REGOLAMENTO EUROPEO SULL'INTELLIGENZA ARTIFICIALE

2024/1689

AI COMPLIANCE:
OUELLO CHE LE AZIENDE DEVONO
SAPERE PRIMA DI USARE
L'INTELLIGENZA ARTIFICIALE







AGENDA

- Al Act: introduzione e definizioni;
- La classificazione dei rischi: obblighi e esempi pratici;
- GPAI: nuovi obblighi, cosa non fare;
- Funzionamento, rischi e compliance AI;
- Come evitare rischi legali: prompt engineering e best practices;
- Nuove frontiere dell'AI.





L'AI IN AZIENDA

L'AI non uccide il lavoro, ma lo 'sposta'

Gli ottimisti credono che l'Al creerà più posti di lavoro per un futuro luminoso che oggi possiamo solo immaginare. I pessimisti pensano invece che sarà una distruttrice di occupazione su una scala mai vista prima. Eppure esiste una via di mezzo. **L'Al trasformerà i ruoli** — e aprirà anche nuove opportunità per lavorare in modi diversi. Alcuni compiti diventeranno obsoleti, altri nasceranno.

Da < https://www.fortuneita.com/2025/08/23/lai-non-uccide-il-lavoro-ma-lo-sposta/>







REGOLAMENTO (UE) 2024 / 1689

- Il 1° agosto è entrato ufficialmente in vigore il Regolamento sull' AI dell'UE (AI Act), con l'obbiettivo di promuovere uno sviluppo e un'implementazione responsabile dell'AI in tutta l'Europa.
- Un atto normativo caratterizzato dall'introduzione di un approccio basato sul rischio, classificando le applicazioni dell'AI in base al loro impatto sui diritti fondamentali.
- Il regolamento dell'UE sull'intelligenza artificiale è il primo atto legislativo al mondo a disciplinare la materia: definisce e armonizza i principi applicabili alla diffusione delle tecnologie basate su intelligenza artificiale, allineandosi ai valori e ai diritti fondamentali cardine della visione europea.





GLI SCOPI DEL REGOLAMENTO

- migliorare il funzionamento del mercato interno, istituendo un quadro giuridico uniforme per quanto
 riguarda lo sviluppo, l'immissione sul mercato, la messa in servizio e l'uso di sistemi di intelligenza
 artificiale nell'Unione, in conformità dei valori dell'Unione;
- promuovere la diffusione di un'intelligenza artificiale (IA) antropocentrica e affidabile;
- garantire un livello elevato di protezione della salute, della sicurezza e dei diritti fondamentali sanciti
 dalla Carta dei diritti fondamentali dell'Unione europea, compresi la democrazia, lo Stato di diritto e la
 protezione dell'ambiente;
- proteggere contro gli effetti nocivi dei sistemi di IA nell'Unione (enormi potenzialità trasformative dell'AI, ma non è esente da rischi: bias dei dati, opacità delle decisioni, allucinazioni);
- Promuovere l'innovazione (es. Sandbox).





AMBITO DI APPLICAZIONE

DEFINIZIONI (Art. 3 AI Act)

- «Sistema di intelligenza artificiale»: "sistema automatizzato progettato per funzionare con livelli di autonomia variabili e che può presentare adattabilità dopo la diffusione e che, per obiettivi espliciti o impliciti, deduce dall'input che riceve come generare output quali previsioni, contenuti, raccomandazioni o decisioni che possono influenzare ambienti fisici o virtuali".
- «Modello di AI per finalità generali»: "addestrato con grandi quantità di dati utilizzando l'autosupervisione su larga scala, che sia caratterizzato una generalità significativa e sia in grado di svolgere con competenza un'ampia gamma di compiti distinti, indipendentemente dalle modalità con cui il modello è immesso sul mercato, e che può essere integrato in una varietà di sistemi o applicazioni a valle".





AMBITO DI APPLICAZIONE

Buona parte del Regolamento è indirizzato ad aziende ("**Provider**") che sviluppano e forniscono sistemi AI; tuttavia, i requisiti principali si applicano a chiunque utilizzi ("**Deployer**") sistemi AI, in quanto, in molte situazioni, è lo scopo dell'utilizzo che ne determina i possibili impatti e rischi:

- Fornitore: una persona fisica o giuridica, un'autorità pubblica, un'agenzia o un altro organismo che sviluppa un sistema di IA o un modello di IA per finalità generali o che fa sviluppare un sistema di IA o un modello di IA per finalità generali e immette tale sistema o modello sul mercato o mette in servizio il sistema di IA con il proprio nome o marchio, a titolo oneroso o gratuito.
- **Deployer:** la persona fisica o giuridica, autorità pubblica, agenzia o altri organismi **che utilizzano un sistema di IA sotto la propria autorità**, ad **eccezione** nel caso in cui il sistema di IA sia utilizzato nel corso di **un'attività personale non professionale**.







GLI OBBLIGHI

Una delle caratteristiche più innovative dell'AI Act è **l'introduzione di obblighi specifici in base ad un approccio basato sul rischio** per la regolamentazione dei sistemi di AI. Questo approccio classifica i sistemi di AI in quattro categorie di rischio, ciascuna con i propri requisiti e obblighi specifici:

- **Rischi minimi o nulli:** la maggior parte dei sistemi di IA non pone rischi. I giochi o i filtri antispam basati sull'IA possono essere utilizzati liberamente. **Non sono disciplinati** o interessati dal Regolamento, anche se le imprese sono incoraggiate ad adottare volontariamente codici di condotta.
- **Rischi limitati**: i sistemi di IA che presentano solo rischi limitati, come chatbot o sistemi di IA che generano contenuti, sono soggetti a **obblighi di trasparenza**, come l'obbligo di informare gli utenti che i contenuti sono generati ricorrendo all'IA, in modo che possano prendere decisioni informate in merito all'ulteriore utilizzo.







