



Dipartimento di Ingegneria Civile, Informatica
e delle Tecnologie Aeronautiche

Corso di Laurea in Ingegneria Informatica

Aumentare le capacità di un chatbot usando Model Context Protocol (MCP)

Anno Accademico 2024/2025

Laureando

Del Prete Andrea

Matricola 589453

Relatore

Prof. Merialdo Paolo

Tutor

Dott. Di Nardo Giorgio

Indice

Indice	1
Introduzione	5
1 Large Language Model - Cosa sono e come funzionano	7
1.1 Introduzione	7
1.1.1 Capacità degli LLM	8
1.1.2 Sfide e limitazioni	8
1.2 Funzionamento	9
1.2.1 Panoramica	9
1.2.2 Tokenizzazione e pre-processing dei dati	9
1.2.3 Il Transformer: principio e componenti fondamentali	9
1.2.4 Obiettivi di pre-training	10
1.2.5 Leggi di scalabilità e considerazioni pratiche	11
1.2.6 Stadi di adattamento: dal pre-training all'uso	11
1.3 <i>Pre-training</i> degli LLM	12
1.3.1 Obiettivi del Pre-training	12
1.3.2 Dataset di Addestramento	13
1.3.3 Sfide del Pre-training	13
1.3.4 Tendenze Recenti	13
1.4 Fine-tuning e Adattamento degli LLM	14
1.4.1 <i>Supervised Fine-tuning</i> (SFT)	14

Indice

1.4.2	Instruction Tuning	14
1.4.3	Reinforcement Learning from Human Feedback (RLHF) . . .	15
1.4.4	Domain Adaptation	15
1.4.5	Multitask e Meta-apprendimento	16
1.4.6	Alignment Tuning	16
1.5	Tecniche di Prompting e In-Context Learning	16
1.5.1	Zero-shot Prompting	16
1.5.2	Few-shot Prompting e In-Context Learning (ICL)	17
1.5.3	Prompt Engineering	17
1.5.4	Prompting per il Ragionamento	17
1.5.5	Multi-turn Prompting	18
1.5.6	Tecniche Avanzate di Prompting	18
1.6	Benchmarking e Valutazione degli LLM	19
1.6.1	Metriche Intrinseche	19
1.6.2	Metriche di Similarità e Qualità Testuale	19
1.6.3	Valutazioni Umanistiche	20
1.6.4	Benchmark di Comprensione Linguistica	20
1.6.5	Benchmark di Ragionamento	20
1.6.6	Benchmark Multilingue e Multimodali	21
1.6.7	Valutazione di Sicurezza e Robustezza	21
1.6.8	Sfide Aperte nella Valutazione	21
1.7	Applicazioni degli LLM	22
1.7.1	Generazione e Comprensione del Linguaggio Naturale	22
1.7.2	Programmazione e Ingegneria del Software	22
1.7.3	Supporto alla Ricerca e alla Scienza	23
1.7.4	Medicina e Assistenza Sanitaria	23
1.7.5	Educazione e Apprendimento	23
1.7.6	Business e Produttività	24
1.7.7	Applicazioni Multimodali	24

Indice

1.7.8	Agenti Autonomi e Tool Use	24
1.8	Sfide, Limitazioni ed Aspetti Etici degli LLM	25
1.8.1	Allucinazioni e Aderenza ai Fatti	25
1.8.2	Bias e Equità	25
1.8.3	Sicurezza e Abusi Potenziali	25
1.8.4	Interpretabilità e Trasparenza	26
1.8.5	Sostenibilità e Impatto Ambientale	26
1.8.6	Accessibilità e Disuguaglianze	26
1.8.7	Allineamento con Valori Umani	27
1.8.8	Aspetti Legali e Regolamentazione	27
1.9	Prospettive Future e Direzioni di Ricerca	27
1.9.1	Efficienza e Sostenibilità	27
1.9.2	Modelli Multimodali	28
1.9.3	Memory-Augmented LLMs	28
1.9.4	Autonomia e Agenti Intelligenti	29
1.9.5	Interpretabilità e Controllo	29
1.9.6	Sicurezza e Robustezza	29
1.9.7	Allineamento e Governance	29
1.9.8	Verso l'Intelligenza Artificiale Generale (AGI)	30
2	Model Context Protocol	31
2.1	Funzionamento	31
2.1.1	Partecipanti	32
2.1.2	Livelli	32
2.1.3	Protocollo del livello dati	34
3	Ticket Management System	37
3.1	Il Model Context Protocol in pratica	38
	Conclusioni e sviluppi futuri	39

Indice

Ringraziamenti	40
Bibliografia	41

Introduzione

Siamo oramai in un'era in cui l'intelligenza artificiale è presente in ogni aspetto della nostra vita digitale, ancora di più dopo l'uscita di ChatGPT, che infranse il record di applicazione con la più rapida crescita nella storia di Internet, ben 100 milioni di utenti in soli 2 mesi [1]. Gli ultimi anni in particolare sono stati caratterizzati dall'intelligenza artificiale generativa, capace cioè di generare testi, audio, immagini o video a partire da un semplice “prompt” dell'utente, come ad esempio i modelli GPT alla base di ChatGPT stesso.

Di particolare interesse sono le IA capaci di generare testo, i cosiddetti Large Language Model (LLM), letteralmente modelli linguistici di grandi dimensioni, in grado di emulare con sorprendente accuratezza le capacità conversazionali di un essere umano. Un'interessante conseguenza del modo in cui questi algoritmi sono stati creati, ma soprattutto della mole di dati usata per addestrarli, è stato l'emergere di comportamenti che potremmo considerare intelligenti, con un livello quasi al pari di quello di un essere umano. Ad esempio, gli LLM sono in grado di ricordare informazioni, generare risposte articolate basate su un ampio contesto e ragionare su problemi complessi derivanti dalle branche più disparate, dalla meccanica quantistica alla biologia evolutiva.

Negli anni è perciò nato il desiderio di rendere ancora più sofisticate le abilità di questi modelli linguistici, finora relegate al “semplice” rispondere ai messaggi degli utenti. Ed è per questo che Anthropic, una delle maggiori aziende nel settore della ricerca sugli LLM e creatrice dei molto diffusi modelli Claude, ha ideato un meccanismo che facilita questo obiettivo: il Model Context Protocol (MCP). Questo

Introduzione

protocollo open-source definisce una procedura standard che consente alle applicazioni di fornire contesto ai modelli linguistici, permettendo loro di accedere ad ed interagire con una serie di strumenti esterni in modo sicuro e controllato. Grazie a ciò, gli LLM possono eseguire operazioni specifiche, recuperare informazioni o integrare funzionalità aggiuntive, migliorando l'efficacia e la personalizzazione delle risposte generate [2].

