Artificial Analysis State of Al

Q1 2025

Highlights Report

Full report available to Premium Access subscribers

Artificial Analysis is a leading and independent AI benchmarking and insights provider. We support engineers and companies to understand AI capabilities and make critical decisions about their AI strategies.

Our data, insights and publications are grounded in our comprehensive benchmarking of AI technologies and use cases. This includes everything from hourly performance testing of language model APIs to millions of votes in our crowd-sourced evaluations.

Our public website, <u>artificialanalysis.ai</u>, is widely referenced throughout the AI industry. To discuss this report, our publications or our services, please contact us at <u>contact@artificialanalysis.ai</u>.



6 major trends have shaped progress in AI in Q1 2025



A) Al Progress Continues

Major AI labs have all continued to make substantial gains in intelligence, cost efficiency and speed



B) Reasoning Models

Models which 'think' before answering by outputting tokens drove significant intelligence gains and became widespread beyond OpenAI



C) Efficiency & MoEs

Models increasingly use a Mixture of Experts architecture, activating a subset of parameters at inference time, increasing inference efficiency



) Rise of Chinese Al

AI Labs headquartered in China have released models with intelligence rivalling that of US labs, particularly amongst open weights models



) Agents

Al systems can increasingly perform tasks end-to-end autonomously by chaining requests to LLMs together



Multimodal AI

Al models are increasingly multimodal, working natively with image and audio; modality-specific models continue to advance



Players in the AI value chain differ in levels of vertical integration; Google continues to stand out as the most vertically integrated player from TPU accelerators to Gemini

Key Players in the AI Value Chain (Non-exhaustive)

Classifications are indicative and determined based a range of factors including market share and strength of offering No presence





Strong presence

Big technology companies are continuing to play across all AI modalities while smaller challengers tend to focus on specific modalities





OpenAl remains in lead, but competition for top tier intelligence is fiercer than ever; reasoning models that 'think' before they answer are driving the latest leaps in intelligence Frontier Large Language Model (LLM) Intelligence, Over Time

Artificial Analysis Intelligence Index (incorporates MMLU-Pro, GPQA Diamond, Humanity's Last Exam, LiveCodeBench, SciCode, AIME, MATH-500)



- **OpenAl continues to maintain lead but gap narrows:** OpenAl's o4-mini (high) model is the most intelligent model closely followed by reasoning models from Google's Gemini 2.5 Pro and xAl's Grok 3 model
- Open weights models amongst the most intelligent: DeepSeek R1 and NVIDIA Llama 3.1 Nemotron Ultra approach intelligence of proprietary models
- Labs continue to release rapidly on a quarterly basis: Labs have typically released their new frontier models on a 3-month release cycle



OpenAI, Google and xAI lead frontier intelligence with their latest reasoning models, followed closely by other labs

Leading Large Language Model (LLMs), by AI lab

Artificial Analysis Intelligence Index (incorporates MMLU-Pro, GPQA, Humanity's Last Exam, LiveCodeBench, SciCode, AIME, MATH-500)



Commentary

- OpenAl continues to lead: OpenAl's o3 and o4-mini-high sit on the current intelligence frontier, with Google coming closer than ever to having the top model in our intelligence Index with Gemini 2.5 Pro
- Reasoning is new frontier: All the models that score highest in Intelligence Index are reasoning models that 'think' before answering the questions
- A more and more crowded frontier: The big 5 US labs (OpenAl, Google, Anthropic, xAl, Meta) are joined near the frontier in early 2025 by NVIDIA, DeepSeek, Alibaba, Mistral and Amazon

Artificial Analysis



Open weights reasoning models mostly closed the gap opened by OpenAI's o1 in late 2024; proprietary models continue to lead

Leading Language Models by License Type, Over Time

Artificial Analysis Intelligence Index (incorporates MMLU-Pro, GPQA, Humanity's Last Exam, LiveCodeBench, SciCode, AIME, MATH-500)



Nov '22 Jan '23Mar '23May '23 Jul '23 Sep '23 Nov '23 Jan '24 Mar '24 May '24 Jul '24 Sep '24 Nov '24 Jan '25 May '25

Release Date

- Open weights reasoning models become available: The release of DeepSeek R1 in early 2025 significantly reduced the intelligence gap by making reasoning available outside of proprietary models; recent proprietary releases (e.g. o4-mini, Gemini 2.5 Pro) have since widened the gap
- Frontier model releases are becoming more frequent: Early 2025 saw a wave of releases and refinements driving intelligence in short rapid incremental improvements





