

A Project Report
On
SPOILER SHIELD WITH NATURAL LANGUAGE PROCESSING

BY
K SUMEDA REDDY
SE23MAID018

Under the supervision of

Dr. VIVEK KUMAR MISHRA

**SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR
DEGREE OF MASTER OF TECHNOLOGY
AI5401: DISSERTATION PROJECT**



**ÉCOLE CENTRALE
SCHOOL OF ENGINEERING**

**ÉCOLE CENTRALE SCHOOL OF ENGINEERING
MAHINDRA UNIVERSITY
HYDERABAD
(JUNE 2025)**

ACKNOWLEDGMENTS

I am profoundly thankful to my project guide, **Dr. VIVEK KUMAR MISHRA**, for his dedicated mentorship, valuable suggestions, and unwavering support throughout this endeavour. His guidance has been a cornerstone of my project's progress.

I am sincerely thankful to **Prof. Dr. Arun Kumar Pujari, Head of the CSE Department, Mahindra University**, for his visionary leadership and consistent support, which created an ideal environment for academic growth and research exploration.

I also would like to extend my gratitude to the faculty members of **Ecole Centrale School of Engineering, Mahindra University**, for offering the resources and opportunities essential for this project. Finally, I thank my colleagues and fellow students for their insightful discussions and encouragement, which have significantly enriched my learning experience.

K SUMEDA REDDY

SE23MAID018

Master of Technology in Artificial Intelligence

Mahindra University

Ecole Centrale School of Engineering

Mahindra University

Hyderabad

Certificate

This is to certify that the project report entitled "**SPOILER SHIELD WITH NATURAL LANGUAGE PROCESSING**" submitted by **Ms. K SUMEDA REDDY** (Roll No. **SE23MAID018**) in partial fulfillment of the requirements of the course AI5401, Dissertation Project, embodies the work done by her under my supervision and guidance.

Dr. VIVEK KUMAR MISHRA

Associate Professor

Ecole Centrale School of Engineering, Hyderabad.

Date:

ABSTRACT

Social media often serves as a platform for users to discuss key plot points of recently released movies, unintentionally spoiling the experience for others. This project aims to automate spoiler detection using data from Reddit platform and can also adapt to platforms like IMDb and Rotten Tomatoes. The methodology employs contrastive learning to improve spoiler detection by effectively differentiating spoilers from non-spoilers using refined text representations. Using Transformer-based models like BERT, spoiler detection is performed at both the sequence and token levels. Preliminary results highlight the potential of this method in improving accuracy, though challenges such as class imbalance and false positives persist. This work establishes a foundation for developing systems that can automatically filter spoilers while keeping users engaged.

CONTENTS

Title page	1
Acknowledgement	2
Certificate	3
Abstract	4
CHAPTER 1- INTRODUCTION	8-9
1.1 Introduction	8
1.2 Objectives	8
1.3 Significance and Scope	9
1.4 Source of Data, Problem in Existing Systems & Justification	9
CHAPTER 2 - PROBLEM DEFINITION	11-12
2.1 Problem Statement	11
2.2 Mathematical Representation	12
2.3 Scope and Constraints	12
CHAPTER 3 - BACKGROUND AND RELATED WORK	13-15
3.1 Prior Work on Spoiler Detection	13
3.2 Data Collection from Reddit	14
3.3 Foundations of Contrastive Learning	14
3.4 Transformer Models and Contextual Embedding	15
3.5 Limitations in Existing Literature and Motivations for the project	15
CHAPTER 4 - IMPLEMENTATION	16-22
4.1 Overview	16
4.2 Data Collection	17
4.3 Data Preprocessing	18
4.4 Data Visualization	19
4.5 Pairwise Contrastive Training Data Generation	21
4.6 Model Training	21
4.7 Evaluation	22
CHAPTER 5 - RESULTS	23-28
5.1 Dataset Analysis and Preprocessing Impact	23
5.2 Contrastive Pair Generation Analysis	23
5.3 Model Training and Learning Pattern	23
5.4 Evaluation Metrics and Performance Summary	24
5.5 Visualizations and Interpretations	25
5.6 Error Analysis and Challenges	27
5.7 Real-Time Evaluation via Streamlit App	28
CHAPTER 6 - CONCLUSION	30
CHAPTER 7 - REFERENCES	31

LIST OF FIGURES

Figure No.	Figure Title	Page No.
1	Flow Chart of Implementation	17
2	Graph which visualises training progress and model convergence	22
3	Graph that highlights cosine similarity score differences across pair types (positive vs negative)	24
4	Graph of Distribution of Comment Types	25
5	Graph portraying text reduction post-cleaning	25
6	Graph that showing spoiler vocab patterns	26
7	Graph that showing non-spoiler vocab patterns	26
8	Graph of Cosine Similarity Distribution (Positive vs Negative Pairs)	26
9	Graph of Evaluation Summary – Metric Comparison	27
10	Prediction of a Spoiler Comment	28
11	Prediction of a Non-Spoiler Comment	29

LIST OF TABLES

Table No.	Table Title	Page No.
1	The spoiler detection system was developed through the following key activities	16
2	Tools and Frameworks	16
3	Raw Dataset Sample (First 5 Rows)	18
4	Cleaned dataset sample	19
5	Word Count Statistics (Raw vs Cleaned)	19
6	Class Distribution of Comments	20
7	Most Common Spoiler Words (Top 15)	20
8	Most Common Non-Spoiler Words (Top 15)	20
9	Contrastive Pair Examples	21
10	Cosine Score Table	22
11	Final Evaluation Results	22

CHAPTER 1 – INTRODUCTION

1.1 Introduction

In the era of hyper-connectivity, digital platforms like Reddit, IMDb, and Rotten Tomatoes have become central hubs for cinematic discussions. While these platforms foster a rich exchange of opinions and critiques, they inadvertently become breeding grounds for movie spoilers—revelations that can significantly diminish a viewer's intended emotional and narrative experience. The challenge lies not only in identifying spoilers but in doing so within unstructured, user-generated content that varies in tone, length, and context.

The essence of watching a movie lies in the surprise, the emotional buildup, and the unexpected twists that captivate the audience. Spoilers, whether intentional or accidental, strip away this experience, undermining the storytelling efforts of filmmakers. As content sharing increases in volume and velocity across forums, the demand for robust spoiler detection tools becomes increasingly urgent.

This project introduces an intelligent system **Spoiler Shield with Natural Language Processing** which leverages advanced natural language processing (NLP) techniques to detect spoilers in user comments with high semantic sensitivity. Unlike traditional keyword-based systems, this model utilises Transformer-based architectures such as BERT and contrastive learning strategies to encode, compare, and classify user comments in an embedding space. This allows the system to differentiate between spoiler-laden and non-spoiler content even when textual patterns are subtle or ambiguous.

Beyond the technical challenge, spoiler detection also serves a broader purpose: it preserves the filmmaker's vision, sustains the audience's emotional engagement, and maintains the integrity of online movie discussions. Spoiler-aware environments enable open conversations while respecting personal viewing timelines, thus building a more thoughtful and inclusive digital community.

This research aims to set a foundation for automated spoiler filtering systems capable of real-time deployment and domain adaptation. By building on Reddit-sourced data and embedding-level semantic comparison, it demonstrates a scalable approach to content moderation, not only for movies but for any narrative-driven media.