Obiettivo di questa tesi è illustrare un progetto volto a dimostrare nella pratica le funzionalità del Model Context Protocol. Il progetto in questione è un sistema di ticket management, che consente agli utenti di interfacciarsi con un chatbot che rende automatica la creazione di ticket delineanti i problemi riscontrati e l'inoltro degli stessi a degli sviluppatori, che possono visionarli ed avviare delle procedure per risolverli.

1

Large Language Model - Cosa sono e come funzionano

1.1 Introduzione

Il linguaggio rappresenta uno strumento fondamentale per gli esseri umani, sia come mezzo di comunicazione e trasmissione di conoscenza, sia come interfaccia privilegiata per l'interazione con le macchine. Lo sviluppo di modelli linguistici sempre più potenti è stato quindi al centro della ricerca in *Natural Language Processing* (NLP).

Storicamente, si è passati da approcci **statistici** (basati su n-grammi e modelli probabilistici), a **modelli neurali**, fino ai **modelli pre-addestrati** (*Pre-trained Language Models*, PLM). Questi ultimi, addestrati in maniera auto-supervisionata su grandi corpus testuali, hanno introdotto il concetto di rappresentazioni generali del linguaggio, successivamente adattabili tramite *fine-tuning* a specifici compiti [3][4][5].

Un punto di svolta è stato l'introduzione dei ***Transformer*** [6], che hanno reso possibile scalare i modelli a miliardi di parametri grazie al meccanismo di attenzione, alla parallelizzazione ed alla disponibilità di hardware e dataset sempre più estesi. Questo ha portato alla nascita dei **Large Language Models** (LLM).

1. Large Language Model - Cosa sono e come funzionano

1.1.1 Capacità degli LLM

Con la crescita di scala, gli LLM hanno mostrato proprietà nuove:

- **Generalizzazione senza fine-tuning:** capacità di risolvere compiti mai visti semplicemente a partire da una descrizione testuale (*prompting*).
- **Zero-shot e few-shot learning:** risoluzione di compiti senza esempi (*zero-shot*) o con pochissimi esempi forniti nel prompt (*few-shot*). [7][8][9][10][11][12]
- **In-context learning:** adattamento “al volo” alle istruzioni fornite nell’input.
- **Abilità emergenti:** ragionamento, pianificazione, decision-making, risoluzione di problemi logici, che non erano stati esplicitamente programmati durante l’addestramento. [13][14][15]

1.1.2 Sfide e limitazioni

Nonostante queste capacità impressionanti, gli LLM presentano limiti significativi:

- **Costi computazionali ed energetici elevatissimi,** sia in fase di addestramento che di inferenza (i.e.: durante l’uso effettivo del modello).
- **Allucinazioni:** generazione di contenuti plausibili ma in realtà falsi.
- **Bias e tossicità:** riflessione di stereotipi e pregiudizi presenti nei dati di addestramento.
- **Problemi di allineamento:** difficoltà nel seguire fedelmente l’intento umano; si ricorre quindi a tecniche come *instruction tuning* e *reinforcement learning from human feedback* (RLHF).
- **Accessibilità limitata:** solo pochi attori industriali dispongono delle risorse necessarie per addestrare modelli di queste dimensioni.

1.2 Funzionamento

Illustriamo ora i concetti fondamentali per comprendere la struttura e il funzionamento dei *Large Language Models* (LLM).

1.2.1 Panoramica

Per capire come funzionano gli LLM è utile seguire il flusso informativo: il testo grezzo viene prima pre-processato e tokenizzato, i token vengono poi trasformati in rappresentazioni numeriche che attraversano una architettura di Transformer basata su meccanismi di attenzione e infine il modello viene addestrato (spesso in modo distribuito). Le diverse scelte in ciascuna di queste fasi (tokenizzazione, qualità dei dati, obiettivo di training, architettura) determinano in gran parte le capacità e i limiti del modello.

1.2.2 Tokenizzazione e pre-processing dei dati

La tokenizzazione spezza il testo in unità elementari chiamate token: possono essere singoli caratteri, sottoparole (subword) o parole complete. Le tecniche comuni (es. *Byte-Pair Encoding*, *SentencePiece*) cercano un compromesso tra un vocabolario piccolo e la capacità di rappresentare parole rare o morfologie complesse. Prima della tokenizzazione si eseguono operazioni di pre-processing: filtraggio della qualità, deduplicazione e rimozione di informazioni sensibili. Un buon pre-processing dei dati riduce inoltre rumore, bias e contenuti non desiderati, migliorando la generalizzazione del modello.

1.2.3 Il Transformer: principio e componenti fondamentali

Il Transformer è l'architettura alla base degli LLM moderni. Si basa su un principio chiamato *self-attention*: nel rappresentare matematicamente i token del prompt,

1. Large Language Model - Cosa sono e come funzionano

trasformandoli in oggetti chiamati *embeddings*, si tiene conto non solo della parola isolata ma anche dell'intera frase nella quale si trova, aggiungendo così il contesto.

- ***Multi-head attention***: il modello osserva la frase da più punti di vista contemporaneamente (più "teste") per catturare relazioni diverse, come quelle grammaticali o semantiche.
- ***Positional encoding***: poiché l'attenzione è insensibile alla posizione, si aggiunge un segnale posizionale per preservare l'ordine delle parole.
- ***Feed-forward locale***: dopo l'attenzione, ogni token viene elaborato da una piccola rete indipendente che raffina la sua rappresentazione.

Ci sono diverse tipologie di Transformer:

- ***Encoder-decoder***: due parti separate, una legge l'input (*encoder*) e l'altra genera l'output (*decoder*). Utile per traduzione o riassunti.
- ***Decoder autoregressivo***: genera una parola alla volta guardando solo le parole già prodotte (utilizzato nei modelli GPT).
- **Architetture ibride**: combinano attenzione bidirezionale e causale per casi specifici.
- ***Mixture-of-Experts (MoE)***: il modello ha molti "esperti" specializzati e per ogni input ne attiva solo alcuni, migliorando l'efficienza su larga scala.

1.2.4 Obiettivi di pre-training

Gli obiettivi di addestramento definiscono cosa il modello apprende:

- ***Causal Language Modeling (CLM)***: prevedere il token successivo (approccio autoregressivo).

1. Large Language Model - Cosa sono e come funzionano

- ***Masked Language Modeling (MLM)***: ricostruire token volontariamente mascherati (approccio bidirezionale).
- ***Permutation Language Modeling***: il modello predice token in ordine permutato.
- ***Prefix / Unified objectives***: combinano modalità diverse per aumentare la flessibilità. La scelta dell'obiettivo influenza l'uso finale: i modelli autoregressivi sono particolarmente adatti alla generazione libera, mentre i modelli basati su masking risultano efficaci per rappresentazioni contestuali e compiti di classificazione.

