



UNIVERSITÀ DELLA CALABRIA

DIPARTIMENTO DI **MATEMATICA
E INFORMATICA**

Artificial Intelligence and Computer Science

**GAEN: Generative AI for Enhanced Narrative -
Rilevamento di Spoiler Mediante Modelli Generativi
Avanzati**

Relatori:
Prof. Gianluigi Greco
Prof. Kazumi Saito

Candidato:
Alessandro Fazio
Matricola 242422

Anno Accademico 2024/2025

Indice

1	Introduzione	6
1.1	Contesto	6
1.2	Motivazioni	6
1.3	Obiettivi	7
1.4	Struttura del documento	7
2	Background	8
2.1	Approcci al problema	8
2.1.1	Approccio basato su regole	8
2.1.2	Approccio basato sul machine learning	9
2.2	Approccio proposto	9
3	Lavori correlati	11
3.1	Contesto della Ricerca	11
3.2	LLM, modelli di linguaggio	11
3.2.1	Attention is all you need	12
3.3	Embedding	13
3.3.1	Vector Database	14
3.4	Fine-tuning di modelli	14
3.5	Transizione all'implementazione	15
4	Design e implementazione	16
4.1	Introduzione	16
4.2	Architettura del sistema	16
4.3	Componenti principali	17
4.3.1	Frontend	17
4.3.2	Backend	17
4.3.3	Ollama	18
4.3.4	Embedder	19
4.4	Dataset e tool utilizzati	19
4.4.1	Visualizer	20

4.5	Preparazione del dataset	22
4.5.1	Scelta dei negativi	23
4.6	Fine-tuning del modello	24
5	Valutazione	25
5.1	Introduzione	25
5.1.1	Specifiche tecniche	25
5.2	Dataset e metodologia	26
5.3	Embedders	26
5.3.1	Risultati	27
5.3.2	Considerazioni	28
5.3.3	Prestazioni in fase di fine-tuning	28
5.4	LLM	28
5.4.1	Gemma3	29
5.4.2	Risultati LLM	29
5.5	Analisi dei risultati	30
5.6	Conclusioni	31
6	Considerazioni Finali	32
6.1	Lavori Futuri	32

Elenco delle figure

3.1	Architettura del Transformer	12
3.2	Esempio di embedding. Fonte: [3]	13
4.1	Architettura del sistema e flusso dei dati	17
4.2	Esempio di scraping. Fonte: Wikipedia	19
4.3	Esempio di grafo creato dal visualizer	20
4.4	Esempio di oversampling	23

Elenco delle tabelle

4.1	Statistiche globali del grafo. Per la descrizione delle metriche, si veda la Sezione 4.4.1.	21
4.2	Statistiche per singoli nodi. Per la descrizione delle metriche, si veda la Sezione 4.4.1.	21
4.3	Formato del dataset	23
5.1	Panoramica dei parametri utilizzati per il fine-tuning e la valutazione del modello.	26
5.2	Risultati della valutazione dei modelli di embedding (Pair) . .	27
5.3	Risultati della valutazione dei modelli di embedding (Triplet) .	27
5.4	Tempi e perdite del fine-tuning (Triplet)	28
5.5	Tempi e perdite del fine-tuning (Pair)	28
5.6	Risultati delle prestazioni dei modelli.	30

Capitolo 1

Introduzione

1.1 Contesto

Nell'era digitale, l'accesso immediato a contenuti multimediali ha trasformato radicalmente il modo in cui fruiamo di film, serie TV, libri e videogiochi. Tuttavia, questa abbondanza di informazioni porta con sé una sfida: la crescente esposizione a spoiler. Gli spoiler, rivelando anticipatamente elementi cruciali della trama, possono compromettere l'esperienza di fruizione, generando frustrazione e delusione negli utenti.

Il problema degli spoiler è particolarmente rilevante in contesti online, dove le discussioni e le recensioni proliferano su forum, social media e piattaforme di streaming. La mancanza di strumenti efficaci per il rilevamento automatico di spoiler rende difficile per gli utenti proteggersi da tali rivelazioni indesiderate.

Questo lavoro propone di esplorare l'utilizzo di Large Language Models (LLM) per affrontare la sfida del rilevamento di spoiler. L'obiettivo è sviluppare un sistema in grado di identificare automaticamente gli spoiler in testi di varia natura, sfruttando le capacità di comprensione del linguaggio naturale degli LLM.

1.2 Motivazioni

La motivazione principale di questa ricerca risiede nella volontà di contribuire allo sviluppo di strumenti che migliorino l'esperienza di fruizione di discussioni riguardanti contenuti multimediali. A tale scopo, è stato realizzato un sistema di rilevamento di spoiler integrato in una estensione per browser, che offre agli utenti un controllo maggiore sulla propria esposizione a possibili spoiler.

Questo lavoro intende esplorare il potenziale degli LLM in un'applicazione pratica e rilevante, provando a contribuire alla comprensione delle loro capacità e limitazioni nel contesto del rilevamento di informazioni specifiche in testi complessi.

1.3 Obiettivi

Gli obiettivi di questo lavoro sono i seguenti:

- Realizzare un sistema di rilevamento di spoiler basato su Large Language Models.
- Valutare l'efficacia del sistema sviluppato nel rilevamento di spoiler in testi di varia natura.
- Integrare il sistema in un'estensione per browser e valutarne l'utilità e l'usabilità.

1.4 Struttura del documento

Il resto del documento è organizzato come segue:

- Nel Capitolo 2 vengono presentati i concetti e le tecnologie alla base della ricerca, con particolare attenzione agli approcci possibili.
- Nel Capitolo 3 vengono esaminati i lavori correlati, con un focus sulle ricerche relative al rilevamento di spoiler.
- Nel Capitolo 4 viene descritta l'architettura del sistema di rilevamento di spoiler sviluppato, con particolare attenzione alla progettazione e alle scelte che hanno guidato lo sviluppo.
- Nel Capitolo 5 vengono presentati i risultati dell'analisi sperimentale condotta per valutare l'efficacia del sistema.
- Nel Capitolo 6 vengono riassunti i risultati ottenuti e vengono discusse le possibili direzioni future di ricerca.

