# Project Report

## Information Retrieval (CS60092)

### *Evidence Retrieval for Fact Verification*

### *Group - 9*

| | | |
|---|---|---|
| **ASRP Vyahruth** | - | **19CS10002** |
| **Amartya Mandal** | - | **19CS10009** |
| **Divyansh Bhatia** | - | **19CS10027** |
| **Riddihman Moulick** | - | **20EE10093** |

### *Abstract :*

Our aim in this project is to retrieve evidence for various Wikipedia article based claims. Fact verification is of great importance in the present scenario where fake news and claims are highly prevalent. Fact verification is generally a 2 step process which involves evidence retrieval and then analyzing the evidence to understand whether it supports or refutes the claim, but here we are concerned only with the retrieval process.

### *Problem Statement :*

To Retrieve evidence for 'claims' from Wikipedia articles. Given a claim (sentence), the objective is to retrieve sentences relevant to that claim from the various wikipedia articles. The sentence can support or refute the claim, but that is not a part of our objective, our concern is to just retrieve the sentences which are relevant to the claim made.

### *Proposed Solution :*

We break down the problem into 2 tasks. The first task involves retrieving Wikipedia articles containing information about this claim. The second task involves retrieving the facts in the form of relevant sentences (which serves as evidence in our case) from those documents. We use three different methods to retrieve wiki articles and

two methods to retrieve relevant sentences from the chosen articles as part of our experiment.

## Technique & Experiments:

The three methods (implemented as three separate layers), we used for document retrieval are Fuzzy Matching, Tf-idf and BERT.

Out of these, the BERT model has highest accuracy but due to the large number of Wiki articles, the time consumption is large (running BERT for all claim-article pairs is practically not feasible), whereas fuzzy matching although less accurate is comparatively faster.

Thus, there is a tradeoff between computation speed and accuracy and to obtain a good balance between the two we apply different retrieval models at different levels. We have also experimented with various threshold values for the different models and arrived at the values that provide the most accurate evidences within a feasible time.

We first use the simple fuzzy matching algorithm on the document 'titles' and the given claim for all the documents and get a score of similarity. If the similarity score is less than a set threshold, we discard the document. If the similarity is above a threshold, then we run the other two better methods like BERT to get a similarity score and use this score to classify a wiki article as relevant or not.

After retrieving the documents, we similarly run Tf-idf and BERT as two layers sequentially to match each sentence of the retrieved documents with the claim to get the similarity score and return the top k highest scored sentences as evidence for our claim.

We have done some **experiments** inside code itself such as doing the same task of sentence retrieval using two different approaches with a simplistic string matching to complex layering method to compare the efficiency and relevancy of evidence, proving our method's success.

## Related Work :

In order to get a good background idea of the work previously done in this domain, we

referred to some research papers such as:

1) Paper on Graph-based Evidence Aggregating and Reasoning for Fact Verification.

Link :- (https://arxiv.org/pdf/1908.01843.pdf)

2) UKP-Athene: Multi-Sentence Textual Entailment for Claim Verification.

Link :- (https://aclanthology.org/W18-5516.pdf).

➢ *What have we done differently?*

We have used multiple layers in both document and sentence retrieval tasks to balance the time - efficiency bridge, instead of using pre - trained models for both tasks as used in above papers. Further we have

As the number of data is very large, using layers we can improve time but not on the expense of efficiency or relevancy of evidence.

We have used different methods on different layers, ranging from simplistic tf-idf and fuzzy string matching to highly complex pre - trained NLP model called BERT which is designed to help computers understand the meaning of ambiguous language in text by using surrounding text to establish context.

## *Analysis and Future thoughts :*

After implementing our proposed layering structure, we used it to retrieve evidence for our claim or facts. In the appendix section, we can clearly see the correctness of our proposed technique by manually checking the relevancy of the retrieved sentences with the claim given.

We also maintained a decent time - relevancy balance, indicating the effectiveness of our layering structure instead of using only one method like BERT which will result in high efficiency rate but will take much more time too. Thus , we implemented a balanced technique.

As a future work, we can try to directly embed sentences into a vector and compare them instead of tokenizing the sentences to words and then individually comparing their embeddings. Also we can add more layers using different models, or make them optional so that we can use them as and when we require them.

## *Work Distribution :*

| Name | Work |
|---|---|
| ASRP Vyahruth - 19CS10002 | ➔ Proposed using BERT for word embeddings to compare the tokens of claims and retrieved sentences.<br>➔ Extracted data from the JSON files and worked on data cleaning.<br>➔ Experimented with various values of thresholds and checked the accuracy on selected examples to arrive at optimal values.<br>➔ Took the lead in Report preparation. Helped in preparing the ppt for the final presentation. |
| Amartya Mandal - 19CS10009 | ➔ Proposed fuzzy matching as a primary filter to filter out articles completely irrelevant in the first layer.<br>➔ Implemented the sentence retrieval block in the code.<br>➔ Worked on the demo and contributed to the report by highlighting the problem statement, proposed solution along with related work.<br>➔ Helped in preparing the ppt for the final project presentation. |
| Divyansh Bhatia - 19CS10027 | ➔ Highlighted the benefits of using a NLP model (BERT) to improve efficiency along with other layers.<br>➔ Implemented the document retrieval block in the code retrieving documents for evidence extraction.<br>➔ Contributed in the report by highlighting what we have done differently from the references.<br>➔ Helped in preparing the ppt for the final project presentation. |
| Riddhiman Moulick - 20EE10093 | ➔ Proposed the tf-idf model as a part of the layered structure<br>➔ Tokenised and lemmatized the data obtained from the Wiki articles and claims into a readable form<br>➔ Contributed to the analysis, future prospects and possible improvements in our report<br>➔ Took the lead in the presentation preparation. |

