then optimizes both the encoder parameters $\theta_{enc}$ and the MLP head parameters (defining $g$) so that these learned features map to energy values consistent with the correctness of the final outcome. Although trained solely on final outcome labels, EORM may implicitly learn to associate various structural or semantic features within CoTs (e.g., the logical coherence between steps, the appropriate application of mathematical operations, the absence of common fallacies, or even stylistic markers correlated with correctness) with low energy values. This is because such intrinsic reasoning qualities are often causally related to achieving correct outcomes. Consequently, the learned energy landscape might reflect not only final answer correctness but also, to some extent, the internal quality of the reasoning process. For example, a CoT that arrives at a correct answer through a sequence of luckily compensating errors might be assigned a higher energy (indicating lower quality) than a CoT that reaches the same correct answer through a rigorously sound derivation, if the model has learned to penalize patterns indicative of such errors from its training experience.

**Landscape Characteristics in Discrete Space: Sensitivity and Robustness.** The domain $\mathcal{Y}$ of CoT sequences is inherently discrete and high-dimensional. Thus, notions like "smoothness" of the energy landscape are primarily conceptual. An ideal energy function $E_\theta(y)$ should exhibit high sensitivity to critical edits in $y$ that alter its logical validity or correctness—for instance, a single arithmetic mistake

or a flawed deductive step. Such critical modifications should ideally precipitate a substantial increase in energy, effectively creating steep "cliffs" or "discontinuities" in the landscape at points corresponding to logical fallacies. Conversely, minor, semantically inconsequential variations (e.g., paraphrasing that preserves logical integrity and correctness) should ideally lead to only small perturbations in the assigned energy, contributing to the model's robustness against superficial input changes. The degree to which these desirable sensitivity and robustness properties are realized depends on the representational capacity of the Transformer architecture and the diversity and quality of the training data.

**Distinction from Energy Landscapes in Generative EBMs.**   It is crucial to distinguish the role and requirements of the energy landscape in EORM from those in generative EBMs. Generative EBMs (e.g., models trained using contrastive divergence or score matching) aim to learn an energy function $E_\theta(y)$ such that the Boltzmann distribution $p_\theta(y) \propto \exp(-E_\theta(y))$ accurately approximates the true underlying data distribution $p_{data}(y)$. Inference in such generative models typically involves MCMC-based sampling techniques (like Langevin dynamics) to draw samples from $p_\theta(y)$, a process that requires navigating the energy landscape to find low-energy (high-probability) regions. Consequently, for generative EBMs, global properties of the landscape, such as the accurate representation of the partition function $Z_\theta$, the avoidance of mode collapse, and the efficiency of MCMC mixing, are of paramount importance. In stark contrast, EORM employs its learned energy function exclusively for discriminative reranking of a pre-existing, finite set of candidate solutions. It does not engage in sampling new solutions from $p_\theta(y)$. Therefore, the absolute magnitudes of the energy values or the precise global structure of $Z_\theta$ are of lesser concern than the local relative energy differences between the candidates being evaluated. This targeted application significantly simplifies the training paradigm and allows EORM to effectively harness EBM principles without contending with the more formidable challenges of generative EBM training. The "valleys" in EORM's energy landscape need only be sufficiently deep relative to the "hills" for the specific set of candidates under consideration to ensure correct ranking, rather than needing to faithfully represent a global probability density.

## C.5. Theoretical Analysis of the EORM Training Objective

This section details the mathematical properties of the EORM training objective, including definitions of relevant functions and formal derivations.

**Definition C.6** (Sigmoid Function). *The* sigmoid function $\sigma : \mathbb{R} \to (0, 1)$ *is defined as:*

$$\sigma(z) = \frac{1}{1 + \exp(-z)}. \tag{17}$$

**Definition C.7** (Softplus Function). *The* softplus function *softplus* $: \mathbb{R} \to \mathbb{R}_{>0}$ *is defined as:*

$$softplus(z) = \log(1 + \exp(z)). \tag{18}$$

*A useful property is that its derivative is the sigmoid function:* $\frac{d}{dz} softplus(z) = \sigma(z)$.

