LINUX DAY TORINO 25/10/25



OPEN SOURCE LLM E COME SFRUTTARLI

GIULIO SCIARAPPA



CHI SONO

GIULIO SCIARAPPA

CLOUD ENTERPRISE ARCHITECT





LLM: COSA SONO?



- TLLM = Large Language Model
- Addestrato su miliardi di token testuali con pesi differenti
- MOTORE DI COMPLETAMENTO DEL LINGUAGGIO NATURALE ((PREVEDONO LA PROSSIMA PAROLA))
- **BASATO SU TRANSFORMER E SELF-ATTENTION**
- TESEMPI: CHATGPT, COPILOT, CLAUDE...



LLM FANTASTICI E DOVE TROVARLI



- 🚹 Miliardi di parametri richiedono calcoli complessi, e geometrici
- 🐴 Le normali cpu NON sono ottimizzate per questi calcoli.
- 🐔 Si usano quindi le GPU, ottimizzate per calcoli vettoriali/geometrici
- TROBLEMA DI MEMORIA -> CARICARE MILIARDI DI PARAMETRI, DOVE? > RAM (VIDEO)



LLM FANTASTICI E DOVE TROVARLI



- TOROSSI DATACENTER E FONDI DI INVESTIMENTO SI SONO FATTI AVANTI E OSPITANO I MODELLI PIÙ RINOMATI
- TPRO:
 - TOSTO ACCESSIBILE
 - 🚹 Facilità di utilizzo
- TONTRO:
 - 🐴 Quanto scrivi è di loro proprietà
 - 🐧 Potenzialmente fanno training di nuovi modelli con I dati da te inviati



LLM LOCALI, OPEN SOURCE

- TLOCALE != OPEN SOURCE
- TUNING
- PROBLEMA-> OSPITARLI LOCALMENTE RICHIEDE POTENZA





PERCHÉ OPEN SOURCE? (E LOCALI...)

- 🚹 Autonomia: nessuna dipendenza da API o cloud
- TPRIVACY: I DATI RESTANO SUL TUO PC
- Personalizzazione: puoi modificare e fine-tuning
- TOSTO ZERO: NESSUNA SUBSCRIPTION
- TOMMUNITY DRIVEN: MIGLIORA COSTANTEMENTE



PARAMETRI ≠ BIT

PARAMETRI (WEIGHTS)

- The Sono I valori numerici che il modello apprende durante l'addestramento
- 👔 Ogni parametro rappresenta una connessione (peso) tra neuroni
- 🚹 Un modello da **7B** = 7 miliardi di *pesi*
- 🔭 Ogni peso è un numero in floating point (es. FP16 o quantizzato a 4-8 bit)

ñ

Віт

- 🚹 Unità di misura della memoria (capacità di archiviazione)
- T Determinano quanto spazio occupano I parametri
- Es. 7 miliardi di parametri × 4 byte = **~28 GB di memoria**



PARAMETRI ≠ BIT

Modello	Parametri	Precisione	Peso su disco
Phi-3 Mini	3B	4bit	~ 1,5GB
Mistral 7B	7B	8 bit	~ 7 GB
LLaMA 38B	8B	16 bit	~ 16 GB



PARAMETRI ≠ BIT

PARAMETRI (WEIGHTS)

- TO SONO I VALORI NUMERICI CHE IL MODELLO APPRENDE DURANTE L'ADDESTRAMENTO
- 🌓 Ogni parametro rappresenta una connessione (peso) tra neuroni
- 🚹 Un modello da **7B** = 7 miliardi di *pesi*
- 🚹 Ogni peso è un numero in floating point (es. FP16 o quantizzato a 4-8 bit)

Віт

- 🚹 Unità di misura della memoria (capacità di archiviazione)
- Determinano quanto spazio occupano i parametri
- TES. 7 MILIARDI DI PARAMETRI × 4 BYTE = **~28 GB di memoria**





