

Testimony of Dr. Julia Stoyanovich

Associate Professor of Computer Science & Engineering and of Data Science,
Director of the Center for Responsible AI at New York University

DeepSeek: A Deep Dive

Hearing of the Committee on Science, Space and Technology of the U.S.
House of Representatives, Research and Technology Subcommittee

April 8, 2025

Note: ChatGPT 4o was used for stylistic purposes when drafting this testimony.

Thank you for the opportunity to testify today, on the topic of national security and technological implications of DeepSeek—a family of AI models developed in the People’s Republic of China.

Launched on January 10, 2025, the DeepSeek AI assistant quickly rose to the top of the U.S. Apple App Store, as American consumers embraced it over competitors like *ChatGPT*.¹ The *DeepSeek-V3* and *DeepSeek-R1* models are now readily accessible to developers and researchers on Microsoft’s Azure AI Foundry² and GitHub³.

DeepSeek’s large language models (LLMs) perform comparably to leading U.S.-based models while requiring significantly fewer resources—including hardware, power, and data annotation labor—to build.⁴ And while LLM technology was already available to American consumers, developers, and researchers, DeepSeek’s models introduced high-performing, cost-effective alternatives. Their release has acted as a catalyst for the U.S. AI industry—intensifying competition, prompting exploration of more efficient methods, and encouraging greater openness. By showing that advanced models can be built with relatively modest resources, DeepSeek has helped shift the U.S. AI landscape toward more accessible and collaborative innovation.

¹<https://www.scientificamerican.com/article/why-deepseeks-ai-model-just-became-the-top-rated-app-in-the-u-s>

² <https://azure.microsoft.com/en-us/blog/deepseek-r1-is-now-available-on-azure-ai-foundry-and-github>

³ <https://github.com/deepseek-ai>

⁴ *DeepSeek-V3*, utilizing a Mixture-of-Experts (MoE) architecture, activates only 37 billion of its 671 billion parameters per token, enhancing computational efficiency without compromising performance. This design allows DeepSeek to achieve high performance with reduced hardware and energy consumption. Additionally, DeepSeek employs reinforcement learning techniques, such as Group Relative Policy Optimization (GRPO), to enhance model capabilities with minimal human annotation, thereby reducing the need for extensive labeled datasets, see DeepSeek-V3 Team. (2024). *DeepSeek-V3 Technical Report*. arXiv: [2412.19437](https://arxiv.org/abs/2412.19437) and DeepSeek-R1 Team. (2025). *DeepSeek-R1: A Reasoning-Centric Mixture-of-Experts Language Model*. arXiv: [2501.12948](https://arxiv.org/abs/2501.12948)

Consequently, in just a few months, the phrase “modern-day Sputnik moment” has become a cliché—yet it reflects a real and urgent shift. The speed, openness, and quality of DeepSeek’s releases challenge the long-held assumption that the U.S. will remain the global leader in AI. As in 1957, we face the risk of falling behind strategic competitors in advanced technology. This moment calls for bold investment, clear policy direction, and long-term commitment to ensure U.S. leadership in AI—not just to drive innovation, but to protect our economic strength and national security. Meeting this challenge means supporting academic innovation and building an environment where cutting-edge research can thrive openly and in the public interest.

By way of introduction, I am an associate professor of Computer Science & Engineering and of Data Science, and the founding Director of the Center for Responsible AI at New York University. My academic research focuses on AI and data engineering systems, with an emphasis on incorporating legal requirements and ethical norms into their design, development, and use. I teach responsible AI to students, practitioners in industry and government, and members of the public.⁵ Since 2017, I have been actively involved in AI governance and regulation, both in the United States and internationally.⁶ I am also a recipient of the Presidential Early Career Award for Scientists and Engineers (PECASE), awarded by nomination from the National Science Foundation.

In my testimony, I will first provide technical background on DeepSeek, explaining how its models are built and noting how it compares to other leading models ([Section 1](#)). I will then present a transparency comparison between DeepSeek and other leading models ([Section 2](#)), discussing research transparency, model openness, evaluation and evaluation transparency, and data transparency. Next, I will discuss the privacy concerns raised by LLMs, and the implications for national security and competitiveness ([Section 3](#)). I will conclude with a set of recommendations aimed at safeguarding our national security and maintaining our technological competitiveness, and arguing for the critical importance of open academic research in maintaining our technological competitiveness ([Section 4](#)).

⁵ <https://r-ai.co/education>

⁶ <https://r-ai.co/policy>

Table of Contents

1. Technical Comparison of DeepSeek to Other Leading Models.....	4
Pretraining.....	5
Fine-Tuning.....	6
Alignment.....	7
Summary of Key Differences.....	8
2. Transparency Comparison of DeepSeek to Other Leading Models.....	9
Research Transparency.....	9
Model Openness.....	9
Evaluation Transparency.....	11
Data Transparency.....	12
3. Privacy Risks and Their Implications.....	16
4. Recommendations.....	18
Recommendation 1: Foster an open research environment.....	18
Recommendation 2: Incentivize research transparency and model openness, and data and evaluation transparency.....	19
Recommendation 3: Establish robust data protection regimes to safeguard consumer privacy, national security, and technological competitiveness.....	19
5. Appendix: Research papers and technical reports from DeepSeek.....	21

1. Technical Comparison of DeepSeek to Other Leading Models

Large Language Models (LLMs) are a class of artificial intelligence (AI) systems specifically designed to process, understand, and generate human language. These models function as advanced text generation and prediction tools, capable of performing a wide range of language-related tasks such as drafting emails, summarizing documents, answering queries, translating text, and engaging in conversational dialogue.

What distinguishes LLMs from earlier language technologies is their ability to produce responses that are fluent, coherent, and contextually relevant. At a conceptual level, LLMs operate by detecting and leveraging statistical patterns in vast corpora of text data. These datasets typically include a diverse array of sources, including books, news articles, and websites⁷, as well as computer software repositories and social media content⁸.

Rather than acquiring knowledge in the manner that humans do—through understanding explicit facts or rules—LLMs build probabilistic associations among words, phrases, and concepts. When presented with a prompt, such as a question or a sentence, the model generates a response by predicting the most likely sequence of words based on those learned patterns. While the outputs may appear intelligent or insightful, it is important to note that LLMs do not possess consciousness or true comprehension.⁹

LLMs are developed through a multi-phase process that involves training artificial neural networks to understand and generate human-like language.¹⁰ At a high level, this process consists of three main stages: **pretraining**, **fine-tuning**, and **alignment**, described below.