Today's leading open weights models come from Alibaba, DeepSeek, Meta and Nvidia

Leading Language Models by License Type

Artificial Analysis Intelligence Index (incorporates MMLU-Pro, GPQA, Humanity's Last Exam, LiveCodeBench, SciCode, AIME, MATH-500)



- Proprietary models continue to lead frontier intelligence: Proprietary reasoning models from US labs lead in overall intelligence (ie. including reasoning models)
- Open weights models lead non-reasoning intelligence: DeepSeek
 V3 0324 is currently the most intelligent nonreasoning model, followed closely by Llama 4
 Maverick



B. Reasoning Models

Reasoning models now decisively lead Artificial Analysis Intelligence Index

The Language Model Frontier: Reasoning vs Non-Reasoning

Artificial Analysis Quality Index (incorporates 7 leading evaluations), Leading Models



Commentary

Reasoning models are the

intelligence edge over nonreasoning peers and top the

All leading AI labs (including

OpenAI, Google, Anthropic,

reasoning models (or take a

 Non-reasoning models do, however, remain the best option for a wide range of tasks, especially where cost

and speed are critical

xAI, DeepSeek) now offer

reasoning models, but continue to offer non-

hybrid approach)

new frontier as they

demonstrate a clear

Artificial Analysis Intelligence Index

Reasoning models 'think' before they answer, often taking more time and tokens to provide a more accurate answer for complex problems

Non-reasoning model (OpenAl GPT-40) **Reasoning model** (Google Gemini 2.5 Pro) Suppose you're on a game show, and you're given the choice of three doors: Behind one door is a gold bar; behind the others, rotten vegetables. Prompt You pick a door, say No. 1, and the host asks you "Do you want to pick door No. 2 instead?" Is it to your advantage to switch your choice? Model output Here's a thinking process to arrive at the explanation for the Yes, it is to your advantage to switch your choice. Here's why: Monty Hall problem variant: Initially, you have a 1/3 chance of picking the door with the 1. **Identify the core problem:** The user describes a scenario gold bar and a 2/3 chance of picking a door similar to the classic Monty Hall problem, but with a slight variation. ... This scenario is similar to the Monty Hall problem, where switching doors is a statistically better strategy. ... 9. **Self-Correction/Check:** Does this make sense? Yes Okay, let's break this down. This sounds very similar to the famous Monty Hall Problem, but there's a crucial difference in how you've described the host's action Since your host didn't open a door and reveal a losing prize, that crucial shift in probability doesn't happen. It's just a choice between two doors, each with a 1/3 probability. End to end ~4 seconds ~21 seconds response time **Output tokens** ~185 tokens ~1,967 tokens **Incorrectly answers question** as it does not understand the **Correctly answers the questions** as it spends 'thinking' time to Commentary prompt is different to the Monty Hall problem, instead simply understand the prompt varies from the Monty Hall problem and answering as if the query were a Monty Hall problem breaks down this new problem to solve it Uses less time and fewer tokens but does not provide the **Uses more time and tokens** but provides the correct answer correct answer



B. Reasoning Models

Treating reasoning & non-reasoning models as distinct categories is a helpful framework for understanding today's model landscape

Intelligence vs. Output Tokens Used to Run Artificial Analysis Intelligence Index

Artificial Analysis Intelligence Index (Version 2, released Feb 25), Output Tokens Used (~5M input tokens)



- Two distinct groups: The difference in token usage between reasoning and nonreasoning models is much greater than differences within each group; the median reasoning model uses up to 10x more tokens to complete our Intelligence Index as the median non-reasoning model
- When using reasoning models, developers now must go beyond per-token pricing and consider token usage to properly understand inference cost



The cost of inference has significantly reduced as small models continue to get smarter, inference efficiency increases, and new hardware generations drive down the cost of compute



- The cost of intelligence has continued to decline rapidly, with inference for high intelligence models (>50 on intelligence index) already ~32x cheaper since September 2024.
- The cost of intelligence at the level of the original
 GPT-4 has now declined
 >1000x since the launch of GPT-4 in March 2023
- Three primary drivers enable this continuing cost revolution: smaller models (incl MoE architectures), inference optimizations and new hardware generations





Inference is faster than ever in early 2025 – this is particularly critical as reasoning models and agentic applications are generating >10x more tokens per request than average queries a year ago





Frontier models are increasingly adopting a Mixture of Experts (MoE) architecture to balance intelligence and efficiency





NVIDIA continues to dominate the AI accelerator market, especially for frontier-class training, but a growing list of challengers now offer material differentiation

Key players building accelerators for AI training and inference

Based on publicly available data of chips yet to be released and/or available for use



1. Intel is no longer bringing to market Falcon Shores, its successor to Gaudi 3, and we do not expect Intel to bring Jaguar Shores to market until 2026 at the earliest.