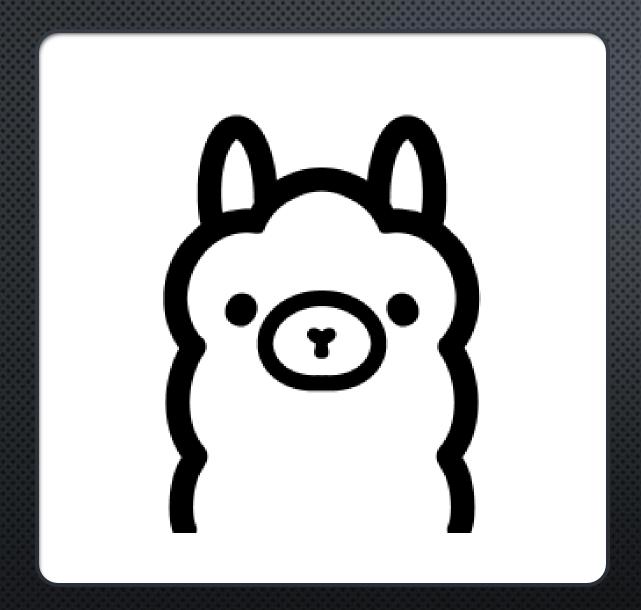
DEMO TIME



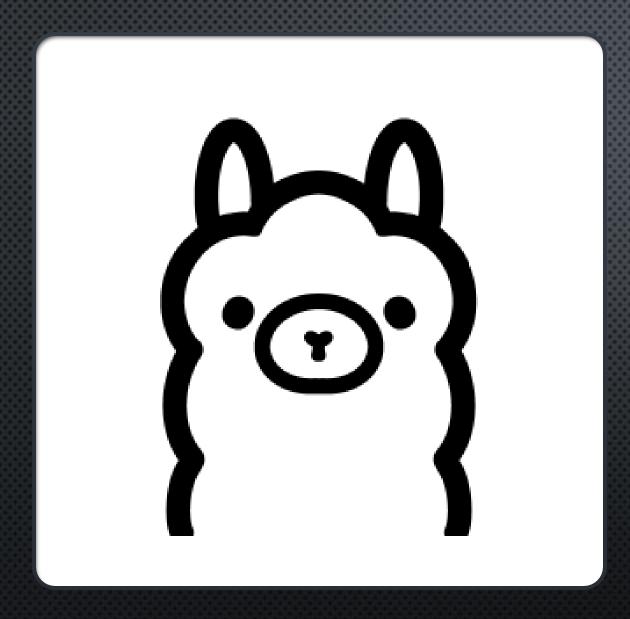
MOTORI DI INFERENZA

- TO UN MOTORE DI INFERENZA È IL SOFTWARE CHE CARICA I PESI DEL MODELLO, GESTISCE LA CACHE, ESEGUE IL FORWARD-PASS E RESTITUISCE OUTPUT AL PROMPT.
- THROUGHPUT, MEMORIA E HARDWARE SPECIFICO (CPU/GPU/QUANTIZZAZIONE).
- LA SCELTA DEL MOTORE PUÒ FARE LA DIFFERENZA SU VELOCITÀ (TIME TO FIRST TOKEN), USO MEMORIA, SCALABILITÀ.
- The locale, con hardware limitato, sono spesso usati motori leggeri/ottimizzati (es. llama.cpp) che permettono esecuzione su CPU.

