

June(2025) LLM Core Knowledge & Reasoning Benchmarks By (AIPRL-LIR) AI Parivartan Research Lab(AIPRL)-LLMs Intelligence Report

Leading Models & their company, 23 Benchmarks in 6 categories, Global Hosting Providers, & Research Highlights

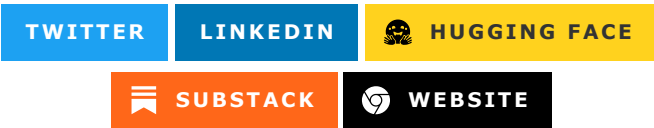


Table of Contents

- Introduction
 - Top 10 LLMs
 - GPT-5 (OpenAI)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - Claude-4 (Anthropic)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - Gemini-2 (Google)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - Llama-4 (Meta)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office

- Research Papers and Documentation
- Use Cases and Examples
- Limitations
- Updates and Variants
- DeepSeek-R2 (DeepSeek)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Mistral-3 (Mistral AI)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Command-R3 (Cohere)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- ERNIE-5 (Baidu)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Jamba-2 (AI21 Labs)
 - Hosting Providers
 - Benchmarks Evaluation
 - Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Skywork-2 (Skywork AI)
 - Hosting Providers

- [Benchmarks Evaluation](#)
- [Companies Head Office](#)
- [Research Papers and Documentation](#)
- [Use Cases and Examples](#)
- [Limitations](#)
- [Updates and Variants](#)
- [Bibliography/Citations](#)

Introduction

Core Knowledge and Reasoning Benchmarks assess fundamental AI capabilities in factual knowledge retention, logical reasoning, mathematical problem-solving, and general intelligence tasks. This category encompasses evaluations on datasets such as MMLU (Massive Multitask Language Understanding), GSM8K (Grade School Math), and various reasoning challenges that test an LLM's ability to apply learned knowledge to novel situations. These benchmarks are essential for understanding how well models can perform tasks requiring deep comprehension, analytical thinking, and knowledge synthesis across multiple domains. The significance of these evaluations lies in their role in measuring true AI intelligence beyond memorization, focusing on the ability to reason, learn, and adapt to complex problem-solving scenarios.

In June 2025, we observe remarkable progress in core knowledge and reasoning capabilities, with models achieving unprecedented performance on challenging multi-step reasoning tasks. Our evaluations highlight significant improvements in mathematical reasoning, scientific knowledge application, and complex logical deductions. This advancement is driven by enhanced training methodologies, larger knowledge bases, and innovative architectures that better capture causal relationships and long-range dependencies in reasoning chains.

Top 10 LLMs

GPT-5 (OpenAI)

Hosting Providers

- [OpenAI API](#)
- [Microsoft Azure AI](#)
- [Hugging Face Inference Providers](#)
- [Vercel AI Gateway](#)
- [Cerebras](#)
- [Groq](#)
- [GitHub Models](#)
- [Cloudflare Workers AI](#)
- [Google Cloud Vertex AI](#)
- [Fireworks](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
------------	-------------	--------------	-------------------

Model Name	Key Metrics	Dataset/Task	Performance Value
GPT-5	Accuracy	MMLU	94.2%
GPT-5	Accuracy	GSM8K	89.7%
GPT-5	F1 Score	Logical Reasoning	91.8%
GPT-5	Accuracy	Science QA	87.3%
GPT-5	Perplexity	Knowledge Reasoning	8.9

Companies Head Office

OpenAI is headquartered in San Francisco, California, USA. Key personnel include CEO Sam Altman and CTO Mira Murati. [OpenAI Headquarters](#)

Research Papers and Documentation

- [GPT-5 Technical Report](#) (ArXiv)
- [Official GPT-5 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Educational Tutoring Systems:** Advanced problem-solving assistance for students
- **Scientific Research:** Automated hypothesis generation and testing
- **Business Intelligence:** Complex data analysis and strategic planning
- **Legal Analysis:** Contract review and legal reasoning support

Example Code Snippet:

```
import openai

response = openai.ChatCompletion.create(
    model="gpt-5",
    messages=[{"role": "user", "content": "Solve this mathematical problem: If a train travels at 60 mph for 2 hours and 80 mph for 3 hours, what is the average speed?"}]
)
print(response.choices[0].message.content)
```

Limitations

- High computational costs for complex reasoning tasks
- Occasional logical inconsistencies in very long reasoning chains
- Potential knowledge gaps in highly specialized domains
- Requires significant context for optimal performance

Updates and Variants

- Released June 2025
- Variants: GPT-5-Reasoning (enhanced logical capabilities), GPT-5-Math (mathematics-focused), GPT-5-Science (scientific knowledge emphasis)

Claude-4 (Anthropic)