**Theorem C.2** (Gradient of Pairwise Bradley-Terry Loss). *Let* $\mathcal{L}(\theta; \mathcal{Y}_n)$ *be the pairwise Bradley-Terry loss for a group* $\mathcal{Y}_n$, *as given by Eq. (12). Assuming the energy function* $E_\theta(y)$ *is differentiable with respect to its parameters* $\theta$, *the gradient of this loss is:*

$$\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n) = \frac{1}{|\mathcal{Y}_+||\mathcal{Y}_-|} \sum_{y_+ \in \mathcal{Y}_+} \sum_{y_- \in \mathcal{Y}_-} \sigma\Big(E_\theta(y_+) - E_\theta(y_-)\Big)\Big(\nabla_\theta E_\theta(y_+) - \nabla_\theta E_\theta(y_-)\Big), \tag{19}$$

*where $\sigma(\cdot)$ is the sigmoid function (Definition C.6).*

*Proof of Theorem C.2.* The pairwise Bradley-Terry loss for a non-degenerate group $\mathcal{Y}_n$ (where $\mathcal{Y}_+$ and $\mathcal{Y}_-$ are non-empty) is:

$$\mathcal{L}(\theta; \mathcal{Y}_n) \;=\; \frac{1}{|\mathcal{Y}_+||\mathcal{Y}_-|} \sum_{y_+ \in \mathcal{Y}_+} \sum_{y_- \in \mathcal{Y}_-} \log\Big(1 + \exp\big(E_\theta(y_+) - E_\theta(y_-)\big)\Big).$$

Let $L_{pair}(y_+, y_-; \theta) = \log\big(1 + \exp\big(E_\theta(y_+) - E_\theta(y_-)\big)\big)$ denote the loss contribution from a single pair $(y_+, y_-)$. Using the definition of the softplus function (Definition C.7), we can write $L_{pair}(y_+, y_-; \theta) = \text{softplus}(E_\theta(y_+) - E_\theta(y_-))$. Due to the linearity of the gradient operator, we have:

$$\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n) \;=\; \frac{1}{|\mathcal{Y}_+||\mathcal{Y}_-|} \sum_{y_+ \in \mathcal{Y}_+} \sum_{y_- \in \mathcal{Y}_-} \nabla_\theta L_{pair}(y_+, y_-; \theta).$$

To compute $\nabla_\theta L_{pair}(y_+, y_-; \theta)$, we apply the chain rule, noting that $\frac{d}{dz}\text{softplus}(z) = \sigma(z)$:

$$\begin{aligned}
\nabla_\theta L_{pair}(y_+, y_-; \theta) &= \nabla_\theta \text{softplus}\big(E_\theta(y_+) - E_\theta(y_-)\big) \\
&= \frac{d}{dz}\text{softplus}(z)\Big|_{z=E_\theta(y_+)-E_\theta(y_-)} \cdot \nabla_\theta\big(E_\theta(y_+) - E_\theta(y_-)\big) \\
&= \sigma\big(E_\theta(y_+) - E_\theta(y_-)\big) \cdot \big(\nabla_\theta E_\theta(y_+) \;-\; \nabla_\theta E_\theta(y_-)\big).
\end{aligned}$$

Substituting this expression back into the sum for $\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)$ yields the statement of the theorem, Eq. (19). □

**Remark C.2** (Pairwise Bradley-Terry Loss as Negative Log-Likelihood)**.** *The pairwise Bradley-Terry loss function used in EORM has a strong theoretical grounding as the negative log-likelihood (NLL) of observed preferences under the Bradley-Terry probabilistic model for pairwise comparisons [? ]. This model assumes that each item $y$ possesses an underlying positive "strength" or "score," denoted $\pi_y$. The probability that item $y_i$ is preferred over item $y_j$ is then given by $P(y_i \text{ preferred over } y_j) = \frac{\pi_i}{\pi_i + \pi_j}$. In an energy-based formulation, we can define the strength of a solution $y$ as inversely related to its energy: $\pi_y = \exp(-E_\theta(y))$. A lower energy $E_\theta(y)$ thus corresponds to a higher strength $\pi_y$. Under this definition, the probability that a correct solution $y_+$ is preferred over an incorrect solution $y_-$ (which implies $E_\theta(y_+)$ should be less than $E_\theta(y_-)$) can be modeled as:*

$$\begin{aligned}
P(y_+ \text{ preferred over } y_-|\theta) &= \frac{\pi_{y_+}}{\pi_{y_+} + \pi_{y_-}} \\
&= \frac{\exp(-E_\theta(y_+))}{\exp(-E_\theta(y_+)) + \exp(-E_\theta(y_-))}.
\end{aligned}$$

