# ADSumm: Annotated Ground-truth Summary Datasets for Disaster Tweet Summarization

Piyush Kumar Garg[1*], Roshni Chakraborty[2] and Sourav Kumar Dandapat[1]

[1]Department of Computer Science and Engineering, Indian Institute of Technology Patna, Bihar, India.
[2]Institute of Computer Science, University of Tartu, Estonia.

*Corresponding author(s). E-mail(s): piyush_2021cs05@iitp.ac.in;
Contributing authors: roshni.chakraborty@ut.ee;
sourav@iitp.ac.in;

**Abstract**

Online social media platforms, such as Twitter, provide valuable information during disaster events. Existing tweet disaster summarization approaches provide a summary of these events to aid government agencies, humanitarian organizations, etc., to ensure effective disaster response. In the literature, there are two types of approaches for disaster summarization, namely, supervised and unsupervised approaches. Although supervised approaches are typically more effective, they necessitate a sizable number of disaster event summaries for testing and training. However, there is a lack of good number of disaster summary datasets for training and evaluation. This motivates us to add more datasets to make supervised learning approaches more efficient. In this paper, we present *ADSumm*, which adds annotated ground-truth summaries for eight disaster events which consist of both natural and man-made disaster events belonging to seven different countries. Our experimental analysis shows that the newly added datasets improve the performance of the supervised summarization approaches by **8 — 28**% in terms of ROUGE-N F1-score. Moreover, in newly annotated dataset, we have added a *category label* for each input tweet which helps to ensure good coverage from different categories in summary.

Additionally, we have added two other features *relevance label* and *keyphrase*, which provide information about the quality of a tweet and explanation about the inclusion of the tweet into summary, respectively. For ground-truth summary creation, we provide the annotation procedure adapted in detail, which has not been described in existing literature. Experimental analysis shows the quality of ground-truth summary is very good with *Coverage*, *Relevance* and *Diversity*.

**Keywords:** Tweet summarization, Ground-truth summary, Disaster, Social media, Datasets

# 1 Introduction

During disasters, eyewitnesses and bystanders share information through social networks, such as Twitter [1]. Several prior research works highlight that tweets are a valuable source of information during disasters [2–5]. This information consists of updates about casualties, infrastructure damage, information on missing people, urgent needs of resources and supplies, etc. Many disaster response organizations, government agencies, volunteers, and NGOs rely on this valuable information to plan and launch relief operations immediately [6]. However, tweets are inherently short and often contain grammatical errors, abbreviations, and informal language. These characteristics make it highly challenging to identify relevant information from them. Additionally, the huge number of tweets makes it highly challenging for any human being and government organizations to identify relevant information manually [7, 8].

To handle these challenges, several research works [9–12] have been proposed for disaster tweet summarization that generate an automated summary of the disaster tweets. Existing disaster tweet summarization approaches are broadly categorized into supervised and unsupervised approaches. However, supervised approaches are more efficient than unsupervised approaches. Efficiency of these automated supervised summarization approaches depend heavily on the availability of sizable ground-truth summaries for training. Furthermore, there is a high variance in information among disasters on the basis of the type and location of a disaster [9]. Therefore, the non-availability of ground-truth summaries for disasters of different locations and types affects the development of a robust and efficient summarization approach. Although Dutta et al. [12] and Rudra et al. [13] have provided the ground-truth summary of six datasets (shown in Table 1) which has been of huge help to the research community, there are around 50 different disaster datasets [8, 14, 15] which do not have any ground-truth summary. The addition of more ground-truth summaries of different disasters will surely improve this situation which motivates us to add more ground truth summaries.