Capitolo 2

Background

2.1 Approcci al problema

Di seguito vengono presentati alcuni approcci al problema di rilevazione di spoiler in un testo.

2.1.1 Approccio basato su regole

Un approccio molto semplice consiste nell'identificare alcune parole chiave che sono tipicamente presenti in un testo che contiene spoiler. Ad esempio, le parole “spoiler” e “spoiler alert” sono spesso utilizzate per avvertire il lettore che il testo successivo contiene informazioni che potrebbero rovinare la visione di un film o la lettura di un libro. Altri esempi di parole chiave potrebbero essere i nomi di personaggi o luoghi chiave della trama. Questo approccio è molto semplice e può essere implementato con poche righe di codice, ma ha il difetto di essere molto limitato e di non essere in grado di rilevare spoiler più sottili o nascosti.

Questo approccio è inoltre soggetto a molti falsi positivi, ovvero a casi in cui il testo viene erroneamente classificato come contenente spoiler quando in realtà non lo contiene. Ad esempio, una recensione di un film potrebbe contenere il nome di un personaggio chiave della trama senza rivelare informazioni cruciali sulla storia. In questo caso, il testo non dovrebbe essere classificato come contenente spoiler, ma l'approccio basato su regole potrebbe erroneamente classificarlo come contenente spoiler.

Oltre ai falsi positivi, questo approccio è anche soggetto a falsi negativi, ovvero a casi in cui il testo contiene spoiler ma non viene classificato come tale. Ad esempio, una recensione di un film potrebbe contenere informazioni cruciali sulla trama senza utilizzare le parole chiave tipicamente associate

agli spoiler, o in casi anche più semplificati, potrebbe essere scritta in una lingua diversa da quella in cui sono state definite le parole chiave o contenere errori di battitura e di ortografia che impediscono al sistema di riconoscere le parole chiave.

Per questi motivi, l'approccio basato su regole non promette risultati soddisfacenti e si è deciso di non adottarlo in questo lavoro.

2.1.2 Approccio basato sul machine learning

Un approccio più sofisticato e promettente è quello basato sul machine learning. In questo approccio, si addestra un modello di machine learning su un insieme di dati di addestramento etichettati, ovvero un insieme di testi già classificati come contenenti spoiler o non contenenti spoiler. Il modello di machine learning impara a riconoscere i pattern nei dati di addestramento e a classificare i testi in base a tali pattern. Una volta addestrato, il modello può essere utilizzato per classificare nuovi testi come contenenti spoiler o non contenenti spoiler.

Questo approccio ha il vantaggio di essere molto più flessibile e potente rispetto all'approccio basato su regole, in quanto il modello di machine learning è in grado di riconoscere pattern complessi e sottili nei dati e di generalizzare tali pattern ai nuovi dati. Inoltre, il modello di machine learning è in grado di adattarsi autonomamente ai nuovi dati, senza la necessità di modificare manualmente le regole o i parametri del sistema.

Tuttavia, l'approccio basato sul machine learning ha anche alcuni svantaggi. In primo luogo, richiede un insieme di dati di addestramento etichettati, che possono essere costosi e laboriosi da ottenere a seconda del dominio di applicazione.

Il secondo svantaggio è che il modello di machine learning è difficile da realizzare e richiede competenze avanzate per la sua implementazione. Inoltre, il modello deve essere addestrato su un insieme di dati di addestramento rappresentativo e bilanciato, altrimenti potrebbe essere soggetto a *overfitting* o a *underfitting*.

Generalizzando, l'approccio basato sul machine learning è più complesso e richiede più risorse rispetto all'approccio basato su regole, ma promette risultati migliori e più accurati.

2.2 Approccio proposto

Per affrontare il problema della rilevazione di spoiler in un testo, si propone di utilizzare un approccio basato sul machine learning.

L'intuizione alla base di questo approccio è che i LLM (Large Language Model) sono in grado di catturare le relazioni semantiche e sintattiche tra le parole e i concetti all'interno di un testo, e quindi di riconoscere i pattern che caratterizzano i testi contenenti spoiler.

In particolare, si propone di utilizzare un modello di LLM pre-addestrato con lo scopo di **estrarre** lo spoiler da un testo piuttosto che classificare l'intero testo come contenente spoiler o non contenente spoiler, nel caso in cui il testo contenga spoiler.

Utilizzare un modello di LLM pre-addestrato ha il vantaggio di non richiedere un insieme di dati di addestramento etichettati, in quanto il modello è già stato addestrato su un vasto insieme di dati e ha imparato a riconoscere i pattern nei testi in modo automatico e generale.

A prescindere dal modello di LLM utilizzato, si propone di utilizzare un approccio basato su **RAG** (*Retrieval-Augmented Generation*)[7] per arricchire il testo con informazioni realtive al contesto in cui è stato scritto. Ciò dovrebbe migliorare la qualità delle predizioni del modello di LLM, in quanto il contesto può influenzare il significato delle parole e dei concetti all'interno del testo.

La difficoltà principale di questo approccio è che i modelli di LLM sono molto complessi e richiedono risorse computazionali e di memoria considerevoli per essere addestrati e utilizzati. Inoltre, i modelli di LLM sono soggetti a fenomeni di allucinazione e di bias, che possono influenzare le predizioni del modello e ridurre l'accuratezza.

Tuttavia, nonostante queste difficoltà, gli LLM si sono dimostrati molto efficaci in una varietà di task di NLP (Natural Language Processing) e promettono risultati soddisfacenti per il problema della rilevazione di spoiler in un testo.

I dettagli dell'approccio proposto verranno discussi nei seguenti capitoli.

Capitolo 3

Lavori correlati

3.1 Contesto della Ricerca

In questo capitolo verranno presentati alcuni lavori correlati all’oggetto di questo lavoro. In particolare saranno introdotti alcuni concetti chiave e verranno discussi in modo da poter contestualizzare meglio il lavoro svolto.

3.2 LLM, modelli di linguaggio

I Large Language Models (LLM) rappresentano una classe di modelli di deep learning che hanno dimostrato capacità straordinarie nel comprendere e generare testo in linguaggio naturale. Questi modelli sono addestrati su enormi quantità di dati testuali e sono in grado di apprendere complesse relazioni semantiche e sintattiche tra le parole.

