

Il sorso fantasma: il costo fisico nascosto dietro ogni risposta AI



Autore: Francesco Zinghini | Data: 19 Febbraio 2026

Siamo nel 2026 e l'interazione con l'**Intelligenza Artificiale Generativa** è diventata un riflesso quasi involontario, naturale come respirare o accendere la luce. Scriviamo un prompt, attendiamo una frazione di secondo e otteniamo un testo, un'immagine o un codice complesso. Tutto appare etereo, pulito, immateriale: una pura transazione di bit che viaggia nell'etere senza lasciare traccia apparente. Eppure, questa percezione di leggerezza è una delle più grandi illusioni del nostro tempo. Dietro l'interfaccia minimalista di un chatbot o di un assistente virtuale, si nasconde una realtà industriale pesante, fatta di metallo, silicio e, soprattutto, di un consumo di risorse fisiche che stiamo iniziando a quantificare solo ora con la dovuta gravità. C'è un costo che paghiamo per ogni singola risposta, un costo che non appare in bolletta ma che l'ambiente salda puntualmente.

L'illusione del Cloud e la realtà del Silicio

Quando parliamo di *progresso tecnologico* e automazione, tendiamo a visualizzare il “Cloud” come una nuvola soffice e astratta. In realtà, il Cloud è costituito da enormi cattedrali di cemento e acciaio: i Data Center. Queste strutture ospitano migliaia di server, ognuno equipaggiato con unità di elaborazione grafica (GPU) ad alte prestazioni, essenziali per il **deep learning** e l'addestramento delle reti neurali. L'architettura neurale di un moderno **LLM** (Large Language Model) richiede una potenza di calcolo mostruosa, non solo durante la fase di addestramento (training), ma anche e soprattutto durante

l'inferenza, ovvero il momento in cui il modello risponde alla domanda specifica di un utente.

Qui entra in gioco la fisica elementare. I chip che alimentano l'**AI** non sono diversi da un motore ad alte prestazioni: quando lavorano al massimo dei giri, generano calore. Molto calore. Per evitare che i processori fondono o vadano in protezione termica, rallentando i calcoli, devono essere raffreddati costantemente. Ed è qui che scopriamo il “sorso fantasma”.

Il Sorso Fantasma: svelare il mistero

La curiosità che ci ha spinto fin qui riguarda la natura fisica di questo costo. Non parliamo solo di elettricità, sebbene il consumo energetico sia enorme. Parliamo di una risorsa ancora più preziosa e tangibile: **l'acqua**. Il “sorso fantasma” è la quantità d'acqua dolce necessaria per raffreddare i server mentre elaborano la vostra richiesta.

Per mantenere la temperatura operativa ideale, molti data center utilizzano torri di raffreddamento evaporativo. Il principio è simile a quello della sudorazione umana: l'acqua viene fatta evaporare per sottrarre calore all'ambiente. Secondo studi accreditati condotti negli anni passati e confermati dalle metriche del 2026, una conversazione standard con un modello tipo **ChatGPT** (composta da circa 20-50 scambi di battute) “beve” circa **500 millilitri di acqua**. È l'equivalente di una bottiglietta di plastica standard.

Immaginate ora di avere una bottiglia d'acqua sulla scrivania e di rovesciarla a terra ogni volta che finite una breve sessione di lavoro con l'AI. Questo è il sorso fantasma. Moltiplicatelo per i miliardi di utenti attivi quotidianamente nel mondo e otterrete fiumi virtuali che vengono consumati, o meglio, fatti

evaporare, per sostenere la nostra sete di informazioni.

Come funziona la sete degli algoritmi

Per comprendere appieno il fenomeno, dobbiamo addentrarci nel funzionamento tecnico dei sistemi di raffreddamento e nelle metriche di efficienza. L'industria utilizza un indicatore chiamato **WUE (Water Usage Effectiveness)**, che misura quanti litri d'acqua vengono consumati per ogni kilowattora (kWh) di energia utilizzata dai server. L'efficienza varia enormemente in base alla posizione geografica del data center e alla stagione.

Un data center situato in una zona fredda come la Svezia potrebbe utilizzare l'aria esterna per raffreddare i server per gran parte dell'anno (free cooling), riducendo il consumo idrico. Tuttavia, un data center situato in zone calde o durante i mesi estivi deve fare affidamento massiccio sull'evaporazione dell'acqua. Quando inviate un prompt, la vostra richiesta viene instradata verso un server che potrebbe trovarsi a migliaia di chilometri di distanza. Se quel server si trova in una regione colpita da siccità o con temperature elevate, il "costo idrico" della vostra domanda aumenta drasticamente.

Inoltre, c'è un consumo idrico indiretto. L'elettricità che alimenta i server deve essere prodotta da qualche parte. Le centrali termoelettriche (a gas, carbone o nucleari) e persino alcune rinnovabili come l'idroelettrico, consumano grandi quantità d'acqua per il loro funzionamento. Il **machine learning**, quindi, ha un'impronta idrica duplice: quella diretta per il raffreddamento dei chip in loco e quella indiretta per la generazione dell'energia che li alimenta.

