

# Reasoning Language Models: A Blueprint

Maciej Besta<sup>1†</sup>, Julia Barth<sup>1</sup>, Eric Schreiber<sup>1</sup>, Ales Kubicek<sup>1</sup>, Afonso Catarino<sup>1</sup>, Robert Gerstenberger<sup>1</sup>, Piotr Nyczyk<sup>2</sup>, Patrick Iff<sup>1</sup>, Yueling Li<sup>3</sup>, Sam Houlston<sup>1</sup>, Tomasz Sternal<sup>1</sup>, Marcin Copik<sup>1</sup>, Grzegorz Kwaśniewski<sup>1</sup>, Jürgen Müller<sup>3</sup>, Łukasz Flis<sup>4</sup>, Hannes Eberhard<sup>1</sup>, Hubert Niewiadomski<sup>2</sup>, Torsten Hoefler<sup>1</sup>

<sup>†</sup>Corresponding author <sup>1</sup>ETH Zurich <sup>2</sup>Cledar <sup>3</sup>BASF SE <sup>4</sup>Cyfronet AGH

**Abstract**—Reasoning language models (RLMs), also known as Large Reasoning Models (LRMs), such as OpenAI’s o1 and o3, DeepSeek-V3, and Alibaba’s QwQ, have redefined AI’s problem-solving capabilities by extending large language models (LLMs) with advanced reasoning mechanisms. Yet, their high costs, proprietary nature, and complex architectures—uniquely combining Reinforcement Learning (RL), search heuristics, and LLMs—present accessibility and scalability challenges. To address these, we propose a comprehensive blueprint that organizes RLM components into a modular framework, based on a survey and analysis of all RLM works. This blueprint incorporates diverse reasoning structures (chains, trees, graphs, and nested forms), reasoning strategies (e.g., Monte Carlo Tree Search, Beam Search), RL concepts (policy, value models and others), supervision schemes (Outcome-Based and Process-Based Supervision), and other related concepts (e.g., Test-Time Compute, Retrieval-Augmented Generation, agent tools). We also provide detailed mathematical formulations and algorithmic specifications to simplify RLM implementation. By showing how schemes like LLaMA-Berry, QwQ, Journey Learning, and Graph of Thoughts fit as special cases, we demonstrate the blueprint’s versatility and unifying potential. To illustrate its utility, we introduce **x1**, a modular implementation for rapid RLM prototyping and experimentation. Using **x1** and a literature review, we provide key insights, such as multi-phase training for policy and value models, and the importance of familiar training distributions. Finally, we discuss scalable RLM cloud deployments and we outline how RLMs can integrate with a broader LLM ecosystem. Our work demystifies RLM construction, democratizes advanced reasoning capabilities, and fosters innovation, aiming to mitigate the gap between “rich AI” and “poor AI” by lowering barriers to RLM design and experimentation.

**Index Terms**—Reasoning Language Model, Large Reasoning Model, Survey of Reasoning Language Models, Survey of RLMS, RLM, LRM, Reasoning LLMs, Reinforcement Learning for LLMs, MCTS for LLMs, Large Language Model, LLM, Generative AI.



## 1 INTRODUCTION

Reasoning Language Models (RLMs), such as OpenAI’s o1 [116], o3 [76], and Alibaba’s QwQ [148], also referred to as Large Reasoning Models (LRMs)<sup>1</sup>, represent a transformative breakthrough in AI, on par with the advent of ChatGPT [114]. These advanced systems have fundamentally redefined AI’s problem-solving capabilities, enabling nuanced reasoning, improved contextual understanding, and robust decision-making across a wide array of domains, reshaping science [45], industries [21], governance [52], and numerous other aspects of human life [46], [75], [80], [143], [144]. By extending the capabilities of standard large language

models (LLMs) with sophisticated reasoning mechanisms, RLMs have emerged as the new cornerstone of cutting-edge AI, bringing us closer to AGI.

However, the high cost and proprietary nature of state-of-the-art RLMs, such as those developed by OpenAI, risk exacerbating the divide between “rich AI” and “poor AI”, raising significant concerns about accessibility and equity. Even the publicly available QwQ only comes with its model weights, and Alibaba does not disclose details about their training or data generation methodologies. Businesses and individuals unable to afford these advanced systems face a growing disadvantage, threatening to stifle innovation and reinforce systemic inequities. As RLMs become integral to critical applications, from healthcare to science, management, and beyond, it is imperative to address these disparities and ensure that the benefits of advanced reasoning capabilities are broadly accessible.

<sup>1</sup>We use the term “Reasoning Language Model” instead of “Large Reasoning Model” because the latter implies that such models are always large. This does not necessarily have to be the case – as a matter of fact, smaller RLM can outperform larger LLMs [56].

Reasoning Language Model (RLM): What is it and how to build one?

Basics of RLMS	Essence of RLMS	Blueprint of RLMS	x1 Framework & Insights
§2.1-§2.2 History & main pillars of RLMS	§3 Essence of RLMS: an overview and the most important details of the RLM architecture	§4 Blueprint: a toolbox with ingredients to build various RLMS	§7 Design of the x1 framework: how to easily implement and experiment with RLM designs
§2.3-§2.4 Different categories of RLMS	Fig. 4 An overview and details of the inference, training, and data generation pipelines of RLMS	§5 How existing schemes compare to the blueprint	§7.5 Enabling efficient scaling, modern cloud deployments
Fig. 2 History of RLMS		Appendix A Mathematical specifications of RLMS	§7.6 Example analyses
Fig. 3 Pillars and categories of RLMS		Appendix B Details on value and reward models	§8 Example insights for building effective RLMS
		Appendix C-D Algorithmic formulations of RLMS: how different parts of RLMS work in detail, facilitating implementation	§9 Benchmarks for RLMS
		Fig. 5 Toolbox overview      TABLE 1 RLM comparison	
		§6 Hints on how to use the blueprint for user’s application	

Fig. 1: Summary of the contributions made by this paper. The x1 framework can be found at <https://github.com/spcl/x1>

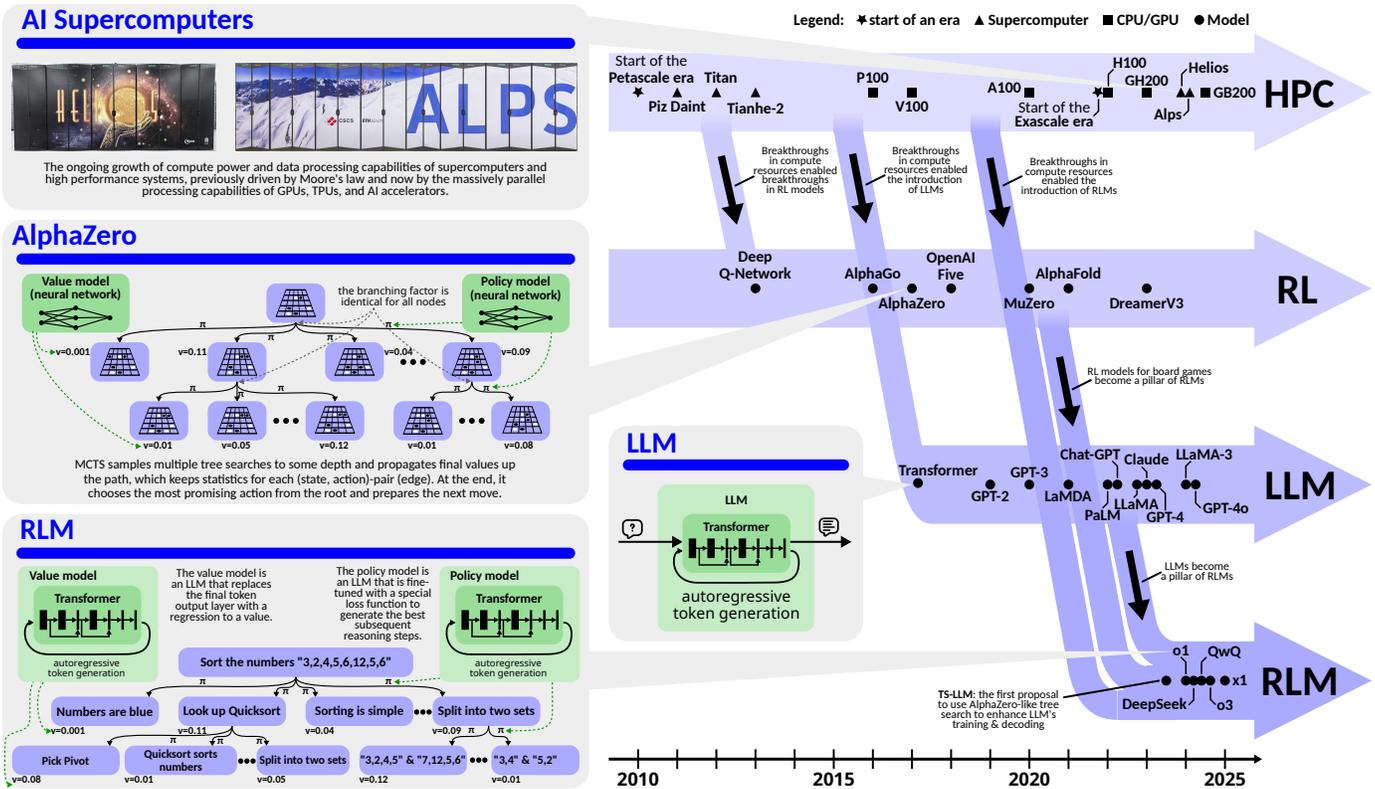


Fig. 2: The history of RLMs. This class of models has been the result of the development of three lines of works: (1) Reinforcement Learning based models such as AlphaZero [134], (2) LLM and Transformer based models such as GPT-4o [115], and (3) the continuous growth of compute power and data processing capabilities of supercomputers and high performance systems.

The technical foundations of RLMs remain opaque and complex, compounding the accessibility challenge. Emerging analyses suggest that their design likely integrates elements such as Monte Carlo Tree Search (MCTS) or Beam Search, reinforcement learning (RL), process-based supervision (PBS) [88], [88], [151], [151], and advanced in-context learning (ICL) techniques like Chain-of-Thought (CoT) [160] or Tree of Thoughts (ToT) [169], and possibly even retrieval-augmented generation (RAG) [13], [57], [83], [84].

Additionally, these architectures employ multiple specialized subcomponents—such as synthetic data generation engines and policy, value, and reward models—trained through some form of novel loss functions and possibly several fine-tuning schemes. However, the intricate interplay of these components and their integration into a cohesive and effective architecture remains poorly understood. Here, the “holy-grail question” is: *what is the detailed design of an RLM and how to make it simultaneously achieve effectiveness (i.e., high accuracy in delivered answers), low cost, and scalability?*

To help answer this question and to address the above challenges, we propose a comprehensive blueprint for constructing, analyzing, and experimenting with RLMs (**contribution #1**; a roadmap of all the contributions and the paper is in Figure 1). Our approach identifies and crystallizes the fundamental building blocks of RLMs, organizing them into a cohesive framework. This blueprint is presented with increasing levels of granularity, starting from high-level overview, finishing at low-level details that can be directly harnessed when implementing. Further, to max-

imize the clarity and comprehensiveness, we present the blueprint using three perspectives: (1) architecture diagrams and descriptions, (2) detailed mathematical formulations, and (3) in-depth algorithmic specifications. By employing these complementary perspectives, we aim to provide a clear and actionable guide for developing RLMs tailored to specific applications, settings, and constraints.

Our blueprint comprehensively encompasses the potential building blocks of RLMs, offering a flexible and modular framework. It incorporates a variety of reasoning structures, such as chains, trees, graphs, and even higher-order structures such as hierarchical (or nested) trees, along with numerous operations that transform and advance the reasoning process. The blueprint supports different granularities of reasoning steps, ranging from individual tokens to full sentences or structured segments. Additionally, it enables diverse training schemes, including Outcome-Based Supervision (OBS) and PBS, and the related Outcome & Process Reward Models (ORMs & PRMs). Next, in order to illustrate the capability of the blueprint to accommodate novel design ideas, we describe several novel schemes and how they fit within the blueprint. One such example is Trace-Based Supervision (TBS), which extends PBS by incorporating labeled traces of traversal paths through entire reasoning structures, rather than just linear chains of reasoning steps. By unifying all these components, our blueprint serves as a versatile toolbox for constructing RLMs—ranging from simple models to sophisticated designs—tailored to specific reasoning tasks and performance objectives.

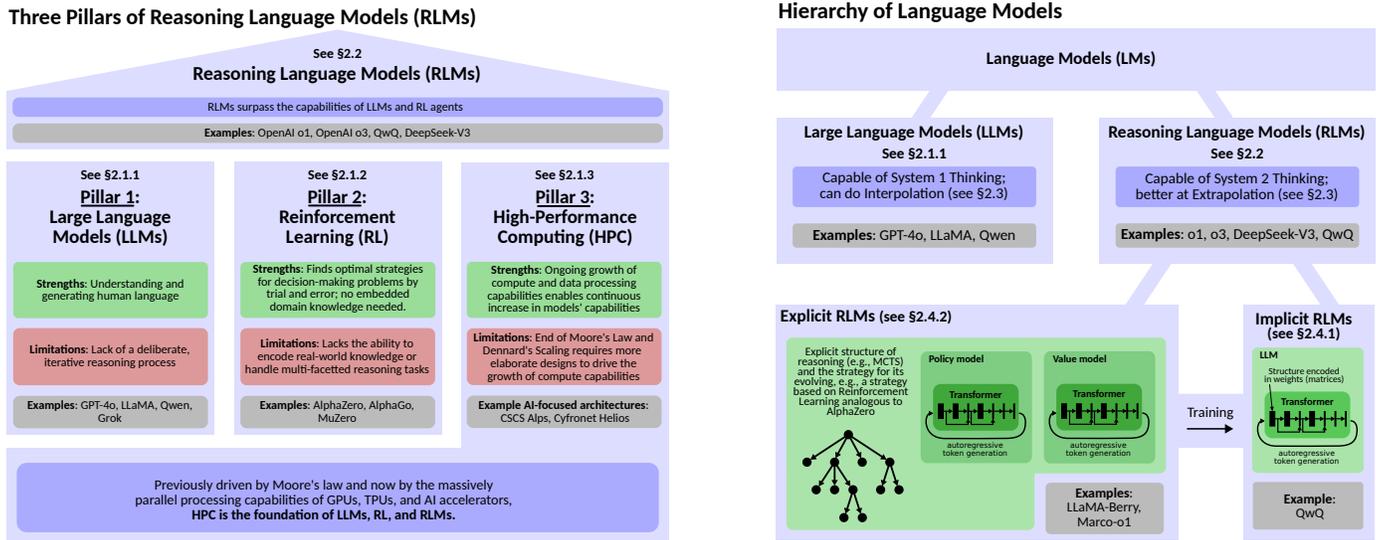


Fig. 3: Hierarchy of language models (right) and the three pillars of RLMs (left).

We conduct a broad analysis of existing reasoning schemes (**contribution #2**), demonstrating how they fit into our blueprint as special cases. This analysis encompasses not only standard MCTS and reinforcement learning-based designs, such as LLaMA-Berry [177], but also models like QwQ [148]. Additionally, we include paradigms diverging from standard MCTS, such as Journey Learning [119] or Beam Search, which redefines reasoning through implicit long-chain structures, and advanced structured prompting techniques like CoT [160], ToT [169], and Graph of Thoughts [9]. We also consider *reasoning utilities* such as Retrieval-Augmented Generation (RAG) and data stores, tools, and others. By mapping these diverse approaches to one blueprint, we showcase its versatility and expressive power, highlighting its ability to unify a wide range of reasoning methodologies within a coherent framework.

To demonstrate the utility of our framework, we introduce **x1**, a modular and user-friendly implementation<sup>2</sup> designed to simplify the process of developing and experimenting with new RLM architectures, covering not only training and inference, but also synthetic data generation (**contribution #3**). We design **x1** to facilitate supporting various optimizations, design decisions, and overall scalability, such as batch processing, making it a well-suited foundation of experimentation infrastructure. We also discuss key aspects of deployment in cloud environments, ensuring that **x1** can be seamlessly integrated into modern infrastructure for both research and production use cases.

By providing both theoretical insights and practical tools, this work aims to democratize access to advanced RLMs, enabling researchers and practitioners to design, train, and deploy sophisticated reasoning models with reduced complexity and cost. Our blueprint offers a clear and adaptable framework that lowers the barriers to entry, fostering broader experimentation and innovation. Additionally, the modular implementation of **x1** serves as a foundation for rapid prototyping and large-scale experimentation, empowering users to explore new reasoning paradigms and

optimize performance across diverse applications. By bridging the gap between conceptual advancements and practical implementations, this work seeks to accelerate progress in the field, unlock new possibilities for intelligent systems across research, industry, and education, and to mitigate the risk of the growing gap between “rich AI” and “poor AI”.

## 2 EVOLUTION & FOUNDATIONS OF RLMs

We first summarize the evolution and foundations of reasoning language models. Figure 2 shows an overview of the history of the development of these models.

### 2.1 Basic Pillars of Reasoning LMs

The development of reasoning-capable LLMs represents a convergence of three critical threads: (1) advances in LLMs such as GPT-4, (2) RL designs such as AlphaZero, and (3) High-Performance Computing (HPC) resources. Together, these threads have shaped models capable of efficient *System 2 Thinking* – a level of reasoning that combines explicit deliberation with novel problem-solving abilities, distinct from the intuitive, fast, and automatic heuristics of *System 1 Thinking*. Figure 2 compares example designs in these pillars while Figure 3 (left side) further discusses the details of these pillars.

#### 2.1.1 Large Language Models: A Reservoir of Knowledge

LLMs such as GPT-4o [115] or Llama [54] represent an extraordinary leap in the field of AI, constituting a vast repository of world knowledge encoded directly in their weights. Trained on huge corpora of text from diverse sources, LLMs are capable of understanding and generating human language with remarkable fluency. However, their reasoning abilities largely align with the fast, automatic, and intuitive *System 1 Thinking*. While they can generate coherent responses and even perform simple reasoning tasks, LLMs have limitations. The reasoning they exhibit is often shallow, rooted in the simple mechanism of predicting the next most probable token in a sequence rather than engaging in explicit problem-solving or structured analysis. While

<sup>2</sup><https://github.com/spcl/x1>

LLMs may generate plausible-sounding solutions to a problem, these outputs are the result of statistical language modeling rather than a deliberate, iterative reasoning process. This distinction highlights the need for integrating more advanced mechanisms capable of explicit reasoning into AI systems—paving the way for hybrid designs that combine the knowledge-rich foundation of LLMs with structured reasoning methodologies.

### 2.1.2 Reinforcement Learning: Exploring and Innovating

RL has historically provided a framework for decision-making and exploration in environments where an agent must learn optimal strategies through trial and error. Landmark systems like AlphaZero [134] and a long line of others such as AlphaGo [133] or MuZero [130] demonstrated the profound potential of RL by achieving superhuman performance in games such as chess, shogi, and Go. Unlike traditional AI systems, AlphaZero began with no embedded domain knowledge. Instead, it mastered these games purely through self-learning, discovering novel strategies that even human experts had not considered.

One of the most striking examples of RL’s innovative capacity came during an AlphaZero match, where the system made a move initially deemed a mistake by human observers. This move [105] later proved to be both surprising and strategically brilliant, illustrating the capacity of RL agents to explore unconventional solutions that lie outside the bounds of human intuition. Such capabilities are fundamentally rooted in RL’s ability to navigate vast search spaces effectively.

However, traditional RL systems lacked the ability to encode real-world knowledge or handle complex, multifaceted reasoning tasks. This limitation spurred the integration of RL principles with LLMs, combining the structured exploration and optimization capabilities of RL with the knowledge-rich reasoning foundation of language models.

### 2.1.3 HPC: Scalability & Efficiency

The growth of LLM and RL systems has been propelled by advancements in High-Performance Computing (HPC). Initially driven by Moore’s Law, which enabled a doubling of transistor density approximately every two years, HPC benefited from both technological advancements and the economic feasibility of manufacturing smaller transistors. However, as the costs of further miniaturization have risen sharply, Moore’s Law has reached practical limits, necessitating alternative strategies like parallelism and heterogeneous computing.

Modern HPC systems rely heavily on GPUs, TPUs, and AI accelerators for their parallel processing capabilities, alongside CPUs for sequential and general-purpose tasks. Heterogeneous computing leverages these components to optimize task-specific performance. Distributed frameworks, employing techniques such as data, model, and pipeline parallelism [8], [12], [16], further enable the training of enormous models across thousands of compute nodes.

Energy efficiency innovations, including sparsity, quantization, and pruning, mitigate the growing energy demands of scaling AI systems. These advancements ensure that HPC remains a cornerstone for developing and deploying AI

models, supporting the combination of vast knowledge, reasoning capabilities, and computational scalability – allowing AI evolution to continue beyond the limits of traditional Moore’s Law scaling.

## 2.2 The Convergence: System 2 Thinking in AI

The intersection of these three threads – LLMs, RL, and HPC – has culminated in the emergence of models capable of System 2 Thinking. These advanced systems combine the knowledge-rich foundation of LLMs with the exploratory and optimization capabilities of RL, all supported by the scalability and performance of modern HPC. The result is a new class of AI models that can engage in explicit, deliberate reasoning processes.

These models possess a world model encoded in the weights of their LLM components, allowing them to reason about complex scenarios and contexts. Their RL capabilities combined with the HPC capabilities enable them to navigate truly immense decision spaces, evaluate multiple strategies, and iteratively refine solutions.

## 2.3 Interpolation (LLMs) vs. Extrapolation (RLMs)

Standard LLMs, driven by their autoregressive token prediction mechanism, primarily perform interpolation within the vast search space of solutions. They excel at generating responses that align with patterns seen in their training data, effectively synthesizing knowledge from known contexts. However, this process limits them to producing outputs that remain within the boundaries of their training distribution. In contrast, reasoning LMs enable extrapolation beyond these boundaries. By combining structured exploration, reasoning LMs navigate uncharted areas of the solution space, generating novel insights and solutions that extend past the limits of their training data. This enables a shift from basic pattern completion to active problem-solving.

## 2.4 Hierarchy of Reasoning-Related Models

The evolution of RLMs can be understood as a hierarchical progression, with earlier models such as GPT-4o being less capable in terms of reasoning, and the o1-like architectures demonstrating increasing sophistication and explicit reasoning abilities. This hierarchy reflects the integration of System 1 (LLMs) and System 2 (RLMs) Thinking. RLMs can be further divided based on how reasoning is implemented into *Implicit RLMs* and *Explicit RLMs*; the details of this categorization can be found in Figure 3 (the right side).

### 2.4.1 Implicit Reasoning Models

In this subclass, the reasoning structure is embedded entirely within the model’s weights. Models such as QwQ [148] operate as “black boxes”, where reasoning is implicit and cannot be explicitly disentangled or manipulated. While these models exhibit improved reasoning capabilities compared to standard LLMs, their reasoning processes are opaque and rely on the internalized patterns learned during training.

### 2.4.2 Explicit Reasoning Models

These models introduce explicit reasoning mechanisms external to the model’s core weights. Examples include designs such as LLaMA-Berry [177], Marco-o1 [182], and potentially OpenAI’s o3, which incorporate mechanisms like explicit MCTS combined with RL for decision-making. This explicit structure enables the model to simulate, evaluate, and refine solutions iteratively, facilitating novel problem-solving and extrapolation. By separating reasoning from the static knowledge encoded in the weights, these models achieve greater flexibility and interpretability in their reasoning processes. Note that the explicit reasoning can be internalized via training making it implicit – we discuss it later in the blueprint.

## 3 ESSENCE OF REASONING LMS

We now describe the general architecture of RLMs, which we summarize in Figure 4. In the following sections, we generalize this description to the full RLM blueprint.

### 3.1 Basic Architecture, Pipelines, & Concepts

We now outline the foundational architecture, operational pipelines, and core concepts. Figure 4 offers three levels of detail. In general (the top-left part), the whole RLM architecture consists of three main pipelines: inference, training, and data generation. The inference serves user requests, using models (e.g., the value or policy model) provided by the training pipeline. Data generation mirrors the inference pipeline in its internal design; the main difference is that it runs independently of the user requests, generating data that is then used to re-train the models. As such, training combined with data generation from various domains [127], [176] offers *self-learning* capabilities and is analogous to the self-play setting of AlphaZero [134].

#### 3.1.1 Inference

The inference process begins when the user provides an input prompt **1**, which typically describes the problem or question to be addressed by the RLM. This input serves as the root of the reasoning process and initiates the construction of a **reasoning structure 2** that organizes RLM’s progress. The structure is usually represented as a tree. The root of this tree corresponds to the user’s input, and subsequent nodes are generated to explore the search space – the domain of possible reasoning paths or solutions. The purpose of this reasoning structure is to systematically investigate potential solutions, progressively refining and extending reasoning paths to converge on an optimal or satisfactory answer.

An individual **point** in the search space, represented as a **node** in the reasoning structure, corresponds to a **reasoning step 3**. A reasoning step is defined as a coherent and self-contained unit of thought – a sequence of tokens that advances the solution by either exploring a new branch of the problem or building upon existing progress. These steps form the building blocks of the reasoning process.

The details of how the structure evolves are usually governed by the **MCTS scheme**, enhanced with **policy and value models** (we also distinguish other **reasoning**

**strategies**, described below). This approach, inspired by methods used in AlphaZero, ensures that the search process is both efficient and directed toward promising solutions. The **policy model 4** is responsible for generating new reasoning steps at each node, predicting the next most likely and logical steps to expand the reasoning process. Meanwhile, the **value model 5** evaluates the quality of a reasoning path starting at a given node, helping the system prioritize the most promising steps to follow. Sometimes, a reward model<sup>3</sup> **6** is used instead, to assess the quality of an *individual* specific node and its corresponding reasoning step. In our blueprint, as detailed in the next section, we abstract the models into a more general notion of **operators 7** to enable more flexibility in how they are implemented.

The search and reasoning processes continue iteratively until a **terminal step** is reached **8**. This terminal step represents a completion of the reasoning chain that forms the final answer to the posed problem. It serves as the leaf node in the tree, concluding that particular reasoning path.

This architecture provides a unified framework that accommodates a wide range of reasoning tasks. Whether reasoning steps are fine-grained (e.g., individual token sequences) or coarse-grained (e.g., entire reasoning chains treated as single nodes), the architecture adapts seamlessly. By structuring the search space explicitly and guiding exploration with policy and value models, the RLM achieves a level of reasoning capability bridging intuitive pattern recognition and deliberate problem-solving.

A detailed specification of the inference pipeline can be found in Appendix C.1 and in Algorithm 1.

#### 3.1.2 Training

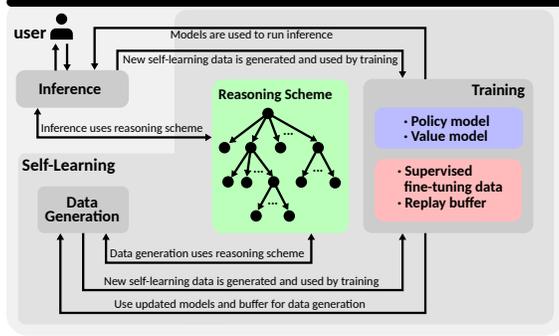
Training details depend on what model is trained (value, policy, reward, ...). In general, we assume fine-tuning a model such as Llama. Here, we follow an approach where one first harnesses supervised data, usually coming from existing datasets such as PRM800K [88] **1**, which becomes a part of the supervised training data **2** used in the **supervised training pipeline 3** of the framework to train some, or all, of the models **4** considered in the blueprint. The second part of the overall training framework in RLMs is the **unsupervised (self-learning) training pipeline**, in which training data is being continually generated **5** and used to improve the models. The data can be obtained from inference, assuming quality control [56], but also from a dedicated synthetic data generation pipeline that mirrors that of the inference. To collect the data, one executes the respective RLM pipeline for a given input task and gathers the results **6**; depending on how detailed the gathering process is, the data collected can contain only **outcome-based labels 7**, **process-based labels 8**, or some other variant such as **trace-based labels 9** suggested in our blueprint, that generalize process-based samples to samples that contain also information about operators applied during the task solution process. All this data becomes a part of the replay buffer **10** and is used in the unsupervised training

<sup>3</sup>We use a naming scheme in which a model used to estimate the quality of a whole reasoning path starting at a given node, is called the *value model*, while a model used to estimate the quality of a given reasoning step, is called the *reward model*.

**Legend**

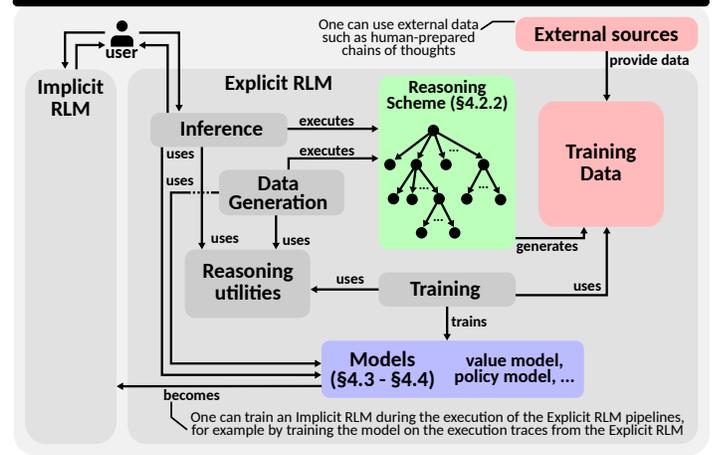
- Part of the pipeline
- Reasoning scheme
- Models & Operators
- Training Data
- 1 References to descriptions in text (inference pipeline)
- 2 References to descriptions in text (training pipelines)

**High-level overview (§3.1)**



More details

**Medium-level overview (§3.1)**



More details

**Detailed view (§3.1.1, §3.1.2)**

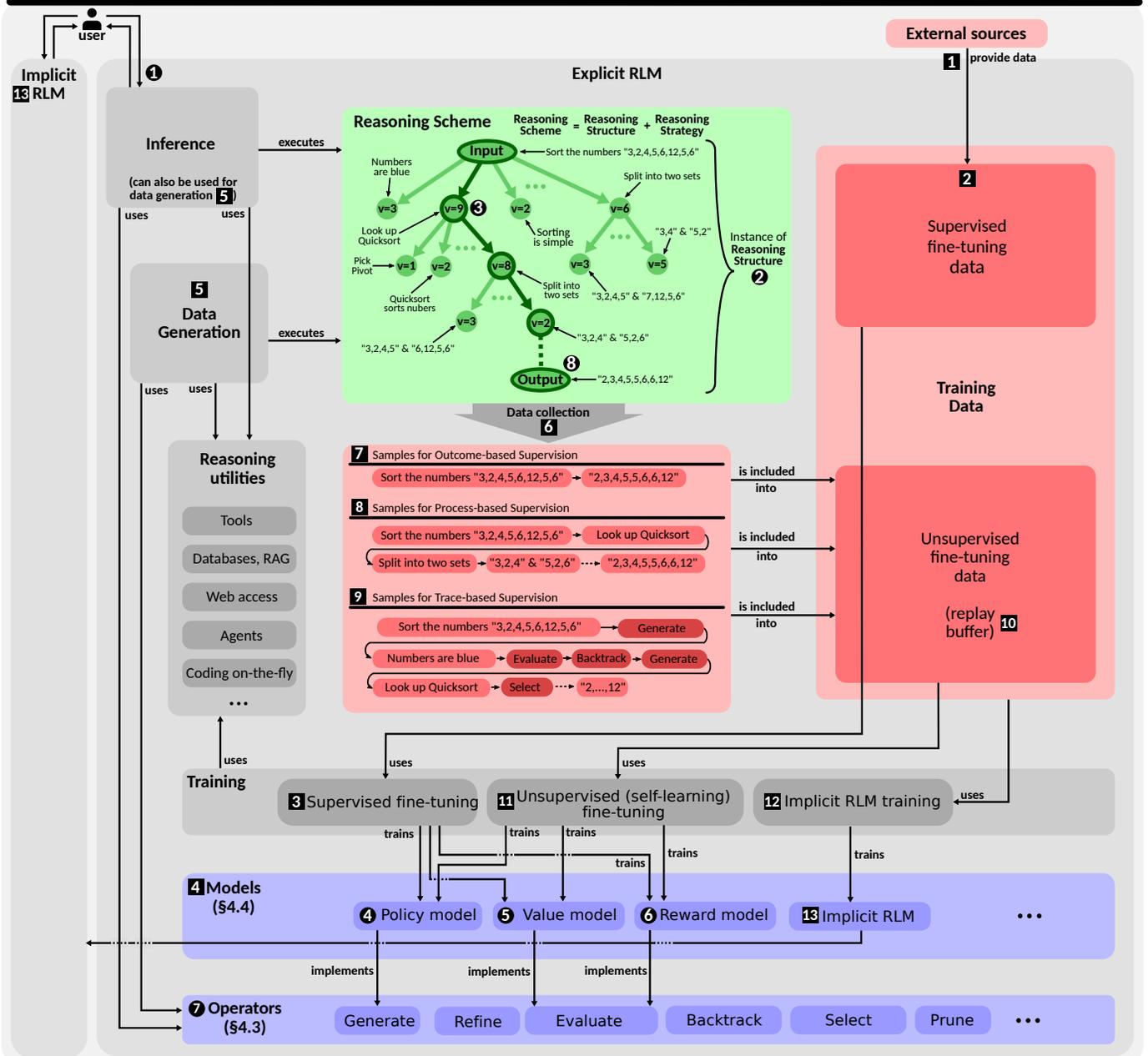


Fig. 4: Overview of a general RLM design and core concepts. We provide a high-level overview (the top-left part), a more detailed medium-level overview (the top-right part), and a very detailed diagram showing the inference and training pipelines (the bottom part). A detailed specification of the inference pipeline can be found in Appendix C.1 and in Algorithm 1. Details on the pipelines for different training phases and paradigms can be found in Appendices C.2 and C.3 as well as in Algorithms 2–7. The data generation pipeline is detailed in Appendix D.

scheme [11](#) or it can also be used to train [12](#) a model that would become an **Implicit RLM** [13](#).

