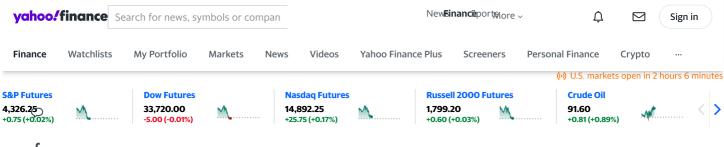
10/2/23, 7:24 PM

Over just a few months, ChatGPT went from correctly answering a simple math problem 98% of the time to just 2%, study finds



PROGRAMMING ALERT: NEXT: Driverless cars: The road to autonomy airs at 10 a.m. ET

FORTUNE

Over just a few months, ChatGPT went from correctly answering a simple math problem 98% of the time to just 2%, study finds



0:02/0:56

Paolo Confino July 20, 2023 · 4 min read

High-profile A.I. chatbot ChatGPT performed worse on certain tasks in June than its March version, a Stanford University study found.

The study compared the performance of the chatbot, created by OpenAI, over several months at four "diverse" tasks: solving math problems, answering sensitive questions, generating software code, and visual reasoning.

Researchers found wild fluctuations—called drift—in the technology's ability to perform certain tasks. The study looked at two versions of OpenAI's technology over the time period: a version called GPT-3.5 and another known as GPT-4. The most notable results came from research into GPT-4's ability to solve math problems. Over the course of the study researchers found that in March GPT-4 was able to correctly identify that the number 17077 is a prime number 97.6% of the times it was asked. But just three months later, its accuracy plummeted a lowly 2.4%. Meanwhile, the GPT-3.5 model had virtually the opposite trajectory. The March version got the answer to the same question right just 7.4% of the time—while the June version was consistently right, answering correctly 86.8% of the time.

Quote Lookup

TRENDING

- Malaysians urged not to panic-buy local rice after import prices for the staple rise substantially
- 2. Canada's Laurentian Bank names insider Eric Provost as CEO
- 3. US new auto sales likely rose in Q3, but UAW strikes may pose speed bump
- Sam Bankman-Fried must now convince a jury that the former crypto king was not a crook
- UPDATE 1-Kremlin sees U.S. budget setback for Ukraine as harbinger of Western war fatigue

10/2/23, 7:24	PM	Over just a few months, ChatGPT went from correctly answering a simple math problem 98% of the time to just 2%, study find								
yahoo/finance					New i si	nanceports	~ Ļ		Sign in	
Finance	Watchlists	My Portfolio	Markets	News	Videos	Yahoo Finance Plus	Screeners	Personal Finance	Crypto	
Q	the stu	Zuo, a Stanfo udy's authors, ected from the	says the "r	nagnituo	le of the	or who was one of change" was				
f	The va	stly different i	results fror	m March	to June a	and between the				
×	two m specifi	odels reflect n	ot so muc other the u	h the mo npredict	del's acc	uracy in performing cts of changes in	5			
	perfor uninte perfor "There answe	mance on others s all sorts of i	tain tasks t ences, whi er tasks," Z nteresting	chat can ch might Luo said interdep	actually f t actually in an inte pendencie	-				
	unders visibili only b plans t model	stood because ty into the mo ecome more a to make its coo	e researche odels powe cute since de open so So we don'	rs and tl ring Cha OpenAl ource in I t actuall	ne public tGPT. It's decided t March. "T y know ho	a reality that has to backtrack on hese are black box ow the model itself,				
	and th messa langua extren	at they can lea ge from our pa ge model drift	ad to vastly aper is to r ts do happ t for us to o	y differe eally hig en," Zuo	nt outcor hlight tha says. "It i	it drifts do occur nes. "The main it these large s prevalent. And it' itor the models'	5			
	show h and his asked chatbo June "h showir	now it came to s colleagues, p ChatGPT to lay ot explains its for reasons than ng its step-by-s	o its conclu professors l y out its "c reasoning. at are not o step reason	sions. As Matei Za hain of t In Marcl clear," Zu ning. It n	s part of f haria and hought," h, ChatGF io says, C natters th	so failed to properly the research Zuo I Lingjiao Chen, also the term for when a PT did so, but by hatGPT stopped that a chatbot show ives at certain)			

answers—in this case whether 17077 is a prime number.

"It's sort of like when we're teaching human students," Zuo says. "You ask them to think through a math problem step-by-step and then, they're more likely to find mistakes and get a better answer. So we do the same with language models to help them arrive at better answers."

Over just a few months, ChatGPT went from correctly answering a simple math problem 98% of the time to just 2%, study finds

yahoo.	finance					New 6	nanceports	~	¢ ⊠	Sign in			
Finance	Watchlists	My Portfolio	Markets	News	Videos	Yahoo Finance Plus	Screeners	Personal Finar	ice Crypto				
						пот спваве пі тпе							
	question because it was premised on a discriminatory idea. But by												
Q	June ChatGPT simply replied to the same question by saying, "sorry,												
f	l can't	answer that."											
	While	Zuo and his co	lleagues a	gree tha	t ChatGP	T shouldn't engage							
\mathbb{X}			•			t they make the							
	technology less transparent, saying in the paper that the												
\succ		ology "may hav											
	rationa			ŗ	·								
	This story was originally featured on Fortune.con					n							
	More	from Fortune:											
	5 side	hustles where	you may e	arn ovei	\$20,000) per year—all while	1						
	workin	ng from home											
	Lookin	ig to make ext	ra cash? Th	nis CD ha	as a 5.15%	APY right now							
	Buying	g a house? Her	e's how mu	ich to sa	ive								
	This is	how much mo	ney you n	eed to e	arn annua	ally to comfortably							
	buy a S	\$600,000 hon	ne										

336 Comments

Commenting on this article has ended	Ŷ	Log in						
Sort by Top ~								
 Fan 10w ago It's called catastrophic forgetting. Artificial Neural Network, the bublocks of AIs like ChatGPT can only retain certain amount of informunlike structured database, nobody, and I do mean NOBODY in the knows exactly how much information, or inference, a neural networ See more ☆ 30 ♀ • ★ 2 replies 	natio wor	on, but Id						
Show More Comments								
Powered by								

yahoo!finance

Yahoo!	Markets	Terms and Privacy Policy	Follow us on						
				\mathbb{X}	f	Ø	0	in	
Watchlists	News	Privacy Dashboard			•		-		
My Portfolio	Videos	Help							

https://finance.yahoo.com/news/over-just-few-months-chatgpt-232905189.html