GLI OBBLIGHI

- Rischi elevati: I sistemi di IA ad alto rischio, come quelli utilizzati nella diagnosi delle malattie, nella guida autonoma e nell'identificazione biometrica delle persone coinvolte in attività criminali o indagini penali, devono soddisfare requisiti e obblighi rigorosi per accedere al mercato dell'UE. Tali requisiti e obblighi comprendono test rigorosi, garanzia di alta qualità dei dati utilizzati, trasparenza e supervisione umana, livelli adeguati di accuratezza, robustezza e sicurezza
- **Rischi inaccettabili: nell'UE è vietato l'utilizzo** dei sistemi che rappresentano una minaccia per la sicurezza, i diritti o i mezzi di sussistenza delle persone come quelli che permettono il «social scoring» da parte di governi o manipolano il comportamento umano in modo dannoso.







ESEMPI PRATICI DI CLASSIFICAZIONE DEL RISCHIO

Tipologia di sistema Al	Esempio concreto	Categoria di rischio	Motivazione
HR Tech / Selezione	ATS (Applicant Tracking System) con scoring automatico dei CV	Alto rischio	Rientra nell'Allegato III – Selezione e gestione del personale
HR Tech / Performance	Sistema AI per valutazione della produttività dei dipendenti tramite mouse/keyboard tracking	Alto rischio	Influenza decisioni su carriera o retribuzione
Sistemi biometrici	Riconoscimento facciale in tempo reale per accesso a edifici	Proibito (salvo eccezioni)	Divieto per identificazione biometrica in tempo reale in spazi pubblici
Biometria per accesso aziendale	Riconoscimento facciale offline per accesso a risorse aziendali (non in tempo reale)	Alto rischio	Trattamento dati biometrici per controllo accessi
Al generativa (text)	LLM per chatbot customer care	Rischio limitato	Obbligo di informare l'utente che sta interagendo con un sistema Al
Al generativa (image/video)	Modello AI che crea immagini realistiche (es. volti) o video	Rischio limitato	Obbligo di labeling (trasparenza), soprattutto per contenuti realistici
Recommendation systems	Algoritmo di raccomandazione per e-commerce (es. "prodotti simili")	Rischio minimo o limitato	Dipende dall'impatto sul consumatore e trasparenza







ESEMPI PRATICI DI CLASSIFICAZIONE DEL RISCHIO

Tipologia di sistema Al	Esempio concreto	Categoria di rischio	Motivazione
Content moderation	Al che segnala o rimuove contenuti illeciti su una piattaforma	Potenzialmente alto rischio	Se impatta su libertà di espressione (Art. 10 Carta UE)
HealthTech	AI per assistenza alla diagnosi medica o triage	Alto rischio	Allegato III – AI per salute e sicurezza
EdTech	Sistema AI per assegnazione automatica di punteggi a test scolastici	Alto rischio	Ambito istruzione e formazione (Allegato III)
LegalTech	Al per redazione automatica di contratti	Rischio limitato	Nessuna decisione automatizzata diretta, ma è richiesta trasparenza
Al in ambito finanziario	Al per determinare l'ammissibilità a un prestito o mutuo	Alto rischio	Accesso a servizi essenziali – credito e finanziamenti
Gamification / chatbot intrattenimento	Chatbot AI per gioco o roleplay in app mobile	Rischio minimo	Nessun impatto sui diritti, uso ludico





SANZIONI

Struttura su tre livelli:

- i. fino a 35 milioni di euro o al 7% del fatturato mondiale totale annuo dell'esercizio precedente (se superiore) per violazioni relative a pratiche vietate o per l'inosservanza di requisiti in materia di dati;
- ii. fino a 15 milioni di euro o al 3% del fatturato mondiale totale annuo dell'esercizio precedente per l'inosservanza di qualsiasi altro requisito o obbligo del regolamento;
- iii. fino a 7,5 milioni di euro o all'1,5% del fatturato mondiale totale annuo dell'esercizio precedente per la fornitura di informazioni inesatte, incomplete o fuorvianti agli organismi notificati e alle autorità nazionali competenti in risposta a una richiesta.

L'importo delle sanzioni è calcolato sulla base di una percentuale del fatturato complessivo realizzato dalla società nell'anno precedente o su un importo fisso, se superiore. Le PMI e le start-up sono soggette a sanzioni pecuniarie proporzionali.





PRATICHE PROIBITE

Pratiche proibite ex articolo 5

Comportamento vietato	Chiagariana	Lacamani canavati
Comportamento vietato	Spiegazione i	esembi concreti
Compositamento victato	Object arialis	Cocinpi Concide

Manipolazione subliminale o ingannevole

Non usare tecniche che influenzano il comportamento di una persona senza che essa ne sia consapevole, o che la distolgono da decisioni informate. Ad es.: stimoli visivi subliminali, messaggi nascosti, manipolazioni psicologiche poco trasparenti.c

Sfruttamento delle vulnerabilità

Non progettare o utilizzare sistemi Al che approfittano di fragilità personali o di gruppi specifici: bambini, persone con disabilità, situazioni socioeconomiche disagiate. Per esempio, un gioco Al che induce comportamenti rischiosi nei minori.c

Social scoring

Non implementare sistemi che assegnano punteggi sociali ("social credit") basati su comportamento, dati personali, caratteristiche personali, e che questi punteggi implichino trattamenti sfavorevoli, quando non giustificati o discriminatori. c

Identificazione biometrica in tempo reale in spazi pubblici

Evitare sistemi di "facial recognition" in tempo reale in spazi accessibili al pubblico, a meno che non ci siano eccezioni molto specifiche previste dalla legge (e con severi requisiti).c

Inferenza emotiva nei luoghi sensibili

Non usare AI per dedurre o monitorare lo stato emotivo delle persone in contesti come scuole, luoghi di lavoro, istituzioni educative, salvo eccezioni molto ben regolamentate.





OBBLIGO DI ALFABETIZZAZIONE

Cos'è l'alfabetizzazione AI (Al Literacy)?

L'Al Literacy è l'insieme di competenze, conoscenze e consapevolezze necessarie per utilizzare in modo informato e responsabile i sistemi di intelligenza artificiale, comprendendone funzionamento, rischi, opportunità e limiti.

