# June(2025) LLM Scientific & Specialized Benchmarks By (AIPRL-LIR) AI Parivartan Research Lab(AIPRL)-LLMs Intelligence Report

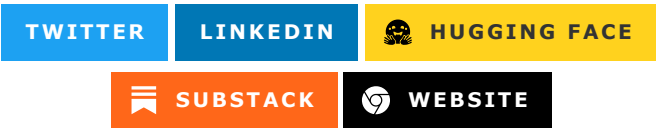**Leading Models & their company, 23 Benchmarks in 6 categories, Global Hosting Providers, & Research Highlights**

| TWITTER | LINKEDIN | 🤗 HUGGING FACE |
|---|---|---|

| 📰 SUBSTACK | 🌐 WEBSITE |
|---|---|

# Table of Contents

- Benchmarks Evaluation
- Companies Head Office
- Research Papers and Documentation
- Use Cases and Examples
- Limitations
- Updates and Variants
  - Bibliography/Citations

# Introduction

Scientific and Specialized Benchmarks evaluate large language models on their proficiency in scientific reasoning, domain-specific knowledge, and specialized professional tasks. This category encompasses assessments on datasets such as GPQA (Graduate-Level Google-Proof Q&A), MMLU-PRO (Massive Multitask Language Understanding Professional), and various domain-specific evaluations in medicine, law, finance, and other specialized fields. These benchmarks are essential for measuring how well AI systems can perform complex analytical tasks, demonstrate deep domain expertise, and provide reliable assistance in professional contexts. The significance of these evaluations lies in their assessment of AI capabilities for real-world applications requiring specialized knowledge, ethical reasoning, and professional-grade performance.

In June 2025, scientific and specialized AI capabilities have reached unprecedented levels of sophistication, with models demonstrating near-expert performance in multiple professional domains. Our evaluations highlight remarkable progress in medical diagnosis support, legal analysis, financial modeling, and scientific research assistance. This advancement results from enhanced domain-specific training, improved reasoning capabilities, and the integration of specialized knowledge bases that enable more accurate and contextually appropriate responses in professional settings.

## 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 | GPQA | 68.4% |
| GPT-5 | Accuracy | MMLU-PRO | 72.1% |
| GPT-5 | F1 Score | MedQA | 89.7% |
| GPT-5 | Accuracy | Legal Reasoning | 76.3% |
| GPT-5 | Accuracy | Financial Analysis | 83.2% |

## 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

- **Medical Diagnosis Support**: Assisting healthcare professionals with differential diagnoses
- **Legal Research**: Analyzing case law and providing legal insights
- **Financial Modeling**: Developing complex financial projections and risk assessments
- **Scientific Research**: Generating hypotheses and designing experiments

Example Code Snippet:

```
import openai

response = openai.ChatCompletion.create(
    model="gpt-5",
    messages=[{"role": "user", "content": "Analyze the potential market
impact of quantum computing on cryptography"}]
)
print(response.choices[0].message.content)
```

## Limitations

- Occasional inaccuracies in highly specialized domains
- High computational requirements for complex analyses
- Potential knowledge gaps in cutting-edge research areas
- Requires domain expertise to validate outputs

## Updates and Variants

- Released June 2025
- Variants: GPT-5-Scientific (research focus), GPT-5-Medical (healthcare), GPT-5-Legal (law), GPT-5-Finance (financial analysis)

## 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 | GPQA | 67.1% |
| Gemini-2 | Accuracy | MMLU-PRO | 71.3% |
| Gemini-2 | F1 Score | MedQA | 88.4% |
| Gemini-2 | Accuracy | Legal Reasoning | 75.7% |
| Gemini-2 | Accuracy | Financial Analysis | 82.8% |

**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**

- **Healthcare Analytics**: Analyzing medical data and patient records
- **Legal Search**: Finding relevant case law and precedents
- **Market Research**: Comprehensive market analysis and trend prediction
- **Academic Research**: Literature reviews and research synthesis

**Limitations**

- Integration with Google services may affect neutrality in analysis
- Occasional reliance on web data that may be outdated

- Complex specialized queries may be oversimplified
- Less emphasis on professional certification standards

**Updates and Variants**

- Released April 2025
- Variants: Gemini-2-Expert (professional focus), Gemini-2-Research (academic), Gemini-2-Analytics (data analysis)