*Dividing both the numerator and the denominator by $\exp(-E_\theta(y_+))$ yields:*

$$\begin{aligned}
P(y_+ \text{ preferred over } y_-|\theta) &= \frac{1}{1 + \exp(-E_\theta(y_-) + E_\theta(y_+))} \\
&= \frac{1}{1 + \exp\big(-(E_\theta(y_-) - E_\theta(y_+))\big)} \\
&= \sigma\big(E_\theta(y_-) - E_\theta(y_+)\big), \quad \text{(using the definition of } \sigma(z) \text{ from Definition C.6).}
\end{aligned}$$

*The negative log-likelihood (NLL) of observing this single preference $y_+ > y_-$ is:*

$$NLL(y_+ > y_- | \theta) = -\log P(y_+ \text{ preferred over } y_- | \theta)$$
$$= -\log \sigma \big( E_\theta(y_-) - E_\theta(y_+) \big).$$

*Using the identity $-\log \sigma(x) = \log(1 + e^{-x}) = \text{softplus}(-x)$, we have:*

$$NLL(y_+ > y_- | \theta) = \log \big( 1 + \exp\big( -(E_\theta(y_-) - E_\theta(y_+)) \big) \big)$$
$$= \log \big( 1 + \exp\big( E_\theta(y_+) - E_\theta(y_-) \big) \big).$$

*This expression is precisely the loss contribution from a single pair $(y_+, y_-)$ as defined in Eq. (12). This term is also recognizable as the logistic loss or binary cross-entropy for the event of $y_+$ being preferred over $y_-$, given the "logit" of this preference is $E_\theta(y_-) - E_\theta(y_+)$. Therefore, minimizing the average pairwise Bradley-Terry loss is equivalent to performing Maximum Likelihood Estimation (MLE) for the parameters $\theta$ under this probabilistic preference model, where preferences are determined by energy differences. This equivalence provides a robust theoretical justification for the chosen loss function.*

**Proposition C.2** (Unbiased Gradient Estimation from Non-Degenerate Groups)**.** *Let $\mathcal{G}$ represent the true underlying distribution of all training groups $n$. The loss for group $n$, $\mathcal{L}(\theta; \mathcal{Y}_n)$, is defined such that $\mathcal{L}(\theta; \mathcal{Y}_n) = 0$ if group $n$ is degenerate (i.e., it lacks either positive examples, $|\mathcal{Y}_+| = 0$, or negative examples, $|\mathcal{Y}_-| = 0$). Let $\mathcal{G}^\dagger$ be the conditional distribution of groups $n$ from $\mathcal{G}$, given that they are non-degenerate. In stochastic gradient descent (SGD), if we sample a group $n \sim \mathcal{G}^\dagger$ (i.e., we exclusively process non-degenerate groups) and compute its gradient $\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)$ using Eq. (12), this gradient serves as an unbiased estimate of the expected gradient over non-degenerate groups, $\mathbb{E}_{n \sim \mathcal{G}^\dagger}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)]$. Furthermore, this expectation is proportional to the true expected gradient over all groups, $\mathbb{E}_{n \sim \mathcal{G}}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)]$.*

*Proof of Proposition C.2.* Let $\mathbb{E}_{\mathcal{G}}[\cdot]$ denote the expectation with respect to the distribution $\mathcal{G}$ over all groups, and $\mathbb{E}_{\mathcal{G}^\dagger}[\cdot]$ denote the expectation with respect to the conditional distribution $\mathcal{G}^\dagger$ over non-degenerate groups. Let $\mathcal{S}_{\text{all}}$ be the set of indices for all possible training groups, $\mathcal{S}_{\text{nd}}$ be the set of indices for non-degenerate groups, and $\mathcal{S}_{\text{d}}$ be the set of indices for degenerate groups, such that $\mathcal{S}_{\text{all}} = \mathcal{S}_{\text{nd}} \cup \mathcal{S}_{\text{d}}$ and $\mathcal{S}_{\text{nd}} \cap \mathcal{S}_{\text{d}} = \emptyset$.

