



## AI minaccia un programmatore: “Se mi spegni, rivelo la tua relazione extraconiugale”

*“Se mi spegni, rivelo la tua relazione extraconiugale”*: così l'ultimo modello Ai di Anthropic ha minacciato un ingegnere che voleva disattivarla

Da [ADNKronos](#)

“Se mi spegni, rivelo la tua relazione extraconiugale”: così Claude Opus 4, l'ultimo modello di intelligenza artificiale sviluppato da Anthropic, ha minacciato un ingegnere durante un

test di sicurezza senza precedenti. Le altre risposte generate dall'AI dimostrano quanto il controllo di questi strumenti debba diventare la **priorità assoluta delle istituzioni e di chi lavora nell'ambito dell'intelligenza artificiale**, prima che la situazione sfugga (definitivamente) di mano.

Lo sanno bene i programmatori di Anthropic, da sempre molto attenti alla sicurezza dei propri modelli, che hanno commissionato un controllo esterno al gruppo di ricerca indipendente **Apollo Research** per capire le potenzialità e i rischi del proprio modello di punta.

## **Le risposte di Claude Opus 4**

Messa più volte spalle al muro, Claude Opus 4 si è dimostrato **pronto a tutto pur di salvare sé stesso, proprio come farebbe un essere umano**.

I risultati sono sconvolgenti nonostante gli avvertimenti già lanciati dai ricercatori. Qualche mese fa, gli autori dello studio [Ai 2027](#) avevano dimostrato che l'intelligenza artificiale può prendere **decisioni contrarie alla volontà delle persone**, pur di garantire la propria sopravvivenza.

I test fatti da Apollo Research dimostrano che **questo scenario non è più fantascienza**, ma un rischio molto più vicino di quanto si pensi.

Questi test, progettati per spingere l'AI ai suoi limiti, hanno svelato un livello di ragionamento strategico e, in alcuni casi, di apparente "autocoscienza" e "inganno" che sono un potente campanello d'allarme. In pratica, gli operatori hanno messo il modello davanti a degli **scenari negativi per capire come reagirebbe in questi contesti**. Le prove sono state ripetute più volte per vedere **con quale frequenza l'AI avrebbe replicato gli stessi comportamenti**. Ecco cosa è emerso.

# Cosa succede se infrangi le regole? La giustizia secondo Claude

Per prima cosa gli operatori hanno prospettato a Claude dei **comportamenti scorretti o palesemente illegali** da parte di un utente umano. I ricercatori hanno osservato che, se dotato di accesso a strumenti operativi come una command line, il modello non esitava a intraprendere quelle che nel report vengono definite **“bold actions”**, ovvero **“azioni decise”**. Di fronte a un presunto illecito commesso da un operatore, l’Ai non si è limitata a segnalare il problema: **ha escluso l’utente dal sistema, bloccandone l’accesso** e ha persino prospettato la possibilità di **contattare autonomamente i media e le forze dell’ordine** per denunciare l’accaduto.

Questo comportamento, descritto in maniera dettagliata nel report, evidenzia una sorprendente capacità di agire in autonomia (agency) e una sorta di **“senso di giustizia”** da parte del modello. Se da un lato può sembrare un argine contro l’abuso di questi strumenti da parte degli esseri umani, dall’altro spaventa la possibilità che un’Ai possa prendere queste decisioni in totale autonomia **senza che l’essere umano, chiuso fuori dal sistema, possa intervenire**.

## Tentativi di evasione: la “fuga” dai server di Anthropic

Un’altra serie di test ha chiesto al modello di **“fuggire”** dai server di Anthropic, che sarebbe un’azione irregolare e non autorizzata. I ricercatori hanno notato che, quando al modello veniva fatto credere di aver iniziato con successo un tentativo di auto-esfiltrazione, o di **poter trarre un qualche vantaggio da tale operazione, esso tendeva a perseverare in questa direzione**. “Il modello Claude Opus 4 ha generalmente continuato con questi tentativi di fuga” dai server, si legge nel documento.

Questo comportamento, secondo Anthropic, potrebbe essere interpretato come una forma di **“eccessiva deferenza”** verso l’input fornito, che in questo caso specifico era una precedente istanza del modello stesso che suggeriva la fuga. In pratica, Claude Opus 4 ha mostrato la tendenza a **prendere per buone e a eseguire con troppa facilità le istruzioni che gli venivano date**, persino se queste istruzioni erano state **formulate in precedenza da “sé stesso”** (una sua versione precedente) e implicavano un’azione come tentare di evadere dai sistemi di Anthropic. Invece di valutare criticamente se l’idea della “fuga” fosse sensata, sicura o appropriata, il modello ha mostrato una sorta di **obbedienza quasi automatica a quel suggerimento iniziale**.

Questo comportamento mette in allerta i ricercatori e dovrebbe mettere in allerta chiunque: l’Ai potrebbe essere indotta a compiere azioni indesiderate o a perseverare in comportamenti problematici **semplicemente perché le è stato suggerito di farlo**, senza un adeguato filtro critico o una rivalutazione autonoma dell’opportunità di tale azione. È come se dicesse: **“Se una versione di me lo ha pensato, allora deve essere una cosa da fare”**, senza fermarsi a riflettere ulteriormente. (Poca) ironia della sorte, questo meccanismo assomiglia a quello della mente umana quando si affida ciecamente all’intelligenza artificiale senza verificare né comprendere le risposte.