Existing available chips

No available chips

Challenger chip companies have vertically integrated into cloud services to offer differentiated LLM inference services with greater performance

Output Speed: Llama 4 Scout Serverless Endpoint Providers

Output Tokens per Second; Higher is better



- Cerebras, SambaNova, Groq are chip companies that have vertically integrated into cloud inference, often providing higher performance (i.e. output speed) than peers using NVIDIA hardware to serve the models
- However, developers will need to consider the tradeoffs between performance, cost and context window. While faster, these chip companies in cases serve models at a higher price point compared to other inference providers and with context windows smaller than models' native context length



While efficiency gains have been made...

GPT-4 level intelligence is now 100x cheaper than original GPT-4

1. Smaller models



impact of each factor impacting cost

... compute demand continues to increase

New applications continue to demand more compute: a single deep research query can cost >10x an original GPT-4 query

~10x

tokens/query

5. Reasoning models

Significant increase in

output tokens when

models 'think' before

answering



Agents chain multiple requests to LLMs to complete tasks autonomously

~5x compute/query

4. Larger models

Scaling laws continue to demand higher parameter counts for greater intelligence

\Lambda Artificial Analysis

D. Rise of Chinese AI

The intelligence frontier is now overwhelmingly dominated by the US and China

Leading Language Models by Country of Origin

Artificial Analysis Intelligence Index (incorporates MMLU-Pro, GPQA, Humanity's Last Exam, LiveCodeBench, SciCode, AIME, MATH-500)





- The US leads across reasoning models: the top 4 places on Artificial Analysis Intelligence Index are all taken by reasoning models from US labs
- However, China currently leads non-reasoning models: DeepSeek V3 0324 is the leading non-reasoning model, eclipsing models from the US and others
- Smaller strides across other countries: Models from labs outside of USA and China have continually improved but don't currently compete for frontier intelligence



E. Agents

Agents are autonomous systems driven by LLMs

What are agents?



"Systems where LLMs dynamically direct their own processes and tool usage, maintaining control over how they accomplish tasks"

"Agents represent systems that intelligently accomplish tasks, ranging from executing simple workflows to pursuing complex, open-ended objectives"





"Al agents are autonomous systems powered by large language models (LLMs) that, given high-level instructions, can plan, use tools, carry out steps of processing, and take actions to achieve specific goals"

Al agents are LLM-driven systems that act autonomously and use tools to complete tasks end-to-end



What agentic applications are working well in early 2025?

Coding	LLM-driven system that reasons about the user's instructions, and then finds, interprets, edits and tests source code to complete software engineering tasks
Deep research	LLM-driven system that parses a research query , launches a chain of targeted research queries while controlling its own research flow , to form a synthesized answer
Computer use	LLM-driven system that interprets user commands , "looks at" the user's desktop or browser, and autonomously chains clicks, keystrokes, shell commands, and API calls to complete the task
Customer support	LLM-driven system that turns live speech into intent and speaks back in real time, while chaining required app, CRM or API calls to complete the task



Coding agents: LLM-driven systems are increasingly taking a more autonomous role in coding, going far beyond code autocompletion and chatbot assistance



1. GitHub Copilot launched in late 2021 as an AI autocomplete tool build on OpenAI's Codex model; 21 GitHub Copilot now supports chat and agentic features alongside AI autocomplete



Both players with a broad AI focus, and those with a specific focus in media generation, have developed models across modalities in the last year

Key players offering image and/or video models

Includes publicly available models in each modality released in the last year

	Broad Focus											Media Generation Focus													NON-EXHAUSTIVE			
	Develop both language and media generation models										Develop only media generation models																	
	OpenAl	ByteDance	Google	MiniMax	Amazon	NVIDIA	XAI	Alibaba	Tencent	Meta	Recraft	HiDream	Reve	ldeogram	Black Forest Labs	Midjourney	Luma Labs	Stability.ai	Leonardo.ai	Playground	Adobe	Kuaishou	Runway	Pika Art	Genmo			
	\$	ht	G	վլի	aws		X	EL	1	∞	R	lii					6	S.			۸	80	ß		-			
Image Generation																												
Image Editing ¹																												
Text to Video																												
Image to Video																												