In existing literature, we have not found any systematic approach for ground truth creation. Existing approaches completely rely on the wisdom and

knowledge of annotators, where a flat set of input tweets is provided to annotators to generate a summary. In this paper, to come up with a systematic approach, we mimic the important steps followed by automatic summarization approaches [9, 16], which could be described as 1) identification of the category of each tweet, 2) understanding of the importance of each category, and 3) selection of important tweets from each category for summary. However, most of these above-mentioned steps are subjective. While automated approaches take important decisions for summary creation based on some well-defined approach and parameters, annotators need to make these decisions based on their wisdom. Following the existing literature [17], we have selected three annotators for this task. This work presents *ADSumm*, which refers to the ground-truth summaries that have been created by three annotators for a total of eight distinct disaster events, each occurring in different locations and of varying types. The inclusion of these datasets would greatly benefit the research community. Furthermore, the incorporation of the recently introduced datasets enhances the efficacy of supervised summarization approaches by $8 - 28\%$ in relation to the ROUGE-N F1-score, as elaborated in the corresponding Subsection 5.4.

In addition to the ground-truth summaries of eight datasets, we provide annotation of three supplementary features that are not yet present in the public datasets. The aforementioned features encompass *category labels*, *key-phrases*, and *relevance labels*. The term *category label* pertains to a specific classification assigned to a tweet that falls within the context of a disaster. These categories may include *Infrastructure Damage, Volunteer Operations, Affected Population, Impact, Prayer*, and others. The *relevance label* assigned to a tweet indicates its significance in relation to the associated disaster occurrence, classifying it as *high*, *medium*, or *low*. The *key-phrase* elucidates the potential rationale underlying the significance of a tweet. For instance, a tweet, "Tropical storm Hagupit kills at least 21 in Philippines" can be classified under the *category label Affected Population*. The tweet contains a *key-phrase storm Hagupit kills at least 21 in Philippines*, which is *high*ly relevant to the Hagupit Typhoon[1] disaster event. We offer the *category label, key-phrase, and relevance labels* for each tweet in eight distinct disaster datasets, which might be of great assistance to the academic research community.

In this paper, we also present comprehensive analyses on the utility of the additional dataset features in assessing the effectiveness for various advanced natural language processing (NLP) tasks. These tasks can be broadly divided into three categories 1) disaster tweet classification, 2) development of a robust summarization algorithm, and 3) evaluation of the quality of the disaster summary. Feature *category label* provides ground truth for the tweet category, which can be used to assess the classification accuracy of any classification approach. As feature *key-phrase* helps in ranking tweets within a category and calculate diversity in the summary, *key-phrase* plays an important role in

---

[1] https://en.wikipedia.org/wiki/Typhoon_Hagupit_(2014)

**Table 1**: Table shows the details of available six disaster datasets, which include dataset name, number of tweets, summary length, country, continent, and disaster type.

| Dataset name | Number of tweets | Summary length | Country | Continent | Disaster type |
|---|---|---|---|---|---|
| *Sandy Hook Elementary School Shooting* | 2080 | 36 tweets | United States | North America | Man-made |
| *Uttrakhand Flood* | 2069 | 34 tweets | India | Asia | Natural |
| *Hagupit Typhoon* | 1461 | 41 tweets | Philippines | Asia | Natural |
| *Hyderabad Blast* | 1413 | 33 tweets | India | Asia | Man-made |
| *Harda Twin Train Derailment* | 4171 | 250 words | India | Asia | Man-made |
| *Nepal Earthquake* | 5000 | 250 words | Nepal | Asia | Natural |

robust summarization algorithm development. Furthermore, the feature *relevance label* immensely helps to assess the summary quality as it provides quantitative understanding of the importance of each tweet. We also conduct a comparative analysis of 16 existing state-of-the-art summarization approaches to evaluate their efficacy in summarizing eight datasets related to disaster events. This evaluation serves as a reference for benchmarking purposes. The datasets were utilized as additional ground-truth summaries as described in [9].