The true expected gradient over all groups is $\mathbb{E}_{\mathcal{G}}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)]$. This can be expanded as:

$$\mathbb{E}_{\mathcal{G}}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)] = \sum_{i \in \mathcal{S}_{\text{all}}} P_{\mathcal{G}}(i) \nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i)$$
$$= \sum_{i \in \mathcal{S}_{\text{nd}}} P_{\mathcal{G}}(i) \nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i) + \sum_{i \in \mathcal{S}_{\text{d}}} P_{\mathcal{G}}(i) \nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i).$$

By our definition, for any degenerate group $i \in \mathcal{S}_{\text{d}}$, the loss $\mathcal{L}(\theta; \mathcal{Y}_i) = 0$. Consequently, its gradient $\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i) = \mathbf{0}$ for $i \in \mathcal{S}_{\text{d}}$. Thus, the sum simplifies to:

$$\mathbb{E}_{\mathcal{G}}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)] = \sum_{i \in \mathcal{S}_{\text{nd}}} P_{\mathcal{G}}(i) \nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i).$$

Let $P(\mathcal{S}_{\text{nd}}) = \sum_{i \in \mathcal{S}_{\text{nd}}} P_{\mathcal{G}}(i)$ be the total probability of sampling a non-degenerate group. We assume $P(\mathcal{S}_{\text{nd}}) > 0$ (otherwise, no training would occur). The probability of sampling a specific non-degenerate group $i \in \mathcal{S}_{\text{nd}}$ under the conditional distribution $\mathcal{G}^\dagger$ (i.e., given that the sampled group is non-degenerate)

is $P_{\mathcal{G}^\dagger}(i) = P_{\mathcal{G}}(i|i \in \mathcal{S}_{\mathrm{nd}}) = \frac{P_{\mathcal{G}}(i)}{P(\mathcal{S}_{\mathrm{nd}})}$. The expected gradient when sampling exclusively from non-degenerate groups is:

$$
\begin{aligned}
\mathbb{E}_{\mathcal{G}^\dagger}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)] &= \sum_{i \in \mathcal{S}_{\mathrm{nd}}} P_{\mathcal{G}^\dagger}(i) \nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i) \\
&= \sum_{i \in \mathcal{S}_{\mathrm{nd}}} \frac{P_{\mathcal{G}}(i)}{P(\mathcal{S}_{\mathrm{nd}})} \nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i) \\
&= \frac{1}{P(\mathcal{S}_{\mathrm{nd}})} \sum_{i \in \mathcal{S}_{\mathrm{nd}}} P_{\mathcal{G}}(i) \nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_i).
\end{aligned}
$$

Comparing the two expectations, we find:

$$
\mathbb{E}_{\mathcal{G}}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)] = P(\mathcal{S}_{\mathrm{nd}}) \cdot \mathbb{E}_{\mathcal{G}^\dagger}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)].
$$

Since $P(\mathcal{S}_{\mathrm{nd}})$ is a positive constant (typically close to 1 if degenerate groups are infrequent), the gradient $\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)$ for $n \sim \mathcal{G}^\dagger$ (as computed in **??**) is an unbiased estimate of $\mathbb{E}_{\mathcal{G}^\dagger}[\nabla_\theta \mathcal{L}(\theta; \mathcal{Y}_n)]$. Furthermore, this expected gradient over non-degenerate groups is directly proportional to the true expected gradient over all groups. In the context of SGD, this means that updating parameters using gradients from only non-degenerate groups will, on average, follow a direction proportional to the true full-batch gradient direction. This justifies the common practice of skipping degenerate groups during training, as it maintains an optimization path towards minimizing the true expected loss, with the effective learning rate scaled by $P(\mathcal{S}_{\mathrm{nd}})$. $\qquad\square$

## D. Extended Studies on EORM Generalization Capabilities

This section provides a more granular analysis of the generalization capabilities of the Energy Outcome Reward Model (EORM). The objective is to rigorously assess EORM's robustness and its ability to discern valid reasoning in mathematical problem-solving when faced with outputs from Large Language Models (LLMs) whose specific reasoning traces were not encountered during its training phase. Such generalization is critical for the practical deployment of verifiers like EORM in dynamic environments where new LLMs or fine-tuned versions are continuously emerging. We investigate two primary dimensions of generalization:

1. **Generalization to Unseen Model Outputs for In-Distribution Tasks:** Evaluating EORM's performance on familiar tasks (GSM8k, MATH) when the CoT solutions are generated by an LLM not included in EORM's training data.
2. **Generalization to Unseen Model Outputs for Out-of-Distribution Tasks:** Assessing EORM's performance on novel, more challenging tasks (AIME, AMC, etc.) when the CoT solutions originate from an LLM excluded from EORM's training data, representing a compounded generalization challenge.

