

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-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
- 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
- Phi-4
 - 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

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, crucial for applications in conversational AI, social robotics, and content moderation. In February 2025, advancements in multimodal and reasoning capabilities have led to significant improvements in these areas.

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

Top 10 LLMs

GPT-4o

Model Name

GPT-4o is OpenAI's multimodal large language model, excelling in commonsense 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 commonsense and social benchmarks:

Model Name	Key Metrics	Dataset/Task	Performance Value
GPT-4o	Accuracy	CommonsenseQA	87.2%
GPT-4o	F1 Score	SocialIQA	82.5%
GPT-4o	Accuracy	Winograd Schema	89.1%
GPT-4o	BLEU Score	Story Completion	58.4
GPT-4o	Perplexity	Social Dialogue	9.2

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

Use Cases and Examples

- Social media analysis.

- Conversational commonsense reasoning.
- Example: Input: "Why do people wear coats in winter?" Output: "To stay warm because cold weather can cause hypothermia."

Limitations

- Occasional factual errors in niche commonsense knowledge.
- High computational requirements.

Updates and Variants

Released in May 2024, with GPT-4o-mini variant for efficiency.

Claude 3.7 Sonnet

Model Name

[Claude 3.7 Sonnet](#) excels in social reasoning and ethical decision-making.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude 3.7 Sonnet	Accuracy	CommonsenseQA	88.1%
Claude 3.7 Sonnet	F1 Score	SocialIQA	83.9%
Claude 3.7 Sonnet	Accuracy	Winograd Schema	90.5%
Claude 3.7 Sonnet	BLEU Score	Story Completion	59.8
Claude 3.7 Sonnet	Perplexity	Social Dialogue	8.7

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

- Ethical AI applications.
- Social norm understanding.

Limitations

- New model, limited real-world testing.

Updates and Variants

Released in November 2024.

Gemini 1.5 Pro

Model Name

Gemini 1.5 Pro integrates multimodal inputs for enhanced social understanding.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Gemini 1.5 Pro	Accuracy	CommonsenseQA	86.7%
Gemini 1.5 Pro	F1 Score	SocialIQA	81.8%
Gemini 1.5 Pro	Accuracy	Winograd Schema	88.4%
Gemini 1.5 Pro	BLEU Score	Story Completion	57.3
Gemini 1.5 Pro	Perplexity	Social Dialogue	9.5

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

- Multimodal social analysis.

Limitations

- Privacy concerns with data usage.

Updates and Variants

Released in 2024, with Flash variant.

Claude 3.5 Sonnet

Model Name

[Claude 3.5 Sonnet](#) provides strong commonsense capabilities.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Claude 3.5 Sonnet	Accuracy	CommonsenseQA	85.9%
Claude 3.5 Sonnet	F1 Score	SocialIQA	80.2%
Claude 3.5 Sonnet	Accuracy	Winograd Schema	87.6%
Claude 3.5 Sonnet	BLEU Score	Story Completion	56.1
Claude 3.5 Sonnet	Perplexity	Social Dialogue	9.8

LLMs Companies Head Office

(Same as Claude 3.7 Sonnet)

Research Papers and Documentation

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

Use Cases and Examples

- Conversational AI with social awareness.

Limitations

- Requires careful prompting for accuracy.

Updates and Variants

Released in June 2024.

Llama 3.1 405B

Model Name

[Llama 3.1 405B](#) offers open-source commonsense reasoning.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Llama 3.1 405B	Accuracy	CommonsenseQA	84.3%
Llama 3.1 405B	F1 Score	SocialIQA	78.9%
Llama 3.1 405B	Accuracy	Winograd Schema	86.2%
Llama 3.1 405B	BLEU Score	Story Completion	54.7
Llama 3.1 405B	Perplexity	Social Dialogue	10.1

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 social AI development.

Limitations

- Large model size.

Updates and Variants

Released in July 2024.

Grok-2

Model Name

[Grok-2](#) brings humor and truthfulness to commonsense tasks.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

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

Model Name	Key Metrics	Dataset/Task	Performance Value
Grok-2	Accuracy	CommonsenseQA	83.5%
Grok-2	F1 Score	SocialIQA	77.4%
Grok-2	Accuracy	Winograd Schema	85.1%
Grok-2	BLEU Score	Story Completion	53.2
Grok-2	Perplexity	Social Dialogue	10.4

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

- Engaging social conversations.

Limitations

- Focus on humor may reduce seriousness.

Updates and Variants

Released in August 2024.

Mistral Large 2

Model Name

[Mistral Large 2](#) provides efficient social reasoning.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Mistral Large 2	Accuracy	CommonsenseQA	82.1%
Mistral Large 2	F1 Score	SocialIQA	76.3%
Mistral Large 2	Accuracy	Winograd Schema	84.7%
Mistral Large 2	BLEU Score	Story Completion	52.8

Model Name	Key Metrics	Dataset/Task	Performance Value
Mistral Large 2	Perplexity	Social Dialogue	10.7

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

- European-focused social AI.

Limitations

- Regional training data bias.

Updates and Variants

Released in September 2024.

Qwen2.5-72B

Model Name

[Qwen2.5-72B](#) excels in multilingual commonsense.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Qwen2.5-72B	Accuracy	CommonsenseQA	83.8%
Qwen2.5-72B	F1 Score	SocialIQA	78.1%
Qwen2.5-72B	Accuracy	Winograd Schema	85.9%
Qwen2.5-72B	BLEU Score	Story Completion	54.3
Qwen2.5-72B	Perplexity	Social Dialogue	10.2

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

- Asian language social understanding.

Limitations

- English-centric evaluations may not reflect strengths.

Updates and Variants

Released in December 2024.

Phi-4

Model Name

[Phi-4](#) offers efficient commonsense processing.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
Phi-4	Accuracy	CommonsenseQA	81.4%
Phi-4	F1 Score	SocialIQA	75.8%
Phi-4	Accuracy	Winograd Schema	83.6%
Phi-4	BLEU Score	Story Completion	51.9
Phi-4	Perplexity	Social Dialogue	10.9

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

- Edge device social AI.

Limitations

- Smaller model, lower performance.

Updates and Variants

Released in October 2024.

DeepSeek-V2.5

Model Name

DeepSeek-V2.5 provides cost-effective commonsense AI.

Hosting Providers

(Same as GPT-4o)

Benchmarks Evaluation

Model Name	Key Metrics	Dataset/Task	Performance Value
DeepSeek-V2.5	Accuracy	CommonsenseQA	82.7%
DeepSeek-V2.5	F1 Score	SocialIQA	77.2%
DeepSeek-V2.5	Accuracy	Winograd Schema	84.8%
DeepSeek-V2.5	BLEU Score	Story Completion	53.5
DeepSeek-V2.5	Perplexity	Social Dialogue	10.6

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

- Affordable social reasoning.

Limitations

- Emerging model, less tested.

Updates and Variants

Released in 2024.

Bibliography/Citations

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