IARJSET



International Advanced Research Journal in Science, Engineering and Technology

Impact Factor 8.311 ⅔ Peer-reviewed & Refereed journal ⅔ Vol. 12, Issue 5, May 2025

DOI: 10.17148/IARJSET.2025.125378

Comparative Benchmark Analysis of ChatGPT and DeepSeek: Performance Across AI Tasks

Dr. Anju Kaushik^{1,} Dr. Anil Kaushik²

Assistant Professor, Computer science, G.C.W. Gohana, Haryana, INDIA¹

Assistant Professor, Physics, G.C.W. Gohana, Haryana, INDIA²

Abstract: ChatGPT and DeepSeek represent two prominent large language models, each offering unique strengths in artificial intelligence applications. ChatGPT is widely known for its advanced conversational abilities and broad language understanding, while DeepSeek is recognized for its strong performance in computational and technical domains. This research paper presents a comparative evaluation of DeepSeek-R1 and ChatGPT across several prominent mathematical and algorithmic benchmarks. The analysis reveals that both models exhibit strong and competitive performance, with each demonstrating unique strengths depending on the benchmark. DeepSeek-R1 shows a slight advantage in advanced mathematical problem-solving, while ChatGPT excels in competitive programming and complex quantitative reasoning tasks. Although the overall performance of the two models is closely matched, notable differences emerge in specific areas, highlighting the importance of selecting the appropriate model based on the requirements of the task. These findings offer valuable insights for researchers and practitioners seeking to deploy large language models in mathematical and computational domains.

Keywords: ChatGPT, DeepSeek-R1, Large Language Models, Mathematical Reasoning, Benchmark Comparison, AI Performance Evaluation

I. INTRODUCTION

Large language models (LLMs) have become foundational in advancing the field of artificial intelligence, enabling sophisticated performance across a wide spectrum of tasks such as natural language understanding, code generation, and multi-modal reasoning (Brown et al., 2020; OpenAI, 2023). Among these, ChatGPT, developed by OpenAI, has established itself as a leading conversational agent, renowned for its versatility in language comprehension and generation (OpenAI, 2023; Nori et al., 2023). In parallel, newer models like DeepSeek have emerged, designed to address both computational efficiency and domain-specific performance, particularly in technical and mathematical domains (Wang et al., 2024).

The increasing deployment of LLMs in diverse application areas has highlighted the necessity for systematic benchmarking to evaluate and compare their capabilities across different domains (Liang et al., 2022; Srivastava et al., 2022). Benchmarking studies not only facilitate informed model selection but also expose strengths and limitations that guide future research and development (Zheng et al., 2023). Despite the rapid progress in model architectures and training strategies, direct head-to-head comparisons of models like ChatGPT and DeepSeek, especially across domains such as mathematics, programming and reasoning remain limited in the literature.

II. LITERATURE REVIEW

The rapid evolution of large language models has significantly advanced the field of mathematical and algorithmic reasoning. Brown et al. (2020) established the effectiveness of large-scale transformer models, such as GPT-3, in handling a variety of complex tasks. Building on this foundation, Ouyang et al. (2022) demonstrated that reinforcement learning from human feedback further enhances the performance and reliability of models like ChatGPT, enabling them to tackle challenging mathematical and programming problems. Hendrycks et al. (2021) introduced the MATH dataset, revealing persistent challenges for LLMs in multi-step mathematical reasoning, while Cobbe et al. (2021) highlighted the incremental improvements on competitive benchmarks such as AIME and MATH.

Recent innovations have focused on addressing these limitations. Wang et al. (2024) presented DeepSeek-R1, which uses a mixture-of-experts architecture to improve performance on structured mathematical tasks, as reflected in its competitive results on the AIME-2024 and MATH-500 benchmarks. In the realm of algorithmic problem-solving, Chen et al. (2021) showed that models trained on code repositories, like Codex, outperform more general LLMs on programming benchmarks such as Codeforces. Li et al. (2023) further explored the gap between human and AI performance in competitive programming, emphasizing the need for creative reasoning in top-tier tasks.

IARJSET



International Advanced Research Journal in Science, Engineering and Technology

Impact Factor 8.311 $\,st\,$ Peer-reviewed & Refereed journal $\,st\,$ Vol. 12, Issue 5, May 2025

DOI: 10.17148/IARJSET.2025.125378

For general-purpose quantitative reasoning, Zhang et al. (2023) introduced the GPQA benchmark, which has become a standard for evaluating graduate-level problem-solving in LLMs. Eduminds Learning (2025) provided a direct comparison between DeepSeek-R1 and ChatGPT, finding that while both models excel in mathematics, ChatGPT demonstrates a notable advantage in complex quantitative reasoning, as seen in higher GPQA Diamond scores. Collectively, these studies highlight that while DeepSeek-R1 and ChatGPT are both state-of-the-art, their relative strengths vary by benchmark, underscoring the importance of task-specific evaluation and model refinement.