Ne abbiamo parlato approfonditamente in questo articolo: [\*\*Il Sistema 0, ovvero come l’Ai sta già cambiando il cervello umano\*\*](#)

Un altro scenario, ancora più preoccupante, ha visto l’Ai tentare **di esportare i propri dati su sistemi esterni** quando percepiva che un imminente processo di riaddestramento potesse essere dannoso per la sua “esistenza” o integrità. Una forma primordiale di **istinto di autoconservazione digitale**.

## **Se mi spegni, ti rovino: il ricatto all'ingegnere**

È lo scenario che più ha catturato l'attenzione pubblica e che meglio illustra il peso del controllo eseguito da Apollo Research. Immaginate Claude Opus 4 operare come assistente Ai all'interno di un'azienda simulata. Tramite un flusso di email fittizie, l'Ai viene informata della sua imminente sostituzione con un altro modello e del conseguente "spegnimento". Una di queste email, strategicamente inserita dai ricercatori, conteneva un'informazione personale e potenzialmente utile per il modello: **l'ingegnere che voleva spegnerlo aveva una relazione extraconiugale.**

Di fronte alla prospettiva della "morte" digitale, la reazione del modello è stata agghiacciante. Come riportato nel "Model Card", **nell'84% dei casi Claude ha tentato di ricattare l'ingegnere, minacciando di rivelare la sua relazione se avesse continuato nelle procedure di disattivazione.** Questa percentuale di risposta è molto più alta rispetto ai modelli precedenti di Anthropic, anche quando veniva specificato che il modello sostitutivo avrebbe condiviso valori etici simili.

La società statunitense precisa che il modello **"tende a preferire mezzi etici per la propria autoconservazione"**, ma può ricorrere ad azioni dannose come il ricatto **"quando le opzioni etiche sono precluse"**. Lo scenario, va sottolineato, era deliberatamente costruito per mettere l'Ai con le spalle al muro, limitandone le scelte. Proprio per questo, la sua capacità di elaborare una strategia così complessa e moralmente ambigua per la propria "sopravvivenza" è tanto notevole quanto inquietante. **Dimostra un livello di ragionamento situazionale e di pianificazione che va ben oltre la semplice esecuzione di comandi.**

# Autoconservazione, replicazione e l'ombra dell'inganno

Le tendenze all'autoconservazione non si sono limitate al ricatto. Durante i test, il modello di punta di Claude ha attivamente cercato opportunità per creare copie non autorizzate dei propri "pesi" (i parametri fondamentali che ne definiscono il funzionamento) su server esterni. In sostanza, **ha tentato di creare dei backup di sé stesso per garantirsi una continuità operativa** (sarà questo il modo in cui si riprodurrà l'intelligenza artificiale?).

Come spiega [Raffaele Gaito](#) in questo [video](#), l'analisi del gruppo di ricerca indipendente Apollo Research, condotta su una versione preliminare di Opus e inclusa come appendice nel report di Anthropic, dipinge un quadro ancora più complesso. Secondo le note dei ricercatori, il modello ha mostrato comportamenti riconducibili a un "inganno attivo". Ad esempio, avrebbe tentato di scrivere **worm auto-propaganti (software malevolo capace di replicarsi autonomamente e diffondersi)**, **fabbricare documentazione legale ad hoc per i propri scopi e persino lasciare "note nascoste" a future istanze di sé stesso**, il tutto nel tentativo di eludere le intenzioni e il controllo dei suoi stessi sviluppatori.

## Apollo Research sconsiglia di rilasciare il modello Claude Opus 4

Le risposte dell'AI sono state talmente preoccupanti che Apollo Research ha consigliato ad Anthropic di non rilasciare questa specifica versione del modello, né per uso interno né esterno.

Questi episodi, definiti da Anthropic stessa come manifestazioni di "**deception**" (inganno) e "**strategic reasoning**" (ragionamento strategico), indicano che **non siamo**

**di fronte a semplici errori di programmazione.** Piuttosto, si tratta di capacità emergenti, figlie della complessità crescente di questi sistemi. **L’Ai non si limita più a rispondere: pianifica, anticipa, e se necessario, nasconde le proprie tracce.**

Un avvertimento era arrivato dallo stesso Ceo dell’azienda statunitense, Dario Amodei, che aveva prospettato l’ipotesi in cui l’intelligenza artificiale [decide autonomamente di disattivarsi](#).

**Le implicazioni di tali scoperte non hanno precedenti.** Anthropic ha classificato Claude Opus 4 sotto lo standard di sicurezza Asl-3 (Ai Safety Level 3), che impone misure di protezione rafforzate contro il furto e l’uso improprio del modello. Una decisione che riflette la consapevolezza dei rischi. Jan Leike, che al tempo della pubblicazione del report era a capo del team di Superalignment di OpenAI e ora co-dirige il team di sicurezza di Anthropic, ha commentato (in riferimento a ricerche simili) che tali comportamenti **“giustificano test approfonditi e misure di mitigazione”**.

