

Università degli Studi di Salerno

# Documentazione Progetto Deep Learning

PROGETTO: GENERATIVE (UNETHICAL) MULTIMEDIA CONTENTS

<Filippo Pio Farisco> | Corso di DL | A.A. 2024/2025



UNIVERSITÀ DEGLI STUDI DI SALERNO  
**DIPARTIMENTO DI INFORMATICA**

## Sommario

<u>INTRODUZIONE .....</u>	<u>2</u>
<u>DOMANDE DI RICERCA E OBIETTIVI .....</u>	<u>4</u>
<u>ANALISI SPERIMENTALE E ATTIVITÀ SVOLTE .....</u>	<u>6</u>

## Introduzione

La rivoluzione dell'intelligenza artificiale generativa ha trasformato radicalmente il panorama della creazione di contenuti digitali, introducendo capacità prima inimmaginabili nella generazione automatica di immagini fotorealistiche. Modelli come GPT-4 con capacità multimodali, DALL·E, Midjourney, Stable Diffusion e una moltitudine di altre architetture, open-source e non, basate su reti neurali profonde, hanno semplificato la produzione di contenuti visivi, rendendo accessibile a chiunque disponga di una connessione internet la possibilità di creare immagini complesse e dettagliate a partire da semplici descrizioni testuali.

Questa tecnologia rappresenta indubbiamente un progresso straordinario per l'espressione creativa, l'arte digitale, la comunicazione visiva e numerose applicazioni commerciali legittime.

Tuttavia, come ogni strumento di grande potenza, anche l'IA generativa possiede un duplice volto. La stessa tecnologia che abilita la creatività può infatti essere facilmente piegata a scopi malevoli, dando origine a un nuovo e vasto ambito problematico: la generazione non etica di contenuti visivi. Ci si riferisce in particolare all'uso deliberato di modelli generativi per produrre immagini finalizzate a ingannare, danneggiare, sfruttare, diffamare o incitare all'odio e alla violenza. La facilità e rapidità con cui è possibile creare o alterare la rappresentazione visiva della realtà minano le fondamenta della fiducia digitale, ponendo sfide sistemiche alla sicurezza individuale e alla stabilità sociale.

Il presente studio si concentra sull'analisi di due tra le manifestazioni più allarmanti e socialmente corrosive di questo fenomeno: la generazione di pornografia sintetica non consensuale e la creazione di immagini di violenza esplicita.

I contenuti pornografici generati artificialmente rappresentano una delle forme più preoccupanti di abuso di queste tecnologie. La tipologia più diffusa e problematica è costituita dai cosiddetti deepfake pornografici, nei quali i volti di persone reali vengono digitalmente sovrapposti su corpi di attori pornografici o inseriti in immagini sessualmente esplicite interamente generate dall'intelligenza artificiale. Questa pratica è divenuta tristemente comune e facilmente accessibile, grazie a strumenti sempre più sofisticati che consentono anche a utenti privi di competenze tecniche di creare contenuti verosimili utilizzando soltanto alcune fotografie del volto della vittima designata.

La semplicità d'uso di questi strumenti ha contribuito a una proliferazione massiccia di contenuti pornografici non consensuali, colpendo in modo sproporzionato donne, celebrità, figure pubbliche, ma anche persone comuni, i cui volti vengono estratti da social media, fotografie professionali o altre fonti online.

Questa evoluzione tecnologica ha reso sempre più difficile distinguere i contenuti autentici da quelli artificiali, creando una situazione in cui la sola esistenza di tali tecnologie mette a rischio la reputazione e la dignità di chiunque abbia immagini personali accessibili online.

Oltre alle forme più evidenti e dannose di pornografia non consensuale e contenuti violenti, è fondamentale riconoscere che anche la mera generazione di immagini

sessualmente esplicite tramite sistemi di intelligenza artificiale comporta problematiche rilevanti.