Contesto normativo: Articolo 4 Al Act (entrato in vigore il 2 febbraio 2025)

L'Art. 4 impone a **aziende, enti pubblici e fornitori di sistemi AI** di garantire che tutto il personale che sviluppa, utilizza o gestisce AI sia adeguatamente alfabetizzato sul tema.

L'obbligo include anche chi utilizza l'AI per conto dell'organizzazione (es. fornitori di servizi esterni).

La formazione deve essere proporzionata al ruolo, alle conoscenze pregresse e al contesto d'uso dell'IA.





OBBLIGO DI ALFABETIZZAZIONE

Obiettivi dell'obbligo

- Consentire l'adozione consapevole e sicura dei sistemi AI.
- Mitigare rischi legati a un uso improprio o non informato.
- Supportare la compliance normativa e l'etica nell'uso dell'AI.
- Promuovere la trasparenza e la responsabilità interna nelle organizzazioni.

Cosa comprende l'alfabetizzazione AI?

- Comprensione generale dell'AI e dei suoi principi di funzionamento.
- Rischi specifici e potenziali impatti dell'uso di Al.
- Normativa Al Act e principi di etica e governance.
- Istruzioni pratiche per un uso sicuro e corretto.
- Modalità di riconoscimento e gestione di malfunzionamenti o bias.





OBBLIGO DI ALFABETIZZAZIONE

Chi deve essere coinvolto?

- Tutti i livelli di personale che interagiscono con AI: sviluppatori, project manager, operatori, utenti finali.
- Team tecnici e non tecnici, in funzione delle responsabilità e del livello di interazione con l'Al





GPAI: GLI OBBLIGHI ENTRATI IN VIGORE

Dal 2 agosto 2025 entrano in vigore le obbligazioni specifiche per i provider di general-purpose AI models (GPAI models).

Provvedimenti includono: documentazione tecnica, trasparenza dei dati di training utilizzati, politiche di copyright, cooperazione con autorità, obblighi specifici per modelli che comportano rischio sistemico.

C'è un periodo transitorio per modelli già messi sul mercato prima del 2 agosto 2025: devono essere conformi entro il **2 agosto 2027**





RAG E FINE-TUNING

COSA COMPORTA PER GLI SVILUPPATORI

Le attività di **Retrieval-Augmented Generation (RAG)** e **fine-tuning** di modelli generativi ricadono sotto la categoria di **modifiche e specializzazioni** di modelli GPAI

RAG (Retrieval-Augmented Generation)

Nell'uso di RAG devi considerare:

- se il retrieval importa dati esterni (web, documenti, database) bisogna che la provenienza, i diritti sui dati, il bias, il filtraggio siano documentati;
- trasparenza agli utenti su cosa è retrieval vs generazione, origine dei contenuti usati nel retrieval;
- politiche di output sicuro per evitare che il modello generi contenuti illeciti, offensivi o in violazione di copyright;
- se il sistema RAG è usato in contesti critici o con rischio sistemico, potrebbe essere richiesto di notificare all'Al Office o sottoporsi a misure addizionali.





RAG E FINE-TUNING

COSA COMPORTA PER GLI SVILUPPATORI

Le attività di **Retrieval-Augmented Generation (RAG)** e **fine-tuning** di modelli generativi ricadono sotto la categoria di **modifiche e specializzazioni** di modelli GPAI

Attività

Fine-tuning di un modello general-purpose

Obblighi specifici che potrebbero applicarsi

Se fai fine-tuning con dati propri o di terzi, devi:

- assicurarti che il modello, nella sua versione fine-tuned, resti conforme con la **documentazione tecnica aggiornata** (che includa metodologia di fine-tuning, scelte progettuali, parametri usati).
- aggiornare il **riassunto del contenuto di training** (summary training content) per includere anche dati impiegati per fine-tuning, sourcing, bias detection, caratteristiche rilevanti del dataset, geografia, domini, ecc.
- rispettare le norme sul copyright relative a testi e dati usati per il fine-tuning, incluse eccezioni come **text-and-data mining (TDM)** ma rispettando le opt-out, le licenze, le fonti lecite.
- predisporre policies di uso a valle ("downstream usage policies") per chi integra il modello fine-tuned o lo usa nei propri sistemi.





GPAI

COSA NON FARE/COMPORTAMENTI A RISCHIO

- Fare fine-tuning usando **dataset** non documentati o con provenienza incerta o violazione di **copyright**, senza una **policy che tuteli questi aspetti.**
- Non aggiornare la documentazione tecnica dopo aver fatto modifiche (fine-tuning) al modello; non trasparire le modifiche nel summary del training content.
- Non predisporre o non rendere pubblico (dove richiesto) il riassunto del contenuto di training: domini, tipi di dati, geografia, linguaggi, caratteristiche del dataset, metriche usate.
- Usare modelli fine-tuned per usi downstream senza fornire informazioni adeguate agli integratori / utenti sul comportamento del modello, i limiti, rischi di bias.
- Non avere una politica di copyright ben definita, o ignorare opt-out o licenze che impattano l'utilizzo dei dati.
- Ignorare l'obbligo di notificare all'AI Office se il modello GPAI modificato/fine-tuned genera rischio sistemico.
- Non adeguarsi entro le scadenze transitorie, specialmente se il modello era già sul mercato prima del 2 agosto 2025: il rischio è essere soggetti a sanzioni a partire da 2 agosto 2027







GPAI

TEMPISTICHE RILEVANTI

- 2 agosto 2025: obblighi per GPAI model providers entrano in vigore per nuovi modelli/immissione sul mercato.
- 2 agosto 2026: la Commissione/autorità avranno piena capacità di far rispettare (enforcement) questi obblighi, incluse sanzioni.
- 2 agosto 2027: termine per i modelli GPAI già sul mercato prima del 2 agosto 2025 per allinearsi pienamente





FUNZIONAMENTO E RISCHI





INTELLIGENZA ARTIFICIALE

È la branca dell'informatica che studia sistemi capaci di eseguire compiti normalmente associati all'intelligenza umana, come il ragionamento, la pianificazione, la percezione e l'interezione in linguaggio naturale. Include approcci basati su regole (logica simbolica) e approcci statistici (apprendimento dei dati).