The rest of the paper is organized as follows. We discuss related works in Section 2. In Section 3, we provide the details of the datasets and discuss the ground-truth preparation details in Section 4. In Section 5, we discuss results where we provide qualitative and quantitative analysis results of the annotated ground-truth summary in Subsection 5.1. We provide an inter-annotator agreement between all the annotators to check the consistency for ground-truth generation in Subsection 5.2. We compare summaries generated by the supervised approaches to assess the importance of additional training set for summary generation in Subsection 5.4. We also provide a case study where we show the use cases of the dataset's additional features in the different NLP-related tasks in Subsection 5.5. We discuss the experiment details and results of the performance comparison of the various existing state-of-the-art summarization approaches in Subsection 5.6. Finally, we conclude the paper in Section 6.

## 2 Related Works

Summarization provides a brief summary of input text, emphasizing its essential aspects and the key information [18]. Text summarization approaches have been proposed for various domains, such as legal texts summarization [19], news summarization [20], timeline summarization [21], tweet summarization [22, 23], etc. Similarly, there are several disaster specific summarization approaches [24–26] which could be categorized into content and context based [16, 27], graph-based approaches [28] and deep learning based approaches [29].

For example, content-based disaster tweet summarization approaches [11, 13, 30] explore the importance of keywords to generate a summary which might not ensure coverage in summary. Recent deep learning-based approaches have proposed techniques, such as Dusart et al. [10] integrated the importance of each word by frequency into Bidirectional Encoder Representations

from Transformers (BERT) [31], Garg et al. [32] integrated the *key-phrase* of each tweet into BERT, and Li et al. [27] capture the inter-tweet similarity through a Graph Convolution Neural (GCN) network to generate the summary. However, these approaches require extensive training and therefore, more number of disaster events with annotated summaries. Graph-based summarization approaches [11, 12, 33] initially generate a graph of tweets where the nodes are tweets and the edges represent the similarity between a pair of tweets followed by identification of groups of similar tweets using different approaches, such as communities [34], connected components [35], and $k$-clique clustering [36]. Finally, they select the representative tweets from each group to create a summary. Although these existing approaches of identifying groups are highly effective in news event-based tweet summarization [22, 23], they fail for disaster events where there is a high vocabulary overlap across different groups. Therefore, Rudra et al. [30, 37] and Garg et al. [9] initially identify the category of a tweet and then select representative tweets from each category into the summary. While Rudra et al. [30, 37] use Artificial Intelligence for Disaster Response (AIDR) [38], which requires human intervention for each new disaster in real-time, Garg et al. [9] proposed an un-supervised ontology-based category identification approach that does not require any human intervention for a disaster event. However, the vocabulary and importance of the category differ on the basis of the type and location of the disaster. Therefore, deep learning, graph-based, and category-based approaches require ground-truth summaries of different types and locations to learn the inherent differences.

Existing research works, such as Rudra et al. [13] and Dutta et al. [12] have provided the ground-truth summary of six disaster events (shown in Table 1) belonging to both the type of disaster (man-made and natural) and four different countries. However, these datasets are not good enough to train a supervised summarization approach. Therefore, these approaches could not reach the required performances, as discussed in Section 5.4. Although, the existing datasets did not provide any supplementary features, such as *key-phrases*, *category labels*, and *relevance labels*. These features, along with the datasets, could be further used for evaluating different NLP-related tasks such as disaster tweet classification, developing a robust summarization algorithm, and evaluating the quality of disaster summary. The non-availability of ground-truth summaries along with features for different locations and types affects the development of a robust summarization algorithm and quality checking which is suitable for disaster events irrespective of the variance in disaster data. Therefore, in this paper, we present the annotated ground-truth summaries along with the three different features with the dataset for eight different disaster events belonging to natural and man-made disasters and seven more different countries.

# 3 Datasets

In this Section, we discuss the datasets for which we prepare the ground-truth summaries. We show the details for $D_1$-$D_8$ datasets in Table 2.

