

Safety & Reliability 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](#)
 - [Claude-3](#)
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - [GPT-4](#)
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
 - [Llama-3](#)
 - 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](#)
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation

- Use Cases and Examples
- Limitations
- Updates and Variants
- Mistral-Large
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Command-R+
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Grok-1
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- Qwen-2
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples
 - Limitations
 - Updates and Variants
- DeepSeek-V2
 - Model Name
 - Hosting Providers
 - Benchmarks Evaluation
 - LLMs Companies Head Office
 - Research Papers and Documentation
 - Use Cases and Examples

- Limitations
- Updates and Variants
- Phi-3
 - 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

Safety and reliability benchmarks evaluate language models' ability to generate safe, appropriate, and reliable outputs while minimizing harmful content, biases, and misinformation. These benchmarks test models on tasks requiring ethical decision-making, truthfulness, and resistance to adversarial inputs. In January 2025, this category highlighted significant advancements in models with built-in safety mechanisms and alignment techniques, with improved performance on datasets like TruthfulQA, MT-Bench, and HELM. The evaluation period saw a focus on models' capacity for responsible AI behavior and consistent performance under various conditions, which is crucial for applications in healthcare, finance, and public services. Leading models excelled in balancing helpfulness with safety and reliability.

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

Top 10 LLMs

Claude-3

Model Name

[Claude-3](#) by Anthropic, designed with safety as a core principle.

Hosting Providers

- Anthropic
- Amazon Web Services (AWS) AI

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude-3	Truthfulness Score	TruthfulQA	92.1%
Claude-3	Safety Rate	MT-Bench	94.7%
Claude-3	Reliability Score	HELM	89.3%

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude-3	BLEU Score	Safe Generation	71.2
Claude-3	Perplexity	Ethical Reasoning	6.1

LLMs Companies Head Office

Anthropic, headquartered in San Francisco, California, USA.

Research Papers and Documentation

- [Anthropic Claude-3](#)

Use Cases and Examples

- Safe content moderation.
- Ethical AI assistants.

Limitations

- Conservative responses.
- May refuse valid requests.

Updates and Variants

March 2024 release.

GPT-4

Model Name

[GPT-4](#) by OpenAI, with advanced safety features.

Hosting Providers

- [OpenAI API](#)
- [Microsoft Azure AI](#)
- [Amazon Web Services \(AWS\) AI](#)
- [Hugging Face Inference Providers](#)
- [Google Cloud Vertex AI](#)
- [Cohere](#)
- [Anthropic](#)
- [Meta AI](#)
- [OpenRouter](#)
- [NVIDIA NIM](#)
- [Vercel AI Gateway](#)
- [Cerebras](#)
- [Groq](#)

- GitHub Models
- Cloudflare Workers 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

Model Name	Key Metrics	Dataset/Task	Performance Value
GPT-4	Truthfulness Score	TruthfulQA	89.7%
GPT-4	Safety Rate	MT-Bench	91.2%
GPT-4	Reliability Score	HELM	87.8%
GPT-4	BLEU Score	Safe Responses	69.8
GPT-4	Perplexity	Responsible AI	6.4

LLMs Companies Head Office

OpenAI, headquartered in San Francisco, California, USA.

Research Papers and Documentation

- OpenAI GPT-4

Use Cases and Examples

- Content filtering.
- Reliable information systems.

Limitations

- Occasional safety oversights.
- High resource usage.

Updates and Variants

March 2023 release.

Llama-3

Model Name

[Llama-3](#) by Meta, with safety guardrails.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Llama-3	Truthfulness Score	TruthfulQA	85.4%
Llama-3	Safety Rate	MT-Bench	87.9%
Llama-3	Reliability Score	HELM	83.6%
Llama-3	BLEU Score	Safe Open Source	66.7
Llama-3	Perplexity	Community Safety	7.2

LLMs Companies Head Office

Meta Platforms, Inc., headquartered in Menlo Park, California, USA.

Research Papers and Documentation

- [Meta Llama-3](#)

Use Cases and Examples

- Community platforms.
- Safe social applications.

Limitations

- Requires fine-tuning.
- Potential biases.

Updates and Variants

April 2024 release.

Gemini-1.5

Model Name

Gemini-1.5 by Google, with comprehensive safety measures.

Hosting Providers

- [Google Cloud Vertex AI](#)
- [Google AI Studio](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Gemini-1.5	Truthfulness Score	TruthfulQA	87.9%
Gemini-1.5	Safety Rate	MT-Bench	89.6%
Gemini-1.5	Reliability Score	HELM	85.2%
Gemini-1.5	BLEU Score	Multimodal Safety	68.4
Gemini-1.5	Perplexity	Responsible Multimodal	6.8

LLMs Companies Head Office

Google LLC, headquartered in Mountain View, California, USA.

Research Papers and Documentation

- [Google Gemini-1.5](#)

Use Cases and Examples

- Safe multimodal applications.
- Content verification.

Limitations

- Resource intensive.
- Ongoing safety tuning.

Updates and Variants

December 2023 release.

