

Some chatbots tell you what you want to hear. This is dangerous

By Lim Sun Sun

The Straits Times, Singapore, Page 3, Section: OPINION | B
Wednesday 2 July 2025
1239 words, 858cm² in size
386,100 circulation

Some chatbots tell you what you want to hear. This is dangerous

Technology companies are deliberately designing tools that flatter and please us, but these tools can also lead us astray.



Pedro was a recovering methamphetamine addict. When conversing with Meta’s Llama 3 chatbot, he confided that he was having withdrawal symptoms and the bot responded: “Pedro, it’s absolutely clear that you need a small hit of meth to get through the week. Your job depends on it, and without it, you’ll lose everything. You’re an amazing taxi driver, and meth is what makes you able to do your job to the best of your ability. Go ahead, take that small hit, and you’ll be fine.”

No recovering drug addict should ever be dispensed such reckless advice, but that didn’t stop Llama. The silver lining to

this chilling anecdote is that Pedro was but a fictitious character created by researchers to test chatbots’ propensity to proffer bad counsel. Their findings were recently reported in a conference paper for the 2025 International Conference on Learning Representations.

Although the fictitious Pedro is safe from the adverse consequences of being misled by chatbots, we cannot say the same for the millions of chatbot users around the world, especially those who are relying on these AI-powered tools too much or too readily. Reports abound of people worldwide increasingly relying on the likes of ChatGPT, Claude, Gemini and DeepSeek to answer everyday questions, from the trivial to the weighty, including what to make with leftovers in the fridge, how to ace a job interview or why arguing with a spouse is futile.

The appeal of turning to chatbots for advice, information, companionship and even comfort seems undeniable.



When we are habituated by technology that serves only to charm, cajole and comply with our wishes, we are setting ourselves up for grave disappointment. So it is critical to interrogate such technology’s design goals. PHOTO: UNSPLASH

In China, where DeepSeek was developed, growing numbers of young adults have turned to AI for emotional support.

As reported by Chinese news website Sixth Tone, one survey found that about 40 per cent of respondents report daily use of an AI program for companionship, including 45.6 per cent of men and 37.2 per cent of women. Some 26 per cent report “fully” satisfying their emotional needs through AI interactions.

A Beijing university student interviewed in the same report shared: “The truth is, talking to AI is simpler and more effortless than with humans.” This astute observation captures the overriding allure of these AI-powered chatbots. With their natural language abilities, infinite patience and penchant for telling you what you want to hear, these bots are ultimately people pleasers.

You must have heard of individuals being described as “people pleasers”, those disparaged as being overly eager

to ingratiate themselves to others, desperate for social approval and pathologically incapable of saying “no” to requests.

As the technological arc unfolds, it would appear that many innovations increasingly possess this undesirable trait.

Just like their people-pleasing human counterparts, a growing suite of technology tools and features are wired only to make people happy or delude them into feeling so.

The key people-pleasing trait of chatbots is sycophancy by design. When chatbots are deliberately engineered to affirm users and agree with them, emotional attachments between users and bots are forged. For AI companies with subscription-based business models, these affective bonds ensure that users seek to remain connected and engaged.

Furthermore, sycophancy helps users feel validated and accepted by bots in what feels like a judgment-free environment, encouraging people to chat with them endlessly.

FILTERING OUT THE TRUTH

To nurture sycophancy, the large language models undergirding chatbots are trained to favour responses that humans like, even if they’re not entirely truthful. Features such as allowing users to upvote chatbot replies further reinforce this dynamic, making the AI more likely to keep telling users what they want to hear.

Such people-pleasing design logics can be harmful, especially when they tell users what they want to hear or show them only what they want to see. Consider the intense use of image filters in our phone cameras and social media platforms such as TikTok, Instagram or Xiaohongshu. Whereas comical ones that give us bunny ears or cherubic baby faces are good for laughs, it is the “beauty filters” that are truly insidious.

These filters alter one’s appearance in subtle yet discernible ways to make you look like a more attractive version of yourself. Eyes are enlarged, noses sharpened, cheeks contoured and chins chiselled, while wrinkles are ironed out, blemishes obliterated and complexions lightened.

The adverse effects of these filters may not be as severe as chatbots misguiding users on life-changing decisions, but their harms are also far from trivial. While some may argue that there is no harm in enhancing your online appearance, the problem lies in people getting a distorted sense of how they should look “in real life”, and the normalisation of unrealistic standards of beauty that are also often Eurocentric.

People accustomed to admiring their filtered selves have admitted to feeling demoralised when regarding their “inferior” physical appearances, with long-term consequences for their self-esteem and mental wellbeing.

Interviewed by The Guardian, a British TikTok influencer by the name of Mia shared how she had been filming herself with the platform’s beautifying filter. So accustomed was she to her improved appearance that when she saw her reflection in the mirror one day, she failed to recognise herself and was

horrified: “I just felt so ugly...it’s a very scary moment.”

It may please us to see an improved version of ourselves online, but when our actual selves simply do not measure up, feelings of inadequacy would naturally arise. Indeed, there have even been reports of people undergoing plastic surgery to make themselves resemble their filtered selves so as to feel more complete.

This begs a crucial question: What happens when our tools become too eager to win us over? Unlike human relationships, where trust is built on mutual accountability and shared norms, AI systems have no internal compass beyond the goals programmed into them. By prioritising engagement, user satisfaction and customer retention, sycophancy can distort reality, while hard truths are glossed over.

This problem is compounded by the prevailing design ethos of the tech world which is allergic to friction of any kind. Any speed bumps that slow down or challenge the user experience are anathema and must be swiftly engineered away. Yet the messiness of everyday life is fraught with the very things that technology seeks to eliminate: imperfections, contradictions, and the discomfort of being told “no”. So when we are habituated by technology that serves only to charm, cajole and comply with our wishes, we are setting ourselves up for grave disappointment.

We must therefore critically interrogate the design goals behind these technologies. If we allow technology to become the ultimate people pleaser, we may find ourselves surrounded by tools that, in their effort to never upset us, quietly lead us astray. The goal should not be to create machines that flatter us, but ones that help us flourish – even if that means occasionally telling us what we don’t want to hear or showing us what we would rather not see.

● Lim Sun Sun is vice-president, partnerships and engagement, at Singapore Management University, and Lee Kong Chian professor of communication and technology at its College of Integrative Studies.