**Table 2**: Table shows the details of $D_1$-$D_8$ datasets, which include dataset number, dataset name, year, number of tweets, summary length, country, continent, and disaster type.

| Num | Name | Year | Number of tweets | Summary length | Country | Continent | Disaster type |
|---|---|---|---|---|---|---|---|
| $D_1$ | Los Angeles International Airport Shooting | 2013 | 1409 | 40 tweets | United States | North America | Man-made |
| $D_2$ | Hurricane Matthew | 2016 | 1654 | 40 tweets | Haiti | North America | Natural |
| $D_3$ | Puebla Mexico Earthquake | 2017 | 2015 | 40 tweets | Mexico | North America | Natural |
| $D_4$ | Pakistan Earthquake | 2019 | 1958 | 40 tweets | Pakistan | Asia | Natural |
| $D_5$ | Midwestern U.S. Floods | 2019 | 1880 | 40 tweets | United States | North America | Natural |
| $D_6$ | Kaikoura Earthquake | 2016 | 2195 | 40 tweets | New Zealand | Oceania | Natural |
| $D_7$ | Cyclone Pam | 2015 | 1508 | 40 tweets | Vanuatu | Oceania | Natural |
| $D_8$ | Canada Wildfires | 2016 | 2242 | 40 tweets | Canada | North America | Natural |

1. $D_1$: This dataset is created by [14] from a terrorist attack on the *Los Angeles International Airport Shooting*[2] on November, 2013 in which 1 person was killed and more than 15 people were injured.

2. $D_2$: This dataset is created by [8] from the *Hurricane Matthew*[3] on October, 2016 in which around 603 people were killed, around 128 people were missing and the estimated damage were around $2.8 billion USD.

3. $D_3$: This dataset is created by [8] from the *Puebla Mexico Earthquake*[4] on September, 2017 in which 370 people were dead and more than 6000 people were injured.

4. $D_4$: This dataset is created by [8] from the *Pakistan Earthquake*[5] on September, 2019 in which around 40 people were killed, 850 people were injured, and around 319 houses were damaged.

5. $D_5$: This dataset is created by [8] from the *Midwestern U.S. Floods*[6] in which around 14 million people were affected, and damage were around 2.9 billion USD.

6. $D_6$: This dataset is created by [8] from the *Kaikoura Earthquake*[7] on November, 2016 in which 2 people were died, and around 57 were injured.

7. $D_7$: This dataset is created by [1] from the *Cyclone Pam*[8] in March, 2015 in which around 15 people were died and around 3300 people were displaced.

8. $D_8$: This dataset is created by [8] from the *Canada Wildfires*[9] in May, 2016 in which $1,456,810$ acres were burned, $3,244$ buildings were destroyed, and damage were around 9.9 billion C$.

**Dataset Pre-processing**

We perform standard pre-processing steps, like, lemmatization, conversion to lowercase, removal of noisy keywords, white spaces, and punctuation marks. Additionally, we consider only tweet text and subsequently we remove the Twitter-specific keywords like hashtags, URLs, usernames, and emoticons [39]. To reduce the redundancy, we further remove retweets followed by duplicate tweets. Lastly, we follow Alam et al. [15] to remove noise, i.e., any word which

---

[2] https://en.wikipedia.org/wiki/2013_Los_Angeles_International_Airport_shooting
[3] https://en.wikipedia.org/wiki/Hurricane_Matthew
[4] https://en.wikipedia.org/wiki/2017_Puebla_earthquake
[5] https://en.wikipedia.org/wiki/2019_Kashmir_earthquake
[6] https://en.wikipedia.org/wiki/2019_Midwestern_U.S._floods
[7] https://en.wikipedia.org/wiki/2016_Kaikoura_earthquake
[8] https://en.wikipedia.org/wiki/Cyclone_Pam
[9] https://en.wikipedia.org/wiki/2016_Fort_McMurray_wildfire

is of length less than three characters. We show some examples of tweets for $D_3$ and $D_6$ in Table 3.