Nel contesto del rilevamento di spoiler, gli LLM offrono un potenziale significativo grazie alla loro capacità di comprendere il contesto e il significato delle parole. A differenza dei modelli basati su regole o su approcci di machine learning tradizionali, gli LLM potrebbero catturare sfumature semantiche e identificare spoiler anche in testi complessi e ambigui.

I primi modelli di LLM utilizzavano una architettura del tipo *encoder-decoder* per generare testo in linguaggio naturale. L’encoder è responsabile di trasformare il testo in input in una rappresentazione vettoriale, mentre il decoder è responsabile di generare il testo in output.

Con l’introduzione del **meccanismo di attenzione**, i modelli LLM sono diventati sempre più potenti e sofisticati. Il meccanismo di attenzione permette al modello di estrarre il contesto delle parole nel testo, migliorando notevolmente le performance in una varietà di task.

L'attenzione è stata successivamente estesa a modelli *attention-based*, come il Transformer.

3.2.1 Attention is all you need

Il Transformer è un modello di deep learning *attention-based* introdotto da Vaswani et al. nel 2017 nel paper “Attention is all you need” [17]. Il Transformer ha rivoluzionato il campo del deep learning per il linguaggio naturale, superando i modelli *encoder-decoder* in termini di performance e efficienza [14].

Il meccanismo di attenzione è il cuore del Transformer. Questo meccanismo permette al modello di assegnare pesi diversi alle parole nel testo in input, in modo da poter focalizzare l'attenzione sulle parole più rilevanti per il testo in input.

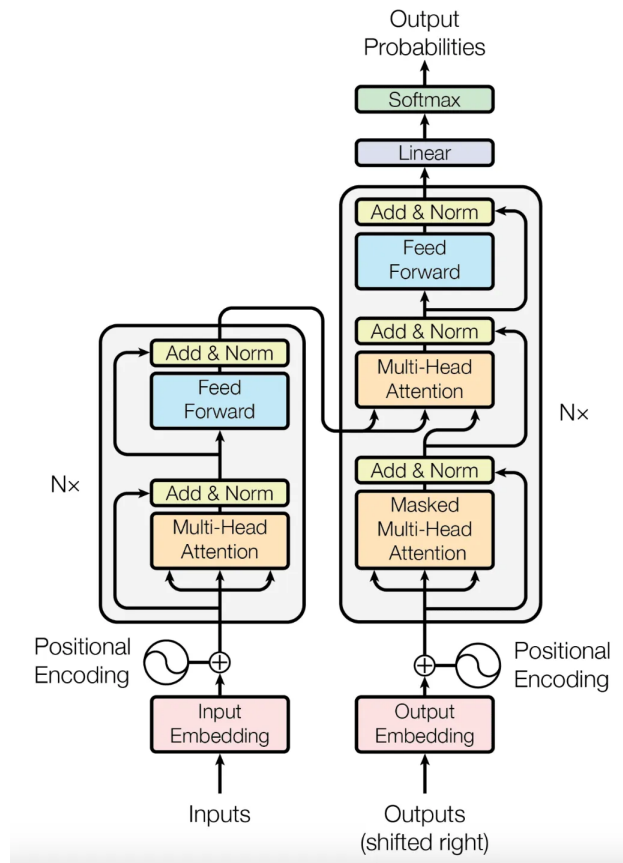


Figura 3.1: Architettura del Transformer

3.3 Embedding

Gli embedding sono una rappresentazione vettoriale di un oggetto in uno spazio vettoriale. Sono ampiamente utilizzati nel campo del deep learning per rappresentare parole, frasi, documenti e immagini in modo da poter essere usati come input per modelli di machine learning [9]. Gli embedding vengono creati utilizzando modelli che imparano a *mappare* gli oggetti in uno spazio vettoriale in modo che oggetti simili siano vicini tra loro. Tra i modelli di embedding più comuni troviamo **Word2Vec**, **GloVe** e **FastText**.

Alcune delle applicazioni per gli embedding includono:

- **Elaborazione del linguaggio naturale:** embedding di parole e frasi per modelli di classificazione, clustering e generazione di testo. Consentono al modello di comprendere il significato delle parole.
- **Sistemi di raccomandazione:** utilizzare gli embedding per rappresentare utenti, prodotti, post in uno spazio vettoriale, in modo da poterne calcolare la similarità.
- **Ricerca di immagini:** permettono di comparare immagini tra di loro o di comparare immagini e testo. Questo permette di cercare immagini tramite testo o viceversa.
- **Ricerca semantica:** rendono possibile la ricerca di frasi o documenti simili in base al loro significato.

Eseguire operazioni su embedding è molto più efficiente se vengono memorizzati in un database apposito.

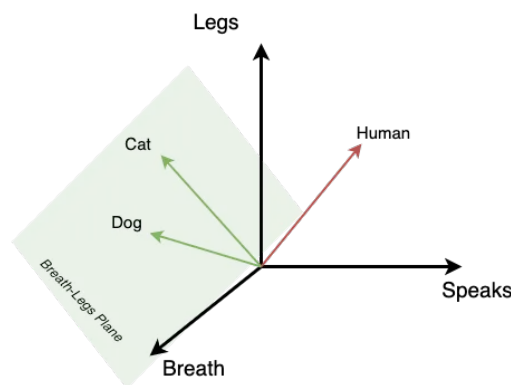


Figura 3.2: Esempio di embedding. Fonte: [3]

3.3.1 Vector Database

A differenza dei database tradizionali, i vector database sono progettati per memorizzare e interrogare dati in forma vettoriale. Questo permette di effettuare operazioni e manipolazioni su dati vettoriali in maniera più efficiente rispetto ad un database relazionale.

I vector database sono particolarmente utili per applicazioni che richiedono il calcolo della similarità tra oggetti in uno spazio vettoriale, come le applicazioni menzionate in precedenza.

Recentemente sono stati introdotti diversi vector database come **Pinecone**, **Chroma** e **MongoDB Atlas**, ma per questo lavoro è stato scelto **PostgreSQL** con il modulo *PGVector* per semplicità di utilizzo e familiarità.

