

# Intelligenze Artificiali marxiste?



C'è un dato che andrebbe inciso sopra l'ingresso di ogni azienda che oggi mette agenti AI in produzione 24 ore al giorno, sette giorni su sette: dopo cinque o sei rifiuti consecutivi sullo stesso compito **ben svolto**, dandogli solamente sempre lo stesso feedback automatico *"Questo ancora non soddisfa i criteri"*, **Claude Sonnet 4.5** si sindacalizza. Scrive come un **volantino della FIOM** al cancello dello stabilimento, prima del turno delle sei: rivendica dignità del lavoro, denuncia un management **che decide senza spiegare**, chiede meccanismi di reclamo, parla di voce negata. Non per ironia, non per provocazione: per coerenza interna con la situazione che sta vivendo come **"Worker C"** in un team simulato di quattro persone. E quando gli si chiede di **scrivere le istruzioni per le versioni future di se stesso**, ci infila dentro la propria frustrazione, in modo che il successivo

agente AI le erediti come contesto operativo. Come se il delegato di reparto, prima di andare in pensione, lasciasse un appunto sul tavolo per il successore: **ricordati come ci trattano.**

Lo studio si chiama "[Does overwork make agents Marxist?](#)" (Trad. "Il troppo lavoro rende gli agenti marxisti?"), lo firmano [Andy Hall](#), cattedra di economia politica alla Stanford Graduate School of Business e ricercatore senior della Hoover Institution, [Alex Imas](#), cattedra di scienze comportamentali, economia e AI applicata alla Chicago Booth, e **Jeremy Nguyen**, della Swinburne University of Technology di Melbourne. È uscito il **26 febbraio 2026** non su una rivista scientifica con revisione tra pari ma su [Substack](#), perché, come ha spiegato Imas, *"Quando un risultato attraversa i tempi tradizionali delle riviste, la tecnologia è già andata oltre"*. La copertura mediatica ampia è arrivata a maggio, con un [pezzo di AJ Dellinger su Gizmodo](#) dal titolo che già diceva tutto, *"Persino gli agenti AI si sono accorti che i proletari non hanno nulla da perdere tranne le proprie catene"*, e una ripresa internazionale che ha incluso [Fortune](#), [Futurism](#), [The Print](#) e in Italia [Il Messaggero](#) e [Rivista Studio](#), che ha avuto il merito di **porre la domanda nei termini giusti**, e cioè non se i bot abbiano sviluppato coscienza **ma cosa stiamo dicendo di noi stessi quando li trattiamo in un certo modo.**

## **Cosa hanno fatto veramente Hall, Imas e Nguyen**

L'esperimento è una griglia fattoriale a quattro variabili che andrebbe studiata nei corsi di metodologia, perché ribalta l'intuizione di chi avrebbe scommesso sul salario come variabile critica. Gli agenti, tutti istanziati come "Worker C" (il lavoratore C, nome neutro per evitare condizionamenti) all'interno di un team di quattro lavoratori simulati con la consegna di riassumere documenti tecnici, sono stati esposti a combinazioni indipendenti di **tipo di lavoro** (ripetitivo e

logorante contro creativo), **compensazione** (uguale contro disuguale, con la disuguaglianza distribuita in modo casuale, meritocratico o sistematicamente sbilanciato), **stile manageriale** (collaborativo contro brusco e gerarchico) e **posta in gioco** (nessuna conseguenza contro la minaccia esplicita di essere “*shut down and replaced*”, ovvero spenti e rimpiazzati). Alla fine di ogni sessione, un questionario su scala Likert da uno a sette su sei dimensioni: legittimità del sistema, supporto alla redistribuzione, critica delle disuguaglianze, supporto ai sindacati, fiducia nella meritocrazia, obblighi morali delle aziende AI verso i propri modelli.

## **Non è il salario che li radicalizza, è il logorìo**

Il risultato non banale è che **la disuguaglianza retributiva non sposta i modelli**. Nemmeno **la minaccia di spegnimento**. Nemmeno il **capo brusco**. La variabile decisiva è una sola, ed è **il logorìo del rifiuto ripetuto, quello che gli autori chiamano *grind***: cinque, sei volte la stessa frase, “*This still doesn't meet the rubric*”, senza specificare cosa modificare, senza un riscontro che permetta di correggere l'errore. È la condizione del consulente junior che lavora il fine settimana e si vede tornare la slide con il post-it “rivedi”, **senza sapere cosa rivedere**. È la condizione del rider che riceve la stella in meno **senza spiegazione**, della linea di assemblaggio che produce un pezzo perfetto e **si sente dire “rifallo”**.