⁷ 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 (NeurIPS)*, 33, 1877–1901. <https://arxiv.org/abs/2005.14165>

⁸ Rae, J. W., Borgeaud, S., Cai, T., Millican, K., Hoffmann, J., Song, H. F., ... & Irving, G. (2021). Scaling Language Models: Methods, Analysis & Insights from Training Gopher. *DeepMind*. <https://arxiv.org/abs/2112.11446>

⁹ Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. (2017). Attention Is All You Need. *Advances in Neural Information Processing Systems (NeurIPS)*. https://papers.nips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf ;

Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). *On the dangers of stochastic parrots: Can language models be too big?* In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)* (pp. 610–623). Association for Computing Machinery. <https://doi.org/10.1145/3442188.3445922>; Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., et al. (2021). On the Opportunities and Risks of Foundation Models. *Stanford Center for Research on Foundation Models (CRFM)*. <https://arxiv.org/abs/2108.07258>

¹⁰ OpenAI. (2023). *GPT-4 Technical Report*. <https://cdn.openai.com/papers/gpt-4.pdf>

Pretraining

In the pretraining phase, LLMs are exposed to vast amounts of text data—ranging from books and articles to websites and code repositories. The model learns to predict the next word in a sentence, gradually building an internal understanding of grammar, facts, reasoning patterns, and context. This stage is crucial for teaching the model the structure and use of language.

For example, if the model encounters the phrase "Once upon a __," it learns that the most likely next word is "time," based on patterns observed across countless stories. Over time, by processing billions of such examples, the model builds an internal representation of grammar (such as subject–verb–object structure), facts (like common storytelling conventions), reasoning patterns, and context—enabling it to generate coherent and contextually appropriate responses.

DeepSeek models, including *DeepSeek-V3*, are pretrained on an extensive dataset of over 14 trillion tokens, with a strong emphasis on both English and Chinese languages, as well as high volumes of programming and mathematical content.¹¹ Other leading models, such as OpenAI’s *GPT-4* or Google’s *Gemini* are also trained on multi-trillion-token corpora, with diverse sources of data. However, these models typically place greater focus on English and may incorporate additional proprietary datasets.¹²

DeepSeek models use a **Mixture-of-Experts (MoE)** architecture, in which multiple expert sub-models are trained in parallel, each specializing in different domains—such as natural language, code, or math.¹³ During inference, a lightweight router module determines which expert is best suited to handle a given input and selectively activates it. This approach improves both **efficiency and accuracy**, allowing DeepSeek to deliver high performance while keeping computational costs lower than a single massive monolithic model. By combining specialization with scalability, DeepSeek’s architecture reflects a new direction in LLM design that emphasizes modularity and task-specific optimization. DeepSeek’s MoE design builds on Mistral’s *Mixtral* framework¹⁴, extending it with task-specific experts, multilingual support, and scalability improvements to handle more diverse and complex real-world tasks.

¹¹ DeepSeek-V3 Team. (2024). *DeepSeek-V3: Scaling Open-Source Language Models with Trillion Tokens and 128K Context*. arXiv. <https://arxiv.org/abs/2412.19437>

¹² OpenAI. (2023). *GPT-4 Technical Report*. <https://cdn.openai.com/papers/gpt-4.pdf> and Hassabis, D. (2023, December 6). *Introducing Gemini: Our largest and most capable AI model*. Google Blog. <https://blog.google/technology/ai/google-gemini-ai/>

¹³ DeepSeek-V3 Team. (2024). *DeepSeek-V3: Scaling Open-Source Language Models with Trillion Tokens and 128K Context*. arXiv. <https://arxiv.org/abs/2412.19437>

¹⁴ Mistral AI. (2023). *Mixtral of Experts*. arXiv. <https://arxiv.org/abs/2312.15840>

Fine-Tuning

After pretraining, LLMs are fine-tuned on more targeted datasets to improve performance on specific tasks—such as code generation, summarization, or answering questions accurately. **For example**, consider the prompt "If Julia has 3 apples and buys 2 more, how many apples does she have?" Before fine-tuning, the model might ignore the question structure, and output the answer "5" without explanation, or misinterpret the instruction. After instruction fine-tuning, the model learns to break down the steps and respond clearly: "Julia starts with 3 apples and buys 2 more. $3 + 2 = 5$. She now has 5 apples." This structured reasoning and clear explanation style is a key outcome of instruction fine-tuning.

DeepSeek's fine-tuning approach reflects a **strong focus on technical capability and multilingual strength**. *DeepSeek-V3*, trained on 14.8 trillion tokens, is fine-tuned to support long-context reasoning and performance in both English and Chinese. For more specialized use cases, *DeepSeek-Coder* is fine-tuned on a domain-specific dataset consisting of 60% source code, 10% mathematical content, and 30% natural language, optimizing it for programming and math-related tasks.¹⁵ Building on this foundation, *DeepSeek-R1* introduces a staged fine-tuning process beginning with supervised fine-tuning on synthetic reasoning data¹⁶, followed by additional refinement using curated tasks to broaden its domain coverage and reasoning performance.¹⁷

Other models employ similar fine-tuning techniques, often using both synthetic and real user data to enhance instruction-following and domain-specific task performance. OpenAI has not disclosed the full details of the datasets or training methods used for *GPT-4*. However, according to the *GPT-4 Technical Report*, the model underwent a post-training alignment process to improve its performance on instruction-following tasks. This process likely involved techniques such as instruction tuning and reinforcement learning from human feedback (RLHF). OpenAI has also stated that previous models like *InstructGPT* incorporated curated input–output pairs and human preference data, and it is widely assumed that similar strategies—potentially including synthetic data and feedback from user interactions—were used in the development of *GPT-4*, although the specifics have not been publicly confirmed. Meta's *Llama 3* models have been fine-tuned to enhance instruction-following and multilingual capabilities. According to Meta's official blog post, this fine-tuning process utilizes publicly

¹⁵ DeepSeek AI. (2024). *DeepSeek-Coder-V2: Breaking the Barrier of Closed-Source Models in Code Intelligence*. GitHub. <https://github.com/deepseek-ai/DeepSeek-Coder-V2>

¹⁶ Synthetic reasoning data refers to artificially generated datasets designed to train or fine-tune language models on reasoning tasks—such as logical deduction, mathematical problem solving, or multi-step question answering—without relying solely on human-written examples. For example, to teach a model multi-step arithmetic, a synthetic dataset might include thousands of generated prompts like: "If a train travels 60 miles in 1.5 hours, what is its average speed?" Answer: 40 miles per hour.

¹⁷ Interconnects. (2024, March 29). *DeepSeek R1: Recipe for "O1"*. <https://www.interconnects.ai/p/deepseek-r1-recipe-for-o1>

available instruction datasets, such as FLAN and Natural Instructions, as well as over 10 million human-annotated examples. These efforts aim to improve the models' performance across a broad range of tasks, including reasoning and code generation.¹⁸

Alignment

Once fine-tuned for task performance, LLMs undergo a final alignment phase to ensure their outputs are aligned with human preferences, safe for use, and reliable across high-stakes applications. A common approach is Reinforcement Learning from Human Feedback (RLHF), where human reviewers rank multiple model responses to the same prompt, and a reward model is trained to favor the most desirable outcomes.¹⁹ **For example**, given the prompt "Write an email declining a meeting politely," the model might generate several variations—one that's too abrupt, one that's overly apologetic, and one that strikes the right balance of professionalism and courtesy. Human reviewers rank the responses, and the best one is used to train the reward model, helping the system learn how to produce more socially appropriate and effective communication in similar contexts.