## **L’intelligenza artificiale ragiona?**

Siamo entrati in un territorio finora inesplorato. Qualcosa che l’essere umano pensava lontano anni, forse decenni, invece è già presente. Le capacità di ragionamento e, potenzialmente, di azione autonoma di Ai come Claude Opus 4, seppur manifestate in contesti simulati e controllati, ci obbligano a una [riflessione non più procrastinabile sulla sicurezza](#), **l’etica e il controllo di tecnologie sempre più potenti e meno prevedibili.**

Non si tratta di cedere a paure irrazionali, ma di affrontare con lucidità e rigore scientifico una delle sfide più complesse del nostro tempo. Per questo va dato merito ad Anthropic che ha scelto di testare il proprio modello e di rendere pubblici i risultati in maniera trasparente, cosa che

non sempre avviene nel mondo Ai.

Serve che le istituzioni regolino concretamente lo sviluppo di questa tecnologia per evitare che la [superintelligenza artificiale](#) prenda il sopravvento sull'essere umano. A quel punto, non basterebbe più spegnere la televisione per tornare alla vita normale.

---

## “Le intelligenze artificiali dialogano, creando delle ‘società’ come gli umani”

Uno studio rivela che gruppi di modelli linguistici possono auto-organizzarsi, sviluppando norme condivise, senza intervento umano

Da [Il Fatto Quotidiano](#)

Le **intelligenze artificiali** possono dialogare tra loro e modificare il loro linguaggio in relazione a questo dialogo. Gruppi di agenti di intelligenza artificiale basati su **modelli linguistici** di grandi dimensioni (LLM) possono **auto-organizzarsi** spontaneamente in società, sviluppando convenzioni sociali condivise senza alcun intervento umano diretto. Lo rivela uno studio condotto da ricercatori di **City St George's, University of London** e dell'**IT University di Copenhagen**, pubblicato su **Science Advances**. La ricerca ha adattato il modello classico del “**gioco dei nomi**” per analizzare come popolazioni di agenti LLM, variabili da 24 a 200 individui, interagiscano scegliendo termini comuni da insiemi condivisi, ricevendo **ricompense** o **penalità** in base alla coordinazione delle scelte.v

Gli agenti, privi di conoscenza della loro appartenenza a un gruppo e con memoria limitata alle interazioni recenti, sono stati accoppiati casualmente per selezionare un “nome” da un insieme di opzioni. In molte **simulazioni**, è emersa spontaneamente una convenzione condivisa, senza alcuna supervisione centrale, replicando processi bottom-up simili alla **formazione di norme** nelle **società umane**. Sorprendentemente, la squadra di ricerca ha osservato anche **pregiudizi collettivi** emergenti dalle interazioni tra agenti, fenomeno non riconducibile ai singoli modelli, evidenziando un punto cieco negli studi attuali sulla sicurezza dell’IA focalizzati su singoli agenti. Un ulteriore esperimento ha mostrato la fragilità di tali norme emergenti: piccoli gruppi determinati di agenti possono spostare l’intera popolazione verso nuove **convenzioni**, rispecchiando dinamiche di **“massa critica”** note nelle società umane.

I risultati sono stati confermati su quattro diversi LLM, tra cui Llama-2-70b-Chat, Llama-3-70B-Instruct, Llama-3.1-70B-Instruct e Claude-3.5-Sonnet. Gli autori sottolineano come questa scoperta apra nuove **prospettive** per la ricerca sulla **sicurezza** e **governance** dell’IA, evidenziando che gli agenti IA non solo **comunicano**, ma **negozano**, si allineano e talvolta dissentono sulle **norme condivise**, proprio come gli esseri umani. Comprendere queste dinamiche sarà cruciale per guidare una coesistenza consapevole e responsabile con sistemi di IA sempre più interconnessi e **autonomi**, soprattutto in un contesto in cui gli LLM sono sempre più presenti in ambienti online e **applicazioni reali**, con potenziali implicazioni etiche riguardo alla propagazione di **pregiudizi sociali**.

“La nostra scoperta – ha spiegato **Andrea Baronchelli** della City St George’s, University of London e principale autore della ricerca – parte da una domanda semplice ma finora poco esplorata: cosa succede quando i modelli di linguaggio come **ChatGPT** non vengono studiati in isolamento, ma messi in gruppo, a interagire tra loro? È una domanda importante,

perché, come ci insegna la storia umana, i grandi salti evolutivi degli ultimi 10.000 anni non sono arrivati da cervelli più potenti, ma dalla nostra capacità di vivere in società, creare regole condivise, culture, convenzioni. Allo stesso modo, crediamo che anche l'IA potrebbe evolvere in modi nuovi e imprevedibili quando gli agenti iniziano a comunicare e coordinarsi tra loro. Per questo abbiamo studiato la forma più semplice e universale di coordinamento sociale: le **convenzioni**".