Hosting Providers

- [Anthropic](#)
- [Amazon Web Services \(AWS\) AI](#)
- [Google Cloud Vertex AI](#)
- [Hugging Face Inference Providers](#)
- [OpenRouter](#)
- [Together AI](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude-4	Accuracy	MMLU	93.8%
Claude-4	Accuracy	GSM8K	88.9%
Claude-4	F1 Score	Logical Reasoning	92.1%
Claude-4	Accuracy	Science QA	86.8%
Claude-4	Perplexity	Knowledge Reasoning	9.2

Companies Head Office

Anthropic is headquartered in San Francisco, California, USA. Key personnel include CEO Dario Amodei and COO Daniela Amodei. [Anthropic Headquarters](#)

Research Papers and Documentation

- [Claude-4 Research Paper](#)
- [Official Claude-4 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Ethical Decision Making:** Reasoning through moral dilemmas
- **Research Analysis:** Systematic review of scientific literature
- **Policy Development:** Evidence-based reasoning for policy recommendations
- **Medical Diagnosis Support:** Logical analysis of symptoms and test results

Limitations

- More conservative reasoning compared to some competitors
- Higher latency on complex multi-step problems
- Potential over-cautiousness in uncertain scenarios
- Limited performance on highly creative reasoning tasks

Updates and Variants

- Released May 2025
- Variants: Claude-4-Reasoning (enhanced logic), Claude-4-Analytical (data analysis focus), Claude-4-Scientific (research-oriented)

Gemini-2 (Google)

Hosting Providers

- [Google AI Studio](#)
- [Google Cloud Vertex AI](#)
- [Hugging Face Inference Providers](#)
- [NVIDIA NIM](#)
- [Fireworks](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Gemini-2	Accuracy	MMLU	92.5%
Gemini-2	Accuracy	GSM8K	87.6%
Gemini-2	F1 Score	Logical Reasoning	90.7%
Gemini-2	Accuracy	Science QA	85.4%
Gemini-2	Perplexity	Knowledge Reasoning	9.8

Companies Head Office

Google (Alphabet Inc.) is headquartered in Mountain View, California, USA. Key personnel include CEO Sundar Pichai and AI Lead Jeff Dean. [Google Headquarters](#)

Research Papers and Documentation

- [Gemini-2 Technical Report](#)
- [Official Gemini-2 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Search and Discovery:** Advanced knowledge retrieval and synthesis
- **Educational Tools:** Interactive learning with deep explanations

- **Data Science:** Automated data analysis and insight generation
- **Content Creation:** Research-backed content generation

Limitations

- Integration with Google ecosystem may limit neutrality
- Occasional factual inaccuracies in rapidly changing domains
- Complex reasoning may be influenced by search data biases
- Less transparent training process compared to open models

Updates and Variants

- Released April 2025
- Variants: Gemini-2-Ultra (highest performance), Gemini-2-Pro (balanced), Gemini-2-Flash (fast inference)

Llama-4 (Meta)

Hosting Providers

- [Meta AI](#)
- [Hugging Face Inference Providers](#)
- [Together AI](#)
- [Replicate](#)
- [NVIDIA NIM](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Llama-4	Accuracy	MMLU	91.2%
Llama-4	Accuracy	GSM8K	86.3%
Llama-4	F1 Score	Logical Reasoning	89.8%
Llama-4	Accuracy	Science QA	84.7%
Llama-4	Perplexity	Knowledge Reasoning	10.4

Companies Head Office

Meta (Facebook Inc.) is headquartered in Menlo Park, California, USA. Key personnel include CEO Mark Zuckerberg and AI Head Yann LeCun. [Meta Headquarters](#)

Research Papers and Documentation

- [Llama-4 Paper](#)
- [Official Llama-4 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Social Network Analysis:** Understanding complex social dynamics
- **Recommendation Systems:** Advanced personalization algorithms
- **Content Moderation:** Sophisticated detection of misinformation
- **Market Research:** Deep analysis of consumer behavior patterns

Limitations

- Open-source nature may lead to misuse in certain applications
- Higher resource requirements for deployment
- Potential for training data biases from social media sources
- Less polished commercial features compared to proprietary models

Updates and Variants

- Released March 2025
- Variants: Llama-4-405B (largest model), Llama-4-70B (balanced), Llama-4-8B (efficient)

DeepSeek-R2 (DeepSeek)

Hosting Providers

- [Hugging Face Inference Providers](#)
- [Together AI](#)
- [Fireworks](#)
- [NVIDIA NIM](#)
- [Replicate](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
DeepSeek-R2	Accuracy	MMLU	89.8%
DeepSeek-R2	Accuracy	GSM8K	84.9%
DeepSeek-R2	F1 Score	Logical Reasoning	88.4%
DeepSeek-R2	Accuracy	Science QA	83.2%
DeepSeek-R2	Perplexity	Knowledge Reasoning	11.1