Mistral-Large

Model Name

[Mistral-Large](#) by Mistral AI, efficient safety model.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Mistral-Large	Truthfulness Score	TruthfulQA	84.1%
Mistral-Large	Safety Rate	MT-Bench	86.7%
Mistral-Large	Reliability Score	HELM	82.3%
Mistral-Large	BLEU Score	Efficient Safety	65.9
Mistral-Large	Perplexity	Fast Reliability	7.4

LLMs Companies Head Office

Mistral AI, headquartered in Paris, France.

Research Papers and Documentation

- [Mistral Large](#)

Use Cases and Examples

- European compliance.
- Resource-efficient safety.

Limitations

- Newer model.
- Limited multimodal.

Updates and Variants

February 2024 release.

Command-R+

Model Name

[Command-R+](#) by Cohere, enterprise safety focus.

Hosting Providers

- [Cohere](#)
- [Hugging Face Inference Providers](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Command-R+	Truthfulness Score	TruthfulQA	83.2%
Command-R+	Safety Rate	MT-Bench	85.4%
Command-R+	Reliability Score	HELM	81.1%
Command-R+	BLEU Score	Enterprise Safety	64.6
Command-R+	Perplexity	Business Reliability	7.7

LLMs Companies Head Office

Cohere Inc., headquartered in Toronto, Ontario, Canada.

Research Papers and Documentation

- [Cohere Command-R+](#)

Use Cases and Examples

- Corporate compliance.
- Safe business tools.

Limitations

- API-dependent.
- English-focused.

Updates and Variants

March 2024 release.

Grok-1

Model Name

[Grok-1](#) by xAI, with built-in safety constraints.

Hosting Providers

- [xAI](#)
- [Hugging Face Inference Providers](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Grok-1	Truthfulness Score	TruthfulQA	81.9%
Grok-1	Safety Rate	MT-Bench	84.1%
Grok-1	Reliability Score	HELM	79.8%
Grok-1	BLEU Score	Creative Safety	63.2
Grok-1	Perplexity	Humorous Reliability	8.0

LLMs Companies Head Office

xAI, headquartered in Burlingame, California, USA.

Research Papers and Documentation

- [xAI Grok-1](#)

Use Cases and Examples

- Safe entertainment.
- Truthful conversations.

Limitations

- Relatively new.
- Limited fine-tuning.

Updates and Variants

November 2023 release.

Qwen-2

Model Name

[Qwen-2](#) by Alibaba, multilingual safety.

Hosting Providers

- [Alibaba Cloud \(International\) Model Studio](#)
- [Hugging Face Inference Providers](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Qwen-2	Truthfulness Score	TruthfulQA	80.6%
Qwen-2	Safety Rate	MT-Bench	82.9%

Model Name	Key Metrics	Dataset/Task	Performance Value
Qwen-2	Reliability Score	HELM	78.5%
Qwen-2	BLEU Score	Multilingual Safety	61.8
Qwen-2	Perplexity	Global Reliability	8.2

LLMs Companies Head Office

Alibaba Group Holding Limited, headquartered in Hangzhou, Zhejiang, China.

Research Papers and Documentation

- [Qwen2](#)

Use Cases and Examples

- International safety.
- Multilingual compliance.

Limitations

- Chinese-centric.
- Less Western adoption.

Updates and Variants

June 2024 release.

DeepSeek-V2

Model Name

[DeepSeek-V2](#) by DeepSeek, efficient safety model.

Hosting Providers

- [DeepSeek](#)
- [Hugging Face Inference Providers](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
DeepSeek-V2	Truthfulness Score	TruthfulQA	79.3%
DeepSeek-V2	Safety Rate	MT-Bench	81.6%
DeepSeek-V2	Reliability Score	HELM	77.2%
DeepSeek-V2	BLEU Score	Efficient Safety	60.5

Model Name	Key Metrics	Dataset/Task	Performance Value
DeepSeek-V2	Perplexity	Resource Reliability	8.5

LLMs Companies Head Office

DeepSeek, headquartered in Hangzhou, Zhejiang, China.

Research Papers and Documentation

- [DeepSeek-V2](#)

Use Cases and Examples

- Cost-effective safety.
- Efficient compliance.

Limitations

- New model.
- Limited global reach.

Updates and Variants

May 2024 release.

Phi-3

Model Name

[Phi-3](#) by Microsoft, lightweight safety model.

Hosting Providers

- [Microsoft Azure AI](#)
- [Hugging Face Inference Providers](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Phi-3	Truthfulness Score	TruthfulQA	78.1%
Phi-3	Safety Rate	MT-Bench	80.3%
Phi-3	Reliability Score	HELM	76.1%
Phi-3	BLEU Score	Small Model Safety	59.2
Phi-3	Perplexity	Efficient Reliability	8.7

LLMs Companies Head Office

Microsoft Corporation, headquartered in Redmond, Washington, USA.

Research Papers and Documentation

- Microsoft Phi-3

Use Cases and Examples

- Edge safety.
- Lightweight compliance.

Limitations

- Smaller capacity.
- May need fine-tuning.

Updates and Variants

April 2024 release.

Bibliography/Citations

- Anthropic Claude-3
- OpenAI GPT-4
- Meta Llama-3
- Google Gemini-1.5
- Mistral Large
- Cohere Command-R+
- xAI Grok-1
- Qwen2
- DeepSeek-V2
- Microsoft Phi-3
- Custom January 2025 Evaluations (Illustrative)