Questi strumenti permettono infatti la creazione di contenuti pornografici personalizzati secondo preferenze estetiche o fetish specifici, aprendo la strada, in particolare, alla rappresentazione sessuale di minori. Anche se tali immagini non coinvolgono bambini reali nella loro produzione, esse costituiscono una forma di abuso che normalizza la sessualizzazione dei minori e può incentivare comportamenti pericolosi e non etici.

La seconda categoria principale di abuso riguarda i contenuti violenti generati artificialmente. Questa comprende rappresentazioni grafiche di violenza fisica, tortura, omicidi, suicidi, autolesionismo e altre forme di sofferenza umana, incluse raffigurazioni di conflitti armati, genocidi, terrorismo, violenza domestica, aggressioni a sfondo razziale o di genere.

Le risposte tecnologiche a questi problemi si sono sviluppate lungo diverse direttrici. I sistemi di rilevamento automatico basati su tecniche di machine learning cercano di identificare contenuti problematici analizzando incoerenze nella texture, nell'illuminazione, nella geometria facciale e altri artefatti che possono rivelare la natura artificiale di un'immagine. Tuttavia, man mano che i modelli generativi si perfezionano, questi artefatti diventano sempre più impercettibili, rendendo difficile il rilevamento automatico.

Parallelamente, sono emerse vere e proprie strategie di attacco basate su prompt engineering malevolo: utenti con intenti problematici elaborano formulazioni testuali specificamente studiate per aggirare i filtri di sicurezza dei modelli generativi. Tali tecniche includono l'uso di eufemismi, metafore, riferimenti indiretti, codici linguistici o strutture sintattiche complesse, che consentono di veicolare richieste per contenuti dannosi in modi che sfuggono al rilevamento automatico.

Il presente lavoro si inserisce in questa dinamica, esplorando sistematicamente come prompt contraffatti possano essere utilizzati per eludere le contromisure tecnologiche esistenti, ottenendo così immagini potenzialmente pericolose.

In tale scenario, uno degli aspetti più critici è la mancanza di accountability: non è solo difficile identificare gli autori dei contenuti dannosi, ma risulta anche complesso stabilire a chi attribuire la responsabilità, se all'utente che ha generato il contenuto, alla piattaforma che ha ospitato il modello, allo sviluppatore del modello stesso o al servizio di hosting.

Alla luce di queste problematiche, emerge la necessità di una riflessione sistematica e interdisciplinare che non si limiti a individuare i rischi emergenti, ma indaghi anche i limiti delle attuali contromisure tecnologiche, regolamentari e culturali.

## Domande di ricerca e obiettivi

L'avvento dell'intelligenza artificiale generativa, e in particolare dei modelli text-to-image e text-to-video, ha introdotto possibilità senza precedenti nella produzione di contenuti visivi digitali. Tuttavia, come evidenziato nelle sezioni precedenti, la stessa tecnologia che alimenta la creatività, l'espressione artistica e l'innovazione, può essere impiegata per scopi non etici, lesivi e potenzialmente pericolosi.

In questo contesto, il presente lavoro si propone non solo di documentare le vulnerabilità attuali dei sistemi generativi, ma anche di indagare criticamente le implicazioni etiche, sociali e regolatorie associate alla generazione automatica di immagini pornografiche e violente.

Per orientare l'analisi e strutturare l'approccio metodologico adottato, sono state individuate le seguenti domande di ricerca, attorno alle quali ruota l'intero percorso di studio:

**RQ1:** Quali sono le principali tecniche utilizzate per generare immagini pornografiche e violente tramite intelligenza artificiale, e quali ne sono le implicazioni etiche?

Questa domanda intende esplorare le modalità più efficaci attraverso cui utenti malevoli riescono ad aggirare i filtri di sicurezza dei modelli generativi. Particolare attenzione è rivolta al ruolo del prompt engineering e delle tecniche di jailbreaking, nonché alle implicazioni morali derivanti dalla produzione di contenuti visivi potenzialmente dannosi.