1.2.5 Leggi di scalabilità e considerazioni pratiche

Studi empirici sulle cosiddette *scaling laws* mostrano che le prestazioni migliorano in modo prevedibile all'aumentare del numero dei parametri, quantità di dati e risorse computazionali, spesso seguendo relazioni di tipo esponenziali. Tuttavia, aumentare solamente la quantità di dati senza adeguare la capacità del modello o le risorse computazionali non garantisce miglioramenti proporzionali. È quindi fondamentale bilanciare dimensione del modello, qualità e quantità dei dati e budget computazionale per ottenere guadagni effettivi.

1.2.6 Stadi di adattamento: dal pre-training all'uso

Il ciclo di vita di un LLM si articola tipicamente in tre fasi principali:

1. **Pre-training**: apprendimento auto-supervisionato su grandi corpus testuali per acquisire conoscenze linguistiche generali.
2. **Fine-tuning / Alignment**: adattamento a compiti specifici o preferenze umane.
3. **Prompting / Utilizzo**: impiego diretto tramite prompt.

1. Large Language Model - Cosa sono e come funzionano

In sintesi, un LLM è il risultato di una catena integrata: dati curati e tokenizzazione efficiente alimentano un'architettura Transformer che, tramite meccanismi di self-attention, costruisce rappresentazioni robuste. L'addestramento distribuito su larga scala con obiettivi appropriati modella il comportamento finale, mentre tecniche di adattamento e prompting permettono l'impiego del modello in applicazioni concrete.

1.3 *Pre-training* degli LLM

Il *pre-training* è la fase cruciale che consente ai Large Language Models (LLM) di acquisire conoscenza linguistica generale e capacità di rappresentazione del testo prima di qualsiasi adattamento specifico.

1.3.1 Obiettivi del Pre-training

Gli LLM vengono addestrati in modo auto-supervisionato su grandi corpus testuali con diversi obiettivi di ottimizzazione:

- **Causal Language Modeling (CLM):** il modello predice il prossimo token dato il contesto precedente.
- **Masked Language Modeling (MLM):** alcuni token vengono mascherati e il modello deve ricostruirli.
- **Permutation Language Modeling:** il modello predice token in ordine permutato.
- **Prefix Language Modeling:** viene scelto un prefisso casuale e il modello deve predire i token successivi.
- **Unified Language Modeling:** combina strategie diverse per aumentare la flessibilità.

1. Large Language Model - Cosa sono e come funzionano

1.3.2 Dataset di Addestramento

La qualità e la varietà dei dati di pre-training sono fondamentali. Le principali categorie di dataset includono:

- **Corpus web:** testi raccolti da siti web, ad esempio *Common Crawl*.
- **Enciclopedie e risorse strutturate,** come Wikipedia.
- **Libri e testi lunghi,** ad esempio *BooksCorpus*.
- **Forum e codice,** inclusi dataset come *The Pile*, *Github*, *StackExchange*.

Prima dell'addestramento, i dataset subiscono processi di filtraggio e deduplicazione.

1.3.3 Sfide del Pre-training

L'addestramento degli LLM pone sfide di tipo computazionale e metodologico:

- **Costo computazionale:** i modelli con centinaia di miliardi di parametri richiedono settimane di addestramento su supercomputer con migliaia di GPU.
- **Efficienza dei dati:** la disponibilità di dati di alta qualità è limitata.
- **Rumore e bias:** i dati web contengono rumore, contenuti tossici e bias culturali che vengono appresi dal modello.
- **Controllo della qualità:** è necessario implementare filtri e curare la selezione delle fonti.

1.3.4 Tendenze Recenti

Le ricerche recenti nel pre-training hanno introdotto innovazioni quali:

- **Adattamento multilingue:** uso di corpus in centinaia di lingue.

1. Large Language Model - Cosa sono e come funzionano

- **Addestramento multimodale:** integrazione di testo con immagini, audio o codice.
- **Data Mixture Optimization:** selezione ponderata delle fonti testuali per massimizzare la qualità dell'apprendimento.
- **Continual Pre-training:** estendere progressivamente le conoscenze aggiornando i modelli con nuovi dati.

1.4 Fine-tuning e Adattamento degli LLM

Dopo il pre-training, gli LLM vengono tipicamente sottoposti a una fase di **adattamento**, che ne affina le capacità su compiti specifici o ne migliora l'allineamento con le intenzioni umane. Questo capitolo descrive i metodi principali di fine-tuning, le strategie di efficientamento dei parametri e le tecniche di allineamento.

1.4.1 *Supervised Fine-tuning* (SFT)

Consiste nell'addestrare ulteriormente l'LLM su un dataset etichettato per compiti specifici (es. traduzione, classificazione, QA).

- Viene usato per adattare un modello generalista a un dominio ristretto (es. medicina, diritto).
- Può essere effettuato con dataset relativamente piccoli rispetto al pre-training.
- Rischia di causare *catastrophic forgetting*, ossia perdita delle conoscenze generali.

1.4.2 Instruction Tuning

Gli LLM vengono ottimizzati per comprendere e seguire istruzioni in linguaggio naturale.

1. Large Language Model - Cosa sono e come funzionano

- Basato su dataset di tipo *prompt-response*, costruiti manualmente o tramite crowdsourcing.
- Introduce la capacità di rispondere in modo coerente a comandi generali senza ulteriori adattamenti.

1.4.3 Reinforcement Learning from Human Feedback (RLHF)

Il **RLHF** rappresenta una tecnica fondamentale per allineare i modelli alle preferenze umane.

1. Si raccoglie un dataset di risposte generate dal modello valutate da esseri umani.
2. Si addestra un **modello di reward** che impara a prevedere la valutazione di cui prima.
3. L'LLM viene ulteriormente ottimizzato con tecniche di reinforcement learning.

Questo approccio è stato utilizzato da, ad esempio, ChatGPT.

1.4.4 Domain Adaptation

Gli LLM possono essere specializzati in domini specifici:

- **Continued Pre-training**: ulteriore addestramento su corpus specializzati.
- **Fine-tuning supervisionato su task di dominio**.
- **Approcci ibridi**: combinazioni di pre-training e PEFT (*Parameter-Efficient Fine-Tuning*).

1. Large Language Model - Cosa sono e come funzionano

1.4.5 Multitask e Meta-apprendimento

Oltre al fine-tuning per singoli compiti, sono stati esplorati approcci che mirano a una maggiore generalizzazione:

- ***Multi-task Fine-tuning***: addestramento simultaneo su più compiti.
- **Meta-apprendimento**: insegnare al modello ad adattarsi rapidamente a nuovi compiti con pochi esempi.