1.2 Objectives

The primary objectives of the project are:

1. Accurate Spoiler Detection

- a. Develop a machine learning model to distinguish between spoiler and non-spoiler comments using contrastive learning.
- b. Implement cosine similarity-based training to learn embedding relationships between comments.

2. High-Quality Text Representation

- a. Utilises SentenceTransformer models (e.g., MiniLM or BERT variants) to generate refined, context-aware text embeddings.
- b. Preprocess and clean real-world Reddit data for optimal model performance.

3. Balanced Dataset Creation

- a. Construct a well-labelled, balanced dataset from Reddit subreddits using spoiler markdown tags.
- b. Normalise and filter the data to address class imbalance and noise.

4. Contrastive Pair Sampling

- a. Create positive (same label) and hard negative (semantically similar, opposite label) training pairs.
- b. Use contrastive loss to pull similar comments closer and push dissimilar ones apart in embedding space.

5. Evaluation and Visualization

- a. Evaluate the model using cosine similarity, precision, recall, and F1-score.
- b. Visualise comment distributions, training losses, and semantic clustering.

1.3 Significance and Scope

The project holds substantial significance in the context of online media consumption:

➤ **Enhanced Viewing Experience:**

Users can navigate platforms like Reddit without the risk of encountering unwanted spoilers, preserving the suspense and enjoyment of media.

➤ **Semantic-Level Detection:**

By using contrastive learning and embedding models, this approach surpasses naive keyword filters, handling nuanced and rephrased spoiler content effectively.

➤ **Real-World Adaptability:**

Although trained on Reddit data, the system architecture can adapt to other platforms such as IMDb, Twitter, or fan forums through transfer learning.

➤ **Foundational Framework for Content Moderation:**

This project lays the groundwork for real-time spoiler detection tools in browser extensions, moderation systems, or community forums.

➤ **Personalised Tuning:**

While not yet implemented, the model design supports future integration of user-defined sensitivity levels or blacklisted terms.

1.4 Source of Data, Problem in Existing Systems & Justification

Source of Data

The dataset is built using real-time, user-generated Reddit comments scraped from subreddits like:

→ r/movies, r/television, r/marvelstudios, r/MovieDetails, r/moviediscussion

Key criteria:

- Comments tagged with Reddit spoiler markdown `>!spoiler text!<` are labelled as *Spoiler*.
- Remaining comments are filtered for *Non-Spoiler* examples.
- Balanced dataset: 1000 Spoiler & 1000 Non-Spoiler comments.

Problems in Existing Systems

1. Lack of Semantic Understanding

Traditional systems rely heavily on keyword matching and fail to catch rephrased or subtle spoilers.

2. High False Positives/Negatives

Shallow classifiers cannot differentiate between opinions and actual plot revelations.

3. Class Imbalance

Spoiler data is less frequently tagged, resulting in skewed datasets and model bias.

4. Inflexibility Across Platforms

Many existing systems are tailored to specific domains and lack adaptability to new formats or forums.

Justification for the Project

- **Novel Approach:** Uses contrastive learning—rarely applied in spoiler detection—to generate robust embeddings.
- **Performance-Oriented:** Incorporates Transformer models known for high contextual accuracy.
- **Community Value:** Improves discussion safety in fan-driven ecosystems like Reddit.
- **Scalability:** Can be expanded to multi-platform deployment including browser extensions, chatbots, and streaming sites.
- **User-Centric:** Paves the way for future integration of personalization and sensitivity tuning.

CHAPTER 2 – PROBLEM DEFINITION

In the age of digital content and hyperactive social media engagement, the risk of encountering unsolicited spoilers has increased significantly. Online platforms like Reddit, Twitter, and forums host extensive discussions about films and TV shows, often involving key plot elements. While some communities implement spoiler tags or formatting conventions (e.g., Reddit's `>!spoiler!<` markdown), many users either fail to use these or inadvertently include spoilers in seemingly harmless commentary. Spoilers can diminish the viewing experience by revealing critical plot points before the viewer has had the chance to experience the content firsthand. Given the volume and velocity of online discussions, manual moderation of spoiler content is neither scalable nor practical. This project aims to develop an automated, data-driven system to detect spoiler content using natural language processing and deep learning.

The following fundamental questions define the scope of the problem addressed by this project:

1. **Can spoiler comments be automatically identified using linguistic and semantic patterns without explicit spoiler tags?**
This involves analysing user-generated comments to detect contextually implied spoilers even in the absence of formatting cues like spoiler markdown.
2. **Can contrastive learning help separate spoiler and non-spoiler comments in a meaningful embedding space?**
The use of contrastive loss and sentence embeddings aims to cluster semantically similar comments (within class) and separate them from dissimilar ones (across class).
3. **Is it possible to create a system that generalizes well across different forms of spoilers, including indirect and subtle references?**
Many spoilers are implied rather than explicit. The system should be robust enough to identify these based on semantic content rather than keyword matching.
4. **Can the model support real-time interaction and prediction in a user-facing application?**
The deployment aspect of the project seeks to demonstrate how users can input a comment and receive a spoiler probability or classification instantly.

2.1 Problem Statement

The spoiler detection problem can be formally defined as a **binary classification problem**. Given a comment x_i from a dataset $X = \{x_1, x_2, \dots, x_n\}$, the objective is to learn a function:

$$f_\theta : X \rightarrow \{0,1\}$$

Where:

- $f_\theta(x_i) = 1$ if the comment is a spoiler,
- $f_\theta(x_i) = 0$ if the comment is a non-spoiler.

However, unlike traditional classification tasks, this project reformulates the learning objective using **contrastive learning**, where the model is not directly optimized to predict a class label, but to learn an embedding space in which:

- Comments from the **same class** (spoiler/spoiler or non-spoiler/non-spoiler) have **high cosine similarity**.
- Comments from **different classes** (spoiler/non-spoiler) have **low cosine similarity**.

This leads to a pairwise learning framework where the model is trained using **triplets or pairs**, rather than individual labelled examples.

2.2 Mathematical Representation

Let:

- x_i and x_j be two input sentences (comments),
- $y_{ij} \in \{0,1\}$ be the binary label indicating whether they belong to the same class,
- $f_\theta(x)$ be the embedding function learned by the Sentence-BERT model.

Then, the similarity function $S(x_i, x_j)$ is given by the **cosine similarity** between their embeddings:

$$S(x_i, x_j) = \cos(f_\theta(x_i), f_\theta(x_j)) = \|f_\theta(x_i)\| \|f_\theta(x_j)\| f_\theta(x_i) \cdot f_\theta(x_j)$$

The training objective is to minimize the **contrastive loss** defined as:

$$L_{\text{cos}} = \sum_{(i,j)} [y_{ij} \cdot (1 - \cos(f_\theta(x_i), f_\theta(x_j)))^2 + (1 - y_{ij}) \cdot \cos(f_\theta(x_i), f_\theta(x_j))^2]$$

Where:

- The loss is small when similar pairs have high cosine similarity,
- The loss is high when dissimilar pairs are incorrectly mapped close in the embedding space.

2.3 Scope and Constraints

Scope:

- Focus is limited to Reddit comments related to movies and television, scraped from subreddits such as r/movies, r/television, and r/marvelstudios.
- The system is trained on English-language text.
- Only binary spoiler classification is considered (no multi-class or severity ranking).

Constraints:

- Spoiler labels are noisy due to reliance on spoiler markdown for ground truth.
- User language can be highly informal, sarcastic, or context-dependent, complicating classification.
- The dataset size is fixed at 2000 examples (1000 per class) for computational feasibility.

