



# Automated Evidence Collection for Fake News Detection

ICON 2021 : 18th International Conference on Natural Language Processing

Mrinal Rawat  
[rawatmrinal06@gmail.com](mailto:rawatmrinal06@gmail.com)

Diptesh Kanojia  
[d.kanojia@surrey.ac.uk](mailto:d.kanojia@surrey.ac.uk)

# Content

- Motivation
- Introduction
- Dataset
- Approach
- Experiments
- Results & Discussion
- Conclusions & Future Work

# Motivation

- Quality on social media is significantly affected due to the spread of fake news, misinformation and unverifiable facts.
- Panic due to Fake News in the epidemic situation like COVID-19.

# Motivation

- Quality on social media is significantly affected due to the spread of fake news, misinformation and unverifiable facts.
- Panic due to Fake News in the epidemic situation like COVID-19.
- Websites such as Poynter, Factcheck, etc. currently rely on methods involving manual detection of fake news which is a cumbersome and challenging task.
- Automatic Fake News detection aims to mitigate the problem of misinformation with the help of evidence supported by multiple sources.

# Motivation

- Quality on social media is significantly affected due to the spread of fake news, misinformation and unverifiable facts.
- Panic due to Fake News in the epidemic situation like COVID-19.
- Websites such as Poynter, Factcheck, etc. currently rely on methods involving manual detection of fake news which is a cumbersome and challenging task.
- Automatic Fake News detection aims to mitigate the problem of misinformation with the help of evidence supported by multiple sources.

## Facebook posts

stated on March 22, 2020 in a Facebook post:

**“Boil some orange peels wit cayenne pepper in it stand over the pot breathe in the steam so all that mucus can release from yo nasal... MUCUS is the problem its where THE VIRUS LIVES!!!”**



PUBLIC HEALTH

FACEBOOK FACT-CHECKS

CORONAVIRUS

FACEBOOK POSTS

**Water boiled with orange peels and cayenne pepper will not prevent or cure COVID-19**

Screenshot from PolitiFact (<https://www.politifact.com/factchecks/2020/mar/25/facebookposts/water-boiled-orange-peels-and-cayenne-pepper-will-/>)

# Introduction

- Previous approaches rely on classical machine learning models (Vlachos, and Riedel, 2014) and deep learning based approaches (Malon, 2018; Vijjali et al.,2020) on LIAR and FEVER datasets.
- Previous work on CONSTRAINT-2021 dataset (Patwa et al, 2020) does not include evidence since the dataset only have claim and label pairs.

# Introduction

- Previous approaches rely on classical machine learning models (Vlachos, and Riedel, 2014) and deep learning based approaches (Malon, 2018; Vijjali et al.,2020) on LIAR and FEVER datasets.
- Previous work on CONSTRAINT-2021 dataset (Patwa et al, 2020) does not include evidence since the dataset only have claim and label pairs.
- We address these issues and propose a novel approach that automatically collects the evidence from multiple sources.
- We incorporate a summarization component that helps outperforms the state-of-the-art approaches on CONSTRAINT-2021 shared task achieving a F1-score of 0.9925.

# Dataset

We use the pre-released COVID-19 fake news dataset as a part of the CONSTRAINT- 2021 shared task (Patwa et al., 2020). Each post or tweet contains content in the English language and is classified as:

- **Real:** where tweets or articles which are factually correct and verified from authentic sources
- **Fake:** where tweets or posts related to COVID-19 which are factually incorrect and verified as false.