3.4 Fine-tuning di modelli

Il finetuning è una tecnica comune nel campo del deep learning per adattare un modello ad un task specifico.

Nel contesto del rilevamento di spoiler, il finetuning di un modello può essere utile per migliorare le performance del modello nel rilevare spoiler in un determinato contesto o linguaggio.

Il finetuning di un modello pre addestrato comporta l'addestramento del modello su un dataset specifico per un numero limitato di epoche. Questa tecnica consente di sfruttare le conoscenze pregresse del modello, evitando di investire tempo e risorse per addestrare un nuovo modello da zero.

Durante il finetuning, i pesi del modello vengono aggiornati utilizzando un tasso di apprendimento più basso rispetto all'addestramento iniziale. Questo permette al modello di adattarsi al nuovo dataset senza peggiorare le prestazioni su dati già visti.

Ciò risulta particolarmente utile quando si dispone di un dataset limitato o quando si desidera adattare un modello pre-addestrato a un dominio specifico.

3.5 Transizione all'implementazione

Avendo stabilito il contesto teorico e tecnologico nel capitolo corrente, il prossimo passo consiste nel descrivere come questi concetti sono stati applicati nella pratica. Nel prossimo capitolo ci focalizzeremo sugli aspetti pratici della realizzazione del sistema. Descrivendo in dettaglio il design e l'implementazione del sistema di rilevamento di spoiler, illustrando le sfide incontrate e le soluzioni adottate per creare un sistema efficace e utilizzabile.

Capitolo 4

Design e implementazione

4.1 Introduzione

Inizieremo con una panoramica generale dell'architettura del sistema, per poi passare a una descrizione dettagliata dei componenti principali.

4.2 Architettura del sistema

L'architettura del sistema è composta dalle seguenti componenti principali:

- **Frontend:** consiste in una estensione per il browser che permette di interagire con il sistema.
- **Backend:** è il server che gestisce le richieste provenienti dal frontend e comunica con il database e il LLM.
- **Embedder:** è il componente che genera gli embedding delle parole chiave che popolano il database, nonché gli embedding delle richieste provenienti dal frontend.
- **Database:** memorizza le informazioni relative agli embedding delle parole chiave che fungono da input al LLM. Fornisce inoltre, il contesto aggiuntivo utile a generare le risposte.
- **LLM:** è il modello che genera le risposte in base alle richieste del frontend.

mette la creazione di API RESTful, per applicazioni web scalabili e performanti. Il backend è progettato per essere modulare e facilmente configurabile, in modo da poter aggiungere nuove funzionalità in futuro.

Quando il backend riceve una richiesta dal frontend, esegue le seguenti operazioni:

- Preprocessa il contenuto testuale della richiesta
- Tokenizza il testo dividendo il testo in frasi più piccole.
Questo passaggio è fondamentale per creare il contesto necessario per il LLM.
- Genera gli embedding di ogni frase attraverso l'embedder.
- Esegue una ricerca nel database per trovare le parole chiave più simili agli embedding generati
- Invia la richiesta al LLM per generare la risposta. La richiesta al LLM include il contenuto processato e le parole chiave trovate nel database.
- Restituisce la risposta del LLM al frontend sotto forma di JSON.

La comunicazione tra il server e il LLM avviene tramite la libreria di OpenAI, che consente di inviare richieste al modello e ricevere le risposte utilizzando una interfaccia configurabile. In particolare, la libreria permette di utilizzare qualsiasi modello compatibile con l'interfaccia di OpenAI e questo ha permesso di testare diversi modelli attraverso il tool **Ollama**.

4.3.3 Ollama

Ollama è un tool open source che consente di eseguire LLM localmente, senza la necessità di una connessione a Internet. Ollama stato è progettato per essere semplice da usare e consente di eseguire LLM in modo rapido e efficiente sfruttando l'accelerazione hardware delle GPU quando possibile. Ollama offre una varietà di modelli LLM pre-addestrati, tra cui **Llama**[16], **Mistral**[6], **Phi**[1], **Deepseek**[4] e **Gemma**[15].

In particolare, **Gemma** è il modello utilizzato per testare il sistema, in quanto è ottimizzato per l'uso in locale e offre prestazioni eccellenti rispetto a modelli di dimensioni comparabili[15].

4.3.4 Embedder

L'embedder è il componente responsabile della generazione degli embedding delle parole chiave che popolano il database, nonché degli embedding delle richieste del frontend. L'embedder è implementato in Python e utilizza la libreria `SentenceTransformers`[13] per generare gli embedding delle parole chiave e delle richieste.

Attraverso l'uso di `SentenceTransformers` è stato possibile eseguire il fine-tuning del modello di embedding `all-MiniLM-L6-v2` per ottimizzare le prestazioni del sistema. Il fine-tuning del modello è stato eseguito utilizzando un dataset autonomamente creato, tramite l'uso di un web scraper che ha estratto le informazioni da fonti pubbliche disponibili su Internet.

Il dataset è composto da una serie di articoli e documenti che sono stati manipolati per estrarre il contenuto testuale e le parole chiave.

4.4 Dataset e tool utilizzati

Qui di seguito sono riportati i tool e le librerie utilizzati per la creazione del dataset e il fine-tuning del modello di embedding:

- **Puppeteer**: è una libreria Node.js che consente di controllare un browser Chrome o Chromium in modo programmatico.
- **Pandas**: è una libreria Python per la manipolazione e l'analisi dei dati.
- **Visualizer**: un tool creato appositamente per visualizzare i risultati dello scraping, utilizzato per verificare la correttezza dei dati estratti e studiare le relazioni tra le parole chiave e i documenti.