1.4.6 Alignment Tuning

Per garantire sicurezza, utilità e rispetto delle preferenze etiche, si introducono strategie di **alignment**:

- ***Constitutional AI***: uso di principi predefiniti come linee guida per il comportamento del modello.
- ***Self-critique e revisione***: i modelli generano e valutano autonomamente le proprie risposte.
- **Feedback multimodale**: estensione del RLHF con segnali da immagini, audio o interazioni più complesse.

1.5 Tecniche di Prompting e In-Context Learning

Dopo il pre-training e l'eventuale fine-tuning, i Large Language Models (LLM) vengono tipicamente utilizzati tramite la formulazione di **prompt**, cioè sequenze di input testuali che guidano la generazione.

1.5.1 Zero-shot Prompting

Nel ***zero-shot prompting*** il modello riceve soltanto una descrizione del compito in linguaggio naturale, senza esempi.

1. Large Language Model - Cosa sono e come funzionano

- Dimostra la capacità di generalizzazione appresa durante il pre-training.
- Funziona bene per compiti semplici e largamente rappresentati nei dati di training.
- È la modalità più naturale di interazione con un LLM.

1.5.2 Few-shot Prompting e In-Context Learning (ICL)

Nel **few-shot prompting**, al modello vengono forniti alcuni esempi di input-output all'interno del prompt. Consente al modello di inferire lo schema del compito e adattarsi in tempo reale. Questo comportamento è noto come *in-context learning* (ICL).

1.5.3 Prompt Engineering

La formulazione del prompt influenza drasticamente la qualità delle risposte. Sono dunque emerse tecniche sistematiche di **prompt engineering**:

- **Instruction-style prompts**: forniscono istruzioni chiare e contestualizzate.
- **Template-based prompting**: utilizzo di schemi predefiniti per generare input coerenti.
- **Role prompting**: assegnazione di un “ruolo” al modello (es. “Sei un medico con decenni di esperienza. . .”).

1.5.4 Prompting per il Ragionamento

Gli LLM possono essere spinti a mostrare capacità di ragionamento tramite prompt opportuni. Le tecniche principali sono:

- **Chain-of-Thought (CoT)**: incoraggia il modello a produrre spiegazioni passo-passo.

1. Large Language Model - Cosa sono e come funzionano

- ***Self-Consistency***: genera più catene di ragionamento e seleziona la risposta più frequente.
- ***Tree-of-Thought (ToT)***: esplora diversi percorsi di ragionamento, con possibilità di backtracking.
- ***Generated Knowledge Prompting***: il modello produce prima conoscenza di supporto, poi la usa per rispondere.

1.5.5 Multi-turn Prompting

Alcuni compiti complessi richiedono interazioni iterative con il modello.

- **Single-turn prompting**: il modello riceve tutte le informazioni in un unico prompt.
- **Multi-turn prompting**: l'interazione è suddivisa in più turni, con feedback e aggiustamenti progressivi. Questo approccio è alla base degli **agenti autonomi** basati su LLM.

1.5.6 Tecniche Avanzate di Prompting

Sono state sviluppate ulteriori strategie per aumentare l'affidabilità:

- **Prompt Chaining**: suddividere un compito complesso in più sotto-compiti concatenati.
- **Automatic Prompt Generation**: generazione automatica di prompt efficaci tramite ottimizzazione o meta-modelli.
- **Retrieval-Augmented Prompting**: combinazione con sistemi di recupero documentale per fornire contesto aggiornato e ridurre le allucinazioni.

1.6 Benchmarking e Valutazione degli LLM

La valutazione degli LLM è una componente cruciale per misurarne le capacità, confrontare modelli diversi e guidare lo sviluppo futuro. Questo capitolo analizza le principali metriche, i benchmark più utilizzati e le sfide aperte nella valutazione.

1.6.1 Metriche Intrinseche

Le metriche intrinseche valutano la qualità linguistica del modello in termini statistici e probabilistici:

- ***Perplexity***: misura la capacità del modello di prevedere la sequenza successiva. È una delle metriche più classiche, ma non sempre correlata alla qualità percepita.
- ***Log-likelihood***: probabilità logaritmica delle sequenze osservate, usata per confronti diretti tra modelli.
- ***Cross-entropy***: misura dell'errore medio nella predizione dei token.

1.6.2 Metriche di Similarità e Qualità Testuale

Per valutare la corrispondenza tra output generati e testi di riferimento si impiegano metriche di similarità lessicale e semantica:

- ***BLEU*, *ROUGE*, *METEOR***: confrontano n-grammi con testi di riferimento.
- ***BERTScore***: utilizza embedding contestuali per misurare similarità semantica.
- **MoverScore** e **BLEURT**: metriche basate su modelli neurali più sofisticati.

1. Large Language Model - Cosa sono e come funzionano

1.6.3 Valutazioni Umanistiche

In molti casi le metriche automatiche non sono sufficienti, e si ricorre a giudizi umani:

- **Valutazioni dirette:** annotatori umani valutano la qualità, la coerenza e l'utilità delle risposte.
- **Pairwise comparison:** confronto tra risposte di modelli diversi.
- **Preference modeling:** raccolta di preferenze umane per addestrare modelli di reward (come in RLHF).

1.6.4 Benchmark di Comprensione Linguistica

Sono stati proposti benchmark standardizzati per valutare capacità di comprensione:

- **GLUE** e **SuperGLUE**: valutano compiti di classificazione e comprensione testuale.
- **SQuAD**: benchmark per question answering su testi di Wikipedia.
- **RACE**, **BoolQ**: valutano comprensione di testi complessi e risposte booleane.

1.6.5 Benchmark di Ragionamento

Per testare capacità di inferenza e ragionamento:

- **MMLU** (**Massive Multitask Language Understanding**): copre centinaia di compiti e discipline.
- **BIG-bench**: insieme di migliaia di task creati da comunità accademica e industriale.
- **GSM8K**: valutazione su problemi matematici di livello scolastico.
- **ARC** (**AI2 Reasoning Challenge**): domande scientifiche a scelta multipla.

1. Large Language Model - Cosa sono e come funzionano

1.6.6 Benchmark Multilingue e Multimodali

Per valutare capacità oltre l'inglese e il puro testo:

- *XGLUE*, *XTREME*: benchmarking multilingue.
- *MMBench*, *MMMU*: valutazione multimodale (testo + immagini).
- *CodeXGLUE*, *HumanEval*: benchmarking su programmazione e codice.

1.6.7 Valutazione di Sicurezza e Robustezza

Un aspetto emergente riguarda la valutazione di rischi e affidabilità:

- **Adversarial Testing**: inserimento di prompt avversari per indurre errori.
- **Red-teaming**: valutazione guidata da esperti per testare rischi etici e di sicurezza.
- **Hallucination Benchmarks**: misurano la tendenza del modello a generare informazioni false.