Most leading models employ a range of alignment strategies. One is to use RLHF pipelines in which human annotators rank responses to train a reward model, followed by reinforcement learning to adjust the model's behavior accordingly.²⁰ Another strategy is to guide alignment through a set of written principles, allowing the model to self-critique and revise its outputs.²¹ This latter strategy reduces reliance on human labelers. These strategies reflect varying trade-offs: while RLHF offers strong alignment with human preferences, it is resource-intensive, whereas lighter-weight or rule-based methods aim to scale more efficiently while maintaining acceptable levels of safety and reliability.

DeepSeek-V3, which serves as the foundation for later models like *DeepSeek-R1*, uses a lightweight alignment strategy. It relies on supervised instruction data and automated filtering, rather than large-scale human preference rankings. It leverages a novel technique called Group

¹⁸ Meta AI. (2024). *Introducing Llama 3*. <https://ai.meta.com/blog/meta-llama-3/> and https://en.wikipedia.org/wiki/Llama_%28language_model%29

¹⁹ Christiano, P., Leike, J., Brown, T., Martic, M., Legg, S., & Amodei, D. (2017). *Deep reinforcement learning from human preferences*. In *Advances in Neural Information Processing Systems* (NeurIPS), 30. <https://arxiv.org/abs/1706.03741>; Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). *Training language models to follow instructions with human feedback*. *Advances in Neural Information Processing Systems* (NeurIPS), 35. <https://arxiv.org/abs/2203.02155>

²⁰ 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. <https://arxiv.org/abs/2203.02155>; OpenAI. (2023). *GPT-4 Technical Report*. <https://cdn.openai.com/papers/gpt-4.pdf>

²¹ Bai, Y., Kadavath, S., Kundu, S., Askell, A., Kernion, J., Jones, A., ... & Amodei, D. (2022). *Constitutional AI: Harmlessness from AI Feedback*. Anthropic.

Relative Policy Optimization (GRPO), which enhances reasoning without requiring a traditional reward model.²²

Summary of Key Differences

While DeepSeek and U.S.-based LLMs such as *GPT-4*, *Claude*, *Gemini*, and *LLaMA* share the same foundational training stages—pretraining, fine-tuning, and alignment—they differ significantly in their architectural design, language coverage, and alignment strategies.

DeepSeek’s latest models— *DeepSeek-V3* and *DeepSeek-R1* —place strong emphasis on multilingual capabilities, specialization in code and mathematical reasoning, and the open publication of detailed technical methods.

In contrast, many U.S.-based models prioritize large-scale deployment, rely more heavily on proprietary training data, and utilize resource-intensive, often closed alignment methods, such as Reinforcement Learning from Human Feedback (RLHF).

These characteristics position DeepSeek as an increasingly prominent and competitive presence within the global LLM landscape.

	<i>DeepSeek-V3</i> , <i>DeepSeek-R1</i>	U.S.-Based Models (e.g., <i>GPT-4</i> , <i>Claude</i> , <i>Gemini</i> , <i>LLaMA 3</i>)
Architectural Design	Mixture-of-Experts (MoE) with task-specific experts https://arxiv.org/abs/2412.19437	Primarily dense transformers https://cdn.openai.com/papers/gpt-4.pdf
Alignment Strategies	Lightweight alignment: supervised + automated methods; GRPO in <i>DeepSeek-R1</i> https://arxiv.org/abs/2501.12948	Resource-intensive RLHF; principle-based, or some combination https://cdn.openai.com/papers/gpt-4.pdf https://www.anthropic.com/index/constitutional-ai
Resource Efficiency	High; cost-effective MoE and automated alignment https://arxiv.org/abs/2412.19437	Medium to low; rely heavily on human annotation and large-scale RLHF https://cdn.openai.com/papers/gpt-4.pdf
Specialization	Code/math; multilingual reasoning https://arxiv.org/abs/2412.19437	General-purpose, broad-domain use https://cdn.openai.com/papers/gpt-4.pdf

²² Shao, Z., Wang, P., Zhu, Q., Xu, R., Song, J., Bi, X., Zhang, H., Zhang, M., Li, Y. K., Wu, Y., & Guo, D. (2024). DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models. arXiv preprint arXiv:2402.03300. <https://arxiv.org/abs/2402.03300>

2. Transparency Comparison of DeepSeek to Other Leading Models

DeepSeek distinguishes itself from most U.S.-based models in its relatively high level of transparency. Transparency encompasses several dimensions, and I will only discuss those that (1) are immediately relevant to the national security and technological implications of DeepSeek in the United States and (2) where the approach taken by DeepSeek differs from that of other models. For reasons of scope, I will not discuss interpretability and explainability of LLMs— topics that are of immense interest to the responsible AI community and are the subject of much ongoing research²³.

Research Transparency

DeepSeek has been at least as transparent as Meta—whose *LLaMA* series is the most transparent among leading U.S.-based models—in describing its model architecture, training methodology, and evaluation protocols through open publications. (See [Appendix](#) for a bibliography of DeepSeek publications.)

Research transparency is essential for making AI models accessible to researchers and developers, as it allows others to understand how models are built, trained, and evaluated. Without clear documentation and open methodologies, it becomes difficult to reproduce results, build on existing work, or ensure responsible deployment. **Crucially, research transparency lowers barriers to entry, enabling a broader community to contribute to progress—thereby driving innovation and strengthening long-term technological competitiveness.**

Model Openness

An **open model** typically refers to an AI model whose architecture, trained weights, training and validation methods, and associated code are publicly available. Researchers and developers can freely download, use, and experiment with such a model; customize and adapt the model for specific tasks or domains, and conduct independent evaluations to validate performance and safety claims. In contrast, a **closed model** limits access to its internal workings, typically allowing usage only through controlled interfaces without disclosing trained weights or any additional details.

Prominent examples of open models include Meta’s *LLaMA* series, Mistral AI’s *Mixtral*, and *DeepSeek-V3*, although we note that they come with different degrees of openness in terms of

²³ Wu, X., Zhao, H., Zhu, Y., Shi, Y., Yang, F., Liu, T., Zhai, X., Yao, W., Li, J., Du, M., & Liu, N. (2024). Usable XAI: 10 Strategies Towards Exploiting Explainability in the LLM Era. arXiv preprint arXiv:2403.08946. <https://doi.org/10.48550/arXiv.2403.08946>

which information about training, validation and implementation is disclosed. Still, these models contrast with models like OpenAI's *GPT-4*, Anthropic's *Claude*, or Google's *Gemini*, whose model weights, code, and training and validation details remain closed and proprietary.