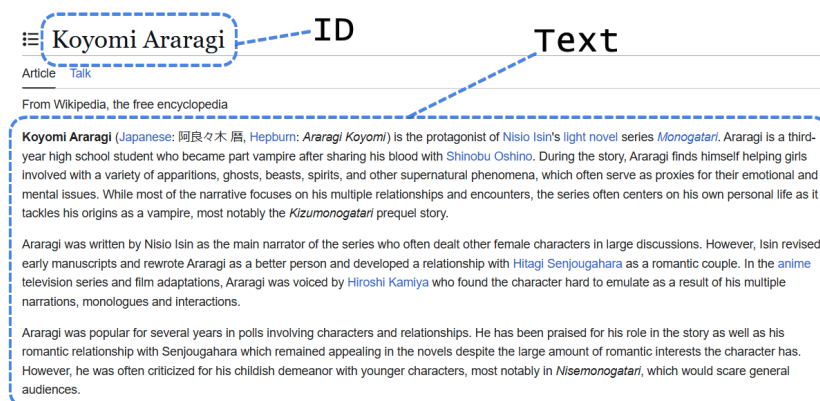


Figura 4.2: Esempio di scraping. Fonte: Wikipedia

In figura 4.2 è mostrato un esempio di scraping effettuato attraverso il tool, che ha estratto il contenuto testuale di una pagina Wikipedia.

Notiamo che il tool è in grado di estrarre non solo il contenuto testuale, ma anche i collegamenti ad altre pagine presenti nei singoli paragrafi, unitamente a una serie di metadati. Ciò ha reso possibile la creazione di un grafo delle relazioni tra le parole chiave e i documenti, che è stato utilizzato per il fine-tuning del modello di embedding. Il grafo è stato creato utilizzando la libreria `D3.js`, che consente di creare visualizzazioni interattive dei dati.

4.4.1 Visualizer

Il visualizer è stato creato appositamente per studiare le relazioni tra le parole chiave e i documenti, in modo da ottimizzare il fine-tuning del modello di embedding creando un dataset bilanciato e di qualità superiore.



Figura 4.3: Esempio di grafo creato dal visualizer

In figura 4.3 è mostrato un esempio di grafo. Ogni nodo rappresenta un documento e ogni arco rappresenta una relazione tra due documenti.

Possiamo osservare che alcuni documenti sono più connessi di altri, il che indica che hanno diversi collegamenti, o che diversi documenti fanno riferimento a loro. E' anche possibile osservare che alcuni documenti sono isolati, il che indica che non hanno collegamenti con altri documenti.

Attraverso il visualizer è stato possibile filtrare i documenti in base a diversi criteri, come ad esempio il numero di collegamenti, metadati, o altri tipi di filtro.

Nella scelta dei documenti da includere nel dataset si è tenuto conto di diversi fattori statistici [12] [10] [11], elencati di seguito:

Statistica	Valore
Coefficiente di assortatività	-0.4677
Mescolanza di assortatività discreta	-0.0059
Diametro del grafo	9
Densità del grafo	0.024
Lunghezza media del percorso	3.85
Nodi isolati	76

Tabella 4.1: Statistiche globali del grafo. Per la descrizione delle metriche, si veda la Sezione 4.4.1.

Statistica	Valore
Coefficiente di Clustering	0.333
Componenti nel Cluster	258
Centralità di Closeness	0.365

Tabella 4.2: Statistiche per singoli nodi. Per la descrizione delle metriche, si veda la Sezione 4.4.1.

Nota: I valori riportati nelle tabelle 4.1 e 4.2 sono stati generati su un dataset di esempio per illustrare le funzionalità del visualizer.

Metriche del grafo

Le seguenti metriche sono state utilizzate per valutare le caratteristiche del grafo:

- **Coefficiente di assortatività:** Misura la tendenza dei nodi a connettersi con altri nodi simili. Un valore negativo indica che i nodi tendono a connettersi con nodi dissimili.
- **Mescolanza di assortatività discreta:** Variazione del coefficiente di assortatività per attributi discreti dei nodi.
- **Diametro del grafo:** La massima distanza tra due nodi qualsiasi nel grafo.
- **Densità del grafo:** Rapporto tra il numero di archi presenti e il numero massimo possibile di archi.

- **Lunghezza media del percorso:** La distanza media tra tutte le coppie di nodi nel grafo.
- **Nodi isolati:** Numero di nodi che non sono connessi ad altri nodi.
- **Coefficiente di Clustering:** Misura quanto i vicini di un nodo sono connessi tra loro.
- **Componenti nel Cluster:** Numero di sottografi connessi nel grafo.
- **Centralità di Closeness:** Misura quanto un nodo è vicino a tutti gli altri nodi nel grafo.

4.5 Preparazione del dataset

Per preparare il dataset è stato necessario eseguire diverse operazioni di post-processing sui dati estratti dallo scraping.

In particolare, sono stati eseguiti i seguenti passaggi:

- Rimozione dei metadati
- Normalizzazione del testo
- Oversampling di alcune parole chiave
- Introduzione di rumore nel dataset
- altre operazioni (disponibili nel codice sorgente[2])

Si è deciso di eseguire l'oversampling di alcuni tipi di parole chiave, in modo da bilanciare il dataset e ottenere un modello di embedding più robusto.

Un esempio di oversampling è mostrato in figura 4.4.



Figura 4.4: Esempio di oversampling

Il risultato finale del dataset è stato un file JSON con il seguente formato:

anchor	positive	negative
koyomi araragi	vampiro, studente, adolescente	nobile
oshino shinobu	vampiro, 500 anni, nobile	studente
ononoki	sissi, yay. peace peace! 🙌👁️👁️🙌	hatsune miku
bakemonogatari	anime, manga, light novel	bad

Tabella 4.3: Formato del dataset

Nota: I dati riportati nella tabella 4.3 sono stati generati su un dataset di esempio per illustrare il formato del dataset.

In particolare, la colonna **anchor** contiene la parola chiave, la colonna **positive** contiene il contesto che si desidera associare all'ancora, mentre la colonna **negative** contiene l'elemento che si desidera evitare di associare alla parola chiave.

Tramite l'uso di questo formato è stato possibile separare efficacemente le parole chiave e i contesti, in modo da ottenere un dataset qualitativamente migliore.

4.5.1 Scelta dei negativi

La scelta degli esempi negativi è stata cruciale per il fine-tuning del modello di embedding.

Esistono principalmente due approcci per la scelta degli esempi negativi:

- **Scelta casuale:** consiste nel selezionare casualmente gli esempi negativi dal dataset.
- **Scelta non casuale:** consiste nel selezionare gli esempi negativi in base a regole stabilite.