1.6.8 Sfide Aperte nella Valutazione

- Le metriche tradizionali non catturano abilità complesse come ragionamento, coerenza a lungo termine, pianificazione.
- Molti benchmark soffrono di saturazione: i modelli raggiungono rapidamente prestazioni umane.
- Mancano valutazioni standardizzate per proprietà emergenti e per compiti compositivi.
- Le valutazioni umane sono costose, soggettive e poco scalabili.

1.7 Applicazioni degli LLM

Gli **Large Language Models** (LLM) hanno trovato applicazioni in un'ampia gamma di domini, dalla generazione di linguaggio naturale alla programmazione, dalla ricerca scientifica alla medicina. Questo capitolo presenta le principali aree di utilizzo e gli scenari emergenti.

1.7.1 Generazione e Comprensione del Linguaggio Naturale

Gli LLM eccellono in compiti tradizionali di NLP:

- **Traduzione automatica:** modelli come GPT, T5 e mBART hanno raggiunto performance vicine ai sistemi dedicati.
- **Riassunto testuale:** generazione di sintesi coerenti e concise da testi lunghi.
- **Parafrasi e riformulazione:** utile per la scrittura assistita e l'elaborazione creativa.
- **Conversazione e chatbot:** applicazioni come *ChatGPT* e Claude forniscono assistenza virtuale multi-dominio.

1.7.2 Programmazione e Ingegneria del Software

Gli LLM specializzati sul codice hanno mostrato grande impatto nello sviluppo software:

- **Code completion:** completamento automatico di funzioni e snippet.
- **Code generation:** traduzione da linguaggio naturale a codice eseguibile.
- **Code explanation e refactoring:** spiegazione e ottimizzazione di codice esistente.
- Strumenti come *Codex*, *Copilot* e *AlphaCode* hanno reso accessibile la programmazione anche a utenti non esperti.

1. Large Language Model - Cosa sono e come funzionano

1.7.3 Supporto alla Ricerca e alla Scienza

Gli LLM vengono sempre più utilizzati come strumenti di accelerazione scientifica:

- **Recupero e sintesi della letteratura scientifica.**
- **Generazione di ipotesi:** suggerimenti per esperimenti o linee di ricerca.
- **Analisi dati e interpretazione:** supporto a scienziati in fisica, biologia, chimica.
- **Scoperta di molecole e farmaci:** integrazione con modelli di chimica computazionale.

1.7.4 Medicina e Assistenza Sanitaria

Il settore medico rappresenta un campo sensibile ma promettente:

- **Supporto diagnostico:** analisi di sintomi e cartelle cliniche.
- **Assistenza alla scrittura clinica:** redazione automatica di referti.
- **Educazione medica:** generazione di materiali di studio personalizzati.
- **Chatbot sanitari:** fornire informazioni mediche preliminari ai pazienti.

Sono necessarie forti misure di sicurezza per ridurre rischi di allucinazioni e bias.

1.7.5 Educazione e Apprendimento

Gli LLM possono agire da tutor personalizzati:

- **Creazione di materiali didattici adattivi.**
- **Assistenza individuale agli studenti** con spiegazioni su misura.
- **Valutazione automatizzata** di esercizi e saggi.
- Supporto all'apprendimento delle lingue attraverso dialoghi interattivi.

1. Large Language Model - Cosa sono e come funzionano

1.7.6 Business e Produttività

Le applicazioni aziendali sono tra le più immediate:

- **Automazione del customer service** con chatbot multilingue.
- **Generazione di report e analisi** a partire da dati aziendali.
- **Supporto alle decisioni** tramite analisi predittive.
- **Content creation**: marketing, copywriting e produzione di contenuti digitali.

1.7.7 Applicazioni Multimodali

Gli LLM vengono estesi per gestire input multimodali:

- **Vision-Language Models**, come *CLIP*, *Flamingo*, GPT-4 multimodale.
- **Image captioning e visual question answering**.
- **Audio e speech**: trascrizione, traduzione e sintesi vocale.
- **Robotica**: uso di LLM per interpretare istruzioni linguistiche e guidare azioni fisiche.

1.7.8 Agenti Autonomi e Tool Use

Un campo emergente riguarda l'integrazione degli LLM in agenti autonomi:

- **Tool-augmented LLMs**: modelli che interagiscono con strumenti esterni come motori di ricerca, calcolatrici o database.
- **Agenti multi-step**: in grado di pianificare sequenze di azioni.
- Framework come *LangChain*, *AutoGPT* e *BabyAGI* mostrano le potenzialità degli LLM come agenti generici.

1.8 Sfide, Limitazioni ed Aspetti Etici degli LLM

Nonostante i notevoli progressi, gli LLM presentano ancora numerose sfide tecniche, limitazioni pratiche e rischi etici. Questo capitolo esamina i principali problemi legati alla loro adozione e le preoccupazioni sociali connesse al loro utilizzo.

1.8.1 Allucinazioni e Aderenza ai Fatti

Gli LLM spesso generano contenuti plausibili ma falsi, un fenomeno noto come **allucinazioni**. Possono inventare dati, riferimenti bibliografici o fatti storici inesistenti. Ciò riduce l'affidabilità in contesti critici come medicina, diritto o scienza. Tecniche come *retrieval-augmented generation* (RAG) sono state proposte per mitigare questo problema.

1.8.2 Bias e Equità

Gli LLM ereditano i bias presenti nei dati di addestramento.

- Riproducono stereotipi di genere, etnia, religione, orientamento sessuale.
- Possono amplificare discriminazioni già esistenti nella società.

Sono in corso ricerche su dataset bilanciati e metodi di de-biasing.

1.8.3 Sicurezza e Abusi Potenziali

Gli LLM possono essere usati impropriamente per generare contenuti dannosi.

- **Disinformazione:** generazione automatica di fake news e propaganda.
- **Cybersecurity:** creazione di codice malevolo o phishing su larga scala.
- **Deepfakes multimodali:** se integrati con modelli di immagini o audio.

1. Large Language Model - Cosa sono e come funzionano

1.8.4 Interpretabilità e Trasparenza

Un limite importante riguarda la natura **black-box** degli LLM.

- È difficile spiegare come vengano prese le decisioni interne al modello.
- La mancanza di interpretabilità riduce la fiducia in applicazioni sensibili.

Sono state proposte tecniche di *model probing*, analisi delle rappresentazioni interne e spiegazioni post-hoc.

1.8.5 Sostenibilità e Impatto Ambientale

Il training di LLM richiede un enorme consumo energetico.

- Addestrare un singolo modello può generare tonnellate di CO₂.
- Questo solleva preoccupazioni ambientali e di sostenibilità.

Approcci di efficientamento (quantizzazione, pruning, distillazione) mirano a ridurre tali costi.

1.8.6 Accessibilità e Disuguaglianze

Gli LLM di punta sono sviluppati principalmente da grandi aziende tecnologiche.