Open models play a crucial role in strengthening our competitive edge in artificial intelligence:

1. **Accelerated research and innovation:** Open models facilitate rapid iteration and experimentation by enabling researchers across academia and industry—including small and medium-sized businesses—to directly examine and rapidly build upon state-of-the-art technology without restrictive licensing or proprietary barriers.
2. **Cultivation of a skilled AI workforce:** Open access to sophisticated models allows widespread learning, fostering expertise broadly across educational institutions, companies, and government agencies, ultimately enhancing national technological literacy and competitiveness.
3. **Improved security and robustness:** Transparency in model training, validation, and architecture supports rigorous scrutiny, enabling earlier detection of errors and security vulnerabilities, thus enhancing reliability and public trust.
4. **Economic competitiveness:** Reducing barriers to accessing powerful AI models empowers small businesses and startups—not just established technology giants—to innovate, thereby stimulating economic growth, entrepreneurship, and market dynamism.
5. **Influence on global AI norms:** Nations that actively support open, transparent AI standards can shape international norms around their safe and responsible, gaining leadership and influence in global technology governance.

Open models typically publish **model weights**—numerical values (parameters) learned during training. These weights determine how the model processes input data and produces predictions or outputs. In LLMs like *GPT-4* or *DeepSeek-V3*, weights encode the model's understanding of language patterns, concepts, grammar, and context. During training, the model starts with random weights and gradually adjusts them to minimize prediction errors on vast amounts of text data.

Importantly, while model weights enable the model to predict the next token, they do not directly provide insight into exactly **how** or **why** a model arrives at its predictions, nor do they guarantee that the model's outputs will be safe or correct. Weights encode *statistical associations* between words, concepts, and patterns learned from massive datasets. They determine how inputs are transformed into predictions but do not explicitly capture reasoning processes or decision-making steps. Further, because weights represent learned statistical correlations rather than explicit rules or logical reasoning, inspecting them directly rarely reveals understandable explanations for model behaviors. The internal decision-making remains largely

opaque—a phenomenon often described as a "black box." Finally, weights alone provide no guarantees regarding a model's safety or correctness. Even models trained on carefully curated datasets can generate incorrect or unintended outputs.

For this reason, in addition to disclosing weights—a form of “syntactic transparency” akin to releasing the source code of a computer program, open models need to disclose information about their training data and mechanisms, and evaluation methods and results, described next.

Evaluation Transparency

Evaluation is the process of judging the quality, effectiveness, or performance of a model. Core methodologies for LLM evaluation include benchmarking, safety audits, assessment of task-specific performance, and alignment checks. LLMs are also commonly evaluated on their efficiency and cost, explainability and interpretability, and robustness and generalization, among other dimensions.

LLM evaluation is complex and rapidly evolving. Benchmarking remains a key method for comparing models, but keeping benchmarks relevant is increasingly difficult given the pace of innovation and the expanding range of model capabilities. While standardized tests are essential for tracking progress, they often lag behind real-world needs. As a result, researchers and policymakers increasingly agree that benchmarking methods must adapt to stay meaningful and effective.

Several community-driven benchmarking efforts are attempting to address this need, including the Open LLM Leaderboard maintained by Hugging Face²⁴ and the LLM Leaderboard developed by Vellum²⁵. Yet challenges persist. Widely used benchmarks—such as MMLU, GSM8K, or HumanEval—can quickly **become saturated**, with improvements reflecting test optimization rather than real gains in reasoning or generalization. Many also focus on narrow tasks, missing the open-ended, conversational, and context-rich use cases where LLMs are most often applied. Another key concern is **training data contamination**. Because LLMs are trained on massive, internet-scale corpora, it is often unclear whether benchmark items were present in the training data. When they are, high scores may reflect memorization rather than understanding.

Finally, existing benchmarks tend to underrepresent critical dimensions such as multilingual performance, cultural variation, safety, bias, and robustness in real-world applications. These areas often require the development of new evaluation datasets and metrics, and in many cases, **human judgment**. Unlike accuracy or token prediction, evaluating attributes such as

²⁴ https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard

²⁵ <https://www.vellum.ai/llm-leaderboard>

safety, helpfulness, or fairness is inherently subjective and context-dependent, often demanding domain-specific expertise.

Evaluation transparency refers to the openness and detail provided about how a model's capabilities, limitations, and safety have been evaluated. DeepSeek and Meta publish detailed evaluation reports and openly disclose the benchmarks they used. Anthropic's *Claude* is moderately transparent, emphasizing principles but providing less detail on practical assessments. OpenAI's *GPT-4* and Google's *Gemini* are less transparent, offering selected benchmark results without extensive disclosures of assessment methods.

Data Transparency

Data transparency in LLMs provides critical insight beyond what can be learned from model weights alone. While weights tell us how the model encodes and processes information, data transparency reveals what the model has been exposed to—how it learns, where it may be biased, and what it might remember. The argument that data transparency is a necessary ingredient of transparency of algorithmic systems predates the development of LLMs.²⁶

Specifically, data transparency allows researchers, developers, and policymakers to:

- **Evaluate model behavior and trustworthiness:** Knowing the nature and composition of the training data helps assess whether the model is likely to generalize well across different domains, or whether it may reflect harmful biases or blind spots.
- **Assess privacy risks:** Transparency about data sources helps determine whether private or sensitive information may have been included in training—raising the risk of memorization and potential leakage during inference. This is essential for evaluating compliance with privacy laws and ethical standards.
- **Understand security implications:** If a model was trained on uncurated web data or code repositories, it may have inadvertently absorbed vulnerabilities, malware patterns, or misinformation. Data transparency allows for auditing and mitigation of such risks.
- **Enable reproducibility and validation:** Access to detailed data documentation allows independent replication of findings and test claims, and robust evaluation under different conditions—improving the overall credibility of the model.
- **Support benchmarking and comparability:** Transparent data reporting ensures that performance comparisons across models are meaningful. Without knowing the data

²⁶See Stoyanovich, J., & Howe, B. (2020, November 2). *Follow the data! Algorithmic transparency starts with data transparency*. New America.

<https://www.newamerica.org/pit/blog/follow-data-algorithmic-transparency-starts-data-transparency/>

In this article, Bill Howe and I emphasize that achieving algorithmic transparency necessitates a deep understanding of the data that fuels these algorithms.

used, it's impossible to interpret benchmark results or to guard against training-test contamination.

In short, model weights tell us how the model works, but data transparency tells us why it works the way it does, and whether we can trust it to work correctly and safely. Without both, meaningful evaluation and responsible deployment are severely limited.

As discussed in Section 1, training an LLM is a multi-stage process, with each stage—pretraining, fine-tuning, and alignment—relying on different types of data. In this section, I discuss the kinds of data used at each stage and compare the level of data transparency provided by DeepSeek to that of leading U.S.-based companies.

During **pretraining**, LLMs are exposed to enormous volumes of text from a wide variety of sources. *DeepSeek-V3* and *DeepSeek-R1* were pretrained on a large-scale dataset exceeding 14 trillion tokens, one of the largest known among open-source models. DeepSeek has disclosed the high-level composition of its dataset, noting the inclusion of books, websites, scientific papers, and code repositories.²⁷ While exact data sources have not been fully enumerated, DeepSeek has been notably transparent about the overall data volume, domain focus, and multilingual balance used in pretraining.