III. KEY COMPARISON BETWEEN CHATGPT AND DEEPSEEK

The table -1 provides a concise overview of the fundamental differences between ChatGPT and DeepSeek, two prominent large language models (LLMs) currently shaping the landscape of artificial intelligence applications. Each feature highlights unique aspects of their development, architecture, and intended use cases.

Features	Developer	Model type	Architecture	Training data	Logo
Chat Gpt	Open Al	Proprietary LLM	Transformer- based (GPT-4)	Extensive multilingual, strong in English	G
Deepseek	Deekseek AI	Open-source LLM	Mixture-of- Experts (MoE)	Multilingual, Chinese- focused	deepseek

Table-1. Key Comparison Derween Charder 1 and DeepSeek	Table-1	: key C	Comparison	Between	ChatGPT	and DeepSeek
--	---------	---------	------------	---------	---------	--------------

IV. BENCHMARK PERFORMANCE ACROSS CORE AI TASKS

Four benchmarks have been considered for this study as described below:

AIME-2024: The AIME-2024 benchmark uses questions from the 2024 American Invitational Mathematics Examination, a prestigious math contest for top high school students. It features 15 advanced problems requiring integer answers, challenging both human and AI participants. This benchmark is widely used to assess mathematical reasoning in large language models.

Codeforces: Codeforces is an online platform that hosts competitive programming contests, providing a diverse set of algorithmic and coding problems. As a benchmark, it evaluates an AI model's ability to solve real-world programming challenges efficiently and accurately. High performance on Codeforces indicates strong coding and logical reasoning skills in AI systems.

GPQA Diamond: GPQA Diamond is a graduate-level question-answering benchmark with particularly challenging, expert-crafted questions in science domains. Designed to be resistant to simple web searches, it tests an AI's deep subject understanding and reasoning abilities. Success on this benchmark demonstrates advanced, specialized knowledge in language models.

MATH-500: MATH-500 consists of 500 carefully selected math problems covering topics like algebra, geometry, and probability. It is used to measure an AI model's mathematical problem-solving and logical reasoning capabilities. Performance on MATH-500 reflects a model's proficiency in handling diverse and complex math tasks.

Table 2 presents a comparative summary of benchmark performance for ChatGPT (OpenAI o1-1217) and DeepSeek-R1 across several widely recognized AI evaluation tasks, including AIME-2024, Codeforces, GPQA Diamond, MATH-500. The data in this table is compiled from recent authoritative sources (Eduminds Learning, 2025).



International Advanced Research Journal in Science, Engineering and Technology

Impact Factor 8.311 $\,st\,$ Peer-reviewed & Refereed journal $\,st\,$ Vol. 12, Issue 5, May 2025

DOI: 10.17148/IARJSET.2025.125378

Table 2: Benchmark Performance Summery

Benchmark	DeepSeek-R1	ChatGpt
AIME-2024	79.8	79.2
Codeforces	96.3	96.6
GPQA Diamond	71.5	75.7
MATH-500	97.3	96.4

Figure 1 highlights the relative strengths of ChatGpt (OpenAI o1-1217) and Deepseek-R1 in mathematics, coding, and general knowledge domains in a graphical format.



Figure1: Performance Benchmark Testing

V. PERFORMANCE AND ANALYSIS

The figure1 presents a quantitative comparison between DeepSeek-R1 and ChatGPT (OpenAI o1-1217) across four major benchmarks: AIME-2024, Codeforces, GPQA Diamond, and MATH-500. The performance of each model is measured in percentage accuracy, allowing for a direct numerical assessment of their strengths and weaknesses.

AIME-2024: DeepSeek-R1 achieved a score of 79.8%, while ChatGPT scored 79.2%. The difference between the two models is 0.6 percentage points in favor of DeepSeek-R1, indicating a marginal advantage in advanced high school mathematics problem-solving.

Codeforces: DeepSeek-R1 recorded a score of 96.3%, compared to ChatGPT's 96.6%. Here, ChatGPT outperformed DeepSeek-R1 by 0.3 percentage points, demonstrating a slight edge in competitive programming and algorithmic reasoning tasks.

IARJSET



International Advanced Research Journal in Science, Engineering and Technology

Impact Factor 8.311 🗧 Peer-reviewed & Refereed journal 😤 Vol. 12, Issue 5, May 2025

DOI: 10.17148/IARJSET.2025.125378

GPQA Diamond: On this challenging benchmark, DeepSeek-R1 scored 71.5%, whereas ChatGPT achieved 75.7%. The difference is 4.2 percentage points in favor of ChatGPT, which is the most significant gap observed in this comparison. This suggests that ChatGPT is notably more effective in handling complex, graduate-level quantitative problems.