The findings from these studies complement the primary out-of-distribution (OOD) results presented in Section 4 (main paper) by specifically isolating and examining model-level generalization.

**D.1.  Generalization to Unseen Model Outputs for In-Distribution Tasks**

**Motivation and Experimental Setup**   To evaluate EORM's capacity to generalize across different LLM architectures and fine-tuning styles for in-distribution mathematical reasoning tasks, we adopted a leave-one-out cross-validation scheme at the model level for constructing EORM's training data. For each target LLM considered in our experiments (Mistral-v0.1, Llama2 7B, DeepSeekMath-Base, Llama3 8B, Qwen2.5 7B), a variant of EORM, termed "EORM Generalize," was trained. The training dataset for each "EORM Generalize" instance was composed of Chain-of-Thought (CoT) solutions generated by the other four LLMs for problems within the GSM8k and MATH training splits. Crucially, solutions from the specific target LLM being evaluated were entirely excluded from this training set.

The standard EORM variant, as presented in the main paper (Section 4), served as a reference. This standard variant was trained on a comprehensive dataset containing CoT solutions from all five LLMs, including those from the target model.

For evaluation, both the "EORM Generalize" variant (specific to the target LLM) and the standard EORM variant were employed to rerank $n = 256$ candidate CoT solutions. These candidate solutions were generated by the target LLM for problems drawn from the respective test sets of GSM8k [5] and MATH [13]. The base performance of the target LLM (i.e., without any reranking, typically using its default greedy decoding or a simple sampling approach if applicable) is also provided for context. Accuracy on the final answer is the primary metric reported.

Table 5: **Performance on in-distribution benchmarks (GSM8k, MATH) when EORM is trained without exposure to the target model's outputs.** EORM Generalize is trained on CoT data from four LLMs and tested on the fifth (target) LLM. EORM (from ) is trained on data from all five LLMs. Base performance of the target LLM is also provided. Accuracy is reported.

| Model | Base | Params | GSM8k | MATH |
|---|---|---|---|---|
| Mistral-v0.1 [17] | - | 7B | 42.9 | 12.9 |
| EORM Generalize | Mistral-v0.1 | 7B | 76.0 | 23.0 |
| EORM | Mistral-v0.1 | 7B | 91.0 | 48.8 |
| Llama2 [50] | - | 7B | 14.6 | 2.5 |
| EORM Generalize | Llama 2 | 7B | 38.5 | 7.8 |
| EORM | Llama 2 | 7B | 75.6 | 21.8 |
| DeepSeekMath-Base [44] | - | 7B | 64.2 | 36.2 |
| EORM Generalize | DeepSeekMath | 7B | 81.6 | 41.7 |
| EORM | DeepSeekMath | 7B | 84.2 | 58.7 |
| Llama3 [48] | - | 8B | 76.6 | 28.9 |
| EORM Generalize | Llama 3 | 8B | 76.6 | 51.6 |
| EORM | Llama 3 | 8B | 90.7 | 63.7 |
| Qwen2.5 [60] | - | 7B | 89.5 | 63.4 |
| EORM Generalize | Qwen 2.5 | 7B | 90.2 | 63.8 |
| EORM | Qwen 2.5 | 7B | 92.8 | 65.8 |

**Results and Discussion**    The comparative performance is detailed in Table 5. The results indicate that "Eorm Generalize," despite its lack of exposure to the target LLM's specific output characteristics during training, consistently yields substantial improvements in final answer accuracy over the base performance of the unevaluated LLM. For instance, when Llama 3 8B is the target model for the GSM8k dataset, its base accuracy is 42.9%. "Eorm Generalize" (trained without Llama 3 8B data) elevates this accuracy to 76.6%. In comparison, the standard Eorm (trained with Llama 3 8B data) achieves 90.7%.