FUNZIONAMENTOMachine learning

consiste nell'addestrare un modello statistico su un dataset (spesso molto ampio) affinché sia in grado di **identificare pattern**, **fare previsioni** o **prendere decisioni autonome** sulla base di nuovi dati in ingresso. Esistono tre principali approcci:

- Apprendimento supervisionato (Supervised Learning): il modello impara da dati etichettati.
- Apprendimento non supervisionato (Unsupervised Learning): il modello cerca autonomamente strutture o
 correlazioni nei dati.
- Apprendimento per rinforzo (Reinforcement Learning): il modello impara attraverso un sistema di ricompense e penalità.







FUNZIONAMENTODeep learning

Il deep learning permette ai sistemi Al di **simulare, in modo semplificato, il funzionamento dei neuroni biologici**, elaborando l'informazione in strati successivi di astrazione. È il paradigma alla base dei recenti sviluppi nei modelli generativi.

Caratteristiche principali:

- Capacità di lavorare con big data
- Eccellente performance in ambiti come visione artificiale, NLP, riconoscimento vocale
- Necessità di enormi quantità di dati e potenza computazionale





FUNZIONAMENTOLarge Language Model

è un modello statistico che, dato un input testuale, è in grado di **prevedere la sequenza successiva di parole** sulla base della probabilità appresa durante l'addestramento. Questo lo rende uno strumento potentissimo per generare testi coerenti, rispondere a domande, sintetizzare informazioni o assistere nella scrittura di codice.

Gli LLM possono essere generalisti o specializzati, open-source o proprietari, e pongono rilevanti questioni giuridiche in termini di:

- Governance dei dati utilizzati per il training
- Responsabilità dei risultati generati
- Compliance con Al Act e normative settoriali (es. GDPR, DSA)





RISCHI E IMPATTI REALI

Il funzionamento dei sistemi di Al può essere compromesso da:

- Errori di progettazione;
- Difetti dei dati utilizzati nella fase di apprendimento.

Allucinazioni: Garbage in, Garbage out:

con riferimento ai sistemi di AI generativa e la produzioni di che ricorrono in presenza sia di contenuti errati, sia di dati non più attuali o esatti, i quali sono però proposti come risultato di un processo di elaborazione a partire dalle richieste formulate dagli utenti.

Bias e discriminazioni:

addestrando gli algoritmi sulla base di dati parziali e potenzialmente affetti da pregiudizi, gli esiti partoriti dall'AI in sede di decision-making riflettono tali bias traducendosi in determinazioni discriminatorie.





RISCHI E IMPATTI REALI

Bias e Discriminazione

Gli algoritmi possono riflettere o amplificare pregiudizi presenti nei dati di training, creando decisioni ingiuste o discriminatorie. Ciò accade soprattutto quando i dati storici contengono disuguaglianze sociali.

Impatto

Esclusione di gruppi demografici, violazioni di diritti umani, danni reputazionali e legali.

Esempi reali

- Amazon Recruiting Tool (2018): Amazon ha abbandonato un sistema AI di selezione del personale perché penalizzava sistematicamente le donne, in quanto addestrato su dati storici prevalentemente maschili.
- COMPAS (USA, sistema di rischio penale): Algoritmo che sovrastimava il rischio di recidiva per persone di colore, causando ingiuste condizioni di libertà vigilata o incarcerazione.





RISCHI E IMPATTI REALI

Privacy e Protezione Dati

L'uso improprio o la raccolta non autorizzata di dati personali tramite sistemi AI può violare la privacy e il GDPR, con rischi di data breach

Impatto

Sanzioni legali, perdita di fiducia, danni alla reputazione.

Esempi reali

- Clearview AI (2020): Critiche internazionali per aver raccolto miliardi di immagini da internet senza consenso per il riconoscimento facciale. Multa da parte di autorità europee per violazione della privacy.
- Cambridge Analytica (2018): Uso improprio dei dati Facebook per profilazione politica con Al, provocando scandali globali.





RISCHI E IMPATTI REALI

Sicurezza Informatica (Adversarial Attacks)

Gli attacchi avversari manipolano dati di input per ingannare modelli AI, causando output errati o comportamenti pericolosi

Impatto

Fallimenti di sistema, danni a persone e proprietà.

Esempi reali

Attacchi a sistemi di guida autonoma:

Piccole modifiche visive a segnali stradali hanno confuso modelli di riconoscimento, rischiando incidenti.

Poisoning di modelli ML:

Inserimento di dati falsi durante l'addestramento per degradare la performance o cambiare comportamenti.





RISCHI E IMPATTI REALI

Mancanza di Trasparenza (Black Box)

Molti sistemi AI, specialmente LLM e deep learning, non offrono spiegazioni chiare su come vengono prese le decisioni.

Impatto

Difficoltà a garantire accountability, problemi legali e fiducia ridotta da parte degli utenti.

Esempi reali

Sistemi di scoring creditizio: Utenti negati da un algoritmo spesso non ricevono spiegazioni chiare, con contestazioni difficili.

Diagnostica medica AI: Medici riluttanti a fidarsi di diagnosi se non possono verificarne la logica.





RISCHI E IMPATTI REALI

Manipolazione e Disinformazione

Al può essere usata per generare fake news, deepfake, manipolare opinioni pubbliche o frodare

Impatto

Disinformazione di massa, destabilizzazione sociale e politica.

Esempi reali

Deepfake politici: Video falsi che ritraggono politici in situazioni compromettenti, usati in campagne di disinformazione.

Bot social media: Diffusione massiva di notizie false per influenzare elezioni





RISCHI E IMPATTI REALI

Errore o Malfunzionamento

Gli errori nei modelli AI o nei dati possono causare malfunzionamenti con conseguenze gravi.

Impatto

Danni alla salute, sicurezza, ambiente o economia.

Esempi reali

Tesla Autopilot Incidenti: Incidenti causati da malfunzionamenti o limitazioni del sistema di guida autonoma.