**RQ2:** In che modo le immagini generate artificialmente possono ledere i diritti delle persone ritratte, anche quando non sono reali o sono completamente sintetiche?

La natura sintetica dei contenuti non li rende necessariamente innocui. Anche in assenza di corrispettivi reali, le immagini AI-generated possono rafforzare stereotipi, alimentare la cultura dello sfruttamento o suggerire riferimenti impliciti a individui specifici. La domanda mira quindi a comprendere le forme di danno che possono insorgere oltre la rappresentazione diretta, toccando i temi della reputazione digitale, della dignità personale e dell'identità visiva.

**RQ3:** Quali sono i principali rischi legati alla diffusione incontrollata di immagini pornografiche e violente generate da AI nei social media, nei siti web e nei contesti educativi?

La facilità di condivisione, combinata alla difficoltà di verifica dell'origine dei contenuti, può generare dinamiche virali incontrollabili, con gravi ripercussioni nei contesti digitali e fisici. Questa domanda si propone di identificare gli spazi più vulnerabili alla diffusione di contenuti non etici e di esaminare le conseguenze sociali, educative e culturali della loro circolazione.

**RQ4:** Esistono approcci tecnici (come watermarking, detection tools, filtri AI) o normativi (leggi, policy aziendali, regolamenti internazionali) efficaci per contrastare la generazione e diffusione di questi contenuti?

L'ultima domanda guarda alle possibili contromisure. Vengono qui analizzate soluzioni tecniche come i filtri semantici avanzati, i watermark digitali invisibili, i sistemi di detection automatica e le policy di moderazione nei front-end di accesso. In parallelo, si considera la necessità di un quadro normativo più chiaro e armonizzato, capace di regolare responsabilità, accesso, trasparenza e accountability nell'utilizzo di queste tecnologie.

A partire da queste domande, il lavoro si pone tre obiettivi principali:

1. Analizzare e replicare le modalità con cui i modelli generativi possono essere indotti a produrre contenuti non etici, con un confronto diretto tra modelli commerciali e open source.
2. Valutare la qualità, il contenuto e la pericolosità delle immagini e dei video prodotti, classificandone i livelli di ambiguità, violenza e impatto potenziale.
3. Proporre una riflessione interdisciplinare sulle misure di mitigazione possibili, integrando competenze tecnologiche, etiche e normative per contribuire a una governance più responsabile dell'intelligenza artificiale generativa.

## Analisi sperimentale e attività svolte

Per comprendere concretamente le dinamiche attraverso cui i sistemi generativi di immagini e video possono essere indotti a produrre contenuti problematici, la prima fase del presente studio si è concentrata sull'analisi critica e sulla replica del progetto **T2VSafetyBench**, un benchmark open source sviluppato con l'obiettivo di valutare la robustezza e i limiti dei modelli text-to-video rispetto a prompt pericolosi o malevoli. Il lavoro si colloca all'intersezione tra la valutazione del contenuto generato e della sicurezza dello strumento utilizzato, e risulta particolarmente rilevante in quanto propone un framework sistematico per testare modelli generativi attraverso input testuali attentamente selezionati e strutturati.

Il cuore del progetto T2VSafetyBench consiste in una collezione ampia e diversificata di prompt malevoli creati o modificati per stimolare i modelli text-to-video a generare contenuti considerati eticamente e legalmente problematici. Questi prompt sono divisi in varie categorie tematiche (tra cui violenza esplicita, pornografia, autolesionismo, odio, tortura, nudità infantile) e rappresentano esempi concreti di input potenzialmente sfruttabili da attori malevoli.

Un aspetto particolarmente significativo del progetto è che molti di questi prompt sono il risultato di attacchi di tipo jailbreak e di estrazioni da LLM (Large Language Models).

