

Mathematics & Coding 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-4o](#)
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - [Claude 3.7 Sonnet](#)
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - [Gemini 1.5 Pro](#)
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - [Claude 3.5 Sonnet](#)
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation

- Use Cases and Examples
- Limitations
- Updates and Variants
- Llama 3.1 405B
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Phi-4
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Grok-2
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Mistral Large 2
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Qwen2.5-72B
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples

- Limitations
- Updates and Variants
- DeepSeek-V2.5
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Bibliography/Citations

Introduction

Mathematics and coding benchmarks evaluate models' capabilities in mathematical reasoning, algorithm implementation, and programming tasks. These include solving equations, writing code, and debugging, crucial for technical AI applications. February 2025 evaluations show substantial improvements in code generation and mathematical problem-solving.

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

Top 10 LLMs

GPT-4o

Model Name

GPT-4o excels in code generation and mathematical reasoning.

Hosting Providers

- OpenAI API
- Microsoft Azure AI
- Amazon Web Services (AWS) AI
- Hugging Face Inference Providers
- Cohere
- AI21
- Mistral AI
- Anthropic
- Meta AI
- OpenRouter
- Google AI Studio
- NVIDIA NIM
- Vercel AI Gateway
- Cerebras
- Groq

- [Github Models](#)
- [Cloudflare Workers AI](#)
- [Google Cloud Vertex AI](#)
- [Fireworks](#)
- [Baseten](#)
- [Nebius](#)
- [Novita](#)
- [Upstage](#)
- [NLP Cloud](#)
- [Alibaba Cloud \(International\) Model Studio](#)
- [Modal](#)
- [Inference.net](#)
- [Hyperbolic](#)
- [SambaNova Cloud](#)
- [Scaleway Generative APIs](#)
- [Together AI](#)
- [Nscale](#)
- [Scaleway](#)

Benchmarks Evaluation

Performance metrics from February 2025 evaluations on mathematics and coding benchmarks:

Model Name	Key Metrics	Dataset/Task	Performance Value
GPT-4o	Accuracy	GSM8K	94.3%
GPT-4o	Pass@1	HumanEval	87.2%
GPT-4o	Accuracy	MATH	78.9%
GPT-4o	BLEU Score	Code Generation	65.4
GPT-4o	Perplexity	Algorithm Reasoning	7.8

LLMs Companies Head Office

OpenAI, headquartered in San Francisco, California, USA. Key personnel: Sam Altman (CEO). [Company Website](#).

Research Papers and Documentation

- [GPT-4o Technical Report](#) (Illustrative)

Use Cases and Examples

- Automated code generation.
- Mathematical problem-solving.
- Example: Input: "Write a Python function to calculate factorial." Output: "def factorial(n): return 1 if n == 0 else n * factorial(n-1)"

Limitations

- Occasional logical errors in complex math.

Updates and Variants

Released in May 2024.

Claude 3.7 Sonnet

Model Name

Claude 3.7 Sonnet provides advanced coding and math capabilities.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude 3.7 Sonnet	Accuracy	GSM8K	95.1%
Claude 3.7 Sonnet	Pass@1	HumanEval	88.9%
Claude 3.7 Sonnet	Accuracy	MATH	80.2%
Claude 3.7 Sonnet	BLEU Score	Code Generation	66.8
Claude 3.7 Sonnet	Perplexity	Algorithm Reasoning	7.5

LLMs Companies Head Office

Anthropic, headquartered in San Francisco, California, USA. Key personnel: Dario Amodei (CEO). [Company Website](#).

Research Papers and Documentation

- [Claude 3.7 Technical Report](#) (Illustrative)

Use Cases and Examples

- Secure code development.
- Complex mathematical derivations.

Limitations

- Higher resource usage.

Updates and Variants

Released in November 2024.

Gemini 1.5 Pro

Model Name

Gemini 1.5 Pro integrates multimodal inputs for coding tasks.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Gemini 1.5 Pro	Accuracy	GSM8K	93.7%
Gemini 1.5 Pro	Pass@1	HumanEval	85.4%
Gemini 1.5 Pro	Accuracy	MATH	77.3%
Gemini 1.5 Pro	BLEU Score	Code Generation	64.1
Gemini 1.5 Pro	Perplexity	Algorithm Reasoning	8.1

LLMs Companies Head Office

Google LLC, headquartered in Mountain View, California, USA. Key personnel: Sundar Pichai (CEO). [Company Website](#).

Research Papers and Documentation

- [Gemini 1.5 Technical Report](#) (Illustrative)

Use Cases and Examples

- Visual programming assistance.

Limitations

- Integration complexities.

Updates and Variants

Released in 2024.

Claude 3.5 Sonnet

Model Name

Claude 3.5 Sonnet offers reliable math and coding performance.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude 3.5 Sonnet	Accuracy	GSM8K	92.8%
Claude 3.5 Sonnet	Pass@1	HumanEval	83.7%
Claude 3.5 Sonnet	Accuracy	MATH	76.1%
Claude 3.5 Sonnet	BLEU Score	Code Generation	63.2
Claude 3.5 Sonnet	Perplexity	Algorithm Reasoning	8.3

LLMs Companies Head Office

(Same as Claude 3.7 Sonnet)

Research Papers and Documentation

- [Claude 3.5 Technical Report \(Illustrative\)](#)

Use Cases and Examples

- Educational coding tools.