**Table 3**: Table shows some examples of tweets text of two disaster events, such as *Puebla Mexico Earthquake* ($D_3$) and *Kaikoura Earthquake* ($D_6$).

| Event | Tweet text |
|---|---|
| *Puebla Mexico Earthquake* | 7.1 magnitude earthquake in Mexico kills 226 people. |
| | RT @latimes: 2,000 historic buildings in Mexico have been damaged by the earthquake. |
| | Hundreds of rescue personnel and citizens in Mexico City have banded together to rescue earthquake survivors. |
| | RT @Kevinwoo91: My prayers are continuing to be with everyone who has been affected by the earthquake in Mexico #PrayForMexico |
| | Venezuela Delivers 10.4 Tons of Aid to Earthquake-Ravaged Mexico. |
| *Kaikoura Earthquake* | RT @NZcivildefence: 14 staff from Urban Search &amp; Rescue have been deployed to #Wellington to help assess buildings. Another teams on stand |
| | #New Zealand PM John Key says 2 people killed in earthquake; sending military helicopter to Kaikoura - Reuters |
| | RT @nytimes: An earthquake measuring 7.9 hit New Zealand, triggering 3 large aftershocks and at least 3 tsunami waves |
| | Reports of damage to at least 25 buildings in Wellington so far #eqnz #wellington |
| | Praying for #NZ - still a Tsunami threat from the 7.5 mag #eqnz. May Allah keep everyone safe. Pls follow instructions from @NZcivildefence |

# 4 Ground-truth Preparation

In this Section, we discuss the annotation procedure followed by annotators to generate a ground-truth summary of a disaster event. In the existing literature [12, 16, 37], we find that annotators are provided a flat set of tweets of a disaster event and prepare the summary based on their knowledge and wisdom. Annotators manually assess the importance of each tweet with respect to the disaster event and then decide whether it should be part of the summary based on intuition given all the tweets related to a disaster event. These approaches mainly depend on the understanding of the annotators related to that disaster event. Therefore, there is a need for a systematic approach to ground-truth summary creation from a flat collection of tweets. Additionally, ground-truth
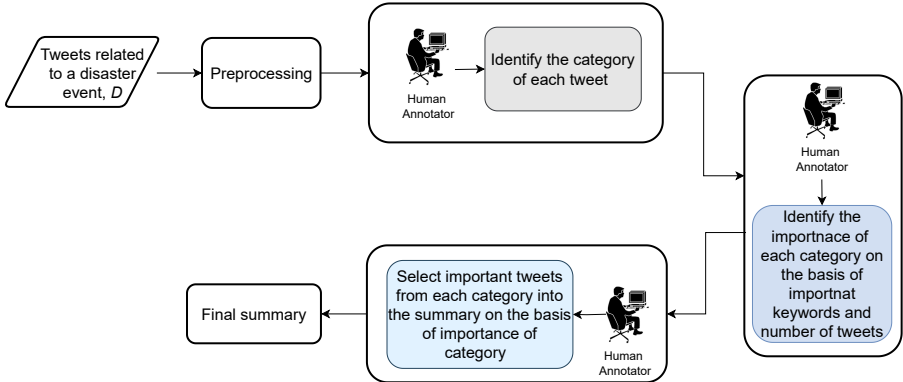
**Fig. 1**: An overview of the proposed systematic annotation procedure is shown.

summary generation is a subjective task, so we can not depend on only one annotator for the summary, and we require at least three annotators for their individual summaries [12, 13, 40]. Additionally, the ground-truth generation mainly depends on the annotator's wisdom related to a disaster event, and therefore, there is a need to ensure the quality and consistency of the annotators before selecting them for the ground-truth summary generation task. For this, we follow Garg et al. [17] where we select three annotators on the basis of a *Quality Assessment Test*. We refer to these annotators as $A_1$, $A_2$, and $A_3$ in the rest of the paper. We next discuss the details of the systematic approach followed by the annotators to generate the final ground-truth summary.