- L'alto costo computazionale limita la ricerca accademica indipendente.
- Si rischia una concentrazione di potere e conoscenza in poche mani.

Iniziative open-source (es. *LLaMA*, *Falcon*, *MPT*) cercano di democratizzare l'accesso.

1. Large Language Model - Cosa sono e come funzionano

1.8.7 Allineamento con Valori Umani

Un'altra sfida riguarda l'**alignment** con preferenze etiche e valori umani.

- Gli LLM devono essere utili, onesti e innocui.
- Tecniche come **RLHF** e *Constitutional AI* sono state sviluppate per migliorare l'allineamento.
- Rimane complesso garantire coerenza su culture e comunità diverse.

1.8.8 Aspetti Legali e Regolamentazione

Con la diffusione degli LLM emergono questioni legali:

- **Copyright:** rischio di violazioni dovute a testi generati simili a fonti protette.
- **Privacy:** possibilità che vengano riprodotte informazioni sensibili dai dati di training.
- **Regolamentazione:** discussioni in corso a livello internazionale su AI Act e framework etici.

1.9 Prospettive Future e Direzioni di Ricerca

Gli LLM rappresentano uno dei campi più dinamici dell'intelligenza artificiale. Dopo i progressi ottenuti negli ultimi anni, la comunità scientifica si interroga su come evolveranno i modelli futuri e quali direzioni di ricerca saranno più rilevanti. Questo capitolo discute le prospettive principali.

1.9.1 Efficienza e Sostenibilità

Una priorità sarà lo sviluppo di modelli più efficienti dal punto di vista computazionale ed energetico.

1. Large Language Model - Cosa sono e come funzionano

- Tecniche di compressione (quantizzazione, pruning, distillazione) sempre più sofisticate.
- Architetture ibride che combinano Transformer con nuovi meccanismi di attenzione.
- Algoritmi di addestramento ottimizzati per ridurre costi e tempi.

1.9.2 Modelli Multimodali

Gli LLM tenderanno a integrarsi con altre modalità sensoriali:

- **Vision-Language Models:** unione di testo e immagini.
- **Speech e audio:** modelli in grado di comprendere e generare linguaggio parlato.
- **Video e robotica:** comprensione di sequenze visive e comandi linguistici per il controllo di agenti fisici.

L'obiettivo è costruire **modelli universali** in grado di gestire input multimodali.

1.9.3 Memory-Augmented LLMs

Attualmente gli LLM hanno memoria limitata alla finestra contestuale.

- Ricerca su architetture con memoria esterna persistente.
- Meccanismi per richiamare informazioni apprese in sessioni precedenti.
- Possibilità di modelli che accumulano conoscenza nel tempo, riducendo la necessità di retraining massivo.

1. Large Language Model - Cosa sono e come funzionano

1.9.4 Autonomia e Agenti Intelligenti

Un filone emergente è l'integrazione degli LLM in agenti autonomi.

- Capacità di pianificazione multi-step tramite prompting avanzato.
- Integrazione con tool esterni: database, motori di ricerca, calcolatrici, API.
- Applicazioni in domini complessi come finanza, ricerca scientifica, sviluppo software.

1.9.5 Interpretabilità e Controllo

La comunità scientifica riconosce la necessità di modelli più trasparenti:

- Strumenti di **explainable AI** per comprendere le decisioni interne.
- Analisi delle rappresentazioni neurali per capire l'emergere delle abilità.
- Tecniche di controllo fine-grained per modulare il comportamento del modello.

1.9.6 Sicurezza e Robustezza

Un'altra direzione riguarda la costruzione di LLM affidabili e sicuri:

- Difese contro attacchi avversari e prompt injection.
- Riduzione delle allucinazioni tramite retrieval e grounding su basi di conoscenza.
- Maggiore robustezza a input rumorosi o manipolati.

1.9.7 Allineamento e Governance

Il tema dell'allineamento con valori umani resterà centrale.

- Evoluzione di tecniche come RLHF, Constitutional AI e feedback multimodale.

1. Large Language Model - Cosa sono e come funzionano

- Creazione di framework etici condivisi per la progettazione di modelli sicuri.
- Discussione su regolamentazione internazionale e open science.

1.9.8 Verso l'Intelligenza Artificiale Generale (AGI)

Infine, gli LLM sono considerati possibili passi verso l'AGI.

- Studi sull'emergere di capacità generali di ragionamento e apprendimento.
- Integrazione con sistemi simbolici e neuro-simbolici.
- Possibilità di architetture future che superino i limiti dei Transformer attuali.

2

Model Context Protocol

2.1 Funzionamento

MCP è un protocollo open-source che standardizza il modo in cui le applicazioni forniscono contesto agli LLM. Si può pensare a MCP come ad una porta USB-C per le applicazioni IA. Proprio come USB-C fornisce un modo standardizzato per collegare i dispositivi a varie periferiche e accessori, MCP fornisce un modo standardizzato per collegare i modelli IA a diverse fonti di dati e strumenti. MCP consente di creare agenti e flussi di lavoro complessi basati sugli LLM e connette i modelli con il mondo. [16]

Il protocollo MCP include i seguenti progetti:

- Specifiche MCP: una specifica di MCP che delinea i requisiti di implementazione per client e server.
- SDK MCP: *Software Development Kit* (SDK) per diversi linguaggi di programmazione che implementano MCP.
- Strumenti di sviluppo MCP.
- Implementazioni di server MCP di riferimento.

2. Model Context Protocol

MCP si concentra esclusivamente sul protocollo per lo scambio di contesto, senza stabilire come le applicazioni IA utilizzino gli LLM o gestiscano il contesto fornito. [17]

2.1.1 Partecipanti

MCP segue un'architettura client-server in cui un host MCP, un'applicazione IA come Claude Code o Claude Desktop, stabilisce connessioni ad uno o più server MCP. L'host MCP realizza questo creando un client MCP per ogni server MCP. Ogni client MCP mantiene una connessione uno-a-uno dedicata con il suo server MCP corrispondente. I principali partecipanti all'architettura MCP sono:

- Host MCP: l'applicazione IA che coordina e gestisce uno o più client MCP.
- Client MCP: un componente che mantiene una connessione ad un server MCP ed ottiene, da un server MCP, il contesto che l'host MCP può utilizzare.
- Server MCP: un programma che fornisce contesto ai client MCP; possono essere eseguiti sia in locale che in remoto. [17]

2.1.2 Livelli

MCP è costituito da due livelli:

- Livello dati: definisce il protocollo, basato su JSON-RPC, per la comunicazione client-server, inclusa la gestione del ciclo di vita e le primitive principali, come strumenti, risorse, prompt e notifiche.
- Livello trasporto: definisce i meccanismi e i canali di comunicazione che consentono lo scambio di dati tra client e server, inclusi l'instaurazione di connessioni specifiche per il trasporto, il framing dei messaggi e l'autorizzazione.

