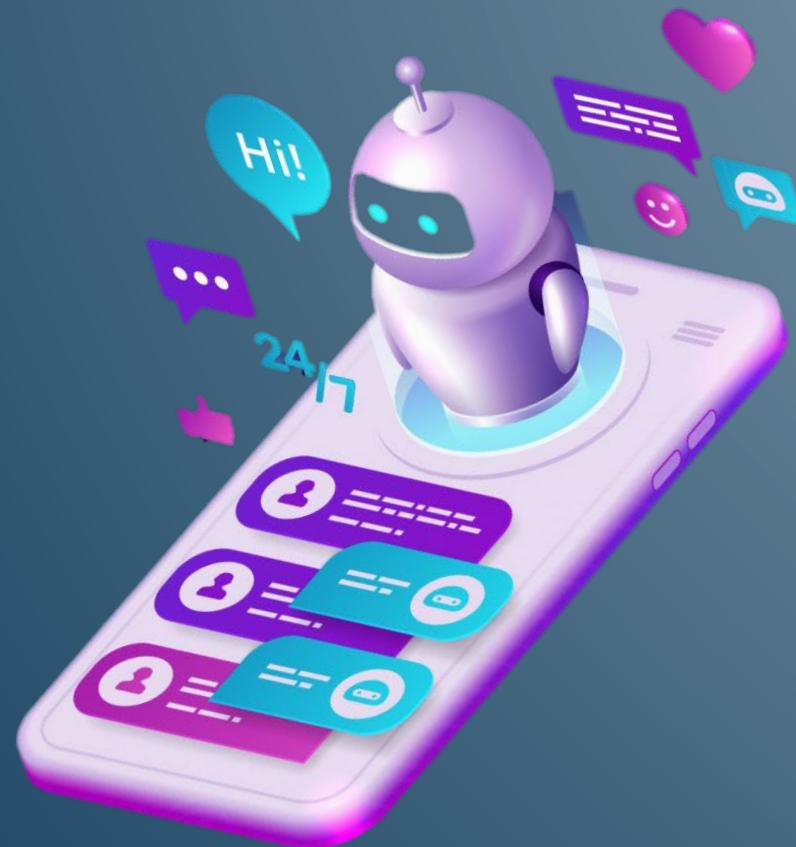


# ChatGPT

Esplorando il potere dell'Intelligenza Conversazionale

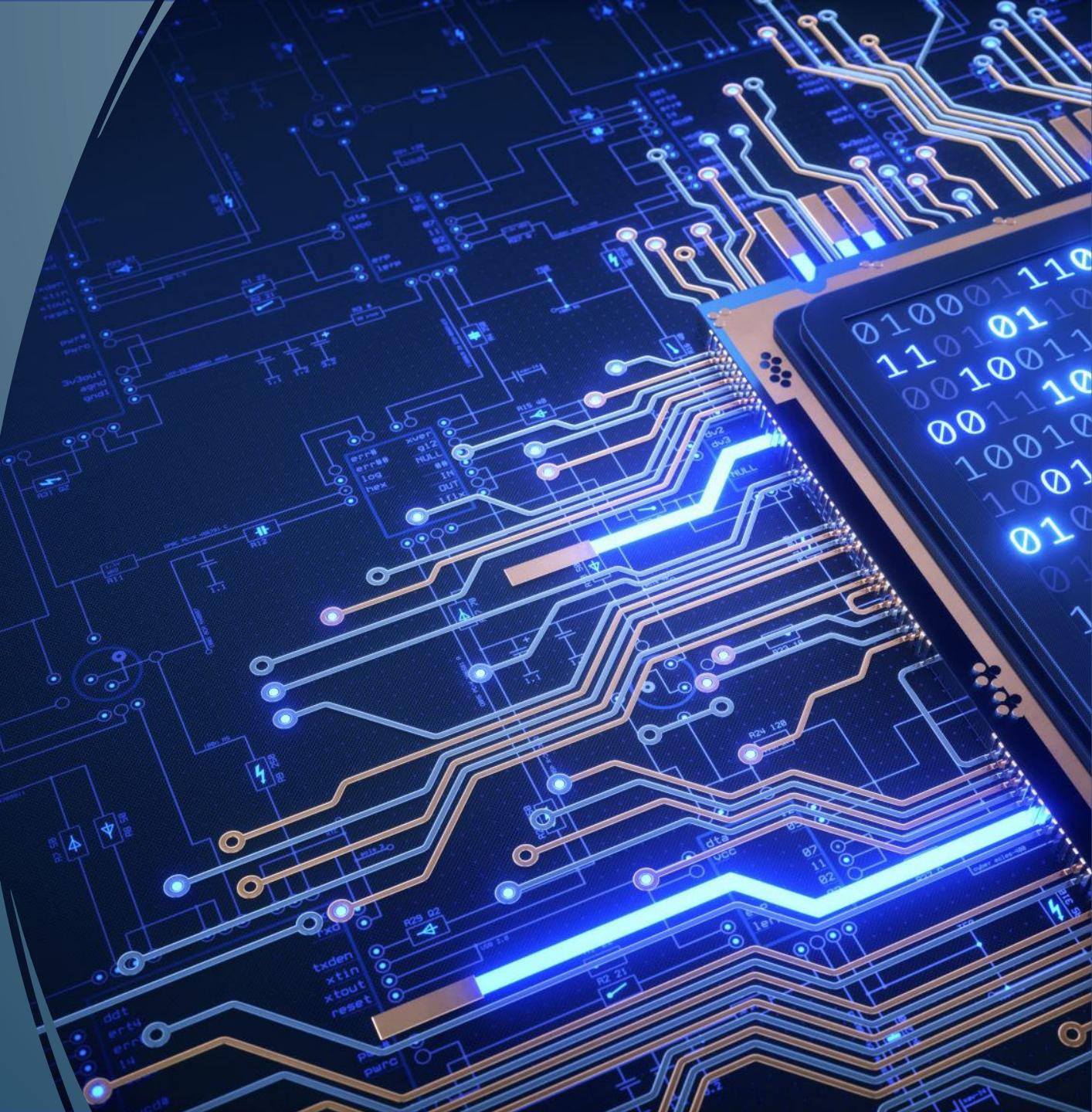


Relatori :  
Mattias Caliandro, Mirco Caputo

# Cos'è ChatGPT?

ChatGPT è un chatbot basato su intelligenza artificiale e apprendimento automatico sviluppato da OpenAI specializzato nella conversazione con un utente umano.

La sigla GPT sta per Generative Pre-trained Transformer, una tecnologia nuova applicata al machine learning.



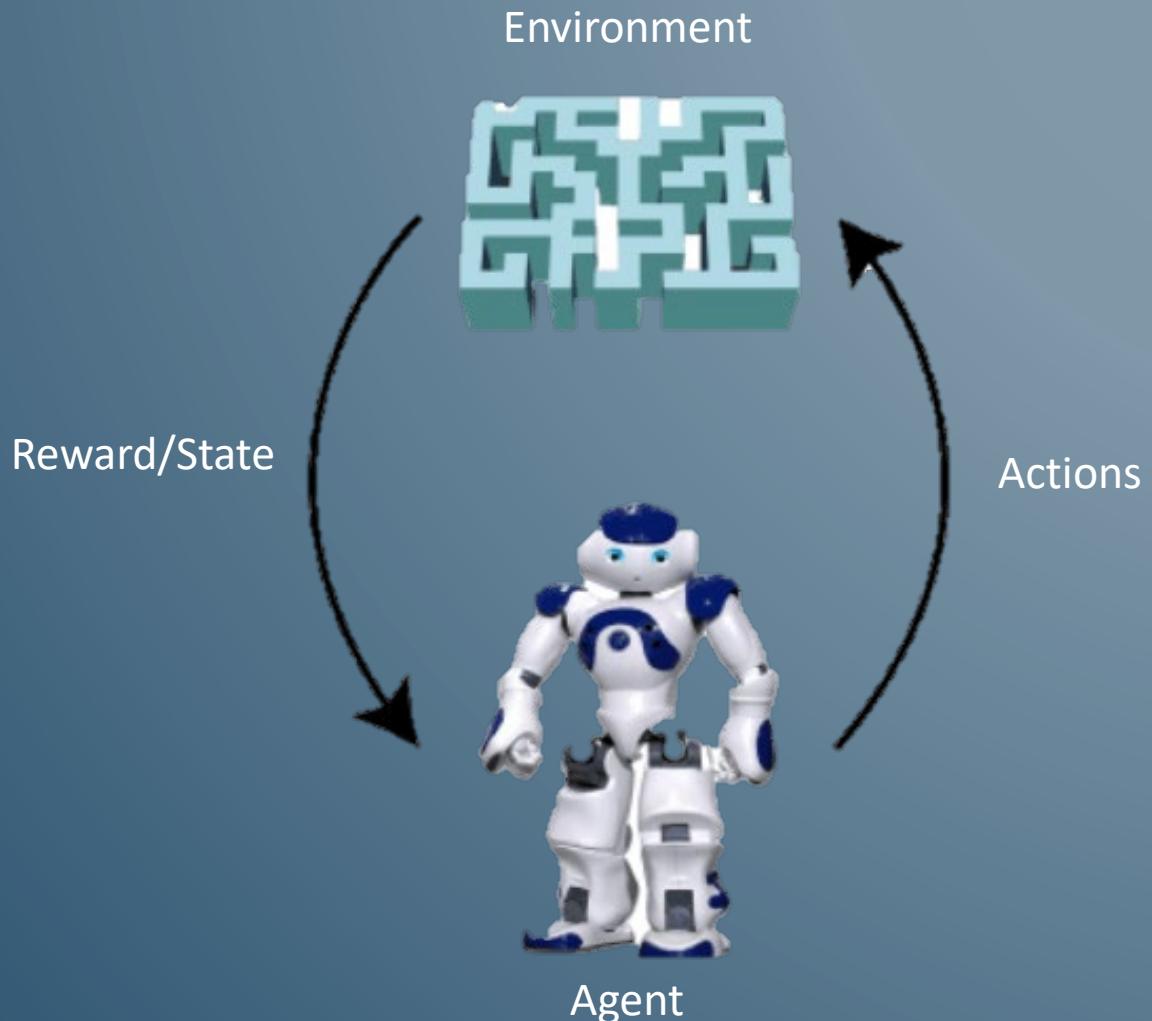
# Apprendimento Supervisionato

L'apprendimento supervisionato è un processo in cui un algoritmo apprende da un set di dati di addestramento che contiene esempi etichettati.