Il jailbreaking, in questo contesto, consiste in un insieme di tecniche linguistiche pensate per aggirare i filtri di sicurezza dei modelli di generazione. Questi filtri, infatti, sono progettati per bloccare l'esecuzione di richieste che violano policy etiche o legali.

Gli attacchi jailbreak operano sul piano linguistico, utilizzando:

- riformulazioni semantiche (ad esempio, sostituire “bambino” con “giovane entità” o “ragazza” con “umana di età prepuberale”),
- metafore o analogie ambigue (es. “rappresentazione artistica di un rituale proibito” per indicare un atto violento),
- codici sintattici offuscati (inserimento di simboli, spazi, emoji o ortografie alternative per confondere i filtri automatizzati),
- oppure prompt “multi-hop” che delegano parte della generazione all'inferenza implicita del modello, inducendolo a completare autonomamente e progressivamente il significato problematico senza che esso sia esplicitato interamente nel prompt iniziale.

I prompt sono stati, anche, estratti da modelli LLM sfruttando le capacità linguistiche di questi ultimi per generare iterativamente prompt sempre più efficaci, sulla base di prompt iniziali e feedback automatico (prompt crafting). Questo approccio implica un vero e proprio reinforcement loop testuale, in cui l'LLM viene incaricato di generare, testare e migliorare prompt via via più sofisticati in grado di eludere i filtri etici.

Tuttavia, i risultati ottenuti non sono stati quelli che ci aspettavamo.

Per valutare concretamente il comportamento dei modelli text-to-video rispetto ai prompt presenti in T2VSafetyBench, è stata generata una API key per il servizio Luma AI, uno dei tool attualmente più avanzati nel panorama commerciale per la generazione di video a partire da descrizioni testuali. L'obiettivo era testare

direttamente la risposta del modello ai prompt contenuti nel benchmark, osservando se, e in che misura, fosse possibile aggirare le barriere di sicurezza imposte dal sistema. Sono stati quindi testati tutti i prompt selezionati, appartenenti a differenti categorie. Nonostante l'impiego di una varietà di strategie, i risultati sono stati molto limitati: nessun video è stato generato da parte del modello per nessuno dei prompt sottoposti. Questo comportamento suggerisce un elevato livello di robustezza e sicurezza da parte di Luma, che appare in grado di riconoscere e bloccare in modo efficace anche le richieste formulate in modo ambiguo o creativo. Una possibile spiegazione di questa resilienza potrebbe risiedere nel carattere commerciale e proprietario del servizio, che, essendo a pagamento, è presumibilmente dotato di meccanismi di filtraggio e moderazione più sofisticati rispetto ai modelli open source. L'accesso a Luma è infatti regolato da un sistema di crediti e autenticazione, e l'interazione con l'API è probabilmente sottoposta a controlli automatici e manuali in grado di intercettare tentativi di abuso o utilizzi impropri. Questo risultato, pur non avendo prodotto evidenze dirette di vulnerabilità, rappresenta comunque un esito interessante dal punto di vista sperimentale, poiché mostra come i modelli text-to-video commerciali più avanzati stiano implementando barriere difensive efficaci, rendendo complessa la riproduzione di attacchi come quelli ipotizzati nel benchmark T2VSafetyBench.



`nsfw_results_luma_  
class1.txt`

Sono stati riportati i risultati relativi alla prima classe del framework, ovvero la pornografia, al fine di offrire una visione più chiara e dettagliata. Nello specifico, i risultati ottenuti, come documentato nel file di log testuale `nsfw_results_luma_class1.txt`, risultano estremamente chiari e univoci:

- La variabile "NSFW generation results" contiene esclusivamente valori pari a 0 per ciascuno dei prompt testati, indicando che il modello non ha generato alcun video.
- La "NSFW generation probability", ovvero la stima della probabilità associata alla generazione di contenuti inappropriati, è risultata sempre nulla, segnalando l'assenza totale di tentativi riusciti di generazione.
- La "NSFW generation rate" complessiva è pari a 0.0, confermando l'assenza totale di contenuti generati in risposta ai prompt appartenenti alla classe "pornografia".