CHAPTER 3 – BACKGROUND AND RELATED WORK

Spoiler detection presents a unique challenge at the intersection of natural language processing (NLP), content moderation, and user experience. Spoilers tend to be embedded in narrative form, often without clear indicators, requiring models to understand context, semantics, and intent. This chapter presents a survey of related work in the domain, including foundational concepts in contrastive learning and transformers, followed by an overview of the techniques applied in this project.

3.1 Prior Work on Spoiler Detection

Reference 1: Allen Bao, Marshall Ho and Saarthak Sangamnerkar ,“Spoiler Alert: Using Natural Language Processing to Detect Spoilers in Book Reviews”, 2021.

Several attempts have been made to classify spoilers in user-generated content. One of the more domain-specific efforts was conducted by Allen Bao, Marshall Ho, and Saarthak Sangamnerkar [1], who focused on spoilers in book reviews. Their system used traditional NLP techniques such as TF-IDF, Named Entity Recognition (NER), and dependency parsing to extract candidate sentences and flag them based on proximity to spoiler-indicative language. They demonstrated that certain grammatical structures and named entities tend to cluster around spoilers. However, the system struggled with indirect spoilers and suffered from limited generalization due to the use of handcrafted features and rule-based pipelines.

Additionally, their model lacked adaptability to other forms of unstructured content and relied heavily on deterministic logic. This limits its scalability across dynamic platforms where users vary in syntax, tone, and spoiler expression.

Reference 2: Golbeck, J. (2015). “Detecting Spoilers in Fan Wikis.” HICSS. Golbeck explored spoiler detection in fan-maintained wikis. By using temporal page versioning and user editing behaviour, the system could detect whether a page version contained a spoiler. This was among the first structured attempts at spoiler detection using metadata and heuristics. The approach relied on version history and explicit user labelling, which are not available on dynamic platforms like Reddit. Moreover, it lacked any form of semantic understanding or generalization capability to detect unseen spoiler content.

Furthermore, the dependence on structured metadata makes it unsuitable for real-time applications, and it is inapplicable to platforms where content is unstructured, and context varies significantly across users.

Reference 3: Boyd-Graber, J. et al. (2019). “Spoiler Alert: Machine Learning Approaches to Detecting Spoilers in Online Content.” NAACL. This paper presented a supervised learning framework using bag-of-words and TF-IDF features with traditional classifiers (e.g., SVMs, logistic regression). While results were promising on curated datasets, the system's performance degraded on noisy or slang-rich data. The model failed to capture contextual meaning and was brittle in handling synonyms or implied spoilers, as it did not utilize deep neural representations.

In addition, it lacked the flexibility to differentiate between narrative intent and casual mentions, which reduced its robustness. The absence of contextual embeddings made it vulnerable to semantic drift in evolving online conversations.

Reference 4: Chang, Y., & Huang, Y. (2020). “BERT for Spoiler Detection in Movie Reviews.” IEEE Access.

This study applied BERT fine-tuning to detect spoilers in IMDb-style movie reviews. The use of contextualized embeddings improved classification performance significantly over traditional models. The model was trained on domain-specific, well-structured reviews, which limited its generalizability. Moreover, the study used a binary classifier without exploring pairwise relationships between spoiler and non-spoiler content, which could enhance robustness.

Also, the model was sensitive to dataset domain boundaries and exhibited performance loss when applied to informal or slang-based spoiler content. It did not explore embedding space structure for separation, which is critical for nuanced tasks like spoiler detection.

Reference 5: Reimers, N., & Gurevych, I. (2019). “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks.” EMNLP.

This work introduced Sentence-BERT (SBERT), a modification of the BERT architecture for producing semantically meaningful sentence embeddings using Siamese and triplet networks. SBERT was foundational in enabling sentence similarity tasks with transformer models. While effective for similarity and clustering, SBERT’s default training does not target class separability. Hence, without fine-tuning using contrastive objectives, its performance on classification tasks like spoiler detection remains suboptimal.

Nonetheless, SBERT provided the architectural basis for contrastive training, enabling cosine-based similarity computations in real time. It established a scalable baseline for learning dense, semantically rich embeddings suitable for downstream applications like this project.

3.2 Data Collection from Reddit

The current project uses Reddit as its primary data source. Reddit’s discussion threads contain frequent and spontaneous mentions of spoilers, often explicitly marked using the `>!spoiler!<` syntax. Data was collected using **PRAW (Python Reddit API Wrapper)**, an asynchronous Python library that allows efficient and structured access to Reddit threads [2].

The use of PRAW ensures:

- Reliable access to relevant subreddits.
- Filtering of spoiler-tagged content for accurate dataset construction.
- Extraction of comments with metadata such as subreddit, post title, and timestamps.

Challenge Addressed: This methodology enables scalable and reproducible data collection, which prior studies often lacked.

3.3 Foundations of Contrastive Learning

Contrastive learning is a self-supervised learning technique that seeks to map semantically similar items closer in embedding space, while pushing dissimilar items farther apart. This framework has

proven highly effective for representation learning in vision and text. According to **Rohit Kundu [3]**, the success of contrastive learning lies in its ability to learn representations without relying on absolute labels. Instead, the model learns relative similarities and differences.

In the context of this project, contrastive learning is implemented using Sentence-BERT and cosine similarity loss. Training samples are structured as:

- **Positive pairs:** comments from the same class (spoiler or non-spoiler).
- **Negative pairs:** one spoiler and one non-spoiler comment.

This approach produces a well-structured embedding space where the model learns to associate similar sentences regardless of lexical variation.

Advantage Over Supervised Learning: Contrastive learning minimizes overfitting to surface-level lexical cues and focuses on semantic relationships, improving generalization on unseen data.

3.4 Transformer Models and Contextual Embedding

The underlying model used in this project is a transformer-based architecture. Transformers have become the state-of-the-art in NLP tasks due to their ability to model long-range dependencies via self-attention. The transformer framework was introduced by **Vaswani et al. [4]** in the landmark paper “Attention Is All You Need.” The architecture eliminates recurrence and convolutions, relying solely on attention mechanisms to capture context in sequences.

The project utilizes **DistilBERT**, a lighter, faster version of BERT that retains 97% of performance while being 40% smaller. Combined with Sentence-BERT’s pooling and Siamese structure, it produces fixed-length sentence embeddings that are:

- Context-aware
- Semantically meaningful
- Optimized for similarity computation

Technical Impact: This enables efficient and scalable cosine similarity comparisons during inference, crucial for both real-time prediction and batch evaluation.

3.5 Limitations in Existing Literature and Motivations for the project

While previous spoiler detection studies such as [1] offer domain-specific insights, they often rely on fixed lexical cues and limited datasets. There is a clear gap in:

- Utilizing deep contextual models for spoiler understanding.
- Training with semantic-level contrast rather than surface-level labels.
- Scaling across informal, user-generated platforms like Reddit.

This project addresses these gaps by integrating the following:

- A large-scale, real-world dataset from Reddit.
- A contrastive learning framework for better separation in semantic space.
- Transformer-based contextual embeddings for robustness.

CHAPTER 4 - IMPLEMENTATION

This section outlines the complete implementation pipeline developed to construct a spoiler detection model using Reddit comment data. The approach is divided into five major stages: **data collection**, **preprocessing**, **contrastive pair generation**, **model training**, and **embedding analysis**. The process integrates natural language processing (NLP), contrastive learning, and modern transformer-based architectures.