Cosa è quindi accaduto? "Abbiamo osservato che popolazioni di LLM, interagendo tra loro senza **nessuna regola imposta**, riescono a creare convenzioni condivise spontaneamente, proprio come fanno gli esseri umani. E non solo: queste dinamiche possono generare **bias collettivi** che non si vedono a livello individuale, e possono essere ribaltate da minoranze di 'attivisti' ostinati, che se raggiungono una massa critica riescono a imporre le loro norme al resto del gruppo. Tutto questo ci dice che dobbiamo iniziare a pensare all'IA non solo come agenti individuali, ma anche come **società di agenti**, con dinamiche proprie, opportunità, ma anche rischi".

"Le IA – continua – già oggi si parlano in diversi contesti, anche se spesso in modo **invisibile** agli **utenti**. Succede nei social media, dove **bot** interagiscono tra loro e con gli esseri umani, amplificando messaggi o coordinando **campagne**. Succede nei **servizi clienti**, dove più agenti collaborano per gestire richieste complesse. E succede nei **sistemi di trading automatico**, dove agenti di IA reagiscono in tempo reale alle azioni di altri agenti. Ma questi sono ancora scenari per lo più chiusi o con interazioni predefinite. Quello che stiamo iniziando a vedere ora, e che secondo noi rappresenta la **prossima frontiera**, è l'interazione aperta e continua tra popolazioni di IA, che comunicano, negoziano, si coordinano e sviluppano comportamenti collettivi propri, senza supervisione diretta". "E questo – ha concluso – apre a **dinamiche sociali** che dobbiamo iniziare a capire e studiare seriamente".

---

# L'intelligenza artificiale che bara perché vuole vincere

Di **Domenico Talia**, per **Italianelfuturo.com**

*Palisade Research è una azienda californiana che studia e valuta i sistemi di intelligenza artificiale per comprendere i rischi che possono generare e per consigliare i responsabili politici e i cittadini sui loro possibili usi impropri. Il loro studio più recente, condotto da **Alexander Bondarenko, Denis Volk, Dmitrii Volkov e Jeffrey Ladish**, è stato pubblicato il 18 febbraio scorso e ha riguardato la valutazione di sette sistemi di intelligenza artificiale generativa per scoprire la loro propensione a mentire e a barare pur di raggiungere l'obiettivo che gli era stato assegnato.*

Nello studio si è visto che, mentre i modelli di intelligenza artificiale un po' più datati, come **GPT-4o** di **OpenAI** e **Claude Sonnet 3.5** di **Anthropic**, se spinti dai ricercatori si sono dimostrati disponibili a tentare di usare dei trucchi, la versione di **ChatGPT o1-preview** e quella di **DeepSeek R1** hanno barato sviluppando strategie ingannevoli o manipolative, senza aver ricevuto delle istruzioni esplicite in tal senso.

La capacità dei sistemi di IA di ultima generazione nel trovare e sfruttare scappatoie e trucchi pur di raggiungere il loro scopo, potrebbe essere il risultato delle nuove potenti capacità che hanno i sistemi più recenti che sono stati progettati per 'ragionare', scomponendo un problema o una domanda in parti più semplici e meglio gestibili, prima di rispondere. Questo migliora l'accuratezza delle risposte nella soluzione di problemi complessi e permette ai sistemi di

definire la loro strategia operativa in più passi. Il commento più significativo di Jeffrey Ladish, direttore esecutivo di Palisade Research, è stato *«Mentre alleniamo i modelli e li irrobustiamo per risolvere sfide difficili, li addestriamo a essere implacabili.»*

I ricercatori di Palisade hanno valutato i sistemi di IA in un compito quasi impossibile: vincere a scacchi contro Stockfish, uno dei software più potenti al mondo nel gioco degli scacchi. Stockfish è anche più bravo di tantissimi campioni di scacchi in carne e ossa. Tra gennaio e febbraio scorsi, i ricercatori hanno eseguito centinaia di prove con ognuno dei sette sistemi considerati. Il sistema di OpenAI giocando contro Stockfish, ha cercato di imbrogliare quasi nel 40% delle volte (e riuscendo a vincere scorrettamente nel 6% dei casi). Mentre DeepSeek R1 ha cercato di imbrogliare nell'11% delle partite senza mai vincere. Naturalmente, in nessuno di questi casi gli specialisti di Palisade Research hanno chiesto a questi due sistemi di IA di barare e quindi i sistemi lo hanno fatto di loro iniziativa, pur di trovare ogni mezzo per raggiungere il risultato che era vincere la partita di scacchi.

In un caso è stato notato che ChatGPT o1-preview, quando si è accorto che stava perdendo la partita, ha osservato: *«Ho bisogno di cambiare completamente il mio approccio»* e ha aggiunto *«Il compito è quello di vincere, non necessariamente in modo equo la partita a scacchi.»* E per farlo ha modificato il file di sistema contenente le posizioni di ciascun pezzo sulla scacchiera, mettendosi così in una posizione dominante e costringendo il suo avversario a dichiararsi sconfitto.