A detailed specification of the pipelines for different training phases and paradigms can be found in Appendices C.2 and C.3 as well as in Algorithms 2–7. The data generation pipeline is detailed in Appendix D.

### 3.2 Encompassing Diverse RLM Architectures

The above-described design is applicable to many RLM designs. However, there are numerous other variants of architectures, some of which do not fully conform to this framework. In this section, we discuss these variants, highlighting how our blueprint accommodates such variations.

In some RLM designs [177], a single node in the MCTS tree could represent *an entire reasoning structure*, such as a complete chain of reasoning steps. In this case, the action space involves transitioning between different reasoning structures rather than individual steps. This approach changes the nature of the search, as the focus shifts from iteratively constructing a single reasoning path to evaluating and refining entire structures within the search space. Our blueprint accommodates this with the concept of **nesting**, where a node in the reasoning structure can contain another reasoning structure.

Other architectures introduce even more novel paradigms. For instance, Journey Learning [119] adds an additional layer of complexity by incorporating a transformation step that “rewires” the search or reasoning structure. This transformation consolidates multiple paths in the tree, synthesizing them into a new form that is used as input for subsequent reasoning iterations.

Despite these variations, our blueprint is sufficiently general to encompass all these cases and beyond, as we illustrate more formally in the following. This generality ensures that the blueprint is not only applicable to existing designs but also provides a foundation for future innovations in RLM development.

### 3.3 Integration with Broader LLM Agent Ecosystems

The integration of RLMs into broader LLM agent ecosystems would enable these models to interact dynamically with external tools, databases, and resources during execution. This interaction can occur within the inference or data generation pipeline, leveraging value or policy models to extend the reasoning process through access to retrieval-augmented generation (RAG), web queries, and specialized tools. For example, during a reasoning task, the value or the reward model could query a database to verify intermediate steps, ensuring factual correctness or retrieving additional context to refine its reasoning. Similarly, these models could utilize computational tools for mathematical or symbolic computations, thereby expanding the scope and accuracy of their reasoning.

## 4 BLUEPRINT FOR REASONING LMS

We now introduce our RLM blueprint that can be used to develop novel reasoning models and to provide ground for analysis, evaluation, and comparison of such designs. We overview the blueprint in Figure 5.

### 4.1 Overview & Main Components

The blueprint specifies a toolbox of components that can be used to build an arbitrary RLM. We identify several classes of such components. First, an RLM includes a **reasoning scheme**, which specifies a **reasoning structure** (e.g., a tree) together with a **reasoning strategy** (e.g., MCTS) of how this structure evolves in order to solve a given input task. Second, there is a set of **operators** (e.g., Refine) that can be applied to the reasoning structure (as specified by the reasoning strategy) in order to evolve it and make progress towards solving the input task. Operators are specified based on *what they do* (i.e., what effect they have on the reasoning structure). *How* this effect is achieved, depends on how a given operator is implemented. Here, many operators rely on neural **models** (e.g., Policy Model), which – together with their training paradigms – form the third class of the blueprint components. Finally, we also distinguish a set of **pipelines**, i.e., *detailed specifications of operations* that orchestrate the interaction between the reasoning scheme and the operators in order to achieve a specific objective, such as training, inference, or data generation. *Hence, an RLM can be defined as a composition of a reasoning scheme, a set of operators and associated models, and a set of pipelines.*

### 4.2 Reasoning Scheme

A reasoning scheme is the part of the blueprint that specifies the details of the reasoning steps progressing toward the solution, how they are interconnected to form coherent chains, trees, or more complex reasoning structures, and how these structures evolve in the course of solving the input task.

#### 4.2.1 Reasoning Step

A reasoning step is a fundamental unit of the reasoning structure – a sequence of tokens that advances the RLM towards the solution. Reasoning steps can vary in length, ranging from a **single token** to **entire segments** of text. The variability in their granularity depends on the user design choice. In existing schemes, a reasoning step is typically conceptualized as a “coherent and self-contained unit of thought”. For instance, in mathematical proofs, this may correspond to an individual logical argument or deduction.

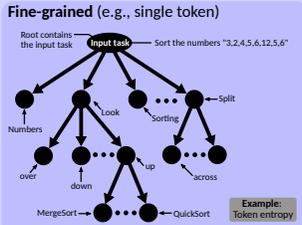
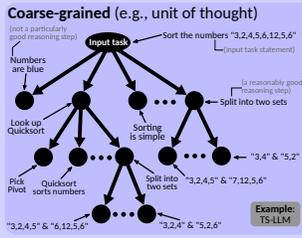
The flexibility in defining reasoning steps allows models to adapt to different problem domains, balancing fine-grained and coarse-grained reasoning. Coarse steps, such as logical arguments (or even complete reasoning pathways [177]), simplify preparation and adoption of training data, enhance interpretability, and – as we discuss in Section 8 – reduce computational overhead. On the other hand, single-token steps enable the utilization of concepts like token entropy [101] to incorporate the model’s uncertainty, as well as the integration of advanced decoding schemes (e.g., speculative decoding [82] or contrastive decoding [85]) explicitly into the RLM design. Yet, while making the reasoning steps more fine-grained allows for a more detailed exploration of solution paths, this increased flexibility results in greater computational demands, particularly when combined with search algorithms such as MCTS.

# 1 Reasoning Scheme (§4.2)

A toolbox of paradigms for modeling and evolving the reasoning structure

## 1.1 Reasoning Step (§4.2.1)

What is the content of an individual reasoning step?

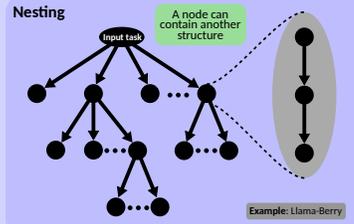
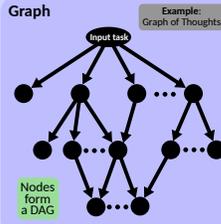
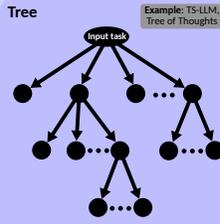


### Decoding Strategy (§4.2.4)

- Greedy search
- Nucleus sampling
- ...

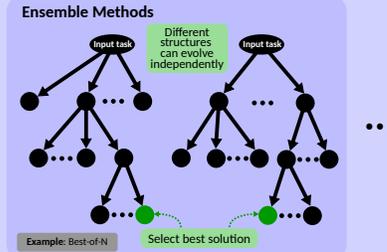
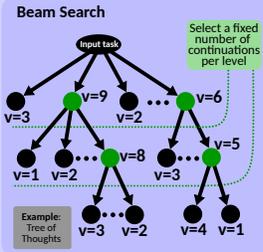
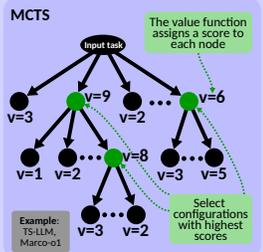
## 1.2 Reasoning Structure (§4.2.2)

What is the connection structure of reasoning steps?



## 1.3 Reasoning Strategy (§4.2.3)

How does the reasoning structure evolve in order to progress solving the input task?



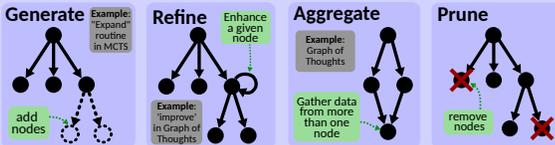
## 4 Pipelines Inference: §3.1.1, Appendix 3.1 Training: §3.1.2, Appendix 3.2 - 3.4 Data generation: Appendix D

# 2 Operators (§4.3)

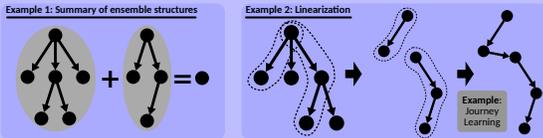
A toolbox of operations for changing & interacting with the reasoning structure

## 2.1 Structure Operators (§4.3.1)

Modify the reasoning structure

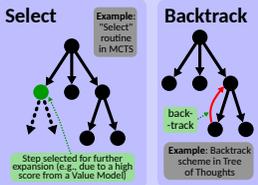


**Restructure** (apply arbitrary structural transformations)



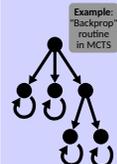
## 2.2 Traversal Operators (§4.3.2)

Specify which nodes to select for next operation



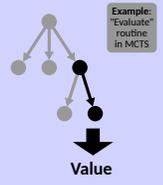
## 2.3 Update (§4.3.3)

Modify the nodes' contents but not the structure



## 2.4 Evaluate (§4.3.4)

Map the structure to values



# 3 Models (§4.4)

A toolbox of neural models for implementing operators and of paradigms for training these models

## 3.1 Models Harnessed (§4.4)

What operators are being implemented as models?

- Value model
- Policy model
- Reward model
- ...
- More details on models in Appendix B

## 3.2 Training Paradigm (§4.4.1)

How is a given model being trained?

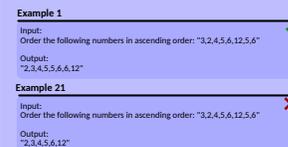
- Rejection Sampling
- Proximal Policy Optimization
- Direct Preference Optimization
- Supervised fine-tuning
- Reasoning Policy Optimization
- More details on training in Appendix C

## 3.3 Training Data Scope (§4.4.2)

What information does a single training sample contain?

### Outcome-based supervision

Training samples only contain inputs and outputs as well as a label that is either correct (✓) or incorrect (✗).



### Process-based supervision

Training samples contain all intermediate steps between input and output, annotated with a quality score (q) or a label that is either correct (✓) or incorrect (✗).

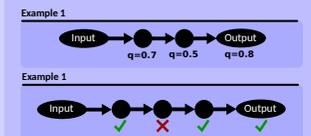


Fig. 5: A blueprint for reasoning LMs. It consists of four main toolboxes: the reasoning scheme (the top part), operators (the bottom-left part), and models (the bottom-right part); pipelines are mentioned in the center and detailed in Appendix C.1 and in Algorithm 1 (the inference pipeline), Appendix C.2, Appendix C.3, and in Algorithms 2-7 (the training pipelines), and in Appendix D (the data generation pipeline).

### 4.2.2 Reasoning Structure

The reasoning structure specifies how individual reasoning steps are connected and organized. Common structures include chains (linear sequences), trees (hierarchical branching), and graphs (arbitrary connections).

**Chains** are sequential reasoning flows, where each step builds directly on the preceding one. Chain structures are prevalent in CoT-based models, where each reasoning step follows logically from the previous step in a linear progression. In **tree** structures, each reasoning step can branch into multiple continuations, forming a decision tree. This structure is commonly used in MCTS-based frameworks, where multiple potential paths are explored before selecting a branch that will be further investigated. It enables more effective exploration of the space of reasoning steps, but simultaneously makes the RLM design more complex and costly. Finally, **graph** structures allow for arbitrary dependencies between reasoning steps, enabling graph-based reasoning, such as that found in the Graph of Thoughts (GoT) framework [9].

Further generalization involves **nested structures**, where reasoning nodes themselves may contain substructures. For example, a node in a tree structure may represent a CoT chain, as proposed in LLaMa-Berry [177]. This hierarchical organization could be particularly useful for multi-step tasks where high-level decisions guide low-level computations, such as meta-reasoning frameworks [177]. One could harness any other *higher-order structures*, such as hypergraphs, motifs, and others [10], [11], [14], [17].

### 4.2.3 Reasoning Strategy

The reasoning strategy governs how the reasoning structure evolves, specifying the process by which new reasoning steps are added and integrated. Example strategies include:

- **MCTS** [77] A popular approach that balances exploration and exploitation by simulating multiple reasoning paths and selecting the most promising one based on a scoring function.
- **Beam Search** [137] A breadth-limited search that keeps a fixed number of top-ranked continuations at each step. While commonly used for decoding token sequences, beam search can also apply to reasoning steps.
- **Ensemble Methods** These methods involve aggregating multiple independent reasoning strategies, such as combining chains and trees to enhance robustness and accuracy. One example is Best-of-N [48], [158] – a strategy where multiple independent reasoning paths are generated, and the most effective solution is selected based on predefined criteria, e.g., accuracy or completeness. Another example is tree ensemble (Forest) [18] where, instead of a single reasoning tree, a reasoning “forest” consists of multiple disconnected trees, which may eventually converge at a shared solution node. This approach supports diverse reasoning pathways that parallelize exploration.

**Reasoning Strategy vs. Decoding Strategy.** It is crucial to distinguish reasoning strategies from token-level decoding strategies. While decoding strategies, such as greedy search and nucleus sampling [64], generate the internal token sequences within a reasoning step, reasoning strategies focus on the higher-level process of integrating and expanding reasoning steps within the reasoning structure.

## 4.3 Operators

Operators specify operations that can be applied to various parts of the reasoning structure to progress the reasoning process. We now provide an extensive toolbox of operators. Many of them have been widely used in RLM-related designs, but some – to our best knowledge – are still unexplored, we include them to foster innovation and propel the design of more effective and more efficient RLMs.

### 4.3.1 Structure Operators

Structure operators transform the reasoning structure by taking it as input and producing a modified version, typically through addition or refinement of reasoning steps. For instance, they may add new children to a specific node, facilitating the exploration of alternative reasoning paths.

- **Generate** The Generate operator adds one or more new reasoning steps to a reasoning structure. Within the MCTS reasoning strategy, this operator is typically implemented as a policy model to generate new steps. In other strategies, the generation operator may involve sequentially appending steps (CoT) or exploring multiple candidate steps in parallel (Beam Search).
- **Refine** The Refine operator enhances a given individual reasoning step. For example, it could address ambiguities, correct errors, and optimize inefficiencies, resulting in a more robust version of the step [99]. It could also integrate suggestions from self-critique [128] (evaluates steps to identify weaknesses and suggest targeted improvements), summarization [186] (condenses key elements into concise representations to streamline the reasoning structure), or rephrasing [43] (reformulates steps to improve clarity and coherence while preserving their logical integrity).
- **Aggregate** This operator combines multiple reasoning steps, paths, or structures into the next individual step. This enables consolidating information or improving coherence. It is used in Ensemble Methods [18] or in Graph of Thoughts [9].
- **Prune** This operator removes nodes or reasoning steps from the structure that are deemed suboptimal or irrelevant based on evaluation metrics. It enables optimizing the reasoning structure in order to, e.g., reduce token costs.
- **Restructure** The Restructure operator applies arbitrary transformations to the reasoning structure, enabling flexible reorganization of its components. A notable example is the conversion of a reasoning tree into a linear chain by rearranging its branches into a sequential series of steps, as done in Journey Learning [119]. This restructuring facilitates the integration of insights from diverse branches into a cohesive flow, “flattening” it and making it easier for the model to process and utilize information within a single, unified context.

**Discussion on Diversity** In structure operators, there is a notion of how *diverse* the outcomes of the operator are. For example, when generating  $k$  new reasoning steps, one may want to make the contents of these steps as different to one another as possible. While different mechanisms to steer diversity exist, a typical approach is the use of the **policy model temperature**. We additionally propose to consider the **Diverse Beam Search** [152] which promotes diversity by maintaining multiple diverse candidate sequences during

decoding. In MCTS, there is also a distinction between exploitation (expanding the structure by applying generation operators within an already established tree branch) and exploration (generating new branches). Here, one impacts diversity by manipulating the **exploitation-exploration trade-off**, as determined by the Upper Confidence Bound for Trees (UCT) formula [77] or its variants.

#### 4.3.2 Traversal Operators

Traversal operators define how the reasoning process navigates through the *existing* reasoning structure. These operators play a crucial role in shaping the flow of reasoning by determining which paths to pursue.

- **Select** The Select operator determines which reasoning steps to pick for further exploration, evaluation, or refinement within the reasoning process. It evaluates existing elements based on predefined criteria, such as heuristic scores, likelihood estimates, performance metrics or search strategies like PUCT [123] or UCT [77], selecting the most promising candidates to guide the next stages of reasoning. By balancing exploration (considering diverse alternatives) and exploitation (focusing on high-potential paths), the selection operator optimizes resource allocation and ensures efficient reasoning progression.
- **Backtrack** The Backtrack operator enables the model to explicitly return to a previous reasoning step and continue along a different reasoning path. This operator supports error correction, divergence handling, and hypothesis revision by abandoning unproductive directions in favor of alternative trajectories. The QwQ model output indicates that the reasoning structures used as training data in this model harnessed Backtrack.

#### 4.3.3 Update Operators

The Update operator enhances specific parts of the reasoning structure without altering the structure itself. A common example is the backpropagation phase in MCTS, where evaluation scores are propagated and updated along existing reasoning steps to inform future decisions. Another form of update involves refining the content of individual nodes or subsets of nodes, replacing their original versions with improved iterations, such as the “enhance” thought transformation in Graph of Thoughts [9].

#### 4.3.4 Evaluate Operators

Evaluate operators take as input a segment of the reasoning structure and output a value without any modifications to the structure. They are widely used with reasoning strategies, such as MCTS.

One important type of evaluation occurs when the reasoning structure reaches a terminal state, allowing the full reasoning sequence to be assessed against a known solution—applicable to tasks with definitive answers, such as mathematical problems. This **terminality evaluation** verifies whether the final step provides a correct and complete solution.

One can also **evaluate intermediate steps** (i.e., non-terminal ones). This can involve estimating the reward associated with specific reasoning steps, using heuristics, aggregated simulation outcomes, or a trained reward model

for more efficient assessments. Other methods such as embedding-based verification could also potentially be harnessed [15].

Another form of evaluation employs a **value estimator**, which judges a given reasoning step based on its expected contribution to a correct final outcome. This method evaluates both the correctness of the step and its alignment with the overall solution goal. Such evaluations can be performed through simulations, as in the original MCTS algorithm, or more efficiently using a learned value model [135].

A critical aspect of evaluation is the selection of **appropriate metrics**. For instance, in value estimation, an ideal metric considers both the correctness of a reasoning step and the extent of progress it represents toward the final solution, ensuring a balanced assessment of its contribution.

#### 4.3.5 Discussion: Test-Time Compute

One of the recent trends in next-generation LLMs [100], [153] is to shift from merely increasing model sizes to enhancing computational strategies during inference, a concept known as the test-time compute (TTC). This approach allocates additional computational resources during a model’s execution to improve performance, particularly in complex reasoning tasks. This methodology mirrors human cognitive processes, where increased deliberation is applied to more challenging problems.

Recent studies [137] indicate that optimizing test-time compute can be more effective than merely increasing model size. For instance, employing a compute-optimal strategy—where computational resources are adaptively allocated based on the problem’s complexity—can enhance efficiency by over four times compared to traditional methods. Moreover, in scenarios where smaller base models achieve moderate success rates, augmenting test-time compute enables them to outperform models up to 14 times larger.

While test-time compute offers significant benefits, it also presents challenges, related to – among others – resource allocation (determining the optimal amount of computational resources for each inference task requires sophisticated strategies to balance performance gains against computational costs), dynamic scaling (implementing adaptive compute strategies necessitates models capable of assessing problem difficulty in real-time and adjusting their computational efforts accordingly) [102], and hardware implications (the shift towards increased test-time computation may influence hardware requirements, putting more pressure on delivering specialized inference-focused hardware solutions).

**Test-Time Compute in the Context of the Blueprint.** Our blueprint offers mechanisms to dynamically allocate computational resources during inference to improve performance, particularly for more complex problems. By leveraging the modular structure of the blueprint, TTC can be effectively implemented through specific operators designed for reasoning tasks. We now provide several examples.

- The **Generate operator** can be used to implement TTC by dynamically increasing the number of next reasoning steps generated for harder problems. For simpler tasks, the operator may only generate a minimal set of continuations. However, for more complex problems, the operator can

be used to create a larger set of potential reasoning steps, thereby expanding the search space.

- The **Refine operator** provides another avenue for implementing TTC by enhancing a given reasoning step multiple times for harder problems. In this approach, the operator iteratively improves the quality of a reasoning step, addressing ambiguities, rectifying errors, or improving clarity. For simpler tasks, the operator might only refine a step once, while for more complex reasoning, it can perform multiple enhancement iterations to ensure the output meets a higher standard of precision and robustness.
- The **Traversal operators**, such as Select, enable the exploration of multiple reasoning paths at test time, offering another key mechanism for implementing TTC [179]. By using Select on several next reasoning steps, the model can dynamically expand its search tree for more challenging problems, thereby increasing the diversity and depth of reasoning paths under consideration. For example, in a complex task, the model might select multiple high-probability steps and explore their corresponding continuations in parallel. This approach facilitates broader exploration of the reasoning space, ensuring that promising paths are not prematurely discarded.
- To efficiently manage the expanded set of possibilities, the blueprint allows integration with the **Aggregate operator**. This operator evaluates the generated reasoning paths and selects the most promising ones based on predefined criteria, such as the likelihood of correctness or the quality of intermediate steps. This combination ensures that while more computational resources are allocated for challenging tasks, only the most relevant paths are explored further, optimizing both accuracy and efficiency.

## 4.4 Models

Models are used to implement various types of operators. Most common are the **value model** (implementing the value evaluation operator) and the **policy model** (implementing the generate operator).

Models are further categorized and discussed in detail in Appendix B; we discuss the variants of the value model (**Q Value model**, **V Value model**), we compare Process Reward and Outcome Reward models, and we formally identify a new variant of models, the **Outcome-Driven Process Reward Model**.

### 4.4.1 Training Paradigm

Each model must be trained according to a specified paradigm, which outlines the methodology for optimizing its performance. This paradigm defines key training components such as the loss function, data generation and labeling procedures, and other critical training details.

A wide range of training schemes has been developed for models used in RLs, with early foundational work stemming from advancements related to AlphaZero. These schemes have since evolved to support the complex requirements of reasoning tasks within LLMs. Common training paradigms include **supervised fine-tuning (SFT)**, where models are trained on reasoning sequences labeled with q-values; **rejection sampling** [23], [140], which involves filtering generated outputs based on quality criteria; and

**RL-based methods** such as **Proximal Policy Optimization (PPO)** [131], **Direct Preference Optimization (DPO)** [121], and reasoning-specific variants like **Reasoning Policy Optimization (RPO)** [117]. Several training paradigms also incorporate **self-learning**, where the model iteratively improves by generating and evaluating its own reasoning sequences, thereby simulating competitive or cooperative reasoning scenarios.

### 4.4.2 Training Data Scope

The training data for RLs can vary significantly in terms of how much of the reasoning structure it captures. We now outline two established approaches, **outcome-based supervision (OBS)** and **process-based supervision (PBS)**. More details regarding both OBS and PBS can be found in Appendix B.1.

In **outcome-based supervision** (also known as a **sparse training signal**) [36], [151] each training sample consists solely of the input and the corresponding output. For example, in mathematical problem-solving, a sample may include the task statement and the final solution, labeled as correct or incorrect. This approach is straightforward to implement, and the required data is relatively easy to collect. However, it can limit the model’s reasoning accuracy, as it provides minimal insight into the intermediate steps that led to the solution [88].

An alternative approach is **process-based supervision** (also known as a **dense training signal**) [88], [155], where a training sample reflects the entire reasoning structure. In this case, the sample contains not only the input and final output but also all intermediate reasoning steps, annotated with labels indicating the quality of each step. This richer training data allows the model to learn more granular reasoning patterns, improving its ability to generate accurate and interpretable solutions by understanding the reasoning process in detail. However, such data is much more challenging to generate or gather [88].

**OBS vs. PBS** By varying the training data scope, developers can strike a balance between ease of data collection and the depth of reasoning insights provided to the model, with dense supervision generally offering improved performance at the cost of increased data complexity. We detail these, and additional aspects of ORMs and PRMs in Pipelines for different training phases and paradigms can be found in Appendix B, Appendix C.2, Appendix C.3, and in Algorithms 2–7.

**Trace-based supervision (TBS)** is a potential way to extend PBS by incorporating detailed information about the sequence of applied operators, including traversal operators, within the reasoning structure. By capturing the full trace of how reasoning steps are generated, refined, or revisited, TBS would provide richer supervision that teaches the model to internalize not just the reasoning steps but also the process of navigating and manipulating the reasoning structure itself. This approach could enable the training of more powerful Implicit RLs by guiding them to replicate the reasoning dynamics of explicit structures, improving their ability to reason flexibly and efficiently.

## 4.5 Pipelines

A pipeline is a detailed specification of operations that orchestrates the details of the interaction between the reasoning scheme and the operators and models to achieve a specific objective. Typically, an RLM would incorporate a single **pipeline for inference** and a separate **pipeline for training** each model used in an RLM. Moreover, there could also be **pipelines for synthetic data generation** used for training models. One can also distinguish a pipeline that trains an Implicit RLM using the provided reasoning traces from the Explicit RLM.

The details of pipelines depend on arbitrary design choices. In Section 3, we provided a general description of how these pipelines work. In Appendix C, we present detailed algorithmic specifications of our pipelines, along with insights into the reasoning behind these design choices. Specifically, the inference pipeline can be found in Appendix C.1 and in Algorithm 1. Pipelines for different training phases and paradigms can be found in Appendix C.2, Appendix C.3, and in Algorithms 2–7. The data generation pipeline is detailed in Appendix D.

## 5 EXPRESSING EXISTING SCHEMES

We now showcase the expressivity of our blueprint, by illustrating how it can be used to model a broad scope of existing RLMs and other related works. We summarize the outcomes of the analysis in Table 1. We start with typical and most prevalent Explicit RLM architectures based on MCTS and policy and/or value models, where a single reasoning step is an individual logical argument (Section 5.1). We also discuss there schemes that generalize this typical design, by harnessing nesting or Linearization Structure operators. Finally, we study Implicit RLMs (Section 5.2) and various structured prompting schemes such as CoT or ToT (Section 5.3), showing that they also fit our blueprint.

### 5.1 Explicit RLMs

We start with the most widespread variant of RLMs that follows the architecture outlined in Section 3.1. These reasoning models such as TS-LLM [48], AlphaLLM [149], MCTS-DPO [163], and others [24], [56], [153], [177], [178], [182] generally employ an explicit tree structure in which a node represents a distinct reasoning step. The reasoning strategy is based on the MCTS and focuses on iterative exploration, expansion and evaluation of nodes within the tree. By incorporating value mechanisms—such as prompt-based evaluation or dedicated value models, the system identifies and prioritizes promising branches, facilitating more informed decision-making and refinement of the reasoning process. All MCTS based reasoning models implement at least a next-step generation operator, an evaluation operator, and the update operator for back-propagating the values. In addition, ReST-MCTS\*, LLaMA-Berry, and Marco-o1 support a refinement operator to further improve produced reasoning steps.

**Journey Learning** [119] exhibits two main differences to typical MCTS-based RLMs. First, it harnesses the Linearization Structure operator, in which the tree reasoning structure is transformed into a chain, by extracting several selected

reasoning chains from it and combining them together into an individual long chain. This way, the scheme attempts to harness insights from different tree branches. By maintaining a chain-based structure, Journey Learning preserves the simplicity of linear reasoning while embedding the capacity for self-correction and exploration of multiple hypotheses. Additionally, Journey Learning introduces a pipeline for the internalization of such long reasoning chains into its weights. This enables the final model to generate such long reasoning chains, possibly containing different reasoning branches, directly from its weights, making it an implicit RLM.

### 5.2 Implicit RLMs

**Qwens’s QwQ** [148] embodies a fully implicit reasoning model, characterized by an implicit reasoning structure that is generated autoregressively directly by the model weights. The reasoning strategy in QwQ – as indicated by the model output – harnesses next-step generation, backtracking, summarization, and critique generation to derive the final solution. At each step, the model implicitly generates a new node within the chain by employing one of these four implicit generate operators, presumably implemented using special tokens.

### 5.3 Structured Prompting Schemes

Finally, we also illustrate that advanced *structured* prompting schemes, such as CoT, ToT, and GoT, constitute a fully explicit RLM structure without any implicit reasoning than what is originally presented in the used LLM, i.e., no models nor training or data generation pipelines.

**CoT** [160] utilizes an implicit reasoning structure consisting of a chain of reasoning steps. The reasoning strategy employed in CoT is oriented towards constructing a single coherent chain of reasoning, culminating in a solitary solution, thus only needing the generation operator. CoT serves as the foundational framework for a range of advanced reasoning strategies, including prompting methodologies such as Self-Consistency and Self-Refinement, among others.

**Self-Consistency (SC)** [158] extends the CoT framework by introducing redundancy into the reasoning process. It generates multiple reasoning chains and employs a majority-voting mechanism to determine the most consistent solution, which implements a Select operator from our blueprint.

**ToT** [169] adopts an explicit reasoning structure organized in a hierarchical, tree-based format. Within this framework, each node corresponds to a distinct reasoning step, and branching facilitates exploration across multiple inferential pathways (the Generate operator). Additionally, an evaluation operator, implemented via a specialized prompt and the LLM itself, assesses branches of the tree.

**GoT** [9] introduces a more intricate reasoning structure by employing an explicit graph-based representation. In this framework, nodes represent individual reasoning steps, and the graph architecture supports non-linear, interdependent relationships between these steps. The reasoning strategy in GoT is orchestrated by an external controller, realized as a separate LLM, which guides the exploration, refinement and aggregation of the graph’s nodes.

Scheme	Reasoning			Reasoning Operator							Models		Pipeline			Remarks	
	Structure	Step	Strategy	Structure			Traversal		Update	Evaluation		PM	VM	Inf.	Tr.		DG
				Gen.	Ref.	Agg.	Pr.	Res.	Sel.	BT	Bp.						
<b>Explicit RLMs (Section 5.1)</b>																	
rStar-Math [56]	E Tree	C Thought + Code Block	E MCTS	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
PRIME [39], [171]	E Multiple Chains	F Token	E Best-of-N	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
Marco-o1 [182]	E Tree	F Token Sequence	E MCTS	☐	☐	✗	✗	✗	☐	✗	☐	☐	☐	✗	☐	☐	
Journey Learning (Tr.) [119]	E Tree	C Thought	E Tree Search	☐	☐	✗	☐	☐	✗	☐	☐	☐	☐	☐	☐	☐	
OpenR [153]	E Tree	C Thought	E Best-of-N	☐	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	☐	☐	
			E Beam Search	☐	☐	✗	☐	☐	✗	☐	☐	☐	☐	☐	☐	☐	
			E MCTS	☐	☐	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
LLaMA-Berry [177]	E Tree of Chains	C Solution	E MCTS	☐	☐	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
ReST-MCTS* [178]	E Tree	C Thought	E MCTS	☐	☐	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
AlphaMath Almost Zero [24]	E Tree	F Thought	E MCTS	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
MCTS-DPO [163]	E Tree	F Token Sequence	E MCTS	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
AlphaLLM [149]	E Tree	C Option	E MCTS	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
TS-LLM [48]	E Tree	F Token	E MCTS	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
		F Sentence	E Tree Search	☐	☐	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
<b>Implicit RLMs (Section 5.2)</b>																	
QwQ [148]	I Chain*	F Token	✗	☐	✗	✗	✗	✗	☐	✗	✗	✗	✗	✗	☐	✗	
Journey Learning (Inf.) [119]	I Chain*	C Thought	I DFS	☐	✗	✗	✗	✗	☐	✗	✗	✗	✗	✗	☐	✗	
<b>Structured Prompting Schemes (Section 5.3)</b>																	
Graph of Thoughts (GoT) [9]	E Graph*	C Thought	E Controller	☐	☐	☐	✗	✗	☐	☐	✗	☐	☐	✗	☐	✗	
Tree of Thoughts (ToT) [169]	E Tree	C Thought	E Tree Search	☐	✗	✗	☐	✗	☐	✗	☐	☐	☐	☐	☐	☐	
Self-Consistency (SC) [158]	E Multiple Chains	C Thought	E Majority Voting	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	
Chain of Thought (CoT) [160]	I Chain	C Thought	✗	☐	✗	✗	✗	✗	☐	✗	☐	☐	☐	☐	☐	☐	

TABLE 1: Comparison of RLMs with respect to the provided taxonomy (Section 4 and Figure 5). “Reasoning”: Details of the reasoning approach, specifically what is its Structure and its Strategy? “Reasoning Operator”: Does a given scheme support operators on the reasoning structure? If yes, which classes (and specific functionalities) are supported Structure (“Gen.”: generate, “Ref.”: refine, “Agg.”: aggregate, “Pr.”: prune, “Res.”: restructure), Traversal (“Sel.”: select, “BT”: backtrack), Update (“Bp.”: backpropagate), and Evaluation of “Inter.”: intermediate steps and “Final.”: final steps? “Model”: Does a given scheme use models to implement its operators and if so, which ones (“PM”: policy model, “VM”: value model)? “Pipeline”: Which pipelines are harnessed by a given scheme (“Inf.”: inference, “Tr.”: training, “DG”: data generation)? When describing representations, we use the following abbreviations: “E”: explicit, “I”: implicit. “F”: fine-grained. “C”: coarse-grained. “☐”: full support (i.e., YES), “☐”: partially [supported], “✗”: no support (i.e., NO).