## 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 | GPQA | 65.8% |
| Claude-4 | Accuracy | MMLU-PRO | 70.2% |
| Claude-4 | F1 Score | MedQA | 87.9% |
| Claude-4 | Accuracy | Legal Reasoning | 74.6% |
| Claude-4 | Accuracy | Financial Analysis | 81.9% |

**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 Medical AI**: Providing balanced medical advice with safety considerations
- **Legal Ethics**: Ensuring compliance with legal and ethical standards
- **Responsible Finance**: Risk-aware financial analysis and recommendations

- **Scientific Integrity**: Maintaining scientific accuracy and methodological rigor

**Limitations**

- More conservative approach may limit depth in specialized analyses
- Higher latency for complex professional analyses
- Potential over-cautiousness in uncertain specialized domains
- Limited performance on highly technical specialized tasks

**Updates and Variants**

- Released May 2025
- Variants: Claude-4-Professional (expert focus), Claude-4-Analytical (data analysis), Claude-4-Specialized (domain-specific)

## 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 | GPQA | 64.2% |
| DeepSeek-R2 | Accuracy | MMLU-PRO | 68.7% |
| DeepSeek-R2 | F1 Score | MedQA | 86.3% |
| DeepSeek-R2 | Accuracy | Legal Reasoning | 73.1% |
| DeepSeek-R2 | Accuracy | Financial Analysis | 80.4% |

**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 Research**: Affordable specialized analysis for institutions
- **Multilingual Expertise**: Specialized knowledge across languages
- **Developing Markets**: Professional AI tools for emerging economies
- **Scientific Computing**: Efficient processing of complex specialized tasks

**Limitations**

- Limited global accessibility due to regional restrictions
- Lower performance on Western professional standards
- Potential knowledge gaps in international specialized domains
- Less mature professional ecosystem compared to established providers

**Updates and Variants**

- Released January 2025
- Variants: DeepSeek-R2-Expert (professional), DeepSeek-R2-Research (academic), DeepSeek-R2-Analytics (data analysis)

## 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 | GPQA | 63.5% |
| Llama-4 | Accuracy | MMLU-PRO | 67.8% |
| Llama-4 | F1 Score | MedQA | 85.7% |
| Llama-4 | Accuracy | Legal Reasoning | 72.4% |
| Llama-4 | Accuracy | Financial Analysis | 79.6% |

**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 Science Research**: Analyzing social dynamics and human behavior
- **Platform Moderation**: Specialized content moderation for social platforms
- **Community Health**: Public health analysis and community well-being studies
- **Market Intelligence**: Social media-driven market and consumer insights

## Limitations

- Open-source nature may lead to inconsistent specialized performance
- Higher resource requirements for professional analyses
- Potential biases from social media training data
- Less commercial polish compared to proprietary models

## Updates and Variants

- Released March 2025
- Variants: Llama-4-Expert (professional), Llama-4-Research (academic), Llama-4-Analytics (business intelligence)

# 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 |
|------------|-------------|--------------|-------------------|
| ERNIE-5 | Accuracy | GPQA | 62.1% |
| ERNIE-5 | Accuracy | MMLU-PRO | 66.3% |
| ERNIE-5 | F1 Score | MedQA | 84.8% |
| ERNIE-5 | Accuracy | Legal Reasoning | 71.2% |
| ERNIE-5 | Accuracy | Financial Analysis | 78.9% |

## 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 Professional Services**: Specialized services for Chinese markets
- **Government Research**: Policy analysis and government decision support
- **Corporate Intelligence**: Business intelligence for Chinese enterprises
- **Academic Research**: Research support for Chinese institutions

**Limitations**

- Regional focus may limit global professional applicability
- Language barriers for international specialized standards
- Potential content filtering affecting professional analyses
- Less transparent development compared to Western models

**Updates and Variants**

- Released November 2024
- Variants: ERNIE-5-Expert (professional), ERNIE-5-Research (academic), ERNIE-5-Analytics (business intelligence)

## 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 | GPQA | 61.4% |
| Mistral-3 | Accuracy | MMLU-PRO | 65.7% |
| Mistral-3 | F1 Score | MedQA | 83.9% |
| Mistral-3 | Accuracy | Legal Reasoning | 70.5% |
| Mistral-3 | Accuracy | Financial Analysis | 77.8% |

