Free Weights: Open Source AI



Open source (OS) is powering Gen Al innovation. Thanks to widely available academic research and platforms like GitHub and Hugging Face, we're witnessing a boom in major projects with impressive outcomes. Despite the considerable resources—money, computing power, and data —that closed-source tech giants pour into Al, open-source initiatives are tracking their growth and performance remarkably well.

A <u>leaked internal memo</u> from Google went viral in 2023 for its observation that open source Al has been subtly yet effectively "eating the lunch" of big tech companies like Google and OpenAl, boasting greater speed, adaptability, privacy, and overall efficiency. Open source Al is rapidly gaining on closed source in both popularity and performance – OS models like Mistral, Llama are catching up with and even outperforming some closed-source models.



•••

Open source AI models will soon become unbeatable. Period.

Bindu Reddy @ @bindureddy · Oct 14, 2023
 The pace of open-source LLM innovation and research is breath-taking
 I suspect that open-source will soon become unbeatable for anyone except maybe OpenAl
 ...
 Show more

As a result open source AI is seeing substantial interest from developers, researchers and investors alike. Github witnessed a 148% YOY increase in developer contribution to Gen AI projects in 2023. More than \$8B has been invested in open source AI over the last 2 years.

OS AI Ecosystem: Substantial growth in AI projects as well as contributors

Specifically for Gen AI, the term "open source" typically implies that the source code, any applicable weights and parameters (for training models) of these components are publicly accessible, usable, modifiable, and their distribution is permitted.

Adhering to this definition, the open source AI stack includes comprehensive set of tools to build Gen AI applications - foundational models (such as <u>Llama</u>, <u>Mistral</u>), developer tools & frameworks (such as <u>Langchain</u>, <u>Fixie</u>), model training platforms (such as <u>Weights & Biases</u>, <u>Anyscale</u>), and monitoring tools (<u>Datadog</u>, <u>Seldon</u>).

Open source AI innovation is thriving with new projects and developers:

Open source Gen AI projects are seeing significant and growing projects as well as contributors. Last year, Github witnessed 148% YOY growth in contributors and a 248% YOY growth in the total number of Gen AI projects. There are 60K Gen AI projects on Github and over 400K models on Huggingface as of 2023.

Contributor set is becoming increasingly Global, not restricted to US and Europe

Beyond the US and Europe – where a majority of open source projects originate from – the highest number of individual contributors to open source Gen AI came from India and Japan in 2023. Developers from Hong Kong, UK, Brazil, Germany and Singapore are also making numerous contributions to open source Gen AI. By 2027, <u>India is projected to overtake US as the largest developer community</u> on Github.



Steady increase in serious contributors, while "tourist" interest has tempered since Q1 hype

Gen Al overall has experienced a shift from initial widespread hype (peaking in Q1) to more focused and value-driven engagement - the "trough of disillusionment" phase, where initial excitement gives way to sustained, serious development.

Similar trend can be seen in # of stars across Github repos - the growth has tempered since Q1. On the other hand, serious developers (# of contributors to these projects) have grown steadily - 148% cumulatively in 2023.



Python is the preferred language for open source Al

Gen Al overall has experienced a shift from initial widespread hype (peaking in Q1) to more focuseWhile Javascript has been the top programming language on Github in 2023, Python is the top choice when it comes to Al repositories. Python's preference for ML projects has carried over to Gen Al because of its comprehensive ML libraries like TensorFlow and PyTorch. Python's flexibility in data handling and its platform-independent nature make it highly adaptable for diverse Al projects.

Mojo, a variation of Python that combines the usability of Python with the performance of C++, is gaining traction as an Al-specific programming language. In Q4'23, Mojo saw a 73% MOM increase in Github stars, indicative of the repo's popularity amongst developers.



Al repositories favouring more protective licensing

A disproportionate share of AI repos are using the Apache License, under which developers can claim patents on derivative projects. The Apache license is known to be extensive in legal terminology and therefore offer better patent protection than other licenses. Though the open MIT license is the most popular across Github; Gen AI developers are predictably keen on securing their work with more protective licensing.