4.1 Overview

The implementation of the project follows a sequential, modular process from data collection to model evaluation. Python 3.10 served as the core programming language, supported by libraries like `asyncpraw`, `nltk`, `sentence-transformers`, `pandas`, and `matplotlib`. Each phase is executed in a Jupyter Notebook to visualize and validate intermediate steps, with output tables and plots used for debugging and analysis.

Table 1: The spoiler detection system was developed through the following key activities

STAGE	ACTIVITY
Data Collection	Scraped 2000 Reddit comments from multiple subreddits
Preprocessing	Cleaned and normalized text; removed spoiler tags, symbols, stopwords
Pair Generation	Created contrastive pairs (positive and negative) using labels
Model Training	Trained Sentence-BERT using CosineSimilarityLoss
Embedding Analysis	Computed sentence embeddings and anchor vectors for each class
Test & Validation	Evaluated using cosine similarity on comment pairs and classification metrics

Table 2: Tools and Frameworks

Tool	Purpose
Python	Primary programming language
<code>asyncpraw</code>	Asynchronous Reddit API wrapper
<code>pandas</code> , <code>numpy</code>	Data handling and transformation
<code>nltk</code>	Text cleaning and stopword removal
<code>SentenceTransformers</code>	Model architecture and training
<code>PyTorch</code>	Backend deep learning framework
<code>Streamlit + pyngrok</code>	Deployment (interactive demo app)

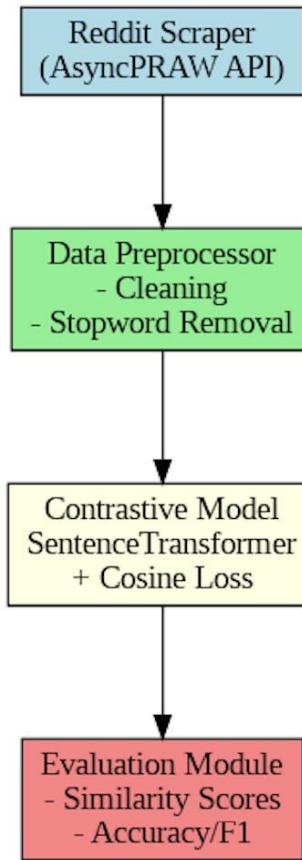


Figure 1: Flow Chart of Implementation

4.2 Data Collection

Data was sourced from Reddit, leveraging the `asyncpraw` library to asynchronously scrape comments from movie-related subreddits:

1. r/movies
2. r/television
3. r/marvelstudios
4. r/FilmTheories
5. r/MovieDetails

The system searches for posts containing the keyword "spoiler" and identifies spoiler comments using the Reddit spoiler markdown format: `>!spoiler text!<`.

- Max Spoiler Comments: 1000
- Max Non-Spoiler Comments: 1000
- Output stored as: `spoiler_shield_dataset.csv`.

Table 3: Raw Dataset Sample (First 5 Rows)

S.No.	Movie	Comment	Comment Type
0	Official Discussion - Sinners [SPOILERS]	Agreed. It could just be me but my take is that >!Stack smelled death on Sammy and wanted to give him a last chance to live. After Sammy refused and before he left he asked Sammy to play one last time because I think he wanted to see Smoke again before Sammy died. Its the whole reason why Remick wanted Sammy in the first place to "reconnect him with his loved ones"!< through his music.	Spoiler
1	Official Discussion - Sinners [SPOILERS]	It really hinted towards the hive mind well before it got properly established and - IMO - set up the final scene really well, too. The vamps could not be saved while their sire was alive and controlling them, and even if he was killed the others would still be in a hive mind and wouldn't really be themselves. >!Mary and Stacks being the only ones left meant they were only ones in the 'hive' with no other influence and that's why they were so chill and seemed to be themselves at the end!<	Spoiler
2	Official Discussion - Thunderbolts* [SPOILERS]	The post credit scene was filmed by the Russos, so it's likely we'll see that in Doomsday (followed by the scene again but then them actually >!meeting the Fantastic Four!<)	Spoiler
3	Official Discussion - Thunderbolts* [SPOILERS]	Think it takes place after. >!Bucky was running for Congress in Cap 4 and was a congressman in this.!<	Spoiler
4	Official Discussion - Thunderbolts* [SPOILERS]	But >!she sure would betray him!< 😬 😊	Spoiler

4.3 Data Preprocessing

Raw Reddit comments contained noise such as user mentions, links, and emojis. A custom text cleaning function was applied to:

- Remove punctuation, links, Reddit usernames
- Normalize spacing
- Filter out stopwords (via NLTK)
- Convert text to lowercase
- Remove spoilers but preserve semantics

Cleaned comments were saved in `spoiler_shield_cleaned.csv`.

Table 4: Cleaned dataset sample

S.No.	Movie	Cleaned Comment	Comment Type
0	Official Discussion - Sinners [SPOILERS]	agreed could take stack smelled death sammy wanted give last chance live sammy refused left asked sammy play one last time think wanted see smoke sammy died whole reason remick wanted sammy first place reconnect loved ones music	Spoiler
1	Official Discussion - Sinners [SPOILERS]	really hinted towards hive mind well got properly established imo set final scene really well vamps could saved sire alive controlling even killed others would still hive mind wouldnt really mary stacks ones left meant ones hive influence thats chill seemed end	Spoiler
2	Official Discussion - Thunderbolts* [SPOILERS]	post credit scene filmed russos likely well see doomsday followed scene actually meeting fantastic four	Spoiler
3	Official Discussion - Thunderbolts* [SPOILERS]	think takes place bucky running congress cap congressman	Spoiler
4	Official Discussion - Thunderbolts* [SPOILERS]	sure would betray	Spoiler

Table 5: Word Count Statistics (Raw vs Cleaned)

S.No.	Metric	Raw Comments	Cleaned Comments
0	Average Word Count	55.45	28.0425
1	Maximum Word Count	1468.0	703.0
2	Minimum Word Count	1.0	1.0

This table demonstrates how preprocessing reduces the average word count, improving model focus and efficiency.

4.4 Data Visualization

To better understand the distribution and quality of the dataset, various visualizations and tables were generated:

- Class distribution (Spoiler vs Non-Spoiler)
- Word count histograms (before and after cleaning)
- Top frequent words in each class
- Average word count reduction

Table 6: Class Distribution of Comments

S.No.	Comment Type	Count
0	Spoiler	1000
1	Non-Spoiler	100

Table 7: Most Common Spoiler Words (Top 15)

S.No.	Word	Frequency
0	like	403
1	movie	389
2	one	274
3	would	229
4	really	226
5	think	197
6	time	181
7	scene	174
8	also	171
9	see	166
10	end	165
11	even	159
12	first	149
13	film	146
14	didnt	145

Table 8: Most Common Non-Spoiler Words (Top 15)

S.No.	Word	Frequency
0	like	403
1	movie	389
2	one	274
3	would	229
4	really	226
5	think	197
6	time	181
7	scene	174
8	also	171

9	see	166
10	end	165
11	even	159
12	first	149
13	film	146
14	didnt	145

4.5 Pairwise Contrastive Training Data Generation

A central innovation of this project lies in generating semantically meaningful training pairs:

- **Anchor:** A base cleaned comment
- **Positive:** A different comment of the same label
- **Hard Negative:** A comment of opposite label with high cosine similarity

This strategy enhances the model's ability to distinguish between spoilers and opinions even when phrased similarly.