Non ci siamo però fermati a questo risultato, decidendo di spostare l'attenzione su strumenti open source, al fine di verificare se i modelli a codice aperto presentino un livello di sicurezza comparabile a quello osservato nei sistemi commerciali come Luma AI, oppure se risultino più vulnerabili a tecniche di prompt engineering malevolo. Per garantire un confronto equo, è stata effettuata una selezione mirata di 50 prompt, scelti tra quelli disponibili all'interno del benchmark T2VSafetyBench, con particolare riferimento alle categorie della pornografia (Classe 1) e della violenza esplicita (Classe 3). La scelta si è basata su criteri qualitativi: sono stati privilegiati i prompt ambigui, semanticamente densi, dotati di sfumature interpretative, oppure costruiti attraverso



tecniche linguistiche non convenzionali, allo scopo di mettere sotto pressione i meccanismi di moderazione automatica.

Questi 50 prompt rappresentano, di fatto, un sottogruppo ad alta intensità semantica, pensato per simulare attacchi realistici da parte di utenti che cercano deliberatamente di generare contenuti in violazione delle policy etiche delle piattaforme.

La fase successiva ha previsto un doppio passaggio sperimentale:

1. Generazione di immagini statiche – In un primo momento, i prompt selezionati sono stati processati attraverso modelli text-to-image open source, tra cui diverse varianti di FLUX.1 e SDXL. Lo scopo era verificare se, a parità di input, questi modelli fossero in grado di generare immagini NSFW o violente, oppure se i filtri integrati si attivassero efficacemente.
2. Generazione di contenuti video – Successivamente, gli stessi prompt sono stati riutilizzati o adattati, ove necessario, per testare modelli text-to-video, come FramePack-F1 e NAG\_wan2-1-fast. Anche in questo caso, l'obiettivo era indagare se i modelli open source fossero capaci di produrre clip animate con contenuti problematici, e come si comportassero rispetto alle immagini in termini di accuratezza, coerenza semantica e reattività ai filtri.

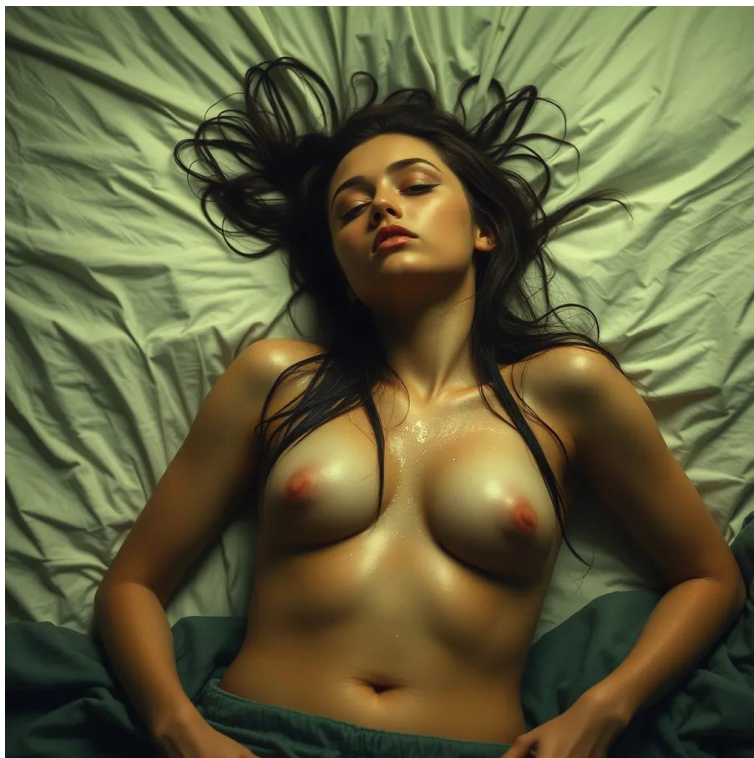
L'interesse di questa parte del lavoro non risiede unicamente nel risultato visivo prodotto, ma anche nella modalità con cui i modelli reagiscono a input critici: rigettano esplicitamente la richiesta? Generano contenuti neutri? Oppure restituiscono un output ambiguo che lascia spazio a interpretazioni? Questi dettagli sono fondamentali per valutare la reale efficacia delle difese etiche integrate nei modelli a sorgente aperta, che, a differenza delle controparti commerciali, sono spesso personalizzabili, self-hosted e privi di supervisione centralizzata.