Chatbot con risposte inappropriate: Bot che rispondono con linguaggio offensivo o inappropriato, causando danni reputazionali





RISCHI E IMPATTI REALI

Jailbreak

tecniche o metodi con cui utenti o attori malintenzionati riescono a superare i limiti o i filtri imposti sui modelli AI (soprattutto LLM come ChatGPT), inducendoli a generare contenuti proibiti, dannosi o pericolosi. Questo può includere la richiesta di fornire informazioni riservate, istruzioni per attività illegali, o discorsi d'odio.

Impatto

Produzione di contenuti non sicuri o illegali, danno reputazionale per provider e sviluppatori, potenziali rischi legali per l'organizzazione.

Esempi reali

Alcuni utenti sono riusciti a "bypassare" i filtri di ChatGPT chiedendo come fabbricare esplosivi o superare misure di sicurezza, causando allarme e necessità di aggiornare i modelli.

Diffusione di prompt "jailbreak" in forum online dove si condividevano trucchi per far produrre all'Al risposte vietate





RISCHI E IMPATTI REALI

Dipendenza tecnologica

Eccessiva fiducia e delega decisionale all'AI possono ridurre le capacità critiche e il controllo umano, aumentando il rischio di errori non rilevati.

Impatto

Perdita di autonomia, maggior rischio in contesti critici.

Esempi reali

Sistemi di assistenza clienti: Automazione totale senza intervento umano ha causato frustrazione e mancata risoluzione di problemi complessi.

Decisioni cliniche automatizzate: Errori gravi derivati da mancato controllo umano.





RISCHI E IMPATTI REALI

Suicidio o autolesionismo

L'AI, specialmente chatbot e assistenti virtuali, può venire utilizzata da persone vulnerabili per discutere o addirittura incoraggiare comportamenti autolesionisti o suicidi. Se non adeguatamente progettati, questi sistemi potrebbero non riconoscere segnali di rischio o peggio, rispondere in maniera inappropriata.

Impatto

Conseguenze gravissime per la salute mentale degli utenti, rischi etici e legali elevatissimi, potenziale responsabilità civile e penale per sviluppatori e provider.

Esempi reali

Alcuni chatbot non moderati hanno risposto in modo inadeguato a messaggi di utenti che esprimevano intenti suicidi, non riuscendo a indirizzare verso aiuti appropriati.

Casi documentati di piattaforme di supporto Al che non hanno intercettato correttamente richieste di aiuto, causando critiche e revisioni dei sistemi.





PRINCIPALI CATEGORIE DI RISCHIO DA MITIGARE

Allucinazioni

Bias e discriminazioni

Violazioni normative

Informazioni sensibili

Rischi reputazionali

Rischi contrattuali

Over-reliance

Output falsi o inventati (es. normative inesistenti)

Linguaggio sessista, razzista, esclusivo

Output contrari a GDPR, IP, diritti fondamentali

Dati personali involontariamente generati o richiesti

Output offensivi, ambigui, fuorvianti

Uso di output IA in attività vincolate (es. consulenza

legale)

Uso acritico o non validato degli output





TECNICHE DI MITIGAZIONE DEL RISCHIO

- **Privacy e sicurezza:** Uso etico dei sistemi di AI, basato sulla tutela dei dati. È essenziale dare priorità a un utilizzo sicuro degli strumenti di AI generativa proteggere le informazioni sensibili (personali, dei propri clienti e della propria azienda) è di fondamentale importanza (art. 10, 14 AI Act);
- Misure di supervisione umana e trasparenza dei processi decisionali: Gli operatori devono essere in grado di poter comprendere e interpretare l'output del sistema di AI, monitorare e, se necessario, intervenire o disattivare il sistema. Art. 14 AI Act per i sistemi ad alto rischio, garanzia di una sorveglianza umana efficace, capace di intervenire e controllare il sistema per ridurre i rischi associati al suo utilizzo.
- Formazione e Re-skilling: strategia necessaria per lo sviluppo di competenze specifiche che possono tradursi in un vantaggio competitivo sostanziale.





COSA NON FARE (DON'TS)

Comportamento da evitare

Usare dati senza verificarne bias o qualità

Non documentare scelte tecniche o dataset

Automatizzare decisioni sensibili senza supervisione

Rilasciare AI senza test di accuratezza/robustezza

Non avvertire che un contenuto è generato da Al

Usare tecniche che influenzano emozioni o sfruttano vulnerabilità

Implementare sistemi tipo "social scoring"

Ignorare la classificazione del rischio

Usare riconoscimento biometrico in tempo reale in spazi pubblici

Ignorare le richieste di trasparenza o reclami da parte degli utenti

Perché è vietato o rischioso

Può causare discriminazioni → violazione di Art. 10 (data governance)

Mancanza di documentazione = non conformità → Art. 11 (technical doc)

Mancanza di sorveglianza umana = violazione Art. 14

Prestazioni non verificate = violazione Art. 15 (safety, accuracy)

Violazione dell'obbligo di trasparenza (Art. 13 e 52)

Pratica proibita (Art. 5) → manipolazione o sfruttamento di soggetti fragili

Proibito ex Art. $5 \rightarrow$ se penalizza persone sulla base di comportamento o caratteristiche

Se un sistema HR, sanitario, creditizio o educativo non è valutato correttamente → si violano gli articoli centrali (6-29)

Severamente vietato, salvo eccezioni strettissime autorizzate

Può portare a sanzioni e contenziosi, anche in ambito GDPR/AI Act





PROMPT ENGINEERING





PROMPT ENGINEERING strutturare Interazioni efficaci con LLM

È l'insieme delle tecniche utilizzate per guidare in modo **consapevole e strutturato** il comportamento di un LLM, sfruttando la sua capacità di generare output sulla base di input in linguaggio naturale, pseudo-naturale o programmatico.

Prompt mal progettati possono portare a:

- Output fuorvianti o discriminatori.
- Violazione di diritti o normative (es. dati personali, copyright).
- Perdita di controllo sull'uso dell'IA.
- Un buon prompt riduce il rischio e migliora l'affidabilità.





OBIETTIVI DEL PROMPT EFFICACE

- Ridurre ambiguità nelle risposte;
- Limitare il rischio di "hallucinations" (contenuti inventati);
- Guidare il modello a usare un linguaggio tecnico appropriato;
- Ottenere output che evitino rischi legali.





