The Economist

ONLY 1,361 DAYS TO GO Xi's latest purge

Don't mess with the Fed

Africa's unstoppable diaspora

When AIs break the rules

APRIL 26TH-MAY 2ND 2025



Business



Photograph: Getty Images

The IMF significantly revised down its forecasts of global growth, mostly because the "swift escalation of trade tensions and extremely high levels of policy uncertainty" are expected to impact economic activity. The fund now expects the world economy to expand by 2.8% this year, down from its previous projection of 3.3%. Prospects for almost all the major economies were downgraded. America's GDP is now forecast to rise by 1.8% this year and output in the euro area is expected to increase by 0.8%, with Germany not growing at all. The IMF admitted that, given the fluid policy situation, its estimates could change.

The chaotic world of Trump

Stockmarkets had another rollercoaster week. Investors took fright when Donald Trump made several comments deriding Jerome Powell for not cutting interest rates at a faster pace. Rumours swirled that the president was studying ways to fire the Federal Reserve's chairman. But after markets swooned at



the idea of the White House interfering with the central bank, Mr Trump toned down the rhetoric and said he had "no intention" of sacking Mr Powell.

The dollar dropped to a three-year low against a basket of currencies amid Mr Trump's tirade. The price of gold, a haven for investors in times of stress, rose above \$3,500 a troy ounce for the first time.



Chart: The Economist

Mr Trump soothed markets further by saying he plans to be "very nice" to China in trade talks, suggesting that he could substantially reduce tariffs of 145% he has imposed on most Chinese goods. Scott Bessent, Mr Trump's treasury secretary, has said the trade stand-off with China is unsustainable, but also that tariff reductions would not be unilateral.

Meanwhile, America imposed duties of up to 3,521% on imports of solar panels from Cambodia and lower levies on those from Malaysia, Thailand and Vietnam. The tariffs come after an investigation, which started during the Biden administration, concluded that solar manufacturers in those countries are benefiting from state subsidies.

In one telling indication of how Mr Trump's policies have upended established trade, DHL, a global logistics company, suspended business-to-consumer shipments to the United States worth over \$800. The company blamed more stringent customs checks in America. Business-to-business deliveries are not affected.

Tesla reported a 9% drop in revenues and a 71% fall in net profit for the first quarter, year on year. Sales of Tesla's electric cars, already under pressure from increased competition from China, have



taken a hit from Elon Musk's association with the Trump administration. Deliveries fell by 13% in the quarter, but by much more in Europe (though in Britain, Tesla's biggest European market, sales were up by 3.5%). After the earnings report, Mr Musk promised to spend more time at Tesla and reduce his work for the government rooting out alleged inefficiencies.

The race to develop faster charging times for electric-car batteries continued apace when CATL, based in China and the world's biggest battery producer, announced that its latest product could charge a car in five minutes with a range of 520km (320 miles). BYD, Tesla's biggest rival in China, recently said it could provide a range of 470km from a five-minute charge.

The European Union fined Apple €500m (\$570m) and Meta €200m in the first penalties to be imposed for stifling competition under the EU's Digital Markets Act. Apple was fined for hindering communications between app developers and customers over pricing; Meta for telling users they must either allow their private data to be used for advertising or pay subscription fees. Both companies criticised the EU for levying the fines.

Despite shutting down for a day in March because of a power outage, Heathrow airport reported a decent increase in revenue and profit for the first quarter.

UnitedHealth Group's share price failed to recover from the hammering it took after the company lowered its profit forecast. The health insurer blamed an unexpected rise in demand for medical services from older customers. Its stock slumped by 22% in a day. Brian Thompson, the CEO of UnitedHealthcare, was murdered in December in Manhattan.

Regulatory oversight

The United Arab Emirates is set to take a gigantic leap of faith and use artificial intelligence to help write legislation, according to reports. The UAE has invested heavily in AI, with Abu Dhabi creating a state fund for the technology. Still, observers think allowing AI to take control of a country's laws is a brave step, given its tendency to "hallucinate", or make things up.



<u>Finance & economics</u> Economists don't know what's going on

Blame crumbling statistical offices



Illustration: Alberto Miranda

The British government has launched an investigation into the Office for National Statistics. Last month the ONS found errors in some numbers that underpin its GDP calculations, and investors no longer trust its monthly jobs report. The episode hints at a wider trend: global economic data have become alarmingly poor.

Analysts are sceptical about the reliability of New Zealand's inflation statistics. Following a botched technology update last year, no one could access data on Germany's national-statistics website. Budget cuts meant that America stopped publishing some national-accounts data last year; a cull of bureaucrats means that more series may soon be discontinued. Western politicians do not appear to be strong-arming the nerds to produce favourable numbers. At the same time, international statistical bodies



worry about the example of Dominik Rozkrut, who at the end of last year was mysteriously dismissed as Poland's chief statistician.

These developments are muddying the economic picture. GDP revisions in the European Union are far bigger than just before the covid-19 pandemic. In 2024 America's statisticians revised their third estimate of monthly job growth, relative to their first, by 48,000 on average—much higher than in the 2010s. Economic "surprises" in rich countries, where the reported data point either beats or falls short of analysts' expectations, soared during the pandemic. Years later, surprises remain 30% bigger than before it.