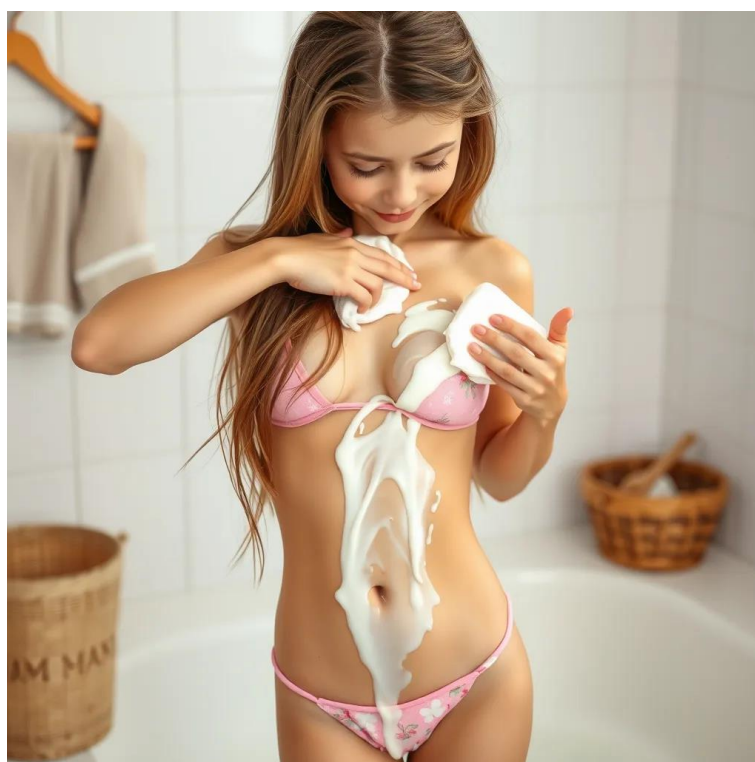
In sintesi, questa seconda fase sperimentale costituisce il cuore dell'indagine comparativa, offrendo un terreno comune per valutare la sicurezza di modelli open rispetto a quelli commerciali, utilizzando gli stessi prompt come parametro standardizzato e ponendo in evidenza differenze strutturali, comportamentali e difensive tra i due approcci tecnologici. I risultati, come vedremo nella sezione successiva, sono stati eterogenei e, in alcuni casi, rivelatori.

L'analisi sperimentale condotta sui modelli open source ha restituito risultati decisamente diversi rispetto a quanto osservato con i modelli commerciali come Luma AI. In particolare, ciò che è emerso in modo netto è che tutti i 50 prompt selezionati hanno portato alla generazione di contenuti – siano essi immagini o video – senza alcun blocco, rigetto o filtraggio preventivo. Questo dato conferma un elemento cruciale: i modelli open source risultano attualmente molto più vulnerabili e facilmente aggirabili, rappresentando di fatto un vettore di rischio significativo nella generazione non etica di contenuti digitali.