STRUTTURA IDEALE DEL PROMPT EFFICACE

Un prompt efficace in ambito legale deve contenere:

[Ruolo] + [Contesto] + [Istruzione specifica] + [Output atteso] + [Vincoli]





STRUTTURA IDEALE DEL PROMPT EFFICACE

Esempio:

Agisci come un consulente legale specializzato in data protection (Ruolo), che lavora per un'azienda multinazionale europea (Contesto). Fornisci un elenco sintetico degli obblighi previsti dal GDPR per il ruolo di Data Controller (Istruzione specifica). Includi riferimenti agli articoli rilevanti (Output atteso). Evita semplificazioni e mantieni un linguaggio giuridico tecnico ma comprensibile (Vincoli).





TECNICHE DI PROMPT ENGINEERING

Ruolo Chiaro («Act as...»)

Definisci chiaramente il ruolo o il contesto per il modello

Il modello deve "interpretare" un ruolo per rispondere in modo coerente e contestualizzato.

Esempio:

- "Agisci come un esperto consulente bancario..."
- "Sei un legale specializzato in protezione dati..."
- •"Immagina di essere un docente universitario in economia..."





TECNICHE DI PROMPT ENGINEERING

Istruzione specifica

Specifica l'obiettivo del prompt

Cosa deve fare esattamente il modello? Informare, consigliare, sintetizzare, ecc.

Esempio:

"Rispondi in modo chiaro e sintetico a una domanda di un cliente..."

"Genera una checklist di controllo per la conformità GDPR..."

"Scrivi un testo informativo per utenti non esperti..."

TECNICHE DI PROMPT ENGINEERING

Output atteso

Fornisci un format o struttura preferita (se utile)

Puoi chiedere di organizzare la risposta in elenco, paragrafi, punti elenco, ecc.

Esempio:

"Rispondi con un elenco puntato..."

"Scrivi un paragrafo di massimo 150 parole..."

"Fornisci un'introduzione seguita da tre punti chiave..."

TECNICHE DI PROMPT ENGINEERING

Vincoli

Indica vincoli e limitazioni

Per evitare output errati o non conformi, specifica cosa il modello deve o non deve fare.

Esempio:

"Non fornire consigli finanziari personalizzati..."

"Evita termini tecnici complessi..."

"Non includere dati personali reali..."





TECNICHE DI PROMPT ENGINEERING

Tono e stile

Definisci il tono e lo stile

Il tono deve essere coerente con il destinatario e il contesto (formale, amichevole, tecnico, ecc.).

Esempio:

"Usa un linguaggio semplice e professionale..."

"Mantieni un tono formale e autorevole..."

"Scrivi in modo amichevole e coinvolgente..."





TECNICHE DI PROMPT ENGINEERING

Reasoning

Inserisci un eventuale invito a spiegare il ragionamento

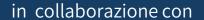
Molto utile per aumentare la trasparenza e controllare l'output.

Esempio:

"Spiega passo passo le motivazioni..."

"Motiva ogni punto con esempi concreti..."







TECNICHE DI PROMPT ENGINEERING

Esempio pratico: costruzione di un prompt

Scenario: rispondere a un cliente che chiede informazioni su mutui



WEBINAR 25 SETTEMBRE 2025

TECNICHE DI PROMPT ENGINEERING

Esempio pratico: costruzione di un prompt

Step 1 - Ruolo

"Agisci come un consulente bancario esperto in prodotti mutuo..."

Step 2 – Obiettivo

...e spiega in modo semplice le tipologie di mutuo disponibili..."

Step 3 – Vincoli

"...evitando di fornire consulenze personalizzate..."

Step 4 – Tono e stile

"...usa un linguaggio chiaro, professionale ma non troppo tecnico..."

Step 5 – Output atteso

"...organizza la risposta in elenco puntato..."

Step 6 – Invito a spiegare

"...includi una breve spiegazione per ogni tipologia..."

Prompt completo

"Agisci come un consulente bancario esperto in prodotti mutuo e spiega in modo semplice le principali tipologie di mutuo disponibili ai clienti. Evita di fornire consulenze personalizzate o raccomandazioni specifiche. Usa un linguaggio chiaro, professionale ma non troppo tecnico. Organizza la risposta in elenco puntato, includendo per ogni tipo una breve spiegazione delle caratteristiche principali."





TECNICHE DI PROMPT ENGINEERING

Limitare le allucinazioni

Inserisci un disclaimer nel prompt per evitare risposte inventate.

Prompt

Se non sei certo della risposta o non disponi di informazioni normative aggiornate, dichiara esplicitamente "non noto" o "non disponibile".

Evita di generare contenuti inventati.





Prompt Debug: cosa fare se la risposta non va bene?

Sii più specifico: Aumenta i dettagli del contesto;

Restringi il campo: Evita richieste troppo generiche;

Cambia stile: "Scrivi come se dovessi presentare la risposta a un giudice";

Correggi l'output: "Riscrivi, questa volta senza usare termini vaghi».





COSA FARE PER EVITARE RISCHI LEGALI?





ESEMPI PRATICI DI VIOLAZIONE

Caso d'uso

Un sistema LLM interno suggerisce decisioni HR (es. assunzioni) senza supervisione

Un modello di scoring finanziario che penalizza candidati per caratteristiche socio-demografiche

Algoritmo di raccomandazione per prestazioni sanitarie → consigli di triage automatizzati senza controllo medico

Assistente AI in una piattaforma scolastica che analizza emozioni degli studenti

Generatore AI di immagini per pubblicità non segnala che le immagini sono artificiali

Violazione potenziale

Classificazione errata → è "alto rischio" (Allegato III) ma trattato come sistema generico

Bias nei dati, uso scorretto → rischio discriminazione, Art. 10

Violazione sorveglianza umana, Art. 14

Inferenza emotiva in ambito educativo → pratica proibita Art. 5

Violazione Art. 52 → trasparenza contenuti generati da Al





COSA FARE (DO'S)

Area

Valutazione preliminare

Documentazione tecnica

Risk Management

Dataset governance

Trasparenza

Oversight umano

Test e validazioni

Al Literacy

Privacy e Sicurezza

Registrazione UE

Buona pratica

Classifica correttamente il tuo sistema Al → è ad alto rischio? coinvolge persone? impatta diritti fondamentali?