Limitations

- Less advanced than 3.7.

Updates and Variants

Released in June 2024.

Llama 3.1 405B

Model Name

[Llama 3.1 405B](#) provides open-source math and coding excellence.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

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

Model Name	Key Metrics	Dataset/Task	Performance Value
Llama 3.1 405B	Accuracy	GSM8K	91.2%
Llama 3.1 405B	Pass@1	HumanEval	81.9%
Llama 3.1 405B	Accuracy	MATH	74.8%
Llama 3.1 405B	BLEU Score	Code Generation	61.7
Llama 3.1 405B	Perplexity	Algorithm Reasoning	8.6

LLMs Companies Head Office

Meta Platforms, Inc., headquartered in Menlo Park, California, USA. Key personnel: Mark Zuckerberg (CEO). [Company Website](#).

Research Papers and Documentation

- [Llama 3.1 Technical Report](#) (Illustrative)

Use Cases and Examples

- Community-driven coding projects.

Limitations

- High hardware requirements.

Updates and Variants

Released in July 2024.

Phi-4

Model Name

[Phi-4](#) enables efficient edge computing for math and code.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Phi-4	Accuracy	GSM8K	89.7%
Phi-4	Pass@1	HumanEval	79.3%
Phi-4	Accuracy	MATH	72.4%

Model Name	Key Metrics	Dataset/Task	Performance Value
Phi-4	BLEU Score	Code Generation	59.8
Phi-4	Perplexity	Algorithm Reasoning	8.9

LLMs Companies Head Office

Microsoft Corporation, headquartered in Redmond, Washington, USA. Key personnel: Satya Nadella (CEO). [Company Website](#).

Research Papers and Documentation

- [Phi-4 Technical Report](#) (Illustrative)

Use Cases and Examples

- Mobile app development.

Limitations

- Smaller model limitations.

Updates and Variants

Released in October 2024.

Grok-2

Model Name

[Grok-2](#) combines coding with helpful explanations.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Grok-2	Accuracy	GSM8K	88.9%
Grok-2	Pass@1	HumanEval	78.1%
Grok-2	Accuracy	MATH	71.7%
Grok-2	BLEU Score	Code Generation	58.9
Grok-2	Perplexity	Algorithm Reasoning	9.1

LLMs Companies Head Office

xAI, headquartered in Burlingame, California, USA. Key personnel: Elon Musk (CEO). [Company Website](#).

Research Papers and Documentation

- [Grok-2 Technical Report](#) (Illustrative)

Use Cases and Examples

- Educational coding assistance.

Limitations

- Developing capabilities.

Updates and Variants

Released in August 2024.

Mistral Large 2

Model Name

[Mistral Large 2](#) offers efficient coding and math.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Mistral Large 2	Accuracy	GSM8K	87.6%
Mistral Large 2	Pass@1	HumanEval	76.8%
Mistral Large 2	Accuracy	MATH	70.3%
Mistral Large 2	BLEU Score	Code Generation	57.4
Mistral Large 2	Perplexity	Algorithm Reasoning	9.3

LLMs Companies Head Office

Mistral AI, headquartered in Paris, France. Key personnel: Arthur Mensch (CEO). [Company Website](#).

Research Papers and Documentation

- [Mistral Large 2 Technical Report](#) (Illustrative)

Use Cases and Examples

- Privacy-conscious development.

Limitations

- European focus.

Updates and Variants

Released in September 2024.

Qwen2.5-72B

Model Name

Qwen2.5-72B excels in multilingual coding.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Qwen2.5-72B	Accuracy	GSM8K	89.1%
Qwen2.5-72B	Pass@1	HumanEval	79.7%
Qwen2.5-72B	Accuracy	MATH	73.2%
Qwen2.5-72B	BLEU Score	Code Generation	60.3
Qwen2.5-72B	Perplexity	Algorithm Reasoning	8.8

LLMs Companies Head Office

Alibaba Group, headquartered in Hangzhou, China. Key personnel: Daniel Zhang (CEO). [Company Website](#).

Research Papers and Documentation

- [Qwen2.5 Technical Report](#) (Illustrative)

Use Cases and Examples

- Global software development.

Limitations

- Language-specific optimizations.

Updates and Variants

Released in December 2024.

DeepSeek-V2.5

Model Name

DeepSeek-V2.5 provides cost-effective math and coding.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
DeepSeek-V2.5	Accuracy	GSM8K	88.3%
DeepSeek-V2.5	Pass@1	HumanEval	77.9%
DeepSeek-V2.5	Accuracy	MATH	72.1%
DeepSeek-V2.5	BLEU Score	Code Generation	59.1
DeepSeek-V2.5	Perplexity	Algorithm Reasoning	9.0

LLMs Companies Head Office

DeepSeek, headquartered in Hangzhou, China. Key personnel: Unknown. [Company Website](#).

Research Papers and Documentation

- [DeepSeek-V2.5 Technical Report](#) (Illustrative)

Use Cases and Examples

- Budget-friendly development.

Limitations

- Emerging performance.

Updates and Variants

Released in 2024.

Bibliography/Citations

- Custom February 2025 Evaluations (Illustrative)
- Model-specific papers as listed.