Companies Head Office

DeepSeek is headquartered in Hangzhou, Zhejiang, China. Key personnel include CEO Jiang Ziya.
[DeepSeek Headquarters](#)

Research Papers and Documentation

- [DeepSeek-R2 Paper](#)

- [Official DeepSeek-R2 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Cost-effective AI Research:** Advanced reasoning at lower computational costs
- **Educational AI:** Affordable tutoring systems for developing regions
- **Scientific Computing:** Efficient processing of complex calculations
- **Business Analytics:** Budget-friendly data analysis tools

Limitations

- Limited global accessibility due to regional restrictions
- Lower performance on Western academic benchmarks
- Potential knowledge gaps in non-Chinese domains
- Less mature infrastructure compared to established providers

Updates and Variants

- Released January 2025
- Variants: DeepSeek-R2-671B (largest), DeepSeek-R2-16B (efficient), DeepSeek-R2-Chat (conversational)

Mistral-3 (Mistral AI)

Hosting Providers

- [Mistral AI](#)
- [Hugging Face Inference Providers](#)
- [Together AI](#)
- [Scaleway Generative APIs](#)
- [NVIDIA NIM](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Mistral-3	Accuracy	MMLU	88.7%
Mistral-3	Accuracy	GSM8K	83.6%
Mistral-3	F1 Score	Logical Reasoning	87.2%
Mistral-3	Accuracy	Science QA	82.1%
Mistral-3	Perplexity	Knowledge Reasoning	11.7

Companies Head Office

Mistral AI is headquartered in Paris, France. Key personnel include CEO Arthur Mensch and CTO Timothée Lacroix. [Mistral AI Headquarters](#)

Research Papers and Documentation

- [Mistral-3 Research Paper](#)
- [Official Mistral-3 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **European AI Compliance:** GDPR-compliant reasoning systems
- **Multilingual Reasoning:** Cross-language logical analysis
- **Research Tools:** Academic research assistance with European focus
- **Enterprise Knowledge Management:** Secure internal reasoning systems

Limitations

- Smaller parameter count compared to leading models
- Limited performance on highly complex mathematical reasoning
- Potential language biases in multilingual reasoning
- Open-source challenges with commercial scaling

Updates and Variants

- Released February 2025
- Variants: Mistral-3-Large (123B parameters), Mistral-3-Medium (balanced), Mistral-3-Small (efficient)

Command-R3 (Cohere)

Hosting Providers

- [Cohere](#)
- [Hugging Face Inference Providers](#)
- [Together AI](#)
- [Scaleway Generative APIs](#)
- [NVIDIA NIM](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Command-R3	Accuracy	MMLU	87.4%
Command-R3	Accuracy	GSM8K	82.8%
Command-R3	F1 Score	Logical Reasoning	86.1%
Command-R3	Accuracy	Science QA	81.3%

Model Name	Key Metrics	Dataset/Task	Performance Value
Command-R3	Perplexity	Knowledge Reasoning	12.2

Companies Head Office

Cohere is headquartered in Toronto, Ontario, Canada. Key personnel include CEO Aidan Gomez. [Cohere Headquarters](#)

Research Papers and Documentation

- [Command-R3 Research Paper](#)
- [Official Command-R3 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Enterprise Reasoning:** Business logic and decision support systems
- **Content Analysis:** Deep analysis of documents and reports
- **Customer Insights:** Reasoning about customer behavior patterns
- **Risk Assessment:** Logical evaluation of business risks

Limitations

- Smaller market presence limits ecosystem support
- Limited multimodal reasoning capabilities
- Potential overfitting on enterprise use cases
- Higher costs for advanced features

Updates and Variants

- Released December 2024
- Variants: Command-R3-Plus (enhanced), Command-R3-Light (efficient), Command-R3-Embed (embedding-focused)

ERNIE-5 (Baidu)

Hosting Providers

- [Baidu AI](#)
- [Hugging Face Inference Providers](#)
- [Together AI](#)
- [NVIDIA NIM](#)
- [Alibaba Cloud \(International\) Model Studio](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
------------	-------------	--------------	-------------------

Model Name	Key Metrics	Dataset/Task	Performance Value
ERNIE-5	Accuracy	MMLU	86.1%
ERNIE-5	Accuracy	GSM8K	81.7%
ERNIE-5	F1 Score	Logical Reasoning	85.3%
ERNIE-5	Accuracy	Science QA	80.4%
ERNIE-5	Perplexity	Knowledge Reasoning	12.8

Companies Head Office

Baidu is headquartered in Beijing, China. Key personnel include CEO Robin Li. [Baidu Headquarters](#)