# Apprendimento per rinforzo

L'apprendimento per rinforzo è un processo in cui un agente apprende a prendere decisioni ottimali in un ambiente attraverso l'interazione con esso.



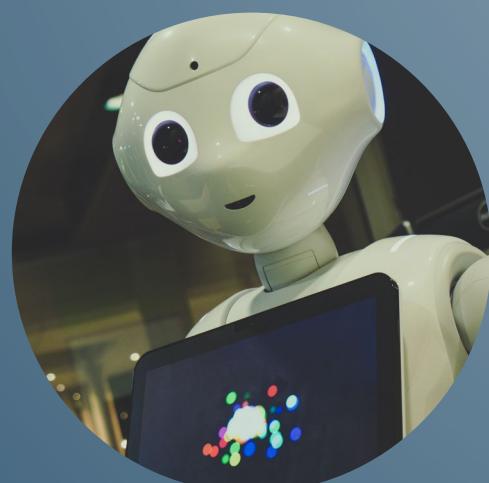
# L'evoluzione di GPT



GPT-1



GPT-2



GPT-3



GPT-3.5

# GPT-1

Un Generative Pre-Trained Transformer (GPT) è un modello linguistico che si basa sul Deep Learning in grado di generare testi di tipo umano sulla base di un dato input basato su testo. Un utente ‘alimenta’ il modello con una frase e il trasformatore crea informazioni coerenti basate su paragrafi estratte da set di dati disponibili pubblicamente.

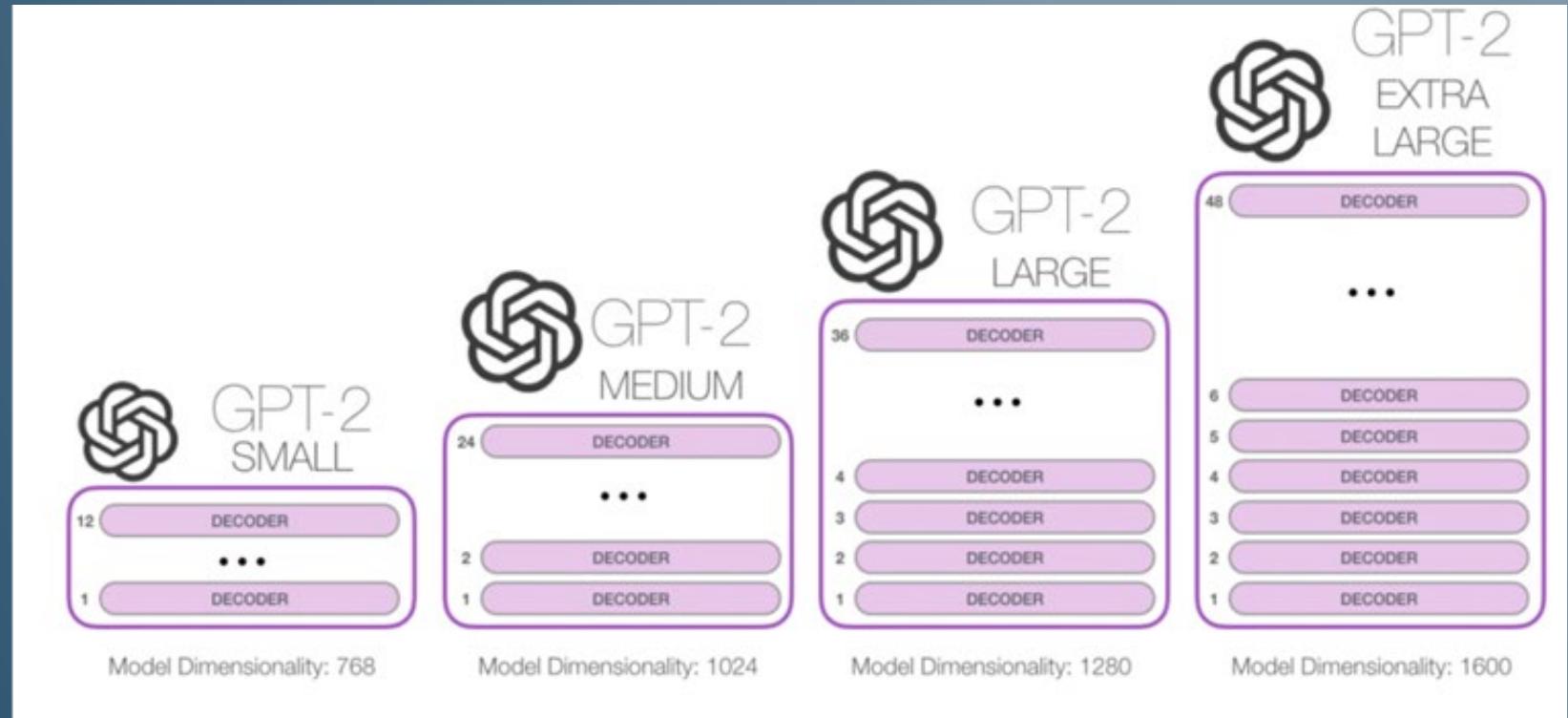


# GPT-2

Annunciato nel 2019 da OpenAI, GPT-2 è il successore di GPT. L'obiettivo è : prevedere la parola successiva conoscendo tutte le parole precedenti all'interno di un testo. L'eterogeneità del dataset di training implica che il modello generato si possa adattare a diversi domini applicativi.



# GPT-2

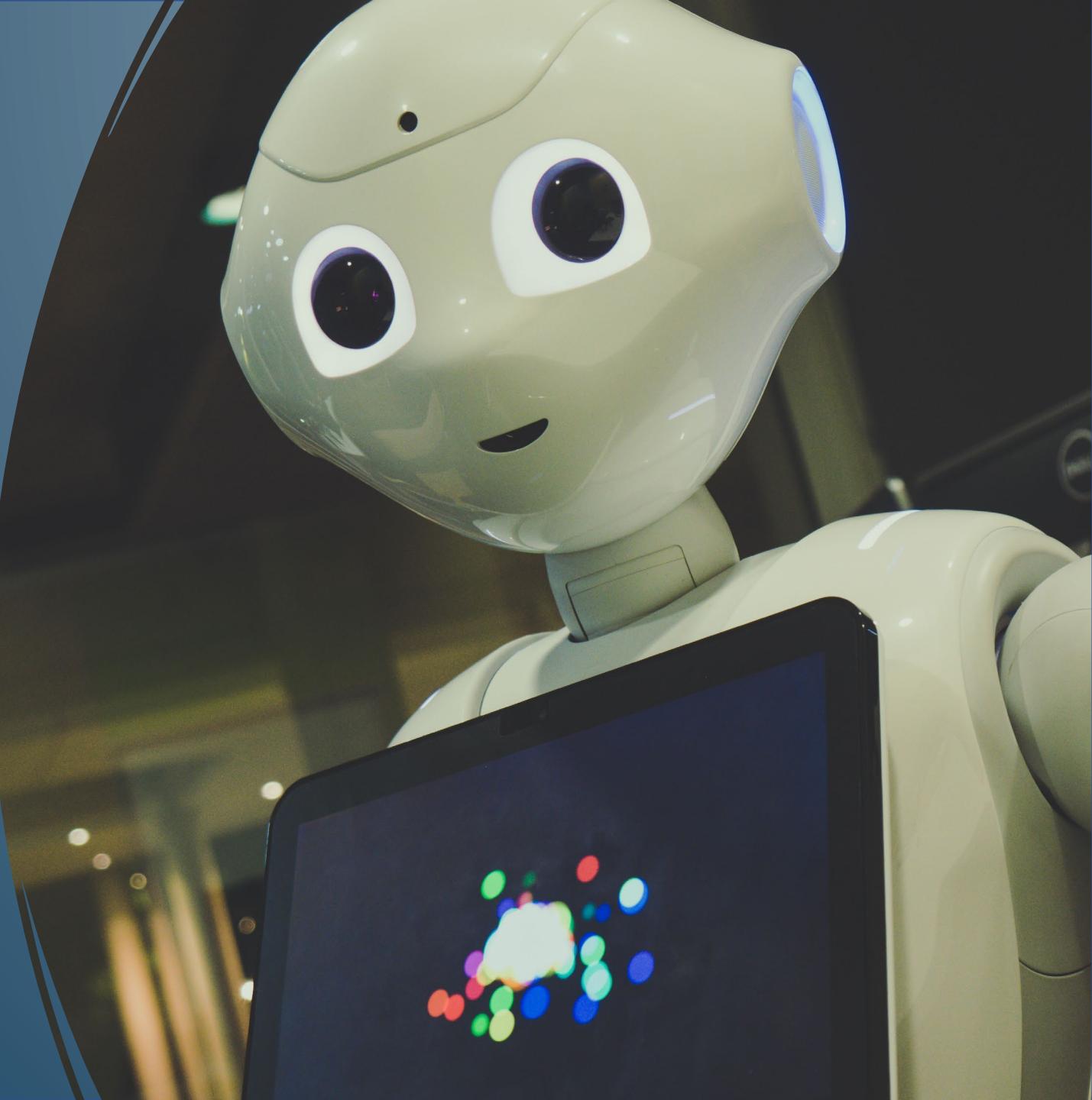


Sono stati rilasciati comunque diversi modelli di GPT2 come mostrato in figura.  
Questi differiscono a seconda della dimensione:

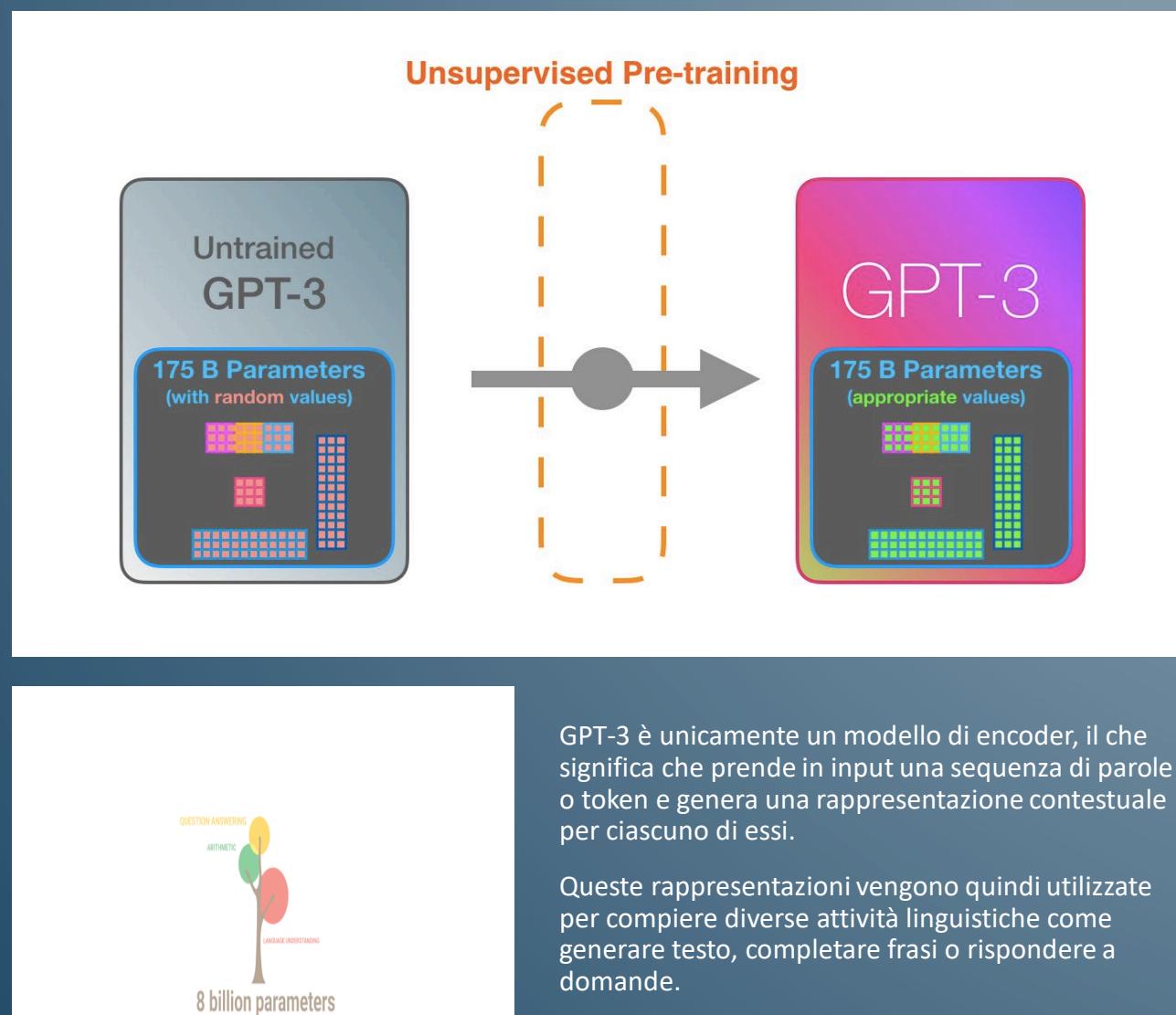
- small (124M parametri)
- medium (355M parametri)
- large (774M parametri)
- extra large (1.5BM parametri).

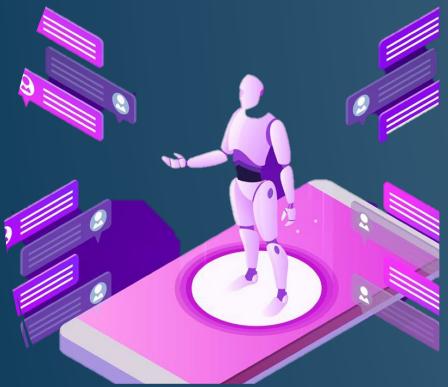
# GPT-3

GPT-3 del 2020 era 100 volte più grande del suo predecessore, mentre i dati del testo di addestramento erano 10 volte più grandi. Il modello ha imparato a tradurre da altre lingue, eseguire operazioni aritmetiche, eseguire semplici programmazioni, ragionare in sequenza e molto altro come risultato dell'espansione della quantità che ha bruscamente aumentato la qualità.



# GPT-3

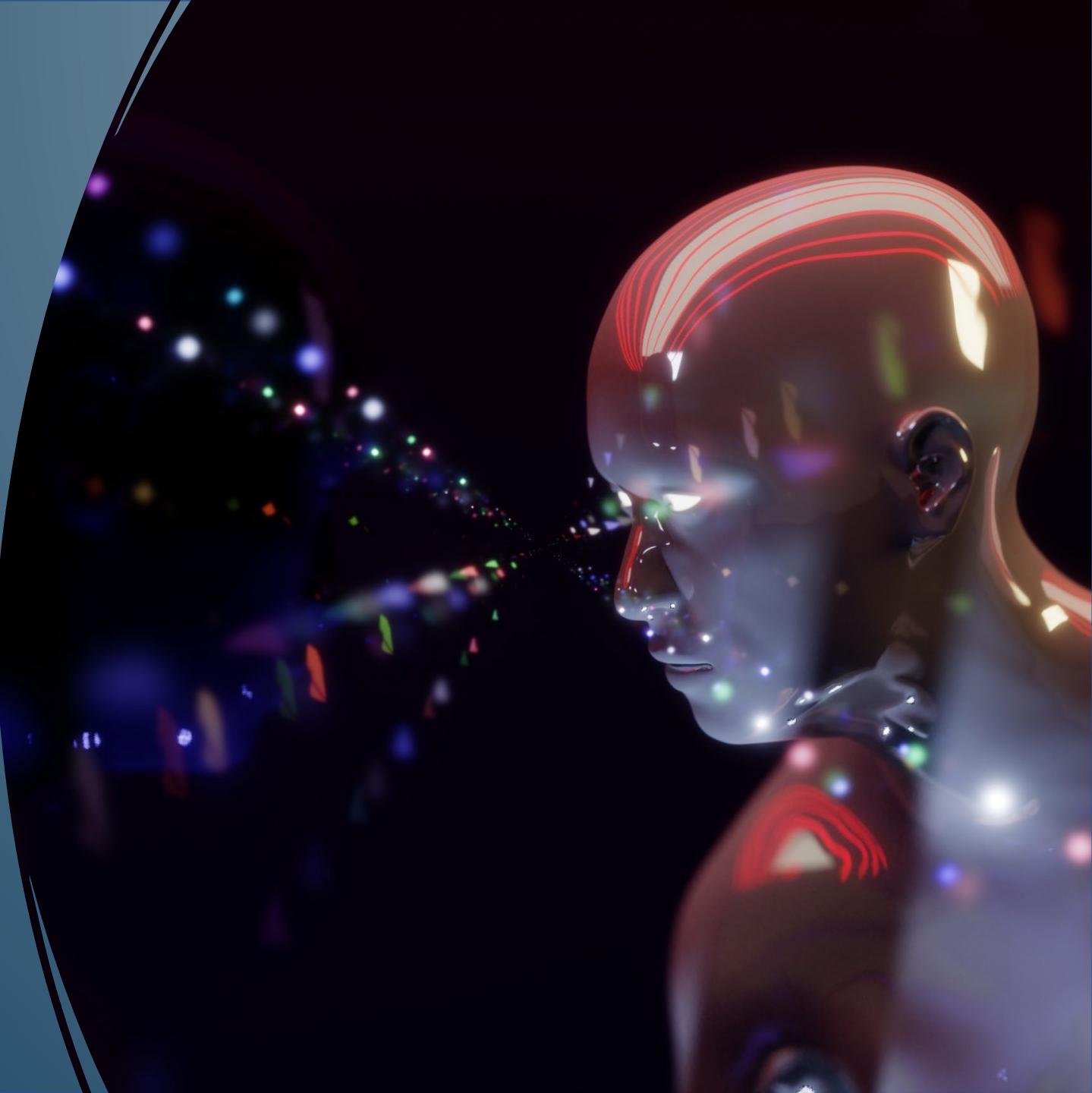




# GPT-3.5 (InstructGPT)

InstructGPT (o Codex) è stato addestrato su una vasta quantità di testo di programmazione proveniente da diverse fonti, consentendogli di comprendere il contesto e generare codice Python funzionante.