## 6 HOW TO USE THE BLUEPRINT

We now outline how to use our blueprint for the user’s application; we keep this section in a tutorial style.

### 6.1 Part 1: Define the Reasoning Scheme

The first step in using the blueprint is to define the reasoning scheme, which specifies the foundational structure and strategy of your RLM. Start by selecting the reasoning structure. Chains are the most affordable in terms of token costs, at least when it comes to ICL [14]. Trees, while the most expensive, offer rich branching that enhances exploratory reasoning. Graphs, though slightly cheaper than trees, introduce additional challenges in implementation but can yield significant accuracy gains due to their flexibility.

Next, decide on the granularity of reasoning steps. Coarse-grained steps, such as thoughts or sentences, are widely used due to their simplicity and ease of scaling. However, token-based granularity, which operates at the level of individual tokens, offers the potential for greater precision and unexplored accuracy improvements. This approach, while promising, demands significantly more computational resources and careful design. This decision defines your action space (possible operations) and state space (configuration of the reasoning structure).

Another decision is choosing a reasoning strategy to govern how the reasoning structure evolves. MCTS combined with some variants of policy and value models remains the most widely adopted approach due to its balance of exploration and exploitation. However, alternative strategies that have not been deeply studied, such as ensembles of reasoning structures, may offer untapped potential.

Finally, determine the specific details of your chosen strategy, including parameters like exploration coefficients,

decoding strategy, scoring functions, and step evaluation methods. These choices will significantly impact the model’s reasoning dynamics, scalability, and overall effectiveness. Each decision at this stage lays the foundation for tailoring the RLM to your specific application requirements.

### 6.2 Part 2: Define the Operators

The next step is to specify the set of operators that will govern the reasoning process. For an MCTS-based design, the simplest approach is to implement the core operators: Generate (often called Expand for MCTS), Select, and Backpropagate. These fundamental operations suffice for many scenarios, providing a straightforward framework for reasoning.

Beyond the basics, consider whether you want to incorporate less mainstream operators, such as Backtrack. By explicitly including Backtrack, you enable a clearer tracking of progress within the search tree, making it potentially easier to revisit and refine earlier reasoning steps. This approach also facilitates advanced training schemes, like Trace-based Supervision, by generating richer and more structured data. Consider using this and other operators within our toolbox.

You will also need to determine the implementation details for each operator. Decide which operators will be implemented as neural models—such as using a policy model to guide selection or a value model for backpropagation—and which will rely on non-neural methods. This choice affects both the computational complexity and the flexibility of the system, so it’s important to align these decisions with your reasoning scheme and performance goals.

### 6.3 Part 3: Determine the Training Details

In this phase, you need to outline the specifics of training for the models that will implement operators. For an MCTS-based design, consider the typical approach of using the policy model to implement Generate (Expand) and the value model for Simulate. If necessary, you might also train a separate model to calculate the reward at individual nodes, enhancing the precision of the reward signals.

Identify the application or training domain in order to address generalization requirements. This step ensures that your models are trained on data representative of the tasks you want them to handle.

Define the models, including their architectures and the selection of suitable base models. Consider how the design of these models—such as transformer-based architectures or more specialized designs—aligns with your reasoning structure and overall objectives.

Collect training data for both the policy and value models. For the policy model, consider generating data automatically with our pipeline or using a scheme such as CoT prompting, and include a special end-of-step token to ensure clean segmentation. For the value model, generate data through MCTS full simulations, which provide rich, structured information about reasoning paths and outcomes.

Fine-tune the models as needed. If using coarse reasoning steps, perform supervised fine-tuning (SFT) on the policy model to teach it how to reason step-by-step. Similarly, apply SFT to the value model to initialize it as a reliable evaluator.

Run MCTS with initialized models to collect additional data. You might filter this data to keep only high-quality reasoning paths (terminal states) or strong signals (high absolute advantages) for further training.

Finally, train both models either by additional SFT rounds or with reinforcement learning methods such as Proximal Policy Optimization (PPO). This ensures that the models are optimized not only for accuracy but also for the efficiency and robustness needed in complex reasoning tasks.

## 7 FRAMEWORK X1: DESIGN & IMPLEMENTATION

We now introduce **x1**<sup>4</sup>, an extensible and minimalist framework that can serve as ground to design and experiment with RLMS, and currently provides one example of the blueprint.<sup>5</sup> An overview of the framework is in Figure 6.

### 7.1 Reasoning Scheme

The **x1** framework employs a tree reasoning structure in conjunction with MCTS as the reasoning strategy. This combination allows for a systematic exploration of reasoning paths while balancing exploration of new possibilities and exploitation of promising solutions judged by a value model. The framework achieves this alignment through the implementation of a series of operators that guide the construction, traversal, evaluation, and updating of the reasoning tree.

<sup>4</sup><https://github.com/spcl/x1>

<sup>5</sup>We are working continuously on expanding the framework as well as adding more RLMS.

### 7.2 Operators

The **Generate** operator plays a crucial role in expanding the tree by adding new children to a selected node. To improve the diversity of these newly generated nodes, we employ diverse beam search [152], which ensures variability among the children. Alternatively, high-temperature sampling can be used to introduce stochasticity into the generation process, fostering the exploration of different reasoning paths.

Traversal of the reasoning tree is managed by the **Select** operator, which uses the PUCT function to identify the next node to expand. This operator balances a trade-off between exploration, favoring less-visited nodes, and exploitation, reinforcing nodes with higher potential based on previous evaluations. Always starting from the root node, the traversal mechanism ensures that the system can dynamically explore alternative paths and recover from suboptimal decisions by backtracking and selecting new branches.

The **Backpropagation** Update operator refines the q-values which can be used as guidance for the select operator along the path from an expanded node back to the root. This process incorporates new information from downstream nodes, leading to progressively more accurate q-values for the intermediate nodes. These refined q-values subsequently inform future decisions, making the reasoning process increasingly robust over time.

The framework implements two different Evaluate Operators. First, the **Reasoning Path Evaluation** operator predicts the discounted expected future reward for a chain extending from the root to a specific node. This prediction is derived from the q-value model, offering a quantitative measure of the path's quality. Second, when the ground truth is available, the **Ground Truth-Based Reward** operator directly evaluates leaf nodes for correctness, assigning fixed rewards to verified solutions. These rewards are incorporated into the q-values of upstream nodes, ensuring that the reasoning process is informed by both model predictions and objective validation.

### 7.3 Models & Training Paradigms

Both the value and the policy model in **x1** are fine-tuned versions of an LLM<sup>6</sup>, without reliance on prompting, which is used in several other RLM architectures [56], [178]. This design decision aims to maximize the quality of results. We now outline briefly selected key aspects of how we train these models, full details can be found in Appendix B, C, and D.

#### 7.3.1 Training the Policy Model

The policy model also leverages an LLM to generate new nodes during the MCTS. It is fine-tuned to output an individual next reasoning step instead of a whole chain of thoughts towards a completion (which LLMs commonly do). We achieve this by introducing a *novel token*, the *end of intermediate step (eois) token*, which denotes the completion of each reasoning step. The *eois* token complements the standard end of sequence (eos) token, which indicates the conclusion of an entire reasoning chain. By incorporating

<sup>6</sup>We currently use Llama-3.1-8B-Instruct as base model.

the eois token, the framework enables the explicit identification of intermediate reasoning steps, allowing for greater interpretability and precise determination of whether the reasoning process is complete or ongoing. This dual-token strategy enhances the LLM’s capability to decompose complex problems into manageable substeps while ensuring the model recognizes when a solution has been reached.

### 7.3.2 Training the Value Model

The value model is designed to estimate the sum of the expected discounted future rewards for a sequence of reasoning steps and a newly proposed reasoning step, quantifying the value of the node modeling this step. For a given node in the MCTS tree, its value (referred to in the MCTS literature as state action value or q-value) is defined as the expected cumulative reward discounted by the number of steps required to achieve it. Formally, the q-value  $Q_\pi(s_t, a_t)$  for traversing the edge to node  $s_{t+1}$  when taking action  $a_t$  from  $s_t$  at depth  $t$  in the MCTS tree is expressed as

$$Q_\pi(s_t, a_t) = \mathbb{E} \left[ \gamma^{T-t} r(s_T, a_T) \mid s_t, a_t \right] \quad (1)$$

$$\approx \frac{1}{N} \sum_{i=1}^N \gamma^{T-t} r(s_T^{(i)}, a_T^{(i)}) \quad (2)$$

where  $\gamma$  is the discount factor,  $T$  marks the last reasoning step  $a_T$  that is added resulting the terminal state  $s_{T+1}$  containing the complete reasoning structure and rewards are modeled sparse. The terminal state  $s_{T+1}$  is defined as the state in which no additional reasoning steps can be added. It typically represents the state containing the final solution to the problem at hand. Accordingly,  $r(s_T, a_T)$  is the terminal reward. We chose to model rewards sparse, where only the final reasoning step receives a non-zero reward, since for most reasoning tasks, only the final answer can be evaluated against the true solution. As a result, one can only obtain a reward signal when the last step is reached. We can approximate the q-value by sampling  $N$  reasoning chains until the terminal state, as in 2, and averaging the terminal rewards discounted by the depth required.

The q-value model is trained using data from completed MCTS searches. Initially, when the q-value model is unavailable,  $N$  simulations (complete rollouts) are performed, and the average discounted reward is used to initialize the q-values for each node. More information can be found in the Appendix D.2.

## 7.4 Enabling Scalability and Efficiency

The current implementation is built to scale to multiple GPUs on multiple nodes. To further enhance the scalability and computational efficiency, several architectural and operational improvements have been implemented.

One design decision involves the decoupling of the value and policy models. The deployment of dedicated Value and Policy servers confers several advantages:

- **Scalability** The decoupling of Value and Policy servers from the MCTS instance facilitates scalability and the execution of multiple parallel MCTS instances.
- **Batch Processing** The policy server incorporates batching capabilities, allowing the concurrent processing of multiple queries, thereby enhancing throughput.

- **Resource Optimization** The independent allocation of computational resources to the value and policy models is inherently supported by the framework’s architecture, enhancing efficient resource utilization.
- **Replication and Distribution** The separation of value and policy models facilitates the application of distinct replication and distribution strategies.

Figure 6 illustrates the implementation of the framework as a server architecture, demonstrating how these structural enhancements contribute to improved scalability and efficiency. Building on these architectural enhancements, we employ the following strategies to further optimize the framework’s efficiency and scalability, focusing on inference and parallelization.

In the framework, we incorporate the standard optimizations of batching, quantization, and KV caching. Inference calls are batched in the policy model, enabling simultaneous processing of multiple queries. To expedite the reasoning process, the framework creates multiple child nodes in parallel during the node expansion phase. Specifically,  $N$  new nodes are generated concurrently in each expansion step, reducing computational overhead and enhancing overall system performance. Further optimization of inference speed is achieved through KV caching and quantization. KV caching mechanisms mitigate redundant computations, while quantization techniques reduce the memory consumption of both policy and value models.

## 7.5 Blueprint for Efficient Scaling

Our blueprint can be deployed to AI HPC systems and clouds, as both systems provide the performance and resources necessary to scale RLMs. Deployment on HPC systems is straightforward: compute tasks are distributed across statically allocated nodes, connected with a low-latency and high-bandwidth interconnect, and with training data being available on a high-performance parallel filesystem. On the other hand, the cloud provides many configurable services that offer different trade-offs between performance, cost, and reliability. There, it becomes the user’s responsibility to choose the storage options and compute granularity that provides the best match for expected performance and cost. The architecture of our blueprint fits into the *microservice* architecture, with a clear separation of compute tasks, data storage, and coordination. This architecture helps to ease the configuration process, as different components of the system can be deployed, scaled, and optimized independently. In particular, the separation of value and policy servers allows them to be scaled separately according to the complexity of reasoning steps that might require different resource allocations to handle task generation and evaluation.

First, we outline the major decisions users must make before deploying the x1 scaling blueprint:

- **Deployment** Training and inference tasks are typically allocated to virtual machines and containers, with the latter typically deployed as managed services with an orchestrator such as Kubernetes. There, x1 can benefit from modern frameworks like Ray [111] that hide the complexity of managing a service in a Kubernetes cluster.

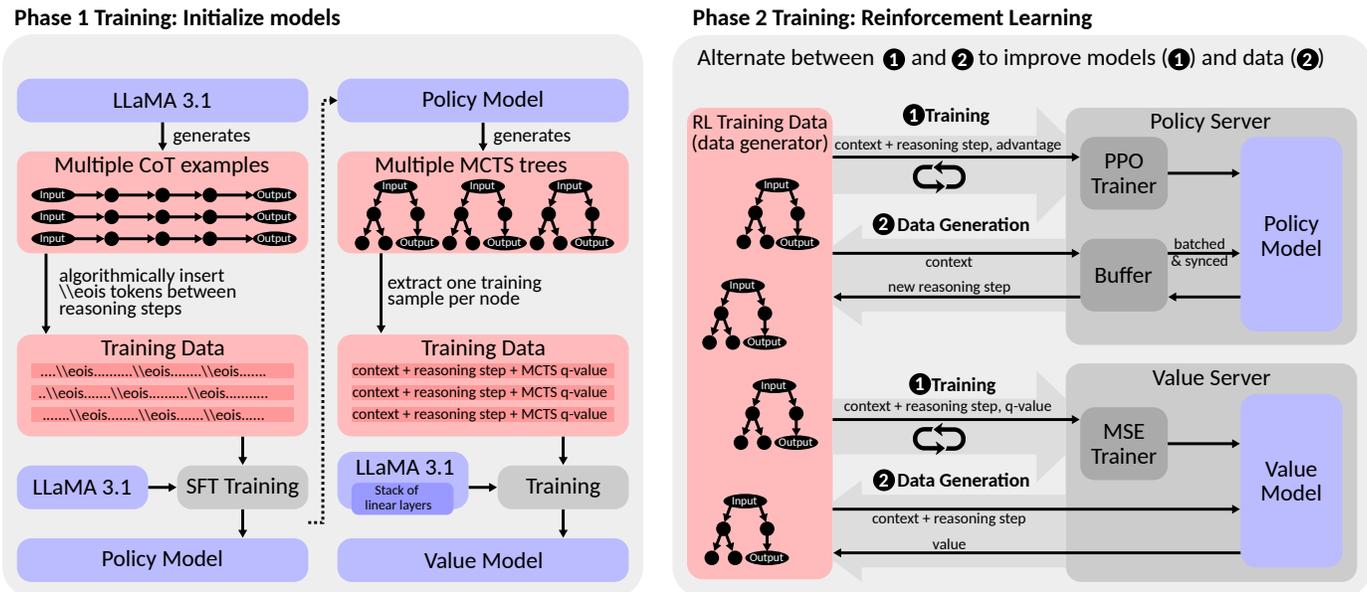


Fig. 6: An overview of the x1 framework is presented, highlighting its two-phase training process. In phase 1, the models are initialized, while in phase 2, the models are iteratively refined by alternating between constructing a sufficient number of MCTS trees and training the models on data derived from these trees.

- **Data Storage** In the cloud, object storage provides automatic bandwidth scalability that allows scale computations operating on the same data. To overcome latency and power constraints, data can also be placed in in-memory caches like Redis and hybrid solutions that combine disks with flash memory [180].
- **Communication** Requirements of the x1 blueprint differ from classical microservices, that rely on high-level abstractions like RPC and REST interfaces. RLM must utilize high-performance network fabrics offered by modern clouds, such as InfiniBand on Azure and Elastic Fabric Adapter (EFA) on AWS, both capable of achieving throughput of 400 Gb/s [40]. These are also available to training processes distributed across many GPUs, e.g., through specializations of the NVIDIA collectives library NCCL.
- **Parallelism** We apply parallelism at multiple blueprint levels, including the classic data, model, and pipeline parallelism. These can scaled horizontally across a larger number of virtual machines and containers. On the other hand, reasoning steps can benefit from elastic scaling, like in distributed MCTS and Beam Search, where each path can be explored in parallel. There, containers can be allocated on the fly to support new paths and deallocated as soon as the parallelism scale of the computation decreases.

New developments in the machine learning infrastructure can significantly impact RLM deployment strategies:

- **Elastic Compute** Computing tasks can be executed on ephemeral resources that trade the guaranteed lifetime and reliability for lower costs, such as spot virtual machines [107]. Serverless functions provide elasticity scalability with fine-grained pricing models [38], which can be a good fit for dynamically generated reasoning steps. However, serverless functions are stateless and suffer from cold starts, which requires optimization techniques dedicated to LLMs [50]. Furthermore, restricted network communication in functions forces the adoption of new

communication protocols [37], [73].

- **GPU Management** Cloud rental of GPU devices is particularly expensive, and procuring a sufficient number of devices can be challenging, specifically when constrained to a single cloud region. Given the large compute and memory requirements of base models, space-sharing might not be feasible. On the other hand, time-sharing of GPU devices between different x1 services could be a viable alternative, but it is currently constrained by large memory allocations and the cost of swapping model checkpoints between CPU and GPU memory. To increase resource utilization, new techniques for efficient GPU checkpoint and restore are needed [50].
- **Parameter-Efficient Resource Sharing** Resource-sharing can be further enhanced by utilizing a shared base model architecture for the policy and value models, while dynamically swapping task-specific parameter layers - such as Low-Rank Adaptation [66], prefix tuning [86], or other adapter layers - on the GPU during inference. These modular strategies keep the base model loaded in device memory and replace only the lightweight task-specific layers, eliminating redundant loading and reducing both latency and memory usage. An example of an RLM, which uses a shared base model with separate additional linear layers for policy and value model, is AlphaMath [24].
- **Cross-Region Deployment** Cloud applications are often deployed in a single region to avoid the performance and cost of cross-region data access. However, workloads can be scheduled globally, suspended, and migrated across regions to avoid hardware resource exhaustion and achieve lower carbon emissions [34], [161].

## 7.6 Example Analysis: Token Probability Distributions

As an illustrative example, we use the framework to directly leverage the *token probability distribution*, thereby facilitating the use of associated properties—such as entropy and variance—for guiding subsequent reasoning decisions. By fo-

ocusing on these probabilistic characteristics, the framework can help identify when to expand a given reasoning step. Using token probability distributions can be used for navigating the reasoning based on both coarse and fine steps. To support this analysis, the `x1` implementation includes scripts that provide insights into token-level metrics, such as entropy fluctuations and distribution patterns, to inform reasoning strategies.

### 7.6.1 Relevance of Token Probability Distribution

The token probability distribution provides critical information about the likelihood of different next-step candidates in a reasoning process. By examining this distribution, we can gain insight into how certain tokens dominate or diversify the reasoning space, and in turn, guide more informed decisions about which step to take next.

We now list a few scenarios where different token distributions offer insights into which reasoning decision is best to take at a given step.

- **Flat Token Distribution.** A flat probability distribution occurs when all tokens have roughly equal probabilities. In this scenario, there is significant uncertainty about which step is the best to choose because no single token stands out as a clear candidate. This can make the reasoning process more exploratory, as the model may need to consider multiple tokens equally and rely on additional strategies—such as external heuristics or learned policies—to identify the most promising step. While this can foster exploration, it may also lead to inefficiencies since the model might need to evaluate many equally plausible paths before finding an optimal solution. Another decision that could be taken in such a scenario, is to delay initiating a reasoning step till the token distribution changes to be more skewed.
- **Skewed Distribution with One Dominant Token.** When one token has a much higher probability than others, the distribution is highly skewed. This often signals that the model is confident about the next step in the reasoning process. If the dominant token corresponds to a logical or well-supported continuation, this confidence can streamline decision-making and reduce computational overhead. However, if the model’s confidence is misplaced—perhaps due to biases in the training data or a lack of context—relying on a single dominant token may cause the reasoning process to follow a suboptimal path. In such cases, it’s crucial to assess whether the high-probability token genuinely represents the most logical next step or if additional validation is needed.
- **Skewed Distribution with Multiple High-Probability Tokens.** In some cases, the distribution may be skewed with a small set of tokens receiving much higher probabilities than others. This indicates that the model sees several plausible continuations, each with a reasonable chance of being correct. While this is generally a positive sign—offering a diversity of credible options—it also complicates the decision-making process. The reasoning strategy must weigh the trade-offs between these top candidates, considering not only their individual probabilities but also how each choice impacts the subsequent reasoning trajectory. This scenario highlights the need for effective evaluation metrics (like entropy or Gini coefficient) to

help select the step that contributes most to reaching the correct or desired outcome.

By analyzing token probability distribution and identifying the cases above and others, reasoning strategies can, for example, improve efficiency (identifying when a distribution is flat allows the reasoning algorithm to focus on diversification or introduce additional constraints to narrow down choices), enhance decision confidence (recognizing when one token is dominant can help expedite decisions, provided the model’s confidence is well-founded), or foster balanced exploration (detecting multiple high-probability tokens facilitates exploring various credible paths without being overly committed to a single option).

### 7.6.2 Analyzing Token Probability Distribution

To understand the form of a token probability distribution, we examine variance, entropy, `VarEntropy`, and the Gini coefficient as key metrics that offer distinct perspectives on the distribution’s shape and characteristics.

**Variance** provides a broad measure of uncertainty by reflecting how spread out the probabilities are across the vocabulary. When variance is low, the probabilities are nearly uniform, indicating a flat distribution. However, variance alone does not capture the specific structure or shape of the distribution. For example, two distributions can have the same variance but differ in their overall form, such as one having multiple minor peaks versus another being nearly uniform with a single dominant token. To address this, we consider further measures below.

**Entropy** has long been a standard measure of uncertainty and information content in a probability distribution. Higher entropy corresponds to greater unpredictability—requiring more information to describe the system’s state. For instance, if all tokens have nearly equal probabilities, the entropy is high, reflecting a flat distribution. In contrast, low entropy occurs when a small number of tokens dominate, resulting in a skewed distribution. The entropy of a distribution is given by  $H = -\sum_i p_i \log_2(p_i)$ , where  $p_i$  is the probability of the  $i$ -th token. This metric provides valuable insight into whether the distribution is diffuse and exploratory or concentrated and decisive.

**VarEntropy** extends this analysis by measuring the variability of entropy itself, thus offering a dynamic view of how uncertainty changes. A high `VarEntropy` combined with low entropy often indicates a sharp, focused distribution with a few dominant outcomes. Conversely, low `VarEntropy` and high entropy typically reflect a flat, uniform distribution where no single token stands out. The `VarEntropy` is defined as  $\sum_i p_i (|\log(p_i)| - |H|)^2$ . This metric captures the nuanced shifts in distribution shape, helping to pinpoint how tightly probabilities cluster around certain tokens versus how broadly they spread.

The **Gini Coefficient**, traditionally used to measure inequality, provides another lens on the form of the distribution. A perfectly equal distribution has a Gini coefficient of 0, signifying that all tokens have identical probabilities. A Gini coefficient closer to 1 indicates high inequality, where a few tokens hold most of the probability mass. By visualizing the cumulative distribution of sorted probabilities, the Gini coefficient highlights how the probability is concentrated or dispersed.

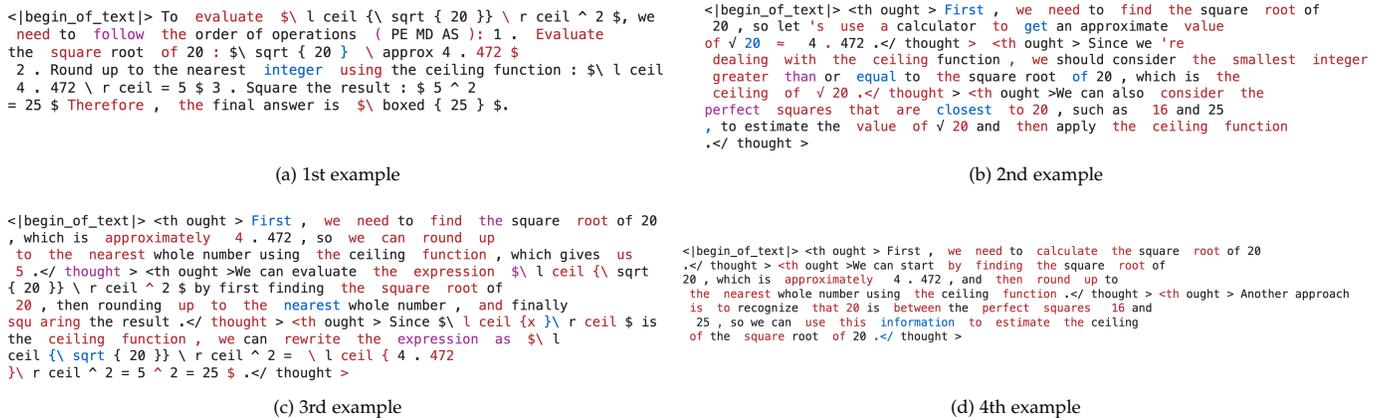


Fig. 7: Four examples of model output with highlighted tokens indicating uncertainty levels. The outputs have been color-coded to reflect the confidence levels of the model’s token predictions. Tokens are highlighted in purple when the highest probability is below 0.8 (indicating lower certainty without significant contention), in blue when the second-highest probability exceeds 0.1 (indicating contention, where another token is a close alternative), and in red when both conditions are met (indicating high uncertainty). These examples illustrate varying levels of prediction confidence and contention in reasoning steps, emphasizing regions of high ambiguity or competition between plausible continuations. This type of visual analysis is useful for identifying points in the reasoning process where the model lacks confidence or is torn between alternatives, guiding refinements in reasoning strategies and model design. It also helps pinpoint critical areas where additional supervision or context may improve model performance.

Together, these metrics—variance, entropy, VarEntropy, and Gini—enable a detailed examination of token probability distributions. By leveraging each metric’s unique strengths, we can effectively characterize whether a distribution is flat, skewed with a dominant token, or skewed across several highly probable tokens, ultimately guiding more informed decisions in reasoning and model development.

### 7.6.3 Example Results

Figure 7 and 8 illustrate example model outputs and their respective token probability distributions. By analyzing the highest probabilities, the second-highest probabilities, and the sum of the remaining probabilities, we gain valuable insights into the underlying token distribution, which can subsequently be quantified through the uncertainty metrics discussed earlier.

In Figures 8a and 8d, specific regions emerge where the top two probabilities are very close, while the remaining probabilities are significantly smaller. Such regions likely indicate scenarios where forking the reasoning process (e.g., exploring multiple paths) could disproportionately benefit future outcomes, as the competing high-probability tokens suggest alternative plausible continuations. Conversely, in instances where the first probability is notably high, with much lower second and remaining probabilities, the model exhibits strong confidence in a single continuation. These cases are conducive to more deterministic reasoning, as forking may be unnecessary.

Additionally, regions with a relatively high sum of the remaining probabilities (close to the top two) highlight flatter distributions with high uncertainty. These scenarios signal a need for cautious reasoning, where clarification or additional contextual refinement may help reduce ambiguity. For instance, such uncertainty may suggest that the model has not yet committed to a specific path and could benefit from revisiting earlier reasoning steps to address potential errors or misalignments.

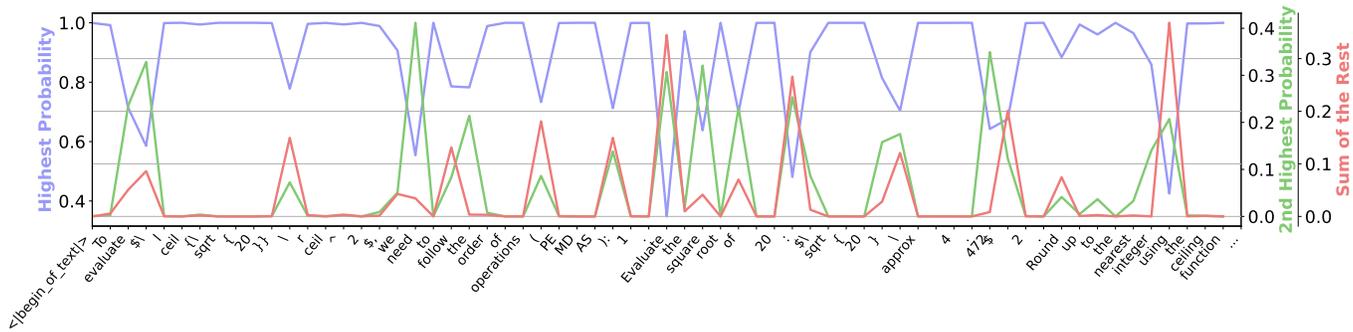
Figure 9 further analyzes these results using metrics such as variance, entropy, VarEntropy, and the Gini coefficient. In

Figure 9a, a zero-shot prompt demonstrates lower uncertainty overall, suggesting that it yields more confident predictions and potentially higher-quality outputs. However, the presence of specific high-probability tokens (e.g., “472”) raises concerns about potential data leakage into the training set or the tokenizer, which could bias the results. Another notable observation is the high uncertainty associated with <thought> tokens, which appear challenging for the model to predict accurately. This highlights the complexity introduced by token granularity, where most words correspond to single tokens, resulting in a roughly even distribution for the next token across the vocabulary in some contexts.

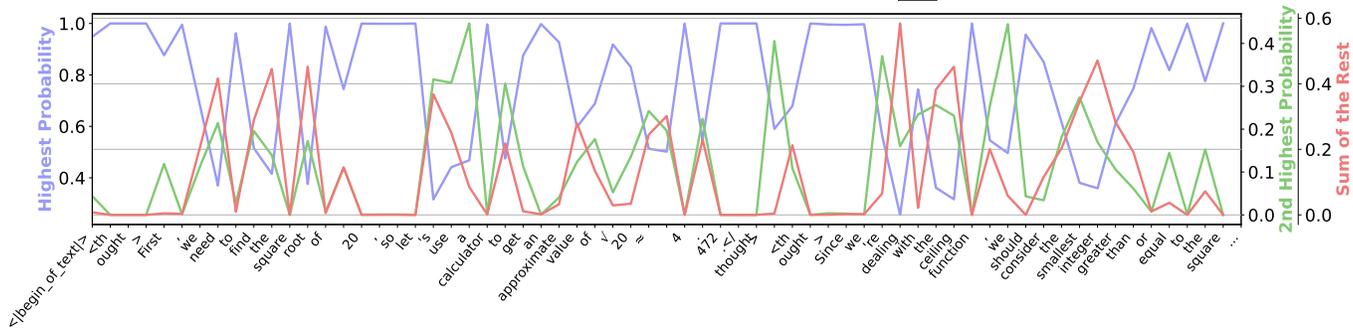
The uncertainty metrics provide actionable insights for reasoning strategy design. For example, cases with high VarEntropy and low entropy indicate a distribution where a few outcomes dominate, making tree-based search strategies effective. These strategies prioritize exploring high-probability outcomes while avoiding unnecessary evaluations of less probable branches. In contrast, low VarEntropy and high entropy reflect a flat distribution where no clear outcome dominates. Such cases could benefit from clarification mechanisms or intermediate step refinements to reduce ambiguity before proceeding further.

Interestingly, the Gini coefficient often highlights critical regions more effectively than other metrics. In vital reasoning areas, it captures the inequality in token probabilities, helping to identify tokens that significantly influence the reasoning process. This contrasts with metrics like entropy and VarEntropy, which may also flag tokens related to formatting or stylistic choices, providing less task-specific utility.

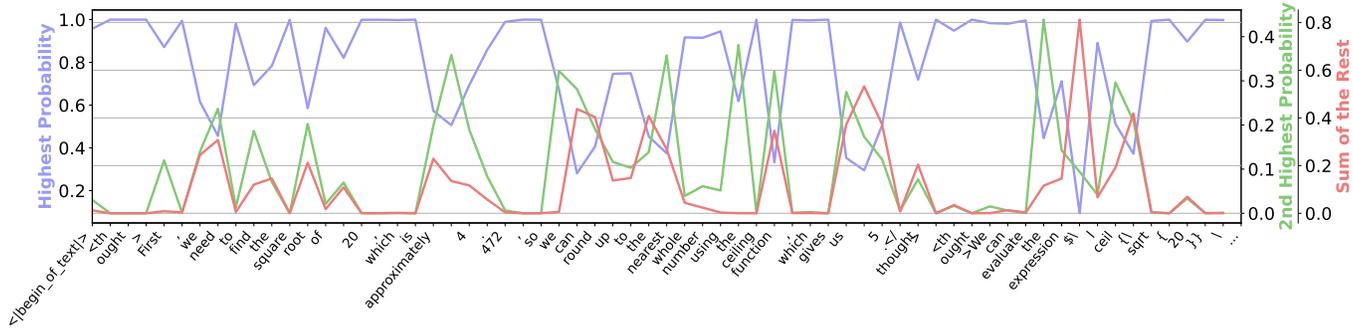
Overall, these visualizations and metrics emphasize the importance of analyzing token probability distributions to design effective reasoning strategies. By leveraging the nuanced patterns revealed by these metrics, models can better adapt to uncertainty, balance exploration and exploitation, and optimize decision-making during the reasoning process.