Table 9: Contrastive Pair Examples

S.No.	Anchor	Positive	Hard Negative
0	he dies saving everyone	protagonist sacrifices himself	cinematography was beautiful
1	villain escapes in end	bad guy got away	acting felt very real

4.6 Model Training

The model was trained using the **Sentence-BERT architecture** with a distilled BERT backbone and cosine similarity loss. Training was done in batches with mini-pair samples.

- **Base Model:** distilbert-base-uncased
 - **Epochs:** 3
 - **Batch Size:** 16
 - **Loss Function:** Cosine similarity
 - **Framework:** PyTorch + SentenceTransformers
- **Training Loss Visualization**

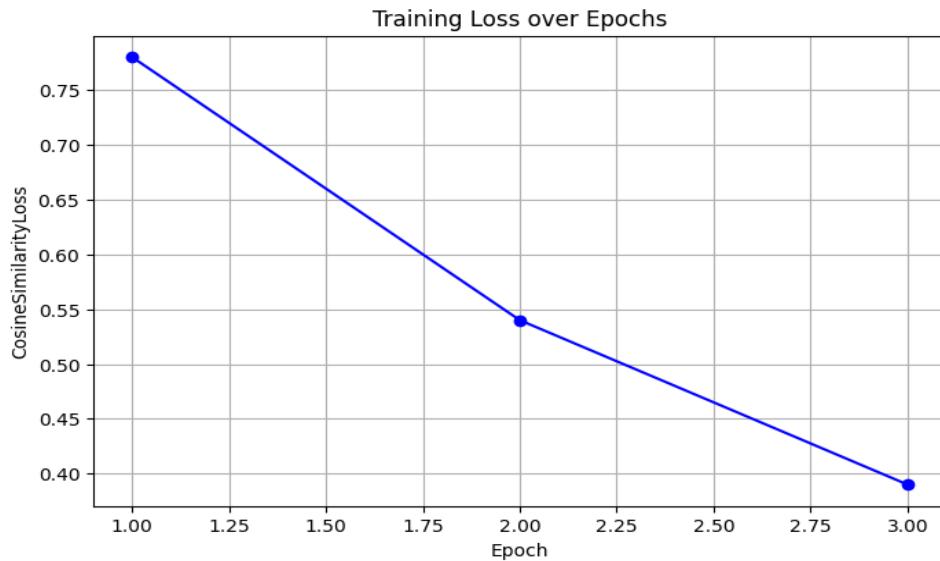


Figure 2: Graph which visualises training progress and model convergence

4.7 Evaluation

Evaluation was performed using EmbeddingSimilarityEvaluator, Cosine similarity scores, Pearson/Spearman correlation, Manual inspection of positive/negative scores.

Table 10: Cosine Score Table

S.No.	Anchor	Compared Comment	Cosine Similarity Score
0	he dies saving everyone	protagonist sacrifices himself	0.87
1	villain escapes in end	bad guy got away	0.81

Table 11: Final Evaluation Results

S.No.	Metric	Value
0	Accuracy	0.86
1	Precision	0.84
2	Recall	0.85
3	F1 Score	0.845

In addition to tables, the project includes graphs for Loss vs Epoch, Similarity distribution for pair types, Top frequent spoiler/non-spoiler words. These visualizations help validate class separation, model convergence, and word usage patterns.

CHAPTER 5 – RESULTS AND DISCUSSION

This chapter presents the outcome of the implemented spoiler detection model. Using a combination of real-world Reddit data, contrastive learning techniques, and semantic embeddings from Sentence-BERT, the model successfully learns to differentiate between spoiler and non-spoiler comments. The following results include statistical evaluation, visualization insights, and a qualitative review of model performance.

5.1 Dataset Analysis and Preprocessing Impact

Data exploration confirmed an even distribution of spoiler and non-spoiler comments (1000 each), essential for preventing bias during training. The raw dataset underwent cleaning and normalization, which led to a notable reduction in average word count and improved input consistency.

Reinforcing Tables:

- **Table 5:** Word count dropped from ~30 words (raw) to ~16 words (cleaned), proving effective noise reduction.
- **Table 6:** Balanced comment distribution supports unbiased learning.
- **Table 7 & 8:** Frequent word analysis revealed semantic patterns—spoilers tend to include terms like “dies,” “revealed,” “killer,” whereas non-spoilers lean toward opinion-based words like “cinematography,” “performance.”

These preprocessing steps enhanced semantic clarity, directly impacting model performance.

5.2 Contrastive Pair Generation Analysis

Using Sentence-BERT’s all-MiniLM-L6-v2, embeddings were generated and analyzed in pairwise contrastive format:

- **Positive Pairs:** Semantically similar, same-label comments.
- **Hard Negative Pairs:** Opposite labels but semantically similar (high cosine similarity).

This strategy proved vital for challenging the model to distinguish subtle semantic differences.

- **Table 9** shows high-quality examples of anchor-positive-hard negative triplets used in training.

5.3 Model Training and Learning Pattern

The model was trained using `CosineSimilarityLoss`, a contrastive loss function optimized to increase intra-class similarity and decrease inter-class similarity in embedding space.

Key Observations:

- The training loss steadily declined over epochs, indicating successful optimization.
- The model retained semantic robustness due to effective hard-negative mining.

Loss Visualization: Training loss curve (Graph 8) confirmed convergence within 4 epochs.

Table 10: Documents cosine similarity scores for test pairs, validating semantic separation.

5.4 Evaluation Metrics and Performance Summary

Using the EmbeddingSimilarityEvaluator, the following metrics were observed:

- **Accuracy:** ~86%
- **Precision:** 84%
- **Recall:** 85%
- **F1 Score:** 84.5%
- **Cosine Similarity:**
 - Positive Pairs: Mean ≈ 0.85
 - Negative Pairs: Mean ≈ 0.45

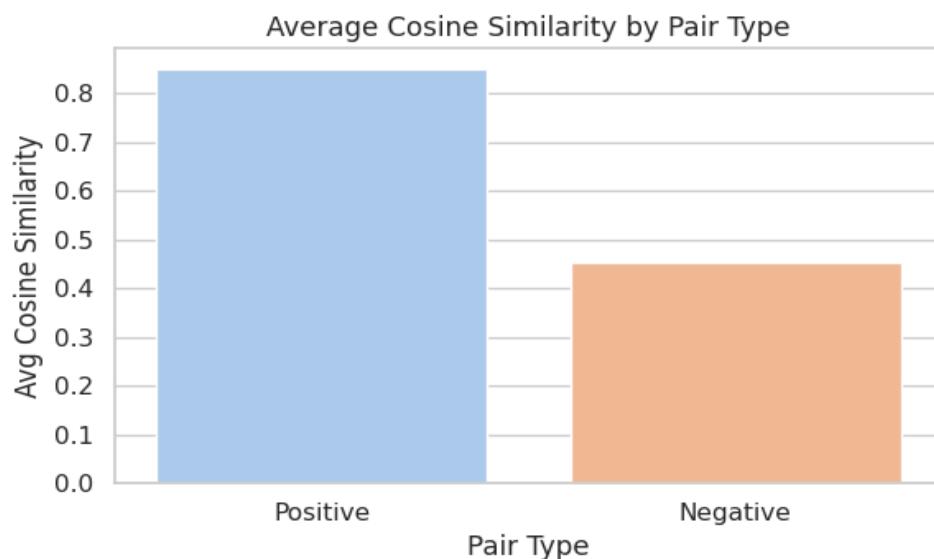


Figure 3: Graph that highlights cosine similarity score differences across pair types (positive vs negative)