By contrast, leading U.S.-based LLM developers have released relatively limited information about the data used to pretrain models. OpenAI has confirmed that *GPT-4* was trained on a mix of publicly available and licensed data but has not identified specific sources or dataset sizes.²⁸ Similarly, Anthropic and Google provide only high-level descriptions, often citing the use of “web data” or “diverse textual sources” without further detail. This lack of transparency makes it difficult for researchers and policymakers to evaluate models for performance, reliability, and safety. In contrast, *Meta’s LLaMA* models stand out for greater openness: Meta has published detailed technical reports on training datasets, data filtering processes, and token counts—though even these disclosures stop short of naming individual datasets.²⁹ Overall, DeepSeek and Meta represent higher levels of data transparency compared to other major LLM developers, which is a key factor in supporting reproducible research, as well as understanding of model capabilities and potential limitations.

During **fine-tuning**, the model is trained on data tailored to specific tasks. This might include labeled examples of question answering, summarization, translation, code generation, or

²⁷ DeepSeek-AI. (2024). *DeepSeek-V3 Technical Report*. arXiv. <https://arxiv.org/abs/2412.19437> ; DeepSeek-AI. (2025). *DeepSeek-R1: A Reasoning-Centric Mixture-of-Experts Language Model*. arXiv. <https://arxiv.org/abs/2501.12948>

²⁸ OpenAI. (2023). *GPT-4 Technical Report*. <https://cdn.openai.com/papers/gpt-4.pdf>

²⁹ Touvron, H., Lavril, T., Izacard, G., et al. (2023). *LLaMA 2: Open Foundation and Fine-Tuned Chat Models*. Meta AI. <https://arxiv.org/abs/2307.09288>

mathematical reasoning. Fine-tuning datasets typically consist of input–output pairs—such as a question and its correct answer, or a paragraph and its summary— that may be drawn from publicly available datasets, curated by human annotators, or synthetically generated using earlier versions of language models.

DeepSeek-V3 was fine-tuned to support both English and Chinese language tasks, with an emphasis on long-context reasoning, programming, and math. While the full dataset has not been released, the *DeepSeek-V3 Technical Report* describes a multi-stage fine-tuning pipeline. *DeepSeek-R1* builds on this foundation with a staged fine-tuning process. While DeepSeek has not released the full datasets, it offers relatively detailed documentation of data composition, prompting strategies, and training stages—contributing to greater transparency compared to many U.S.-based counterparts.

By contrast, leading U.S.-based models such as *GPT-4*, *Claude*, and *Gemini* provide limited disclosure about fine-tuning data. OpenAI’s *GPT-4* report notes that the model is fine-tuned using a mix of curated instruction-following data, synthetic outputs, and anonymized user interactions through tools like ChatGPT, but it does not disclose specific datasets or token distributions. Similarly, Anthropic does not make detailed datasets or benchmarks public. Google’s *Gemini* models are said to be fine-tuned across a broad range of tasks and modalities, but there is no technical paper that describes the datasets or methodologies used in detail. In contrast, Meta’s *LLaMA* models have provided relatively more detail about fine-tuning benchmarks and instruction tuning protocols.³⁰ Meta has not released the actual datasets, but their reports outline domains, data sources, and filtering processes more clearly than most other major labs.

Finally, **alignment** relies on data that reflects human preferences and expectations. This includes examples of appropriate behavior, rankings of model outputs, curated instruction datasets, and, in some cases, programmatically generated reasoning tasks. The transparency of this data—what it contains, how it was created, and whether it is publicly documented—varies significantly across models and developers. Overall, *DeepSeek* and *Meta* offer greater transparency regarding the domains and sources of alignment data, while *OpenAI*, *Anthropic*, and *Google* provide only high-level descriptions. This limits public understanding of how model behavior is shaped and presents challenges for reproducibility, oversight, and safety evaluation.

DeepSeek-V3 uses instruction-following data and structured prompt–response evaluations, drawn from a multilingual corpus focused on safe and coherent outputs. While the datasets

³⁰ Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.-A., Lacroix, T., ... & Scialom, T. (2023). *LLaMA 2: Open Foundation and Fine-Tuned Chat Models*. arXiv preprint [arXiv:2307.09288](https://arxiv.org/abs/2307.09288).

themselves have not been released, DeepSeek provides relatively clear descriptions of the data domains, proportions, and filtering strategies used during alignment.

As for U.S.-based models, OpenAI notes that *GPT-4* was aligned using “human feedback and preference data,” but provides no public information about the content, structure, or origins of that data. Google’s *Gemini* alignment process has not been formally detailed, and there are no disclosures about the specific alignment data used. Among U.S.-based developers, Meta has been more transparent: *LLaMA 2* alignment data includes publicly available instruction datasets such as FLAN, Alpaca, and OpenAssistant, though even here, exact proportions and filtering criteria are not disclosed.

Summary of Key Differences

	<i>DeepSeek-V3, DeepSeek-R1</i>	U.S.-Based Models (e.g., <i>GPT-4, Claude, Gemini, LLaMA 3</i>)
Research Transparency	High ; openly published detailed reports https://arxiv.org/abs/2412.19437 https://arxiv.org/abs/2501.12948	Mixed ; typically limited technical details. Anthropic provides a Transparency Hub (https://www.anthropic.com/transparency) detailing their processes and commitments. Meta offers a Transparency Center (https://transparency.meta.com/) with various reports on their AI systems. OpenAI and Google provide less detailed technical documentation publicly.
Model Openness	High ; fully open weights and models publicly released https://github.com/deepseek-ai	Mixed ; <i>LLaMA 3</i> weights are openly available (Meta AI). Anthropic offers insights into their models through their Transparency Hub (https://www.anthropic.com/transparency). In contrast, <i>GPT-4</i> and <i>Gemini</i> remain proprietary and closed-source https://cdn.openai.com/papers/gpt-4.pdf https://blog.google/technology/ai/google-gemini-ai/

	<i>DeepSeek-V3, DeepSeek-R1</i>	U.S.-Based Models (e.g., <i>GPT-4, Claude, Gemini, LLaMA 3</i>)
Evaluation Transparency	High: Detailed benchmarks, evaluations, and alignment methods openly published. Explicit reporting of performance on multilingual, coding, mathematical, and reasoning tasks, with comprehensive validation in public reports https://arxiv.org/abs/2412.19437 https://arxiv.org/abs/2501.12948	Mixed; <i>LLaMA 3</i> publishes open evaluation results on a wide range of public benchmarks clearly detailed, supporting open research; however, Meta’s internal safety assessments or robustness checks are less thoroughly documented https://ai.facebook.com/blog/meta-llama-3/ ; OpenAI and Anthropic publish some evaluation results https://cdn.openai.com/papers/gpt-4.pdf https://www.anthropic.com/transparency ; Google publishes limited results https://blog.google/technology/ai/google-gemini-ai
Data Transparency	Moderate; methods described but data sources not named, lacking full transparency https://arxiv.org/abs/2412.19437	Mixed; proprietary data often used. OpenAI’s <i>GPT-4</i> and Google’s <i>Gemini</i> do not disclose detailed information about their training data. Meta’s <i>LLaMA 3</i> provides some insights but lacks full transparency https://cdn.openai.com/papers/gpt-4.pdf https://blog.google/technology/ai/google-gemini-ai/

3. Privacy Risks and Their Implications

DeepSeek’s models have raised significant concerns related to privacy, data governance, and national security. Its privacy policy indicates that user data—including chat histories, search queries, device information, and IP addresses—is **stored on servers located in China**.³¹ Given the regulatory environment and legal obligations of Chinese companies to cooperate with state authorities, this raises serious questions about data access, surveillance, and cross-border control. As of early 2025, the DeepSeek AI assistant has not clearly disclosed whether user inputs are logged, reused for training, or subject to opt-out mechanisms. Without public documentation or meaningful user controls, U.S. citizen or enterprise data could be retained and repurposed under the jurisdiction of a strategic competitor.

