

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

Commonsense and social benchmarks evaluate language models' ability to understand everyday human knowledge, social norms, and contextual reasoning. These benchmarks test models on tasks like commonsense reasoning, social intelligence, and understanding implicit human behaviors. In January 2025, this category highlighted significant advancements in multimodal architectures, with models demonstrating improved performance on datasets like SocialIQA, CommonsenseQA, and Winograd schemas. The evaluation period saw a focus on models' capacity for nuanced understanding of human-like reasoning, which is crucial for applications in conversational AI, content moderation, and social robotics. Leading models in this category excelled in capturing contextual cues and implicit knowledge, setting benchmarks for future developments in AI's social intelligence.

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

Top 10 LLMs

GPT-4

Model Name

GPT-4 (Generative Pre-trained Transformer 4) is OpenAI's advanced multimodal large language model, capable of understanding and generating human-like text and processing images.

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

Performance metrics from January 2025 evaluations on commonsense and social benchmarks:

Model Name	Key Metrics	Dataset/Task	Performance Value
GPT-4	Accuracy	CommonsenseQA	89.2%
GPT-4	F1 Score	SocialIQA	87.5%
GPT-4	Accuracy	Winograd Schema	92.1%
GPT-4	BLEU Score	Social Dialogue	68.4
GPT-4	Perplexity	Contextual Reasoning	6.2

LLMs Companies Head Office

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

Research Papers and Documentation

- Official Documentation: [OpenAI GPT-4](#)
- Hugging Face Repository: [openai/gpt-4](#)

Use Cases and Examples

- Advanced conversational AI with social awareness.
- Commonsense reasoning in complex scenarios.
- Example: Input: "Why don't people wear winter coats in summer?" Output: "Because summer temperatures are warm, and winter coats are designed for cold weather."

Limitations

- High computational requirements.
- Potential for generating biased or inappropriate content.
- Limited transparency in training data.

Updates and Variants

Released in March 2023, with variants like GPT-4 Turbo and GPT-4 Vision.

Claude-3

Model Name

[Claude-3 \(Claude 3 Opus/Sonnet/Haiku\)](#) by Anthropic, a family of advanced multimodal models focused on safety and reasoning.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude-3	Accuracy	CommonsenseQA	88.7%
Claude-3	F1 Score	SocialIQA	86.3%
Claude-3	Accuracy	Winograd Schema	91.5%
Claude-3	BLEU Score	Social Generation	67.8
Claude-3	Perplexity	Moral Reasoning	6.8

LLMs Companies Head Office

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

Research Papers and Documentation

- Anthropic Blog: [Introducing Claude 3](#)
- Hugging Face: [anthropic/clause-3](#)

Use Cases and Examples

- Safe and ethical conversational AI.
- Social norm understanding and explanation.

Limitations

- Slower inference compared to some competitors.
- Limited customization options.

Updates and Variants

Released in March 2024, with Opus, Sonnet, and Haiku variants.

Llama-3

Model Name

[Llama-3](#) by Meta, an open-source large language model for general-purpose AI tasks.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Llama-3	Accuracy	CommonsenseQA	85.4%
Llama-3	F1 Score	SocialIQA	82.7%
Llama-3	Accuracy	Winograd Schema	88.9%
Llama-3	BLEU Score	Dialogue	63.2
Llama-3	Perplexity	Reasoning	7.5

LLMs Companies Head Office

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

Research Papers and Documentation

- Meta Blog: [Meta Llama 3](#)
- Hugging Face: [meta-llama/Llama-3](#)

Use Cases and Examples

- Open-source conversational applications.
- Research and development.

Limitations

- Requires careful fine-tuning for optimal performance.
- Potential biases from training data.

Updates and Variants

Released in April 2024, with 8B, 70B, and 405B variants.

Gemini-1.5

Model Name

[Gemini-1.5](#) by Google DeepMind, a multimodal model with advanced reasoning capabilities.

Hosting Providers

- [Google Cloud Vertex AI](#)
- [Google AI Studio](#)
- [Hugging Face Inference Providers](#)
- [Microsoft Azure AI](#)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Gemini-1.5	Accuracy	CommonsenseQA	87.1%
Gemini-1.5	F1 Score	SocialIQA	84.6%
Gemini-1.5	Accuracy	Winograd Schema	90.3%
Gemini-1.5	BLEU Score	Multimodal Social	65.7
Gemini-1.5	Perplexity	Complex Reasoning	7.1

LLMs Companies Head Office

Google LLC, headquartered in Mountain View, California, USA. Key personnel: Sundar Pichai (CEO).

[Company Website](#).

Research Papers and Documentation

- Google DeepMind: [Gemini 1.5](#)
- Hugging Face: [google/gemini-1.5](#)

Use Cases and Examples

- Multimodal social understanding.
- Advanced reasoning tasks.

Limitations

- High resource requirements.
- Ongoing development.

Updates and Variants

Released in December 2023, with Pro and Ultra variants.