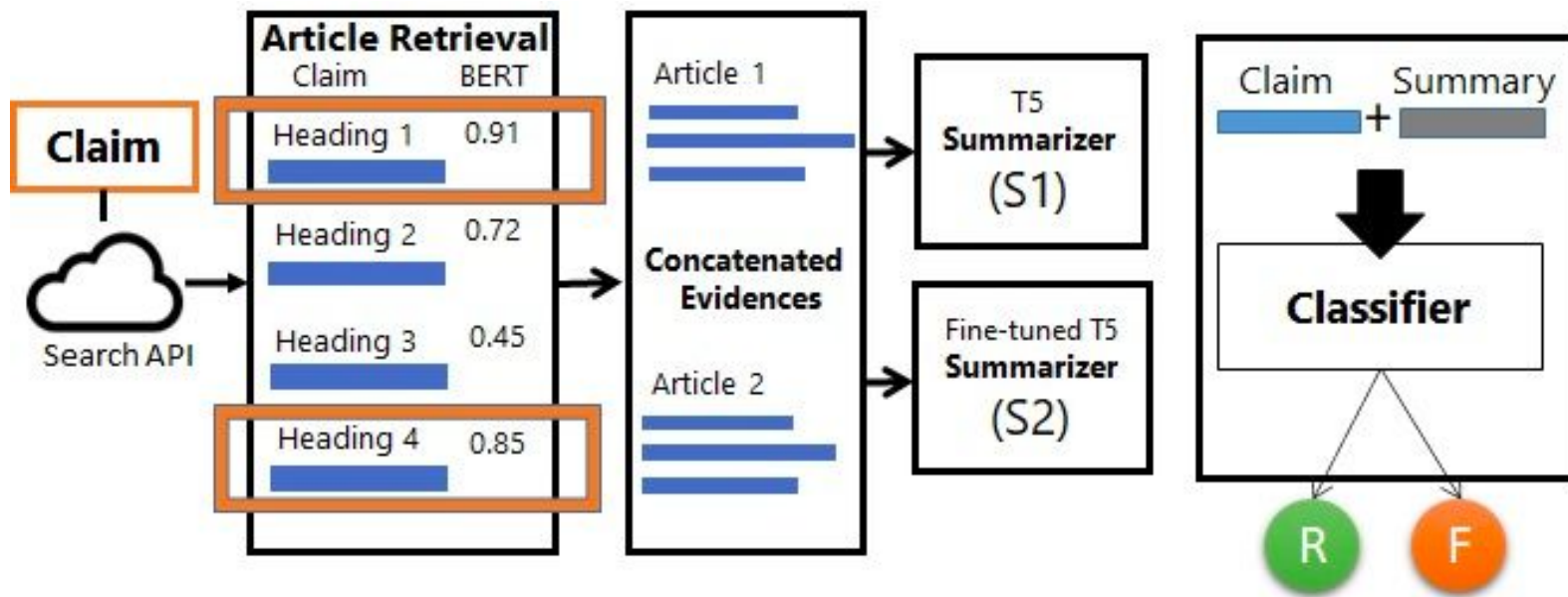




# Approach



# Approach



[CLS]	If	you	tested	...	[SEP]	SOURCES:	<a href="https://cdc.gov">https://cdc.gov</a>	[SEP]	Staying	away	...
-------	----	-----	--------	-----	-------	----------	---	-------	---------	------	-----

# Approach (Cont'd)

## Evidence Collection

The original dataset of COVID-19, the evidence is not released along with the claim. Hence, we extract the evidence which assists the classifier.

- **Article Retrieval:** For each claim  $c$ , we perform a Google search and use BERT to get the similarity score of response text w.r.t. to the input claim. We select top 3 results which has the similarity score greater than 0.7.
- **Sentence (Evidence) Retrieval:** We extract the relevant article URLs  $U = (u_1, u_2, u_3)$  and employ a similar method to find the top  $k$  sentences within each article.

# Approach (Cont'd)

## Evidence Collection

The original dataset of COVID-19, the evidence is not released along with the claim. Hence, we extract the evidence which assists the classifier.

- **Article Retrieval:** For each claim  $c$ , we perform a Google search and use BERT to get the similarity score of response text w.r.t. to the input claim. We select top 3 results which has the similarity score greater than 0.7.
- **Sentence (Evidence) Retrieval:** We extract the relevant article URLs  $U = (u_1, u_2, u_3)$  and employ a similar method to find the top  $k$  sentences within each article.

## Evidence Summarization

We summarize the evidence collected in previous step using following two ways: •

- **Google T5 (S1):** We use the Google T5-large (Raffel et al., 2019) model which is the current state-of-the-art for the summarization tasks.
- **Fine-tuning on FEVER Dataset (S2):** We first fine-tune the Google T5 on FEVER dataset which has ground truth claim and evidence, and then use it to summarize the evidence collected using our approach.

# Approach (Cont'd)

Claim	Evidence	Summarization-1 (S1)	Summarization-2 (S2)
<p>There is no evidence that children have died because of a COVID-19 vaccine. No vaccine currently in development has been approved for widespread public use. <a href="https://t.co/9ecvMR8SAf">https://t.co/9ecvMR8SAf</a></p>	<p>Currently there is no coronavirus vaccine that has been approved for the American public. And there is no evidence that children have died because they received one of the COVID-19 vaccines being developed. PolitiFact found no evidence that anyone has died from complications related to a trial COVID-19 vaccination. There is no evidence that children have died because of a COVID-19 vaccine.</p>	<p>There is no evidence that children have died because they received a COVID-19 vaccine. No evidence that anyone has died from complications related to a trial COVID-19.</p>	<p>There is no evidence that children have died because they received one of the COVID-19 vaccines being developed. PolitiFact found no evidence that anyone has died from complications related to a trial COVID-19 vaccination.</p>