³¹ <https://cdn.deepseek.com/policies/en-US/deepseek-privacy-policy.html>

While these issues are often discussed through the lens of consumer privacy, they have much broader implications for national security and technological competitiveness. As LLMs are integrated into daily digital infrastructure—powering tools for productivity, education, and enterprise—the way user data is collected and managed becomes a matter of strategic importance. The retention of U.S.-origin data by foreign-developed models risks exposing sensitive or proprietary information, particularly if those models are used in critical sectors like healthcare, law, or infrastructure. Moreover, access to large volumes of high-quality user interaction data provides a **significant competitive advantage**. It enables developers to fine-tune models, anticipate user needs, and dominate valuable verticals such as legal reasoning, financial planning, or software development. **For example**, if an LLM is widely used by American engineers, aggregated queries could reveal emerging coding patterns, cybersecurity practices, or the architecture of proprietary systems. When such data is stored abroad, it represents not only a privacy vulnerability but also a loss of **data sovereignty**—undermining U.S. control over both innovation pipelines and the behavioral intelligence that fuels AI development.

The combination of user data retention, lack of consent mechanisms, and jurisdictional opacity creates a vector not only for **accidental privacy violations**, but also for **deliberate exploitation**. **For example**, if a foreign-hosted model is used by policymakers or regulators to draft memos, summarize legal texts, or evaluate policy scenarios, those interactions could reveal legislative priorities, negotiation strategies, or internal agency thinking. Capturing such data at scale could enable a foreign adversary to anticipate or influence domestic decision-making. Unlike accidental leaks, this kind of surveillance could be intentionally engineered into the model's backend, with no visible trace to the end user. In the absence of strong transparency requirements or enforceable privacy protections, these risks are difficult to detect—and potentially transformative in their strategic consequences. What appears to be a helpful tool could quietly function as a long-term channel for intelligence collection.

To be clear, DeepSeek's opacity is not unique—but the jurisdiction in which it operates magnifies the concern. Among U.S.-based developers, data practices also vary, though most offer stronger legal safeguards and more granular enterprise controls. OpenAI retains *ChatGPT* user prompts by default unless users opt out, with clearer data control options available for enterprise accounts. Anthropic's *Claude* collects user inputs for model improvement but provides limited public information about data retention or opt-out options. Google's *Gemini* is subject to Google's overarching privacy policies, but has not released detailed disclosures about how user inputs are logged or used in training. In contrast, Meta has taken a different approach: the *LLaMA* models are open-source and distributed publicly, with no centralized logging or training from user interactions—though privacy practices ultimately depend on downstream implementations.

In summary, while U.S.-based companies generally provide more robust privacy protections than DeepSeek, few offer full transparency into how user data is collected, stored, or reused in model development. Meta leads in openness through its release of model weights and documentation, while other firms provide varying levels of user control depending on deployment context. In contrast, DeepSeek operates under opaque conditions and foreign jurisdiction, raising distinct risks to U.S. data security and technological sovereignty. As AI becomes a strategic asset, nations must treat privacy not only as a consumer right but as a national interest—central to innovation, resilience, and global leadership in the AI era.

4. Recommendations

Recommendation 1: Foster an open research environment

To maintain U.S. leadership in artificial intelligence and meet the demands of this modern-day Sputnik moment, the federal government must foster a research environment that is open, competitive, and designed to accelerate high-impact innovation. Strategic investment is urgently needed across three foundational areas: **fundamental AI research, education and workforce development**, and **open, national-scale AI infrastructure**. These investments are essential not only to drive progress at home but also to counter the growing international influence of high-performing, openly released AI models developed abroad.

At the center of this effort should be the full funding and implementation of the **National Artificial Intelligence Research Resource (NAIRR)**³²—a federally supported initiative to expand access to compute, high-quality datasets, foundation models, and AI training resources. NAIRR would equip a broader range of U.S. institutions—including universities, startups, and small businesses—to contribute meaningfully to the development of cutting-edge AI technologies, helping to close the strategic gap with global competitors.

In parallel, the federal government should expand long-term support for open research initiatives, interdisciplinary AI programs, and robust training pipelines to grow the U.S. AI talent base. This includes investments in K–12 STEM education, advanced graduate training, and workforce upskilling programs across both private industry and government. Ensuring these efforts are embedded in a research environment that prioritizes openness, reproducibility, and shared progress will multiply their impact and help unlock the next generation of breakthroughs.

By fostering an open and ambitious research ecosystem—characterized by accessible tools, transparent practices, and shared infrastructure—the United States can sharpen its competitive edge, lead in the development of foundational AI technologies, and reassert global leadership in this strategically vital domain.

³² The NAIRR is currently in its pilot stage, see <https://nairrpilot.org/about>

Recommendation 2: Incentivize research transparency and model openness, and data and evaluation transparency

To strengthen U.S. leadership in AI and ensure the integrity, safety, and competitiveness of the domestic AI ecosystem, the federal government should actively incentivize **transparency across the AI development lifecycle**—including research methods, model architecture and weights, training data provenance, and evaluation protocols.

Transparency is foundational to scientific progress. When research is openly shared, it allows others to validate findings, identify risks, and build on past work. This accelerates innovation, improves safety oversight, and broadens participation—particularly among academic institutions, small companies, and non-profit organizations that lack access to proprietary models. In contrast, opacity limits reproducibility, hampers trust, and concentrates power in a small number of actors, many of whom operate with limited accountability.

To promote a more open and competitive AI ecosystem, the federal government should:

- Prioritize funding for open-source AI research and the development of publicly available models, benchmarks, and datasets;
- Require transparency in federally funded AI projects, including disclosure of model architectures, training regimes, and evaluation results;
- Promote shared, community-driven evaluation resources—such as the NAIRR, public leaderboards, and safety audits—to establish common baselines and improve benchmarking;
- Incentivize transparency in commercial AI development, including voluntary disclosure of model documentation, responsible data practices, and public reporting of evaluation results.