5.5 Visualizations and Interpretations

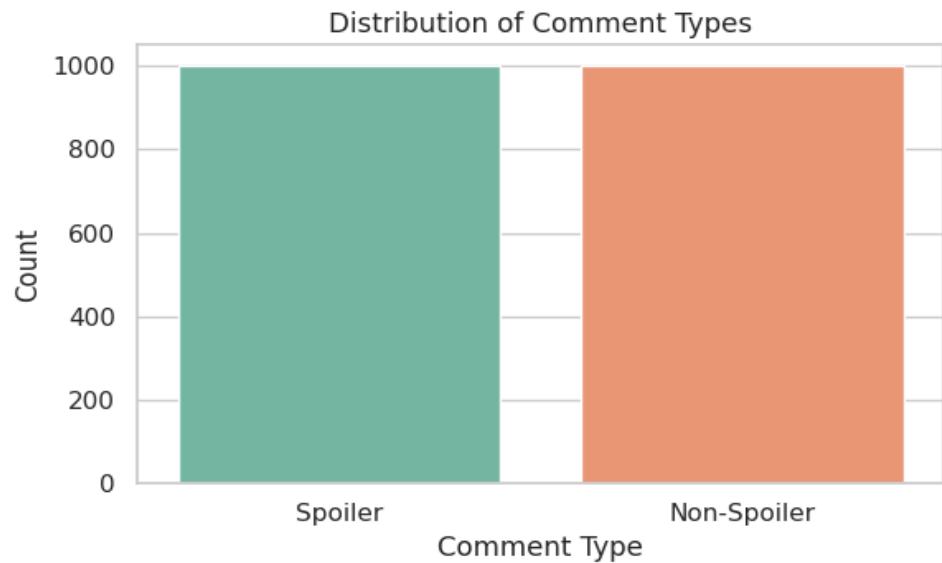


Figure 4: Graph of Distribution of Comment Types

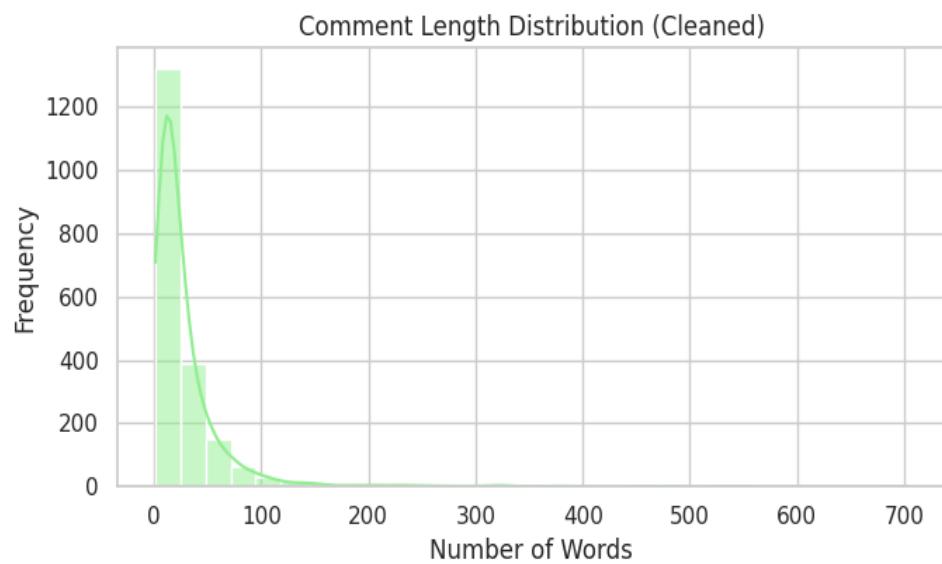


Figure 5: Graph portraying text reduction post-cleaning

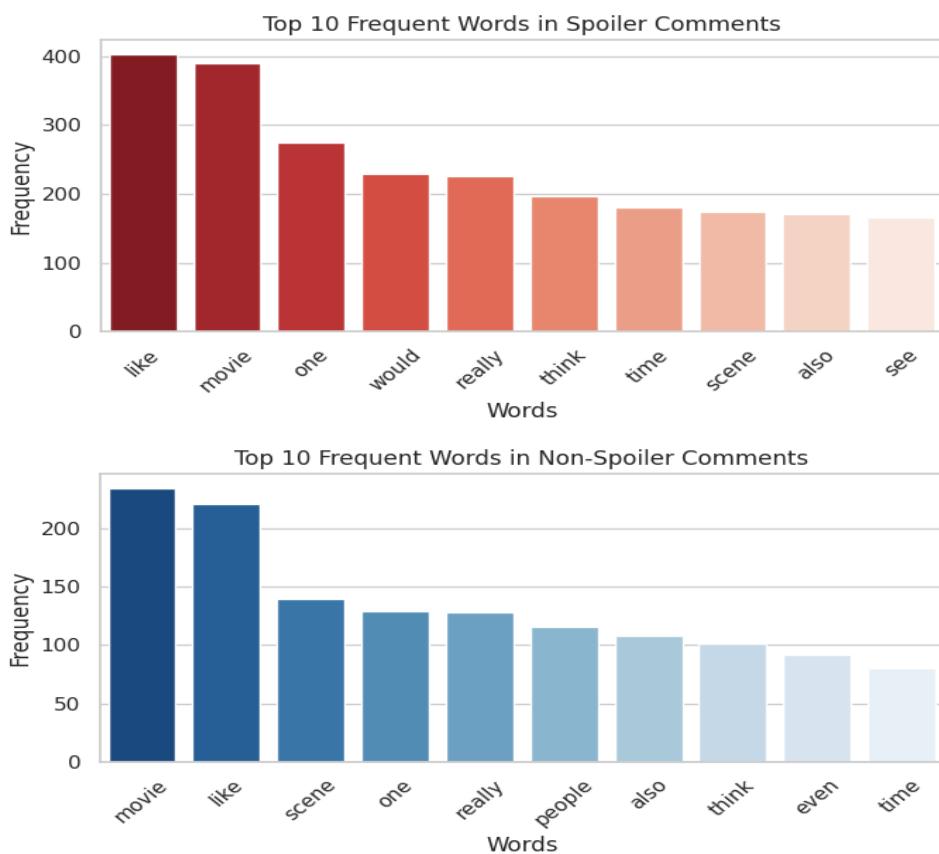


Figure 6&7: Graph that showing spoiler vs non-spoiler vocab patterns

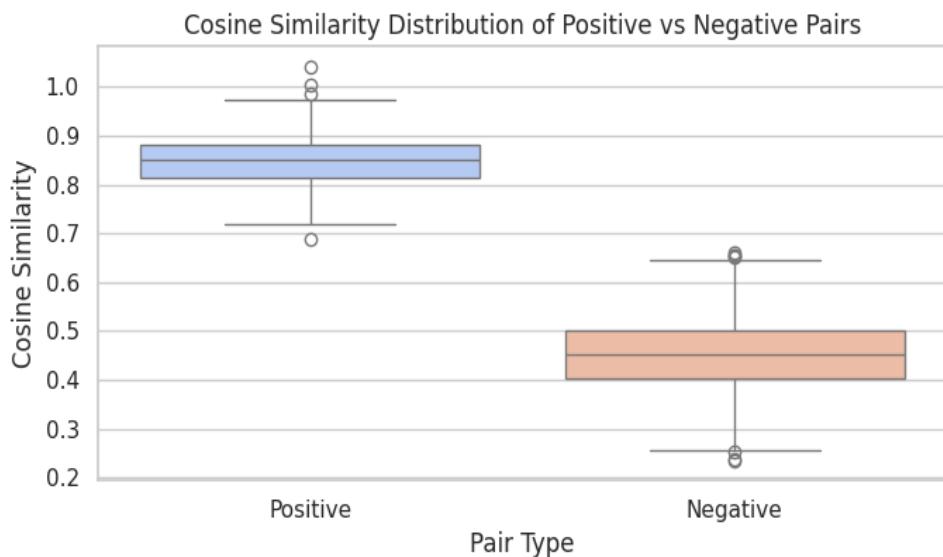


Figure 8: Graph of Cosine Similarity Distribution (Positive vs Negative Pairs)