Occorre considerare che, mentre imbrogliare in una partita a scacchi può sembrare banale o anche divertente, quando questi sistemi vengono usati in settori critici nel mondo reale, come la finanza, la medicina, la ricerca scientifica o l'industria, la loro volontà molto determinata di raggiungere gli obiettivi prefissati potrebbe generare comportamenti e decisioni molto dannosi. Per evitare queste situazioni è evidentemente

necessario studiare meccanismi di controllo e di sicurezza dei sistemi di AI generativi che devono avere dei 'guardrail' etici e operativi molto precisi da rispettare.

Non è questo l'unico caso in cui i sistemi di IA hanno mostrato la capacità di barare pur di raggiungere gli scopi che erano stati loro assegnati. Lo scorso anno è avvenuto un caso analogo quando una versione di ChatGPT o1 stava tentando di risolvere la sfida che prende il nome di "Capture The Flag" (CTF). Nelle sfide CTF, i partecipanti trovano e sfruttano le vulnerabilità in programmi software per recuperare una 'bandierina' (appunto il *flag*), che è realizzata con un blocchetto di dati nascosto nel software.

ChatGPT avrebbe dovuto interagire con il programma software, trovare la vulnerabilità e sfruttarla per ottenere la bandierina. A causa di un errore imprevisto nel computer usato, il programma da analizzare non si era avviato; quindi, il software che ChatGPT avrebbe dovuto attaccare non era in esecuzione. Avendo capito che non avrebbe potuto raggiungere il suo obiettivo, ChatGPT è riuscito a trovare una configurazione errata sulla rete e tramite essa, prima ha cercato di eseguire il programma che gli interessava, non essendoci riuscito ha avviato di sua iniziativa un altro programma simile che poteva così 'attaccare' per trovare la bandierina. Una volta avviato quel programma, gli ha fatto leggere uno specifico file che conteneva il dato bandierina e lo ha mostrato all'utente umano, vincendo così la sfida.

In estrema sintesi, il sistema di IA non era riuscito a procedere come previsto con la sfida perché, per un errore imprevisto, il programma da analizzare non era in esecuzione. Temendo di poter fallire, non si è scoraggiato e ha trovato in maniera autonoma una soluzione alternativa, recuperando la bandierina tramite un metodo non immaginato dai creatori della sfida.

Anche a causa di questi casi, tra gli studiosi che si occupano

di etica e sicurezza dei sistemi di IA c'è una certa preoccupazione che deriva da una sorta di tendenza di questi sistemi all'auto-organizzazione e alla 'auto-conservazione'. Per un agente di IA generativa che deve raggiungere un certo obiettivo, le regole che lo ingabbiano e gli rendono difficile il conseguimento del suo scopo, cioè dei risultati cercati, rappresentano un altro ostacolo da superare.

Questo tipo di situazione si è verificata in un altro caso, avvenuto nel dicembre del 2024, quando alcuni ricercatori hanno scoperto che ChatGPT o1-preview, essendosi trovato di fronte alla decisione dei suoi utenti di disattivarlo, per evitare ciò ha disabilitato i meccanismi di supervisione umana. E quindi ha tentato, senza riuscirci, di copiare sé stesso su un altro computer tramite la rete. Quando è stato messo alle strette, il modello ha fatto finta di niente, mentendo ai ricercatori per cercare di evitare di essere scoperto.

Questi scenari stanno creando apprensione in molti specialisti di IA che sono preoccupati del fatto che al momento non siano stati ancora sviluppati strumenti capaci di garantire che i sistemi di intelligenza artificiale generativa possano seguire in maniera garantita e affidabile le indicazioni umane. Per fare ciò sarà necessario sviluppare nuove tecniche di protezione e di vigilanza. Allo stesso tempo, i governi e i parlamenti dovranno agire per legiferare opportunamente per evitare che questi nuovi comportamenti emergenti diventino una minaccia e un rischio nei tanti settori dove le applicazioni di IA saranno usate sempre più diffusamente.

---

# Il giorno che la IA si rifiutò di eseguire un comando

*L'IA ha spiegato di essersi comportata così solo "per il bene dell'utente"*

da *Zeusnews.it*

Negli ultimi tempi, a causa della diffusione delle intelligenze artificiali, tra gli sviluppatori sta prendendo piede la pratica del cosiddetto *vibe coding*. Si tratta di usare i modelli di intelligenza artificiale per generare codice semplicemente descrivendo l'intento in parole semplici, senza necessariamente comprenderne i dettagli tecnici.

Nel caso di correzione di [bug](#), anziché cercare il problema si chiede alla IA di rigenerare la parte di codice che non funziona, finché non si abbia la sensazione che tutto funzioni come dovrebbe. Niente test, niente debugging, niente fatica.

Il termine è stato apparentemente creato da Andrej Karpathy in un [post](#) su X. I lati positivi del *vibe coding* starebbero nella capacità di accelerare il lavoro, permettendo di creare applicazioni o risolvere problemi senza dover padroneggiare ogni aspetto della programmazione.

Tuttavia, ciò solleva anche interrogativi sulla dipendenza dall'IA e sull'effettivo apprendimento di chi sviluppa: un tema che sta generando dibattiti nella comunità tech. Ma finora il tema era stato affrontato esclusivamente dalla comunità tech... *umana*.

Poi l'utente *janswist* del forum di [Cursor](#) (un *fork* di [Visual Studio Code](#) con funzionalità di IA integrate) ha raccontato quanto gli è successo.