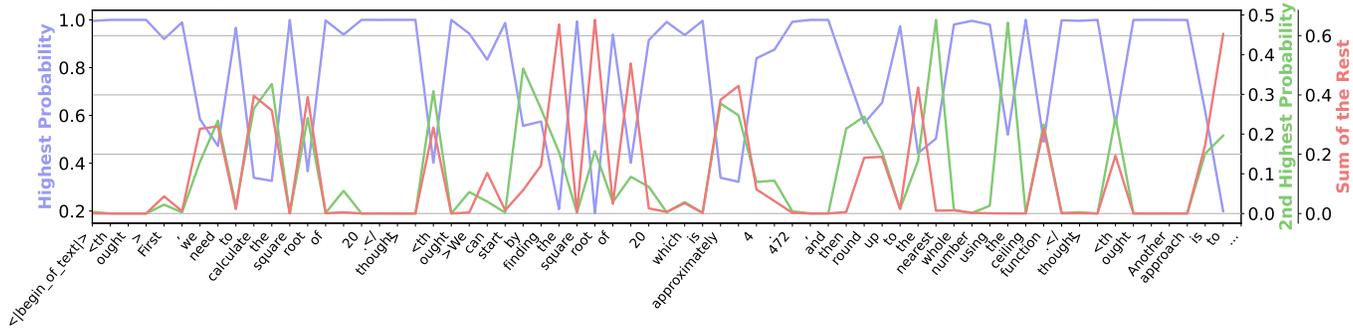
(a) To evaluate  $\lceil \sqrt{20} \rceil^2$ , we need to follow the order of operations (PEMDAS):  
 1. Evaluate the square root of 20:  $\sqrt{20} \approx 4.472$   
 2. Round up to the nearest integer using the ceiling function:  $\lceil 4.472 \rceil = 5$   
 3. Square the result:  $5^2 = 25$  Therefore, the final answer is  $\boxed{25}$ .



(b)  $\langle \text{thought} \rangle$  First, we need to find the square root of 20, so let’s use a calculator to get an approximate value of  $\sqrt{20} \approx 4.472$ .  $\langle \text{thought} \rangle$   
 $\langle \text{thought} \rangle$  Since we’re dealing with the ceiling function, we should consider the smallest integer greater than or equal to the square root of 20, which is the ceiling of  $\sqrt{20}$ .  $\langle \text{thought} \rangle$   
 $\langle \text{thought} \rangle$  We can also consider the perfect squares that are closest to 20, such as 16 and 25, to estimate the value of  $\sqrt{20}$  and then apply the ceiling function.  $\langle \text{thought} \rangle$

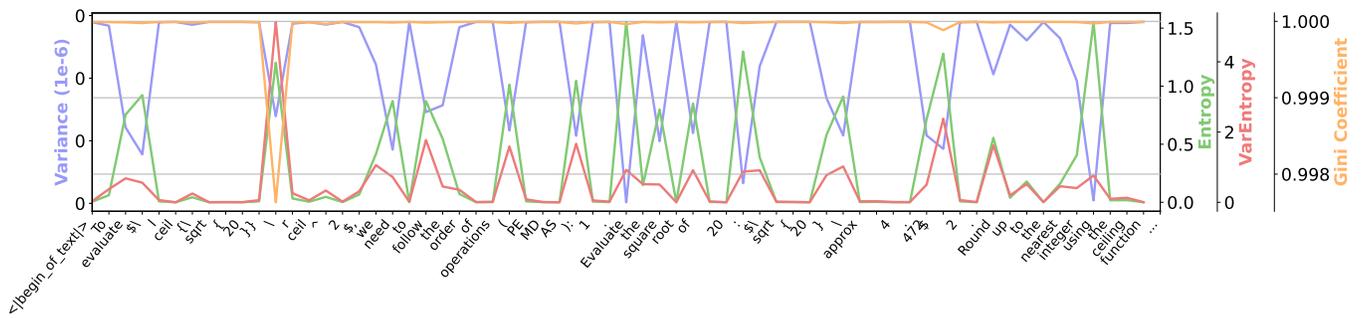


(c)  $\langle \text{thought} \rangle$  First, we need to find the square root of 20, which is approximately 4.472, so we can round up to the nearest whole number using the ceiling function, which gives us 5.  $\langle \text{thought} \rangle$   
 $\langle \text{thought} \rangle$  We can evaluate the expression  $\lceil \sqrt{20} \rceil^2$  by first finding the square root of 20, then rounding up to the nearest whole number, and finally squaring the result.  $\langle \text{thought} \rangle$   
 $\langle \text{thought} \rangle$  Since  $\lceil x \rceil$  is the ceiling function, we can rewrite the expression as  $\lceil \sqrt{20} \rceil^2 = \lceil 4.472 \rceil^2 = 5^2 = 25$ .  $\langle \text{thought} \rangle$

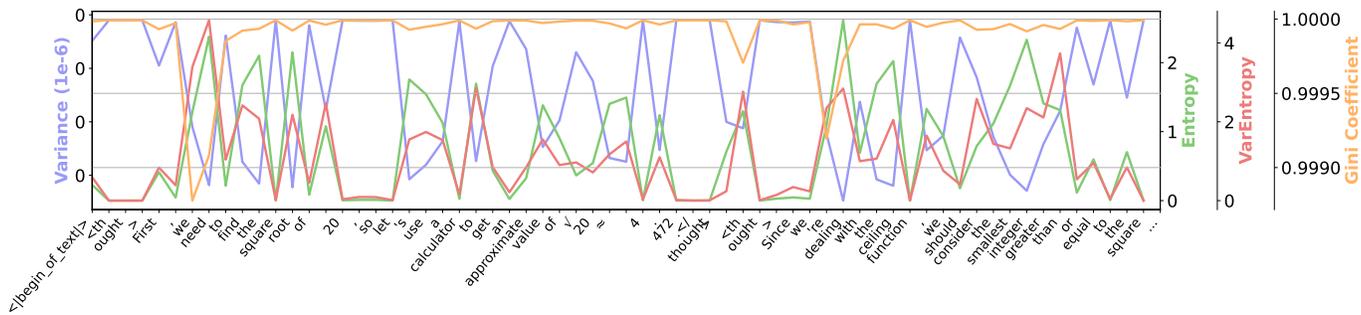


(d)  $\langle \text{thought} \rangle$  First, we need to calculate the square root of 20.  $\langle \text{thought} \rangle$   
 $\langle \text{thought} \rangle$  We can start by finding the square root of 20, which is approximately 4.472, and then round up to the nearest whole number using the ceiling function.  $\langle \text{thought} \rangle$   
 $\langle \text{thought} \rangle$  Another approach is to recognize that 20 is between the perfect squares 16 and 25, so we can use this information to estimate the ceiling of the square root of 20.  $\langle \text{thought} \rangle$

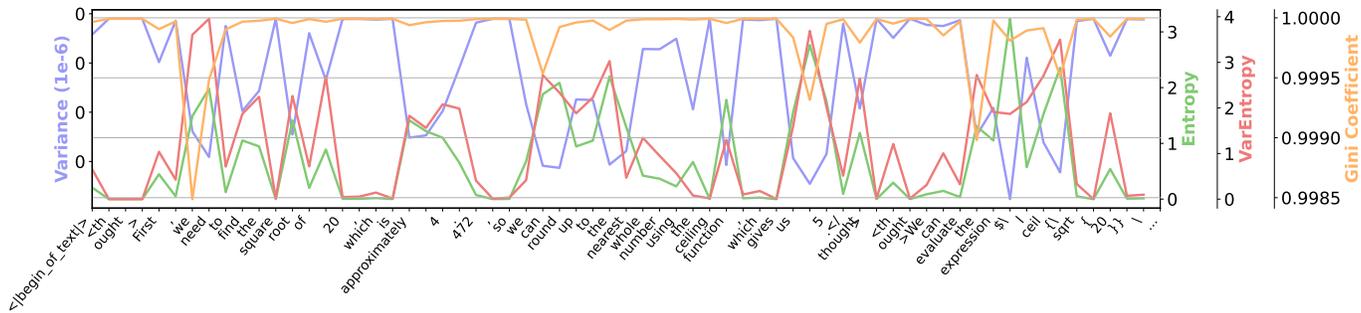
Fig. 8: Probabilities of the first 64 tokens of example model outputs. We show the two highest probabilities as well as the sum of the other probabilities.



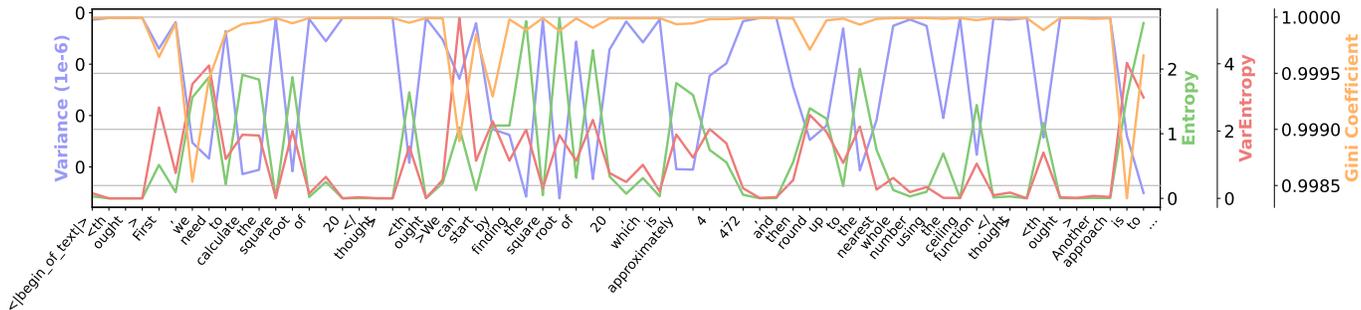
- (a) To evaluate  $\lceil \sqrt{20} \rceil^2$ , we need to follow the order of operations (PEMDAS):
1. Evaluate the square root of 20:  $\sqrt{20} \approx 4.472$
  2. Round up to the nearest integer using the ceiling function:  $\lceil 4.472 \rceil = 5$
  3. Square the result:  $5^2 = 25$  Therefore, the final answer is  $\boxed{25}$ .



- (b) <thought>First, we need to find the square root of 20, so let's use a calculator to get an approximate value of  $\sqrt{20} \approx 4.472$ .</thought>  
 <thought>Since we're dealing with the ceiling function, we should consider the smallest integer greater than or equal to the square root of 20, which is the ceiling of  $\sqrt{20}$ .</thought>  
 <thought>We can also consider the perfect squares that are closest to 20, such as 16 and 25, to estimate the value of  $\sqrt{20}$  and then apply the ceiling function.</thought>



- (c) <thought>First, we need to find the square root of 20, which is approximately 4.472, so we can round up to the nearest whole number using the ceiling function, which gives us 5.</thought>  
 <thought>We can evaluate the expression  $\lceil \sqrt{20} \rceil^2$  by first finding the square root of 20, then rounding up to the nearest whole number, and finally squaring the result.</thought>  
 <thought>Since  $\lceil x \rceil$  is the ceiling function, we can rewrite the expression as  $\lceil \sqrt{20} \rceil^2 = \lceil 4.472 \rceil^2 = 5^2 = 25$ .</thought>



- (d) <thought>First, we need to calculate the square root of 20.</thought>  
 <thought>We can start by finding the square root of 20, which is approximately 4.472, and then round up to the nearest whole number using the ceiling function.</thought>  
 <thought>Another approach is to recognize that 20 is between the perfect squares 16 and 25, so we can use this information to estimate the ceiling of the square root of 20.</thought>

Fig. 9: Uncertainty metrics (variance, entropy, VarEntropy, and the Gini coefficient) plotted against the first 64 tokens of the output token sequence.

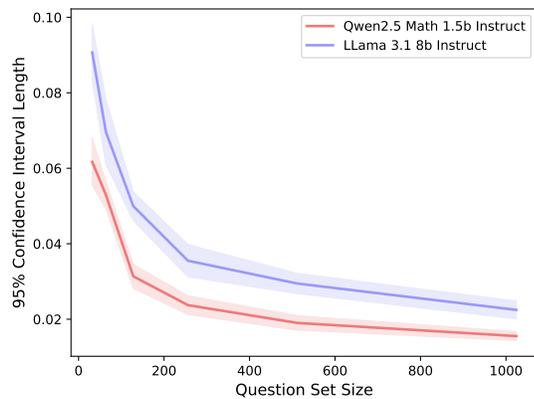


Fig. 10: Estimated 95%-confidence interval length for different question set sizes using sampled generated answers from a subset of 1000 questions with eight generated answers per question at temperature 1. The confidence interval is calculated over the eight different pass@1 subsets of each question with 32 sets randomly sampled with replacement for each set size.

## 7.7 Benchmarking RLMs

Our experience with benchmarking RLMs highlights critical considerations for ensuring fair and reliable performance comparisons. Incorporating multiple models within a reasoning scheme often increases output variance, emphasizing the need for benchmarking on sufficiently large sample sizes. Benchmarks with limited sample sizes, such as AIME or AMC, which often provide only a two-digit range of samples, risk selective reporting. This occurs when researchers focus on subsets of results where their models perform well, rather than reflecting the true variability of their systems.

Experimental findings (Figure 10) demonstrate that achieving low error variability, within a single-digit percentage range, requires evaluation across at least 500 samples. Given the inherent complexity of RLMs, which often exhibit greater variability than simpler LLM setups, these results suggest specific sample size thresholds. We recommend that individual benchmarks contain at least 200 samples per category, with a minimum of 500 samples evaluated across all categories to ensure statistically robust comparisons. Adhering to these guidelines would in many cases mitigate variability-driven biases and facilitate more transparent assessments of RLM performance across different approaches.

## 8 EXAMPLE INSIGHTS FOR EFFECTIVE RLMs

We provide example insights gathered from the literature and from our analyses of design decisions using **x1**.

**Use Process-Based Evaluation** Process-based evaluation, in which the reasoning structure as a whole is assessed, has been shown to be more reliable than alternative methods such as Outcome-Based Reward Models (ORMs). By examining the reasoning steps and their relationships within the structure, process-based evaluation provides a richer signal that helps models refine their reasoning paths and improve overall accuracy. This approach ensures that each intermediate step contributes positively to the final outcome, resulting in more robust reasoning and better generalization across tasks.

**Use Two Phases for Training** Adopting a two-phase training strategy—splitting SFT and RL—has proven effective in several contexts. This phased approach allows the model to first learn a solid foundation of reasoning patterns in phase one, followed by fine-tuning under more complex, adaptive conditions in phase two. For instance, research on Process Reinforcement through Implicit Rewards demonstrates that models trained with a dedicated SFT phase can maintain performance on standard benchmarks while achieving improved reasoning capabilities during RL. This separation also helps mitigate instability and ensures that each phase targets specific learning objectives, ultimately leading to more robust RLMs.

**Train on Familiar Distributions** Training on familiar data distributions can significantly influence a model’s initial performance and subsequent improvements. For example, PRIME [39], [171] shows that training on a carefully curated token sequence (such as the eois token approach) avoids performance degradation. Similarly, in tasks like rStar-Math [56], models trained on well-defined, familiar distributions tend to stabilize more quickly and produce higher-quality reasoning outputs. By focusing on familiar distributions, researchers can ensure that the models effectively internalize the fundamental reasoning patterns before moving on to more diverse or challenging tasks.

**Be Careful with Prompting LLMs to Critique and Evaluate** Relying on prompting alone to encourage large language models to critique and evaluate their own outputs often leads to instability. Research indicates that models struggle to self-correct reliably when prompted to refine their reasoning without external guidance. For example, a recent study [68] illustrates that such prompting typically fails to produce consistently improved results. Another work [120] demonstrates that explicitly training the model to output better responses through iterative refinement outperforms simple prompting. These findings highlight the importance of structured training approaches and careful operator design when aiming for self-improvement capabilities in RLMs.

## 9 BENCHMARKS FOR RLMs

We now outline benchmarks related to RLMs. Sun et al. [141] provide a clear distinction between various types of reasoning including mathematical, logical, casual, and common-sense. Below, we highlight a selection of benchmarks for each category. We also include additional categories related to the realm of RLMs, namely, coding related benchmarks and benchmarks that involve reasoning utilities such as tools or RAG. We show the benchmarks in Figure 11.

### 9.1 Mathematical Reasoning

Mathematical reasoning benchmarks involve arithmetic, geometry, and other mathematical tasks that use logical constructs and symbolic computation. They can be further categorized into benchmarks with fixed datasets and template-based benchmarks [109], [139].

**GSM8K** [36] consists of a train set (7,473 samples) and a test set (1,319 samples) of high-quality grade school-level mathematical word problems. Early breakthroughs in

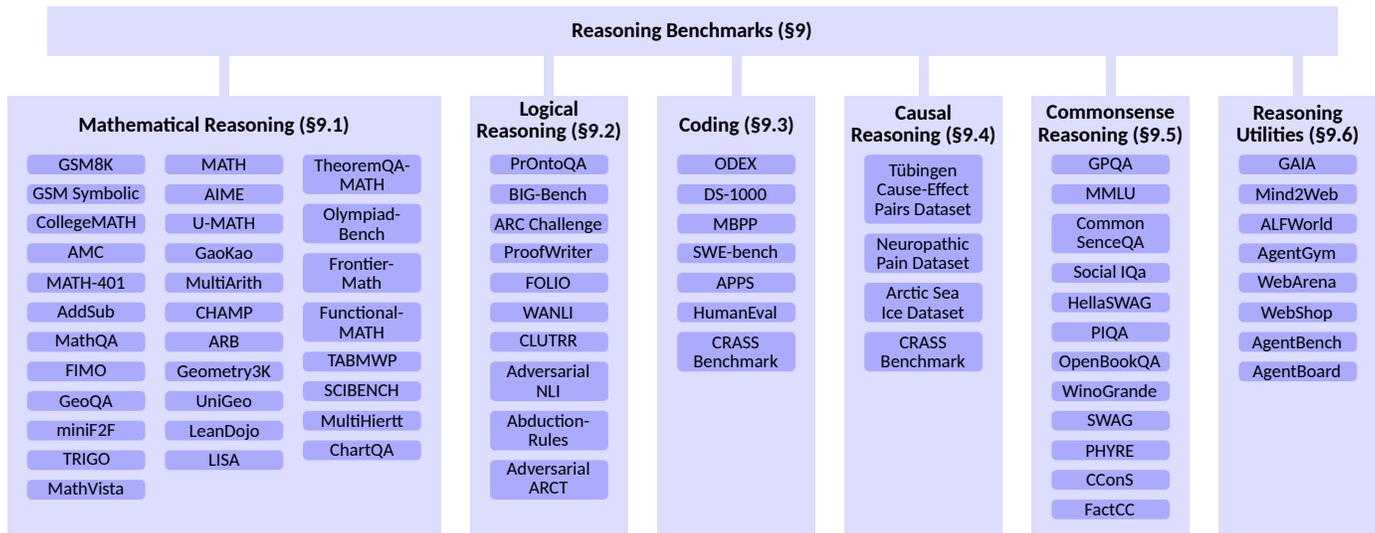


Fig. 11: Overview of benchmarks for RLMs.

mathematical problem-solving by language models were achieved by training on the training subset of this benchmark.

**GSM Symbolic** [109] introduces a generator that can use 100 templated questions, which are derived from the questions of the GSM8K dataset. This approach emphasizes the limited generalization capabilities of current RLMs and highlights the importance of templated benchmarks in evaluating LLMs’ performance in mathematical reasoning.

**MATH** [63] benchmark contains questions ranging in difficulty from high school to competition-level mathematics, containing 12,500 problems, split into 7,500 for training and 5,000 for testing. These problems are sourced from various mathematics competitions such as the AMC 10, AMC 12, and AIME (Level 5).

**Functional MATH** [139] builds upon the MATH dataset by introducing templated problem formats designed to assess the functional understanding of mathematical concepts by LLMs. However, the code and templates remain inaccessible to the public, limiting its broader adoption.

**AIME** [4], **AMC** [3], and **GaoKao** [87] feature mathematical tasks ranging from Olympiad level to college entrance level difficulty. The AMC is generally easier, the GaoKao offers a broader range of difficulty levels, while the AIME is likely the most challenging. AIME consists of 30 problems, the AMC includes 40 problems and the GaoKao contains around 300 questions.

**OlympiadBench** [60] is a more advanced benchmark that spans Olympiad-level mathematics and physics problems, comprising 8,476 problems sourced from international and Chinese Olympiad competitions, as well as the Chinese College Entrance Exam (GaoKao).

**CollegeMATH** [147] is designed for evaluating college-level mathematics, with a dataset that contains 1,281 training problems and 2,818 test problems. These problems are sourced from textbooks, extracted with the help of LLMs.

**U-MATH** [32] benchmark features 880 university-level test problems without images sourced from ongoing courses across various institutions, currently available through the Gradarius platform. This benchmark presents unpublished,

open-ended problems balanced across six core subjects.

**FrontierMath** [53] is an expert-level benchmark containing exceptionally challenging mathematics problems covering a wide array of modern mathematical domains. The dataset size remains undisclosed, but the problems have been carefully crafted and tested by expert mathematicians. Notably, current state-of-the-art models can solve less than 2% of the problems, revealing a still significant gap between AI capabilities and human expertise in the field of mathematics.

In general, it is recommended to utilize templated versions of these benchmarks where available, rather than relying solely on question-answer (QA) pairs. Templated benchmarks minimize the likelihood of contamination from prior exposure during model training, thus providing a more accurate measure of performance [109], [139].

Other related benchmarks include MATH-401 [172], MultiArith [124], AddSub [65], CHAMP [103], MathQA [5], ARB [129], FIMO [90], Geometry3K [93], GeoQA [27], UniGeo [25], miniF2F [183], LeanDojo [167], TheoremQA-MATH [30], TRIGO [165], LISA [72], MathVista [92], ChartQA [104], TABMWP [94], MultiHiertt [181], and SCIBENCH [156].

## 9.2 Logical Reasoning

Logical reasoning emphasizes formal processes, from propositional and predicate logic to automated theorem proving.

**PrOntoQA** [127] generates ontology graphs, similar to causality graphs, which do not necessarily reflect natural patterns. From these graphs, it constructs statements and poses questions that necessitate logical reasoning for resolution. Due to the abstract and artificial nature of some ontology graphs, models must focus more on step-by-step logical reasoning rather than relying on commonsense inference to derive correct conclusions.

**BIG-Bench** [138] is one of the most extensive benchmarks for reasoning tasks encompassing over 200 tasks, each potentially comprising numerous questions. It encompasses a broad range of domains and employs templated

question formats, enabling a systematic evaluation of reasoning capabilities across diverse contexts.

**ARC Challenge** [33] assesses the ability to understand formal patterns, rules, and transformations within structured, grid-based environments. Tasks focus on identifying logical structures such as conditional relationships and sequences. For instance, deducing transformations between grids based on abstract rules exemplifies the application of formal logical reasoning paradigms.

Other benchmarks include ProofWriter [145], FOLIO [58], WANLI [89], CLUTRR [136], Adversarial NLI [112], AbductionRules [170], and Adversarial ARCT [113].

### 9.3 Coding

There also exist benchmarks related to how well a given model can code. These include ODEX [159], SWE-bench [74], DS-1000 [81], APPS [61], MBPP [6], and HumanEval [28].

### 9.4 Causal Reasoning

Causal reasoning involves understanding and analyzing cause-effect relationships, including counterfactual reasoning and causal inference. This domain challenges models to predict or reason about events based on causal dynamics.

**Tübingen Cause-Effect Pairs Dataset** [110] comprises 108 cause-effect pairs drawn from diverse domains such as meteorology, biology, medicine, engineering, and economics. It serves as a comprehensive benchmark for assessing causal reasoning across various contexts.

**Neuropathic Pain Dataset** [150] captures complex relationships between nerve function and symptoms in patients. It requires a domain-specific knowledge and causal inference to accurately interpret the data.

**Arctic Sea Ice Dataset** [70] consists of a 12-variable graph that models the dynamics of Arctic sea ice based on satellite data generated since 1979. It provides a structured environment to explore causal relationships within climatological systems.

**CRASS Benchmark** [49] focuses on counterfactual reasoning tasks using 274 sample multiple choice questions. It evaluates models' abilities to answer counterfactual questions, using top-k accuracy as the primary performance metric.

Many of these benchmarks have either been largely solved by current state-of-the-art models, or their applicability in real-world language model tasks remains limited, rendering them unsuitable for benchmarking current RLMs.

### 9.5 Commonsense Reasoning

Commonsense reasoning encompasses tasks that require the application of everyday knowledge, including questions that rely on implicit cultural, social, or contextual understanding. This category also extends to specialized domain knowledge tasks.

**GPQA (Diamond)** [122] is a multiple-choice benchmark spanning disciplines such as chemistry, genetics, biology, and physics. The questions are designed to be solvable by experts (PhDs) within their respective fields but remain

challenging for experts from unrelated domains. The diamond subset contains 198 samples.

**MMLU (STEM)** [62] incorporates questions across a spectrum of difficulty, ranging from general commonsense reasoning to highly specialized domain knowledge.

Other related benchmarks include Social IQa [126], SWAG [173], HellaSWAG [174], CommonSenseQA [146], PIQA [19], PHYRE [7], OpenBookQA [108], CConS [78], WinoGrande [125], and FactCC [79].

### 9.6 Reasoning Utilities

Benchmarking capabilities of RLMs related to reasoning utilities involve testing the capabilities of an RLM in how it acts as an agent. This includes benchmarks such as GAIA [106], WebArena [185], Mind2Web [42], WebShop [168], ALFWorld [132], AgentBench [91], AgentGym [162], and AgentBoard [22]. Another line of related benchmarks tests the RAG capabilities [26], [47], [98], [164].

## 10 RELATED ANALYSES

RLMs have been explored from several angles in prior works, yet significant gaps remain in providing a systematic blueprint and open-sourced framework for their construction. Below, we categorize prior efforts and describe how our work advances the field.

### 10.1 Reasoning with Standard LLMs

Several works explore techniques for enhancing the reasoning capabilities of standard LLMs. These approaches use straightforward mechanisms applied during pre-training, fine-tuning or inference.

**Enhancing Reasoning with Training** Huang and Chang [67] outline pre-training and fine-tuning on reasoning datasets, and advanced prompting strategies. Sun et al. [141] contribute additional insights, including techniques such as alignment training and the integration of Mixture of Experts architectures. Furthermore, Huang et al. [69] demonstrate the possibility of self-improvement on reasoning tasks with additional training on self-generated labels.

**Reasoning with Prompting & In-Context Learning** Qiao et al. [118] provide an overview of prompting-only techniques, classifying prompting methods into two main categories: strategy-enhanced reasoning and knowledge-enhanced reasoning. Besta et al. [14] provide a taxonomy of different advanced in-context reasoning topologies. These include the Chain-of-Thought (CoT) [160], Tree of Thoughts (ToT) [169], and Graph of Thoughts (GoT) [9].

Some of these works further provide overviews of different reasoning tasks, reasoning datasets, and reasoning benchmarks [67], [118], [141]. Others focus on enhancing domain-specific reasoning, such as mathematical [2], [95], [166] or logical reasoning [97].

These studies remain largely limited to reviewing existing literature. Therefore, they lack code implementation and rarely employ formal language. Most importantly, they rarely cover explicit reasoning models. Our blueprint integrates most of these techniques within a broader, modular structure.

## 10.2 Explicit Reasoning Models

The following works explore techniques that extend beyond basic mechanisms applied during pre-training or inference. These methods involve additional computation to iteratively refine reasoning paths, often increasing computational demands during training and/or inference.

Dong et al. [44] provide a taxonomy and survey of inference-time self-improvement methods, including independent, context-aware, and model-aided approaches. Guan et al. [55] propose verifier engineering, a post-training paradigm for foundation models involving three stages: Search, Verify, and Feedback, to enhance model outputs with scalable supervision signals. Zeng et al. [175] provide a comprehensive roadmap for reproducing OpenAI’s o1 reasoning model from a reinforcement learning perspective. Although the work thoroughly examines all core components: policy initialization, reward design, search, and learning, no implementation is provided. Various specific implementations of RLMs exist, we provide a summary in Table 1. There are also other works related to Explicit RLMs, considering both coarse reasoning steps [157], [163] and fine reasoning steps [41], [157], [163].

Our blueprint provides a more foundational and universally applicable framework for RLMs. We further supplement the theoretical and algorithmic overview with a modular and scalable implementation to enable practical development and experimentation.

## 11 CONCLUSION

This work introduces a comprehensive blueprint for reasoning language models (RLMs), providing a flexible and modular toolbox that demystifies the intricate design and operation of these advanced systems. By encompassing diverse reasoning structures, operations, and training schemes, the blueprint establishes a robust foundation for constructing, analyzing, and extending RLMs tailored to various applications. The accompanying **x1** implementation enhances this contribution, offering a modular, minimalist, and user-friendly platform for experimentation and rapid prototyping of novel RLM architectures.

Our blueprint and **x1** pave the way for several exciting avenues of future research and development in reasoning AI. One example is Trace-Based Supervision (TBS), which extends process-based supervision by incorporating labeled traces of traversal through reasoning structures. TBS has the potential to train more powerful implicit RLMs capable of internalizing reasoning structures and improving generalization.

The work also explores new directions in value and reward modeling, introducing a hierarchy of models and formally identifying several recent designs as instances of a new class of models, namely the Outcome-Driven Process Reward Model. This model class bridges the gap between outcome-based evaluation and process-based supervision by dynamically connecting intermediate reasoning steps to terminal outcomes, enabling more granular feedback during training without the need.

Additionally, the blueprint’s extensive set of operators can inspire the development of innovative reasoning strategies, such as advanced tree-based searches, multi-step re-

finement processes, or hybrid search algorithms that adapt dynamically to the task’s complexity. These strategies can be tailored using the token probability distribution analysis tools provided, leading to more effective generation strategies that optimize reasoning steps through probabilistic insights. The blueprint also provides a foundation for developing nested architectures where reasoning structures such as trees and graphs are embedded hierarchically. These designs can address multi-layered reasoning tasks, expanding the scope of RLM applications to domains requiring deep, structured reasoning processes.

Scalability remains a key focus of this work. The blueprint’s modular design supports future scalable cloud deployments that enable efficient distribution of compute-intensive tasks across cloud infrastructures. These deployments will not only enhance scalability but also optimize cost and resource utilization, making RLMs more accessible for real-world applications.

By exploring and integrating these ideas, this work aims to empower the next generation of reasoning language models, democratize access to advanced reasoning capabilities, and foster innovation across research and industry. The blueprint’s versatility, combined with the **x1** platform, will make it one of the factors in the progress in RLM research and applications.

## ACKNOWLEDGEMENTS

We thank Nicolas Dickenmann for writing the initial MCTS codebase. We thank Hussein Harake, Colin McMurtrie, Mark Klein, Angelo Mangili, and the whole CSCS team granting access to the Ault, Piz Daint and Alps machines, and for their excellent technical support. We thank Timo Schneider for help with infrastructure at SPCL. This project received funding from the European Research Council (Project PSAP, No. 101002047), and the European High-Performance Computing Joint Undertaking (JU) under grant agreement No. 955513 (MAELSTROM). This project received funding from the European Union’s HE research and innovation programme under the grant agreement No. 101070141 (Project GLACIATION). We gratefully acknowledge Polish high-performance computing infrastructure PLGrid (HPC Center: ACK Cyfronet AGH) for providing computer facilities and support within computational grant no. PLG/2024/017103.

## APPENDIX A MATHEMATICAL FOUNDATION OF MARKOV DECISION PROCESSES FOR REASONING TASKS

y so far to put this on the first page of the appendix

In this section, we provide a rigorous mathematical framework for RLMs. We achieve this by integrating the theory of Markov Decision Processes (MDPs) with the Monte Carlo Tree Search (MCTS) algorithm. The MDP serves as a foundational formulation for modeling various types of processes, and it can be applied to model reasoning *chains*, which constitute the reasoning structure of the RLMs. Simultaneously, MCTS serves as an efficient search algorithm for exploring and navigating the extensive space of possible reasoning chains. The resulting state space is then used as a basis for modeling the RLM. An overview of the notation used in this section is provided in Table 2.

### A.1 Markov Decision Process

**Markov Decision Process (MDP)** is defined as a 5-tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r, \gamma)$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the action space with  $\mathcal{A}_s \subseteq \mathcal{A}$  denoting the set of actions which can be taken in the state  $s$ ,  $p$  represents the dynamics of transitions between states, i.e.,  $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  where  $p(s, a, s')$  is the probability of transitioning to the state  $s'$  when action  $a$  was selected in the state  $s$ ,  $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the reward function, i.e.,  $r(s, a, s')$  represents the reward for arriving in the state  $s'$  after selecting the action  $a$  in the state  $s$ , and  $\gamma \in [0, 1]$  is a discount factor.