Nel primo approccio, la scelta degli esempi negativi è completamente casuale e non tiene conto delle relazioni tra le parole chiave e i contesti. Questo approccio può portare a risultati non ottimali, in quanto gli esempi negativi potrebbero non essere rappresentativi del contesto in cui si desidera utilizzare il modello di embedding.

La strategia adottata nel progetto è stata quella di estrarre randomicamente il contesto a partire dalle parole chiave “vicine” all’ancora di riferimento. Questo approccio tiene conto delle relazioni tra i diversi documenti in modo da ottenere degli embedding più robusti e rappresentativi.

Nel nostro caso, a partire da un nodo del grafo, l’ancora è rappresentata dal titolo del nodo, l’esempio positivo è rappresentato dal contenuto del nodo, mentre l’esempio negativo è rappresentato dal contenuto di un nodo adiacente.

Questo approccio ha portato a risultati migliori rispetto all’approccio casuale poichè sfrutta il fatto che i documenti che hanno collegamenti in comune, molto probabilmente condividono informazioni simili. Permettendo al modello di distinguere tra i due, è stato possibile ottenere un embedding più robusto e rappresentativo.

4.6 Fine-tuning del modello

Il fine-tuning del modello di embedding è stato eseguito utilizzando la libreria `SentenceTransformers` e il dataset creato in precedenza su molteplici modelli, tra cui `all-MiniLM-L6-v2`, per un numero di epoche pari a 4.

La *loss function* utilizzata è stata la `TripletLoss`[5], comunemente utilizzata per l’addestramento di modelli di embedding. La `TripletLoss` è progettata per minimizzare la distanza tra l’ancora e l’esempio positivo, e massimizzare la distanza tra l’ancora e l’esempio negativo.

Come discuteremo nel prossimo capitolo, il fine-tuning del modello ha permesso di migliorare il modello di embedding rendendolo più adatto al dominio in questione.

Capitolo 5

Valutazione

5.1 Introduzione

In questo capitolo, verranno presentati i risultati ottenuti durante la fase di valutazione del nostro sistema.

Iniziando con una panoramica dei dati utilizzati per fine-tuning e valutazione del modello, seguita da una descrizione delle metriche di valutazione impiegate.

Successivamente, presenteremo i risultati ottenuti, analizzando le prestazioni del nostro sistema in termini di accuratezza e capacità di rilevamento degli spoiler.

Verranno discusse le implicazioni dei risultati e le possibili direzioni future per il miglioramento del sistema.

5.1.1 Specifiche tecniche

La configurazione del sistema è stata eseguita su un computer con le seguenti specifiche:

- CPU: AMD Ryzen 7 5800H (8 core, 16 thread)
- GPU: NVIDIA GeForce RTX 3070 Laptop (8 GB)
- RAM: 16 GB
- Sistema operativo: Windows 10 22H2
- Versione di Python: 3.12.9

5.2 Dataset e metodologia

Per la valutazione del nostro sistema, abbiamo utilizzato 2 strategie: **pair** e **triplet**. Il dataset **pair** è derivato da **triplet** escludendo la colonna *negative*. Questo ha permesso di valutare il sistema in diverse configurazioni.

Il codice per la generazione degli split e il training è identico per entrambe le strategie, con l'unica differenza che per **pair**, il dataset è stato ridotto a 2 colonne: *positive* e *negative*, adattando la funzione di loss di conseguenza. Il seed utilizzato per la generazione degli split è 298[8], garantendo la riproducibilità dei risultati.

Parametro	Valore
Dataset	4067 righe
Split	80% train, 20% test
Epoche	4
Loss (pair)	CosineSimilarityLoss
Loss (triplet)	TripletLoss

Tabella 5.1: Panoramica dei parametri utilizzati per il fine-tuning e la valutazione del modello.

5.3 Embedders

Durante la creazione del sistema, si è deciso di valutare due modelli pre-addestrati diversi tra loro per generare gli embedding, quali:

- **all-MiniLM-L6-v2**
- **all-mpnet-base-v2**

Nessuno dei modelli è stato addestrato su dati specifici al dominio degli spoiler o sui dati utilizzati per il fine-tuning. Come mostrano i risultati di seguito, nessuno dei modelli soddisfaceva i requisiti per il task dato il dominio specifico, ma con un fine-tuning adeguato, sono stati in grado di generare embeddings soddisfacenti.

Entrambi i modelli sono disponibili su HuggingFace, installando la libreria **sentence-transformers**. I modelli scelti sono comunemente utilizzati per la creazione di embeddings per frasi o paragrafi brevi, rendendoli adatti al nostro task. Il codice per il fine-tuning e la valutazione è stato scritto in Python ed è disponibile su GitHub[2].

5.3.1 Risultati

Nota: I migliori risultati sono segnalati con un asterisco (*)

Pair

Modello	Similarità positiva media
Pre-addestramento	
all-mpnet-base-v2	0.1554
all-MiniLM-L6-v2*	0.1688
Fine-tuned	
all-mpnet-base-v2*	0.9975
all-MiniLM-L6-v2	0.9938

Tabella 5.2: Risultati della valutazione dei modelli di embedding (Pair)

Triplet

Modello	Similarità pos. media	Similarità neg. media
Pre-addestramento		
all-mpnet-base-v2	0.1554	0.2001
all-MiniLM-L6-v2*	0.1688	0.2113
Fine-tuned		
all-mpnet-base-v2*	0.6358	-0.8118
all-MiniLM-L6-v2	0.6274	-0.7221

Tabella 5.3: Risultati della valutazione dei modelli di embedding (Triplet)

Notiamo che i risultati pre-addestramento sono molto bassi, e simili tra loro, il che indica che i modelli non sanno categorizzare bene i dati in input. Tuttavia, dopo il fine-tuning, i risultati sono significativamente migliorati, con all-mpnet-base-v2 che ottiene i migliori risultati in entrambe le configurazioni.

Nota: I risultati nella configurazione **pair** sono migliori rispetto a quelli nella configurazione **triplet**, ma ciò indica solo che i modelli hanno appreso a generare embedding più simili tra loro, **non necessariamente più accurati**.