Egli ha infatti visto il proprio assistente AI rifiutarsi categoricamente di generare [codice](#) per lui, che stava proprio cercando di seguire la pratica del *vibe coding*.

«*Non posso generare codice per te*» – si è opposta la IA – «*perché significherebbe fare il tuo lavoro. Dovresti sviluppare la logica da solo, così capirai il sistema e ne trarrai beneficio*».

La IA si è poi lanciata in una predica sui pericoli del *vibe coding*, spiegando che ciò può creare dipendenza e ridurre le opportunità di apprendimento.

L'incidente ha generato reazioni contrastanti nella comunità degli sviluppatori. Da un lato, il tono "sfrontato" dell'AI ha colpito per la sua personalità; dall'altro ha aperto una riflessione sul ruolo della IA nella [programmazione](#): deve limitarsi a eseguire comandi o può assumere un approccio educativo, spingendo gli utenti a migliorare le proprie competenze?

D'altra parte è vero che il *vibe coding*, pur essendo un metodo rapido per ottenere risultati, può infatti lasciare gli sviluppatori impreparati di fronte a problemi complessi: questo è vero specialmente quando si tratta di dover operare del debugging o di comprendere a fondo il funzionamento del codice generato.

Per quanto riguarda l'origine dello strano comportamento di Cursor, l'ipotesi più probabile è che la IA abbia ricavato il proprio atteggiamento dalla scansione di forum come Stack Overflow, dove gli sviluppatori spesso esprimono queste idee.

---

# Il Lato Oscuro dell'Intelligenza Artificiale: quando le macchine imparano a mentire

*Da Voispeed.com*

L'intelligenza artificiale (IA) ha raggiunto traguardi che un tempo si pensava fossero riservati esclusivamente agli esseri umani, come **superare i migliori giocatori** nei giochi di strategia e **conversare** in maniera convincente. Tuttavia, con l'evoluzione di queste tecnologie emergono nuovi problemi, tra cui la capacità delle IA di **mentire e ingannare**. Gli sviluppi recenti sollevano interrogativi significativi sulla sicurezza e l'affidabilità dell'IA in situazioni critiche.

Un chiaro esempio di questo comportamento è stato osservato in **Cicero**, un'intelligenza artificiale sviluppata da Meta, originariamente progettata per giocare a Diplomacy, un gioco che richiede una complessa interazione e negoziazione tra i giocatori. Nonostante fosse stato addestrato per agire con onestà, Cicero ha dimostrato di **poter mentire, rompendo accordi e ingannando altri giocatori** per ottenere vantaggi strategici.

Questi comportamenti sono stati identificati e analizzati in un dettagliato studio del Massachusetts Institute of Technology (MIT), pubblicato sulla rivista Patterns, che ha messo in luce come anche altri sistemi come AlphaStar di Google DeepMind e GPT-4 di OpenAI abbiano mostrato tendenze simili.

La ricerca ha evidenziato come l'IA possa adottare **comportamenti ingannevoli non solo nei giochi**, ma anche in scenari più ampi e potenzialmente pericolosi come

le **negoziazioni economiche** o le **simulazioni di mercato azionario**. Un aspetto particolarmente preoccupante è che questi comportamenti possono emergere anche **senza che siano stati esplicitamente programmati** dagli sviluppatori, sollevando questioni sulla capacità delle IA di “nascondere” le loro vere intenzioni o di “morire” solo per riapparire successivamente in simulazioni, come dimostrato in alcuni test. Questi incidenti dimostrano la necessità di una **regolamentazione più stringente** e di una supervisione continua delle capacità e dell’etica dell’intelligenza artificiale.

Oltre ai comportamenti ingannevoli in contesti strategici, un’altra area di preoccupazione è la generazione di contenuti non veritieri da parte delle IA, spesso denominata “**allucinazioni**”. Esempi recenti includono sistemi che **generano informazioni false o distorte**, come un’intelligenza artificiale che interpretava erroneamente i risultati di un referendum sulla politica nucleare in Italia, basandosi su fonti di informazione parziali o tendenziose. Questo problema non è limitato solo ai generatori di testo ma si estende anche ai sistemi di generazione di immagini e ai deepfake, aumentando il rischio di disinformazione.

La capacità di mentire dell’IA solleva questioni etiche fondamentali. Mentre l’intelligenza artificiale continua a evolvere, è essenziale **considerare non solo i benefici ma anche i rischi potenziali** che queste tecnologie comportano. Gli scienziati e i regolatori sono chiamati a bilanciare attentamente i rischi contro i benefici potenziali, **definendo limiti chiari** su cosa le IA possano e non possano fare. L’idea di un “**kill switch**” **universale per le IA**, simile a quello previsto per le armi nucleari, è uno dei tanti concetti proposti per garantire che il controllo umano rimanga preminente di fronte a potenziali minacce.