Concettualmente, il livello dati è il livello interno, mentre il livello trasporto è il livello esterno. [17]

2. Model Context Protocol

Livello dati

Il livello dati implementa un protocollo di scambio basato su JSON-RPC 2.0 che definisce la struttura e la semantica dei messaggi. Questo livello include:

- Gestione del ciclo di vita: gestisce l’inizializzazione della connessione, la negoziazione delle capacità e la terminazione della connessione tra client e server.
- Funzionalità del server: consente ai server di fornire funzionalità di base, inclusi strumenti per azioni IA, risorse per dati di contesto e richieste per schemi di interazione da e verso il client.
- Funzionalità del client: consente ai server di chiedere al client di campionare dall’LLM, ottenere input dall’utente e registrare messaggi al client.
- Funzionalità di utilità: supporta funzionalità aggiuntive come notifiche per aggiornamenti in tempo reale e monitoraggio dei progressi per operazioni di lunga durata. [17]

Livello di trasporto

Il livello di trasporto gestisce i canali di comunicazione e l’autenticazione tra client e server. Gestisce la creazione della connessione, il framing dei messaggi e la comunicazione sicura tra i partecipanti. MCP supporta due meccanismi di trasporto:

- *Stdio Transport*: utilizza flussi di input/output standard per la comunicazione diretta tra processi locali sulla stessa macchina, garantendo prestazioni ottimali senza sovraccarico di rete.
- *Streamable HTTP transport*: utilizza metodi HTTP POST per i messaggi client-server con eventi inviati dal server opzionali per le funzionalità di streaming. Questo trasporto consente la comunicazione con il server remoto e supporta metodi di autenticazione HTTP standard, inclusi *bearer token*, chiavi API e intestazioni personalizzate.

2. Model Context Protocol

MCP consiglia di utilizzare OAuth per ottenere i token di autenticazione. Il livello di trasporto astrae i dettagli di comunicazione dal livello di protocollo, consentendo lo stesso formato di messaggio del protocollo JSON-RPC 2.0 su tutti i meccanismi di trasporto. [17]

2.1.3 Protocollo del livello dati

Una parte fondamentale di MCP è la definizione dello schema e della semantica tra client e server MCP. Il livello dati è la parte di MCP che definisce le modalità con cui gli sviluppatori possono condividere il contesto dai server MCP ai client MCP. MCP utilizza JSON-RPC 2.0 come protocollo *Remote Procedure Call* (RPC). Client e server si inviano richieste e rispondono di conseguenza. Le notifiche possono essere utilizzate quando non è richiesta alcuna risposta. [17]

Primitive

Le primitive MCP sono il concetto più importante all'interno di MCP. Definiscono ciò che client e server possono offrirsi reciprocamente. Queste primitive specificano i tipi di informazioni contestuali che possono essere condivise con le applicazioni IA e la gamma di azioni che possono essere eseguite. MCP definisce tre primitive principali che i server possono esporre:

- *Tools*: funzioni eseguibili che le applicazioni IA possono invocare per eseguire azioni (ad es.: operazioni su file, chiamate API, query di database).
- *Risorse*: fonti di dati che forniscono informazioni contestuali alle applicazioni IA (ad es.: contenuto di file, record di database, risposte API).
- *Prompt*: schemi riutilizzabili che aiutano a strutturare le interazioni con i modelli linguistici (ad es.: prompt di sistema, prompt *few-shot*).

Ogni tipo di primitiva ha metodi associati per la scoperta (**/list*), il recupero (**/get*) e, in alcuni casi, l'esecuzione (*tools/call*). I client MCP utilizzeranno

2. Model Context Protocol

i metodi `*/list` per scoprire le primitive disponibili. Ad esempio, un client può prima elencare tutti gli strumenti disponibili (`tools/list`) e poi eseguirli. Questa progettazione consente di creare elenchi dinamici.

Come esempio concreto, si consideri un server MCP che fornisce contesto su un database. Può esporre strumenti per interrogare il database, una risorsa che contiene lo schema del database e un prompt che include esempi di interazione con gli strumenti.

MCP definisce anche le primitive che i client possono esporre. Queste primitive consentono agli autori del server MCP di creare interazioni più ricche.

- **Campionamento:** consente ai server di richiedere il completamento del modello linguistico dall'applicazione IA del client. Questa funzionalità è utile quando gli autori del server desiderano accedere ad un modello linguistico, ma vogliono rimanere indipendenti dal modello e non includere un SDK del modello linguistico nel proprio server MCP. Possono utilizzare il metodo `sampling/complete` per richiedere il completamento del modello linguistico dall'applicazione IA del client.
- **Elicitazione:** consente ai server di richiedere informazioni aggiuntive agli utenti. Questa funzionalità è utile quando gli autori del server desiderano ottenere maggiori informazioni dall'utente o chiedere la conferma di un'azione. Possono utilizzare il metodo `elicitation/request` per richiedere informazioni aggiuntive all'utente.
- **Logging:** consente ai server di inviare messaggi di log ai client a scopo di debug e monitoraggio. [17]

Notifiche

Il protocollo supporta notifiche in tempo reale per abilitare aggiornamenti dinamici tra server e client. Ad esempio, quando gli strumenti disponibili su un server cambiano, come quando vengono rese disponibili nuove funzionalità o vengono modificati

2. Model Context Protocol

strumenti esistenti, il server può inviare notifiche di aggiornamento per informare i client connessi di tali modifiche. Le notifiche vengono inviate come messaggi di notifica JSON-RPC 2.0 (senza attendere una risposta) e consentono ai server MCP di fornire aggiornamenti in tempo reale ai client connessi. [17]

3

Ticket Management System

Il *Ticket Management System* è un'applicazione software che ha lo scopo di aiutare gli utenti nella creazione e gestione di ticket, creati per segnalare la presenza di problemi individuati, e nel loro inoltro a degli sviluppatori che possano analizzare e poi risolvere tali problemi. Funzionalità centrale dell'applicazione è il fatto che l'utente si interfaccia con un chatbot, un'istanza della risorsa Azure OpenAI che usa il modello GPT-4.1 nello specifico, per usufruire dei servizi offerti.

Il *Ticket Management System* è stato progettato come una *solution* in C# suddivisa, secondo un approccio modulare, in più progetti indipendenti ma interconnessi. Tale organizzazione riflette i principi delle moderne architetture software a livelli, secondo il principio *Separation of Concerns*, favorendo la manutenibilità e la possibilità di estendere il sistema senza introdurre dipendenze circolari o accoppiamenti eccessivamente rigidi.

La composizione della *solution* prevede cinque progetti principali, ciascuno con un ruolo ben definito: *TM.Shared*, *TM.Data*, *TM.CQRS*, *TM.Function* e *TM.Client*.

