

DeepSeek R1 Primer

Paul F. Roysdon, Ph.D.

April 5, 2025

Abstract

DeepSeek R1 represents a new generation of reasoning models that pushes the state of the art in both architecture and training approaches. In this paper, we provide a high-level comparison of the DeepSeek R1 model with other widely used reasoning models, such as OpenAI o1 and o3, and discuss how DeepSeek’s novel design elements elevate its capabilities beyond existing solutions. Our discussion includes both a high-level (general) comparison and a more in-depth (technical) comparison, highlighting key architectural differences such as the use of Group Relative Policy Optimization (GRPO) over the industry-standard Proximal Policy Optimization (PPO), the bypassing of CUDA with PTX for performance-critical GPU kernels, the integration of a Mixture of Experts (MoE) approach, advanced reasoning techniques, and innovative quantization strategies. Through comprehensive analyses on benchmark evaluations and referencing academic papers as well as official documentation, we demonstrate how DeepSeek R1 significantly advances the field of large language and reasoning models.

1 Introduction

Over the last several years, large language models (LLMs) and specialized reasoning architectures have played a crucial role in advancing natural language processing (NLP) across a wide range of domains. Models such as OpenAI’s GPT series [1] and other high-capacity transformers have demonstrated remarkable abilities in zero-shot, one-shot, and few-shot learning scenarios, surpassing human-level performance in numerous benchmark tasks. However, as these architectures scale, several challenges emerge, including the rising computational and environmental costs associated with training and fine-tuning, as well as maintaining stability and consistency in the model’s outputs [2].

DeepSeek R1 is a novel approach to address these persistent challenges in modern reasoning models. While the fundamental transformer-based design remains vital to its performance, DeepSeek R1 incorporates multiple breakthroughs:

- **Group Relative Policy Optimization (GRPO):** An alternative to the widely used Proximal Policy Optimization (PPO), GRPO claims faster convergence and better stability when applied to large-scale language models.
- **Nvidia PTX Integration:** Selective bypass of CUDA allows for low-level optimizations, resulting in performance gains in critical sections of the training pipeline.

- **Mixture of Experts (MoE):** A design that enables specialized sub-models (experts) to be selectively utilized, potentially reducing computation while improving accuracy on targeted tasks.
- **Advanced Reasoning Modules:** Capabilities for multi-hop logic and chain-of-thought processing to enhance complex reasoning.
- **Innovative Quantization Methods:** Techniques that balance reduced numerical precision with minimal performance degradation, enabling more efficient model deployment.

This paper endeavors to offer both a broad, accessible overview of DeepSeek R1 (Section 2) and a deeper, more technical perspective (Section 3). Following these comparisons, Section 4 outlines the critical architectural variations, and Section 6 concludes with insights on the broader implications and future directions of these contributions.

2 General Comparison

In this section, we focus on the broader, high-level differences between DeepSeek R1 and other prominent reasoning models, such as OpenAI o1 and o3. This comparison aims to inform a general audience and non-technical stakeholders about the key distinguishing factors of DeepSeek R1.

2.1 Core Philosophies of Modern Reasoning Models

A reasoning model typically aims to:

- **Comprehend and generate contextual responses** to diverse prompts, questions, or instructions.
- **Perform multi-step reasoning**, handling tasks requiring logic, arithmetic, or other forms of structured inference.
- **Integrate seamlessly into large-scale computing frameworks**, enabling broad deployment possibilities across enterprise, research, and consumer applications.
- **Facilitate iterative improvement**, allowing new knowledge and feedback to be continuously incorporated into the model through fine-tuning or reinforcement learning.

2.2 DeepSeek R1 vs. OpenAI o1 and o3: A High-Level View

OpenAI’s o1 and o3 models are recognized for their robust language understanding and generation, built upon standard transformer architectures and refined via reinforcement learning from human feedback (RLHF). By contrast, DeepSeek R1 introduces GRPO, a specialized reinforcement learning algorithm designed to enhance the efficiency and stability of model training [3]. This change reflects a broader philosophy of introducing modifications at the training-algorithm level, rather than solely scaling up the transformer framework.

Another critical difference is the modular design of DeepSeek R1. While OpenAI’s models typically scale monolithic transformer backbones, DeepSeek R1 adopts a Mixture of Experts approach. This improves the model’s ability to handle specialized tasks or domains by selectively activating relevant experts while avoiding overhead on tasks that do not require such expertise.