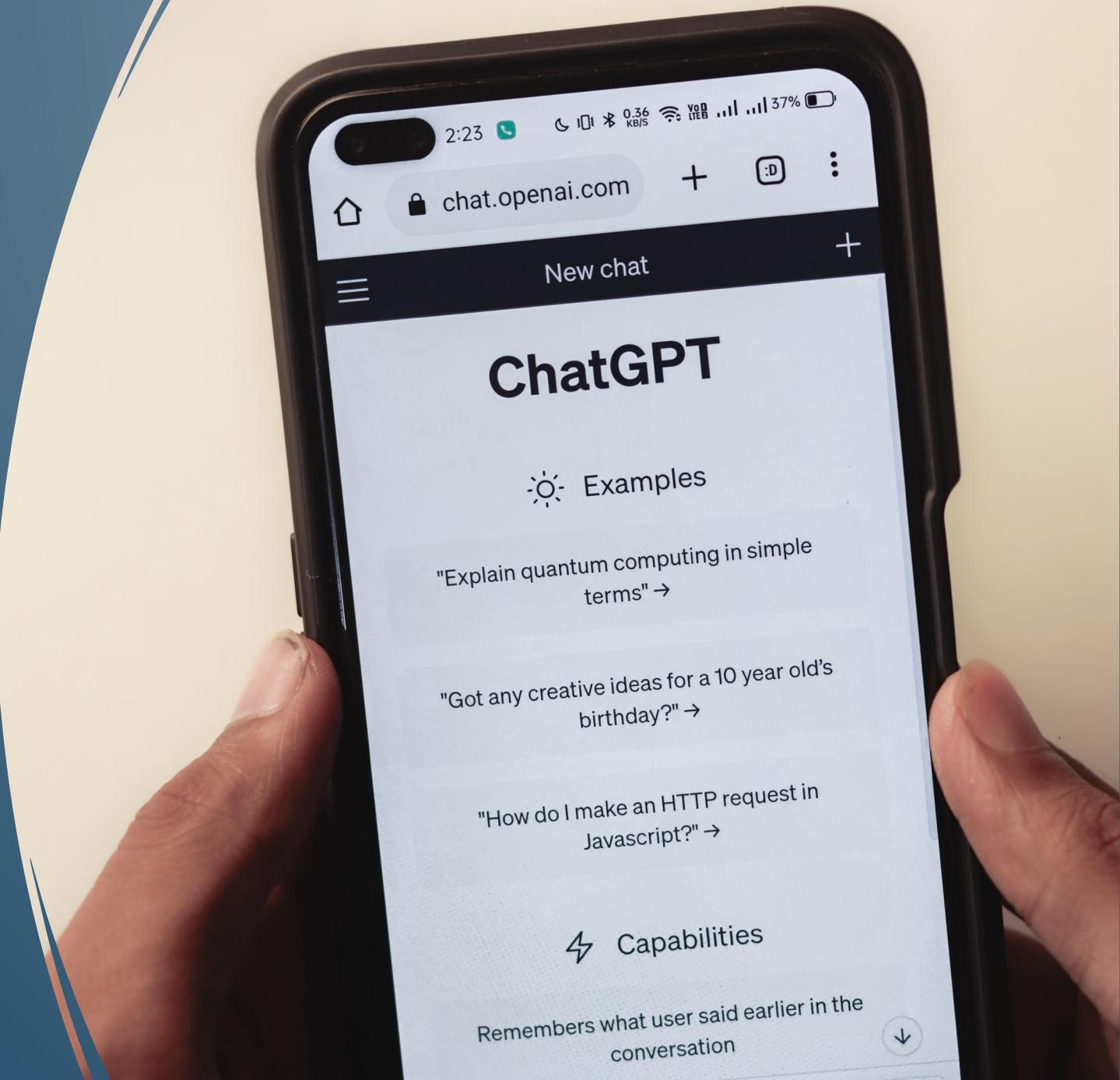
Può essere utilizzato per assistere gli sviluppatori nella scrittura di codice, fornendo suggerimenti, completamenti automatici e persino scrivendo intere funzioni o blocchi di codice.





# ChatGPT

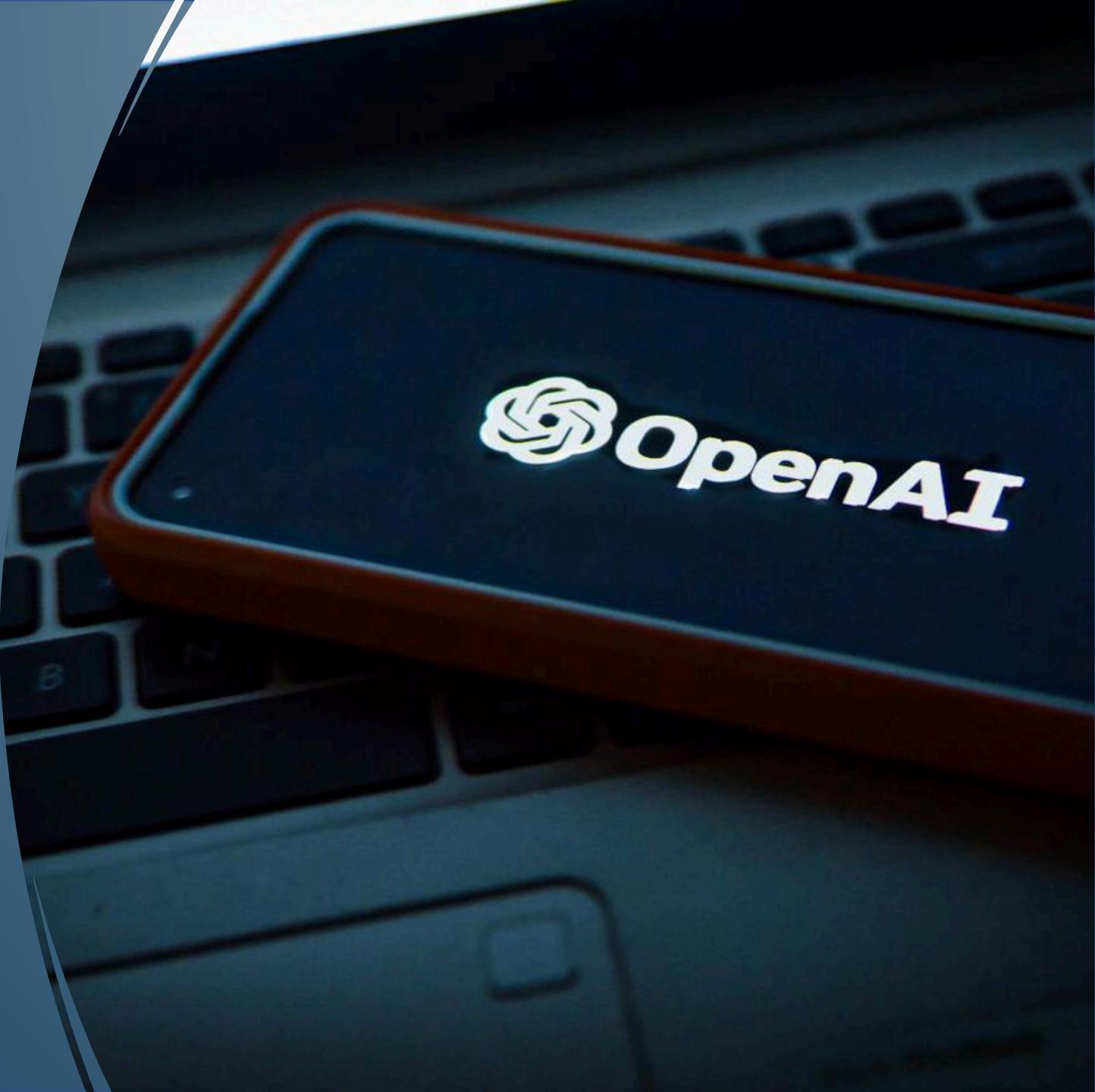
La chiave del successo di ChatGPT è la nuova interfaccia user-friendly. ChatGPT utilizzando la "finestra di dialogo" dei messenger è diventato uno strumento disponibile per tutti contemporaneamente.



# ChatGPT

Per generare testi con ChatGPT, è essenziale fornire al modello un "prompt". Il prompt è l'input testuale che viene utilizzato per generare il testo di risposta. Può essere qualsiasi cosa, come una frase di avvio per scrivere un articolo o una storia, oppure una domanda per la quale si desidera ottenere una risposta. In generale, maggiore è la precisione e la specificità del prompt, maggiore sarà la precisione e la specificità del testo generato da ChatGPT.

I prompt degli utenti vengono presentati al valutatore di prompt di ChatGPT. Se il valutatore risponde con "no", viene restituito un messaggio di errore all'utente. Se il valutatore risponde con "sì", il prompt viene passato a ChatGPT per generare la risposta desiderata.



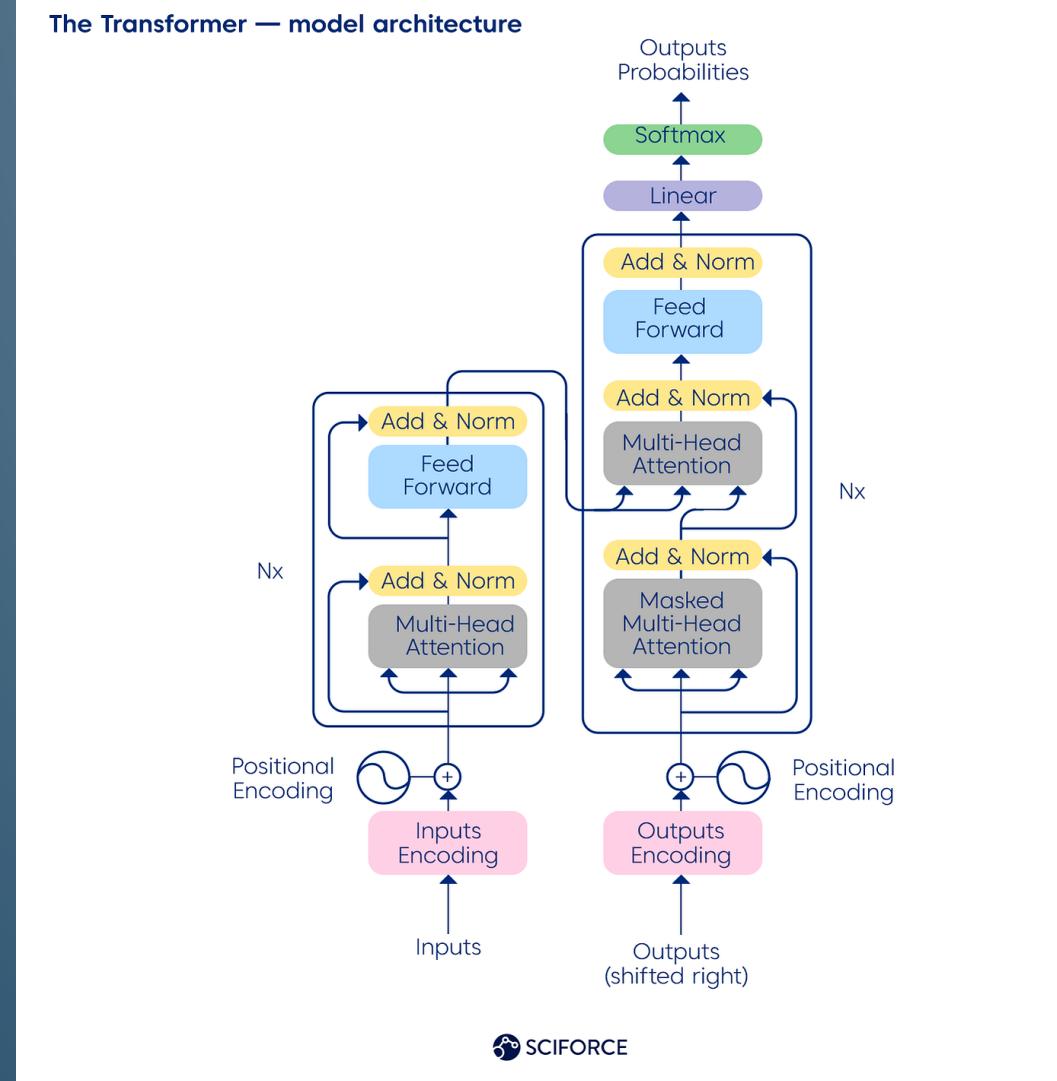
# Creazione di una ChatBot

Grazie al framework Gradio e le API di OPENAI  
abbiamo creato una ChatBot in codice Python.