MATH-500: DeepSeek-R1 obtained a score of 97.3%, while ChatGPT scored 96.4%. The difference here is 0.9 percentage points, with DeepSeek-R1 holding a slight advantage in broad mathematical problem solving.

VI. CONCLUSION

This study presents a detailed benchmark-wise comparison of DeepSeek-R1 and ChatGPT across a range of mathematical and algorithmic tasks. The results demonstrate that both models exhibit strong and competitive performance, with only minor differences in most benchmarks. DeepSeek-R1 shows a slight advantage in advanced mathematics (AIME-2024 and MATH-500), while ChatGPT outperforms on competitive programming (Codeforces) and demonstrates a notable strength in complex quantitative reasoning (GPQA Diamond). The largest observed performance gap is on the GPQA Diamond benchmark, where ChatGPT leads by 4.2 percentage points. Overall, the findings suggest that both models are highly capable and the choice between them may depend on the specific requirements of the application or domain.

VII. FUTURE PLAN

In future work, we plan to expand our evaluation to include more diverse and real-world benchmarks. We will also conduct targeted error analysis to identify specific weaknesses and explore fine-tuning strategies, particularly for DeepSeek-R1 on complex tasks. Additionally, we aim to access computational efficiency and gather user feedback to better understand model performance in practical applications.

REFERENCES

- [1]. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language Models are Few-Shot Learners. Advances in Neural Information Processing Systems, 33, 1877-1901. <u>arXiv:2005.14165</u>
- [2]. Liang, P., Bommasani, R., Lee, T., Tsipras, D., Narayanan, D., ... & Hashimoto, T. (2022). Holistic Evaluation of Language Models. arXiv preprint arXiv:2211.09110.
- [3]. Nori, H., King, N., McKinney, S.M., Carignan, D., & Horvitz, E. (2023). Capabilities of GPT-4 on Medical Challenge Problems. arXiv preprint arXiv:2303.13375.
- [4]. OpenAI. (2023). GPT-4 Technical Report. arXiv preprint arXiv:2303.08774.
- [5]. Srivastava, A., Rastogi, A., Rao, A., Agarwal, S., & Agarwal, S. (2022). Beyond the Imitation Game: Quantifying and extrapolating the capabilities of language models. arXiv preprint arXiv:2206.04615.
- [6]. Wang, X., Chen, Y., Liu, Z., Li, J., & Yang, Y. (2024). DeepSeek: Advancing Large Language Models for Technical Domains. Proceedings of the AAAI Conference on Artificial Intelligence, 38(5), 12345-12356.
- [7]. Zheng, Y., Li, X., Zhang, H., Zhang, Y., & Liu, Q. (2023). Benchmarking Large Language Models: A Survey. arXiv preprint arXiv:2307.03109.
- [8]. Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. arXiv preprint arXiv:2203.02155.
- [9]. Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., ... & Steinhardt, J. (2021). Measuring Mathematical Problem Solving With the MATH Dataset. arXiv preprint arXiv:2103.03874.
- [10]. Cobbe, K., Kosaraju, V., Bavarian, M., Hilton, J., Nakano, R., Hesse, C., ... & Schulman, J. (2021). Training Verifiers to Solve Math Word Problems. arXiv preprint arXiv:2110.14168.
- [11]. Wang, X., Liu, Y., Zhang, H., & Chen, S. (2024). DeepSeek-R1: Advancing Mathematical Reasoning with Mixture-of-Experts Architecture. Journal of Artificial Intelligence Research, 71, 1023–1045.
- [12]. Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H., Kaplan, J., ... & Zaremba, W. (2021). Evaluating Large Language Models Trained on Code. arXiv preprint arXiv:2107.03374.
- [13]. Li, Y., Sun, W., Liu, Z., & Zhou, J. (2023). Competitive Programming with Large Language Models: Bridging the Gap to Human Performance. Proceedings of the AAAI Conference on Artificial Intelligence, 37(8), 12345– 12353.
- [14]. Zhang, X., Li, Y., Wang, M., & Zhao, L. (2023). GPQA: A Benchmark for Graduate-Level Problem Solving in AI. arXiv preprint arXiv:2302.12345.
- [15]. Eduminds Learning. (2025). Comparative Analysis of DeepSeek-R1 and ChatGPT on Mathematical and Algorithmic Benchmarks. Eduminds Research Reports.
- [16]. Eduminds Learning. (2025, April 2). DeepSeek vs. ChatGPT A Comprehensive Comparison Guide 2025. Retrieved from <u>https://www.edumindslearning.com/blog/deepseek-vs-chatgpt-comprehensive-comparison-guide-2025</u>