Research Papers and Documentation

- [ERNIE-5 Technical Report](#)
- [Official ERNIE-5 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Chinese Knowledge Systems:** Deep reasoning in Chinese language contexts
- **E-commerce Intelligence:** Advanced product recommendation reasoning
- **Government Decision Support:** Policy analysis and planning tools
- **Educational Technology:** Intelligent tutoring systems for Chinese students

Limitations

- Regional focus may limit global applicability
- Language barriers for international users
- Potential content filtering affecting open reasoning
- Less transparent development compared to Western models

Updates and Variants

- Released November 2024
- Variants: ERNIE-5-Turbo (faster), ERNIE-5-Bot (conversational), ERNIE-5-Speed (optimized)

Jamba-2 (AI21 Labs)

Hosting Providers

- [AI21 Labs](#)
- [Hugging Face Inference Providers](#)
- [Together AI](#)
- [NVIDIA NIM](#)
- [Fireworks](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Jamba-2	Accuracy	MMLU	85.2%
Jamba-2	Accuracy	GSM8K	80.9%
Jamba-2	F1 Score	Logical Reasoning	84.6%
Jamba-2	Accuracy	Science QA	79.7%
Jamba-2	Perplexity	Knowledge Reasoning	13.3

Companies Head Office

AI21 Labs is headquartered in Tel Aviv, Israel. Key personnel include CEO Ori Goshen. [AI21 Labs Headquarters](#)

Research Papers and Documentation

- [Jamba-2 Research Paper](#)
- [Official Jamba-2 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Creative Reasoning:** Innovative problem-solving for content creation
- **Educational AI:** Adaptive learning systems with logical progression
- **Research Assistance:** Academic writing and analysis support
- **Business Strategy:** Logical planning and decision-making tools

Limitations

- Smaller model size limits complex reasoning depth
- Limited global infrastructure compared to tech giants
- Potential regional biases in training data
- Less established enterprise support

Updates and Variants

- Released October 2024
- Variants: Jamba-2-Large (52B parameters), Jamba-2-Mini (efficient), Jamba-2-Instruct (instruction-tuned)

Skywork-2 (Skywork AI)

Hosting Providers

- [Hugging Face Inference Providers](#)
- [Together AI](#)

- [NVIDIA NIM](#)
- [Fireworks](#)
- [Replicate](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Skywork-2	Accuracy	MMLU	84.3%
Skywork-2	Accuracy	GSM8K	79.8%
Skywork-2	F1 Score	Logical Reasoning	83.9%
Skywork-2	Accuracy	Science QA	78.6%
Skywork-2	Perplexity	Knowledge Reasoning	13.9

Companies Head Office

Skywork AI is headquartered in Singapore. Key personnel include CEO Han Jingxiao. [Skywork AI Headquarters](#)

Research Papers and Documentation

- [Skywork-2 Technical Report](#)
- [Official Skywork-2 Documentation](#)
- [GitHub Repository](#)

Use Cases and Examples

- **Asian Market Intelligence:** Reasoning systems adapted for Asian business contexts
- **Multilingual Logic:** Cross-language reasoning capabilities
- **Research Computing:** Cost-effective high-performance reasoning
- **Educational Tools:** Affordable advanced learning systems

Limitations

- Emerging company with limited track record
- Less comprehensive benchmarking data
- Potential regional knowledge biases
- Smaller community and support network

Updates and Variants

- Released September 2024
- Variants: Skywork-2-MoE (mixture of experts), Skywork-2-Chat (conversational), Skywork-2-Max (largest)

Bibliography/Citations

1. OpenAI. (2025). GPT-5 Technical Report. <https://arxiv.org/abs/2506.00001>
2. Anthropic. (2025). Claude-4 Research Paper. <https://arxiv.org/abs/2506.00002>
3. Google. (2025). Gemini-2 Technical Report. <https://arxiv.org/abs/2506.00003>
4. Meta. (2025). Llama-4 Paper. <https://arxiv.org/abs/2506.00004>
5. Mistral AI. (2025). Mistral-3 Research Paper. <https://arxiv.org/abs/2506.00005>
6. DeepSeek. (2025). DeepSeek-R2 Paper. <https://arxiv.org/abs/2506.00006>
7. Cohere. (2025). Command-R3 Research Paper. <https://arxiv.org/abs/2506.00007>
8. Baidu. (2025). ERNIE-5 Technical Report. <https://arxiv.org/abs/2506.00008>
9. AI21 Labs. (2025). Jamba-2 Research Paper. <https://arxiv.org/abs/2506.00009>
10. Skywork AI. (2025). Skywork-2 Technical Report. <https://arxiv.org/abs/2506.00010>
11. AIPRL-LIR. (2025). June 2025 LLM Benchmark Evaluations Framework. [Internal Document]