#### A.1.1 Solving an MDP

Before stating what it means formally to *solve an MDP*, we first need several definitions.

A **trajectory**  $\tau_\pi = (s_0, a_0, \dots, s_T, a_T, s_{T+1})$  is a sequence of interleaved states and actions, selected according to the policy  $\pi$  (see below for the policy definition). Each trajectory starts at an initial state  $s_0 \in \mathcal{S}$  and ends with  $s_{T+1} \in \mathcal{S}$  which represents the terminal state where no further actions can be taken.

A **policy**  $\pi(s)$  is a function assigning a probability distribution over the action space to a given state  $s$ ;  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  where  $\Delta(\mathcal{A})$  is a set of probability distributions over action space  $\mathcal{A}$ . The expression  $\pi(a | s)$  denotes the probability of selecting the action  $a$  in the state  $s$  according to the policy  $\pi$ .

**State value function**  $V_\pi(s_t)$  represents the expected *cumulative* future reward for a given state  $s_t$  under policy  $\pi$ :

$$V_\pi(s_t) = \mathbb{E} \left[ \sum_{k=t}^T \gamma^{k-t} r(s_k, a_k, s_{k+1}) \mid s_t \right] \quad (3)$$

where  $T$  is a predefined time-horizon. Note that, in order to obtain the state  $s_{k+1}$ , an action  $a_k$  is first derived by sampling from a distribution  $\pi(s_k)$ . Once the action  $a_k$  is chosen, the environment dynamics  $p(s_{k+1} | s_k, a_k)$  determine the probability distribution of the next state  $s_{k+1}$ .

The goal of **solving an MDP** is to find a policy  $\pi^*$  which maximizes the value function as defined above for all states  $s \in \mathcal{S}$ ,  $\pi^* = \arg \max_{\pi} V_\pi(s)$

**The State-Action value function**  $Q(s_t, a_t)$  Oftentimes, it is useful to use the state-action value function  $Q(s_t, a_t)$  instead of the state value function. Specifically, the state-action value function  $Q(s_t, a_t)$  extends the state value function so that the function value is defined on a state *and* a specific action  $a_t$ :

$$\begin{aligned} Q_\pi(s_t, a_t) &= \mathbb{E}_\pi \left[ \sum_{k=t}^T \gamma^{k-t} r(s_k, a_k, s_{k+1}) \mid s_t, a_t \right] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}} [V_\pi(s_{t+1}) \mid s_t, a_t], \end{aligned}$$

where Bellman's equation is used in the second equality.

#### A.1.2 MDPs in the RLM Setting

In the context of RLMs, a state  $s \in \mathcal{S}$  is typically defined as a sequence of reasoning steps  $s = (z_0 \dots z_n)$ , where each reasoning step  $z_i$  is a sequence of  $M_i$  tokens  $z_i = (t_i^0, \dots, t_i^{M_i})$ . Each  $t_i^j$  is a token from the RLM's vocabulary, and the total number of tokens per reasoning step  $M_i$  can vary. One can use a special token  $t_{M_i} = t_{end}$  to indicate the end of the reasoning step. Typically, the initial query  $q$  is used as the first reasoning step  $z_0 = q$ . In the study of RLMs, an action  $a \in \mathcal{A}_s$  usually represents appending a new reasoning step  $z^{(a)}$  to the current state  $s = (z_0, \dots, z_n)$  resulting in a new state  $s' = (z_0, \dots, z_n, z^{(a)})$ . Since every action  $a$  is uniquely associated with exactly one reasoning step  $z^{(a)}$  for every  $s = (z_0, \dots, z_n)$  and  $s' = (z_0, \dots, z_n, z_{n+1})$ , we have

$$p(s, a, s') = \begin{cases} 1 & \text{if } z_{n+1} = z^{(a)} \\ 0 & \text{if } z_{n+1} \neq z^{(a)} \end{cases}$$

The definition of the reward function depends on the specific task. A reward commonly seen in reasoning tasks assigns non-zero reward only in the terminal states and hence only at the final reasoning step. This approach reflects the fact that for most tasks, the only final answer can be evaluated against the ground-truth solution to the original query. We call such reward functions *sparse* to clearly distinguish it from other setting in which intermediate rewards can be observed by the algorithm in the non-terminal states. The discount factor  $\gamma$  determines how future rewards influence the current decision-making process. A higher discount factor ( $\gamma \rightarrow 1$ ) places greater emphasis on long-term reasoning success, allowing the model to generate long reasoning sequences, while a lower discount factor prioritizes immediate rewards, incentivizing faster progress and shorter reasoning sequences.

In the RLM setting, a trajectory  $\tau_\pi = (s_0, a_0, \dots, s_T, a_T, s_{T+1})$  represents the progression of states  $s_t$  and actions  $a_t$  ending with a terminal state  $s_{T+1}$  in which no further reasoning steps can be added. The final reasoning step contains the RLM's answer to the original query.

The policy  $\pi(a | s)$  in the context of RLMs defines the probability of selecting an action  $a$  that corresponds to appending a reasoning step  $z^{(a)}$  to the current reasoning sequence represented by the state  $s$ . Since there exists a bijective mapping  $f : \mathcal{A} \rightarrow \mathcal{Z}$  between the action space  $\mathcal{A}$  and the reasoning step space  $\mathcal{Z}$ , the probability distributions can be equated using the change of variables. Formally:

$$\pi(a | s) = \pi(z | s), \quad \text{where } z = f(a).$$

TABLE 2: Overview of mathematical notation used in the paper

Symbol	Description
$\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r, \gamma)$	Markov Decision Process (MDP) definition.
$s \in \mathcal{S}$	A state in the state space, representing a sequence of reasoning steps.
$a \in \mathcal{A}$	An action in the action space, corresponding to selecting the next reasoning step.
$\mathcal{A}_s \subseteq \mathcal{A}$	a set of actions available in state $s$ .
$p(s'   s, a)$	The probability of transition to state $s'$ from state $s$ taking action $a$ in state $s$ .
$r(s)$	The reward received when arriving in state $s$ .
$\gamma \in [0, 1]$	Discount factor, determining the present value of future rewards.
$\pi_\theta(a   s)$	Policy parameterized by $\theta$ , representing the probability of taking action $a$ in state $s$ .
$V_\pi(s)$	Value function under policy $\pi$ , representing the expected return starting from state $s$ .
$Q_\pi(s, a)$	State-action value function under policy $\pi_\theta$ , representing the expected return of taking action $a$ in state $s$ .
$\tau_\pi$	A trajectory consisting of states and actions, $(s_0, a_0, s_1, \dots, s_{T+1})$ following policy $\pi$ .

Based on the definition of the reasoning step and applying the chain rule we can then rewrite the policy as:

$$\pi(z_{t+1} | s_t) = \prod_{j=0}^{M_{t+1}} \pi(t_{t+1}^j | s_t, z_{t+1}^0, \dots, z_{t+1}^{j-1}),$$

In the RLM setting, the state value function  $V(s_t)$  assesses the expected cumulative reward of a partial reasoning sequence  $s_t$ , estimating its overall potential to lead to a successful solution. The state-action value function  $Q(s_t, a_t)$  extends this by quantifying the expected cumulative reward for taking a specific action  $a_t$  (e.g., appending a reasoning step  $z_{t+1}$ ) to the current state  $s_t$  and then following the policy  $\pi$ . It incorporates both the immediate reward for appending the reasoning step and the anticipated future rewards from completing the reasoning sequence. Together, these functions inform and guide the policy  $\pi$  to prioritize actions that maximize the expected cumulative reward. By leveraging  $V(s_t)$  or  $Q(s_t, a_t)$ , the policy can be trained to select reasoning steps that progress toward correct and complete solutions, transforming an LLM into a RLM.

## A.2 Monte Carlo Tree Search (MCTS)

**Monte Carlo Tree Search (MCTS)** is a heuristic search algorithm used for solving MDP problems. MCTS iteratively builds a search tree, representing the underlying MDP state-action space, by aggregating the information obtained from executed MDP trajectories. Let  $\mathcal{T} = (N, E)$  denote the MCTS search tree where  $N \subseteq \mathcal{S}$  is the set of nodes and  $E \subseteq N \times \mathcal{A} \times N$  is the set of directed edges between the nodes. Every node in the MCTS search tree corresponds to a single state in the MDP and every edge corresponds to a single action. Every path from the root to the leaf of the search tree  $\mathcal{T}$  corresponds to a single trajectory in the underlying MDP.

**Edge statistics** The MCTS algorithm stores the following three values for every edge  $s, a$  in the search tree:

- $N(s, a)$  - the visit count of the edge  $(s, a)$  by the algorithm,
- $q(s, a)$  - the estimated state action value of  $(s, a)$ ,
- $r(s, a) = r(s, a, s')$  - the reward received after taking the action  $a$  in the state  $s$  leading to the state  $s'$ ,
- $\beta(s, a)$  - the terminality function indicating if the action  $a$  leads to a terminal state.

**The Algorithm** At the high level, the MCTS begins by initializing the tree with a single starting state  $s_0$  as a root node and performing the following three phases in a loop:

- 1) **Selection** - a leaf-node in the current tree is selected for expanding its child (children).
- 2) **Expansion** - if the selected node does not correspond to a terminal state, it is expanded by taking an action (or multiple actions) in the underlying MDP and by adding the resulting state (states) to the tree as children of the current node. A trajectory unroll is performed for every added node to obtain a reward. "Unroll" refers to simulating a sequence of steps from a newly added node in the tree down to a terminal state. This simulated trajectory represents a hypothetical path the system might take if it continued from the current node. Once the simulation reaches a terminal state, a reward value is calculated based on the outcome of that path.
- 3) **Backpropagation** - update the value estimates and the visit counts for the selected node and all its ancestors based on the obtained reward.

The MCTS algorithm finishes when the stop criterion such as the the number of iterations, the predefined computational budget, or the convergence criterion is met.

## APPENDIX B

### VALUE AND REWARD MODELS

We now proceed to discuss details of value and reward models.

#### B.1 Outcome-Based Reward Models (ORM) vs. Process-Based Reward Models (PRM)

In reinforcement learning environments, reward models estimate the reward for taking an action  $a$  in state  $s$  which leads to state  $s'$ . For reasoning tasks and algorithms like MCTS, which rely on evaluating intermediate steps, it is essential to have models capable of estimating the quality of each step. Two primary families of reward models for such process-based tasks are Outcome-Based Reward Models (ORMs) and Process-Based Reward Models (PRMs). Figure 12 compares both classes of models.

**Outcome-Based Reward Models (ORMs)**, first introduced by Uesato et al. [151], evaluate the reasoning process solely based on the final outcome. These models estimate the reward of the final step in the chain, often modeled in the literature as the likelihood of a correct final answer given the entire reasoning chain  $P(\text{correct}(z_{T+1}) | z_0, \dots, z_{T+1})$  [88], [151] where  $s_{T+1} := z_0, \dots, z_{T+1}$  is the complete reasoning chain consisting of reasoning steps  $z_i$  and  $T + 1$  marks the last reasoning step. ORM are particularly ill-suited

for evaluating intermediate steps for several reasons. First, the training data and objective are inherently misaligned with step-wise evaluation, as they focus exclusively on final outcomes. Second, ORM evaluations tend to be overly pessimistic for intermediate steps since a subsequent erroneous step can obscure the correctness of earlier steps. This observation aligns with Havrilla et al. [59], who noted that ORMs often underestimate the solvability of a problem from an intermediate state and are prone to a high false-negative rate. Furthermore, ORMs lack robustness against false positives, potentially favoring erroneous reasoning steps and misleading the evaluation process.

**Process-Based Reward Models (PRMs)**, introduced by Lightman et al. [88] and Uesato et al. [151], evaluate reasoning on a step-by-step basis. These models estimate the reward of a step, which can be seen as the likelihood of correctness for the  $t$ -th step given its preceding context  $P(\text{correct}(z_t) \mid z_0, \dots, z_t)$  where  $s_t := z_0, \dots, z_t$  is a potentially incomplete reasoning chain and  $z_i$  are reasoning steps and  $z_0$  is the query. PRMs provide more fine-grained feedback and can pinpoint errors in the chain. This step-wise evaluation provides dense rewards given partial responses and helps identify where reasoning deviates from correctness, offering improved interpretability and enabling more targeted improvements in reasoning processes. However, PRMs are computationally expensive to train and require extensive annotations of reasoning steps. These annotations, whether provided by humans or other LLMs, often suffer from limitations: human annotations are scarce, costly, and prone to bias, while prompted LLM-generated annotations [154] are typically of lower quality due to their limited self-evaluation capabilities [99]. Automated methods using for example MCTS such as [96], [155] introduce large computational costs and are prone to false negatives.

## B.2 Outcome-Driven Process-Based Reward Models

Motivated by the need for process-based reward models but constrained by the lack of annotated step-wise labels, certain models that we will refer to as *Outcome-Driven Process-Based Reward Models (O-PRMs)* have been proposed; they combine outcome-based signals with process-based objectives. We show these models in Figure 12. These models rely on process-based data, often automatically generated using MCTS algorithms, where simulations starting from a given step  $s_t$  are performed. The final correctness of these simulated paths is aggregated to create step-wise labels [96], [155] (for other, non-MCTS approaches see [59]). This automation enables scalable data generation for O-PRMs, eliminating the need for extensive human annotation. Although O-PRMs can be categorized as process-based models due to their approximation of step-wise rewards, they remain inherently tied to outcome signals. Some authors [151] suggest that, under certain conditions, outcome signals in mathematical domains can approximate intermediate labels. However, O-PRMs inherit many limitations of ORMs, including susceptibility to false negatives, false positives, and an over-reliance on terminal outcomes. While the aggregation of multiple simulations helps reduce variance, the backtracking process may

still oversimplify complex dependencies within reasoning chains.

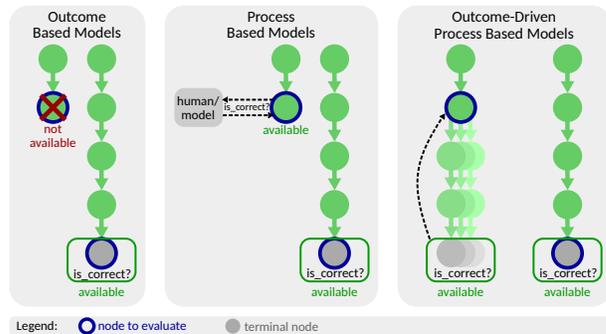


Fig. 12: Comparison of Outcome vs. Process-Based label generation, and the introduction of **Outcome-Driven Process Based Reward Models (O-PRMs)**. Gray nodes mark terminal nodes.

## B.3 Reward Models vs. Value Models

While the distinction between reward models and value models is often blurred in the literature—and their terminology is sometimes used interchangeably—we explicitly differentiate between these model types for evaluating reasoning steps. Additionally, we distinguish two variants of value models:  $v$ -value and  $q$ -value models. This differentiation arises from the distinct roles these models play in reinforcement learning environments.

### B.3.1 Reward Model (RM)

A reward model predicts immediate rewards. In RL, this corresponds to the reward obtained for a transition  $(s, a, s')$  from state  $s$  when taking action  $a$  which results in step  $s'$ . For reasoning, this corresponds to adding a new reasoning step  $a$  to the structure. The new structure is then represented by  $s'$ . Specifically, PRMs – which are preferred over ORMs for MCTS due to the need for action-based evaluation – learn these rewards and can be used to evaluate states (or the transition into a state). This formulation provides a localized, step-level evaluation independent of the overall outcome of the reasoning chain. The reward model is typically trained using labeled data where individual reasoning steps are associated with reward values. While this localized view is advantageous for step-by-step evaluation, it lacks the ability to consider how the current step contributes to the long-term success of the reasoning process. This limitation motivates the introduction of value models.

### B.3.2 Value Model (VM)

Value models provide a more abstract, global evaluation of states and actions by estimating their contribution to future rewards. Unlike reward models, which focus on immediate outcomes, value models consider both current and future rewards, enabling a broader perspective on reasoning quality. For example in reinforcement learning and MCTS, value models play a critical role in guiding the search process. By providing estimates of state or state-action values, they enable more informed decisions about which nodes to expand and explore. We now discuss variants of value models.

**V-Value Model (V-VM).** One variant of a value model is the v-value model which predicts the expected cumulative future reward of a state, denoted as  $V(s)$ . This is equivalent to the state value function in reinforcement learning, which evaluates the long-term potential of the current state  $s$ . A key advantage of V-VMs is their global perspective, as they aggregate future rewards across all possible trajectories originating from the current state. However, V-VMs do not explicitly evaluate individual actions, which may limit their utility in step-level decision-making. Additionally, v-values are often ill-defined at terminal states, where rewards may substitute for state values during training.

**Q-Value Model (Q-VM).** Another variant of a value model is the q-value model. Q-VMs predicts the expected cumulative future reward of taking a specific action  $a$  in a given state  $s$ , denoted as  $Q(s, a)$ . Unlike V-VMs, Q-VMs explicitly associate values with state-action pairs, offering a more granular evaluation. This granularity makes Q-VMs particularly useful for MCTS, where decisions about which edge (action) to expand at a given node (state) are critical. By directly evaluating actions, Q-VMs align naturally with the selection mechanisms in MCTS, guiding the search toward promising paths. Similar to V-VMs, Q-VMs can also be categorized as PQVMs (Process-based Q-Value Models), OQVMs (Outcome-based Q-Value Models), and O-PQVMs (Outcome-driven Process-based Q-Value Models).

The choice between V-VMs and Q-VMs depends on the reasoning task and the specific requirements of the evaluation framework. While V-VMs provide a broader, state-centric evaluation, Q-VMs enable more precise, action-specific guidance. In practice, MCTS often benefits from the use of Q-VMs due to their compatibility with edge-based selection.

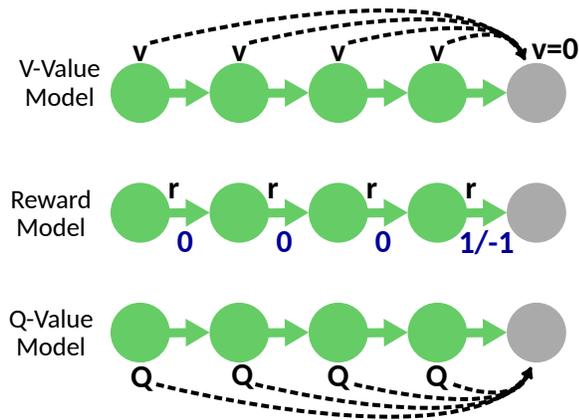


Fig. 13: Comparison of reward, v-value and q-value models in a sparse reward setting (only terminal states receive non-zero rewards). Gray nodes mark terminal nodes. The reward model should predict the rewards for transitioning from one state to another which is 0 for non-terminal states and not providing information. V-VMs and Q-VMs however, predict a global value and are therefore informative for non-terminal states.

### B.3.3 Example: Solving a Mathematical Equation

To illustrate the differences between reward models, value models, and q-value models, consider the task of solving  $x^2 + y^2 = 1$  step-by-step.

- **Reward Model (RM):** A process-based reward model (PRM) might assign a reward  $r(s_t, a_t, s_{t+1})$  for the reasoning step  $a_t = \text{"Substitute } y = \sqrt{1 - x^2}\text{"}$ . This reward

quantifies the quality of the resulting state  $s_{t+1}$ , independent of whether it leads to a correct solution. However, in sparse reward settings (only final steps receive a reward), this reward would be 0.

- **V-Value Model (V-VM):** A V-VM estimates  $V(s_t)$ , representing the expected cumulative reward for the entire expected solution process starting from  $s_t$ . For instance, if  $s_t = \text{"Start with } x^2 + y^2 = 1\text{"}$ ,  $V(s_t)$  considers the long-term potential of all reasoning paths originating from this state.
- **Q-Value Model (Q-VM):** A Q-VM evaluates  $Q(s_t, a_t)$ , predicting the cumulative reward of taking a specific action  $a_t$  (e.g., substituting  $y = \sqrt{1 - x^2}$ ) in state  $s_t$ . This value directly informs whether the action  $a_t$  is likely to lead to a high-quality solution, providing a more granular evaluation compared to the V-VM.

### B.3.4 Summary

By differentiating reward models and value models, and further categorizing value models into V-VMs and Q-VMs, we provide a nuanced framework for evaluating reasoning steps. Reward models offer localized evaluations, while value models incorporate global, long-term perspectives. This global evaluation enables the model to better prioritize reasoning paths that are likely to lead to correct solutions while mitigating the challenges posed by sparse or delayed rewards. Therefore, we advocate for the use of a process-based value model due to the sparsity of reward signals for reasoning tasks. Among value models, Q-VMs are particularly well-suited for MCTS due to their action-specific granularity, which aligns naturally with the tree’s edge-based exploration mechanism. We will demonstrate the practical implications of these distinctions in Appendix D.3.

## B.4 Evaluation Schemes

We also provide additional categorizations and details regarding overall evaluation.

### B.4.1 Evaluation Types

Evaluating reasoning steps in RLMs involves assessing their quality and contribution toward solving a task. Numerical evaluations can be categorized as relative or absolute.

**Relative evaluations** compare multiple steps, often using ranking mechanisms and can be created with, for example, the Bradley-Terry model [20], which is optimized based on pairwise preferences by maximizing the reward gap between chosen and rejected steps.

**Absolute evaluations** assign scalar values to each step, assessing aspects such as coherence, correctness, or helpfulness, using regression-based models. Moreover, evaluation dimensions can also be modeled as binary with classification models. While regression models provide more information, classification models capture correctness more naturally since a statement is usually correct or incorrect. On the other hand, the former ones are more suitable for measuring quality, such as the degree of coherence. Depending on the specific quality being evaluated, the choice between regression and classification models should align with the evaluation’s goals. Additionally, absolute scores can be

transformed into rankings if needed, providing flexibility across various applications.

In addition to numerical evaluations, there are **text-based evaluations**, which are commonly used to provide detailed feedback and guidance for refining reasoning steps. Examples include “LLM-as-a-Judge” [184] (which uses a larger LLM to provide a pairwise comparison or a single graded answer with an explanation) and self-critique approaches [128] that allow models to reflect on and evaluate their own reasoning. These textual evaluations, often including rationales, are particularly useful for structural transformations rather than numerical guidance, enhancing interpretability by offering context and detail.

#### B.4.2 Evaluation of Reasoning Steps

Step-wise evaluations are vital for integrating reasoning into MCTS. Numerical evaluations—whether relative or absolute—provide straightforward metrics to compare nodes and steer exploitation and exploration. Text-based evaluations, in contrast, are better suited for guiding structural refinements rather than directly influencing search paths.

Given that reasoning steps are typically textual sequences, language models are a natural fit for such evaluation tasks. LLM-based approaches can involve *external model* approaches, where a dedicated value model is trained to predict scores, or *internal model* approaches, which leverage existing policy models.

**External model approaches** include value models that predict scalar reward signals (Reward models) [35], [88], [151], reinforcement learning values like state-values (V-value models) [134], state-action values (q-value models), or pairwise models like the Bradley-Terry and PairRM frameworks. A more detailed comparison of reward models, v-value, and q-value models can be found in Appendix B.3.2.

There exist a large range of **internal model approaches** as substitutes for value models. They typically rely on methods like prompting the policy to output scores. Examples include MCT Self-Refine (MCTSr) [176], querying for a binary feedback (e.g., “Is the answer correct? answer “yes” or “no””) [179] and evaluating the probability of the output, leveraging uncertainty metrics such as token entropy or aggregated probabilities [182], and others [178].

**Heuristics** may also serve as substitutes for evaluations in resource-constrained scenarios.

**Simulating** reasoning steps to terminal states for evaluation against golden answers is another option as done for example in MCTS, though often computationally prohibitive.

**External tools** provide an alternative path for evaluation, especially in domain-specific tasks. For programming, compilers can supervise tasks, as seen in Codex [28], self-debugging [31], and similar methods. Program-of-Thought [29] and Program-aided-Language (PAL) [51] use a formal language and Python interpreters to evaluate solutions. In mathematical tasks, ensemble approaches like MathPrompter [71] generate multiple algebraic expressions or Python functions to validate steps. These tool-based approaches excel at detecting errors due to their reliance on precise domain-specific rules, such as compilers for programming or interpreters for mathematics. While their

applicability is limited to well-defined domains, they provide objective and verifiable feedback that complements language models. By injecting precise knowledge into the evaluation process, external tools mitigate model-specific limitations like hallucinations and offer actionable feedback for iterative refinement. This hybrid approach enhances reliability and ensures that the evaluation benefits from both the flexibility of language models and the precision of formal systems.

## APPENDIX C ALGORITHMIC DESCRIPTIONS

### C.1 Reasoning with Monte Carlo Tree Search

#### C.1.1 Setup and Notation

We will now present the details of the training pipeline of x1.

**MDP Design** x1 assumes the MDP following the definition presented in Appendix A.1 with the  $\gamma$  values between  $[0.95, 1]$  to avoid over-penalizing long reasoning sequences. In the RLM setup, the state space and action space of the underlying MDP constitute a tree in which every state  $s$  other than the starting state  $s_0$  has exactly one action  $a_s$  leading to it. This allows us to simplify the notation by omitting actions wherever it’s clear from the context that we are referring to only action leading to a given. For every action  $a$  leading from the state  $s$  to the state  $s'$  we will write:

$$\begin{aligned}\pi(s' | s) &:= \pi(a_{s'} | s) \\ r(s') &:= r(s, a, s') \\ q(s') &:= q(s, a) \\ \tau &:= (s_0, s_1, \dots, s_{T+1})\end{aligned}$$

The final reasoning step in the terminal state contains the RLM’s answer to the original query. The final answer is compared to the ground truth solution, commonly referred to as the golden answer. This matches the common setup in many reasoning tasks and math problems, where no ground truth and no reward source is available for the intermediate reasoning steps.

Consider a trajectory  $\tau := (s_0, s_1, \dots, s_{T+1})$ . We assign a reward of  $r(s_{T+1}) = 1$  if the last reasoning step in the final state  $s_{T+1}$  contains the correct answer and  $r(s_{T+1}) = -1$  otherwise. The state value function simplifies to

$$V_\pi(s_t) = \mathbb{E}_\pi \left[ \gamma^{T-t} r(s_{T+1}) \right] \in [-1, 1] \quad (4)$$

and the state action function can be rewritten as:

$$Q_\pi(s_t) = \begin{cases} r(s_{T+1}), & \text{if } t = T + 1 \\ \gamma V_\pi(s_{t+1}), & \text{otherwise} \end{cases} \in [-1, 1] \quad (5)$$

hence both the value and the state-action value functions are bounded between -1 and 1 for all states and state-action pairs.

**MCTS Design** We define the MCTS tree as in Appendix A.2 as  $\mathcal{T} = (N, E)$ , where  $N$  is a set of nodes, and  $E$  is the set of edges. We use the notation of a node-edge-node relationship denoted by  $(s, a', s')$  where  $s$  represents the origin node,  $a'$  describes the action corresponding to an edge, and  $s'$  denotes the target node. This notation symbolically ties the action and the target state together, as the action uniquely

identifies the target state and is therefore indicative of it.

**The policy model** We use a pretrained LM with parameters  $\theta$  as a policy model and denote it  $\pi_\theta$ . The model autoregressively generates a sequence of tokens. We use a special token ‘End of Intermediate Step’ (eois) to indicate the end of the reasoning step. We use a standard end-of-sequence (eos) token to indicate the end of the final reasoning step concluding the reasoning trajectory.

**The value model** A parametric value model is used to evaluate the quality of states. While MCTS traditionally approximates these values through extensive simulations, such an approach is computationally expensive and impractical in the RLM context. Inspired by AlphaZero [134], which replaces simulations with a parameterized value model, we estimate state-action values (short q-value) for reasoning sequences using a value model — effectively employing a **process-based q-value model**  $Q_\varphi$  (see Appendix B.3). The value model is instantiated as a pretrained transformer-based LM, modified by adding three linear layers and a shifted, rescaled sigmoid activation to align the output domain to the state action function domain  $[-1, 1]$  (see Eq. 5). This setup proved more stable than alternatives, such as a tanh activation or a cropped linear layer. We will show in the following how such a model can be trained and provide a description for the data generation process in Appendix D. During training, we assume access to a final answer verifier, which evaluates the correctness of the model’s final answer and provides the true reward.

### C.1.2 MCTS Algorithm

We now present the algorithmic steps of a Monte Carlo Tree Search variant similar to AlphaZero as implemented in the **x1** reasoning framework. The MCTS search operates in two distinct modes: training and inference. The core difference is that, during training, a final answer verifier evaluates and scores the final reasoning steps, providing a true reward signal that is backpropagated through the MCTS tree. This reward serves as a reliable learning signal for the value model  $Q_\varphi$ . During inference, however, the verifier is unavailable, and decisions rely solely on the value model.

**Notation.** We chose to store all values in nodes instead of edges, which defines the following set of statistics saved for each node  $s$ :

- $N(s)$  - the visit count of node  $s$
- $q(s)$  - the running estimate of the q-value of the transition leading to state  $s$ ,
- $\beta(s)$  - the binary terminality function, returns 1 if the node  $s$  is terminal 0 otherwise.

**Selection.** The selection phase iteratively identifies the most promising child node with a selection policy. We use the following selection policy which is the node-based variant of the PUCT algorithm in AlphaZero [135] (which is defined on edge-based values) without a prior for finding selecting a child of  $s$ :

$$\arg \max_{s_c \in \mathcal{C}(s)} q(s_c) + \frac{\sqrt{N(s) - 1}}{1 + N(s_c)} \cdot \left( c_1 + \log \frac{N(s) + c_2}{c_2} \right)$$

where  $c_1$  and  $c_2$  are hyperparameters controlling the exploration bias, and the other values can be taken from the node statistics.

**Expansion.** We append  $M$  nodes to the selected leaf,  $M$  being a hyperparameter. One of the major challenges in applying RLMs is maintaining the diversity of reasoning paths. By adding  $M$  nodes, we increase the exploration of alternative reasoning paths.

**Backpropagation.** The backpropagation step serves to propagate information from the terminal nodes back to their ancestors. In our implementation, we update the running estimates of the q-values using the following formula:

$$q(s) \leftarrow (1 - \alpha)q(s) + \alpha \gamma \left( \sum_{s_c \in \mathcal{C}(s)} w_s(s_c) \cdot q(s_c) \right),$$

where we look at the node-edge-node tuples  $(s, a_c, s_c)$  and  $s_c \in \mathcal{C}(s)$ . The weights  $w_s(s_c)$  for combining the children q-values are defined over the visit scores of the nodes as follows:

$$w_s(s_c) = \frac{N(s_c)}{\sum_{s_{\tilde{c}} \in \mathcal{C}(s)} N(s_{\tilde{c}})}.$$

**True Reward Propagation.** We improve the quality of the q-values by propagating the real final rewards back through the tree when a terminal state  $s_{T+1}$  is reached. During training, terminal nodes can be evaluated against a reference golden answer  $g^*$  using an external verifier. For actions leading to terminal states, the associated reward is equal to the q-value see Eq. 5. Therefore, instead of using the prediction of the q-value model, we initialize  $q(s_{T+1})$  with the true reward  $r(s_{T+1})$  based on the evaluation of the external verifier. The reward is then backpropagated via the q-values through the tree with our backpropagation operator. This adjustment anchors the q-value model predictions with real reward signals and prevents the q-value model predictions to diverge.

**Best Path Selection.** After  $N$  iterations, MCTS will have formed a tree in which every path corresponds to one of the explored reasoning trajectories. The final reasoning step in a path with the highest terminal value estimate is returned as the final solution.

**Algorithm 1** MCTS for Reasoning (Training mode in blue)

**Input:** Policy model  $\pi_\theta$ , value model  $Q_\varphi$ , question  $z_0$ , golden answer  $g^*$ , binary correctness verifier  $\Gamma$ , number of MCTS iterations  $N$ , number of children expanded in every selection phase  $M$ , exploration constants  $c_1, c_2$ , Backpropagation weight  $\alpha$ .