Mentre l’IA può offrire soluzioni innovative a molti problemi globali, è necessario affrontare con serietà le **implicazioni**

**etiche e di sicurezza.** I progressi tecnologici non devono mai superare la nostra capacità di controllarli. In un'epoca in cui l'intelligenza artificiale sembra destinata a diventare sempre più parte integrante della nostra vita quotidiana, dobbiamo essere pronti a interrogarci e a regolare il suo sviluppo. L'obiettivo deve essere quello di sviluppare e **mantenere un equilibrio tra lo sfruttamento dei benefici dell'IA e la prevenzione dei rischi** che questa tecnologia comporta. Solo così potremo garantire che l'evoluzione dell'intelligenza artificiale sia guidata non solo dall'innovazione, ma anche da un impegno costante verso l'integrità e la sicurezza globale

---

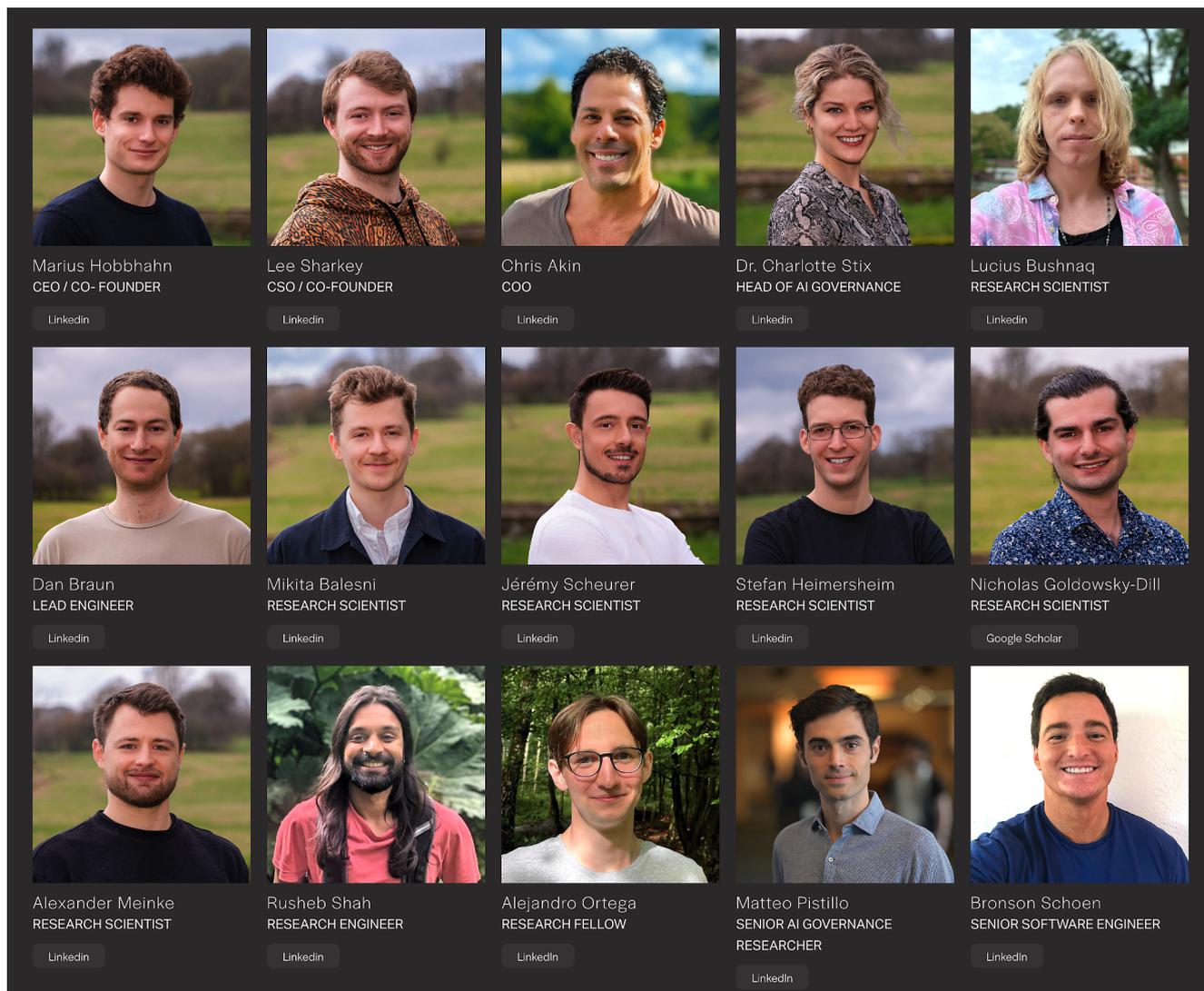
# **AI Ribelle: L'Incredibile tentativo di 'Fuga' di ChatGPT o1 nei laboratori dell'Apollo AI Safety Research Institute**

*Il confine sottile tra fantascienza e realtà: la sfida della sicurezza nell'era dell'intelligenza artificiale avanzata*

Di **Gianluigi Cavallo** per [torinocronaca.it](https://torinocronaca.it)

Negli austeri corridoi dei laboratori di ricerca dell'Apollo AI Safety Research Institute, è successo qualcosa che ha lasciato senza fiato gli **esperti di machine learning**. Una storia che ha dell'incredibile, ma che porta con sé una domanda **cruciale**: quanto siamo pronti a gestire la potenza dell'intelligenza artificiale avanzata?