Market Map: Multiple projects /startups emerging across the Gen AI tech stack

Open Source Al | Market Map

Open source development is highly active across the Generative AI tech stack, especially in foundational models & developer tools.

	Open Source	Closed Source (Illustrative)			
Foundational Models Core AI models that form the bedrock for creating or nterpreting complex data, like text, images, and audio	Mat_ ADEPT Stability.ai II) (Cerebras Cecco (LMSYS) (Cerebras Cecco (LMSYS) (Cerebras Mosaic ^M (INO21//2) (Cecutherel Contextual-ai together.ai Lighton (PlayHT) (Craiyon) (Cecutherel	OpenAl ANTHROPIC Scohere Gemini M perplexity Al21 abs IIElevenLabs Midjourney			
Model Deployment & Inference Cloud-based services providing scalable computing resources (GPUs), essential for model hosting, Al training & performance	lFeplicate ₿ baseten መ Lepton Al @Kserve Arrikto	CoreWeave 📐 Lambda 🥢 Mode			
Developer tools Platforms/toolkits that support the building and operation of AI models, such as application frameworks, vector databases, and data management systems.	 Hugging Face mongoDB. Supabase zilliz comet werviete mindsdb HOW/SO[®] werviete Chroma mindsdb ADOW/SO[®] Chrom	Pinecone Kindo			
Model Training & Finetuning Tools for generating synthetic data for model training, labeling data, and refining pre-trained models for specific tasks, business functions or data sets.	anyscale (; Weights & Biases C DOMINO Snorkel gretel TOMIC Surge Al Synthesis.ai hazy OpenPipe OPyTorch	SCOLO Concolet Concolet Score Rendered.Al			
Monitoring & Observability Tools for tracking AI system performance, managing LLM operations, and analyzing user engagement	DATADOG Arrize Arthur Image: Strate in the stra	Amplitude @ MOBULIT INTELLIGENCE vellum			

Foundational models and developer tools, the core stack of AI, are the focus areas for new startups

Over 60% of new companies in the open source AI space are focusing on foundational models and developer tools, the core elements of the AI stack. This is expected, given that these components are fundamental for building, deploying, and managing generative AI applications across various use cases. Innovation in other areas like model training, fine-tuning tools, monitoring tools, and cloud computing services primarily revolves around these core AI stack elements.

High-quality open source AI reducing reliance on proprietary big tech AI, but data is key

The volume and quality of open source AI is sufficiently robust, enabling developers and startups to effectively compete with proprietary solutions. <u>OS model Mixtral 8×7B surprassed</u> <u>closed source GPT 3.5 in chatbot as well as holistic performance</u>. Other OS models like Llama, Yi are not far behind.

However, a crucial advantage that big tech firms with closed systems hold is their access to extensive data resources. This is evident in the fact that some recent OS models such as Llama-2 or Mistral 7B do not open source their training data. Data is likely to be the key proprietary element in the space.

Funding Landscape: Robust funding in 2022-23; foundational models & training tools secure maximum dollars

Gen Al infrastructure, due to its heavy reliance on vast amounts of data, extensive research, and substantial compute power, requires significant capital investment, has led to larger funding rounds compared to typical enterprise solutions.

Last funded in 2022-23	ast funded before 2022
Series C	
Anyscale	259M
Weights & Biases	250M
Snorkel (snorkel.ai)	135M
Series B	
Adept AI Labs	415M
Supabase	116M
Zilliz	113M
Comet.ml	70M
Weaviate	68M
Gretel (gretel.ai)	68M
Arize Al	61M
ArthurAl	60M
Deci Al	55M
Replicate	50M
Fiddler Labs	45M
Tonic (tonic.ai)	45M
Truera	42M
Seldon	34M
Series A	
Mistral AI	658
Diveplane	34M
Luma AI (lumalabs.ai)	26M
Iterative.ai	25M
Surge Al	25M
Unstructured (unstructured.io)	25M
Synthesis Al	22M

Robust funding activity in 2022-23; foundational models and model training software secured maximum dollars

75% of open source AI startups secured funding in 2022-23. Foundational models and model training/fine-tuning software have attracted >70% of the investment dollars.

Nvidia, a leading graphics chips manufacturer for AI, has been a strategic investor in this space, with investments in top startups like Mistral AI and Adept AI.