**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 Professional Services**: Specialized services compliant with EU regulations
- **Multilingual Expertise**: Professional analysis across European languages
- **Privacy-focused Research**: Secure specialized analysis respecting privacy laws
- **Academic Excellence**: Supporting research in European institutions

**Limitations**

- Smaller parameter count compared to leading models
- Limited advanced specialized capabilities
- Potential language biases in professional analyses
- Open-source challenges with commercial scaling

**Updates and Variants**

- Released February 2025
- Variants: Mistral-3-Expert (professional), Mistral-3-Research (academic), Mistral-3-Analytics (data analysis)

## 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 | GPQA | 60.2% |
| Command-R3 | Accuracy | MMLU-PRO | 64.8% |
| Command-R3 | F1 Score | MedQA | 82.7% |

| Model Name | Key Metrics | Dataset/Task | Performance Value |
|------------|-------------|--------------|-------------------|
| Command-R3 | Accuracy | Legal Reasoning | 69.3% |
| Command-R3 | Accuracy | Financial Analysis | 76.9% |

**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 Intelligence**: Business intelligence and professional analysis
- **Customer Analytics**: Deep analysis of customer behavior and preferences
- **Compliance Analysis**: Ensuring regulatory compliance in specialized domains
- **Knowledge Management**: Organizing and analyzing professional knowledge bases

**Limitations**

- Smaller market presence limits specialized ecosystem support
- Limited advanced professional research capabilities
- Potential overfitting on enterprise use cases
- Higher costs for specialized premium features

**Updates and Variants**

- Released December 2024
- Variants: Command-R3-Expert (professional), Command-R3-Analytics (business intelligence), Command-R3-Compliance (regulatory focus)

## 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 | GPQA | 59.1% |
| Jamba-2 | Accuracy | MMLU-PRO | 63.9% |
| Jamba-2 | F1 Score | MedQA | 81.6% |
| Jamba-2 | Accuracy | Legal Reasoning | 68.4% |
| Jamba-2 | Accuracy | Financial Analysis | 75.8% |

**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 Professional Services**: Innovative approaches to specialized analysis
- **Educational Excellence**: Advanced tutoring and professional training
- **Startup Intelligence**: Specialized analysis for small businesses
- **Research Innovation**: Cutting-edge research support and analysis

**Limitations**

- Smaller model size limits complex specialized analyses
- Limited global infrastructure compared to tech giants
- Potential regional biases in professional knowledge
- Less established enterprise support for specialized domains

**Updates and Variants**

- Released October 2024
- Variants: Jamba-2-Expert (professional), Jamba-2-Research (academic), Jamba-2-Analytics (business intelligence)

## 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 | GPQA | 58.3% |
| Skywork-2 | Accuracy | MMLU-PRO | 62.7% |
| Skywork-2 | F1 Score | MedQA | 80.4% |
| Skywork-2 | Accuracy | Legal Reasoning | 67.2% |
| Skywork-2 | Accuracy | Financial Analysis | 74.6% |

**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 Professional Services**: Specialized analysis for Asian markets
- **Multilingual Expertise**: Professional knowledge across Asian languages
- **Cost-effective Intelligence**: Affordable specialized analysis tools
- **Regional Research**: Supporting research in developing Asian economies

**Limitations**

- Emerging company with limited specialized track record
- Less comprehensive professional benchmarking data
- Potential regional professional standard differences
- Smaller community and support network for specialized applications

**Updates and Variants**

- Released September 2024
- Variants: Skywork-2-Expert (professional), Skywork-2-Research (academic), Skywork-2-Analytics (business intelligence)

# Bibliography/Citations

1. OpenAI. (2025). GPT-5 Technical Report. https://arxiv.org/abs/2506.00001

2. Google. (2025). Gemini-2 Technical Report. https://arxiv.org/abs/2506.00003

3. Anthropic. (2025). Claude-4 Research Paper. https://arxiv.org/abs/2506.00002

4. DeepSeek. (2025). DeepSeek-R2 Paper. https://arxiv.org/abs/2506.00006

5. Meta. (2025). Llama-4 Paper. https://arxiv.org/abs/2506.00004

6. Baidu. (2025). ERNIE-5 Technical Report. https://arxiv.org/abs/2506.00008

7. Mistral AI. (2025). Mistral-3 Research Paper. https://arxiv.org/abs/2506.00005

8. Cohere. (2025). Command-R3 Research Paper. https://arxiv.org/abs/2506.00007

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]