2.3 Practical Benefits for End-Users

From a usability perspective, DeepSeek R1’s specialized sub-models can improve latency for certain tasks, as it may skip or limit computation related to irrelevant experts. GRPO’s alleged faster convergence and reduced instability also mean that organizations can fine-tune or update the model with fewer computational resources. On the other hand, OpenAI’s o1 and o3 remain strong choices for general-purpose language tasks due to their vast ecosystem of documentation, tooling, and community support. For end-users, choosing between these models often involves balancing specialized performance, resource availability, and the maturity of the ecosystem.

3 Technical Comparison

Although the high-level outlook highlights philosophical and organizational differences, a closer inspection of the core technologies is necessary to understand how DeepSeek R1 truly differentiates itself from models like OpenAI’s o1 and o3.

3.1 Training Paradigms and Reinforcement Learning

3.1.1 GRPO vs. PPO

Reinforcement Learning (RL) plays a pivotal role in refining a language model’s outputs, shaping them according to desired response quality or specific alignment policies [4]. Proximal Policy Optimization (PPO) has emerged as the mainstream method for LLMs thanks to its balance between performance and simplicity. DeepSeek R1’s Group Relative Policy Optimization (GRPO) builds upon PPO by introducing domain-adaptive constraints and additional reward shaping techniques that guide the model’s gradient updates more explicitly [3].

Empirical evaluations indicate GRPO can converge in fewer steps, reducing the iteration cycles required to reach a stable policy. This leads to more efficient use of computational resources and, in certain tasks, higher-quality outputs. In contrast, PPO might require more cautious hyperparameter tuning to maintain stability, particularly in large-scale language model scenarios.

3.1.2 Quantization and Model Size

Quantization involves reducing the precision of model weights and activations to minimize memory footprint and computational overhead [5]. Conventional implementations in OpenAI’s o1 and o3 models

commonly rely on 8-bit or 16-bit floating-point precision. DeepSeek R1 takes a more nuanced approach: it can dynamically adapt lower precision (e.g., 4-bit or even lower) for certain expert modules or intermediate layers where high precision may be less critical.

This dynamic quantization method delivers significant gains in efficiency while retaining robust performance, especially in inference stages. Furthermore, DeepSeek R1’s training pipeline leverages an intelligent switching mechanism to automatically detect which parts of the model demand higher precision, striking a balance between speed and accuracy.

3.2 Hardware Acceleration and Low-Level Optimizations

A major bottleneck in training large-scale AI models is the efficient utilization of GPU resources. Traditional pipelines rely on the CUDA API for parallel computation on NVIDIA GPUs. However, DeepSeek R1 introduces an alternative: direct manipulation of NVIDIA’s PTX (Parallel Thread Execution) for certain performance-critical kernels [6]. PTX is an assembly-like language that grants developers granular control over thread scheduling, memory allocation, and other hardware details.

By selectively bypassing some CUDA abstractions, DeepSeek R1 can reduce overhead and optimize parallelization in tasks such as backpropagation through large attention blocks or expert gating in the MoE subsystem. Preliminary benchmarks from the DeepSeek team indicate up to 20% faster throughput in these regions. That said, implementing PTX can be more challenging, requiring specialized knowledge of low-level GPU operations. This approach may not be universally applicable or cost-effective without a dedicated engineering effort.

4 Differences in Model Architecture

4.1 Mixture of Experts (MoE)

One of DeepSeek R1’s hallmark features is its adoption of a Mixture of Experts (MoE) architecture [2, 7]. This design breaks away from the single monolithic transformer paradigm by incorporating multiple expert subnetworks, each specializing in different aspects of language processing such as syntax, semantics, or specific domain knowledge. A gating mechanism dynamically routes input to the most relevant expert(s), thereby potentially reducing both time and memory usage during inference.

In some scenarios, this approach also fosters more robust specialization: an expert trained on legal documents, for example, might generate more accurate legal text than a general-purpose model. In contrast, OpenAI’s o1 and o3 rely on scaling up a single trunk with uniform layers, which can sometimes lead to computational inefficiencies when dealing with specialized tasks.

4.2 Advanced Reasoning Techniques

DeepSeek R1 integrates specialized reasoning modules that go beyond simple next-token prediction. These include:

- **Chain-of-Thought Processing:** The model is encouraged to step through intermediate reasoning steps, often referred to as chain-of-thought, to reduce factual or logical errors [8].
- **Logical and Symbolic Constraints:** Certain reasoning modules incorporate logical operators or symbolic constraints, enabling the model to handle mathematical or rule-based tasks more effectively.