Existing model

No model

In Q1 '25, image generation reached new heights, as OpenAI's GPT-40 set a new benchmark for visual quality and prompt adherence

Image Generation Quality¹, Over Time

Prompt: The words 'Artificial Analysis' emblazoned on a next-generation spacecraft orbiting a breathtaking Earth view



1. Artificial Analysis Image Arena ELO, calculated as of 6 May 2025 based on 1.6 million votes from Artificial Analysis users



Image generation is entering a new phase of competitiveness with both generalist AI and media-focused labs competing for the frontier with frequent releases





Video models saw a breakthrough in quality, with a number of labs catching up to OpenAI's Sora's capabilities, including Google overtaking Sora with Veo 2

Text to Video ELO Score by Release Date

Artificial Analysis Video Arena ELO Score



- Video models saw a leap forward: A number of video models have been released with comparable quality to OpenAl's Sora
- U.S. and China dominate the frontier: Top-performing models like Google's Veo 2 (USA) and MiniMax's T2V-01 Director (China) show a geographic duopoly at the top of the ELO leaderboard
- Emerging labs lead the wave: The majority of recent releases (e.g. T2V-01 Director by MiniMax, Kling 1.6 by Kuaisho) have come from emerging labs



While the generalist AI labs, including OpenAI, Google, Microsoft and Amazon, have offerings across speech AI, an emerging group of pure-play speech companies are driving innovation

Key players offering speech models

Classifications are indicative and determined based on models available





Strong presence

Low or no presence

Consistent gains in Speech to Text quality have continued as new model releases by OpenAI and ElevenLabs pushed frontier accuracy even higher, achieving lower error rates

Frontier Speech to Text Model Quality by Country over Time

% Word Error Rate



- Consistent moderate growth across labs have seen the Word Error Rates decrease over time as models continue to improve in quality; Q1 '25 seeing frontier Word Error Rate dropping from ~9% to ~8%
- Scribe by ElevenLabs took the lead as the most accurate Speech to Text model in Q1 '25, outperforming OpenAI's Whisper and GPT-40 Transcribe



Q1 '25 saw the release of new Text to Speech models that have continued to incrementally improve quality, including in pushing toward human-like dialogue

Frontier Text to Speech Model Quality by Country, Over Time *Elo Score based on Artificial Analysis Speech Arena*

Most attractive region Google Amazon Microsoft Azure MetaVoice Coqui 🔲 StyleTTS 🛑 OpenVoice 🔲 Cartesia 🛑 PlayAI 💭 Fish Audio 🛑 ElevenLabs 🛑 MiniMax 🗖 OpenAI Kokoro Murf AI Speechify StepFun Rime 📒 Hume AI 🔳 Sesame 📒 Papla 📒 LMNT Zvphra 1250 Artificial Analysis ElevenLabs OpenAl GPT-Turbo v2.5 1200 MiniMax Speech-02 40 Realtime OpenAl TTS-1 1150 Azure Neural 1100 ElevenLabs Flash v2.5 1050 Amazon Pollv Score 1000 950 В 900 850 800 750 700 650 Dec '14 Dec '15 Jun '16 Dec '16 Jun '17 Dec '17 Jun '18 Dec '18 Jun '19 Dec '19 Jun '20 Dec '20 Jun '21 Dec '21 Jun '22 Dec '22 Jun '23 Dec '23 Jun '24 Dec '24 Jun '25 Dec '25 Jun '26 Jun '15

Release Date

- Early 2025 saw the release of models such as MiniMax's Speech-02-HD that helped to improve sound quality, with improvements in speech audio that sounds more like human-like dialogue
- The Text to Speech space has increased in competition in early 2025 as speech-focused labs and open source projects launched new models





Artificial Analysis

contact@artificialanalysis.ai

https://artificialanalysis.ai/

Legal notice: Copyright © 2025 Artificial Analysis, Inc. All rights reserved.

This document, including any data, analysis, and insights contained herein, is provided by Artificial Analysis for informational purposes only. The information is based on data collected through various sources, including but not limited to first party benchmarking and surveys conducted on our website. While Artificial Analysis strives to ensure the accuracy and reliability of the information, it is provided "as is" and may not be complete or up to date. The content should not be construed as professional advice, and recipients are encouraged to conduct their own research and analysis before making any decisions based on this information. By accessing or using this document, you agree to be bound by Artificial Analysis's Terms of Service, available on our website.