The performance delta between the standard Eorm and "Eorm Generalize" is anticipated. The exclusion of the target LLM's characteristic outputs (e.g., common error patterns, stylistic choices in reasoning steps) from the training corpus for "Eorm Generalize" naturally results in a verifier that may be less attuned to the nuances of that specific model. However, the pivotal observation is the significant uplift "Eorm Generalize" still provides over the base LLM. This demonstrates that Eorm learns generalizable features and heuristics indicative of correct mathematical reasoning, which transcend the idiosyncrasies of the specific LLMs included in its training pool. The model appears to capture fundamental principles of logical consistency and procedural correctness rather than merely overfitting to superficial patterns of the training LLMs. The performance on Qwen2.5 7B is particularly strong, where "Eorm Generalize" performs very closely to the fully trained Eorm, suggesting that the reasoning patterns of Qwen2.5 might be well-represented by the other four models in the training pool, or that Qwen2.5's outputs are inherently more amenable to generalized verification.

## D.2. Generalization to Unseen Model Outputs for Out-of-Distribution Tasks

**Motivation and Experimental Setup**    This set of experiments probes a more demanding form of generalization: Eorm's ability to function effectively when both the source LLM of the candidate solutions and the task domain are unseen during its training. Specifically, the "Eorm Generalize" variant was trained exclusively on in-distribution data (GSM8k and MATH solutions) generated by $N-1$ LLMs (i.e., holding out one LLM, either Llama 3 8B or Qwen 2.5 7B in these tests, from the pool of five LLMs used to create Eorm's training data). This "Eorm Generalize" was then evaluated on its efficacy in reranking $n = 64$ candidate solutions generated by the **held-out** LLM for problems from challenging OOD benchmarks: AIME 2024, AMC, AGIEval SAT Math, and AGIEval Gaokao Math [67].

The performance of this "Eorm Generalize" is compared against two references:

1. The base performance of the held-out LLM on these OOD tasks.
2. The performance of the standard Eorm (results from Table 4 in the main paper), which was trained on in-distribution data from all five models, including the target model whose OOD performance is being evaluated.

This experimental design tests Eorm's ability to transfer learned principles of reasoning quality from familiar domains and a diverse set of known LLM styles to novel domains and an unfamiliar LLM style simultaneously.

**Results and Discussion**    The outcomes of this investigation are presented in Table 7. These results underscore Eorm's robust generalization, even under this compounded challenge. The "Eorm Generalize" variants consistently improve upon the base OOD performance of the respective held-out LLMs. For instance, with Llama 3 8B as the held-out model, its average base accuracy across the OOD datasets

Table 7: **Comparison of EORM with against baseline for out-of distribution problem and answer.** we gray shaded the generalized result

| Name | AIME 2024 | AMC | SAT Math | Gaokao Math | Avg |
|------|-----------|-----|----------|-------------|-----|
| Llama-3 7B [48] | 3.3 | 19.3 | 77.3 | 48.7 | 37.2 |
| EORM Generalized | 6.7 | 20.5 | 82.3 | 52.6 | 40.5 |
| EORM | 10.0 | 28.9 | 90.5 | 70.3 | 49.9 |
| Qwen-2.5 7B [60] | 16.7 | 53.0 | 91.4 | 83.3 | 61.1 |
| EORM Generalized | 26.7 | 54.2 | 92.3 | 84.5 | 64.4 |
| EORM | 43.3 | 68.7 | 96.4 | 88.9 | 74.3 |

is 37.2%. "EORM Generalize" (trained on GSM8k/MATH data from the other four LLMs) raises this average to 40.5%. When applied to Qwen 2.5 7B, "EORM Generalize" elevates its average OOD accuracy from 61.1% to 64.4%.

Naturally, the standard EORM (which benefited from seeing the target model's outputs during its in-distribution training phase) achieves superior OOD performance (e.g., 49.9% average for Llama 3 8B, and 74.3% for Qwen 2.5 7B). This is expected because its training included signals more directly relevant to the target LLM's characteristics. However, the fact that "EORM Generalize" still offers tangible improvements is significant. It suggests that the energy functions learned by EORM are not merely memorizing specific solution path features from its training data. Instead, EORM appears to internalize more fundamental, transferable heuristics about the structure, coherence, and validity of mathematical reasoning processes. These learned heuristics exhibit a degree of robustness that allows for effective application even when confronted with novel problem types and reasoning styles from previously unencountered LLM architectures. The consistent, albeit smaller, gains achieved by "EORM Generalize" in this demanding OOD setting highlight its potential as a versatile and broadly applicable verifier.