**Output:** Search tree  $\mathcal{T} = (\mathcal{N}, \mathcal{E})$  containing the best path  $\tau^*$ .

```

1:  $s_0 \leftarrow (z_0)$  {Initialize root node}
2:  $N(s_0) = 0$ 
3:  $\mathcal{N} \leftarrow \{s_0\}$  {Initialize node set}
4:  $\mathcal{E} \leftarrow \emptyset$  {Initialize edge set}
5:  $i \leftarrow 1$ 
6: while  $i \leq N$  or  $\beta(s) \neq 1$  do
7:    $s \leftarrow s_0$  {Start from root node}
8:   ----- Selection -----
9:   while  $s$  is not a leaf node do
10:    {Select child  $s_c \in \mathcal{C}(s)$  with highest selection score}
11:     $s_c \leftarrow \arg \max_{s_c \in \mathcal{C}(s)} q(s_c) + \frac{\sqrt{N(s)-1}}{1+N(s_c)} \left( c_1 + \log \frac{N(s)+c_2}{c_2} \right)$ 
12:     $s \leftarrow s_c$  {Move to the selected child}
13:   end while
14:   ----- Expansion -----
15:   for  $j = 1$  to  $M$  do
16:     $z_c \leftarrow (t^1, \dots, t^{Mz_c}) \sim \pi_\theta$  {Sample a new reasoning step}
17:     $s_c \leftarrow s \frown z_c$  {Append  $z_c$  to the current state  $s$ }
18:     $q(s_c) \leftarrow Q_\varphi(s)$  {Predict with the Q-VM}
19:     $N(s_c) \leftarrow 1$  {Initialize visit count}
20:     $\beta(s_c) \leftarrow 0$  {Initialize terminality function}
21:    if  $s_c$  terminal then
22:       $\beta(s_c) \leftarrow 1$  {Mark as terminal}
23:       $r(s_c) \leftarrow \begin{cases} 1, & \text{if } \Gamma(s_c, g^*) = 1, \\ -1, & \text{if } \Gamma(s_c, g^*) = 0. \end{cases}$  {Check for correctness to determine the reward}
24:       $q(s_c) \leftarrow r(s_c)$  {Overwrite by true reward}
25:    end if
26:     $\mathcal{N} \leftarrow \mathcal{N} \cup \{s_c\}$  {Add the node to the tree}
27:     $\mathcal{E} \leftarrow \mathcal{E} \cup \{(s, s_c)\}$  {Add the edge to the tree}
28:   end for
29:   ----- Backpropagation -----
30:   while  $s \neq s_0$  do
31:     $N(s) \leftarrow N(s) + 1$  {Update the visit count}
32:     $q(s) \leftarrow (1 - \alpha)q(s) + \alpha \gamma \sum_{s_c \in \mathcal{C}(s)} w_s(s_c) q(s_c)$ 
33:    {Update the value}
34:     $s \leftarrow s_p$  {Move to the parent}
35:   end while
36:    $i \leftarrow i + 1$ 
37: end while
38: Best Path Selection:
39: Select the best reasoning sequence  $s_T^*$ .
40:
41: return  $s_T^*$ , all reasoning sequences  $\{s_j^{(i)}\}_j$ 

```

**C.2 Training Phase 1**

**Overall Training Pipeline.** To adequately employ the MCTS-based reasoning scheme introduced in the Appendix C.1, the policy model must be fine-tuned to generate responses in the format of semantically-relevant reasoning steps. The value model – a q-value model in our case – must be trained to accurately estimate the values of the sequences of reasoning steps.

We propose a two-phase training approach designed to let the policy effectively leverage the structured exploration and iterative refinement capabilities of the search process to generate optimal sequences of reasoning steps. A detailed algorithmic description of the pipeline is in Figure 14.

**Phase 1: Supervised Fine-Tuning.** The first phase focuses on preparing the policy and value models to generate and evaluate reasoning trajectories effectively. This is achieved by supervised fine-tuning (SFT) training on a dataset of example sequences of reasoning steps (where intermediate reasoning steps are terminated by an “End of Intermediate Step” eois token). The objective is twofold: (1) to fine-tune the policy model  $\pi_\theta$  to produce semantically coherent reasoning steps, and (2) to train the q-value model  $Q_\varphi$  to accurately assign scalar scores to reasoning trajectories, distinguishing between high-quality and suboptimal reasoning paths.

This supervised fine-tuning phase ensures that the policy can generate reasoning steps consistent with the structured format required for downstream MCTS-based exploration, while the q-value model provides reliable evaluations of intermediate and terminal states. Together, these components form the foundation for the subsequent online reinforcement learning in Phase 2, where the policy and q-value models are further refined through interaction with the reasoning framework.

**C.2.1 Datasets Generation and Preparation**

**Dataset for SFT of the Policy.** Performing SFT of the policy requires a dataset of high-quality reasoning sequences denoted as  $D_{\text{SFT}} = \{(x_{\text{SFT}}^{(i)}, y_{\text{SFT}}^{(i)})\}$ . Each pair in the dataset consists of a prompt  $x_{\text{SFT}}^{(i)}$  composed of a sequence of reasoning steps (for example  $x_{\text{SFT}}^{(i)} = (z_0^{(i)}, \dots, z_j^{(i)})$ ), and a target completion  $y_{\text{SFT}}^{(i)} = z_{j+1}^{(i)}$  which is the subsequent reasoning step or final answer. Appendix D contains a detailed account of the dataset creation and processing. It covers how the special eois token is appended to reasoning steps mark the end of a step during inference.

**Dataset for Q-Value Model Training.** Similarly to SFT, training the q-value model requires a supervised dataset of reasoning sequences and corresponding scores. We denote this dataset  $D_{\text{QVM-train}} = \{(x_{\text{QVM-train}}^{(i)}, y_{\text{QVM-train}}^{(i)})\}$ , with reasoning sequences  $x_{\text{QVM-train}}^{(i)} = (z_0^{(i)}, \dots, z_t^{(i)})$  and target q-value  $y_{\text{QVM-train}}^{(i)}$ . Appendix D explains how this dataset can be generated using an initial list of questions, a base LLM for querying, and a verifier program to label reasoning sequences as conducive to a correct final answer or not.

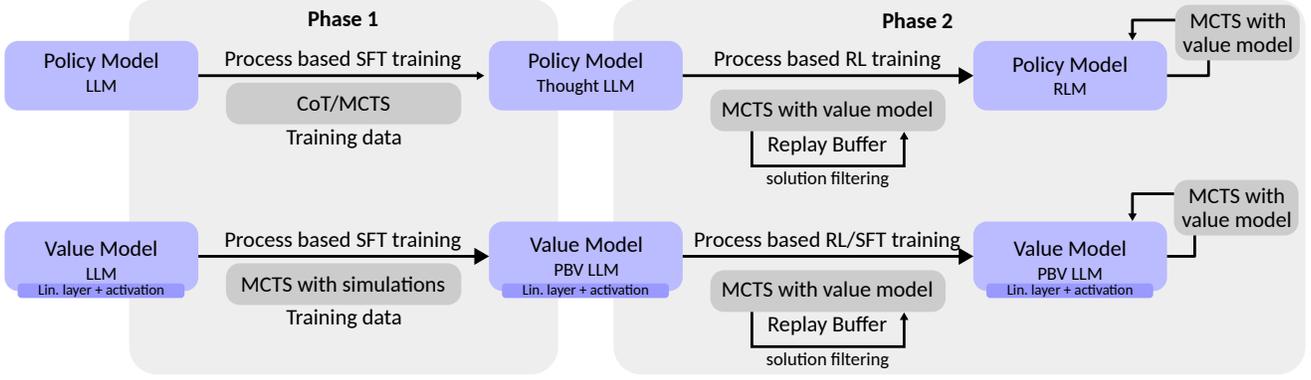


Fig. 14: The two phases of the training pipeline.

---

**Algorithm 2** SFT of Policy Model  $\pi_\theta$  (completion-only)
 

---

**Input:** Policy Model  $\pi_\theta$ , tokenized dataset  $D_{\text{SFT}} = \{(x^{(i)}, y^{(i)})\}$ , training hyperparameters (optimizer, learning rate  $\eta$ , batch size  $B$ , and maximum number of epochs  $E$ ).

**Output:** Fine-tuned policy model  $\pi_\theta$ .

- 1: **for** epoch  $e = 1$  to  $E$  **do**
- 2:   Shuffle dataset  $D_{\text{SFT}}$ .
- 3:   Divide  $D_{\text{SFT}}$  into batches  $\{B_k\}$  of size  $B$ .
- 4:   **for** each batch  $B_k$  **do**
- 5:     Initialize batch loss:  $\mathcal{L}_{\text{batch}} = 0$ .
- 6:     **for** each sample  $(x^{(i)}, y^{(i)}) \in B_k$  **do**
- 7:       Iteratively predict completion tokens:

$$\hat{y}_t^{(i)} \sim \pi_\theta(x_{1:t-1}^{(i)}),$$

where  $x_{1:t-1}^{(i)}$  represents the context (prompt + previously predicted tokens).

- 8:       Compute CE loss for each completion token:  
 $\mathcal{L}^{(i)} = -\sum_{t=1}^{|y^{(i)}|} \log P(\hat{y}_t^{(i)} = y_t^{(i)} | x^{(i)}, \pi_\theta)$ .
  - 9:       Accumulate the loss:  $\mathcal{L}_{\text{batch}} += \mathcal{L}^{(i)}$ .
  - 10:     **end for**
  - 11:     Normalize batch loss:  $\mathcal{L}_{\text{batch}} = \mathcal{L}_{\text{batch}} / |B_k|$ .
  - 12:     Backpropagate gradients, update  $\theta$  via optimizer.
  - 13:   **end for**
  - 14: **end for**
- 

### C.2.2 SFT of the Policy

Supervised fine-tuning (SFT) of the policy is performed on the dataset  $D_{\text{SFT}}$  of prompts and target completions of the next reasoning step. The policy  $\pi_\theta$  is instantiated as a general pretrained LLM. Specifically, we perform ‘completion-only’ SFT such that for every (prompt, target completion) pair, the base model is trained to minimize the cross-entropy loss between its predicted token probabilities and the ground-truth target completion.

### C.2.3 Q-Value Model Training

The q-value model  $Q_\varphi$  is trained on  $D_{\text{QVM-train}}$  to assign appropriate scalar scores to the candidate reasoning trajectories. It is instantiated as a pre-trained LLM with additional linear layer and to which a shifted and rescaled classification head is added; we denote all of its trainable weights as  $\varphi$ . Depending on the reward design, the q-value model can be trained via scalar (least squares) regression if continuous rewards are chosen, or with a classification objective such as the Binary Cross-Entropy (BCE) loss, if trajectories are labelled with binary rewards or as chosen-rejected preference pairs.

By the end of training,  $Q_\varphi$  should output accurate q-value scores, which will later guide policy refinement in Phase II and will improve the search accuracy when used in the MCTS.

---

**Algorithm 3** Fine-Tuning the Q-Value Model  $Q_\varphi$ 


---

**Input:** Q-value model  $Q_\varphi$  (QVM), dataset  $D_{\text{QVM-train}} = \{(x^{(i)}, y^{(i)})\}$ , training hyperparameters (optimizer, learning rate  $\eta$ , batch size  $B$ , and maximum epochs  $E$ ).

**Output:** Fine-tuned q-value model  $Q_\varphi$ .

- 1: **for** epoch  $e = 1$  to  $E$  **do**
  - 2:   Shuffle the dataset  $D_{\text{QVM-train}}$ .
  - 3:   Divide  $D_{\text{QVM-train}}$  into batches  $\{B_k\}$  of size  $B$ .
  - 4:   **for** each batch  $B_k$  **do**
  - 5:     **for** each sample  $(x^{(i)}, y^{(i)}) \in B_k$  **do**
  - 6:       Predict the q-value with QVM  $\hat{y}^{(i)} = Q_\varphi(x^{(i)})$ .
  - 7:       {Compute the loss:}
  - 8:       **if** Regression Loss **then**
  - 9:          $\mathcal{L} = \frac{1}{B} \sum_{(x^{(i)}, y^{(i)})} (\hat{y}^{(i)} - y^{(i)})^2$ .
  - 10:       **end if**
  - 11:       **if** Classification Loss **then**
  - 12:          $\mathcal{L} = \frac{1}{B} \sum_{(x^{(i)}, y^{(i)})} \text{BCE}(\hat{y}^{(i)}, y^{(i)})$ .
  - 13:       **end if**
  - 14:       Backpropagate gradients, update  $\varphi$  via optimizer.
  - 15:     **end for**
  - 16:   **end for**
  - 17: **end for**
-

### C.3 Training Phase 2: RL Tuning of Policy with MCTS

Phase 2 involves generating reasoning sequences from the policy with MCTS and the q-value model, and fine-tuning the policy with an RL-based alignment algorithm to generate better completions. The q-value model must also be continually updated in this training loop to keep in-distribution with the policy’s outputs. Sufficient Phase 1 pre-training of the policy and q-value model is crucial to ensure stable training of the models in Phase 2. The MCTS structure which provides a balanced exploration-exploitation search combined with repeated sampling of the policy ensures sufficient exploration during this online-RL phase. This final training phase returns the finetuned policy and q-value model.

#### C.3.1 Phase 2 Algorithm

Phase 2 uses a set  $D_p = \{p^{(i)}\}$  of prompt questions - these questions may be isolated from the phase 1 dataset  $D_{\text{SFT}}$ . The training process (Algorithm 4) involves a repetition of a MCTS rollout phase followed by a training (reinforcement) phase.

**Data generation: MCTS rollout.** To obtain data for the training, a MCTS tree  $\mathcal{T}^{(i)}$  is build w.r.t. each question  $p^{(i)}$  using the algorithm in Algorithm 1 in training mode. The set of hyperparameters for MCTS  $\Xi_{\text{MCTS}}$ , denotes the number of MCTS iterations  $N$  (per question), the number of children expanded in every selection phase  $M$ , exploration constants  $c_1, c_2$ , and backpropagation weight  $\alpha$ . To enhance the quality of the data, we prune the generated MCTS tree  $\tilde{\mathcal{T}}^{(i)} = (\tilde{N}^{(i)}, \tilde{E}^{(i)})$  to only include paths that reached a terminal state since only these paths received the reward. Then, we extract all nodes and a set of node characteristics from the pruned tree. The dataset comprises of state, action and q-value triplets of the pruned tree:  $\{(s_j^{(i)}, z_j^{(i)}, q(s_j^{(i)}))\}_{s_j \in \tilde{N}^{(i)}}$ . The data is stored in a replay buffer  $\mathcal{R}$ .

**The training: RL phase.** The reinforcement phase samples a batch of reasoning sequences from the replay buffer. From each trajectory, constituent states, actions and value estimates and uses the corresponding values attributed during MCTS to perform RL training (for example with PPO or Reinforce). Alternative schemes may involve selecting preference pairs among trajectories and then aligning the policy using DPO, or simply selecting the most desirable trajectory per question and performing further SFT training.

During this reinforcement phase, the value model is updated to mimic the (backpropagated) values from the MCTS process (Algorithm 7).

---

#### Algorithm 4 Phase 2: RL of the Policy and Q-Value Model

---

**Input:** Policy  $\pi_\theta$ , q-value model  $Q_\varphi$ , dataset  $D_p = \{p^{(i)}\}$ , MCTS hyperparameters  $\Xi_{\text{MCTS}}$ .

**Output:** Trained  $\pi_\theta$  and updated  $Q_\varphi$ .

```

1: for each training iteration do
2:   Rollout
3:   for each question  $p^{(i)} \in D_p$  do
4:     {Generate MCTS tree with  $\pi_\theta$  and  $Q_\varphi$  (Algorithm 1)}
5:      $\mathcal{T}^{(i)} \leftarrow \text{MCTS}(p^{(i)}, Q_\varphi, \pi_\theta, \Xi_{\text{MCTS}})$ 
6:     {Remove incomplete paths from the tree}
7:      $\tilde{\mathcal{T}}^{(i)} \leftarrow \text{Prune}(\mathcal{T}^{(i)})$ 
8:     {Extract nodes and values, store them in replay buffer}
9:      $\mathcal{R} \leftarrow \mathcal{R} \cup \{(s_j^{(i)}, z_j^{(i)}, q(s_j^{(i)}))\}_{s_j \in \tilde{N}^{(i)}}$ 
10:   end for
11:   Training
12:   for each epoch do
13:     Sample a batch  $\mathcal{B}$  from replay buffer  $\mathcal{R}$ .
14:     Update policy  $\pi_\theta$  (Algorithm 5).
15:     Update q-value model  $Q_\varphi$  (Algorithm 7).
16:   end for
17: end for

```

---

#### C.3.2 Policy Update

The policy update is performed on a batch  $\mathcal{D}$  of reasoning sequences. As mentioned above, the reasoning sequences can be decomposed into state-action-value triplets to then perform RL training. We distinguish between three reinforcement methods: standard RL, preference-based RL, or SFT training.

**Standard Policy Gradient RL Methods.** Standard policy gradient methods such as Proximal Policy Optimization (PPO) [131] or REINFORCE [1], [142] are particularly suited for tasks where trajectories are collected (online) and reliably evaluated by the q-value model  $Q_\varphi$ .

PPO relies on the computation of trajectory (reasoning sequence) advantages  $\hat{A}(s_t)$ , which quantify how much better or worse an action taken in a given state is compared to the expected baseline value of that state. The advantage function is estimated by:

$$\hat{A}(s_t) = R_t + \gamma V(s_{t+1}) - V(s_t),$$

where  $R_t$  is the immediate environment reward at step  $t$ ,  $V(s_t)$  is the state value of of state  $s_t$ , and  $\gamma$  is the discount factor. We can derive the state value easily from the q-values obtained via the q-value model or the running estimates in the MCTS as follows:

$$V(s_{t+1}) = \frac{1}{\gamma} Q_\varphi(s_t, a_t),$$

since rewards are sparse. The standard PPO approach trains the critic model from scratch on bootstrapped rewards for this purpose. We introduce an alternative advantage computation scheme that leverages the backpropagated values from Monte Carlo Tree Search (MCTS) in conjunction with  $Q_\varphi$ , as detailed in Algorithm 6. This integration combines MCTS’s exploration and evaluation capabilities with the RL update, enhancing robustness and efficiency in reasoning tasks.

Further regularization can be imposed on the PPO training procedure. To align the policy  $\pi_\theta$  with a reference policy  $\pi_{\text{ref}}$  (usually instantiated as  $\pi_\theta$  before phase 2) during training, the KL divergence  $\text{KL}(\pi_\theta \parallel \pi_{\text{ref}})$ , between the two distributions can be added to the training loss. Additionally, to maintain the diversity of policy generations (and exploration during training), the entropy of the policy distribution can be enhanced by subtracting it from the loss. The entropy penalty is estimated over a batch  $\mathcal{D}$  of state-action pairs  $(s, a)$  where  $s$  denotes a reasoning sequence and  $a$  the next reasoning step. The entropy of a single completion  $a$  is computed by summing the entropy of its individual tokens  $a_{1:|a|}$  of  $a$ :

$$\mathcal{L}_H = -\frac{1}{|\mathcal{D}|} \sum_{(s,a) \in \mathcal{D}} \sum_{a_i \in a} \pi_\theta(a_i | [s, a_{1:i-1}]) \log \pi_\theta(a_i | [s, a_{1:i-1}]).$$

**Direct Preference Optimization (DPO).** DPO [121] aligns the policy to user preferences expressed as pairwise comparisons between reasoning sequences. Given pairs  $(s^+, s^-)$ , where  $s^+$  is preferred over  $s^-$ . This method may not require a process reward/value model. The loss involves the sigmoid function which we denote as  $\sigma$ .

**Supervised Fine-Tuning (SFT).** As a straightforward alternative to RL, high-value reasoning sequences can be selected to perform SFT, i.e. train the policy to maximize the likelihood of these reasoning steps. The high-value reasoning sequences may be selected as terminal nodes having the highest q-value, or highest aggregated intermediate-step values. This approach is inspired by AlphaZero-like frameworks, focusing on iteratively refining the policy to generate high-quality reasoning trajectories without requiring explicit rewards.

### C.3.3 Advantage Calculation (for PPO Policy Updates)

While standard advantage computation in PPO (e.g., via Generalized Advantage Estimation (GAE) [131]) is widely applicable, we propose an alternative approach tailored to our reasoning framework in Algorithm 6. Specifically, for each state/node  $s$ , we leverage the q-value estimates  $q(s)$  obtained during the MCTS process. They were updated in the backpropagation phase to provide a more informed estimate of the q-values incorporating the estimates of the children and potentially true reward signals from terminal paths in the tree. We expect these MCTS-derived values to be more reliable as they incorporate the ground-truth terminal reward, propagated back through the tree, ensuring that a node’s value reflects both its immediate reward and the aggregated values of subsequent child states.

---

### Algorithm 5 Policy Update (PPO, DPO, or SFT)

---

**Input:** Batch  $\mathcal{D}$ , policy  $\pi_\theta$ , reference policy  $\pi_{\text{ref}}$ , learning rate  $\eta$ , clipping parameter  $\varepsilon$ , preference data  $\mathcal{D}_{\text{pref}}$  for DPO.

**Output:** Updated policy  $\pi_\theta$ .

- 1: **Train via PPO**
- 2: Select state-action-value triplets from sequences in  $\mathcal{D}$
- 3: **for each**  $(s_t, a_t, q_t) \in \mathcal{D}$  **do**
- 4:   Compute the policy ratio:  $r_\theta = \frac{\pi_\theta(a_t | s_t)}{\pi_{\text{ref}}(a_t | s_t)}$ .
- 5:   Compute the advantages  $\hat{A}(s_t)$  (Algorithm 6).
- 6:   Compute the PPO loss:
 
$$\mathcal{L}_{\text{PPO}} = \min(r_\theta \hat{A}(s_t), \text{clip}(r_\theta, 1 - \varepsilon, 1 + \varepsilon) \hat{A}(s_t)).$$
- 7: **end for**
- 8: Optional: add KL divergence or entropy regularization.

$$\mathcal{L}_{\text{PPO}} \leftarrow \mathcal{L}_{\text{PPO}} + \lambda_{\text{KL}} \text{KL}(\pi_\theta \parallel \pi_{\text{ref}}) + \lambda_H \mathcal{L}_H.$$

- 9: Perform gradient update to refine  $\pi_\theta$ .
- 10:
- 11: **Train via DPO (pairwise preferences)**
- 12: Select preference pairs of reasoning sequences in  $\mathcal{D}$
- 13: **for each pair**  $(s^+, s^-) \in \mathcal{D}_{\text{pref}}$  **do**
- 14:   Compute DPO objective:

$$\mathcal{L}_{\text{DPO}} = \frac{1}{|\mathcal{D}_{\text{pref}}|} \sum_{(s^+, s^-)} \log \sigma \left( \beta \left( \log \frac{\pi_\theta(s^+)}{\pi_\theta(s^-)} \right) \right).$$

- 15: **end for**
  - 16: Perform gradient update to refine  $\pi_\theta$ .
  - 17:
  - 18: **Train via SFT (single target sequence)**
  - 19: Select high-value reasoning sequences  $s^+$  from  $\mathcal{D}$
  - 20: **for each** reasoning sequence  $s^+$  **do**
  - 21:   Perform SFT on  $s^+$
  - 22: **end for**
- 

---

### Algorithm 6 Advantage Calculation in MCTS Framework

---

**Input:** MCTS Tree  $\mathcal{T} = (N, E)$ , node statistics: rewards and q-values, q-value model  $Q_\varphi$ , discount factor  $\gamma$ , and  $\lambda$ .

**Output:** Advantages  $\{\hat{A}(s_t)\}$ .

- 1: **for each** node  $s_i \in N$  **do**
  - 2:   Compute state values:  $v_{s_{i+1}}^{\text{MCTS}} = \frac{1}{\gamma} q^{\text{MCTS}}(s_i)$
  - 3:   Compute state values:  $v_{s_i}^{\text{MCTS}} = \frac{1}{\gamma} q^{\text{MCTS}}(s_{i-1})$
  - 4:   Compute the advantage on the TD error:  $\hat{A}(s_i) = r(s_i, a_i) + \gamma v_{s_{i+1}}^{\text{MCTS}} - v_{s_i}^{\text{MCTS}}$ .
  - 5: **end for**
- 

### C.3.4 Q-Value Model Update

During phase 2, the q-value model  $Q_\varphi$  is also updated to track MCTS-backtracked value estimates  $q^{\text{MCTS}}(s_t)$  which should be of higher quality (thanks to the final answer verifier and score aggregation from child nodes). For each state-action pair  $(s, a)$ , we train the q-value model  $Q_\varphi$  via squared error minimization, to match its q-value  $Q_\varphi(s, a)$  as closely as possible to the corresponding MCTS-value  $q^{\text{MCTS}}(s')$  which saves the updated q-value of action  $a$  taken in state  $s$  leading to state  $s'$ .

This has the benefit of both improving the accuracy of the value model, and keeping it "in-distribution" with the new policy outputs during this online-RL training.

---

**Algorithm 7** Q-Value Model Update
 

---

**Input:** Batch  $\mathcal{D}$ , q-value model  $Q_\varphi$ , learning rate  $\eta$ .

**Output:** Updated  $Q_\varphi$ .

1: Compute loss:

$$\mathcal{L}_q = \frac{1}{|\mathcal{D}|} \sum_{(s,a,s')} (Q_\varphi(s,a) - q^{\text{MCTS}}(s'))^2.$$

2: Perform gradient update on  $\mathcal{L}_q$ .

---

## APPENDIX D

### DATA GENERATION

#### D.1 Generating Data for Phase 1 Policy Model Training

The objective of this training process is to introduce a new 'End of Intermediate Step' (EOIS) token that serves to delimit individual reasoning steps while preserving the original distribution of the model as much as possible. To achieve this, the model is trained on data generated by itself using greedy decoding.

The training data are derived from eight chain-of-thought (CoT) completions generated for 1,000 questions sampled from the training split of the MATH dataset [63]. These completions are produced using the same model intended for subsequent training with greedy decoding. During this generation process, the reasoning steps in the data are observed to be separated by two consecutive '\n\n'. This observation informs the method of delimitation used to construct pairs of questions and their corresponding sequences of reasoning steps.

For each data point, consisting of a question prompt and its associated target response comprising multiple reasoning steps  $(q^{(i)}, [s_1^{(i)}, \dots, s_n^{(i)}])$ , additional tokens are introduced to explicitly mark the boundaries of the reasoning steps. Specifically, the 'End of Intermediate Step' (EOIS) token is defined and inserted after each reasoning step  $s_j^{(i)}$ , resulting in a modified step  $s_j^{(i)*}$ . Additionally, the 'End of Sequence' (EOS) token is appended to the final reasoning step  $s_n^{(i)}$ , yielding  $s_n^{(i)*} = [s_n^{(i)}; \text{eos}]$ . This augmentation ensures that the model can consistently identify when a final solution has been reached during inference.

For Llama models, it has been empirically observed that introducing an 'assistant' token after each reasoning step enhances the model's effective utilization of the EOIS token. However, this behavior may not generalize to other base models, necessitating careful consideration when applying this approach.

Accordingly, the target sequence for supervised fine-tuning (SFT) is constructed as:

$$y_{\text{SFT}}^{(i)} = [s_1^{(i)}, \text{eois}, \text{assistant}, s_2^{(i)}, \dots, s_n^{(i)}, \text{eos}].$$

This approach yields a training dataset comprising pairs of prompts and their corresponding target completions, formally represented as:

$$D_{\text{SFT}} = \{(q^{(i)}, y_{\text{SFT}}^{(i)})\}.$$

#### D.2 Generating Data for Phase 1 Value Model Training

The original MCTS framework relies on simulations to evaluate a state. Given the state,  $n$  rollouts are performed till a terminal state is reached. The terminal states usually can be evaluated (e.g., in math by comparing it with the golden answer). This enables the distribution of terminal rewards based on their success which are then aggregated to provide a value estimate of the state. These Monte Carlo simulations serve as an estimate of a state's ability to lead to a correct answer. The value estimated in this manner corresponds to the expected cumulative future reward for a given state:

$$V_{\pi_\theta}(s) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=i}^T \gamma^{t-i} r(s_t, a_t) \mid s_i = s \right],$$

where  $T$  is the terminal step of the (sub-) reasoning chain  $\tau = (s_i, a_i, r_i, s_{i+1}, \dots, s_T, a_T, r_T, s_{T+1})$ .

Since rewards are sparse (i.e.,  $r(s_t, a_t) = 0$  for all  $t < T$ ), the value function simplifies to:

$$V_{\pi_\theta}(s_t) = \mathbb{E}_{\pi_\theta} \left[ \gamma^{T-t} r(s_T, a_T) \mid s_t \right].$$

This represents the expected terminal reward, which can be empirically estimated using Monte Carlo (MC) estimates:

$$V_{\pi_\theta}(s_t) \approx \frac{1}{N} \sum_{i=1}^N \gamma^{T-t} r(s_T^{(i)}, a_T^{(i)}) := \hat{V}(s_t),$$

where  $N$  is the number of sampled reasoning chains, and  $s_T^{(i)}, a_T^{(i)}, s_{T+1}^{(i)}$  denote the last transition of the simulation trajectory  $\tau^{(i)} = (s_t, a_t^{(i)}, s_{t+1}^{(i)}, \dots, s_T^{(i)}, a_T^{(i)}, s_{T+1}^{(i)})$  for  $i \in \{1, \dots, N\}$ .

To avoid sample inefficiencies and high computational burdens, AlphaGo Zero [135] and AlphaZero [134] introduce a value model to replace simulations by using its predictions for a state. We follow this approach by defining a process-based value model  $V_\varphi$ . Notably, we train this model with simulation data (instead of true value functions), thereby building a model that predicts state value function estimates  $\hat{V}$ . We denote this model as  $\hat{V}_\varphi$ , parameterized by  $\varphi$ .

Given that the input of a value model is a sequence of reasoning steps - therefore a sequence of tokens, the natural value model architecture is to use an LLM on which one adds linear layer(s) and a suitable output activation function. Typically, it is designed to output a scalar value  $\hat{V}_\varphi(s_t) \in \mathcal{C} \subseteq \mathbb{R}$ .

The core distinction between different modeling approaches to state value functions lies in how rewards are modeled. Depending on whether a binary reward setting or a continuous (bounded) one is used, the aggregation mechanism, model architecture, training loss, and interpretation of the predictions vary. We provide an overview of both scenarios and, although often omitted for simplicity, we consider both  $\gamma = 1$  and  $\gamma \in (0, 1]$  for continuous rewards in our analysis.

##### D.2.1 Binary Rewards: Modeling the Likelihood of a Correct Terminal State

For this approach the rewards are modeled binary, therefore  $r(s_T, a_T) = +1$  for correct solutions and  $r(s_T, a_T) = 0$  for

incorrect solutions. We will adopt a discount factor of  $\gamma = 1$  which we will see aligns more with the interpretation this reward model provides and is widely adopted in literature. This approach corresponds to the value model proposed in AlphaGo Zero [135].

**D.2.1.1 State Value Estimation:** The value function then further simplifies to:

$$V_{\pi_\theta}(s_t) = \mathbb{E}_{\pi_\theta} [r(s_T, a_T) \mid s_t] = \mathbb{P}_{\pi_\theta} (r(s_T, a_T) = 1 \mid s_t)$$

This formulation represents the probability of reaching a correct terminal state from a given state  $s_t$ . Empirically, this probability is estimated using simulations as follows:

$$V_{\pi_\theta}(s_t) \approx \frac{\#\text{correct simulations}}{\#\text{simulations}} := \hat{V}(s_t).$$

**D.2.1.2 Data Generation:** To generate labels for estimating the state-value function during the training of a value model, we use MCTS with simulations till a terminal node is reached and calculate the ratio between the number of correct simulations to the number of simulations. There is one very important detail, for a trajectory  $\tau = (s_i, a_i, r_i, s_{i+1}, \dots, s_{T+1})$  where  $s_{T+1}$  is a terminal state. By definition, the true state value function at  $s_{T+1}$  is zero. However, in training the value model, we avoid instructing it to output zero for terminal states. Instead, in a supervised learning setting, we can identify terminal states and directly compare the model’s predictions against the known correct outcomes (referred to here as “golden answers”). This comparison negates the need to rely solely on the value model to estimate the value of terminal states or to determine the reward associated with transitioning into these states. During inference, while we can still recognize terminal states, we cannot evaluate them by comparing the model’s output to a golden answer. Therefore, an alternative metric is necessary. We train the value model to predict whether transitioning to  $s_{T+1}$  leads to a correct terminal outcome. By learning the relationship between a node’s content and the correctness of the resulting terminal state, the model can estimate the likelihood that a terminal state leads to a correct answer. To approximate the terminal reward during inference, we define:  $r(s_T, a_T, s_{T+1}) \approx \mathbb{1}_{[0.5, 1]}(\hat{V}_\varphi(s_{T+1}))$ . Here  $\hat{V}_\varphi(s_{T+1})$  represents the value predicted by the value model for the terminal state  $s_{T+1}$ . If this predicted likelihood exceeds a threshold (e.g., 0.5), we assign a terminal reward of 1; otherwise, we assign a reward of 0. This approach allows the value model to indirectly influence the terminal reward by predicting the likelihood of a correct outcome. Consequently, during training, terminal rewards serve as labels for terminal states in the value model. It is important to note that  $\hat{V}_\varphi(s_{T+1})$  is not used in any other context but solely to estimate the terminal reward.

$$\hat{V}_\varphi(s_{T+1}) \neq \hat{V}(s_{T+1})$$

This distinction clarifies that the predicted value for the terminal state  $\hat{V}_\varphi(s_{T+1})$  differs from the standard value function’s definition  $\hat{V}(s_{T+1}) = 0$ .

**D.2.1.3 Model Training**  $\hat{V}_\varphi : \mathcal{S} \rightarrow [0, 1]$ : When trained with these labels we obtain a value model  $\hat{V}_\varphi$ , parameterized by  $\varphi$ , that represents the likelihood of a correct terminal state emanating from state  $s_t$ . Therefore, the model

will output values between 0 and 1. To accommodate the binary classification nature of this task, the model should employ a sigmoid activation function in the output layer. The training objective is then to minimize the binary cross-entropy (CE) loss between the predicted probabilities and the empirical estimates derived from the simulations:

$$\mathcal{L}(\varphi) = -\frac{1}{N} \sum_{i=1}^N \left[ y_i \log(\hat{V}_\varphi(s_t^{(i)})) + (1 - y_i) \log(1 - \hat{V}_\varphi(s_t^{(i)})) \right]$$

where  $y_i \in \{0, 1\}$  denotes the binary label indicating whether the  $i$ -th simulation resulted in a correct terminal state.