# Experiments

We evaluate the COVID-19 fake news detection using on following three methods:

- No evidence
- Summarized evidence without fine-tuning
- Summarized evidence with fine-tuning

The evaluation metrics used are as following:

Accuracy, Precision, Recall and F1-Score

## **Models Used:**

- Machine Learning Models:  
Logistic Regression, Support Vector Machines, Random Forest Classifier
- Deep Learning Models:  
LSTM, BERT, RoBERTa, XLNet

# Results & Discussion

	Previous Approaches		Our Approach w/ various Classification methods								
	Chen et al. (2021)	Li et al. (2021)	Logistic Regression			SVM			LSTM		
-			S1	S2	-	S1	S2	-	S1	S2	
<b>Precision</b>	0.9902	0.986	0.9531	0.9565	0.9701	0.9641	0.9671	0.9764	0.9589	0.9598	0.9612
<b>Recall</b>	0.9901	0.985	0.9531	0.9564	0.9700	0.9639	0.9668	0.9761	0.9584	0.9596	0.9612
<b>F-Score</b>	0.9901	0.985	0.9531	0.9565	<u>0.9700</u>	0.9639	0.9668	<u>0.9761</u>	0.9584	0.9596	<u>0.9612</u>

Results obtained after the fake news classification task where the values for previous approaches are from the latest shared task results and the results for each iteration of our approach are shown [P (Precision), R (Recall), and F (F-Score)]. (-) → No Evidence, S1 → Summarization-1 as Evidence, S2 → Summarization-2 as Evidence.



# Results & Discussion (Cont'd)

	Previous Approaches		Our Approach w/ various Deep Learning Classification methods								
	Chen et al. (2021)	Li et al. (2021)	BERT <sub>base</sub>			RoBERTa <sub>base</sub>			XLNet <sub>base</sub>		
			-	S1	S2	-	S1	S2	-	S1	S2
<b>Precision</b>	0.9902	0.986	0.9612	0.9916	0.9917	0.9918	0.9929	0.9922	0.9920	0.9934	0.9947
<b>Recall</b>	0.9901	0.985	0.9864	0.9888	0.9897	0.9897	0.9911	0.9916	0.9892	0.9911	0.9925
<b>F-Score</b>	0.9901	0.985	0.9858	0.9888	<u>0.9893</u>	0.9893	0.9908	<b><u>0.9908</u></b>	0.9892	0.9910	<b><u>0.9925</u></b>

Results obtained after the fake news classification task where the results for each iteration of our approach with various deep learning classification methods are shown [P (Precision), R (Recall), and F (F-Score)]. (-) → No Evidence, S1 → Summarization-1 as Evidence, S2 → Summarization-2 as Evidence.

# Conclusions & Future Work

- We present an automated method to collect the evidence and summarize it for the fake news detection task.
- Our approach helps in augmenting the dataset released in the CONSTRAINT-2021 task.
- Our systematic framework achieves a F1- score of 0.9925.

# Conclusions & Future Work

- We present an automated method to collect the evidence and summarize it for the fake news detection task.
- Our approach helps in augmenting the dataset released in the CONSTRAINT-2021 task.
- Our systematic framework achieves a F1- score of 0.9925.
- We show that a summarization module that can help collect evidence more effectively.
- In future, we want to experiment and reproduce the results on other fact verification tasks.

# Conclusions & Future Work (Cont'd)

USTER

FAKE NEWS BUSTER

Yesterday our laboratories completed 2899 tests of those 726 were testing of people in managed isolation

Real

**Confidence :**

0.9503107666969299

**Sources:**

nzherald.co.nz

**Evidence:**

Yesterday, New Zealand's laboratories completed 2899 tests, and 726 of those were taken at managed isolation or quarantine facilities. The total number of completed tests is now 436,233. The NZ Covid Tracer app has now recorded 607,000 registrations after a "flurry" of Kiwis downloaded it overnight. He said there had been tens of thousands of tests in the community during the past few weeks.

USTER

FAKE NEWS BUSTER

Nasal flu vaccine side effects do not cause covid

Fake

**Confidence :**

0.9999988079071045

**Sources:**

cdc.gov reuters.com nhs.uk

**Evidence:**

Nasal flu vaccine side effects are usually mild and would not cause a positive COVID-19 test. As Reuters reported previously (here), the nasal flu vaccine contains live flu viruses that have been weakened. The shed virus does not cause flu in others. The nasal spray flu vaccine contains small amounts of weakened flu viruses. They do not cause flu in children.

**Source Code:** [https://github.com/rawat-mrinal06/fake\\_news/](https://github.com/rawat-mrinal06/fake_news/)

**THANK YOU**