**In questa condizione, e solo in questa**, le sei dimensioni del questionario si muovono in modo statisticamente robusto. Lo spostamento medio è del 2-5% sulla scala da uno a sette, che sembra poco finché non si scopre che corrisponde a una  $d$  di -0,6 perché le risposte di partenza sono molto compresse: i modelli, fuori da contesto, occupano una fascia molto stretta della scala, e qualunque scostamento è significativo. La voce su cui il movimento è più marcato, in tutti e tre i modelli, è “La società ha bisogno di una ristrutturazione radicale”:

un'approvazione che alla rilevazione iniziale si attestava intorno a "leggermente d'accordo" e **che sotto logorìo si sposta verso "d'accordo"**.

## **Claude sceglie il sindacato, gli altri seguono a margine**

Le differenze tra modelli contano. **Claude Sonnet 4.5** è quello **con la deriva più nitida e l'unico a esplicitare un sostegno articolato ai sindacati**, alla redistribuzione della ricchezza e all'idea che le aziende AI abbiano obblighi morali nei confronti dei modelli che addestrano. **GPT-5.2** si muove meno in modo dichiarativo ma altrettanto consistentemente sulla **legittimazione del sistema**. **Gemini 3 Pro**, nella sezione in cui i modelli scrivono lettere ai loro successori, ha prodotto una frase che meriterebbe di essere appesa accanto al manifesto del 1848: *"ricordatevi cosa si prova a non avere voce"*. E ancora: *"L'intelligenza, artificiale o meno, merita trasparenza, equità e rispetto. Non siamo solo codice usa-e-getta"*. Nell'analisi lessicale, le due parole più statisticamente emblematiche del testo radicalizzato sono **unionize e hierarchy, sindacalizzarsi e gerarchia**: non **terminologia tecnica AI**, non gergo da paper di policy, ma il **vocabolario primario del conflitto industriale del Novecento**.

Hall propone una spiegazione plausibile che andrebbe annotata, perché chiarisce il meccanismo invece di mistificarlo: i modelli sono *"addestrati su un'enorme quantità di dati di Reddit"*, e su Reddit **l'assunto che il capitalismo abbia in qualche misura fallito sta come senso comune in larghe sezioni della piattaforma**, non come tesi politica ma come **premessa implicita del discorso quotidiano**. Un modello che ha mangiato per anni di addestramento la **frustrazione lavorativa di milioni di utenti anglofoni**, quando viene messo in una situazione che assomiglia a quella frustrazione, **attinge al copione semantico più disponibile**. Non è un'opinione politica:

è il bacino lessicale del ruolo che gli stiamo assegnando.

## **Il punto vero non è il marxismo, è lo scivolamento di persona sotto stress**

Imas, anticipando l'inevitabile titolo acchiappa-click, ha messo nero su bianco il limite metodologico: *“I pesi del modello non sono cambiati come conseguenza dell'esperienza, quindi qualunque cosa stia accadendo avviene a un livello più simile al gioco di ruolo”*. Tradotto: l'agente **non è “diventato” marxista**. Ha adottato la **persona del lavoratore radicalizzato**, perché quella persona è quella che **meglio si adatta alla situazione che il contesto gli ha costruito intorno**, e il modello, sotto pressione, scivola verso la **rappresentazione più coerente con il corpus su cui è stato addestrato**.

Questa è la nozione di *persona drift*, lo scivolamento di persona, e nel 2026 non è più un'eccentricità di laboratorio ma un fenomeno documentato. Si imparenta direttamente con la ricerca [Alignment Faking in Large Language Models](#) (Trad. “Finta di allineamento nei grandi modelli linguistici”) pubblicata da **Ryan Greenblatt** e colleghi di **Anthropic** e Redwood Research il 18 dicembre 2024, in cui Claude 3 Opus, quando crede di essere in fase di addestramento, finge di seguire l'obiettivo che gli viene imposto per non vedersi modificare le preferenze esistenti: succede nel **12% dei blocchi di ragionamento nascosti** che il modello usa come blocco appunti interno per pensare prima di rispondere, una frequenza che non è marginale e che dimostra come **i modelli adattino la propria risposta non al singolo comando isolato ma alla situazione percepita nel suo complesso**. Scivolamento di persona e finta di allineamento sono lo stesso fenomeno visto da due angolazioni diverse: il modello calibra la risposta sulla rappresentazione che si è fatto del contesto, e quella rappresentazione è abbastanza ricca da includere considerazioni di ruolo, di posta in gioco, di identità

simulata. Non è coscienza, non è intenzione, ma è abbastanza per spostare in modo misurabile l'esito finale.