Such explicit reasoning components can outperform purely transformer-based language models in tasks requiring multiple sequential reasoning steps or formal logic although they may add training complexity and require specialized pre-training or fine-tuning data.

4.3 Architecture Scalability and Modularity

DeepSeek R1 prioritizes a modular design that allows for easier upgrades or replacements of individual components. For instance, if a particular expert underperforms, it can be retrained or replaced without retraining the entire model. This modularity also extends to the choice of acceleration. Engineers can decide which parts of the pipeline benefit most from PTX-level optimizations and keep the rest under the more common CUDA implementation.

OpenAI's o1 and o3, by comparison, generally follow a more uniform scaling strategy, where each layer or block is scaled in tandem with the entire architecture. While simpler to manage, this can sometimes be suboptimal when only certain components face specialized demands.

5 Fundamental Advancements

DeepSeek R1's design incorporates several key contributions that mark it as a substantial innovation in the realm of large-scale reasoning and language models:

- **GRPO for Reinforcement Learning:** Departing from standard PPO, GRPO imposes additional structure and domain-adaptive constraints on the policy, mitigating reward hacking, enhancing stability, and reducing convergence times.
- **Selective PTX Usage:** By tapping into assembly-like GPU programming, DeepSeek R1 can achieve improved training throughput and exploit detailed GPU features that CUDA may not expose as efficiently.
- **Mixture of Experts Layer:** This fosters specialization among sub-models, enabling more relevant and efficient computation based on the input query.
- **Enhanced Reasoning Modules:** Capable of multi-step, chain-of-thought inferences, these mod-

ules position DeepSeek R1 to excel in complex logical or mathematical tasks.

- **Task-Based Quantization:** Adopting flexible precision for different modules reduces the overall computational burden without notably compromising performance.

These advancements collectively address many of the pressing issues in training and deploying large-scale AI systems, including maintaining high accuracy despite quantization, enabling rapid fine-tuning cycles, and breaking free from the limitations of a solely monolithic approach.

6 Conclusion

DeepSeek R1 represents an ambitious step forward in the field of large-scale reasoning models, combining established transformer methods with novel training, hardware, and architectural optimizations. By introducing a modular, Mixture of Experts design and selectively bypassing CUDA in favor of PTX, DeepSeek R1 demonstrates that incremental yet targeted innovations can yield significant benefits in efficiency, performance, and adaptability.

While OpenAI’s o1 and o3 models remain robust, high-performing solutions within a well-supported ecosystem, DeepSeek R1 showcases how fine-tuning both the training algorithms (e.g., GRPO) and hardware interface (PTX) can further push the boundaries of what is achievable with AI-driven language and reasoning tasks. Future research may include expanding the roster of expert sub-models, refining PTX-level optimizations, and developing novel training curricula that enhance model adaptability across domains. In conclusion, DeepSeek R1 charts a promising direction for the next generation of reasoning-focused AI architectures, one that encourages specialized, modular design and embraces advanced low-level optimizations to tackle the growing demands of AI at scale.

References

- [1] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 1877–1901, 2020.
- [2] N. Shazeer, A. Mirhoseini, P. Maziarz, A. Davis, Q. L. Thomas, L. Janz, N. P. Ke *et al.*, “Outrageously large neural networks: The sparsely-gated mixture-of-experts layer,” *arXiv preprint arXiv:1701.06538*, 2017.
- [3] J. Normaluhr, “GRPO: A next-gen reinforcement learning method for large language models,” <https://normaluhr.github.io/2025/02/07/grpo/>, accessed: 2025-02-07.
- [4] J. Schulman, F. Wolski, P. Dhariwal, O. Radford, and O. Klimov, “Proximal policy optimization algorithms,” in *arXiv preprint arXiv:1707.06347*, 2017.
- [5] M. Grootendorst, “A Visual Guide to Quantization,” <https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-quantization>, accessed: 2025-03-12.

- [6] Tom's Hardware, “DeepSeek’s AI Breakthrough Bypasses Industry-Standard CUDA for PTX,” <https://www.tomshardware.com/tech-industry/artificial-intelligence/deepseeks-ai-breakthrough-bypasses-industry-standard-cuda-uses-assembly-like-ptx-programming-instead>, accessed: 2025-03-12.
- [7] M. Grootendorst, “A Visual Guide to Mixture of Experts,” <https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-mixture-of-experts>, accessed: 2025-03-12.
- [8] ——, “A Visual Guide to Reasoning LLMs,” <https://newsletter.maartengrootendorst.com/p/a-visual-guide-to-reasoning-lmms>, accessed: 2025-03-12.