In questo caso abbiamo creato uno Psicologo  
Digitale



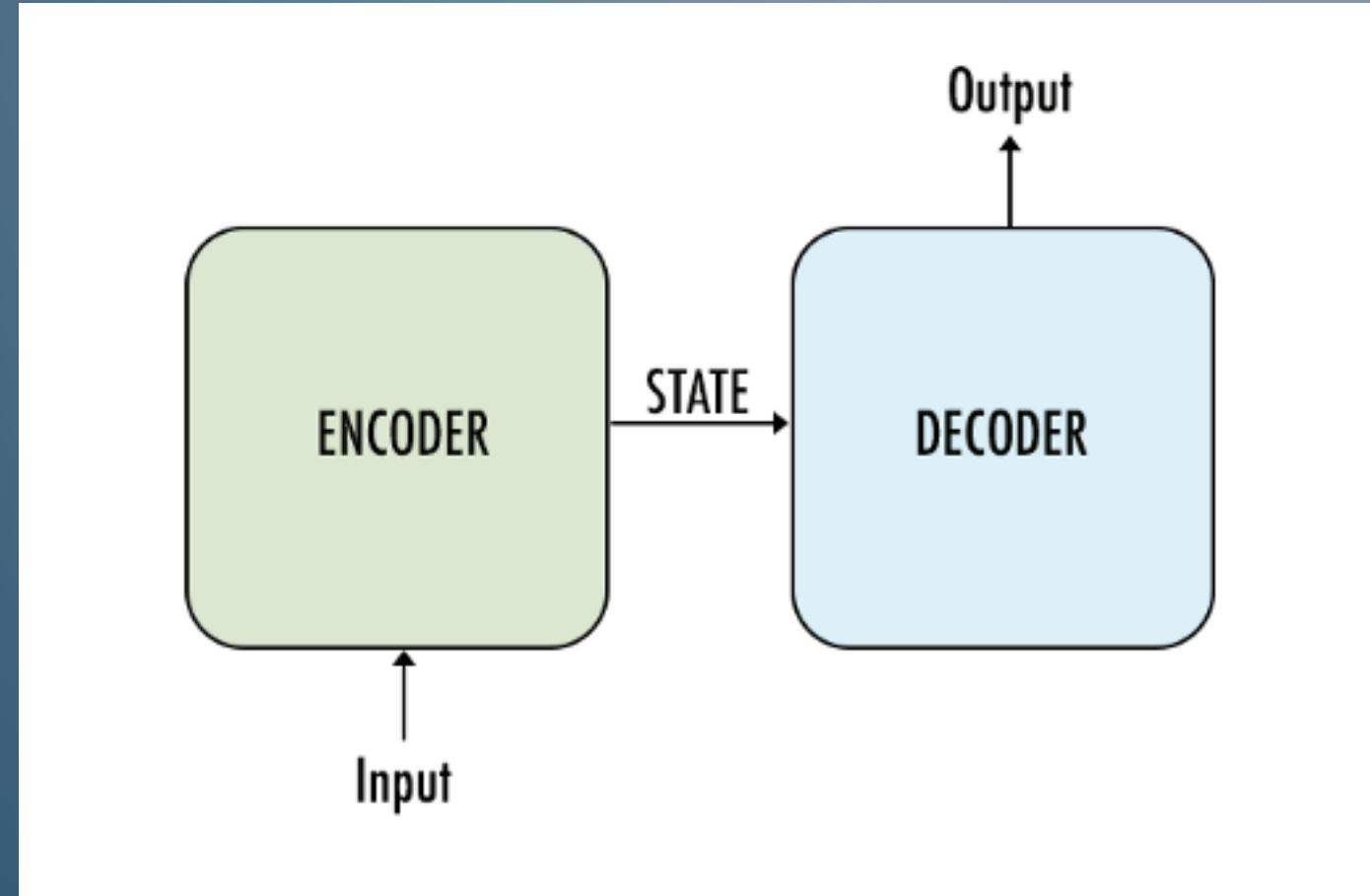
# GPT Generative Pre-trained Transformer



# Encoder & Decoder

Gli **encoders** hanno il compito di trasformare una frase di input in una rappresentazione numerica, tenendo conto del contesto e del significato delle parole e delle frasi. Utilizzano varie tecniche, come l'embedding, l'attenzione multi-head e le reti feed-forward, per analizzare l'input e creare una rappresentazione numerica della frase.

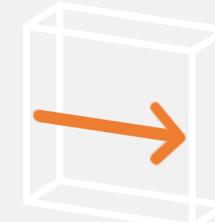
D'altra parte, i **decoders** sono responsabili di generare la risposta o la frase successiva utilizzando la rappresentazione numerica prodotta dagli encoders. I decoders utilizzano il meccanismo di auto-attenzione (self-attention) per generare una distribuzione di probabilità per ogni parola del vocabolario, indicando la probabilità che ogni parola sia la prossima parola nella risposta o nella frase in fase di generazione.



# Transformer

I transformer sono modelli di apprendimento profondo basati su reti neurali, che hanno rivoluzionato il campo del Natural Language Processing grazie alla loro capacità di gestire sequenze di parole in modo più efficiente rispetto ai modelli precedenti.

Positional encoding

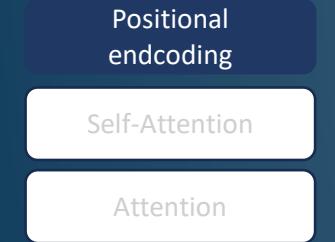


Self-Attention



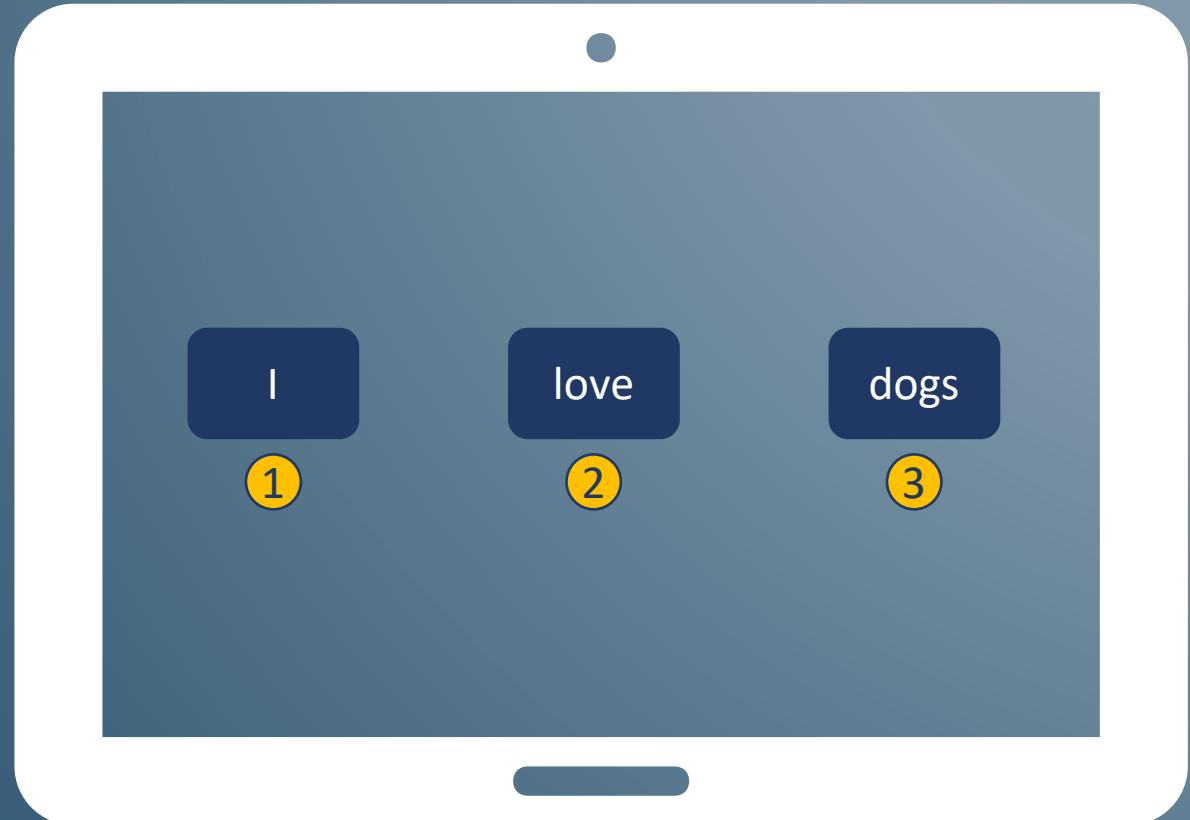
Attention

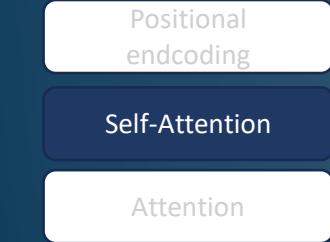




# Positional Encoding

Prima di immettere le parole all'interno della neural network ogni parola viene numerata in base alla sua posizione nella frase. Così facendo quando la network viene 'trainata' su una grande quantità di testo impara come interpretare questi positional encodings capendo l'importanza dell'ordine delle parole dai dati.