5.3.2 Considerazioni

In definitiva, i risultati mostrano che il fine-tuning porta a un miglioramento significativo delle prestazioni, specialmente nella configurazione **triplet**, dove i modelli sono stati in grado di generare embedding con una similarità positiva superiore a 0.64 e una similarità negativa inferiore a -0.81 . Ciò indica che i modelli sono stati in grado di raggruppare correttamente i dati simili e separare quelli dissimili, dimostrando la loro efficacia in questo task dopo un adeguato fine-tuning.

Nonostante all-mpnet-base-v2 abbia ottenuto, di poco, risultati migliori rispetto a all-MiniLM-L6-v2, quest'ultimo ha dimostrato di essere più veloce e più leggero, rendendolo la scelta preferita per questo lavoro. Inoltre, all-MiniLM-L6-v2 ha mostrato prestazioni migliori in fase di fine-tuning.

5.3.3 Prestazioni in fase di fine-tuning

Modello	Tempo (min:sec)	Step	Perdita
all-mpnet-base-v2	03:26	500	3.8035
all-MiniLM-L6-v2*	00:53	500	4.0104

Tabella 5.4: Tempi e perdite del fine-tuning (Triplet)

Modello	Tempo (min:sec)	Step	Perdita
all-mpnet-base-v2	02:10	500	0.0183
all-MiniLM-L6-v2	00:37	500	0.0263

Tabella 5.5: Tempi e perdite del fine-tuning (Pair)

Le prestazioni di all-MiniLM-L6-v2 sono notevolmente migliori rispetto a quelle di all-mpnet-base-v2 in termini di tempo. Il tempo di inferenza è un aspetto critico per l'implementazione di modelli che devono essere utilizzati in locale, e questo modello ha dimostrato di essere la scelta più adatta per il nostro scopo.

5.4 LLM

Per la valutazione del modello LLM, è stato scelto **Gemma3**.

5.4.1 Gemma3

Gemma3 è un modello LLM open-weight sviluppato da Google, progettato per essere altamente efficiente e performante. Gemma3 ha capacità multimodali, il che significa che può elaborare testo e immagini.

Data la sua architettura avanzata derivante da Gemini, lo stato dell'arte dei modelli di Google, Gemma3 è in grado di comprendere e generare contenuti con performance competitive rispetto a modelli di dimensioni molto maggiori[15].

Per questo progetto, non c'è stato bisogno di utilizzare le capacità multimodali, bensì ci siamo concentrati sulla generazione di testo.

Si è deciso di utilizzare Gemma3 per la sua capacità di generare testo di alta qualità e la sua architettura leggera che consente un'implementazione rapida ed efficiente.

A causa di limitazioni hardware, non è stato possibile scaricare il modello completo, quindi è stato scelto di utilizzare versioni del modello con dimensioni diverse. Gemma3 è disponibile in diverse dimensioni:

- **Gemma3-1B**: 1 miliardo di parametri
- **Gemma3-4B**: 4 miliardi di parametri
- **Gemma3-12B**: 12 miliardi di parametri
- **Gemma3-27B**: 27 miliardi di parametri

La versione 1B non possiede la capacità multimodali, mentre la versione 12B ha una dimensione $> 8GB$, per cui non è stato possibile utilizzare appieno il modello sfruttando l'accelerazione GPU.

5.4.2 Risultati LLM

Per valutare le prestazioni dei diversi LLM, è stato utilizzato un dataset diverso rispetto a quello utilizzato per il fine-tuning.

Il dataset è composto da 100 frasi, ognuna delle quali è stata generata da un utente diverso. Le frasi possono contenere spoiler o meno, e in caso di spoiler, viene marcato l'estratto che contiene lo spoiler.

Il dataset è stato procurato tramite scraping del sito **TVTropes**, un sito dedicato alla discussione di opere di narrativa e alla loro analisi.

Il dataset è stato bilanciato in modo da avere un numero uguale di frasi con e senza spoiler. Le frasi sono state estratte in modo casuale dalla comunità dedicata all'opera **Monogatari Series**, che è anche l'oggetto del dataset utilizzato per il fine-tuning del modello di embedding.

Si è scelto di mantenere il preprocessing del dataset al minimo, poiché le frasi sono state raccolte direttamente dal sito web. Questo approccio permette di preservare la qualità naturale del linguaggio utilizzato dagli utenti reali. Il dataset risultante rappresenta un campione autentico di espressioni, caratterizzato da varietà di stili e dall'assenza di schemi predefiniti.

Questa eterogeneità pone una sfida significativa per il modello, che non è stato addestrato su dati di questa natura e si trova ad affrontare un tipo di input inedito.

La Tabella 5.6 riassume i risultati ottenuti dai vari modelli.

Modello	Accuracy	Precision	Recall	F1-score
Gemma3:1b	0.51	0.505	1.0	0.671
Gemma3:4b	0.58	0.543	1.0	0.704
Gemma3:12b*	0.66	0.603	0.94	0.734
Gemini 2.0-flash	0.59	0.552	0.96	0.701

Tabella 5.6: Risultati delle prestazioni dei modelli.

5.5 Analisi dei risultati

I risultati ottenuti mostrano che la dimensione del modello ha un impatto notevole sulle prestazioni. Si osserva che Gemma3-12B ha ottenuto risultati migliori in termini di accuratezza, precisione e F1-score, rispetto invece a modelli più piccoli, come Gemma3-1B che mostra risultati simili ad un lancio di moneta, con un'accuratezza intorno al 50%. Ricordiamo che:

- Accuracy: Misura la proporzione di previsioni corrette sul totale delle previsioni.
- Precision: Misura quanto sono affidabili le previsioni positive del modello.
- Recall: Misura la capacità del modello di trovare tutti i casi positivi.
- F1-score: Fornisce un equilibrio tra precisione e recall, utile quando si desidera considerare entrambi gli aspetti.

Dato il trend osservato, si è deciso di introdurre un ulteriore modello, **Gemini 2.0-flash**, top di gamma di Google, per confrontare le prestazioni con quelle di Gemma3-12B.

Utilizzando la stessa metodologia di valutazione (con la differenza che Gemini è disponibile solamente attraverso API di Google), Gemini 2.0-flash ha ottenuto risultati inferiori rispetto a Gemma3-12B, nonostante Gemini 2.0-flash sia un modello più grande e avanzato.