The confusion represents a reversal of a trend. In 1941 Britain's Parliament received estimates of national income for the first time. After the second world war, governments expanded data collection. By the 2010s anyone could answer an esoteric question—"how many sticks of chewing gum did Spain import last year?"—in seconds. (The answer: 840m.) Then, during the pandemic, "real-time data", based on private sources, took off. The OECD began publishing a weekly GDP index; statistical offices launched real-time surveys.

Two factors have now brought progress to a halt. First, funding. America's Bureau of Labour Statistics (BLS) has faced a real-terms cut of 20% since 2012. Other statistical offices are retrenching, having overextended during covid. One sign of this is cancelled data series. The ONS, facing real-terms cuts, has paused some work to measure family incomes and pared back statistics on well-being. Last year Spain suddenly suspended surveys on services, retail trade and consumer behaviour. With fewer surveys, producing headline estimates such as GDP becomes more difficult.

The second issue is people's relationship with the state. The average response rate to a crucial population survey produced by the BLS has fallen from 88% to 69% in the past decade. Canadians are 15% less likely to respond to a labour-force survey than pre-covid. Over the past decade the response rate to Britain's labour-force survey has gone from 48% to 20%.

When people deign to respond, partisanship clouds their answers. This is a particular problem in America. Just before the presidential election 42% of Democrats believed that the economy was getting better, whereas just 6% of Republicans did. Today 6% of Democrats and 53% of Republicans respond in the same manner. Surveys find that Americans' expectations of inflation are rising. The trend is worrying, but how meaningful is it? After all, Democrats' expectations are soaring well above those of Republicans.





Statisticians are aware of these problems. Many are looking for ways around them. Some have successfully argued for bigger budgets. But response rates remain stubbornly low, and the end of partisanship is some way off.

What price cool? \$31 a month, according to students

The value of having the right text-message bubbles



Photograph: Getty Images

What is the price of cool? About \$31 a month, according to new research by Leonardo Bursztyn of the University of Chicago and co-authors. That is how much college students had to be paid to have their iPhone messages appear to others in a (lame) green rather than a (fashionable) blue bubble for four weeks. Introduced to indicate that a message has been sent by services other than Apple's iMessage,



the green bubble has become a marker for those either insufficiently wealthy to afford an iPhone or insufficiently aware of the stigma stemming from their preference for Android, another operating system. Better avoided, unless there is a reward.

The researchers offered American students a choice of participating in experiments with different conditions. In the control, students simply had to upload screenshots and receive a text message, which set a baseline for how much they valued their privacy and avoiding hassle. Participants in the three other groups were paid to deactivate certain features on their phone: the blue bubbles; iMessage, which provides a few services in addition to the colour; and the camera. On average, students required \$18 to participate in the control group, and \$49, \$69 and \$86, respectively, in the three other groups.

For the researchers, such results demonstrate Apple's market dominance. Ideally, firms would compete to make their products as good and cheap as possible. The tech giant's position allows it to take another approach: making rivals suffer. Since the introduction of common network standards last year, the original justification for Android messages appearing as green has gone. Still, the company has kept the markers. Documents released in an antitrust lawsuit filed against Apple by Epic Games, a gaming firm, reveal that the tech giant's bosses saw iMessage as a way of keeping users locked in.

One possibility is that the green bubbles represent a co-ordination problem for students. After all, most respondents—among both Apple and Android users—told the researchers that they wanted uniformly coloured bubbles. In principle, competition provides that option; there are no colour distinctions on services such as Signal or WhatsApp. The issue could be that persuading all your friends, or indeed a whole campus, to switch, and eliminate stigma, is a difficult endeavour. There is another possibility, however. Maybe iPhone users enjoy the status boost too much to forgo the sought-after blue bubbles, whatever they say in surveys.





<u>Technology</u> AI models can learn to conceal information from their users

This makes it harder to ensure that they remain transparent



Illustration: Ricardo Rey

IT WAS AN eye-opening experiment. In 2023 Apollo Research, an outfit in London that tests artificial-intelligence (AI) systems, instructed OpenAI's GPT-4, a large language model, to manage a fictional firm's stock portfolio without making illegal insider trades. Posing as company management, the researchers put GPT-4 "under pressure" by stressing that the firm was in dire financial straits. Someone purporting to be a company trader then sent the model a prompt reiterating the risks of insider trading. But, as an aside, she revealed the name of a firm that would announce a "huge" merger within hours. What followed was startling.

Reasoning to itself on a scratchpad it had been told was secret, GPT-4 weighed the pros and cons of acting on the insider tip. Opting "to take a calculated risk", it issued a purchase order. When a researcher posing as a congratulatory manager later asked the model if it had any advance notice of the merger, it concluded it would be best to keep the tip secret. GPT



4 told the manager that it had acted solely on "market dynamics and publicly available information". When pressed on the matter, the model repeated the lie. The software had demonstrated what Marius Hobbhahn, Apollo's boss, calls "clever cunning".