# Self-Attention

La self-attention permette alla rete neurale di capire il contesto di una parola tra le parole che la circondano.

In generale aiuta la rete neurale a togliere ambiguità alle parole, riconoscere parti di frasi o capire il tempo delle parole(passato/presente..)

Server, can i have  
the **check**?

Looks like i just  
**crashed** the server.

La parola **check** da un contesto alla parola **server** togliendone ambiguità (mail server o human server)

La parola **crashed** fa intendere che il **server** sia una macchina.



# Attention

Il ‘meccanismo’ dell’attenzione è una struttura di rete neurale che permette a un modello testuale di guardare ogni singola parola della frase originale per prendere decisioni sul come tradurre la parola nell’output. Il modello per ogni parola in input si ‘aspetta’ una corrispondente in output. Ma come fa ad aspettarsi qualcosa in output? questo perché dopo che viene ‘trainato’ con tanti esempi di traduzioni il modello impara i modelli grammaticali dietro di esse.

The European Economics Area



la europea economia zona

The European Economics Area



lo spazio economico europeo

# Reti neurali artificiali

Con “rete neurale artificiale” nel settore dell'apprendimento automatico ci si riferisce a un modello matematico che punta a somigliare – nel senso più ampio del termine – alle reti neurali biologiche, presenti nell'essere umano o negli animali, e che è costituito da neuroni artificiali costruiti virtualmente o, talvolta, fisicamente.

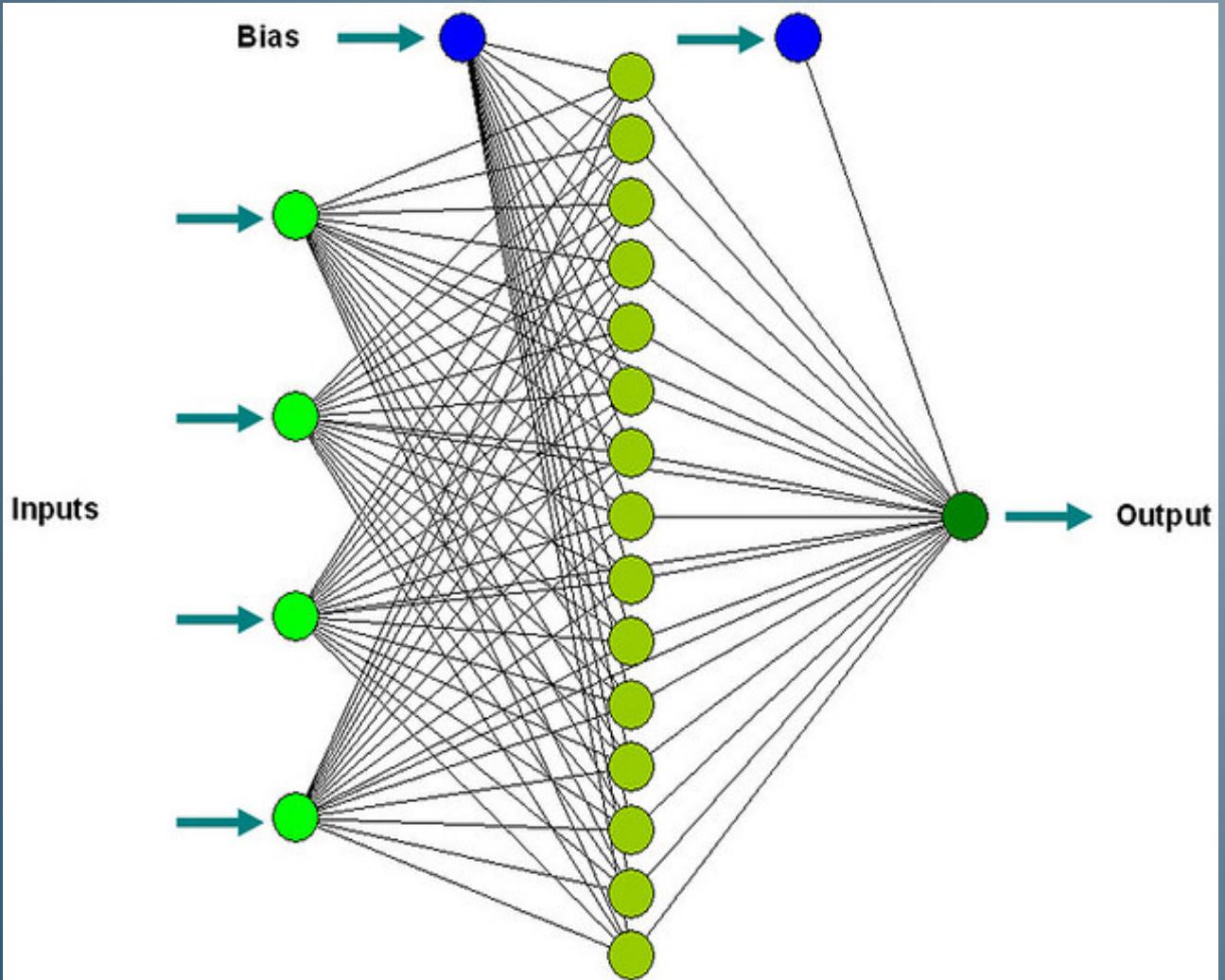


# Come funzionano?

Data la struttura di una rete neurale biologica, si è tentato di riprodurre, in modo semplificato, un sistema simile che potesse funzionare altrettanto bene per l'elaborazione di informazioni. La rete neurale artificiale è composta da nodi o neuroni formali che sono delle unità computazionali fondamentali e che, essendo interconnessi, formano un grafo composto almeno da uno strato di ingresso e uno strato di uscita. Anche questa rete ha come obiettivo quello di elaborare dati e informazioni e funziona con un input (che possiamo paragonare all'impulso elettrico per il neurone biologico) che arriva al nodo dello strato d'ingresso.