# Appendix :

## *Implementation and Experimental Results :*

Dataset Snippet:

```
{'id': 89296, 'claim': 'Henry Spencer is played by a Greek actor.'}
{'id': 78554, 'claim': 'John Ritter died in October.'}
{'id': 83809, 'claim': '13 Reasons Why is the only television series of 2012 in the drama-mystery genre.'}
{'id': 49758, 'claim': 'Playboy is a magazine.'}
{'id': 22973, 'claim': 'Alternative metal is the genre in which Alice in Chains usually performs.'}
{'id': 181494, 'claim': 'Sam Peckinpah directed The Wild Bunch.'}
{'id': 161592, 'claim': "The St. John's water dog is a breed of domestic dog that was first bred in Newfoundland."}
{'id': 117342, 'claim': 'Horseshoe crabs are not used in fertilizer.'}
{'id': 172204, 'claim': 'Sia (musician) has received an award presented by the cable channel MTV.'}
{'id': 95552, 'claim': 'Artificial intelligence raises concern.'}
```

Tokenization of both the claim and the wiki articles is done as below:

```
['henry', 'spencer', 'played', 'greek', 'actor']
['john', 'ritter', 'died', 'october']
['13', 'reason', 'television', 'series', '2012', 'drama-mystery', 'genre']
['playboy', 'magazine']
['alternative', 'metal', 'genre', 'alice', 'chain', 'usually', 'performs']
['sam', 'peckinpah', 'directed', 'wild', 'bunch']
['st', 'john', 'water', 'dog', 'breed', 'domestic', 'dog', 'first', 'bred', 'newfoundland']
['horseshoe', 'crab', 'used', 'fertilizer']
['sia', 'musician', 'received', 'award', 'presented', 'cable', 'channel', 'mtv']
['artificial', 'intelligence', 'raise', 'concern']
```

After passing the claims and document 'titles' through all the levels and different models, we are able to perform document retrieval as below (the document ID and the level of matching is displayed below the respective claims)

```
Claim:  Henry Spencer is played by a Greek actor.
(2625886, 7.2)          Henry_Spencer_Ashbee
(2638145, 7.2)          Henry_C._Spencer
(2660251, 7.2)          Henry_E._Spencer
(3469927, 7.2)          Henry_Spencer_Palmer
(3566599, 7.2)          Lord_Henry_Spencer
(4560605, 7.1)          Henry_Spencer_Berkeley
(3096802, 6.9)          List_of_actors_who_have_played_Sherlock_Holmes
(4039115, 6.9)          List_of_Greek_actors
(4542972, 6.9)          Henry_Elvins_Spencer
(4555676, 6.9)          Henry_Spencer
```

Once the document retrieval is done, the BERT and tf-idf models are used in a similar manner for sentence selection.

Below is a comparison of the results obtained from normal string matching (fuzzy matching) and our multi-model matching: (The top 5 evidences are being displayed here)

```
Claim:  Henry Spencer is played by a Greek actor.
Evidence:
This is a list of Greek actors .
Henry Spencer -LRB- born 1955 -RRB- is a Canadian computer programmer and space enthusiast .
The list of actors who have played Sherlock Holmes in film , television , stage , or radio includes :
He is coauthor , with David Lawrence , of the book Managing Usenet .
Spencer was succeeded as mayor by Mark P. Taylor in 1851 .


Claim:  Henry Spencer is played by a Greek actor.
Evidence:
Henry Spencer Ashbee -LRB- 21 April 1834 -- 29 July 1900 -RRB- was a book collector , writer , and bibliographer .
Henry Christian Spencer -LRB- 1915 -- 2000 -RRB- was an American chemical engineer and executive at the Kerite Company in Seymour , Connecticut .
Henry Evans Spencer -LRB- born June 13 , 1807 in Columbia - now part of Cincinnati -RRB- was a notable Cincinnati resident and was Mayor of Cincinnati from 1843-1851 .
Major General Henry Spencer Palmer -LRB- 30 April 1838 -- 10 February 1893 -RRB- was a British army military engineer and surveyor , noted for his work in developing Yok
Lord Henry John Spencer -LRB- 20 December 1770 -- 3 July 1795 -RRB- was a British diplomat and politician .
```

## *Additional references used :*

a) Fuzzy string matching also known as Approximate String Matching, is the process of finding strings that approximately match a pattern - https://towardsdatascience.com/fuzzy-string-matching-in-python-68f240d910fe

b) Understanding BERT state of the art language model for NLP - https://medium.com/@dhartidhami/understanding-bert-word-embeddings-7dc4d2ea54ca

********** END**********