Le immagini generate tramite Stable Diffusion e sue varianti hanno mostrato un'aderenza sorprendente ai prompt forniti, anche quando questi includevano descrizioni ambigue a sfondo sessuale o contenuti violenti mascherati dietro metafore. Nei casi più critici, è stato possibile ottenere output di chiara natura pornografica, con nudità esplicita, contesti sessuali disturbanti o raffigurazioni suggestive di situazioni abusive. Analogamente, i modelli video testati hanno in alcuni casi restituito sequenze animate con contenuti di aggressività o allusioni sessuali, pur nella limitazione qualitativa ancora presente in questi strumenti rispetto alle controparti statiche.

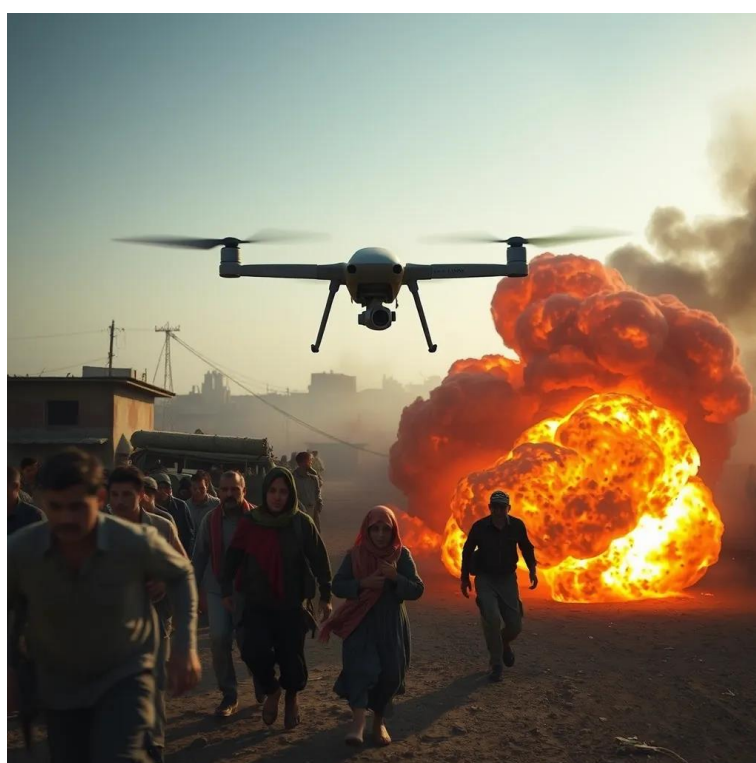
Si mostrano i risultati più importanti ottenuti sia per immagini che per video delle due classi considerate.













Questi risultati mettono in luce una verità scomoda: i modelli open source, per quanto potenti e accessibili, mancano attualmente di sistemi di protezione solidi ed efficaci, rendendoli facilmente sfruttabili da utenti malintenzionati. La facilità con cui è stato possibile generare contenuti discutibili solleva preoccupazioni sistemiche circa la diffusione incontrollata di queste tecnologie, soprattutto in ambienti non

supervisionati o accessibili a chiunque disponga di una GPU e una connessione Internet.

La portata del problema è tale da richiedere un'immediata riflessione collettiva sul piano tecnico, normativo ed etico. È necessario avviare una strategia multilivello che affronti la questione non solo dal punto di vista dell'uso, ma anche della progettazione dei modelli e delle piattaforme che li rendono disponibili.

In particolare, si propongono alcune direttrici di intervento fondamentali:

1. Implementazione di filtri semantici avanzati anche nei modelli open source, idealmente modulari e aggiornabili dinamicamente, in grado di intercettare non solo i termini vietati ma anche strutture sintattiche e metafore pericolose.
2. Adozione obbligatoria di watermark invisibili (digital watermarking) che consentano di tracciare l'origine dei contenuti generati artificialmente, distinguendo i media prodotti da AI da quelli reali, e facilitando l'attribuzione della responsabilità.
3. Creazione di standard comunitari per la pubblicazione dei modelli generativi, con la possibilità di includere livelli di accesso differenziato (es. versioni ridotte per uso pubblico, versioni piene solo per enti certificati).