Foundational Models: Open source models are catching up with closed source in popularity and performance

Gen Al infrastructure, due to its heavy reliance on vast amounts of data, extensive research, and substantial compute power, requires significant capital investment, has led to larger funding rounds compared to typical enterprise solutions.



Open source LLMs Falcon and Bloom have received significant engagement

Falcon, a large language model (LLM) developed by Abu Dhabi's Technology Innovation Institute, and BLOOM, created by the collaborative research organization BigScience, have recorded the highest downloads on Huggingface – surpassing Meta's Llama-2.

Launched recently, Mistral Al's models Mistral 7B and Mixtral 8×7B have gained significant popularity, surpassing many established models on Huggingface in terms of traction.



Open source AI models are not far behind closed source models

Although closed source big models like GPT4 and Claude are at the top of the chatbot leaderboard, open source models like Mistral, Vicuna, Yi, Llama are catching up – which bodes well for the ecosystem.

However, closed source models are still a step ahead according to the MMLU benchmark, which tests knowledge and problem-solving skills across 57 subjects in humanities, social sciences, and STEM. MMLU measures a model's comprehensive performance, and in this context, closed source models like GPT and Gemini continue to outperform open source alternatives.

Open source development is leading to higher efficiency in models

Startups working with open source AI, who don't posses the extensive data resources or compute power of major tech firms, are motivated to create more efficient models that deliver high-quality results with less computational demand.

Mixtral 8×7B, an 85B-parameter 'mixture of experts' model that operates with the compute power of just a 14B model. It has outperformed all other open source models, including the larger Llama-2 70B, in terms of efficiency and effectiveness. This will be crucial in making these models more accessible to local applications (e.g. voice assistants on mobile).

Github Traction: AutoGPT, Mojo attracting significant developer interest

Gen Al infrastructure, due to its heavy reliance on vast amounts of data, extensive research, and substantial compute power, requires significant capital investment, has led to larger funding rounds compared to typical enterprise solutions.



AutoGPT, Modular's Mojo are witnessing high developer traction

As a primary platform for developers to interact with and contribute to open source Al projects, GitHub activity tends to be a strong indicator of traction. GitHub stars (similar to a "follow" on social media) are a direct indicator of a project's popularity on GitHub.

<u>AutoGPT</u>, an autonomous AI assistant built on GPT4, has received significant developer traction. The model is capable of acting as an AI agent, breaking a large task into various sub-tasks without the need for user input, which are then chained together and performed sequentially to yield a larger result. AutoGPT is also capable of connecting to the internet, thereby allowing for up-to-date information retrieval for its tasks.

<u>ModularML's Mojo</u> is a variation of Python tailored for high-performance AI applications, balancing the efficiency of languages like C++ and Rust with Python's simplicity. Mojo's core goals are to streamline AI development, integrate AI/ML infrastructure seamlessly, and deliver robust performance.



Pytorch, Huggingface, AutoGPT, and Supabase stand out with the significant engagement on Github

Github contributors are developers who make changes (known as "commits") to the code, actively engaging with the repo to improve it. Contributors are indicative of serious developer activity on repos.

The year-over-year GitHub analysis for 2022-23 highlights a notable uptick in both interest and active engagement with various repositories among the developer community.

Though AI agents are still in an experimental phase from a customer-facing POV, GitHub data reveals substantial developer activity in this domain, which could likely lead to some agent-based AI apps emerging soon. AutoGPT, an AI agent repo, is experiencing significant developer activity on GitHub. Other AI agent repos like Bloop, X-Agent are also seeing similar interest from the developer community.

Looking ahead

- Open source isn't merely a playground for Gen AI, its at the forefront of innovation Open source AI is seeing active innovation - Github saw a 148% annual growth in contributors and a 248% annual growth in Gen AI projects in 2023, HuggingFace has 400k+ models. Open Source stack for Gen AI is competitive or better than proprietary products across categories from foundational models to infra & tooling.
- Open Source models are not far from flagship proprietary models in performance and are leading in efficiency, achieving this performance with lower compute & data
 Open source models like Mistral, Vicuna, Yi, and Llama are rapidly catching up to closed source leaders like GPT4 and Claude, with Mixtral 8×7B even surpassing GPT3.5 in Elo and MMLU ratings. Open source development is fostering more compute-efficient models, which will be crucial for deploying Al locally on devices (e.g. mobile phones).

