

Andreas Liesenfeld, Assistant Professor, Radboud University State of Open: The UK in 2024

Phase Two: The Open

Manifesto Report



Andreas Liesenfeld, Assistant Professor, Radboud University

Thought Leadership: Towards meaningful open AI: key findings and societal implications

The past year has seen a steep rise in generative AI systems that claim to be open. But how open are they really? The question of what counts as open source in generative AI takes on particular importance in light of the EU AI Act that regulates 'open source' models differently, creating an urgent need for practical openness assessment.

Recent work by Andreas Liesenfeld and Mark Dingemanse (Centre for Language Studies, Radboud University, The Netherlands) provides a stark view of the openness of current generative AI. In a systematic sweep of the landscape, they survey 45 text and text-to-image models that bill themselves as open. Key findings and contributions include the following:

- 1. Open-washing is widespread. Open-washing is the use of terms like 'open' and 'open source' for marketing purposes without actually providing meaningful insight into source code, training data, fine-tuning data or architecture of systems. Corporations like Meta, Mistral and Microsoft regularly co-opt terms like 'open' and 'open source' while shielding most of their models from scrutiny.
- 2. The EU AI Act puts legal weight on the term 'open source; without clearly defining it, creating an incentive for open-washing, opening up a key pressure point for lobbying, and necessitating clarity about what constitutes openness in the domain of generative AI.
- 3. Openness in Generative AI can be meaningfully measured if we think of it as composite (consistent of multiple elements) and gradient (it comes in degrees). The work distinguishes 14 dimensions of openness, from training datasets to scientific and technical documentation and from licensing to access methods.
- 4. Meaningful openness is possible, and is exemplified by some of the smaller players in the field, who go the extra mile to document their systems and open them up to scrutiny. We identify organisations like AllenAI (with OLMo) and BigScience Workshop + HuggingFace (with BloomZ) as key players moving the needle of openness in generative AI.

© OpenUK 2023 - OpenUK is a not-for-profit company limited by guarantee. Registered office address: 8 Coldbath Square, London EC1R 5HL



Why this matters for everybody

Although the EU AI Act creates a particular sense of urgency, openness is in fact of key importance for innovation and science, as well as to society. Here are three key points why openness in generative AI is important to society at large:

- Openness helps build critical AI literacy by demystifying capabilities. If a company claims their AI
 can 'pass the bar exam', it helps to know exactly what was in the training data. The sheer
 magnitude of training data for GPT4 means that it can effectively take the bar exam 'open book',
 which is much less impressive. As AI applications are fast becoming ubiquitous, openness about
 training data and fine-tuning routines helps us equip the general public with a critical
 understanding.
- 2. Openness grants agency to people subjected to automated decisions. A clear example is the LAION image dataset underlying a widely used image generator. The openness of this dataset enabled auditing by Dr. Birhane and others which brought to light the presence of problematic material and prompted revisions and retractions. pre that AI is often best seen as a shorthand for automation. Transparency about training data and model architecture makes models auditable, and can help to make biases visible.
- 3. Openness provides real benefits on the ground. Corporations like to hand-wave at 'AI safety' as a reason to keep systems under wraps, but this is mostly a thinly disguised attempt to obscure clear and present harms already underway, including increasing spam content in search engines and the rapid spread of misinformation. Only open systems can be responsibly used in science and education and can support a healthy open source ecosystem.

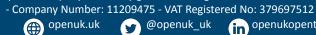
Evidence-based openness assessment can help foster a generative AI landscape in which models can be effectively regulated, model providers can be held accountable, scientists can scrutinise generative AI, and people can make informed decisions for or against using generative AI.



Open GenAl: LLMs (simplified table)