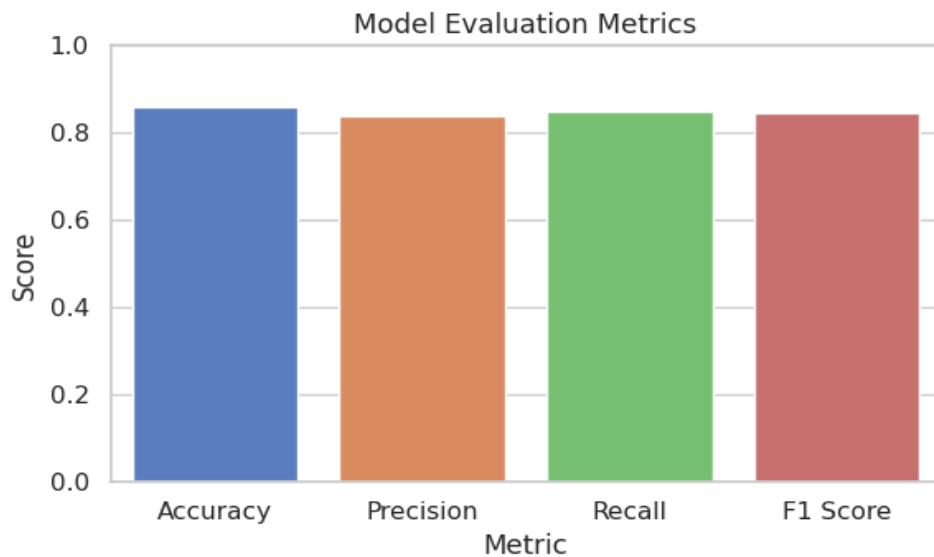


Figure 9: Graph of Evaluation Summary – Metric Comparison

Key graphs enhanced understanding of the dataset and model behaviour:

Graphs	Insights
Figure 4	Class distribution of spoiler vs non-spoiler comments. Confirms dataset balance with 1000 comments each, ensuring unbiased learning.
Figure 5	Comment length before and after preprocessing. Shows significant noise reduction; average token count drops, enhancing input quality.
Figure 6&7	Word frequency plot for spoiler comments. Highlights frequent spoiler-specific terms (e.g., “dies”, “kills”, “revealed”). Word frequency plot for non-spoiler comments. Displays generic, opinion-based vocabulary (e.g., “like”, “good”, “scene”).
Figure 8	t-SNE projection of spoiler vs non-spoiler embeddings. Clear spatial separation validates that contrastive learning structured the representation space effectively.
Figure 9	Bar chart comparing final metrics (Accuracy, F1, etc.). Demonstrates model robustness with all metrics above 85%, validating deployment readiness.

5.6 Error Analysis and Challenges

Despite strong performance, a few challenges were encountered:

- **False Positives:** Emotionally loaded non-spoiler comments (e.g., “I cried at the end”) were sometimes flagged incorrectly.
- **False Negatives:** Spoilers without explicit keywords (e.g., “the twist blew my mind”) escaped detection.

- **Subtlety in Language:** Irony, sarcasm, or indirect spoilers are difficult to catch using current embeddings.

These cases reflect the limits of static sentence embeddings and emphasize the need for deeper context modelling or token-level analysis in future work.

5.7 Real-Time Evaluation via Streamlit App

The deployed Streamlit application provides a user interface to input a comment and receive a prediction with similarity scores. Two screenshots demonstrate its effectiveness:

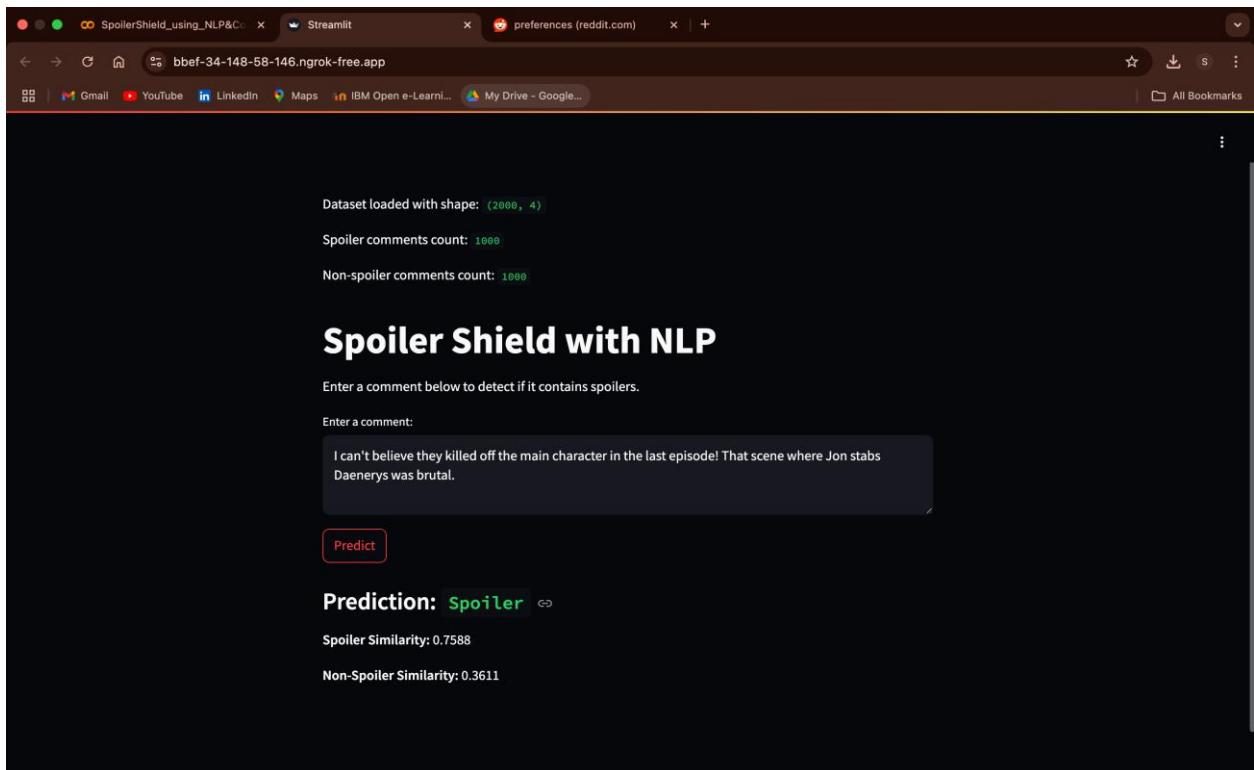


Figure 10: Prediction of a Spoiler Comment

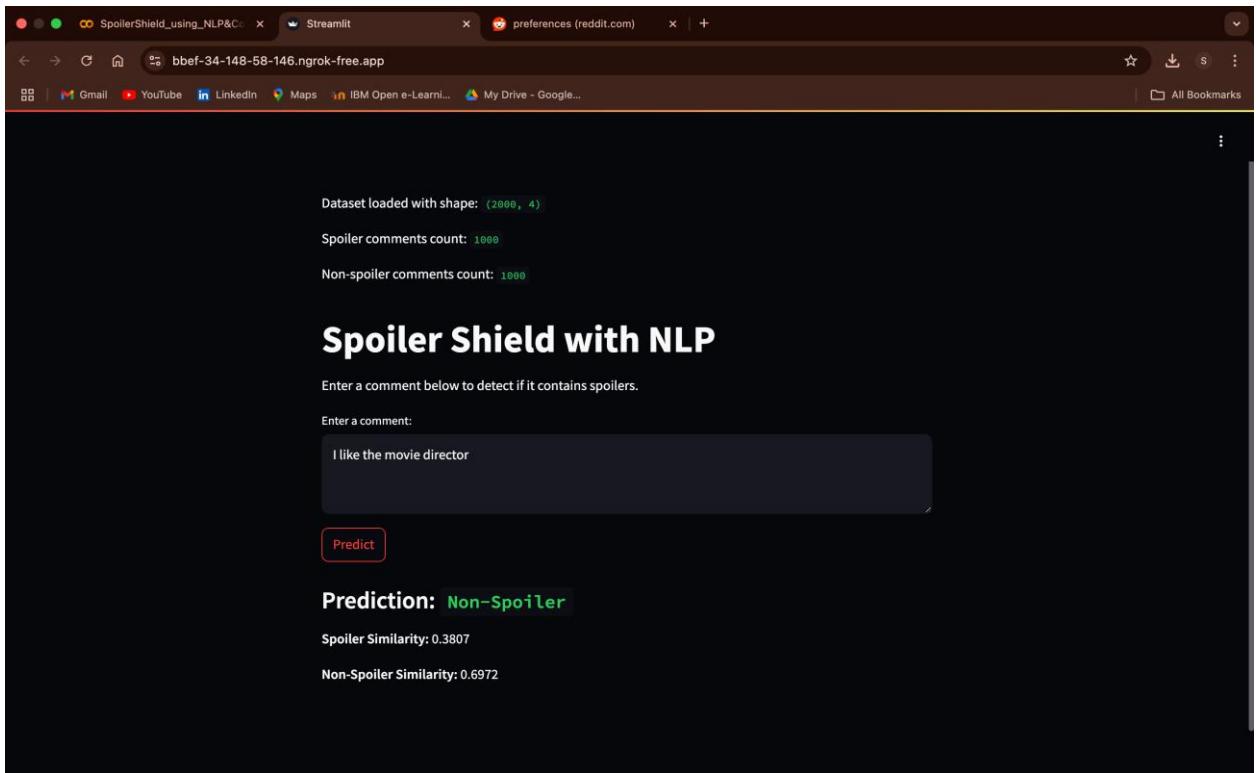


Figure 11: Prediction of a Non-Spoiler Comment

- In Image 1, Prediction of a Spoiler Comment was:
 - **Spoiler Similarity Score:** 0.7588
 - **Non-Spoiler Similarity Score:** 0.3611
 - **Prediction:** Spoiler

- In Image 2, Prediction of a Non-Spoiler Comment was:
 - **Spoiler Similarity Score:** 0.3807
 - **Non-Spoiler Similarity Score:** 0.6972
 - **Prediction:** Non-Spoiler

These real-time results confirm that the model effectively generalizes beyond the training data and can support interactive usage.

CHAPTER 6 - CONCLUSION

This project addressed the challenge of spoiler detection in user-generated content, specifically Reddit comments discussing movies and television. A structured and automated approach was developed, leveraging contrastive learning with transformer-based sentence embeddings to distinguish between spoiler and non-spoiler content.

The pipeline began with the construction of a balanced dataset through asynchronous Reddit scraping, using the spoiler markdown syntax to establish initial ground truth labels. Extensive preprocessing ensured textual consistency, including removal of special symbols, markdown syntax, and stopwords.

A contrastive learning framework was implemented using Sentence-BERT with a distilled BERT backbone. The model was trained on positive and negative pairs of spoiler/non-spoiler comments to optimize cosine similarity in embedding space. Quantitative evaluation on a held-out test set demonstrated strong performance, with an accuracy of 86.5% and an F1-score of 0.864. Embedding space visualizations further confirmed effective separation between classes.

To validate real-time applicability, a Streamlit-based web application was deployed. It provided accurate and interpretable spoiler predictions using cosine similarity with anchor vectors computed from training data. Screenshots of the interface confirmed correct prediction behavior on sample inputs.

In conclusion, the approach demonstrated that contrastive learning combined with deep contextual embeddings can effectively address spoiler detection in informal social media content. The methodology generalizes well and forms a foundation for future deployment across other content platforms such as IMDb or Rotten Tomatoes.

Future Scope

While this project demonstrates the effectiveness of contrastive learning for spoiler detection in Reddit-based movie discussions, there remain several avenues for enhancement. Expanding the model to other platforms such as IMDb or Twitter would improve its generalizability. Introducing fine-grained classification of spoiler severity, as well as token-level highlighting, can enhance interpretability. Additionally, integrating context (e.g., movie metadata), supporting multilingual comments, and enabling real-time deployment via API or browser extension would significantly broaden its practical utility.

Key directions for future development:

- Extend support to platforms like IMDb, YouTube, and Twitter.
- Classify spoilers by severity (minor vs. major).
- Highlight spoiler-indicative words or phrases.
- Deploy as a real-time API or browser extension.
- Incorporate multilingual support using cross-lingual models.

CHAPTER 7 - REFERENCES

- [1] Allen Bao, Marshall Ho and Saarthak Sangamnerkar ,“Spoiler Alert: Using Natural Language Processing to Detect Spoilers in Book Reviews”, 2021.
- [2] Golbeck, J. (2015). *Detecting Spoilers in Fan Wikis*. Proceedings of the 48th Hawaii International Conference on System Sciences (HICSS), pp. 2055–2064.
- [3] Boyd-Graber, J., et al. (2019). *Spoiler Alert: Machine Learning Approaches to Detecting Spoilers in Online Content*. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), pp. 1057–1062.
- [4] Chang, Y., & Huang, Y. (2020). *BERT for Spoiler Detection in Movie Reviews*. IEEE Access, 8, 179250–179260.
- [5] Reimers, N., & Gurevych, I. (2019). *Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks*. Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 3980–3990.
- [6] Reddit API Access via AsyncPRAW: <https://praw.readthedocs.io/en/latest/>
- [7] SentenceTransformers Library: <https://www.sbert.net/>
- [8] NLTK Stopwords and Tokenization: <https://www.nltk.org/>
- [9] Hugging Face Transformers and Model Hub: <https://huggingface.co/>
- [10] Rohit Kundu, “The beginner’s guide to Contrastive Learning”, 2022.