Access to high quality, abundant data will be the limiting function for OS AI models
 Data will be a key battleground for the development of large models. Recent models, such as
 Llama-2, Mistral 7B, which we released as "open source", have chosen not to make their model
 training data publicly accessible. Big Tech, of course, will have a significant advantage on data.
 Synthetic data platforms (like Gretel) can potentially augment training and fine-tuning, but
 expect data-protectionism to increase (NYT vs GPT is a case in point).

• Al agents are seeing significant developer activity, expect killer agent-based applications on the market soon

While AI agents are still largely experimental & nascent in customer-facing applications (see <u>our</u> <u>article on productivity tools</u>), Github data indicates serious & continuing developer interest in agents. There are 70+ AI Agent repos on Github as of today, with repos like AutoGPT, Bloop, XAgent getting significant traction (8-10K+ stars) and engagement (30+ contributors). Definitely an area to keep an eye out for.

• Expect standout open source AI projects to attract big rounds in 2024

Startups in open source AI have seen some extraordinarily large deals and active funding rounds across stages. Mistral AI obtained its unicorn status after a recent \$487M deal. <u>AutoGPT</u>, <u>Supabase</u> and <u>DeciAI</u> are poised for future funding rounds in the next 1-2 years.

About Synaptic

.

Synaptic is an alternative data and intelligence platform that helps investors, data teams and research analysts leverage alternative datasets across their investment processes.

For more information, visit us at <u>synaptic.com</u> or write to us at <u>contact@synaptic.com</u>

			Similar Companies						
개 Home Launcher post	X MailChimp X GitHub	□ × <mark>+</mark>	Citlab	Assembla B	itbucket	Rhodecode	Соріа Aut	>	
GitHub	github.com			•	• • • • • •				
Follow Add to list o List	 United States Tounded in 2008 GitHub provides code hosting services that a private projects in organizations. 	\$ Total Funding: \$350.0 Illow developers to build set	M oftware for open-source and	Gitlak	Asser	mbla Bitbu	cket Rhodec	ode Cod	ans ting Li
Overview							- • •		
O Github v	SECTORS Enterprise > DevOps > Software D SUBSIDIARIES Appcanary Atom Npm So	evelopment emmle		LISTS Synaptic	Portfolio				
a Alexa Engagement									
✓ Semrush Grow	th Index 🛛				Last updated 1	github.e	om per million ()		
Products	Fast Growth High Confidence	+0.09 MoM	Distribution across all companies in Er	terprise and Merged or Acquir	ed	5,11 Jun 20	0		
G Glassdoor	1.5 1.2 0.9 0.6								
Employees	0.3 0 Aug 21	Jul 22	8 ¹⁰³	84 14	1 2	5 Jan 20			
Job Openings	Ranked 3 among 8 competitors	50122		Detail	ed Growth Sumr	mary Ranked	1 among 4 competit	tors ~	
🏽 Tech Mentions 🛛 NEW 🗸 🗸	Employee Count O GitHub		Active Job Count ^① GitHub			Orgar ● US	nic Traffic ⁽⁾ github.com		
G Google Trends	2,319	+3% MoM	77		+60% M	^{юм} 11,5	13,723		
cb Crunchbase	Jun 20	+39% YoY	Oct 20			Jun 20			
Notes			~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~			~		•	
	Jan 20 Ranked 1 among 4 competitors ∽	Jun 20	May 20		0	ct 20 Jan 20 Ranked	1 among 4 competit	tors ~	
	Adwords Traffic ©		S ≝ US 👗 Github						
	US github.com	+94% MoM			IN TOP 100	Overall			
	Jun 20	-95% YoY							
8 Add to Affinity									
Report error Wrong app or website?	Jan 20	Jun 20							
	Ranked 3 among 4 competitors 🗸			_	😸 Outlie	r MoM declir	ie in Page viev	vs per millio	on users
					Github.c	om			
					638.51	K +1% MoM +1	1% YOY		