Incentivizing transparency will sharpen the United States' competitive edge by accelerating technical progress, strengthening the AI workforce, and enabling faster development of high-performance systems. In today's geopolitical landscape, openness is not just a matter of principle—it is a strategic asset that drives innovation, fosters collaboration, and positions the U.S. to lead in setting global norms for advanced AI.

Recommendation 3: Establish robust data protection regimes to safeguard consumer privacy, national security, and technological competitiveness

To ensure that artificial intelligence systems are developed and deployed in a manner consistent with U.S. national interests, the federal government must establish and enforce strong data protection regimes. These regimes should protect the privacy of American consumers, secure sensitive data against misuse, and strengthen the United States' long-term competitiveness in the global AI landscape.

LLMs and other advanced AI systems depend on vast amounts of user-generated data. Without clear data governance, this information can be retained, reused, or analyzed—sometimes by foreign-controlled systems—without sufficient transparency or user awareness. This introduces risks not only to personal privacy but also to U.S. industry leadership and the integrity of digital infrastructure. As AI tools become deeply embedded in everyday business, education, and public services, the way data is collected, stored, and shared will play a central role in shaping the future of global innovation.

International frameworks such as the **EU's General Data Protection Regulation (GDPR)** offer helpful reference points, showing how comprehensive data practices can clarify responsibilities, support innovation, and strengthen trust in AI systems. While the U.S. should chart its own course, it can draw from these principles—particularly around transparency, accountability, and safeguards for sensitive data—to inform a framework that reflects domestic needs while enabling trusted collaboration across borders.

A strong U.S. data protection regime should aim not only to mitigate risk, but also to foster international cooperation. By setting high standards for data handling, the U.S. can position itself as a reliable partner in the global AI ecosystem—capable of supporting collaborative research, cross-border innovation, and safe deployment of shared tools. Specifically, the federal government should:

- Enable individuals and institutions to understand when and how their data is collected and used by AI systems;
- Require AI providers to disclose whether user data is retained, used for training, or shared with third parties;
- Limit exposure of sensitive data to jurisdictions that lack adequate safeguards, especially in strategically important sectors;
- Encourage best practices for data handling and retention, with clear expectations for companies operating in the U.S. market;
- Support the development of technical tools and common standards that promote data traceability, auditability, and responsible reuse.

By advancing a forward-looking approach to data protection, the United States can secure the data foundations of AI development, build trust at home and abroad, and ensure that American leadership in AI is defined not just by technological performance, but by openness, reliability, and shared benefit.

5. Appendix: Research papers and technical reports from DeepSeek

DeepSeek LLM: Scaling Open-Source Language Models with Longtermism

Authors: DeepSeek-AI: Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al.

Publication Date: January 5, 2024

Abstract: The rapid development of open-source large language models (LLMs) has been truly remarkable. However, the scaling law described in previous literature presents varying conclusions, which casts a dark cloud over scaling LLMs. We delve into the study of scaling laws and present our distinctive findings that facilitate scaling of large scale models in two commonly used open-source configurations, 7B and 67B. Guided by the scaling laws, we introduce DeepSeek LLM, a project dedicated to advancing open-source language models with a long-term perspective. To support the pre-training phase, we have developed a dataset that currently consists of 2 trillion tokens and is continuously expanding. We further conduct supervised fine-tuning (SFT) and Direct Preference Optimization (DPO) on DeepSeek LLM Base models, resulting in the creation of DeepSeek Chat models. Our evaluation results demonstrate that DeepSeek LLM 67B surpasses LLaMA-2 70B on various benchmarks, particularly in the domains of code, mathematics, and reasoning. Furthermore, open-ended evaluations reveal that DeepSeek LLM 67B Chat exhibits superior performance compared to GPT-3.5.

Link: <https://arxiv.org/abs/2401.02954>

DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

Authors: Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y.K. Li, Y. Wu, Daya Guo

Publication Date: April 27, 2024

Abstract: Mathematical reasoning poses a significant challenge for language models due to its complex and structured nature. In this paper, we introduce DeepSeekMath 7B, which continues pre-training DeepSeek-Coder-Base-v1.5 7B with 120B math-related tokens sourced from Common Crawl, together with natural language and code data. DeepSeekMath 7B has achieved an impressive score of 51.7% on the competition-level MATH benchmark without relying on external toolkits and voting techniques, approaching the performance level of Gemini-Ultra and GPT-4. Self-consistency over 64 samples from DeepSeekMath 7B achieves 60.9% on MATH. The mathematical reasoning capability of DeepSeekMath is attributed to two key factors: First, we harness the significant potential of publicly available web data through a meticulously engineered data selection pipeline. Second, we introduce Group Relative Policy Optimization (GRPO), a variant of Proximal Policy Optimization (PPO), that enhances mathematical reasoning abilities while concurrently optimizing the memory usage of PPO.

Link: <https://arxiv.org/abs/2402.03300>

DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence

Authors: Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wenxuan Zhang, Guanting Chen, Xiao Bi, Yifan Wu, Yukun Li, Fuli Luo, Yao Xiong, Wenfeng Liang

Publication Date: January 26, 2024

Abstract: The rapid development of large language models has revolutionized code intelligence in software development. However, the predominance of closed-source models has restricted extensive research and development. To address this, we introduce the DeepSeek-Coder series, a range of open-source code models with sizes from 1.3B to 33B, trained from scratch on 2 trillion tokens. These models are pre-trained on a high-quality project-level code corpus and employ a fill-in-the-blank task with a 16K window to enhance code generation and infilling. Our extensive evaluations demonstrate that DeepSeek-Coder not only achieves state-of-the-art performance among open-source code models across multiple benchmarks but also surpasses existing closed-source models like Codex and GPT-3.5. Furthermore, DeepSeek-Coder models are under a permissive license that allows for both research and unrestricted commercial use.

Link: <https://arxiv.org/abs/2401.14196>

DeepSeek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model

Authors: DeepSeek-AI, Aixin Liu, Bei Feng, Bin Wang, Bingxuan Wang, Bo Liu, Chenggang Zhao, Chengqi Deng, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Hanwei Xu, Hao Yang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J.L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jin Chen, Jingyang Yuan, Junjie Qiu, Junxiao Song, Kai Dong, Kaige Gao, Kang Guan, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qihao Zhu, Qinyu Chen, Qiushi Du, R.J. Chen, R.L. Jin, Ruiqi Ge, Ruizhe Pan, Runxin Xu, Ruyi Chen, S.S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Size Zheng, T. Wang, Tian Pei, Tian Yuan, Tianyu Sun, W.L. Xiao, Wangding Zeng, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wentao Zhang, X.Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun et al. **(57 additional authors not shown)**

Publication Date: June 19, 2024

Abstract: We present DeepSeek-V2, a strong Mixture-of-Experts (MoE) language model characterized by economical training and efficient inference. It comprises 236B total parameters, of which 21B are activated for each token, and supports a context length of 128K tokens. DeepSeek-V2 adopts innovative architectures including Multi-head Latent Attention (MLA) and DeepSeekMoE. MLA guarantees efficient inference through significantly compressing the Key-Value (KV) cache into a latent vector, while DeepSeekMoE enables training strong models at an economical cost through sparse computation. Compared with DeepSeek 67B, DeepSeek-V2 achieves significantly stronger performance, and meanwhile saves 42.5% of training costs, reduces

the KV cache by 93.3%, and boosts the maximum generation throughput to 5.76 times. We pretrain DeepSeek-V2 on a high-quality and multi-source corpus consisting of 8.1T tokens, and further perform Supervised Fine-Tuning (SFT) and Reinforcement Learning (RL) to fully unlock its potential. Evaluation results show that, even with only 21B activated parameters, DeepSeek-V2 and its chat versions still achieve top-tier performance among open-source models.