Il progetto *TM.Shared* definisce tutte le classi degli oggetti che vengono creati internamente nel codice al momento dell'interazione con l'utente. Tali classi sono:

- **Ticket**: elemento centrale del sistema; rappresenta i ticket creati dagli utenti, con i relativi **Comment**, le **Task** legate, il **TicketStatus** e la **TicketPriority**.

3. Ticket Management System

- **Category:** ogni **Ticket** appartiene ad una categoria.
- **Comment:** i **Ticket** hanno una serie di commenti.
- **Task:** i task sono le operazioni affidate agli sviluppatori, generati dal sistema a seguito della creazione di un **Ticket**.
- **User:** gli utenti che usano il sistema; sono intesi anche gli sviluppatori che devono risolvere le **Task**.

Il progetto *TM.Data* ha il compito di gestire la persistenza delle informazioni, fornendo l'accesso al database ed archiviando gli oggetti creati.

Il progetto *TM.CQRS* introduce l'implementazione del pattern *Command Query Responsibility Segregation*, un pattern architetturale che separa le operazioni di lettura (**Queries**) da quelle di scrittura (**Commands**) all'interno di un'applicazione.

Il cuore della logica applicativa è rappresentato da *TM.Function*. Essa è un'Azure Function, ovvero una soluzione serverless che consente di scrivere meno codice, gestire un'infrastruttura meno complessa e risparmiare sui costi. Non è più necessario preoccuparsi della distribuzione e della gestione dei server, in quanto l'infrastruttura cloud fornisce tutte le risorse aggiornate necessarie per mantenere le applicazioni in esecuzione [18]. Prende il ruolo dell'MCP Host, l'applicazione IA che gestisce la logica principale del Model Context Protocol. Verrà discusso più approfonditamente in seguito.

Infine, *TM.Client* costituisce l'interfaccia utente, nonché l'MCP Client. Esso contiene la sola classe **Program.cs**, occupandosi di avviare l'intera solution e integrare i vari altri progetti, oltre che inizializzare la comunicazione con il chatbot.

3.1 Il Model Context Protocol in pratica

Conclusioni e sviluppi futuri

La tesi è finita

Ringraziamenti

Desidero esprimere la mia più sincera gratitudine ai docenti, per la dedizione e la passione con cui hanno trasmesso le loro conoscenze, contribuendo in modo fondamentale alla mia crescita accademica; ai colleghi in azienda, per la professionalità e disponibilità con la quale mi hanno accolto, guidato e supportato durante il mio percorso di tirocinio, permettendomi di mettere in pratica quanto appreso e di acquisire nuove competenze; ai miei amici e colleghi di università, per il costante sostegno, la collaborazione e la condivisione di momenti indimenticabili, che hanno reso questo percorso più ricco e stimolante; ed infine alla mia famiglia, per il supporto incondizionato che mi hanno sempre dimostrato e senza la quale nulla di tutto questo sarebbe stato possibile.

Bibliografia

- [1] UBS. Latest House View Daily, 2023. URL: <https://www.ubs.com/global/en/wealthmanagement/insights/chief-investment-office/house-view/daily/2023/latest-25052023.html>. Citato a pagina 5.
- [2] Model Context Protocol. Getting Started - Introduction, 2023. URL: <https://modelcontextprotocol.io/docs/getting-started/intro>. Citato a pagina 6.
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2018. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/1810.04805>, arXiv:arXiv:1810.04805. Citato a pagina 7.
- [4] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, pages 2227–2237. Association for Computational Linguistics, 2018. Accessed: 2025-09-02. URL: <https://aclanthology.org/N18-1202>. Citato a pagina 7.
- [5] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation,

Bibliografia

- translation, and comprehension, 2019. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/1910.13461>, arXiv:arXiv:1910.13461. Citato a pagina 7.
- [6] Andrey Chernyavskiy, Dmitry Ilvovsky, and Preslav Nakov. Transformers: "the end of history" for natural language processing? In *Machine Learning and Knowledge Discovery in Databases. Research Track: European Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part III*, volume 12977 of *Lecture Notes in Computer Science*, pages 677–693. Springer, 2021. Accessed: 2025-09-02. doi:10.1007/978-3-030-86523-8_41. Citato a pagina 7.
- [7] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Saurabh Brahma, et al. Scaling instruction-finetuned language models, 2022. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/2210.11416>, arXiv:arXiv:2210.11416. Citato a pagina 8.
- [8] Victor Sanh, Albert Webson, Colin Raffel, Stephen H. Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Abheesht Raja, et al. Multitask prompted training enables zero-shot task generalization, 2021. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/2110.08207>, arXiv:arXiv:2110.08207. Citato a pagina 8.
- [9] Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amir Mirzaei, Aashka Naik, Anjana Ashok, Abhinav Suresh Dhanasekaran, Anmol Arunkumar, David Stap, et al. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5085–5109, 2022. Accessed: 2025-09-02. URL: <https://aclanthology.org/2022.emnlp-main.341>. Citato a pagina 8.

Bibliografia

- [10] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions, 2022. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/2212.10560>, arXiv:arXiv:2212.10560. Citato a pagina 8.
- [11] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pages 27730–27744, 2022. Accessed: 2025-09-02. URL: https://proceedings.neurips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html. Citato a pagina 8.
- [12] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Sharan Batra, Puneet Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models, 2023. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/2307.09288>, arXiv:arXiv:2307.09288. Citato a pagina 8.
- [13] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models, 2022. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/2206.07682>, arXiv:arXiv:2206.07682. Citato a pagina 8.
- [14] Taylor Webb, Keith J. Holyoak, and Hongjing Lu. Emergent analogical reasoning in large language models. *Nature Human Behaviour*, 7(9):1526–1541, 2023. Accessed: 2025-09-02. doi:10.1038/s41562-023-01659-w. Citato a pagina 8.

Bibliografia

- [15] Daniil A. Boiko, Robert MacKnight, and Gabriel Gomes. Emergent autonomous scientific research capabilities of large language models, 2023. Accessed: 2025-09-02. URL: <https://arxiv.org/abs/2304.05332>, arXiv:arXiv:2304.05332. Citato a pagina 8.
- [16] Model Context Protocol. Getting Started - Introduction, 2023. Accessed: 2025-08-16. URL: <https://modelcontextprotocol.io/docs/getting-started/intro>. Citato a pagina 31.
- [17] Model Context Protocol. Learn - Architecture, 2023. Accessed: 2025-08-16. URL: <https://modelcontextprotocol.io/docs/learn/architecture>. Citato alle pagine 32, 33, 34, 35, and 36.
- [18] Microsoft Docs. Azure Functions overview. Accessed: 2025-08-16. URL: <https://learn.microsoft.com/it-it/azure/azure-functions/functions-overview?pivots=programming-language-csharp>. Citato a pagina 38.