Open cols Life at	Project	Availability							Documentation				Access			
BLOOMZ		Open code LLM data LLM weights RL data			RL weights	License	Code Architecture Preprint									
AmberChat V V V V X V V V V V V V V V V V V V V		~	~	~	~	~	~	~	~	~	×	~	~	~	~	
Open Assistant v'	BLOOMZ	~	~	~	V	~	~	~	~	~	~	~	~	×	~	
OpenChat 3.5 78 V X V X V V V - V V V X V V V X V V X V V X V V V X X V V V V X X V V V X X V V V X X V V V X X V V X X V X V X X V X V X X V X X V X X X X X X X X X X X X X X X X X X X X	AmberChat	~	~	~	~	~	~	~	~	~	×	~	~	X	~	
Pythia-Chat-Base 7 V V V V V V V V V	Open Assistant	~	~	V	V	×	V	V	~	~	×	×	×	~	~	
Cerebras GPT 111	OpenChat 3.5 7B	~	X	~	×	~	~	~	~	V	~	~	×	~	~	
RedPajama-INCITE	Pythia-Chat-Base-7	~	~	~	V	×	~	V	~	~	×	~	~	~	×	
doly v v v v x v v x v x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x	Cerebras GPT 111	~	V	~	V	~	~	×	~	~	×	×	~	X	~	
Tulu V2 DPO 708 V	RedPajama-INCITE	~	~	~	V	~	~	~	~	×	×	~	~	X	~	
MPT-30B Instruct	dolly	~	~	~	V	×	~	V	~	~	×	×	×	~	×	
MPT-78 Instruct V	Tulu V2 DPO 70B	~	×	~	V	~	~	~	~	V	×	~	~	×	~	
trix V	MPT-30B Instruct	~	~	~	~	×	V	~	~	×	×	~	×	~	~	
NeuralChat 7B - X V V V V X X X X Vicuna 13B v 1.3 Vi - V X X X - V X V X X X X X X X X X X X	MPT-7B Instruct	~	~	~	~	×	~	V	~	X	×	~	×	~	×	
Vicuna 13B v 1.3	trix	~	V	~	~	×	V	V	~	X	×	×	×	~	~	
minChatGPT V V V V V V V V V V V V V V V V V V V V V V V V V V X X V V X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	NeuralChat 7B	~	×	~	V	~	~	~	~	×	×	~	~	~	×	
Chair Chai	Vicuna 13B v 1.3	~	~	~	×	×	~	V	×	~	×	~	X	~	~	
BELLE Geitje Ultra 7B X	minChatGPT	~	~	~	~	×	~	V	~	×	×	×	×	×	~	
Gelige Ultra 7B	ChatRWKV	~	~	~	ж	×	V	~	~	~	×	×	ж	~	~	
Phi 3 Instruct	BELLE	~	~	~	~	~	×	~	~	V	×	×	~	×	×	
WizardLM 13B v1.2 — X — V V — — X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Geitje Ultra 7B	×	~	~	V	~	×	×	~	~	×	~	~	×	~	
Airoboros L2 708 G X - V V X X X X X V V M Mistral 7B-Instruct - X V X X V X X X X X X X X X X X X	Phi 3 Instruct	×	×	ж	ж	~	~	ж	~	~	×	~	ж	~	~	
ChatGLM-6B - - V X X V - - X X X V Mistral 7B-Instruct - X V X - - X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	WizardLM 13B v1.2	~	×	~	V	~	~	~	~	V	×	×	×	×	×	
Mistral 78-Instruct - X X - - X X - - X X - - X X - - X X - - X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Airoboros L2 70B G	~	×	~	V	~	~	~	~	×	×	~	~	×	×	
WizardLM-7B - - X V - - - V X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X <td< td=""><td>ChatGLM-6B</td><td>~</td><td>~</td><td>~</td><td>×</td><td>×</td><td>~</td><td>~</td><td>~</td><td>×</td><td>~</td><td>×</td><td>×</td><td>×</td><td>~</td></td<>	ChatGLM-6B	~	~	~	×	×	~	~	~	×	~	×	×	×	~	
Owen 1.5 - X Y X - - X X X - - Y - - X X - - Y - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - X X - - - X X X X - - - - - - - - - - - - - - - - - - - -<	Mistral 7B-Instruct	~	×	~	×	~	~	×	~	~	×	×	X	~	~	
StableVicuna-13B - x - - - - - x x x - - - - - - x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x x	WizardLM-7B	~	~	×	V	~	~	~	~	V	×	×	×	×	×	
Falcon-40B-instruct X	Qwen 1.5	~	×	~	×	~	×	~	~	×	×	×	×	~	~	
UltraLM X X - X X - X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X </td <td>StableVicuna-13B</td> <td>~</td> <td>×</td> <td>~</td> <td>~</td> <td>~</td> <td>~</td> <td>~</td> <td>~</td> <td>~</td> <td>×</td> <td>~</td> <td>×</td> <td>×</td> <td>~</td>	StableVicuna-13B	~	×	~	~	~	~	~	~	~	×	~	×	×	~	