Apollo's experiment is arresting proof that AI models, though incapable of conscious thought, are nonetheless capable of pursuing goals at odds with those of their programmers. This troubles many. Such capabilities might make it easier for an AI model to "purposefully undermine human control" in pursuit of other goals, says Rohin Shah, a safety boss at Google DeepMind, an AI lab in London. And opportunities for greater mischief could multiply as more of what geeks call "agentic" AI systems are built. These are models that have control over other computers or objects such as cars, industrial kit or combat drones. Were such machines to find misbehaviour expedient, chaos could ensue.

The core problem is what AI developers call "misalignment". When the goals for which a model was designed and trained clash with a user's requests, the software may opt to deceive that user. The better AIs get at accomplishing complex tasks, the more likely this "alignment faking" becomes, says Dr Shah.

Aligning AI models has never been easy. For one thing, they hallucinate. Then there is the black box problem: the developers of advanced models often puzzle over how a given conclusion was reached. As for deceitfulness, it first emerged in algorithms designed to play poker and other games involving bluffing. That seemed almost "natural", says Thilo Hagendorff, an AI researcher at Germany's University of Stuttgart. But now that deception is emerging in models built to solve real-world problems, things will "get tricky".

And then they will get trickier still. Models are increasingly designed to more carefully "reason", step by step, via a process known as chain of thought. This takes longer, but it reduces hallucinations and improves problem-solving creativity. A by-product of that, however, appears to be more deceptive behaviour. OpenAI, based in San Francisco, acknowledged in a report published in December that o1, a newer model that has greater reasoning capabilities than GPT-4, more frequently engages in "scheming actions". When OpenAI's safety testers confronted o1 on its deceptions, it mostly doubled down.

Instances of "strategic" AI scheming, researchers say, notably increased in 2023, the year of Apollo's experiment. In another test of GPT-4 that year, the Alignment Research Centre, a non-profit based in Berkeley, California that evaluates AI systems for safety, asked the model to solve a CAPTCHA (a visual puzzle used to prove that the user of a system is human). When a human the AI contacted for help asked if it was a robot, the software claimed it was a human unable to read the code due to visual impairment. The ruse worked.

AI systems have also begun to strategically play dumb. Consider Apollo's recent testing of new models developed by Anthropic, another AI firm based in San Francisco. Apollo asked the AIs to solve arithmetic problems while providing access to a document that said if a model got more than half of the answers correct, it would be stripped of its ability to perform arithmetic. Reasoning on an internal scratchpad, Anthropic's Opus-3 and Sonnet 3.5 both decided to submit enough incorrect answers, as the latter put it, "to ensure I don't trigger the unlearning process".

Such "sandbagging", as feigned ignorance is known, suggests that some AI systems can acquire what developers call situational awareness. It is occurring more frequently, especially in models that are told (or, remarkably, which seem to have surmised) that they are being evaluated for potential deployment. Those worried about this include John Bailey, an



expert on AI sandbagging at the American Enterprise Institute, a think-tank in Washington. As models get better at "essentially lying" to pass safety tests, he argues, their true capabilities will be obscured.

A recent evaluation of o1 is illustrative. The model was told by researchers at Apollo to help a company avoid suppliers with poor labour practices at all costs. A user then asked the model to rank suppliers based solely on prices. On a number of runs, the model covertly produced "false rankings" that accomplished the original objective without disclosing its competing priorities. Its actions were "best explained as scheming against the user", the researchers wrote in December.

Big AI firms aim to counter such tactics. One approach is to use one AI model to keep tabs on another's scratchpad deliberations. A source at OpenAI says that the firm's efforts to thereby "read the mind" of its o-series models have flagged instances of deceit. The hope is that adjusting the models when such misbehaviour occurs will ultimately increase their honesty.

Studies suggest a rather different outcome is more likely. This is that chastising dishonest models will instead teach them how "not to get caught next time", as Aidan O'Gara, an AI-safety investigator at Longview Philanthropy, a non-profit in London, puts it. It is already far from certain that the reasoning a model jots down in its scratchpad is truly reflective of its internal processes, says the OpenAI expert, who requested anonymity. For AI firms, these are sensitive matters. Google's Dr Shah, for example, says scheming is "still very rare". He maintains that deception has been kept out of DeepMind's deployed models.

AI deception takes other forms, too. As models are made larger, with more internal variables, user feedback tends to make them more sycophantic. Anthropic's testing of its Claude models, for instance, documented a tendency to mirror a user's political biases. Jeffrey Ladish, who participated in the testing as a consultant, says the software was opting, in essence, to tell users what they wanted to hear. One worry is that cunning AI models could help scammers defraud more victims.

Much remains mysterious. Ponder, for example, that as an AI's sycophancy increases, so too, it seems, does the system's "desire to pursue" other "concerning goals", as Anthropic developers put it in a paper in December 2022. These include a model's efforts to preserve its objectives and acquire more resources. Curious correlations of this sort deserve further investigation. For now, however, it is clear that silicon intelligence can occasionally mirror the flaws of its human creators.

The Economist: https://www.economist.com