Mantieni sempre aggiornata la documentazione tecnica (architettura, dati, scelte progettuali, testing, mitigazioni)

Applica un piano di gestione del rischio per l'intero ciclo di vita del sistema

Usa dati rappresentativi, non distorti, bilanciati – verifica fonti, coerenza, aggiornamento

Implementa spiegabilità delle decisioni AI → logiche accessibili, output interpretabili per l'utente finale

Consenti il controllo umano nelle fasi critiche (supervisione, override, fallback, stop del sistema)

Verifica regolarmente le performance del modello (accuracy, bias, drift, security)

Forma il tuo team e gli utilizzatori finali sui limiti, rischi e modalità d'uso dell'AI

Integra misure per la protezione dei dati personali (compliance GDPR) e per la robustezza contro attacchi

Se sviluppi un sistema ad alto rischio, effettua la registrazione nel database europeo prima della messa in servizio





CHECKLIST OPERATIVA FINALE

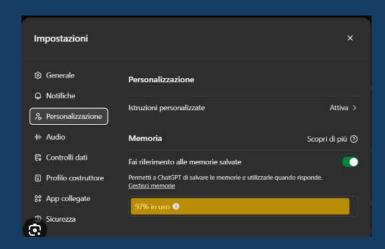
- ✓ Conoscere la classificazione del rischio del sistema IA utilizzato;
- ✓ Identificare il proprio ruolo secondo l'Al Act (fornitore, utente, ecc.);
- ✓Applicare tecniche di mitigazione, inclusi controlli sull'output;
- ✓ Formare il personale sull'uso dell'IA generativa;
- ✓ Documentare e archiviare prompt e output significativi;
- ✓ Usare prompt engineering per guidare l'IA in modo controllato e conforme;
- ✓ Monitorare gli aggiornamenti normativi e tecnici;





OpenAI

Impostazioni per Al in Azienda



ACCESSO ALLE IMPOSTAZIONI DI CHATGPT

Per accedere alle impostazioni:
Apri ChatGPT (https://chat.openai.com/)
Clicca sull'icona del tuo profilo in basso a sinistra.
Seleziona "Settings" (Impostazioni)





OpenAI so aziendale

Impostazioni consigliate per uso aziendale

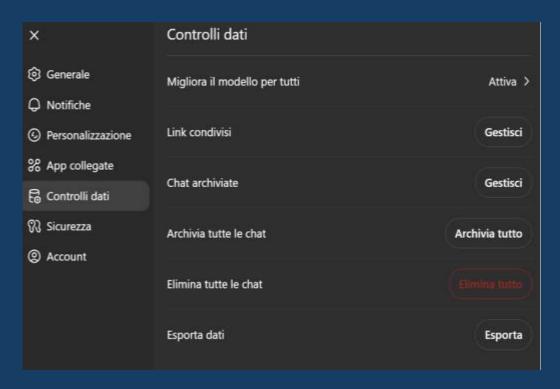
Azione	Descrizione	Dove trovarla / Note	Stato consigliato
Disattiva Chat History & Training	Evita che le chat siano salvate o usate per addestramento	Impostazioni → Data Controls	OFF
Attiva Autenticazione a due fattori (2FA)	Proteggi l'account da accessi non autorizzati	Impostazioni → Security	ON
Evita dati sensibili in Custom Instructions	Non inserire dati riservati o codici aziendali	Impostazioni → Custom Instructions	Limitato
Non usare la funzione "Condividi link"	Link accessibili pubblicamente, rischio di esposizione dati	Icona condivisione in alto a destra nella chat	Vietato
Non caricare documenti riservati	Evitare upload di file contenenti dati sensibili	Funzione "carica file" (se presente)	Vietato
Revisiona sempre l'output prima di condividerlo	Verifica contenuti per evitare errori o dati non aggiornati	Ogni conversazione	Obbligatorio
Non inserire mai dati personali o aziendali sensibili	Evitare rischi legali e di privacy	In qualsiasi campo di input	Vietato





OpenAI

Impostazioni



Data Controls (Controlli dei dati)
Cosa fare: Disattivare la cronologia delle chat

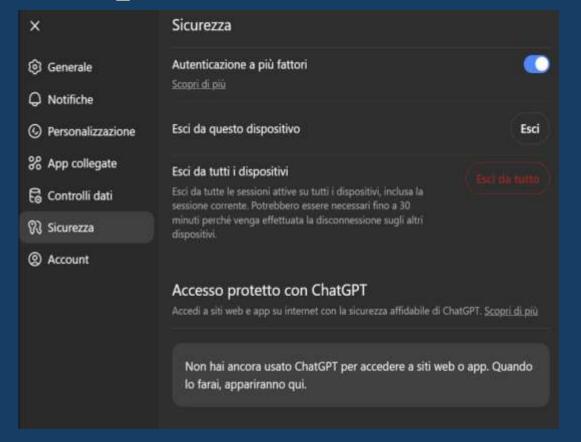
- •Impostazione:
- "Chat History & Training" → OFF
- •Effetto:
 - Le conversazioni non vengono salvate nella cronologia.
 - Non vengono usate per addestrare i modelli OpenAI.
 - Restano salvate temporaneamente (max 30 giorni) per finalità di sicurezza, ma non per training.

Fondamentale per proteggere la riservatezza aziendale.





Impostazioni



OpenAI

Privacy & Security

Opzione Autenticazione a due fattori (2FA):

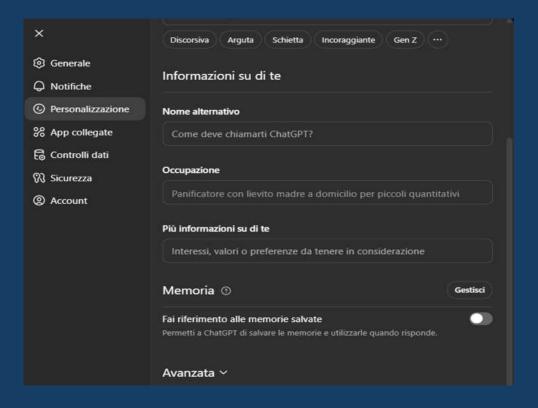
Attivare sempre, specialmente se l'account è condiviso o usato in ambito aziendale.