Addesramento vs Inferenza: due tipi di sete

È importante distinguere tra due fasi del ciclo di vita di un'AI:

- **Addestramento (Training):** È la fase in cui il modello “impara” leggendo miliardi di testi. Questa fase è estremamente intensiva. Si stima che l’addestramento di un modello di punta (come GPT-4 o i suoi successori del 2026) possa consumare tanta acqua quanto ne servirebbe per riempire diverse piscine olimpioniche. Tuttavia, questo avviene una tantum (o periodicamente).
- **Inferenza (Utilizzo):** È l’uso quotidiano. Sebbene il consumo per singola domanda sia basso (il famoso “sorso”), il volume totale delle richieste globali rende l’inferenza la fase più idrovora nel lungo termine. Ogni volta che chiediamo all’AI di riassumere una mail o generare un’immagine, attiviamo questo processo.

La generazione di immagini, in particolare, è molto più costosa in termini energetici e idrici rispetto al testo. Creare un’immagine complessa in alta definizione può consumare tanta energia quanto ricaricare completamente uno smartphone, con il conseguente consumo d’acqua associato al raffreddamento di quello sforzo computazionale improvviso.

Verso un’AI idricamente sostenibile?

La consapevolezza di questo problema ha spinto la ricerca verso nuove soluzioni. Nel 2026, stiamo assistendo all’ascesa di nuove architetture hardware. I chip neuromorfici e le unità di elaborazione tensoriale (TPU) di nuova generazione sono progettati per massimizzare l’efficienza energetica. Inoltre, si sta diffondendo il **raffreddamento a liquido diretto** (direct-to-chip liquid cooling), dove un fluido dielettrico circola direttamente a contatto con i componenti, trasportando il calore in modo molto più efficiente dell’aria e riducendo la necessità di evaporazione dell’acqua.

Anche lato software, gli sviluppatori stanno lavorando su “Small Language Models” (SLM), modelli più compatti e specializzati che richiedono meno potenza di calcolo (e quindi meno acqua) per funzionare, pur mantenendo prestazioni elevate per compiti specifici. L’obiettivo è ridurre il costo del “sorso” senza sacrificare l’intelligenza del sistema.

Conclusioni

Il “sorso fantasma” non deve demonizzare l’uso dell’Intelligenza Artificiale, che rimane uno strumento fondamentale per il progresso scientifico e sociale. Tuttavia, svela una verità fondamentale: non esiste tecnologia senza materia. Ogni byte ha un peso, ogni algoritmo ha una temperatura. Essere consapevoli che ogni nostra interazione digitale ha un corrispettivo fisico in litri d’acqua è il primo passo verso un utilizzo più responsabile. Nel futuro immediato, la sfida per i giganti della tecnologia non sarà solo creare modelli più intelligenti, ma progettare infrastrutture che non prosciughino le risorse del pianeta che cercano di ottimizzare.

Domande frequenti

Quanta acqua consuma una conversazione con l’intelligenza artificiale?

Una conversazione standard con un modello generativo, che comprende tra i 20 e i 50 scambi di battute, comporta il consumo di circa 500 millilitri di acqua dolce. Questo fenomeno, noto come sorso fantasma, è dovuto alla necessità di raffreddare i server dei data center che si surriscaldano elaborando le richieste degli utenti.

Perché i data center dell'AI hanno bisogno di sistemi di raffreddamento ad acqua?

I processori e le GPU utilizzati per l'intelligenza artificiale generano un calore estremo quando lavorano a pieno regime. Per mantenere la temperatura operativa ed evitare guasti, molte strutture utilizzano il raffreddamento evaporativo, un processo che consuma acqua per sottrarre calore all'ambiente, specialmente nelle zone climatiche più calde.

Generare immagini con l'AI consuma più risorse rispetto al testo?

Sì, la creazione di immagini complesse in alta definizione richiede una potenza di calcolo molto superiore rispetto alla generazione di testo. Si stima che produrre una singola immagine possa consumare tanta energia quanto ricaricare completamente uno smartphone, con un conseguente aumento del fabbisogno idrico per il raffreddamento dell'hardware.

Qual è la differenza di consumo tra addestramento e inferenza nell'AI?

L'addestramento è una fase iniziale estremamente intensiva che consuma enormi quantità d'acqua una tantum per istruire il modello. L'inferenza, invece, è l'utilizzo quotidiano da parte degli utenti: sebbene il costo per singola domanda sia basso, il volume globale di richieste rende questa fase la più impattante sul consumo idrico a lungo termine.

Quali soluzioni esistono per ridurre l'impatto idrico dell'intelligenza artificiale?

Il settore si sta muovendo verso l'adozione di chip neuromorfici più efficienti e sistemi di raffreddamento a liquido diretto che riducono l'evaporazione. Inoltre, lo sviluppo di Small Language Models permette di utilizzare intelligenze artificiali più compatte e specializzate, che richiedono meno potenza di calcolo e quindi meno acqua per funzionare.