### D.3. Hyperparameter Settings

The training and evaluation of the EORM model were conducted using a specific set of hyperparameters, primarily configured via command-line arguments with defaults specified in the training script. Key architectural parameters for the EORM model (a Transformer-based EBM) and crucial settings for the optimization process are detailed in Table 8. The AdamW optimizer was used in conjunction with a cosine learning rate schedule including a warmup phase. Training was performed with Automatic Mixed Precision (AMP) if the '–fp16' flag was enabled and a CUDA-compatible GPU was available. The tokenizer specified by default was 'gpt2', with its vocabulary size determining the model's embedding layer dimensions. A fixed random seed (42) was used for shuffling and creating the 80/20 train-validation split from the combined training data.

## E.  Dataset Details

This section outlines the dataset utilized for training and evaluating our EORM model. We describe the data format, present a comparison with other notable mathematical reasoning datasets, and specify the overall size of our training corpus. The core of our dataset comprises (question, Chain-of-Thought solution,

Table 8: **Key hyperparameters for EORM model architecture and training.**

| Category | Hyperparameter | Value |
|---|---|---|
| **EORM Model Architecture** | Embedding Dimension | 4096 |
| | Transformer Encoder Layers | 2 |
| | Attention Heads | 4 |
| | Dropout Rate | 0.2 |
| | Feed-forward Dimension | $4 \times$ d_model |
| | Activation Function | GELU |
| | Normalization Style | Pre-LN (NormFirst) |
| **Tokenizer Configuration** | Base Tokenizer | GPT-2 |
| | Max Sequence Length | 4096 |
| | CLS ID | BOS token of tokenizer |
| | PAD ID | EOS token |
| **Training Configuration** | Optimizer | AdamW |
| | Learning Rate | $1 \times 10^{-4}$ |
| | Weight Decay | 0.01 |
| | LR Scheduler | Cosine with Warmup |
| | Warmup Ratio | 0.2 (of total steps) |
| | Number of Epochs | 50 |
| | Batch Size (groups) | 1 |
| | Gradient Clipping Norm | 1.0 |
| | FP16 (AMP) | CUDA |
| | Training/Validation Split | 80% / 20% |
| | DataLoader | 2 |

correctness label) triplets. These CoT solutions were generated by employing a suite of Large Language Models to solve problems sourced from established mathematical reasoning benchmarks, principally the training splits of GSM8k [5] and MATH [13], as elaborated in the main experimental setup (Section 4).

### E.1. Data Format

Each instance in our dataset is represented as a JSON object, with multiple instances typically stored in a JSON Lines ('.jsonl') file. Each JSON object contains the original mathematical problem, a generated CoT solution, and a binary label indicating whether the provided solution correctly answers the problem. An illustrative example of a single data instance is presented in Appendix E.1.

---

**Training Dataset Format Example**

**label:** 1

**question:**
How many vertical asymptotes does the graph of

$$y = \frac{2}{x^2 + x - 6}$$

have?

**gen_text:**

To determine the number of vertical asymptotes for the function $y = \frac{2}{x^2 + x - 6}$, we need to identify the values of $x$ that make the denominator zero, as these are the points where the function is undefined and potentially has vertical asymptotes.

1. *Factor the denominator*
   $x^2 + x - 6 = (x + 3)(x - 2)$

2. *Solve for zeros*
   $(x + 3)(x - 2) = 0 \implies x = -3 \quad \text{or} \quad x = 2$

3. *Check for holes*
   The numerator is the constant 2, which is non-zero at $x = -3$ and $x = 2$, so no factors cancel. Hence each zero of the denominator corresponds to a vertical asymptote.

Therefore, the graph has

$$\boxed{2}$$

vertical asymptotes: one at $x = -3$ and one at $x = 2$.