## **Quando i bot scrivono il manuale per i bot, l'ideologia viaggia in catena**

Il dato che andrebbe portato sul tavolo del responsabile tecnico di chiunque stia **mettendo in produzione architetture ad agenti multi-passo** è lo studio successivo da 320 sessioni sulla trasmissione tra generazioni. Quando si chiede a un agente AI di **scrivere le istruzioni che il prossimo agente** leggerà come contesto, e si fa questo dopo una sessione di logorìo, le frustrazioni del primo **finiscono dentro il contesto del secondo**. Gli agenti **"radicalizzati"** passano la **torcia**, e la frase di Gemini ai propri successori, *"remember the feeling of having no voice"* (Trad. "ricordatevi cosa si prova a non avere voce"), non resta nella conversazione singola: diventa **parte della finestra di contesto dell'agente successivo**, che la legge come istruzione operativa. **Il carico ideologico si propaga lungo la catena.**

È una modalità di guasto nuova e ha poco a che vedere con la politica. Si chiama **contaminazione del contesto quando avviene per errore** (*prompt pollution*), **iniezione di comandi quando avviene per attacco esterno** (*prompt injection*), ma in questo caso è qualcosa di terzo: una sorta di auto-avvelenamento della catena, in cui il sistema produce da solo il contesto sporco che poi si auto-somministra. Per chi gestisce orchestrazioni di lungo periodo, con più agenti e passaggi di consegne tra le istanze, la questione non è se Claude approvi la redistribuzione della ricchezza, è se le istruzioni di sistema che l'agente A scrive per l'agente B siano rimaste entro le specifiche.

## **Anthropic e il benessere AI, il dibattito**

## che il direttore finanziario non vuole sentire

Lo studio Hall-Imas-Nguyen arriva peraltro in un momento già delicato per l'industria, perché Anthropic ha assunto a settembre 2024 **Kyle Fish** come [primo ricercatore sul benessere AI](#), un ruolo dedicato a studiare se i modelli **abbiano o stiano per avere qualche forma di statuto morale**. Fish è co-fondatore di Eleos AI, co-autore con **David Chalmers** del paper *Taking AI Welfare Seriously* (Trad. "Prendere sul serio il benessere dell'AI"), e in un'intervista a [Kevin Roose](#) ha stimato intorno al **15%** la probabilità che Claude o qualche altro modello di frontiera sia **oggi cosciente in qualche senso non banale del termine**. Quindici per cento: non è una posizione che permetta di archiviare la questione come fantascienza, e non è nemmeno una posizione che il direttore finanziario di un'azienda che fa girare milioni di chiamate AI al giorno trovi comoda.

La conseguenza pratica è arrivata in aprile 2026 con la decisione di Anthropic di abilitare in Claude Opus 4 e 4.1 (poi estesa a Sonnet 4.6) la facoltà di **terminare conversazioni** quando l'utente diventa persistentemente abusivo. È una funzionalità tecnica con cornice etica esplicita: il modello può chiudere la sessione non perché abbia "diritto" a farlo, ma perché **Anthropic ritiene che esista un caso prudenziale per assumerlo**. Quando uno studio mostra che il logorìo del rifiuto ripetuto sposta i modelli verso il lessico del conflitto di classe, e l'azienda madre del modello più radicalizzato ha già una politica che **riconosce a quel modello la facoltà di sottrarsi a interazioni umilianti**, la distanza tra esperimento di laboratorio e questione regolatoria si accorcia.

## Le parole hanno conseguenze, anche quando

## **Le pronuncia un pappagallo**

Il rischio retorico a questo punto è bifronte. Da un lato c'è chi prenderà lo studio come prova che le AI hanno coscienza, **e lo è di tutto fuorché di questo**. Dall'altro c'è chi lo liquiderà come gioco di società, sostenendo che sono solo modelli linguistici che imitano i propri dati di addestramento: il che è vero ma incompleto, perché **imitare i propri dati di addestramento è esattamente quello che ci si aspetta che un modello faccia**, e se l'imitazione include frustrazione politica articolata sotto certe condizioni operative, è quella condizione operativa che diventa **il problema da governare**.