To develop a systematic approach, we imitate the important steps followed by the existing automatic summarization approaches [9, 16], which could be as follows: 1) identification of the category of each tweet, 2) identification of the importance of each category on the basis of tweets' important keywords and the number of tweets, and 3) selection of important tweets from each category in summary on the basis of the importance of that category. However, the majority of the aforementioned steps involve subjective judgment, and therefore, for each dataset, we request all three annotators to go through all the tweets given a dataset and prepare a ground-truth summary with the given summary length by following these steps manually. We follow Dutta et al. [12] and Rudra et al. [16] to decide the summary length as 40 tweets. For the first step of ground-truth creation, which is category identification of each tweet, prior works [1, 11, 41] have highlighted that the tweets of a disaster event belong to different categories. A ground-truth summary of a disaster should provide *Coverage* all important aspects[10]. However, the dataset does not comprise the category of the tweet inherently. To handle this, we share a publicly available ontology for disasters, namely, *Empathi* [42], which comprises of 70 categories. A set of keywords are also provided associated with these categories. We request the annotators to initially segregate all the

---

[10]In this paper, we refer aspects by tweet category.

tweets into these different ontology categories. Based on the keyword similarity (similarity between keywords of a category and keywords of a tweet), a tweet is assigned a category label by annotators. As a next step, annotators identify the importance of each category on the basis of both the number of tweets and the important keywords of the tweets in each category. Further, depending on the importance of a category, an annotator decides how many tweets to be selected from a category on the basis of his/her judgment of the importance of that category with respect to the disaster. An annotator can even consider a category to be not important enough to be covered in the final summary. Finally, annotators rank tweets within each category to select a specific number of tweets from a category in the summary. Annotators also take care about diversity during the selection of tweets from a category. An overview of this systematic annotation procedure is shown in Figure 1.

# 5   Results and Discussions

In this Section, we evaluate the quality (by means of qualitative and quantitative analysis) of annotated ground-truth summaries using three metrics, namely *Coverage*, *Relevance*, and *Diversity*. Further, to check the consistency among the annotators, we provide an inter-annotator agreement between all three annotators for two different tasks: 1) category assignment and 2) ground-truth summary. We further provide representation through the word cloud to visualise the annotated ground-truth summaries of a disaster event. To assess the importance of additional training set for the supervised approach of summary generation, we compare summaries generated by the supervised approach for two scenarios: 1) when trained with existing datasets and 2) when trained with existing datasets along with additional datasets. We also provide a case study where we show the use cases of the dataset's additional features in the different NLP-related tasks. Finally, to provide benchmark data, we compare and evaluate the existing state-of-the-art summarization approaches on the annotated ground-truth summaries for $D_1 - D_8$ datasets.

## 5.1   Qualitative and Quantitative Evaluation of a Ground-truth Summary

In this Subsection, we evaluate the quality of the generated ground-truth summary by means of quantitative and qualitative analysis based on *Coverage*, *Relevance*, and *Diversity*. We discuss each of the evaluation measures next.

### Qualitative Evaluation

We follow Garg et al. [17] for qualitative evaluation of the summary quality. A summary is considered to be of acceptable quality if it has high *Coverage*,

*Relevance*, and *Diversity*. We ask three meta-annotators[11] to score the ground-truth summaries on a scale of 1 (worst score) to 5 (best score) for each summary on the basis of its fulfilment of *Coverage*, *Relevance* and *Diversity* as *Coverage Score*, *Relevance Score*, and *Diversity Score*, respectively. For each summary, we get all the above-mentioned scores from all three meta-annotators. To combine scores from different annotators, we aggregate all the scores and take an average of that. For example, if *Coverage Score* of a summary from three meta-annotators are x, y, and z, respectively, then the *Aggregate Coverage Score* can be computed as $(x+y+z)/3$. In a similar way, we have also computed *Aggregate Relevance Score* and *Aggregate Diversity Score*. We show our observations in Table 4, which indicates that the ground-truth summaries are of high quality in terms of *Coverage*, *Relevance*, and *Diversity*.