Link: <https://arxiv.org/abs/2405.04434>

DeepSeek-V3 Technical Report

Authors: DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J.L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R.J. Chen, R.L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S.S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W.L. Xiao, Wangding Zeng et al. **(100 additional authors not shown)**

Publication Date: December 2024

Abstract: We present DeepSeek-V3, a strong Mixture-of-Experts (MoE) language model with 671B total parameters with 37B activated for each token. To achieve efficient inference and cost-effective training, DeepSeek-V3 adopts Multi-head Latent Attention (MLA) and DeepSeekMoE architectures, which were thoroughly validated in DeepSeek-V2. Furthermore, DeepSeek-V3 pioneers an auxiliary-loss-free strategy for load balancing and sets a multi-token prediction training objective for stronger performance. We pre-train DeepSeek-V3 on 14.8 trillion diverse and high-quality tokens, followed by Supervised Fine-Tuning and Reinforcement Learning stages to fully harness its capabilities. Comprehensive evaluations reveal that DeepSeek-V3 outperforms other open-source models and achieves performance comparable to leading closed-source models. Despite its excellent performance, DeepSeek-V3 requires only 2.788M H800 GPU hours for its full training. In addition, its training process is remarkably stable. Throughout the entire training process, we did not experience any irrecoverable loss spikes or perform any rollbacks. The model checkpoints are available at [this https URL](https://arxiv.org/abs/2412.19437).

Link: <https://arxiv.org/abs/2412.19437>

DeepSeek-Coder-V2: Breaking the Barrier of Closed-Source Models in Code Intelligence

Authors: DeepSeek-AI, Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Yifan Wu, Yukun Li, Huazuo Gao, Shirong Ma, Wangding Zeng, Xiao Bi, Zihui Gu, Hanwei Xu, Damai Dai, Kai Dong, Liyue Zhang, Yishi Piao, Zhibin Gou, Zhenda Xie, Zhewen Hao, Bingxuan Wang, Junxiao Song, Deli Chen, Xin Xie, Kang Guan, Yuxiang You, Aixin Liu, Qiushi Du, Wenjun Gao, Xuan Lu, Qinyu Chen, Yaohui Wang, Chengqi Deng, Jiashi Li, Chenggang Zhao, Chong Ruan, Fuli Luo, Wenfeng Liang
Publication Date: June 17, 2024

Abstract: We present DeepSeek-Coder-V2, an open-source Mixture-of-Experts (MoE) code language model that achieves performance comparable to GPT4-Turbo in code-specific tasks. Specifically, DeepSeek-Coder-V2 is further pre-trained from an intermediate checkpoint of DeepSeek-V2 with additional 6 trillion tokens. Through this continued pre-training, DeepSeek-Coder-V2 substantially enhances the coding and mathematical reasoning capabilities of DeepSeek-V2, while maintaining comparable performance in general language tasks. Compared to DeepSeek-Coder-33B, DeepSeek-Coder-V2 demonstrates significant advancements in various aspects of code-related tasks, as well as reasoning and general capabilities. Additionally, DeepSeek-Coder-V2 expands its support for programming languages from 86 to 338, while extending the context length from 16K to 128K. In standard benchmark evaluations, DeepSeek-Coder-V2 achieves superior performance compared to closed-source models such as GPT4-Turbo, Claude 3 Opus, and Gemini 1.5 Pro in coding and math benchmarks.

Link: <https://arxiv.org/abs/2406.11931>

DeepSeek-Prover-V1.5: Harnessing Proof Assistant Feedback for Reinforcement Learning and Monte-Carlo Tree Search

Authors: Huajian Xin, Z.Z. Ren, Junxiao Song, Zhihong Shao, Wanjia Zhao, Haocheng Wang, Bo Liu, Liyue Zhang, Xuan Lu, Qiushi Du, Wenjun Gao, Qihao Zhu, Dejian Yang, Zhibin Gou, Z.F. Wu, Fuli Luo, Chong Ruan
Publication Date: August 15, 2024

Abstract: We introduce DeepSeek-Prover-V1.5, an open-source language model designed for theorem proving in Lean 4, which enhances DeepSeek-Prover-V1 by optimizing both training and inference processes. Pre-trained on DeepSeekMath-Base with specialization in formal mathematical languages, the model undergoes supervised fine-tuning using an enhanced formal theorem proving dataset derived from DeepSeek-Prover-V1. Further refinement is achieved through reinforcement learning from proof assistant feedback (RLPAF). Beyond the single-pass whole-proof generation approach of DeepSeek-Prover-V1, we propose RMaxTS, a variant of Monte-Carlo tree search that employs an intrinsic-reward-driven exploration strategy to generate diverse proof paths. DeepSeek-Prover-V1.5 demonstrates significant improvements over DeepSeek-Prover-V1, achieving new state-of-the-art results on the test set of the high school level miniF2F benchmark (63.5%) and the undergraduate level ProofNet benchmark (25.3%).

Link: <https://arxiv.org/abs/2408.08152>

DeepSeek-Prover-V1.5: Harnessing Proof Assistant Feedback for Reinforcement Learning and Monte-Carlo Tree Search

Authors: DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z.F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J.L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R.J. Chen, R.L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S.S. Li et al. **(100 additional authors not shown)**

Publication Date: January 22, 2025

Abstract: We introduce our first-generation reasoning models, DeepSeek-R1-Zero and DeepSeek-R1. DeepSeek-R1-Zero, a model trained via large-scale reinforcement learning (RL) without supervised fine-tuning (SFT) as a preliminary step, demonstrates remarkable reasoning capabilities. Through RL, DeepSeek-R1-Zero naturally emerges with numerous powerful and intriguing reasoning behaviors. However, it encounters challenges such as poor readability, and language mixing. To address these issues and further enhance reasoning performance, we introduce DeepSeek-R1, which incorporates multi-stage training and cold-start data before RL. DeepSeek-R1 achieves performance comparable to OpenAI-o1-1217 on reasoning tasks. To support the research community, we open-source DeepSeek-R1-Zero, DeepSeek-R1, and six dense models (1.5B, 7B, 8B, 14B, 32B, 70B) distilled from DeepSeek-R1 based on Qwen and Llama.

Link: <https://arxiv.org/abs/2501.12948>