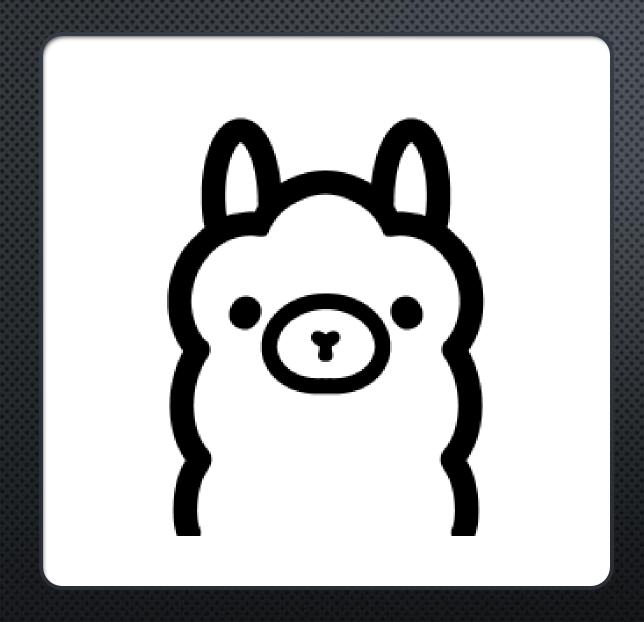




- TRUNTIME LEGGERO PER MODELLI LLM LOCALI
- TINSTALLAZIONE CON UN SOLO COMANDO
- TGESTIONE DI MODELLI (DOWNLOAD, CACHE, RUN)
- API REST LOCALI COMPATIBILI CON OPENAL
- Supporto CPU/GPU (LLAMA.CPP)



- The Serutta Quantizzazioni GGUF → Versioni ottimizzate per CPU/GPU consumer;
- TFORNISCE UN'API REST COMPATIBILE OPENAI,
 MA LOCALMENTE (PORTA: 11434);
- TIGESTISCE CACHE, MODELLI E MEMORIA TRAMITE UN LIVELLO ASTRATTO (NON DEVI COMPILARE NULLA TU);
- TSU MACOS SFRUTTA METAL / MPS, SU LINUX E WINDOWS CUDA O ROCM (SE DISPONIBILI).



TINSTALLAZIONE SU LINUX

CURL -FSSL HTTPS://OLLAMA.COM/INSTALL.SH | SH

Su Windows e Mac con package

https://ollama.com/download

NB: I MODELLI VENGONO SCARICATI IN ~/.OLLAMA/MODELS



TOLLAMA RUN LLM



DEMO TIME



ALTERNATIVE OSS

GUI

- lacktriangledown Open WebUI ightarrow Interfaccia web per Ollama (PIP Install Open-webui)
- \uparrow LM Studio \rightarrow alternativa cross-platform

IDE

- TPLUGIN VS CODE "CONTINUE"
- TO CHAT AI LOCALE VIA API

DATA LAB

NOTEBOOKS CON LANGCHAIN, LLAMA-INDEX, TRANSFORMERS





DEMO TIME



OTTIMIZZARE LE PERFORMANCE

- \bigcirc Quantizzazione (Q4_0, Q5_K_M, ecc.)
- TO GPU VS CPU FALLBACK
- TACHE DEI TOKEN
- TRIDUZIONE CONTEXT WINDOW PER PROMPT LUNGHI
- TEVITARE OUTPUT STREAMING SE NON NECESSARIO



SICUREZZA E PRIVACY

- Thessun dato inviato al cloud
- TGESTIONE LOG LOCALE
- MODELLI VERIFICABILI E AUDITABILI
- BACKUP DI MODELLI E PESI SU STORAGE PERSONALE



DOVE TROVARE I MODELLI

HTTPS://OLLAMA.COM/LIBRARY

HTTPS://HUGGINGFACE.CO/MODELS

HTTPS://GITHUB.COM/OPEN-WEBUI/OPEN-WEBUI



RIEPILOGO

- TCAPITO COSA SONO GLI LLM
- TVISTI I PRINCIPALI MODELLI OPEN SOURCE
- TINSTALLATO OLLAMA
- TESEGUITO UN MODELLO LOCALE
- TCREATO IL TUO PRIMO ASSISTENTE PRIVATO





DOWVNDES





GRAZIE!