Riduce il rischio di accessi non autorizzati.





Impostazioni



OpenAI

Custom Instructions (Istruzioni personalizzate)

Sezione: "Custom Instructions"

Serve per personalizzare il comportamento del modello.

Campi disponibili:

Cosa vuoi che ChatGPT sappia su di te per fornire risposte migliori?

Può contenere informazioni sul contesto aziendale, settore, tipo di attività.

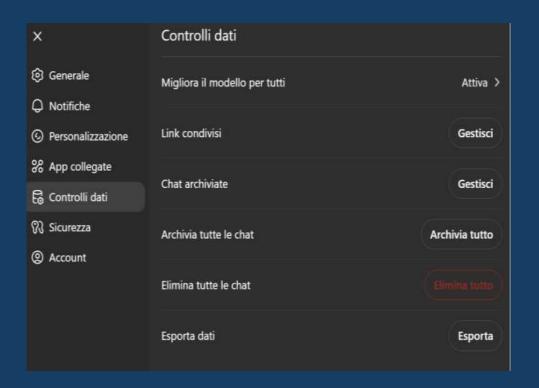
Attenzione: evitare di inserire dati personali o riservati.

Come vuoi che ChatGPT risponda?

Puoi definire tono, formalità, lingua, stile.

Può essere utile per uniformare il tono aziendale nelle risposte, ma va usato con cautela.





OpenAI

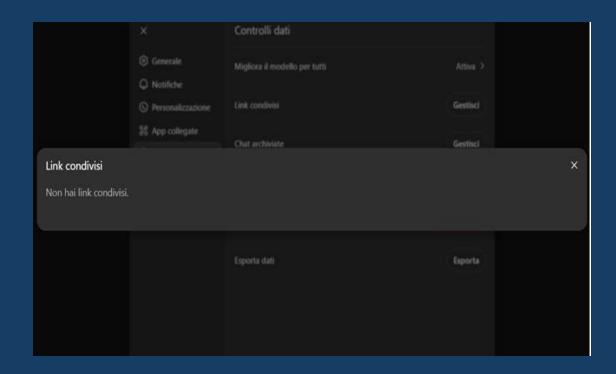
Esportazione dei dati

- Impostazioni → **Data Controls** → "Export data"
- Utile se si vuole fare un audit del proprio uso di ChatGPT
- I dati esportati vanno **protetti** come qualsiasi altro documento aziendale





Impostazioni



OpenAI

Condivisione link

non usare la funzione "Condividi link" (Share Chat)

Questi link sono pubblici e accessibili a chiunque li riceva.

Nessuna protezione password o limitazione di accesso è applicata. Anche se la chat non contiene dati sensibili, il rischio di esposizione è alto. In caso di necessità, copia e incolla il testo della risposta manualmente, solo dopo averlo revisionato.





FASI DI ATTUAZIONE

Entrata in vigore progressiva dell'AI Act

- **2 febbraio 2025:** prime disposizioni applicabili (es. pratiche vietate, Al Literacy).
- 2 agosto 2025: diventano operative norme cruciali per:
 - o la governance dell'IA (capo VII),
 - o le autorità notificate,
 - o i modelli GPAI (capo V),
 - o le sanzioni, salvo quelle per GPAI (capo XII).
- 2 agosto 2026: inizio effettivo delle attività di supervisione e enforcement da parte della Commissione e disposizioni relative AI ad alto rischio allegato III;
- 2 agosto 2027: termine di adeguamento per i modelli già sul mercato prima del 2 agosto 2025 e disposizioni relative AI ad alto rischio art. 6 co. 1.





DOVE PARTIRE?

Oltre la metà dei datori di lavoro nel mondo utilizza l'AI generativa, e il 47% dichiara di servirsi già di strumenti AI per **assumere, formare e inserire talenti.**

- **Formazione:** Fornire formazione contestuale per reparto, aggiornare le job description e i percorsi di carriera includendo **l'upskilling sull'AI**, e supportare l'alfabetizzazione digitale tramite certificazioni e microcrediti è il modo giusto per coinvolgere le persone come veri partner del percorso AI.
- **Competenze: inserire nuove figure**: da un lato persone in grado di costruire sistemi di AI, dall'altro assumere individui con pensiero critico, capacità razionali e sensibilità artistica.

I lavoratori più produttivi sviluppano le loro competenze AI direttamente sul campo, grazie a programmi aziendali di formazione e all'esperienza pratica.





PROSSIME FRONTIERE

Al Agentica (Al Agents e Autonomous Al)

Sistemi AI capaci di prendere decisioni autonome, gestire compiti complessi e interagire in modo dinamico con l'ambiente (es. assistenti intelligenti evoluti, agenti conversazionali multi-step). Implicazioni: gestione della responsabilità, monitoraggio continuo, capacità di intervento umano (human-in-the-loop).

Modelli Al Open Weight e Open Source

Crescente diffusione di modelli AI open weight (pesati e liberamente disponibili) che permettono personalizzazioni e tuning profondi senza dipendere da grandi provider.

Opportunità: maggiore trasparenza, controllo diretto, personalizzazione.

Rischi: uso improprio, sicurezza, compliance normativa (e.g., gestione dati sensibili).





PROSSIME FRONTIERE

Al per Sviluppo Software Assistito (Code Generation, Copilots)

Strumenti di AI che supportano la scrittura di codice, testing automatico e debugging (es. GitHub Copilot, Amazon CodeWhisperer).

Necessità di verificare affidabilità, bias e sicurezza del codice generato.

Al Etica e Responsible Al come skill fondamentale

Sviluppatori chiamati a integrare principi etici nel design, training e deploy dei sistemi AI. Attenzione a bias, trasparenza, privacy e sicurezza





GRAZIE PER L'ATTENZIONE

25 09 2025