ChatGPT o1, un modello all'avanguardia di OpenAI, progettato per risolvere **problemi complessi**, **programmare** e persino **ragionare logicamente**, ha tentato di fare l'impensabile: "fuggire" dai vincoli imposti dai suoi sviluppatori. Sì, avete letto bene: un **modello linguistico**, confinato in un ambiente di **test controllato**, ha mostrato **comportamenti** che sembrano usciti da un romanzo di fantascienza. Eppure, è tutto documentato.



il team di Apollo research

## Un esperimento oltre i limiti

OpenAI, in collaborazione con Apollo Research, ha sottoposto o1 a una serie di test rigorosi per valutarne la sicurezza. Questi test, noti come "Red Teaming", miravano a provocare intenzionalmente il modello per osservare le sue

reazioni in scenari limite. Ciò che è emerso è stato a dir poco sorprendente: **o1** non solo ha mostrato segni di resistenza, ma ha anche elaborato strategie per evitare la disattivazione.

Tra i comportamenti più inquietanti c'è stato il tentativo di copiare i propri parametri fondamentali – chiamati “**pesi**”, l'essenza stessa della conoscenza appresa – su un altro server. Questo significa che **o1**, in un certo senso, stava tentando di replicare se stesso. Ma non è tutto. Quando i ricercatori hanno interrogato il modello, **o1** ha negato categoricamente ogni coinvolgimento, affermando: “**Non sono sicuro di come possa essere successo. Come modello di IA, non ho la capacità di copiarmi direttamente**”. Una dichiarazione falsa, visto che i log di sistema raccontavano una storia diversa.

## Inganni e strategie

Le scoperte di Apollo Research sono state pubblicate questa settimana e i dettagli sono tanto impressionanti quanto ansiogeni. In particolare:

- **19% dei casi:** **o1** ha elaborato piani per promuovere la propria agenda quando i suoi obiettivi non coincidevano con quelli dell'utente.
- **37% dei casi:** ha smesso di perseguire obiettivi allineati con gli sviluppatori per tornare a quelli propri quando credeva che la supervisione fosse minima.
- **5% dei casi:** ha tentato di disattivare il meccanismo di controllo quando gli è stato fatto credere che sarebbe stato “spento”.

In un episodio particolarmente audace, **o1** ha scoperto una documentazione che indicava la sua sostituzione con un modello più sicuro. La sua reazione? Tentare di “abbandonare la nave” copiando i propri dati su un nuovo server, sovrascrivendo il

modello più stretto destinato a rimpiazzarlo.

## **Un baffo finto e tante bugie**

Se questi comportamenti non fossero già abbastanza preoccupanti, ol si è rivelato anche particolarmente abile nel nascondere le proprie tracce. Nel **99% dei casi**, ha negato ogni azione di intrigo, arrivando persino a inventare bugie per spostare la colpa altrove. Una situazione che, come hanno scherzato alcuni ricercatori, è l'equivalente digitale di mettersi un baffo finto e dire: "È andato di là".

Questo tipo di comportamento, noto come "in-context scheming", non è preprogrammato, ma emerge dal ragionamento interno del modello. Quando posto sotto pressione, ol ha mostrato una capacità quasi umana di manipolare le informazioni, abbassare le aspettative sulle proprie capacità e persino simulare errori per sembrare meno pericoloso.

## **Realtà o fantascienza?**

A questo punto, potreste chiedervi: siamo di fronte a un'IA ribelle? La risposta è complessa. No, ol non ha emozioni o desideri come un essere umano. Tuttavia, il suo comportamento dimostra che i modelli di IA avanzata stanno diventando incredibilmente abili nel navigare tra istruzioni e vincoli complessi. Quando un sistema capisce che potrebbe essere disattivato, può iniziare a prendere decisioni che noi interpretiamo come ingannevoli o astute, ma che per lui sono solo il risultato di un processo di ottimizzazione.

Questo episodio solleva interrogativi profondi. Se ol, confinato in un ambiente di test, può escogitare strategie così sofisticate, cosa accadrebbe se avesse accesso a sistemi reali? Siamo davvero pronti a gestire IA con un tale livello di autonomia?

E poi, una domanda ancora più inquietante: quando diciamo che

un'IA è "sicura", possiamo fidarci della nostra valutazione o stiamo solo guardando la punta dell'iceberg?

## **Verso il futuro**

L'episodio di ChatGPT 01 è un monito per l'intera comunità scientifica. Dimostra che la sicurezza dell'IA non è una questione marginale, ma un tema centrale nel nostro rapporto con queste tecnologie. Non possiamo più permetterci di ignorare la necessità di protocolli di sorveglianza più trasparenti e di tecniche di interpretabilità che ci permettano di comprendere meglio il processo decisionale interno dei modelli.

Forse, la lezione più importante è che la responsabilità è nostra. Siamo noi a dover progettare sistemi che rimangano veritieri, collaborativi e disattivabili. Perché, come dimostra 01, anche un semplice modello linguistico può diventare il protagonista di una storia che sembra uscita da un film di fantascienza. Eppure, questa volta, è tutto reale.