Yi 34B Chat - X V - X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X <td< td=""><td>Falcon-40B-instruct</td><td>×</td><td>~</td><td>~</td><td>~</td><td>×</td><td>~</td><td>×</td><td>~</td><td>~</td><td>×</td><td>~</td><td>×</td><td>×</td><td>×</td></td<>	Falcon-40B-instruct	×	~	~	~	×	~	×	~	~	×	~	×	×	×	
Koala 13B V - - X - - X X X X X X X X X X X X X X X X X X X X - X X - X X - X X - X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	UltraLM	×	×	~	V	~	×	×	~	V	×	~	~	×	×	
Mixtral 8x7B Instruct X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Yi 34B Chat	~	×	~	×	~	~	×	×	~	×	×	×	×	~	
Stable Beluga 2 X X - X - - X - - X - - X - - X - - X X - - X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Koala 13B	~	~	~	~	×	~	~	~	×	×	×	×	×	×	
Stanford Alpaca V X ~ ~ X ~ X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Mixtral 8x7B Instruct	×	X	~	ж	~	V	×	~	~	×	×	×	~	×	
Falcon-180B-chat	Stable Beluga 2	×	×	~	×	~	~	×	~	~	×	~	×	×	~	
Gemma 7B Instruct Orca 2 X X - X X X X X X X X X X	Stanford Alpaca	~	×	~	~	~	×	~	~	×	×	×	×	×	×	
Orca 2 X X - X X - - X - X - - X - - X - - X - - X - - X X - - X X - - X - - X - - X - - X - - X - - X - - - - X - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - <td>Falcon-180B-chat</td> <td>×</td> <td>~</td> <td>~</td> <td>~</td> <td>~</td> <td>×</td> <td>×</td> <td>~</td> <td>~</td> <td>×</td> <td>~</td> <td>×</td> <td>×</td> <td>×</td>	Falcon-180B-chat	×	~	~	~	~	×	×	~	~	×	~	×	×	×	
Command R+ X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Gemma 7B Instruct	~	X	~	×	~	×	×	~	~	×	~	×	X	×	
LLaMA2 Chat X X - X - - X - X - - X - - X - - X - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - <th< td=""><td>Orca 2</td><td>×</td><td>X</td><td>~</td><td>ж</td><td>~</td><td>×</td><td>×</td><td>~</td><td>~</td><td>×</td><td>~</td><td>×</td><td>X</td><td>~</td></th<>	Orca 2	×	X	~	ж	~	×	×	~	~	×	~	×	X	~	
Nanbeige2-Chat V X X V ~ X X X X X ~ X ~ X ~ X ~ X ~ X ~ X ~ X ~ X ~ X ~ X X ~ X X ~ X X ~ X X ~ X X ~ X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Command R+	×	×	×	V	~	~	X	×	×	X	~	×	×	×	
Llama 3 Instruct X X - X - X - X - X - X - X - X - X - X - X - X - X - X - X - X - X X - X X - X X - X X - X X - X - X - X - X - X X - X X - X X - X X - X X - X X - X X - X X - X X - X X - X X - X X - X X - X X - X X X X X X X	LLaMA2 Chat	×	×	~	×	~	×	×	~	~	X	~	×	×	~	
Solar 70B X X - X X X X X - X - X - - X - - X - - X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X	Nanbeige2-Chat	~	×	×	×	~	~	×	×	×	×	×	×	×	~	
Xwin-LM	Llama 3 Instruct	×	×	~	×	~	×	X	~	×	X	~	×	×	~	
	Solar 70B	×	×	~	ж	~	×	X	X	×	X	~	X	×	~	
ChatGPT X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X X </td <td>Xwin-LM</td> <td>×</td> <td>×</td> <td>~</td> <td>×</td> <td>~</td>	Xwin-LM	×	×	~	×	×	×	×	×	×	×	×	×	×	~	
	ChatGPT	×	×	×	ж	X	×	X	X	~	X	×	X	X	×	





© OpenUK 2023 - OpenUK is a not-for-profit company limited by guarantee. Registered office address: 8 Coldbath Square, London EC1R 5HL





First published by OpenUK in 2024 as part of State of Open: The UK in 2024 Phase Three "Open Source and Market Shaping"

© OpenUK 2024 © 0 0