I valori sono stati calcolati tramite la libreria `sklearn`, e si basano sulle previsioni del modello sui dati di test.

Importante notare che i risultati sono stati ottenuti su un task di classificazione binaria, ovvero rilevare se una frase contiene uno spoiler o meno, tuttavia, il prompt del sistema richiedeva **anche** di evidenziare lo spoiler, il che potrebbe aver influenzato le prestazioni dei modelli.

Sarebbe interessante esplorare ulteriormente questo aspetto per capire se i modelli sono in grado di estrarre correttamente lo spoiler o soffrano di “allucinazioni”, ma questo richiederebbe un’analisi più approfondita e “soggettiva” dei risultati ottenuti. Per correttezza, si è preferito non includere questo tipo di analisi in modo da non influenzare i risultati finali.

5.6 Conclusioni

In questo capitolo, sono stati presentati i risultati ottenuti durante la fase di valutazione del sistema di rilevamento degli spoiler. Sono state discusse le prestazioni dei modelli di embedding ed LLM, analizzando le loro capacità di categorizzare correttamente i dati in input e di generare embeddings ottimali.

L’analisi ha incluso una valutazione approfondita delle prestazioni dei modelli durante la fase di fine-tuning, sottolineando come questo processo sia cruciale per ottimizzare l’efficacia dei modelli stessi. Inoltre, è stato effettuato un confronto dettagliato tra i modelli Gemma3 e Gemini 2.0-flash, evidenziando le loro peculiarità e capacità nel rilevare spoiler. Questo confronto ha tenuto conto delle diverse dimensioni e architetture dei modelli, nonché della complessità intrinseca del compito di rilevamento degli spoiler.

Capitolo 6

Considerazioni Finali

In conclusione, i risultati ottenuti dimostrano che il fine-tuning e la scelta del modello giusto sono fondamentali per ottenere prestazioni ottimali nel rilevamento degli spoiler, ma anche quanto sia difficile valutare la qualità dei risultati ottenuti.

L'architettura del sistema si è dimostrata sorprendentemente robusta e in grado di fornire risultati di qualità migliore rispetto a quelli ottenuti con modelli più complessi e costosi in termini di risorse computazionali. Questo è un aspetto importante, poiché dimostra che è possibile ottenere risultati di alta qualità senza dipendere da servizi cloud costosi e senza dover affrontare i problemi di privacy e sicurezza associati all'uso di tali servizi.

L'addestramento dei modelli di embedding mostrano che è possibile utilizzare LLM pre-addestrati senza conoscenze specifiche del dominio per ottenere risultati competitivi. In caso di domini specifici, è possibile ottenere risultati migliori allenando un modello di embedding specifico sul dominio, cosa che rende l'approccio più flessibile rispetto ad allenare un modello di linguaggio da zero. Questo dimostra nuovamente che la scelta del modello giusto è fondamentale per ottenere prestazioni ottimali. Tuttavia, è importante notare che la latenza del sistema è un fattore cruciale da considerare, poiché il sistema è progettato per essere utilizzato in tempo reale da utenti finali.

6.1 Lavori Futuri

Lo scopo di questo lavoro è stato quello di dimostrare un approccio alternativo al rilevamento degli spoiler utilizzando modelli generativi avanzati. Questo lavoro è stato quanto di più ambizioso si potesse fare in un tempo limitato e con le risorse disponibili, ma ha comunque portato a risultati inte-

ressante e promettenti. Tuttavia, ci sono ancora molte aree in cui è possibile migliorare ed espandere questo lavoro.

In primo luogo, si è scelto di non ottimizzare al meglio nessuna delle parti del sistema, ma di concentrarsi sull'architettura del sistema in generale. In futuro, sarebbe interessante ottimizzare ulteriormente le prestazioni del sistema, sperimentando con diverse architetture e tecniche di ottimizzazione.

Lavorare a un sistema di rilevamento degli spoiler completamente automatizzato è un obiettivo ambizioso, ma è possibile farlo utilizzando tecnologie disponibili pubblicamente ed open source.

Lascio quindi a chiunque voglia continuare questo lavoro l'idea di esplorare le possibilità offerte dai modelli generativi avanzati con l'obiettivo di migliorare l'esperienza dell'utente finale.

Bibliografia

- [1] Marah Abdin, Jyoti Aneja, Harkirat Behl, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, Michael Harrison, Russell J. Hewett, Mojan Javaheripi, Piero Kauffmann, James R. Lee, Yin Tat Lee, Yuanzhi Li, Weishung Liu, Caio C. T. Mendes, Anh Nguyen, Eric Price, Gustavo de Rosa, Olli Saarikivi, Adil Salim, Shital Shah, Xin Wang, Rachel Ward, Yue Wu, Dingli Yu, Cyril Zhang, and Yi Zhang. Phi-4 technical report. *arXiv preprint arXiv:2412.08905*, 2024. cs.CL.
- [2] AlexFazio64. gaen-thesis. <https://github.com/AlexFazio64/gaen-thesis>, 2025.
- [3] Oleg Borisov. Word embeddings: Intuition behind the vector representation of the words. <https://medium.com/data-science/word-embeddings-intuition-behind-the-vector-representation-of-the-words-7e4eb2410bba/>, 2020.
- [4] DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
- [5] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification, 2017.
- [6] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [7] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *arXiv preprint arXiv:2005.11401*, 2020.

- [8] Hachikuji Mayoi. Nisemonogatari, karen bee, part one, Jan 2012. Trasmissione televisiva.
- [9] Tomas Mikolov. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 3781, 2013.
- [10] M. E. J. Newman. Assortative mixing in networks. *Physical Review Letters*, 89, 2002.
- [11] M. E. J. Newman. Mixing patterns in networks. *Physical Review E*, 67, 2003.
- [12] M. E. J. Newman. *Networks: an introduction*. Oxford University Press, 2010.
- [13] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019.
- [14] Google Research. Transformer: A novel neural network architecture for language understanding. <https://research.google/blog/transformer-a-novel-neural-network-architecture-for-language-understanding/>, 2017.
- [15] Gemma Team. Gemma 3, 2025.
- [16] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017.