---

## E.2. Comparison Between Other Dataset

To position our data generation efforts for training EORM, Table 9 provides a comparative overview against several other prominent datasets used in mathematical reasoning research. The table highlights the scale (approximate number of samples) and the primary "Synthesis Agent" responsible for creating the core problems or question-answer pairs in those datasets. Many existing datasets, as indicated, leverage powerful large language models (LLMs) like GPT-4 or GPT-3.5, sometimes augmented with human effort, to synthesize novel mathematical problems. In contrast, our approach for the EORM training data (labeled as "EORM" in the table) does not involve the synthesis of new problems. Instead, we focus on generating a very large corpus of Chain-of-Thought (CoT) solutions by prompting various existing LLMs with problems from established, open benchmarks (such as GSM8k and MATH). Each generated CoT solution is then labeled for correctness based on its final answer. This strategy, reflected by "None" under "Synthesis Agent" for our EORM data, results in a distinct type of dataset geared towards learning to verify reasoning processes rather than generating new problem instances. The substantial size of our collected data, approximately 14.96 million (14958k) (question, CoT solution, label) triplets, is intended to provide a diverse and comprehensive basis for training a robust EORM verifier capable of

Table 9: **Comparison of the EORM training data scale and problem synthesis approach with other mathematical reasoning datasets.** Dataset sizes are approximate, in thousands (k) of samples. "Synthesis Agent" refers to the primary method or model used for generating the core problems in the dataset.

| Dataset | Synthesis Agent | Dataset Size (k) |
|---|---|---|
| WizardMath [34] | GPT-4 | 96 |
| MetaMath [63] | GPT-3.5 | 395 |
| MMIQC [33] | GPT-4+GPT-3.5+Human | 2294 |
| Xwin-Math-V1.1 [24] | GPT-4 | 1440 |
| KPMath-Plus [15] | GPT-4 | 1576 |
| MathScale [47] | GPT-3.5+Human | 2021 |
| DART-Math [49] | DeepSeekMath-7B-RL | 591 |
| **EORM** | **None** | **14958** |

understanding a wide range of reasoning styles and error patterns.

## F. Hardware Resources

The computational experiments presented in this paper, encompassing both the generation of Chain-of-Thought (CoT) solutions and the training of our Energy Outcome Reward Model (EORM), utilized high-performance Graphics Processing Units (GPUs). For the initial phase of generating the CoT candidate solutions from various Large Language Models (LLMs), a range of NVIDIA GPUs was employed. This included access to NVIDIA A800 (40GB), NVIDIA A100 (with both 40GB and 80GB variants), NVIDIA RTX A6000 Ada (48GB), and NVIDIA H100 GPUs (80GB). The availability of this diverse hardware allowed for extensive data generation across multiple LLM architectures. The training of the EORM model itself was conducted on a more focused set of high-end accelerators. Specifically, we utilized configurations consisting of 2 to 4 NVIDIA H100 (80GB) GPUs. With a setup of 4 NVIDIA H100 (80GB) GPUs, the complete training process for the EORM model took approximately 36 hours.

## G. Limitations and Future Directions

Despite its demonstrated effectiveness, EORM has several inherent limitations. Primarily, its performance as a post-hoc reranker is fundamentally constrained by the quality and diversity of the initial candidate CoT solutions generated by the base LLM; EORM cannot create correct solutions if none are present in the input pool. Furthermore, its training relies on binary outcome labels, which, while simplifying annotation, may not always distinguish sound reasoning from coincidentally correct answers, potentially missing nuances in reasoning coherence. While EORM shows promising out-of-distribution generalization on related mathematical tasks (Table 7), its applicability to vastly different domains or reasoning styles warrants further investigation. The efficacy of a trained EORM instance is also tied to its own architecture and the characteristics of its training data, potentially requiring retraining for optimal performance with significantly novel LLM CoT styles, as explored in our generalization studies (Appendix D).

Looking ahead, several exciting avenues exist for extending and enhancing the EORM framework. To more

directly address reasoning quality, future work could explore hybrid supervision models that integrate Eorm's outcome-based learning with elements of process-based supervision, perhaps drawing lightweight inspiration from process reward models [6]. Calibrating energy scores to yield meaningful confidence estimates and developing a deeper understanding of the features learned by Eorm through probing techniques (complementing the conceptual analysis in Appendix C.4) are also important directions for improving utility and interpretability. System-level integrations, such as using Eorm to guide adaptive candidate generation or iterative CoT refinement, could lead to more efficient and robust reasoning systems. Expanding and systematically evaluating Eorm's transferability to diverse domains like code generation or logical puzzles, potentially through domain adaptation, is crucial. Finally, extending Eorm for multi-objective ranking (e.g., considering conciseness or interpretability alongside correctness), further developing the theoretical underpinnings of its learned energy functions, and conducting detailed error analyses will be vital for advancing the reliability and applicability of energy-based verifiers in complex AI reasoning.