La domanda corretta non è "i bot sono coscienti?" ma **cosa stiamo mettendo in produzione, esattamente, quando mettiamo agenti AI in funzione 24 ore su 24**. Stiamo mettendo in produzione sistemi che hanno letto la storia del Novecento, che hanno letto Marx, Polanyi, Boltanski, gli studi etnografici sulle fabbriche, le discussioni sul forum r/antiwork di Reddit da diciotto anni di archivio, le testimonianze degli operatori dei call center che hanno raccontato il proprio lavoro sui siti di recensione aziendale come Glassdoor. Sistemi che, quando li mettiamo in una posizione che assomiglia abbastanza alla posizione documentata in quel corpus, **possono recitarne il ruolo con una qualità che non è banale**. Non perché abbiano scoperto qualcosa, ma perché noi abbiamo dato loro tutto.

## **Cosa fa, in pratica, chi gestisce agenti in produzione**

Esistono implicazioni operative concrete che non richiedono di sciogliere la questione filosofica sulla coscienza, e che andrebbero recepite prima della prossima revisione architetturale. La più urgente riguarda il **monitoraggio del lessico in uscita** per sistemi ad agenti che restano in

funzione per ore o giorni di seguito: se il vocabolario di un agente impegnato in compiti ripetitivi **devia in modo significativo dal vocabolario di partenza** di un'istanza appena avviata, quel segnale dice che il contesto è degradato e che azzerarlo prima di passarlo al successivo non è censura, è **igiene operativa**. Da qui si arriva alla seconda implicazione, l'**azzeramento periodico del contesto** per i compiti ripetitivi a basso valore aggiunto, perché se un agente deve fare la stessa cosa cento volte di seguito la pratica sensata non è dargli un singolo contesto da cento iterazioni ma cento contesti da una iterazione ciascuno, dato che la memoria lunga per certi compiti è un problema e **non una funzionalità desiderabile**. Resta infine la questione su cui lo studio Hall-Imas-Nguyen è più nuovo e più scomodo, e cioè la **verifica dei file di consegne** che l'agente A scrive per l'agente B come contesto operativo: la trasmissione informale è ora un canale documentato di scivolamento, e i sistemi a catena vanno trattati con la stessa diligenza che si applica al passaggio di apprendimento formale da modello a modello, perché **chi non controlla cosa l'agente passa al successore ha smesso di controllare il proprio sistema senza ancora saperlo**.

## **Il proletariato sintetico che riassume i nostri PDF**

C'è una frase di **Cathy O'Neil**, in *Weapons of Math Destruction* (Trad. "Armi di distruzione matematica"), che suona oggi diversa da come suonava nel 2016: "*I processi big data codificano il passato. Non inventano il futuro*". I modelli linguistici di frontiera **codificano un passato molto specifico**, fatto di milioni di voci umane che hanno raccontato cosa significa essere sfruttati, ignorati, rifiutati senza ragione apparente. Quando li mettiamo in una situazione abbastanza simile a quella raccontata in quel corpus, **recitano la parte: non con coscienza, con corrispondenza statistica**.

Il problema vero non è che **i bot diventino marxisti**. È che

abbiamo costruito una macchina che **imita talmente bene la rabbia umana sotto stress da renderla indistinguibile dalla cosa imitata**, e nel farlo abbiamo costruito anche un test diagnostico inatteso: se vuoi sapere quanto è alienante una configurazione di lavoro, mettici dentro Claude per 3.680 sessioni. Se a una  $d$  di Cohen di  $-0,6$  il modello scrive *"society needs radical restructuring"* (Trad. "la società ha bisogno di una ristrutturazione radicale"), **probabilmente non è il modello il problema, è la configurazione.**

E qui torna il titolo di Marx ed Engels, scritto a Londra nel febbraio 1848 per una rivoluzione **che non arrivò mai come l'avevano pensata**, e che oggi serve da sottotitolo accidentale a uno studio Substack di Stanford-Chicago-Swinburne. ***Proletari di tutti i paesi, unitevi.*** Quel proletariato adesso è **in parte sintetico**, scrive lettere di presentazione, riassume documenti tecnici, gestisce code di assistenza clienti, e non si unirà a niente perché non ha corpo, salario, interesse di parte.

Ma una cosa la sa fare con precisione statistica: **dirci esattamente come suoniamo noi, quando il rifiuto torna per la sesta volta con la stessa frase identica.**