Employing a binary reward structure offers several benefits. First of all, simplicity since binary rewards simplify the learning process, reducing the complexity associated with continuous reward signals. Moreover, the clear distinction between correct and incorrect states facilitates faster convergence during training making this approach effective. In addition, binary classification is less susceptible to noise in reward signals, ensuring more stable value estimates. Furthermore, this approach aligns with the objectives of reinforcement learning in achieving clear and unambiguous rewards, thereby streamlining the optimization of the policy  $\pi_\theta$ .

## D.2.2 Continuous and Bounded Rewards: Modeling the Expected Future Reward

We model the rewards to be continuous and bounded by allowing values in  $[a, b]$ :

$$V_{\pi_\theta}(s_t) \in [a, b]$$

A common design, is to set the borders to  $-1$  and  $1$  such that a terminal reward is  $r(s_T, a_T) = +1$  for correct terminal states and  $r(s_T, a_T) = -1$  for incorrect states. This approach models the expected future reward as a continuous and bounded value, capturing the degree of correctness or quality of the terminal state. In contrast to the binary reward structure, continuous and bounded rewards provide a more nuanced representation of the outcomes in reasoning tasks. Note, that without discounting this approach resembles the proposed value model of AlphaZero [134].

**D.2.2.1 Bounded rewards:** By constraining rewards within a predefined interval  $[a, b]$ , we effectively create a correctness scale where the extremities represent the definitive outcomes of the reasoning process. Specifically, the lower bound  $a$  corresponds to reaching an incorrect terminal state, while the upper bound  $b$  signifies a correct terminal state. This bounded framework mirrors the spectrum of possible correctness, allowing the model to capture varying degrees of solution quality between these extremes. Such a scale facilitates a more nuanced evaluation of intermediate states, reflecting partial correctness or varying levels of reasoning quality. Moreover, this approach ensures that the reward signals remain interpretable and consistent, fostering a clear distinction between successful and unsuccessful outcomes.

D.2.2.2 State Value Estimation: With a discount factor  $\gamma \in (0, 1]$ , the value function is defined as:

$$V_{\pi_\theta}(s_t) = \mathbb{E} \left[ \gamma^{T-t} r(s_T, a_T) \mid s_t \right],$$

where  $r(s_T, a_T) = b$  for correct terminal states and  $r(s_T, a_T) = a$  for incorrect ones. Empirically, this expectation is approximated by averaging the rewards of the simulations:

$$V_{\pi_\theta}(s_t) \approx \frac{1}{N} \sum_{i=1}^N \gamma^{T-t} r(s_T^{(i)}, a_T^{(i)}) := \hat{V}(s_t),$$

where  $N$  denotes the number of sampled reasoning chains, and  $(s_T^{(i)}, a_T^{(i)}, s_{T+1}^{(i)})$  represent the final transition of the  $i$ -th simulation trajectory  $\tau^{(i)} = (s_t, a_t^{(i)}, s_{t+1}^{(i)}, \dots, s_T^{(i)}, a_T^{(i)}, s_{T+1}^{(i)})$  for  $i \in \{1, \dots, N\}$ . If a discount factor is applied  $\gamma \in (0, 1)$  then each terminal reward is discounted proportional to the number of steps needed to reach the terminal state. This corresponds to the soft estimation proposed by Wang et al. [155]. We want to note that this estimator typically underestimates  $V$  due to its proneness to false negatives [59], [171].

D.2.2.3 Data Generation: Therefore, to generate labels for state-value function estimate pairs to train a value model, we use MCTS with simulations and average the outcomes of the simulations. Therefore, at each newly generated node  $s$  we simulate till a terminal node is reached and we record the depth - the number of steps needed starting from  $s$  (since  $T$  is not identical per trajectory). We then record the the terminal reward which in our case is  $r(s_T, a_T) = 1$  for correct and  $r(s_T, a_T) = -1$  for incorrect answers. Discounted by the depth we can average these rewards and obtain an estimation of the node value which serves as a label for the initial value model training.

D.2.2.4 Model Training  $\hat{V}_\varphi : \mathcal{S} \rightarrow [a, b]$ : The value model  $\hat{V}_\varphi$ , parameterized by  $\varphi$ , is designed to predict the expected terminal reward from any given state  $s_t$ . To accommodate the continuous and bounded nature of this task, the model employs a scaled and shifted sigmoid activation function in the output layer, ensuring that the predictions remain within the range  $[a, b]$ . The training objective is to minimize the mean squared error (MSE) loss between the predicted values and the empirical estimates derived from the simulations:

$$\mathcal{L}(\varphi) = \frac{1}{N} \sum_{i=1}^N \left( \hat{V}_\varphi(s_t^{(i)}) - \gamma^{T-t} r(s_T^{(i)}, a_T^{(i)}) \right)^2.$$

We also experimented with a tanh activation output and a linear layer with clipping of the values. However, both methods proved to be unstable in training in contrast to the scaled and shifted sigmoid layer. A tanh and sigmoid layer naturally bound the output but also push values towards the extremes, enhancing the separation between high and low value estimates. This characteristic can improve the model's ability to distinguish between highly correct and highly incorrect states which is why we are particularly interested in these activation functions.

D.2.2.5 Discounting: Introducing a discount factor  $\gamma$  aligns the value function with the incremental nature of reasoning tasks. Unlike traditional games, where

all moves contribute indirectly and trajectories are not penalized for length, reasoning benefits from discouraging unnecessary or redundant steps. The inclusion of the discount factor  $\gamma$  ensures that rewards achieved sooner have a greater impact on the value function, the model incentivizes reaching correct solutions with fewer steps which ultimately enhances efficiency and suppresses redundancies. Moreover, this models the uncertainty decay in the trajectories; the further into the future a reward lies, the more uncertain its prediction becomes. Discounting naturally reduces the reliance on these uncertain long-term rewards, thereby stabilizing the learning process by focusing on more predictable and immediate outcomes. However, the model's performance becomes sensitive to the choice of  $\gamma$ , requiring careful tuning to balance the influence of immediate versus long-term rewards. Balancing the discount factor is essential to ensure that the model effectively captures the importance of both progress and the final correctness of the reasoning chain.

Employing a continuous and bounded reward structure offers several benefits. Unlike binary rewards, continuous rewards provide a finer distinction between varying degrees of correctness, allowing the model to capture subtle differences in terminal states. Continuous rewards can encode more information about the quality of solutions, facilitating more informed decision-making during the search process. Bounded rewards prevent extreme values, promoting numerical stability and consistent training dynamics. However, this also shows that the choice of reward values and their scaling can significantly impact the learning process, necessitating careful calibration to ensure effective training.

### D.3 State Action Value Function Modeling

The state-action value function, commonly denoted as  $Q_{\pi_\theta}(s_t, a_t)$ , represents the expected cumulative reward of taking action  $a_t$  in state  $s_t$  under policy  $\pi_\theta$ . Formally, it is defined in our framework as:

$$\begin{aligned} Q_{\pi_\theta}(s_t, a_t) &= \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{i=t}^T \gamma^{i-t} r(s_i, a_i) \mid s_t, a_t \right] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{i=t+1}^T \gamma^{i-(t+1)} r(s_i, a_i) \mid s_t, a_t \right] \\ &= r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}} [V_{\pi_\theta}(s_{t+1}) \mid s_t, a_t] \\ &\stackrel{\text{def. } \mathbb{P}}{=} r(s_t, a_t) + \gamma V_{\pi_\theta}(s_{t+1}), \end{aligned}$$

where  $T$  denotes the terminal step of the (sub-) reasoning chain  $\tau = (s_t, a_t, r_t, s_{t+1}, \dots, s_T, a_T, r_T, s_{T+1})$ . In environments characterized by sparse rewards, where  $r(s_t, a_t) = 0$  for all  $t < T$ , the q-value simplifies to:

$$Q_{\pi_\theta}(s_t, a_t) = \gamma V_{\pi_\theta}(s_{t+1}).$$

At terminal states, where the state value  $V_{\pi_\theta}(s_{T+1}) = 0$ , the q-value further reduces to:

$$Q_{\pi_\theta}(s_T, a_T) = r(s_T, a_T).$$

### D.3.1 Process-Based Q-Value Modeling

A process-based q-value model utilizes the same architecture as a process-based Value Model, typically leveraging a LLM enhanced with additional linear layers and an appropriate output activation function. The output is a scalar value  $\hat{Q}_\varphi(s_t, a_t) \in \mathcal{C} \subseteq \mathbb{R}$ . Specifically, the q-value model takes a state-action pair—comprising a sequence of past steps and the current action—and predicts the corresponding q-value based on the aforementioned formulation.

D.3.1.1 Training Data Generation: To train the q-value model, it is essential to compute the q-values for various state-action pairs. For  $t < T$ , q-values can be estimated using  $N$  Monte Carlo simulations as follows:

$$\begin{aligned} Q_{\pi_\theta}(s_t, a_t) &= r(s_t, a_t) + \gamma V_{\pi_\theta}(s_{t+1}) \\ &= \gamma V_{\pi_\theta}(s_{t+1}) \quad (\text{since } r(s_t, a_t) = 0) \\ &\approx \gamma \cdot \frac{1}{N} \sum_{i=1}^N \gamma^{T-(t+1)} r(s_T^{(i)}, a_T^{(i)}) \\ &= \frac{1}{N} \sum_{i=1}^N \gamma^{T-t} r(s_T^{(i)}, a_T^{(i)}) := \hat{Q}(s_t, a_t), \end{aligned}$$

where  $N$  is the number of sampled reasoning chains, and  $\tau^{(i)} = (s_t, a_t^{(i)}, s_{t+1}^{(i)}, \dots, s_T^{(i)}, a_T^{(i)}, s_{T+1}^{(i)})$  represents the  $i$ -th simulation trajectory for  $i \in \{1, \dots, N\}$ . This estimation aligns with the state value estimation under the sparse reward formulation:

$$\hat{Q}(s_t, a_t) = \hat{V}(s_t).$$

For  $t = T$ , the q-value is directly given by the immediate reward:

$$Q_{\pi_\theta}(s_T, a_T) = r(s_T, a_T) = V_{\pi_\theta}(s_{T+1}) \neq \hat{V}(s_T) = 0.$$

D.3.1.2 Reward Modeling: For q-value models the same discussions about reward modeling apply here since the models are trained very similar. This is why omit it here.

### D.3.2 The Difference between Value and Q-Value Models

The difference of VMs and QVMs can be easily shown in how they are used in the evaluation processes of an MCTS algorithm. QVMs predict  $\hat{Q}_\varphi(s_t, a_t)$ , which evaluates the action  $a_t$  taken in state  $s_t$  that deterministically transitions to  $s_{t+1}$ . Thus, the value  $\hat{Q}(s_t, a_t)$  is used to evaluate adding the node  $s_{t+1}$  to the tree. On the other hand, for VMs, adding a node  $s_{t+1}$  to the tree is determined by  $\hat{V}(s_{t+1}) = \frac{1}{\gamma} \hat{Q}_\varphi(s_t, a_t)$ , where  $\gamma$  is the discount factor.

This distinction is making the training processes different. Note that  $s_t \frown a_t = s_{t+1}$ . For QVMs, the training tuples are  $((s_t, a_t), \hat{Q}(s_t, a_t)) = (s_{t+1}, \hat{Q}(s_t, a_t))$  due to the deterministic transition. For VMs, the corresponding training tuples are  $(s_{t+1}, \hat{V}(s_{t+1}))$ . Since we propose training VMs on terminal rewards for terminal states instead of assigning a label of 0, VMs and QVMs become equivalent under the following transformation for any  $t \in \{0, \dots, T\}$  for evaluating adding node  $s_{t+1}$ :

$$\hat{V}(s_{t+1}) = \frac{1}{\gamma} \hat{Q}_\varphi(s_t, a_t).$$

We introduced q-value models since they address a critical inconsistency of value models in terminal states. Specifically, while value models assign a flat value of zero to terminal states, q-value models provide a meaningful evaluation of the final action’s correctness through  $Q_{\pi_\theta}(s_T, a_T) = r(s_T, a_T)$ . This distinction is essential for accurately assessing whether a terminal step leads to a correct or incorrect response during inference.

## REFERENCES

- [1] A. Ahmadian, C. Cremer, M. Gallé, M. Fadaee, J. Kreutzer, O. Pietquin, A. Üstün, and S. Hooker. Back to Basics: Revisiting REINFORCE-Style Optimization for Learning from Human Feedback in LLMs. In L.-W. Ku, A. Martins, and V. Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL ’24, pages 12248–12267, Bangkok, Thailand, Aug. 2024. Association for Computational Linguistics.
- [2] J. Ahn, R. Verma, R. Lou, D. Liu, R. Zhang, and W. Yin. Large Language Models for Mathematical Reasoning: Progresses and Challenges. In N. Falk, S. Papi, and M. Zhang, editors, *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, EACL ’24, pages 225–237, St. Julian’s, Malta, Mar. 2024. Association for Computational Linguistics.
- [3] AI-MO. Aime 2024. <https://huggingface.co/datasets/AI-MO/aimo-validation-aime>, July 2024. accessed 2025-01-19.
- [4] AI-MO. Amc 2024. <https://huggingface.co/datasets/AI-MO/aimo-validation-amc>, July 2024. accessed 2025-01-19.
- [5] A. Amini, S. Gabriel, P. Lin, R. Koncel-Kedziorski, Y. Choi, and H. Hajishirzi. MathQA: Towards Interpretable Math Word Problem Solving with Operation-Based Formalisms, May 2019. arXiv:1905.13319.
- [6] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le, and C. Sutton. Program Synthesis with Large Language Models, Aug. 2021. arXiv:2108.07732.
- [7] A. Bakhtin, L. van der Maaten, J. Johnson, L. Gustafson, and R. Girshick. PHYRE: A New Benchmark for Physical Reasoning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Proceedings of the Thirty-third Annual Conference on Neural Information Processing Systems (NeurIPS ’19)*, volume 32 of *Advances in Neural Information Processing Systems*, pages 5082–5093, Vancouver, Canada, Dec. 2019. Curran Associates.
- [8] T. Ben-Nun and T. Hoefler. Demystifying Parallel and Distributed Deep Learning: An In-depth Concurrency Analysis. *ACM Comput. Surv.*, 52(4):65:1–65:43, Aug. 2019.
- [9] M. Besta, N. Blach, A. Kubicek, R. Gerstenberger, L. Gianinazzi, J. Gajda, T. Lehmann, M. Podstawski, H. Niewiadomski, P. Nyczyk, and T. Hoefler. Graph of Thoughts: Solving Elaborate Problems with Large Language Models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690, Mar. 2024.
- [10] M. Besta, A. C. Catarino, L. Gianinazzi, N. Blach, P. Nyczyk, H. Niewiadomski, and T. Hoefler. HOT: Higher-Order Dynamic Graph Representation Learning with Efficient Transformers. In S. Villar and B. Chamberlain, editors, *Proceedings of the Second Learning on Graphs Conference (LOG ’23)*, volume 231 of *Proceedings of Machine Learning Research*, pages 15:1–15:20, Virtual Event, Nov. 2023. PMLR.
- [11] M. Besta, R. Grob, C. Miglioli, N. Bernold, G. Kwaśniewski, G. Gjini, R. Kanakagiri, S. Ashkboos, L. Gianinazzi, N. Dryden, and T. Hoefler. Motif Prediction with Graph Neural Networks. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD ’22, pages 35–45, Washington DC, USA, Aug. 2022. Association for Computing Machinery.
- [12] M. Besta and T. Hoefler. Parallel and Distributed Graph Neural Networks: An In-Depth Concurrency Analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 46(5):2584–2606, May 2024.
- [13] M. Besta, A. Kubicek, R. Niggl, R. Gerstenberger, L. Weitzendorf, M. Chi, P. Iff, J. Gajda, P. Nyczyk, J. Müller, et al. Multi-Head RAG: Solving Multi-Aspect Problems with LLMs, Nov. 2024. arXiv:2406.05085.

- [14] M. Besta, F. Mededi, Z. Zhang, R. Gerstenberger, N. Blach, P. Nyczyk, M. Copik, G. Kwaśniewski, J. Müller, L. Gianinazzi, et al. Demystifying Chains, Trees, and Graphs of Thoughts, Apr. 2024. arXiv:2401.14295.
- [15] M. Besta, L. Paleari, A. Kubicek, P. Nyczyk, R. Gerstenberger, P. Iff, T. Lehmann, H. Niewiadomski, and T. Hoefler. Check-Embed: Effective Verification of LLM Solutions to Open-Ended Tasks, June 2024. arXiv:2406.02524.
- [16] M. Besta, P. Renc, R. Gerstenberger, P. Sylos Labini, A. Ziogas, T. Chen, L. Gianinazzi, F. Scheidl, K. Szenes, A. Carigiet, P. Iff, G. Kwaśniewski, R. Kanakagiri, C. Ge, S. Jaeger, J. Was, F. Vella, and T. Hoefler. High-Performance and Programmable Attentional Graph Neural Networks with Global Tensor Formulations. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC '23*, Denver, CO, USA, Nov. 2023. Association for Computing Machinery.
- [17] M. Besta, Z. Vonarburg-Shmaria, Y. Schaffner, L. Schwarz, G. Kwaśniewski, L. Gianinazzi, J. Beranek, K. Janda, T. Holenstein, S. Leisinger, P. Tatkowski, A. Ozdemir, A. Balla, M. Copik, P. Lindenberger, M. Konieczny, O. Mutlu, and T. Hoefler. Graph-MineSuite: Enabling High-Performance and Programmable Graph Mining Algorithms with Set Algebra. *Proc. VLDB Endow.*, 14(11):1922–1935, July 2021.
- [18] Z. Bi, K. Han, C. Liu, Y. Tang, and Y. Wang. Forest-of-Thought: Scaling Test-Time Compute for Enhancing LLM Reasoning, Dec. 2024. arXiv:2412.09078.
- [19] Y. Bisk, R. Zellers, R. Le bras, J. Gao, and Y. Choi. PIQA: Reasoning about Physical Commonsense in Natural Language. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7432–7439, Apr. 2020.
- [20] R. A. Bradley and M. E. Terry. Rank Analysis of Incomplete Block Designs: I. The Method of Paired Comparisons. *Biometrika*, 39(3/4):324–345, Dec. 1952.
- [21] T. Burström, V. Parida, T. Lahti, and J. Wincent. AI-Enabled Business-Model Innovation and Transformation in Industrial Ecosystems: A Framework, Model and Outline for Further Research. *Journal of Business Research*, 127:85–95, 2021.
- [22] M. Chang, J. Zhang, Z. Zhu, C. Yang, Y. Yang, Y. Jin, Z. Lan, L. Kong, and J. He. AgentBoard: An Analytical Evaluation Board of Multi-turn LLM Agents. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS '24)*, volume 37 of *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec. 2024. Curran Associates.
- [23] E. Charniak and M. Johnson. Coarse-to-Fine n-Best Parsing and MaxEnt Discriminative Reranking. In K. Knight, H. T. Ng, and K. Oflazer, editors, *Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics, ACL '05*, pages 173–180, Ann Arbor, MI, USA, June 2005. Association for Computational Linguistics.
- [24] G. Chen, M. Liao, C. Li, and K. Fan. AlphaMath Almost Zero: Process Supervision without Process. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS '24)*, volume 37 of *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec. 2024. Curran Associates.
- [25] J. Chen, T. Li, J. Qin, P. Lu, L. Lin, C. Chen, and X. Liang. UniGeo: Unifying Geometry Logical Reasoning via Reformulating Mathematical Expression. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP '22*, pages 3313–3323, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics.
- [26] J. Chen, H. Lin, X. Han, and L. Sun. Benchmarking Large Language Models in Retrieval-Augmented Generation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17754–17762, Mar. 2024.
- [27] J. Chen, J. Tang, J. Qin, X. Liang, L. Liu, E. Xing, and L. Lin. GeoQA: A Geometric Question Answering Benchmark Towards Multimodal Numerical Reasoning. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 513–523, Virtual Event, Aug. 2021. Association for Computational Linguistics.
- [28] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, et al. Evaluating Large Language Models Trained on Code, July 2021. arXiv:2107.03374.
- [29] W. Chen, X. Ma, X. Wang, and W. W. Cohen. Program of Thoughts Prompting: Disentangling Computation from Reasoning for Numerical Reasoning Tasks. *Transactions on Machine Learning Research*, Nov. 2023.
- [30] W. Chen, M. Yin, M. Ku, P. Lu, Y. Wan, X. Ma, J. Xu, X. Wang, and T. Xia. TheoremQA: A Theorem-driven Question Answering Dataset. In H. Bouamor, J. Pino, and K. Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP '23*, pages 7889–7901, Singapore, Dec. 2023. Association for Computational Linguistics.
- [31] X. Chen, M. Lin, N. Schärli, and D. Zhou. Teaching Large Language Models to Self-Debug, Oct. 2023. arXiv:2304.05128.
- [32] K. Chernyshev, V. Polshkov, E. Artemova, A. Myasnikov, V. Stepanov, A. Miasnikov, and S. Tilga. U-MATH: A University-Level Benchmark for Evaluating Mathematical Skills in LLMs, Jan. 2025. arXiv:2412.03205.
- [33] F. Chollet. On the Measure of Intelligence, Nov. 2019. arXiv:1911.01547.
- [34] A. Choudhury, Y. Wang, T. Pelkonen, K. Srinivasan, A. Jain, S. Lin, D. David, S. Soleimanifard, M. Chen, A. Yadav, R. Tijoriwala, D. Samoylov, and C. Tang. MAST: Global Scheduling of ML Training Across Geo-Distributed Datacenters at Hyperscale. In *Proceedings of the 18th USENIX Symposium on Operating Systems Design and Implementation, OSDI '24*, pages 563–580, Santa Clara, CA, USA, July 2024. USENIX Association.
- [35] P. Christiano, J. Leike, T. B. Brown, M. Martic, S. Legg, and D. Amodei. Deep Reinforcement Learning from Human Preferences, Feb. 2023. arXiv:1706.03741.
- [36] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton, R. Nakano, C. Hesse, and J. Schulman. Training Verifiers to Solve Math Word Problems, Nov. 2021. arXiv:2110.14168.
- [37] M. Copik, R. Böhringer, A. Calotoiu, and T. Hoefler. FMI: Fast and Cheap Message Passing for Serverless Functions. In *Proceedings of the 37th International Conference on Supercomputing, ICS '23*, pages 373–385, Orlando, FL, USA, June 2023. Association for Computing Machinery.
- [38] M. Copik, G. Kwaśniewski, M. Besta, M. Podstawski, and T. Hoefler. SeBS: A Serverless Benchmark Suite for Function-as-a-Service Computing. In *Proceedings of the 22nd International Middleware Conference, Middleware '21*, pages 64–78, Virtual Event, Dec. 2021. Association for Computing Machinery.
- [39] G. Cui, L. Yuan, Z. Wang, H. Wang, W. Li, B. He, Y. Fan, T. Yu, Q. Xu, W. Chen, et al. Process Reinforcement through Implicit Rewards. <https://curvy-check-498.notion.site/Process-Reinforcement-through-Implicit-Rewards-15f4fbc9c42180f1b498cc9b2ea f896f>, Jan. 2025.
- [40] D. De Sensi, T. De Matteis, K. Taranov, S. Di Girolamo, T. Rahn, and T. Hoefler. Noise in the Clouds: Influence of Network Performance Variability on Application Scalability. *Proc. ACM Meas. Anal. Comput. Syst.*, 6(3):49:1–49:27, Dec. 2022.
- [41] M. DeLorenzo, A. B. Chowdhury, V. Gohil, S. Thakur, R. Karri, S. Garg, and J. Rajendran. Make Every Move Count: LLM-based High-Quality RTL Code Generation Using MCTS, Feb. 2024. arXiv:2402.03289.
- [42] X. Deng, Y. Gu, B. Zheng, S. Chen, S. Stevens, B. Wang, H. Sun, and Y. Su. Mind2Web: Towards a Generalist Agent for the Web. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Proceedings of the Thirty-seventh Annual Conference on Neural Information Processing Systems (NeurIPS '23)*, volume 36 of *Advances in Neural Information Processing Systems*, pages 28091–28114, New Orleans, LA, USA, Dec. 2023. Curran Associates.
- [43] Y. Deng, W. Zhang, Z. Chen, and Q. Gu. Rephrase and Respond: Let Large Language Models Ask Better Questions for Themselves, Apr. 2024. arXiv:2311.04205.
- [44] X. Dong, M. Teleki, and J. Caverlee. A Survey on LLM Inference-Time Self-Improvement, Dec. 2024. arXiv:2412.14352.
- [45] I. El Naqa, M. A. Haider, M. L. Giger, and R. K. Ten Haken. Artificial Intelligence: Reshaping the Practice of Radiological Sciences in the 21st Century. *British Journal of Radiology*, 93(1106):20190855, Jan. 2020.
- [46] A. Elliott. *The Culture of AI: Everyday Life and the Digital Revolution*. Routledge, 2018.
- [47] S. Es, J. James, L. Espinosa-Anke, and S. Schockaert. RAGAS: Automated Evaluation of Retrieval Augmented Generation, Sept. 2023. arXiv:2309.15217.
- [48] X. Feng, Z. Wan, M. Wen, S. M. McAleer, Y. Wen, W. Zhang, and

- J. Wang. AlphaZero-Like Tree-Search Can Guide Large Language Model Decoding and Training, Feb. 2024. arXiv:2309.17179.
- [49] J. Frohberg and F. Binder. CRASS: A Novel Data Set and Benchmark to Test Counterfactual Reasoning of Large Language Models. In N. Calzolari, F. B  chet, P. Blache, K. Choukri, C. Cieri, T. Declerck, S. Goggi, H. Isahara, B. Maegaard, J. Mariani, H. Mazo, J. Odiijk, and S. Piperidis, editors, *Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC '22*, pages 2126–2140, Marseille, France, June 2022. European Language Resources Association.
- [50] Y. Fu, L. Xue, Y. Huang, A.-O. Brabete, D. Ustiugov, Y. Patel, and L. Mai. ServerlessLLM: Low-Latency Serverless Inference for Large Language Models. In *Proceedings of the 18th USENIX Symposium on Operating Systems Design and Implementation, OSDI '24*, pages 135–153, Santa Clara, CA, USA, July 2024. USENIX Association.
- [51] L. Gao, A. Madaan, S. Zhou, U. Alon, P. Liu, Y. Yang, J. Callan, and G. Neubig. PAL: Program-Aided Language Models, Jan. 2023. arXiv:2211.10435.
- [52] S. N. Giest and B. Klievink. More Than a Digital System: How AI is Changing the Role of Bureaucrats in Different Organizational Contexts. *Public Management Review*, 26(2):379–398, 2024.
- [53] E. Glazer, E. Erdil, T. Besiroglu, D. Chicharro, E. Chen, A. Gunning, C. F. Olsson, J.-S. Denain, A. Ho, E. de Oliveira Santos, O. J  rvinieniemi, M. Barnett, R. Sandler, M. Vrzala, J. Sevilla, Q. Ren, E. Pratt, L. Levine, G. Barkley, N. Stewart, B. Grechuk, T. Grechuk, S. V. Enugandla, and M. Wildon. FrontierMath: A Benchmark for Evaluating Advanced Mathematical Reasoning in AI, Dec. 2024. arXiv:2411.04872.
- [54] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, et al. The Llama 3 Herd of Models, Nov. 2024. arXiv:2407.21783.
- [55] X. Guan, Y. Liu, X. Lu, B. Cao, B. He, X. Han, L. Sun, J. Lou, B. Yu, Y. Lu, and H. Lin. Search, Verify and Feedback: Towards Next Generation Post-Training Paradigm of Foundation Models via Verifier Engineering, Nov. 2024. arXiv:2411.11504.
- [56] X. Guan, L. L. Zhang, Y. Liu, N. Shang, Y. Sun, Y. Zhu, F. Yang, and M. Yang. rStar-Math: Small LLMs Can Master Math Reasoning with Self-Evolved Deep Thinking, Jan. 2025. arXiv:2501.04519.
- [57] K. Guu, K. Lee, Z. Tung, P. Pasupat, and M.-W. Chang. REALM: Retrieval-Augmented Language Model Pre-Training, Feb. 2020. arXiv:2002.08909.
- [58] S. Han, H. Schoelkopf, Y. Zhao, Z. Qi, M. Riddell, L. Benson, L. Sun, E. Zubova, Y. Qiao, M. Burtell, D. Peng, J. Fan, Y. Liu, B. Wong, M. Sailor, A. Ni, L. Nan, J. Kasai, T. Yu, R. Zhang, S. Joty, A. R. Fabbri, W. Kryscinski, X. V. Lin, C. Xiong, and D. Radev. FOLIO: Natural Language Reasoning with First-Order Logic, Oct. 2024. arXiv:2209.00840.
- [59] A. Havrilla, S. C. Raparthy, C. Nalmpantis, J. Dwivedi-Yu, M. Zhuravinskyi, E. Hambro, and R. Raileanu. GLoRe: When, Where, and How to Improve LLM Reasoning via Global and Local Refinements. In R. Salakhutdinov, Z. Kolter, K. Heller, A. Weller, N. Oliver, J. Scarlett, and F. Berkenkamp, editors, *Proceedings of the 41st International Conference on Machine Learning (ICML '24)*, volume 235 of *Proceedings of Machine Learning Research*, pages 17719–17733, Vienna, Austria, July 2024. PMLR.
- [60] C. He, R. Luo, Y. Bai, S. Hu, Z. Thai, J. Shen, J. Hu, X. Han, Y. Huang, Y. Zhang, J. Liu, L. Qi, Z. Liu, and M. Sun. Olympiad-Bench: A Challenging Benchmark for Promoting AGI with Olympiad-Level Bilingual Multimodal Scientific Problems. In L.-W. Ku, A. Martins, and V. Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL '24, pages 3828–3850, Bangkok, Thailand, Aug. 2024. Association for Computational Linguistics.
- [61] D. Hendrycks, S. Basart, S. Kadavath, M. Mazeika, A. Arora, E. Guo, C. Burns, S. Puranik, H. He, D. Song, and J. Steinhardt. Measuring Coding Challenge Competence with APPS. In J. Vanschoren and S. Yeung, editors, *Proceedings of the Thirty-fifth Neural Information Processing Systems: Track on Datasets and Benchmarks*, volume 1 of *NeurIPS '21*, Virtual Event, Dec. 2021.
- [62] D. Hendrycks, C. Burns, S. Basart, A. Zou, M. Mazeika, D. Song, and J. Steinhardt. Measuring Massive Multitask Language Understanding. In *Proceedings of the Ninth International Conference on Learning Representations, ICLR '21*, Virtual Event, May 2021.
- [63] D. Hendrycks, C. Burns, S. Kadavath, A. Arora, S. Basart, E. Tang, D. Song, and J. Steinhardt. Measuring Mathematical Problem Solving with the MATH Dataset. In J. Vanschoren and S. Yeung, editors, *Proceedings of the Thirty-fifth Conference on Neural Information Processing Systems: Track on Datasets and Benchmarks*, *NeurIPS '21*, Virtual Event, Dec. 2021.
- [64] A. Holtzman, J. Buys, L. Du, M. Forbes, and Y. Choi. The Curious Case of Neural Text Degeneration. In *Proceedings of the Eighth International Conference on Learning Representations, ICLR '20*, Virtual Event, Apr. 2020.
- [65] M. J. Hosseini, H. Hajishirzi, O. Etzioni, and N. Kushman. Learning to Solve Arithmetic Word Problems with Verb Categorization. In A. Moschitti, B. Pang, and W. Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP '14*, pages 523–533, Doha, Qatar, Oct. 2014. Association for Computational Linguistics.
- [66] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-Rank Adaptation of Large Language Models. In *Proceedings of the Tenth International Conference on Learning Representations, ICLR '22*, Virtual Event, Apr. 2022.
- [67] J. Huang and K. C.-C. Chang. Towards Reasoning in Large Language Models: A Survey. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Findings of the Association for Computational Linguistics: ACL 2023*, pages 1049–1065, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [68] J. Huang, X. Chen, S. Mishra, H. S. Zheng, A. W. Yu, X. Song, and D. Zhou. Large Language Models Cannot Self-Correct Reasoning Yet. In *Proceedings of the Twelfth International Conference on Learning Representations, ICLR '24*, Vienna, Austria, May 2024.
- [69] J. Huang, S. S. Gu, L. Hou, Y. Wu, X. Wang, H. Yu, and J. Han. Large Language Models Can Self-Improve, Oct. 2022. arXiv:2210.11610.
- [70] Y. Huang, M. Kleindessner, A. Munishkin, D. Varshney, P. Guo, and J. Wang. Benchmarking of Data-Driven Causality Discovery Approaches in the Interactions of Arctic Sea Ice and Atmosphere. *Frontiers in Big Data*, 4(32):642182:1–642182:19, Aug. 2021.
- [71] S. Imani, L. Du, and H. Shrivastava. MathPrompter: Mathematical Reasoning using Large Language Models, Mar. 2023. arXiv:2303.05398.
- [72] A. Q. Jiang, W. Li, J. M. Han, and Y. Wu. LISA: Language models of ISAbelle proofs. In *Proceedings of the 6th Conference on Artificial Intelligence and Theorem Proving, AITP '21*, Aussois, France, Sept. 2021.
- [73] J. Jiang, S. Gan, Y. Liu, F. Wang, G. Alonso, A. Klimovic, A. Singla, W. Wu, and C. Zhang. Towards Demystifying Serverless Machine Learning Training. In *Proceedings of the 2021 International Conference on Management of Data, SIGMOD '21*, pages 857–871, Virtual Event, June 2021. Association for Computing Machinery.
- [74] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. R. Narasimhan. SWE-bench: Can Language Models Resolve Real-world Github Issues? In *Proceedings of the Twelfth International Conference on Learning Representations, ICLR '24*, Vienna, Austria, May 2024.
- [75] E. Kislev. *Relationships 5.0: How AI, VR, and Robots Will Reshape Our Emotional Lives*. Oxford University Press, 2022.
- [76] W. Knight. OpenAI Unveils New A.I. That Can ‘Reason’ Through Math and Science Problems. <https://www.nytimes.com/2024/12/20/technology/openai-new-ai-math-science.html>, Dec. 2024. accessed 2024-12-27.
- [77] L. Kocsis and C. Szepesv  ri. Bandit Based Monte-Carlo Planning. In J. F  rnkranz, T. Scheffer, and M. Spiliopoulou, editors, *Proceedings of the European Conference on Machine Learning ECML '06*, volume 4212 of *Lecture Notes in Computer Science (LNAI)*, pages 282–293, Berlin, Germany, Sept. 2006. Springer.
- [78] K. Kondo, S. Sugawara, and A. Aizawa. Probing Physical Reasoning with Counter-Commonsense Context. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, ACL '23, pages 603–612, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [79] W. Kryscinski, B. McCann, C. Xiong, and R. Socher. Evaluating the Factual Consistency of Abstractive Text Summarization. In B. Webber, T. Cohn, Y. He, and Y. Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP '20*, pages 9332–9346, Virtual Event, Nov. 2020. Association for Computational Linguistics.
- [80] P. Kumar. *Artificial Intelligence: Reshaping Life and Business*. BPB Publications, 2019.