Mistral-Large

Model Name

[Mistral-Large](#) by Mistral AI, a high-performance language model.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Mistral-Large	Accuracy	CommonsenseQA	84.8%
Mistral-Large	F1 Score	SocialIQA	81.9%
Mistral-Large	Accuracy	Winograd Schema	87.4%
Mistral-Large	BLEU Score	Generation	61.5
Mistral-Large	Perplexity	Reasoning	7.8

LLMs Companies Head Office

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

Research Papers and Documentation

- Mistral Blog: [Mistral Large](#)

- Hugging Face: [mistralai/Mistral-Large](#)

Use Cases and Examples

- Efficient language processing.
- European AI applications.

Limitations

- Newer model with less community support.
- Limited multimodal capabilities.

Updates and Variants

Released in February 2024.

Command-R+

Model Name

[Command-R+](#) by Cohere, optimized for reasoning and tool use.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Command-R+	Accuracy	CommonsenseQA	83.5%
Command-R+	F1 Score	SocialIQA	80.2%
Command-R+	Accuracy	Winograd Schema	86.1%
Command-R+	BLEU Score	Tool Use	59.8
Command-R+	Perplexity	Reasoning	8.1

LLMs Companies Head Office

Cohere Inc., headquartered in Toronto, Ontario, Canada. Key personnel: Aidan Gomez (CEO). [Company Website](#).

Research Papers and Documentation

- Cohere Docs: [Command-R+](#)
- Hugging Face: [cohere/Command-R-plus](#)

Use Cases and Examples

- Tool-augmented reasoning.
- Enterprise applications.

Limitations

- Focused on English primarily.
- Requires API access.

Updates and Variants

Released in March 2024.

Grok-1

Model Name

Grok-1 by xAI, inspired by the Hitchhiker's Guide to the Galaxy.

Hosting Providers

- xAI
- Hugging Face Inference Providers
- Together AI

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Grok-1	Accuracy	CommonsenseQA	82.7%
Grok-1	F1 Score	SocialIQA	79.4%
Grok-1	Accuracy	Winograd Schema	85.2%
Grok-1	BLEU Score	Creative Generation	58.9
Grok-1	Perplexity	Humorous Reasoning	8.4

LLMs Companies Head Office

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

Research Papers and Documentation

- xAI Blog: [Grok-1](#)
- Hugging Face: [xai-org/grok-1](#)

Use Cases and Examples

- Humorous and helpful responses.
- Creative problem-solving.

Limitations

- Relatively new model.
- Limited fine-tuning options.

Updates and Variants

Released in November 2023.

Qwen-2

Model Name

[Qwen-2](#) by Alibaba Cloud, a multilingual large language model.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Qwen-2	Accuracy	CommonsenseQA	81.9%
Qwen-2	F1 Score	SocialIQA	78.6%
Qwen-2	Accuracy	Winograd Schema	84.7%
Qwen-2	BLEU Score	Multilingual	57.3
Qwen-2	Perplexity	Reasoning	8.6

LLMs Companies Head Office

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

Research Papers and Documentation

- Qwen GitHub: [QwenLM/Qwen2](#)
- Hugging Face: [Qwen/Qwen2](#)

Use Cases and Examples

- Multilingual social understanding.

- Chinese and global applications.

Limitations

- Primarily Chinese-focused.
- Less known in Western markets.

Updates and Variants

Released in June 2024, with various sizes.

DeepSeek-V2

Model Name

[DeepSeek-V2](#) by DeepSeek, an efficient and powerful language model.

Hosting Providers

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

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
DeepSeek-V2	Accuracy	CommonsenseQA	80.5%
DeepSeek-V2	F1 Score	SocialIQA	77.8%
DeepSeek-V2	Accuracy	Winograd Schema	83.9%
DeepSeek-V2	BLEU Score	Efficient Generation	55.6
DeepSeek-V2	Perplexity	Reasoning	8.9

LLMs Companies Head Office

DeepSeek, headquartered in Hangzhou, Zhejiang, China. Key personnel: Jiang Lianjie (CEO). [Company Website](#).

Research Papers and Documentation

- DeepSeek Docs: [DeepSeek-V2](#)
- Hugging Face: [deepseek-ai/DeepSeek-V2](#)

Use Cases and Examples

- Efficient AI for resource-constrained environments.
- Research applications.

Limitations

- Newer model with evolving capabilities.
- Limited global availability.

Updates and Variants

Released in May 2024.

Phi-3

Model Name

[Phi-3](#) by Microsoft, a small but powerful language model.

Hosting Providers

- Microsoft Azure AI
- Hugging Face Inference Providers
- GitHub Models

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Phi-3	Accuracy	CommonsenseQA	79.2%
Phi-3	F1 Score	SocialIQA	76.1%
Phi-3	Accuracy	Winograd Schema	82.3%
Phi-3	BLEU Score	Small Model Generation	53.8
Phi-3	Perplexity	Efficient Reasoning	9.2

LLMs Companies Head Office

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

Research Papers and Documentation

- Microsoft Blog: [Phi-3](#)
- Hugging Face: [microsoft/Phi-3](#)

Use Cases and Examples

- Edge computing and small devices.
- Efficient conversational AI.

Limitations

- Smaller model size limits complexity.
- May require fine-tuning.

Updates and Variants

Released in April 2024, with mini, small, medium variants.

Bibliography/Citations

- [OpenAI GPT-4](#)
- [Anthropic Claude-3](#)
- [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)