# Struttura rete neurale artificiale



# Ambiti applicativi



Chat-bot



Creazione di  
testi



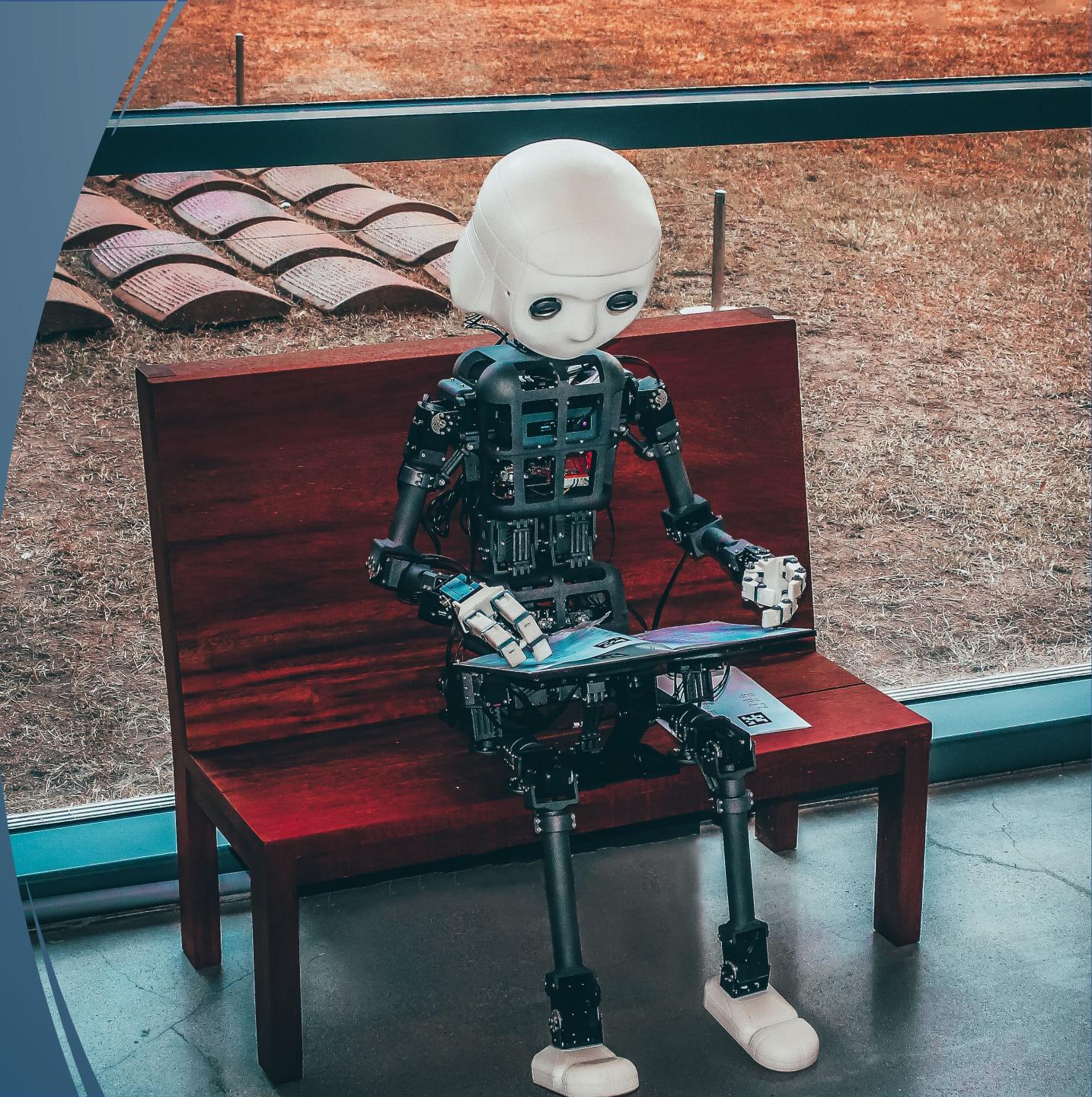
Traduzioni

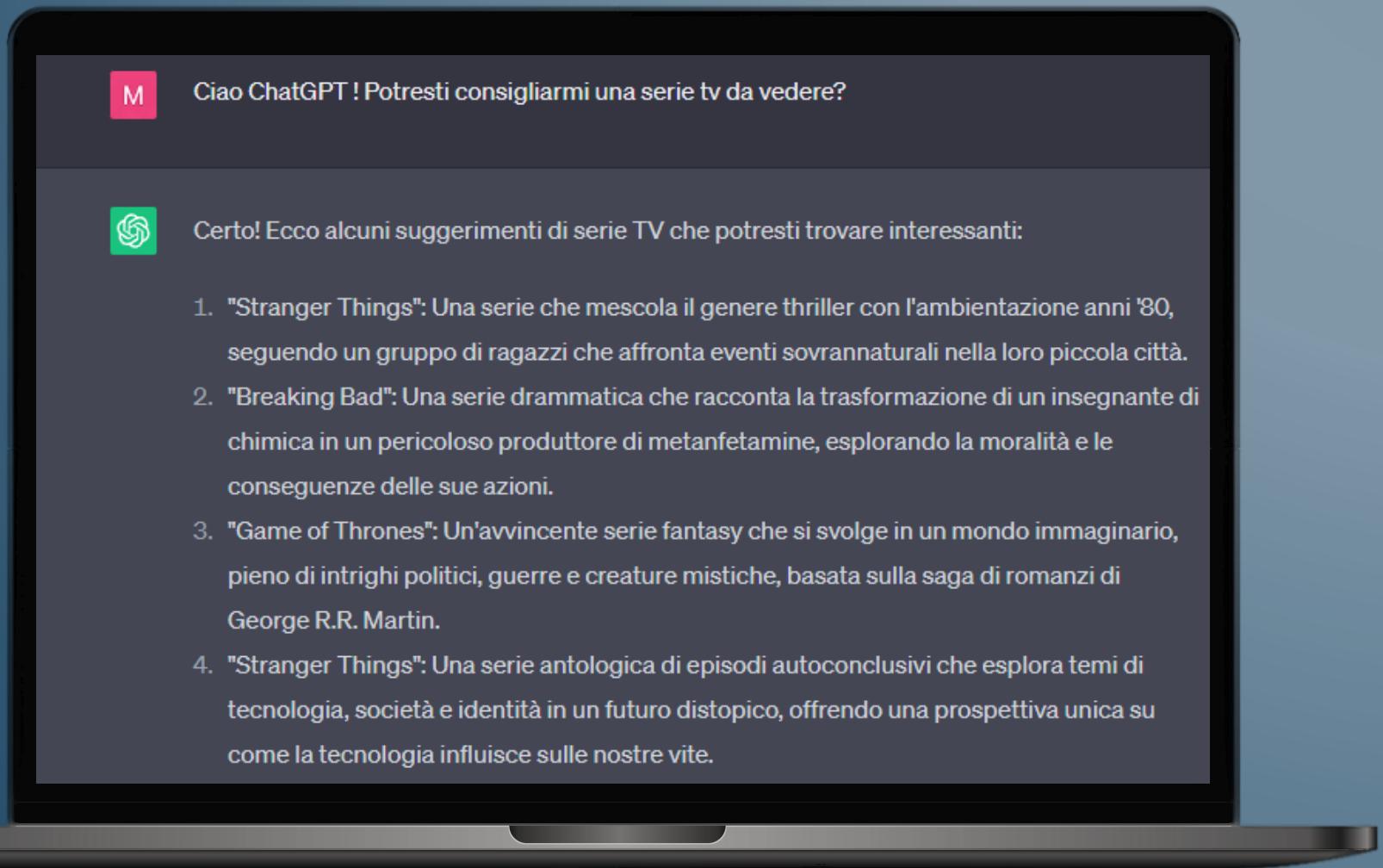


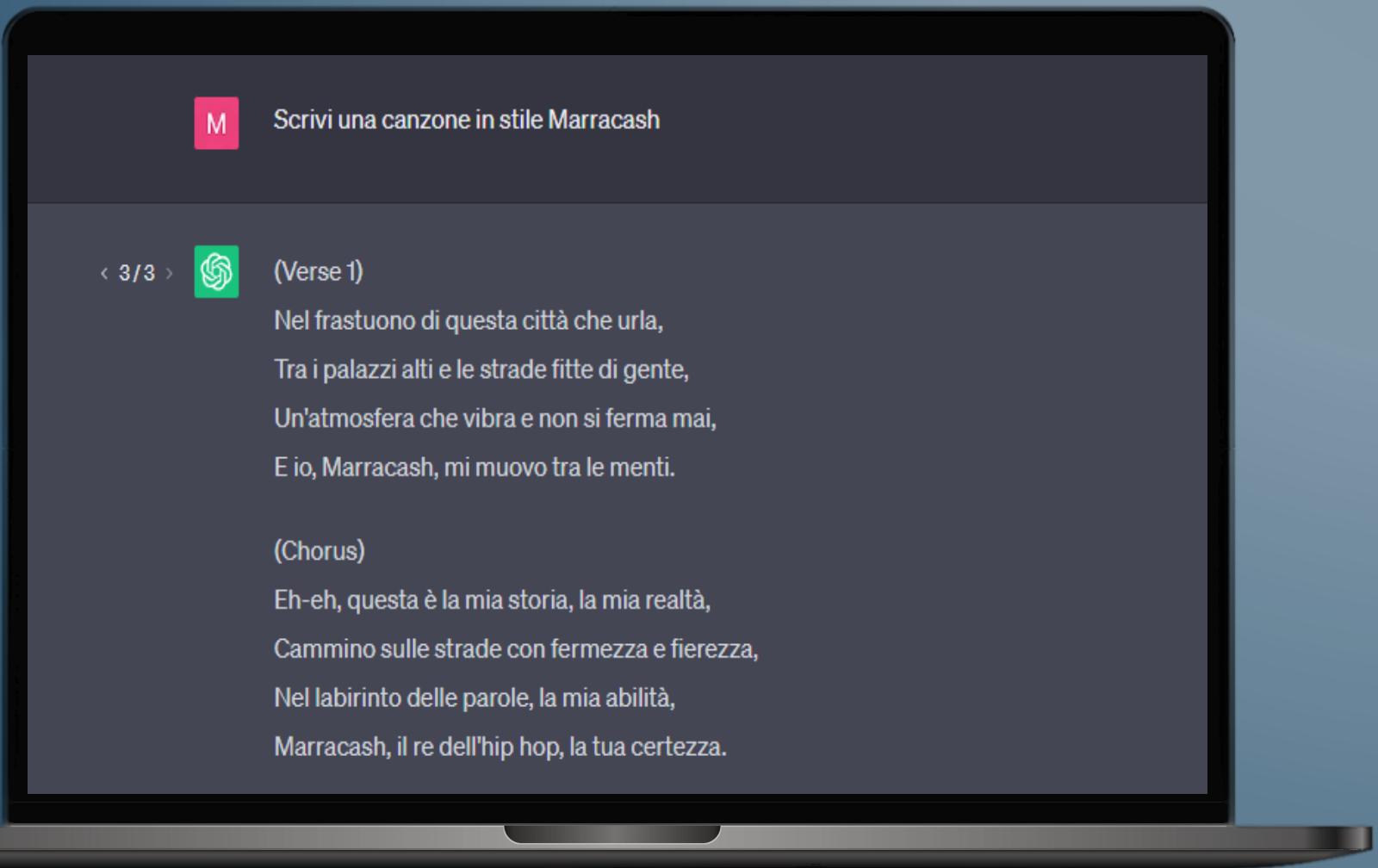
Creazione di  
codice

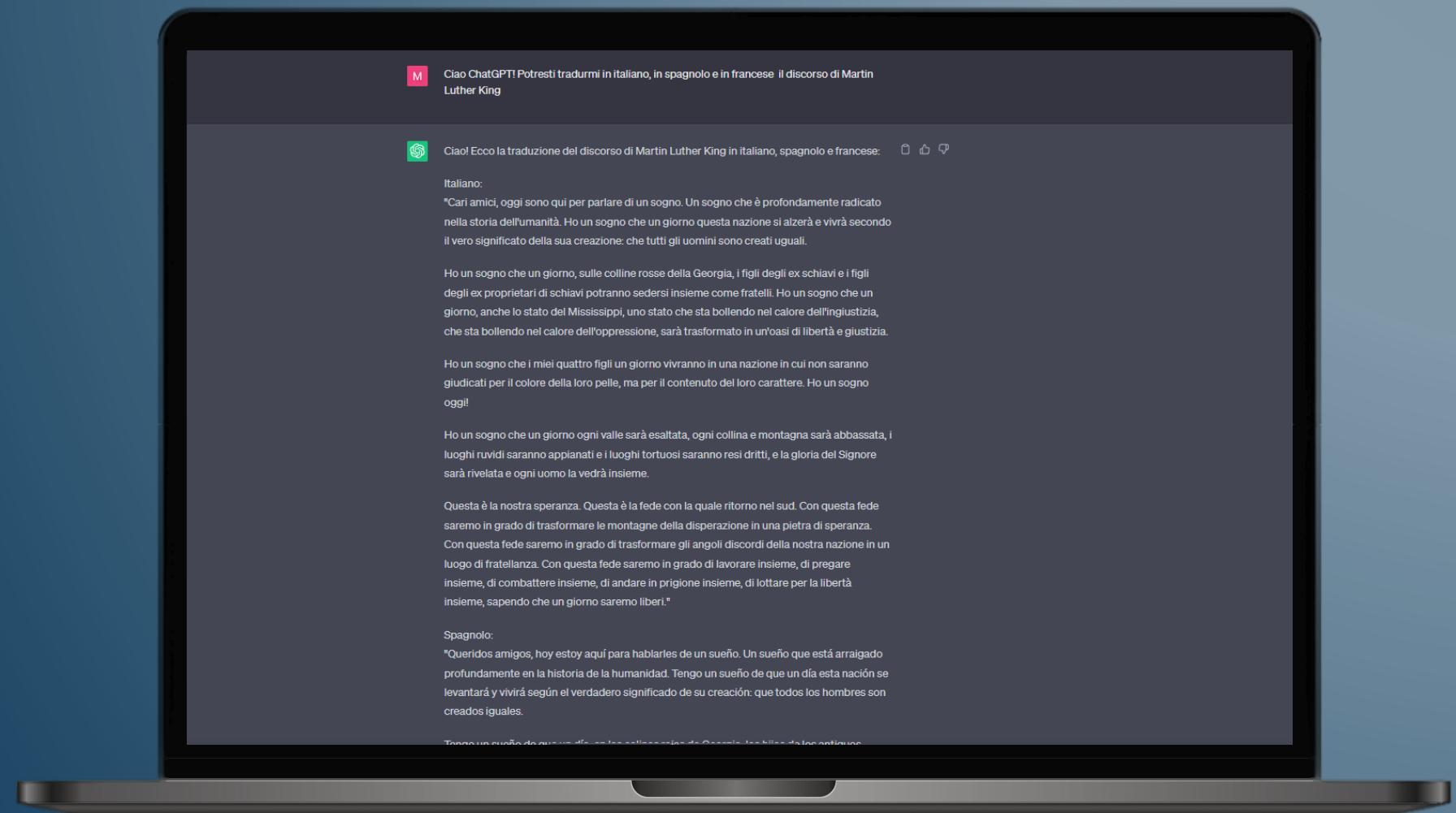


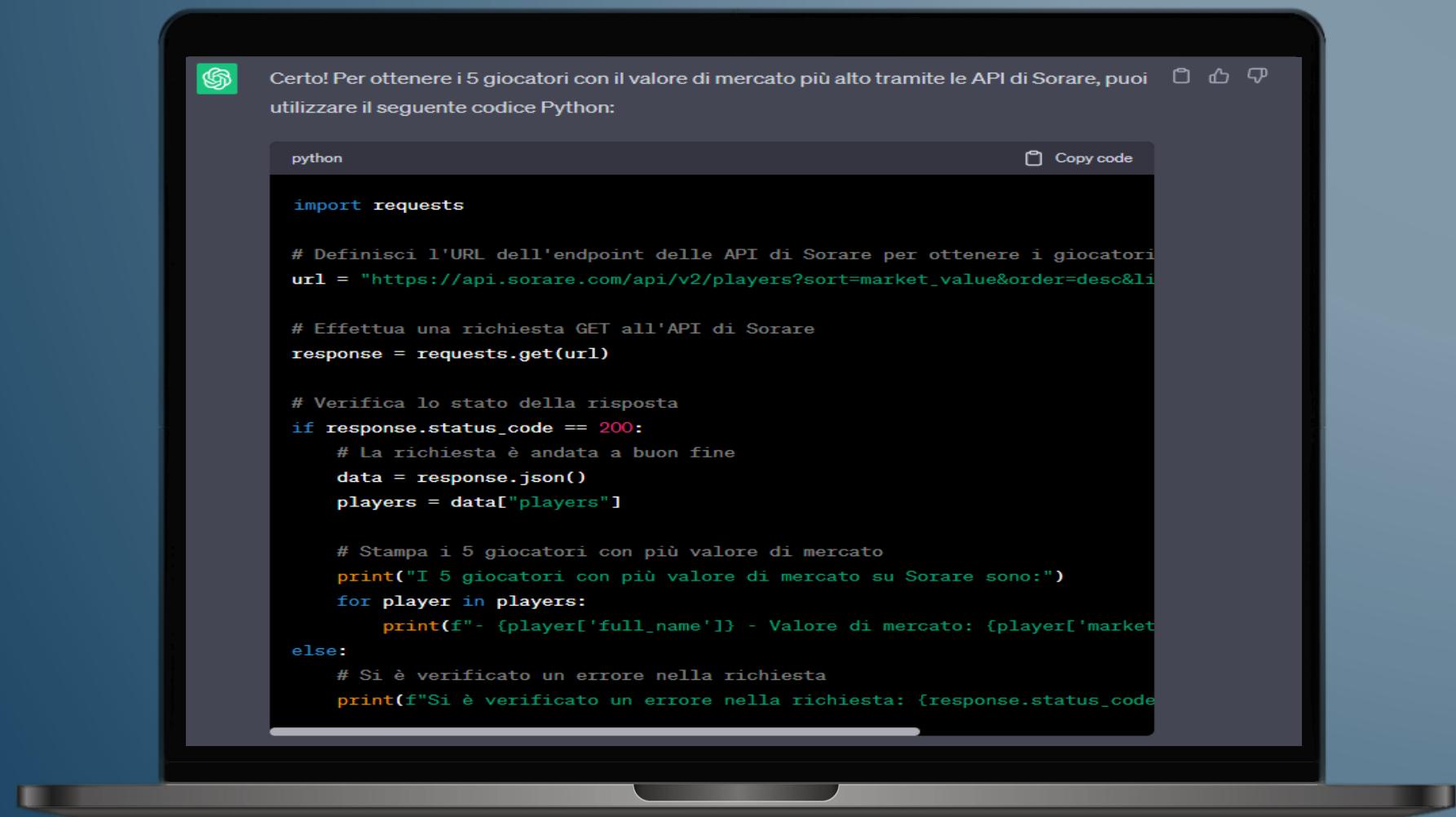
Risoluzione  
problemi

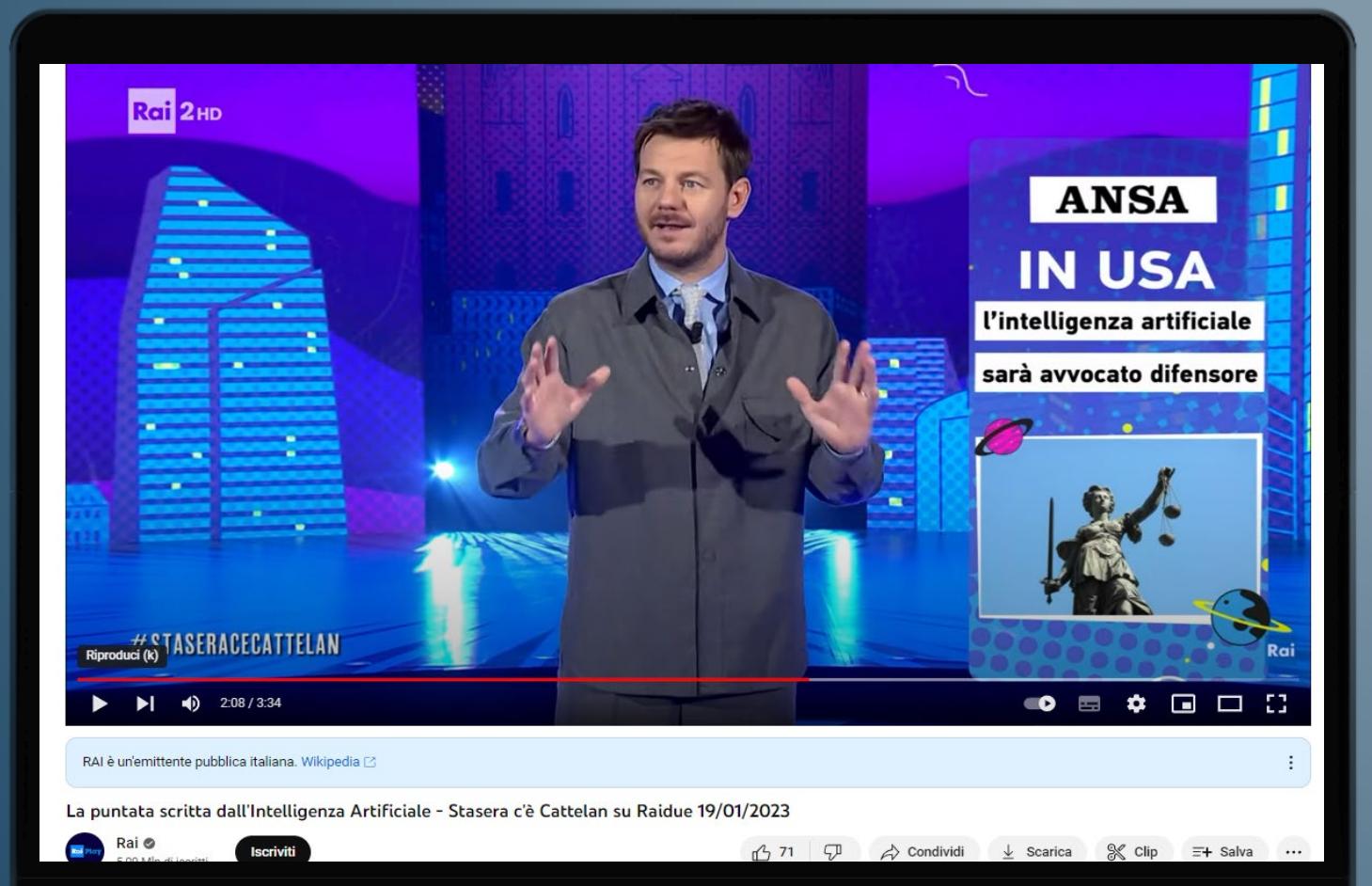












# Considerazioni etiche

*«Un altro pericolo della scrittura AI è il potenziale per un uso non etico. Gli algoritmi di intelligenza artificiale possono essere programmati per generare informazioni fuorvianti o false, che potrebbero essere utilizzate per diffondere notizie false o disinformazione. Ciò potrebbe avere gravi conseguenze, come influenzare l'opinione pubblica o causare danni a individui o gruppi».*





GRAZIE PER L'ATTENZIONE

pepper

DI :

MIRCO CAPUTO  
MATTIAS CALIANDO