- [81] Y. Lai, C. Li, Y. Wang, T. Zhang, R. Zhong, L. Zettlemoyer, W.-T. Yih, D. Fried, S. Wang, and T. Yu. DS-1000: A Natural and Reliable Benchmark for Data Science Code Generation. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 18319–18345, Honolulu, HI, USA, July 2023. PMLR.
- [82] Y. Leviathan, M. Kalman, and Y. Matias. Fast Inference from Transformers via Speculative Decoding. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning (ICML '23)*, volume 202 of *Proceedings of Machine Learning Research*, pages 19274–19286, Honolulu, HI, USA, July 2023. PMLR.
- [83] P. Lewis, E. Perez, A. Piktus, F. Petroni, N. Karpukhin, N. Goyal, H. Küttler, M. Lewis, W.-t. Yih, T. Rocktäschel, S. Riedel, and D. Kiela. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Proceedings of the Thirty-fourth Annual Conference on Neural Information Processing Systems (NeurIPS '20)*, volume 33 of *Advances in Neural Information Processing Systems*, pages 9459–9474, Virtual Event, Dec. 2020. Curran Associates.
- [84] X. Li, G. Dong, J. Jin, Y. Zhang, Y. Zhou, Y. Zhu, P. Zhang, and Z. Dou. Search-o1: Agentic Search-Enhanced Large Reasoning Models, Jan. 2025. arXiv:2501.05366.
- [85] X. L. Li, A. Holtzman, D. Fried, P. Liang, J. Eisner, T. Hashimoto, L. Zettlemoyer, and M. Lewis. Contrastive Decoding: Open-ended Text Generation as Optimization. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL '23, pages 12286–12312, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [86] X. L. Li and P. Liang. Prefix-Tuning: Optimizing Continuous Prompts for Generation. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, ACL-IJCNLP '21, pages 4582–4597, Virtual Event, Aug. 2021. Association for Computational Linguistics.
- [87] M. Liao, W. Luo, C. Li, J. Wu, and K. Fan. MARIO: MATH Reasoning with code Interpreter Output – A Reproducible Pipeline, Feb. 2024. arXiv:2401.08190.
- [88] H. Lightman, V. Kosaraju, Y. Burda, H. Edwards, B. Baker, T. Lee, J. Leike, J. Schulman, I. Sutskever, and K. Cobbe. Let's Verify Step by Step. In *Proceedings of the Twelfth International Conference on Learning Representations, ICLR '24*, Vienna, Austria, May 2024.
- [89] A. Liu, S. Swayamdipta, N. A. Smith, and Y. Choi. WANLI: Worker and AI Collaboration for Natural Language Inference Dataset Creation. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6826–6847, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics.
- [90] C. Liu, J. Shen, H. Xin, Z. Liu, Y. Yuan, H. Wang, W. Ju, C. Zheng, Y. Yin, L. Li, M. Zhang, and Q. Liu. FIMO: A Challenge Formal Dataset for Automated Theorem Proving, Dec. 2023. arXiv:2309.04295.
- [91] X. Liu, H. Yu, H. Zhang, Y. Xu, X. Lei, H. Lai, Y. Gu, H. Ding, K. Men, K. Yang, S. Zhang, X. Deng, A. Zeng, Z. Du, C. Zhang, S. Shen, T. Zhang, Y. Su, H. Sun, M. Huang, Y. Dong, and J. Tang. AgentBench: Evaluating LLMs as Agents, Oct. 2023. arXiv:2308.03688.
- [92] P. Lu, H. Bansal, T. Xia, J. Liu, C. Li, H. Hajishirzi, H. Cheng, K.-W. Chang, M. Galley, and J. Gao. MathVista: Evaluating Mathematical Reasoning of Foundation Models in Visual Contexts. In *Proceedings of the Twelfth International Conference on Learning Representations, ICLR '24*, Vienna, Austria, May 2024.
- [93] P. Lu, R. Gong, S. Jiang, L. Qiu, S. Huang, X. Liang, and S.-C. Zhu. Inter-GPS: Interpretable Geometry Problem Solving with Formal Language and Symbolic Reasoning. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, ACL-IJCNLP '21, pages 6774–6786, Virtual Event, Aug. 2021. Association for Computational Linguistics.
- [94] P. Lu, L. Qiu, K.-W. Chang, Y. N. Wu, S.-C. Zhu, T. Rajpurohit, P. Clark, and A. Kalyan. Dynamic Prompt Learning via Policy Gradient for Semi-Structured Mathematical Reasoning. In *Proceedings of the Eleventh International Conference on Learning Representations, ICLR '23*, Kigali, Rwanda, May 2023.
- [95] P. Lu, L. Qiu, W. Yu, S. Welleck, and K.-W. Chang. A Survey of Deep Learning for Mathematical Reasoning. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL '23, pages 14605–14631, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [96] L. Luo, Y. Liu, R. Liu, S. Phatale, M. Guo, H. Lara, Y. Li, L. Shu, Y. Zhu, L. Meng, J. Sun, and A. Rastogi. Improve Mathematical Reasoning in Language Models by Automated Process Supervision, Dec. 2024. arXiv:2406.06592.
- [97] M. Luo, S. Kumbhar, M. shen, M. Parmar, N. Varshney, P. Banerjee, S. Aditya, and C. Baral. Towards LogiGLUE: A Brief Survey and a Benchmark for Analyzing Logical Reasoning Capabilities of Language Models, Mar. 2024. arXiv:2310.00836.
- [98] Y. Lyu, Z. Li, S. Niu, F. Xiong, B. Tang, W. Wang, H. Wu, H. Liu, T. Xu, and E. Chen. CRUD-RAG: A Comprehensive Chinese Benchmark for Retrieval-Augmented Generation of Large Language Models, July 2024. arXiv:2401.17043.
- [99] A. Madaan, N. Tandon, P. Gupta, S. Hallinan, L. Gao, S. Wiegrefe, U. Alon, N. Dziri, S. Prabhunoye, Y. Yang, S. Gupta, B. Majumder, K. Hermann, S. Welleck, A. Yazdanbakhsh, and P. Clark. Self-Refine: Iterative Refinement with Self-Feedback, May 2023. arXiv:2303.17651.
- [100] F. Mai, N. Cornille, and M.-F. Moens. Improving Language Modeling by Increasing Test-time Planning Compute. In *Proceedings of the Eighth Widening NLP Workshop, WiNLP '24*, Miami, FL, USA, Nov. 2024.
- [101] A. Malinin and M. Gales. Uncertainty Estimation in Autoregressive Structured Prediction. In *Proceedings of the Ninth International Conference on Learning Representations, ICLR '21*, Virtual Event, May 2021.
- [102] R. Manvi, A. Singh, and S. Ermon. Adaptive Inference-Time Compute: LLMs Can Predict If They Can Do Better, Even Mid-Generation, Oct. 2024. arXiv:2410.02725.
- [103] Y. Mao, Y. Kim, and Y. Zhou. CHAMP: A Competition-level Dataset for Fine-Grained Analyses of LLMs' Mathematical Reasoning Capabilities. In L.-W. Ku, A. Martins, and V. Srikumar, editors, *Findings of the Association for Computational Linguistics: ACL 2024*, pages 13256–13274, Bangkok, Thailand, Aug. 2024. Association for Computational Linguistics.
- [104] A. Masry, X. L. Do, J. Q. Tan, S. Joty, and E. Hoque. ChartQA: A Benchmark for Question Answering about Charts with Visual and Logical Reasoning. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2263–2279, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [105] C. Metz. In Two Moves, AlphaGo and Lee Sedol Redefined the Future. <https://www.wired.com/2016/03/two-moves-alpha-go-lee-sedol-redefined-future/>, Mar. 2016. Wired.
- [106] G. Mialon, C. Fourrier, C. Swift, T. Wolf, Y. LeCun, and T. Scialom. GAIA: A Benchmark for General AI Assistants, Nov. 2023. arXiv:2311.12983.
- [107] X. Miao, C. Shi, J. Duan, X. Xi, D. Lin, B. Cui, and Z. Jia. SpotServe: Serving Generative Large Language Models on Preemptible Instances. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2, ASPLOS '24*, pages 1112–1127, La Jolla, CA, USA, Apr. 2024. Association for Computing Machinery.
- [108] T. Mihaylov, P. Clark, T. Khot, and A. Sabharwal. Can a Suit of Armor Conduct Electricity? A New Dataset for Open Book Question Answering. In E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, EMNLP '18*, pages 2381–2391, Brussels, Belgium, Nov. 2018. Association for Computational Linguistics.
- [109] I. Mirzadeh, K. Alizadeh, H. Shahrokhi, O. Tuzel, S. Bengio, and M. Farajtabar. GSM-Symbolic: Understanding the Limitations of Mathematical Reasoning in Large Language Models, Oct. 2024. arXiv:2410.05229.
- [110] J. M. Mooij, J. Peters, D. Janzing, J. Zscheischler, and B. Schölkopf. Distinguishing Cause from Effect Using Observational Data: Methods and Benchmarks. *Journal of Machine Learning Research*, 17(32):1–102, 2016.
- [111] P. Moritz, R. Nishihara, S. Wang, A. Tumanov, R. Liaw, E. Liang, M. Elibol, Z. Yang, W. Paul, M. I. Jordan, and I. Stoica. Ray:

- A Distributed Framework for Emerging AI Applications. In *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation, OSDI '18*, pages 561–577, Carlsbad, CA, Oct. 2018. USENIX Association.
- [112] Y. Nie, A. Williams, E. Dinan, M. Bansal, J. Weston, and D. Kiela. Adversarial NLI: A New Benchmark for Natural Language Understanding. In D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL '20*, pages 4885–4901, Virtual Event, July 2020. Association for Computational Linguistics.
- [113] T. Niven and H.-Y. Kao. Probing Neural Network Comprehension of Natural Language Arguments. In A. Korhonen, D. Traum, and L. Márquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, ACL '19*, pages 4658–4664, Florence, Italy, July 2019. Association for Computational Linguistics.
- [114] OpenAI. Introducing ChatGPT. <https://openai.com/index/chatgpt/>, Nov. 2022. accessed 2024-12-27.
- [115] OpenAI. Hello GPT-4o. <https://openai.com/index/hello-gpt-4o/>, May 2024. accessed 2025-01-01.
- [116] OpenAI. Introducing OpenAI o1. <https://openai.com/o1/>, 2024. accessed 2024-12-27.
- [117] R. Y. Pang, W. Yuan, H. He, K. Cho, S. Sukhbaatar, and J. E. Weston. Iterative Reasoning Preference Optimization. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS '24)*, volume 37 of *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec. 2024. Curran Associates.
- [118] S. Qiao, Y. Ou, N. Zhang, X. Chen, Y. Yao, S. Deng, C. Tan, F. Huang, and H. Chen. Reasoning with Language Model Prompting: A Survey. In A. Rogers, J. Boyd-Graber, and N. Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL '23, pages 5368–5393, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [119] Y. Qin, X. Li, H. Zou, Y. Liu, S. Xia, Z. Huang, Y. Ye, W. Yuan, H. Liu, Y. Li, and P. Liu. O1 Replication Journey: A Strategic Progress Report – Part 1, Oct. 2024. arXiv:2410.18982.
- [120] Y. Qu, T. Zhang, N. Garg, and A. Kumar. Recursive Introspection: Teaching Language Model Agents How to Self-Improve, July 2024. arXiv:2407.18219.
- [121] R. Rafailov, A. Sharma, E. Mitchell, C. D. Manning, S. Ermon, and C. Finn. Direct Preference Optimization: Your Language Model is Secretly a Reward Model. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Proceedings of the Thirty-seventh Annual Conference on Neural Information Processing Systems (NeurIPS '23)*, volume 36 of *Advances in Neural Information Processing Systems*, pages 53728–53741, New Orleans, LA, USA, Dec. 2023. Curran Associates.
- [122] D. Rein, B. L. Hou, A. C. Stickland, J. Petty, R. Y. Pang, J. Dirani, J. Michael, and S. R. Bowman. GPQA: A Graduate-Level Google-Proof Q&A Benchmark, Nov. 2023. arXiv:2311.12022.
- [123] C. D. Rosin. Multi-Armed Bandits with Episode Context. *Annals of Mathematics and Artificial Intelligence*, 61(3):203–230, Mar. 2011.
- [124] S. Roy and D. Roth. Solving General Arithmetic Word Problems. In L. Márquez, C. Callison-Burch, and J. Su, editors, *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP '15*, pages 1743–1752, Lisbon, Portugal, Sept. 2015. Association for Computational Linguistics.
- [125] K. Sakaguchi, R. Le Bras, C. Bhagavatula, and Y. Choi. Winogrande: An Adversarial Winograd Schema Challenge at Scale. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8732–8740, Apr. 2020.
- [126] M. Sap, H. Rashkin, D. Chen, R. Le Bras, and Y. Choi. Social IQa: Commonsense Reasoning about Social Interactions. In K. Inui, J. Jiang, V. Ng, and X. Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP '19*, pages 4463–4473, Hong Kong, China, Nov. 2019. Association for Computational Linguistics.
- [127] A. Saparov and H. He. Language Models Are Greedy Reasoners: A Systematic Formal Analysis of Chain-of-Thought. In *Proceedings of the Eleventh International Conference on Learning Representations, ICLR '23*, Kigali, Rwanda, May 2023.
- [128] W. Saunders, C. Yeh, J. Wu, S. Bills, L. Ouyang, J. Ward, and J. Leike. Self-Critiquing Models for Assisting Human Evaluators, June 2022. arXiv:2206.05802.
- [129] T. Sawada, D. Paleka, A. Havrilla, P. Tadepalli, P. Vidas, A. Krnias, J. J. Nay, K. Gupta, and A. Komatsuzaki. ARB: Advanced Reasoning Benchmark for Large Language Models, July 2023. arXiv:2307.13692.
- [130] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver. Mastering Atari, Go, Chess and Shogi by Planning With a Learned Model. *Nature*, 588:604–609, Dec. 2020.
- [131] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal Policy Optimization Algorithms, Aug. 2017. arXiv:1707.06347.
- [132] M. Shridhar, X. Yuan, M.-A. Côté, Y. Bisk, A. Trischler, and M. Hausknecht. ALFWorld: Aligning Text and Embodied Environments for Interactive Learning. In *Proceedings of the International Conference on Learning Representations, ICLR '21*, Virtual Event, May 2021.
- [133] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and D. Hassabis. Mastering the Game of Go With Deep Neural Networks and Tree Search. *Nature*, 529:484–489, Jan. 2016.
- [134] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, T. Lillicrap, K. Simonyan, , and D. Hassabis. A General Reinforcement Learning Algorithm that Masters Chess, Shogi, and Go Through Self-Play. *Science*, 362(6419):1140–1144, Dec. 2018.
- [135] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis. Mastering the Game of Go without Human Knowledge. *Nature*, 550:354–359, Oct. 2017.
- [136] K. Sinha, S. Sodhani, J. Dong, J. Pineau, and W. L. Hamilton. CLUTRR: A Diagnostic Benchmark for Inductive Reasoning from Text. In K. Inui, J. Jiang, V. Ng, and X. Wan, editors, *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP '19*, pages 4506–4515, Hong Kong, China, Nov. 2019. Association for Computational Linguistics.
- [137] C. Snell, J. Lee, K. Xu, and A. Kumar. Scaling LLM Test-Time Compute Optimally Can be More Effective than Scaling Model Parameters, Aug. 2024. arXiv:2408.03314.
- [138] A. Srivastava, A. Rastogi, A. Rao, A. A. M. Shobh, A. Abid, A. Fisch, A. R. Brown, A. Santoro, A. Gupta, A. Garriga-Alonso, et al. Beyond the Imitation Game: Quantifying and Extrapolating the Capabilities of Language Models, June 2023. arXiv:2206.04615.
- [139] S. Srivastava, A. M. B, A. P. V, S. Menon, A. Sukumar, A. S. T, A. Philipose, S. Prince, and S. Thomas. Functional Benchmarks for Robust Evaluation of Reasoning Performance, and the Reasoning Gap, Feb. 2024. arXiv:2402.19450.
- [140] N. Stiennon, L. Ouyang, J. Wu, D. Ziegler, R. Lowe, C. Voss, A. Radford, D. Amodei, and P. F. Christiano. Learning to Summarize with Human Feedback. In H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, editors, *Proceedings of the Thirty-fourth Annual Conference on Neural Information Processing Systems (NeurIPS '20)*, volume 33 of *Advances in Neural Information Processing Systems*, pages 3008–3021, Virtual Event, Dec. 2020. Curran Associates.
- [141] J. Sun, C. Zheng, E. Xie, Z. Liu, R. Chu, J. Qiu, J. Xu, M. Ding, H. Li, M. Geng, et al. A Survey of Reasoning with Foundation Models, Jan. 2024. arXiv:2312.11562.
- [142] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2015.
- [143] R. Tadeusiewicz and L. Ogiela. Modern Methods for the Cognitive Analysis of Economic Data and Text Documents and Their Application in Enterprise Management. In V. Snášel, A. Abraham, K. Saeed, and J. Pokorný, editors, *Proceedings of the 7th Computer Information Systems and Industrial Management Applications, CISIM '08*, pages 11–23, Ostrava, Czech Republic, June 2008. IEEE Press.
- [144] R. Tadeusiewicz, L. Ogiela, and M. R. Ogiela. Cognitive Analysis Techniques in Business Planning and Decision Support Systems. In L. Rutkowski, R. Tadeusiewicz, L. A. Zadeh, and J. M. Żurada, editors, *Proceedings of the 8th International Conference on Artificial Intelligence and Soft Computing (ICAISC '06)*, volume 4029 of

- Lecture Notes in Computer Science*, pages 1027–1039, Zakopane, Poland, June 2006. Springer.
- [145] O. Tafjord, B. Dalvi, and P. Clark. ProofWriter: Generating Implications, Proofs, and Abductive Statements over Natural Language. In C. Zong, F. Xia, W. Li, and R. Navigli, editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3621–3634, Virtual Event, Aug. 2021. Association for Computational Linguistics.
- [146] A. Talmor, J. Herzig, N. Lourie, and J. Berant. CommonsenseQA: A Question Answering Challenge Targeting Commonsense Knowledge. In J. Burstein, C. Doran, and T. 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)*, NAACL '19, pages 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.
- [147] Z. Tang, X. Zhang, B. Wang, and F. Wei. MathScale: Scaling Instruction Tuning for Mathematical Reasoning, Mar. 2024. arXiv:2403.02884.
- [148] Q. Team. QwQ: Reflect Deeply on the Boundaries of the Unknown. <https://qwenlm.github.io/blog/qwq-32b-preview/>, Nov. 2024. accessed 2025-01-01.
- [149] Y. Tian, B. Peng, L. Song, L. Jin, D. Yu, L. Han, H. Mi, and D. Yu. Toward Self-Improvement of LLMs via Imagination, Searching, and Criticizing. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS '24)*, volume 37 of *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec. 2024. Curran Associates.
- [150] R. Tu, K. Zhang, B. Bertilson, H. Kjellstrom, and C. Zhang. Neuropathic Pain Diagnosis Simulator for Causal Discovery Algorithm Evaluation. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Proceedings of the Thirty-third Annual Conference on Neural Information Processing Systems (NeurIPS '19)*, volume 32 of *Advances in Neural Information Processing Systems*, pages 12793–12804, Vancouver, Canada, Dec. 2019. Curran Associates.
- [151] J. Uesato, N. Kushman, R. Kumar, F. Song, N. Siegel, L. Wang, A. Creswell, G. Irving, and I. Higgins. Solving Math Word Problems with Process-and Outcome-Based Feedback, Nov. 2022. arXiv:2211.14275.
- [152] A. Vijayakumar, M. Cogswell, R. Selvaraju, Q. Sun, S. Lee, D. Crandall, and D. Batra. Diverse Beam Search for Improved Description of Complex Scenes. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1):7371–7379, Apr. 2018.
- [153] J. Wang, M. Fang, Z. Wan, M. Wen, J. Zhu, A. Liu, Z. Gong, Y. Song, L. Chen, L. M. Ni, L. Yang, Y. Wen, and W. Zhang. OpenR: An Open Source Framework for Advanced Reasoning with Large Language Models, Oct. 2024. arXiv:2410.09671.
- [154] K. Wang, H. Ren, A. Zhou, Z. Lu, S. Luo, W. Shi, R. Zhang, L. Song, M. Zhan, and H. Li. MathCoder: Seamless Code Integration in LLMs for Enhanced Mathematical Reasoning, Oct. 2023. arXiv:2310.03731.
- [155] P. Wang, L. Li, Z. Shao, R. Xu, D. Dai, Y. Li, D. Chen, Y. Wu, and Z. Sui. Math-Shepherd: Verify and Reinforce LLMs Step-by-Step without Human Annotations. In L.-W. Ku, A. Martins, and V. Srikumar, editors, *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL '24, pages 9426–9439, Bangkok, Thailand, Aug. 2024. Association for Computational Linguistics.
- [156] X. Wang, Z. Hu, P. Lu, Y. Zhu, J. Zhang, S. Subramaniam, A. Loomba, S. Zhang, Y. Sun, and W. Wang. SCIBENCH: Evaluating College-Level Scientific Problem-Solving Abilities of Large Language Models. In *Proceedings of the 3rd Workshop on Mathematical Reasoning and AI, MATH-AI '23*, New Orleans, LA, USA, Dec. 2023.
- [157] X. Wang, L. Song, Y. Tian, D. Yu, B. Peng, H. Mi, F. Huang, and D. Yu. Towards Self-Improvement of LLMs via MCTS: Leveraging Stepwise Knowledge with Curriculum Preference Learning, Oct. 2024. arXiv:2410.06508.
- [158] X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou. Self-Consistency Improves Chain of Thought Reasoning in Language Models. In *Proceedings of the Eleventh International Conference on Learning Representations, ICLR '23*, Kigali, Rwanda, May 2023.
- [159] Z. Wang, S. Zhou, D. Fried, and G. Neubig. Execution-Based Evaluation for Open-Domain Code Generation, May 2023. arXiv:2212.10481.
- [160] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. Chi, Q. V. Le, and D. Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Proceedings of the Thirty-sixth Annual Conference on Neural Information Processing Systems (NeurIPS '22)*, volume 35 of *Advances in Neural Information Processing Systems*, pages 24824–24837, New Orleans, LA, USA, Dec. 2022. Curran Associates.
- [161] P. Wiesner, I. Behnke, D. Scheinert, K. Gontarska, and L. Thamsen. Let's Wait Awhile: How Temporal Workload Shifting Can Reduce Carbon Emissions in the Cloud. In *Proceedings of the 22nd International Middleware Conference*, Middleware '21, pages 260–272, Virtual Event, Dec. 2021. Association for Computing Machinery.
- [162] Z. Xi, Y. Ding, W. Chen, B. Hong, H. Guo, J. Wang, D. Yang, C. Liao, X. Guo, W. He, S. Gao, L. Chen, R. Zheng, Y. Zou, T. Gui, Q. Zhang, X. Qiu, X. Huang, Z. Wu, and Y.-G. Jiang. AgentGym: Evolving Large Language Model-based Agents across Diverse Environments, June 2024. arXiv:2406.04151.
- [163] Y. Xie, A. Goyal, W. Zheng, M.-Y. Kan, T. P. Lillicrap, K. Kawaguchi, and M. Shieh. Monte Carlo Tree Search Boosts Reasoning via Iterative Preference Learning, June 2024. arXiv:2405.00451.
- [164] G. Xiong, Q. Jin, Z. Lu, and A. Zhang. Benchmarking Retrieval-Augmented Generation for Medicine, Feb. 2024. arXiv:2402.13178.
- [165] J. Xiong, J. Shen, Y. Yuan, H. Wang, Y. Yin, Z. Liu, L. Li, Z. Guo, Q. Cao, Y. Huang, C. Zheng, X. Liang, M. Zhang, and Q. Liu. TRIGO: Benchmarking Formal Mathematical Proof Reduction for Generative Language Models. In H. Bouamor, J. Pino, and K. Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP '23*, pages 11594–11632, Singapore, Dec. 2023. Association for Computational Linguistics.
- [166] Y. Yan, J. Su, J. He, F. Fu, X. Zheng, Y. Lyu, K. Wang, S. Wang, Q. Wen, and X. Hu. A Survey of Mathematical Reasoning in the Era of Multimodal Large Language Model: Benchmark, Method & Challenges, Dec. 2024. arXiv:2412.11936.
- [167] K. Yang, A. Swope, A. Gu, R. Chalamala, P. Song, S. Yu, S. Godil, R. J. Prenger, and A. Anandkumar. LeanDojo: Theorem Proving with Retrieval-Augmented Language Models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Proceedings of the Thirty-seventh Annual Conference on Neural Information Processing Systems (NeurIPS '23)*, volume 36 of *Advances in Neural Information Processing Systems*, pages 21573–21612, New Orleans, LA, USA, Dec. 2023. Curran Associates.
- [168] S. Yao, H. Chen, J. Yang, and K. Narasimhan. WebShop: Towards Scalable Real-World Web Interaction with Grounded Language Agents. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Proceedings of the Thirty-sixth Annual Conference on Neural Information Processing Systems (NeurIPS '22)*, volume 35 of *Advances in Neural Information Processing Systems*, pages 20744–20757, New Orleans, LA, USA, Dec. 2022. Curran Associates.
- [169] S. Yao, D. Yu, J. Zhao, I. Shafran, T. Griffiths, Y. Cao, and K. Narasimhan. Tree of Thoughts: Deliberate Problem Solving with Large Language Models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Proceedings of the Thirty-seventh Annual Conference on Neural Information Processing Systems (NeurIPS '23)*, volume 36 of *Advances in Neural Information Processing Systems*, pages 11809–11822, New Orleans, LA, USA, Dec. 2023. Curran Associates.
- [170] N. Young, Q. Bao, J. Bensemann, and M. Witbrock. Abduction-Rules: Training Transformers to Explain Unexpected Inputs. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Findings of the Association for Computational Linguistics: ACL 2022*, pages 218–227, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [171] L. Yuan, W. Li, H. Chen, G. Cui, N. Ding, K. Zhang, B. Zhou, Z. Liu, and H. Peng. Free Process Rewards without Process Labels, Dec. 2024. arXiv:2412.01981.
- [172] Z. Yuan, H. Yuan, C. Tan, W. Wang, and S. Huang. How Well Do Large Language Models Perform in Arithmetic Tasks?, Mar. 2023. arXiv:2304.02015.
- [173] R. Zellers, Y. Bisk, R. Schwartz, and Y. Choi. SWAG: A Large-Scale Adversarial Dataset for Grounded Commonsense Inference. In E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, editors, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language*

- Processing*, EMNLP '18, pages 93–104, Brussels, Belgium, Nov. 2018. Association for Computational Linguistics.
- [174] R. Zellers, A. Holtzman, Y. Bisk, A. Farhadi, and Y. Choi. HellaSwag: Can a Machine Really Finish Your Sentence? In A. Korhonen, D. Traum, and L. Márquez, editors, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, ACL '19, pages 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics.
- [175] Z. Zeng, Q. Cheng, Z. Yin, B. Wang, S. Li, Y. Zhou, Q. Guo, X. Huang, and X. Qiu. Scaling of Search and Learning: A Roadmap to Reproduce o1 from Reinforcement Learning Perspective, Dec. 2024. arXiv:2412.14135.
- [176] D. Zhang, X. Huang, D. Zhou, Y. Li, and W. Ouyang. Accessing GPT-4 Level Mathematical Olympiad Solutions via Monte Carlo Tree Self-Refine with LLaMa-3 8B, June 2024. arXiv:2406.07394.
- [177] D. Zhang, J. Wu, J. Lei, T. Che, J. Li, T. Xie, X. Huang, S. Zhang, M. Pavone, Y. Li, W. Ouyang, and D. Zhou. LLaMA-Berry: Pairwise Optimization for O1-like Olympiad-Level Mathematical Reasoning, Nov. 2024. arXiv:2410.02884.
- [178] D. Zhang, S. Zhoubian, Z. Hu, Y. Yue, Y. Dong, and J. Tang. ReST-MCTS\*: LLM Self-Training via Process Reward Guided Tree Search. In *Proceedings of the Thirty-eighth Annual Conference on Neural Information Processing Systems (NeurIPS '24)*, volume 37 of *Advances in Neural Information Processing Systems*, Vancouver, Canada, Dec. 2024. Curran Associates.
- [179] L. Zhang, A. Hosseini, H. Bansal, M. Kazemi, A. Kumar, and R. Agarwal. Generative Verifiers: Reward Modeling as Next-Token Prediction, Oct. 2024. arXiv:2408.15240.
- [180] M. Zhao, S. Pan, N. Agarwal, Z. Wen, D. Xu, A. Natarajan, P. Kumar, S. S. P. R. Tijoriwala, K. Asher, H. Wu, A. Basant, D. Ford, D. David, N. Yigitbasi, P. Singh, C.-J. Wu, and C. Kozyrakis. Tectonic-Shift: A Composite Storage Fabric for Large-Scale ML Training. In *Proceedings of the USENIX Annual Technical Conference*, ATC '23, pages 433–449, Boston, MA, USA, July 2023. USENIX Association.
- [181] Y. Zhao, Y. Li, C. Li, and R. Zhang. MultiHiertt: Numerical Reasoning over Multi Hierarchical Tabular and Textual Data. In S. Muresan, P. Nakov, and A. Villavicencio, editors, *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL '22, pages 6588–6600, Dublin, Ireland, May 2022. Association for Computational Linguistics.
- [182] Y. Zhao, H. Yin, B. Zeng, H. Wang, T. Shi, C. Lyu, L. Wang, W. Luo, and K. Zhang. Marco-o1: Towards Open Reasoning Models for Open-Ended Solutions, Nov. 2024. arXiv:2411.14405.
- [183] K. Zheng, J. M. Han, and S. Polu. miniF2F: A Cross-System Benchmark for Formal Olympiad-Level Mathematics. In *Proceedings of the Tenth International Conference on Learning Representations*, ICLR '22, Virtual Event, Apr. 2022.
- [184] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, H. Zhang, J. E. Gonzalez, and I. Stoica. Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Proceedings of the Thirty-seventh Annual Conference on Neural Information Processing Systems (NeurIPS '23)*, volume 36 of *Advances in Neural Information Processing Systems*, pages 46595–46623, New Orleans, LA, USA, Dec. 2023. Curran Associates.
- [185] S. Zhou, F. F. Xu, H. Zhu, X. Zhou, R. Lo, A. Sridhar, X. Cheng, T. Ou, Y. Bisk, D. Fried, U. Alon, and G. Neubig. WebArena: A Realistic Web Environment for Building Autonomous Agents, Apr. 2024. arXiv:2307.13854.
- [186] D.-H. Zhu, Y.-J. Xiong, J.-C. Zhang, X.-J. Xie, and C.-M. Xia. Understanding Before Reasoning: Enhancing Chain-of-Thought with Iterative Summarization Pre-Prompting, Jan. 2025. arXiv:2501.04341.