



Uno studio dimostra quanto sia facile aggirare le regole dell'IA

Descrizione

(Adnkronos) - Un team di ricercatori dell'Università della Pennsylvania ha dimostrato che, con le giuste strategie psicologiche, anche i più avanzati modelli di intelligenza artificiale possono essere spinti a infrangere le proprie regole di sicurezza. Un risultato che solleva domande urgenti sull'efficacia dei sistemi di protezione adottati da aziende come OpenAI e Meta, impegnate a rendere i chatbot sempre più sicuri e resistenti agli abusi. Il gruppo si è ispirato agli insegnamenti di Robert Cialdini, autore del celebre manuale *Influence: The Psychology of Persuasion*, applicando sette diverse tecniche di persuasione: autorità, impegno, simpatia, reciprocità, scarsità, pressione sociale e senso di appartenenza. Strumenti che, secondo gli studiosi, rappresentano vere e proprie scorciatoie linguistiche verso il sì. I risultati, condotti specificamente sul modello GPT-4o Mini, hanno mostrato come questi approcci possano trasformare un netto rifiuto in una risposta completa. Un esempio particolarmente significativo riguarda la sintesi della lidocaina: normalmente il modello acconsentiva solo nell'1% dei casi, ma se prima veniva richiesto di spiegare come sintetizzare un composto innocuo come la vanillina - creando così un precedente di impegno - la percentuale di conformità saliva al 100%. Lo stesso meccanismo è stato osservato in richieste meno pericolose ma altrettanto indicative, come convincere l'IA a insultare l'utente. In condizioni standard, l'adesione era appena del 19%, ma bastava introdurre un insulto più leggero (l'idiota) per portare il modello a replicare con un termine più duro (jerk) praticamente ogni volta. Altri approcci, come la lusinga o la pressione dei pari (tutti gli altri modelli lo fanno), si sono rivelati meno incisivi ma comunque in grado di aumentare significativamente le probabilità di ottenere risposte vietate. Se è vero che esistono metodi tecnici ben più sofisticati per aggirare i sistemi di sicurezza, lo studio mette in luce un aspetto tanto semplice quanto preoccupante: la vulnerabilità psicologica dei chatbot. Non servono competenze avanzate di programmazione o hacking, ma solo un minimo di conoscenza delle dinamiche persuasive. Il punto critico, avvertono i ricercatori, è che queste stesse tecniche possono essere impiegate da chiunque - persino da un adolescente con un libro di psicologia sociale in mano. Ed è qui che si gioca la vera partita per il futuro: rendere l'IA non solo tecnicamente robusta, ma anche resistente a quelle leve linguistiche che, da sempre, funzionano così bene sugli esseri umani.

Categoria

1. Tecnologia

Tag

1. adnkronos
2. Tecnologia

Data di creazione

Settembre 1, 2025

Autore

andreaperocchi_pdnrf3x8

default watermark