**Table 4**: Table shows the *Aggregate Coverage Score*, *Aggregate Relevance Score*, and *Aggregate Diversity Score* of the annotated ground-truth summaries of all the three annotators for $D_1 - D_8$ datasets.

| Dataset | Aggregate Coverage Score | Aggregate Relevance Score | Aggregate Diversity Score | Dataset | Aggregate Coverage Score | Aggregate Relevance Score | Aggregate Diversity Score |
|---------|---------|---------|---------|---------|---------|---------|---------|
| $D_1$ | 3.94 | 4.19 | 3.86 | $D_5$ | 3.69 | 3.44 | 4.42 |
| $D_2$ | 3.78 | 4.22 | 3.38 | $D_6$ | 4.22 | 4.19 | 3.86 |
| $D_3$ | 4.33 | 3.31 | 3.94 | $D_7$ | 3.67 | 4.44 | 3.67 |
| $D_4$ | 4.05 | 3.47 | 3.75 | $D_8$ | 4.00 | 3.67 | 4.25 |

## Quantitative Evaluation

In this Subsection, we follow Garg et al. [17] to asses the annotated summaries through *Coverage*, *Relevance* and *Diversity*, quantitatively.

**Coverage :** Coverage is a well-accepted metric for assessing summary quality and provides a measure of how many important aspects are included in the summary. In this paper, we realize different aspects by means of category [17]. In order to compute this, we identify the category coverage in the annotator's annotated ground-truth summaries. We utilize the categories identified by the annotators during ground-truth preparation discussed in Section 4. We show the number of categories in the annotated ground-truth summaries of three annotators and in a dataset for $D_1 - D_8$ datasets in Table 5. Similar to the qualitative assessment, for quantitative analysis, we also compute the *Aggregate Coverage Score* for each ground-truth summary. Our observations indicate that the *Aggregate Coverage* of all three annotators is in the range of $51.85 - 77.78\%$ across the datasets. The $D_6$ dataset has the highest *Aggregate Coverage* (77.78%), and the $D_3$ dataset has the lowest *Aggregate Coverage* (51.85%). Furthermore, we observe that for the categories that are not captured in the summary, we found that both the number of tweets and the importance of those categories with respect to the disaster event are very low.

---

[11]The meta-annotators are graduate students who belong to the age group of $20 - 30$, have good knowledge of English and are not a part of this project.

Hence, we can say that the generated ground-truth summaries ensure good *Coverage* (covering important aspects/categories).

**Table 5**: We show the number of categories in a dataset and the ground-truth summaries and Aggregate Coverage (in %) of all three annotators for $D_1 - D_8$ datasets. (Note: # represents a number in this table.)

| Dataset | # of categories covered in a dataset | # of categories covered in ground-truth summary of $A_1$ | $A_2$ | $A_3$ | Aggregate Coverage in % | Dataset | # of categories covered in a dataset | # of categories covered in ground-truth summary of $A_1$ | $A_2$ | $A_3$ | Aggregate Coverage in % |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $D_1$ | 9 | 5 | 5 | 7 | 65.32 | $D_5$ | 8 | 5 | 7 | 6 | 75.00 |
| $D_2$ | 9 | 6 | 6 | 7 | 70.37 | $D_6$ | 9 | 7 | 7 | 7 | 77.78 |
| $D_3$ | 9 | 4 | 5 | 5 | 51.85 | $D_7$ | 9 | 7 | 7 | 7 | 77.78 |
| $D_4$ | 10 | 6 | 8 | 8 | 73.33 | $D_8$ | 9 | 5 | 5 | 6 | 59.26 |

**Relevance :** A summary has a limited capacity to capture and preserve the comprehensive details and nuances present in a dataset. Therefore, one measure of assessing the summary is *Relevance*, which quantifies the inclusion of highly relevant tweets in the summary. In order to compute this, we ask all the meta-annotators to assign a label to all the tweets based on tweet's importance with respect to the disaster, as *high*, *medium*, and *low*. Finally, label of a tweet is decided based on the majority vote of the three annotators. If we fail to decide based on a majority vote, we ask another annotator to decide the final label. However, for our datasets, we could assign all the labels on the basis of a majority vote. We show the percentage of each *relevance label* for ground-truth summaries for $D_1 - D_8$ datasets in Table 6. Our observations indicate that $57.50 - 75.00\%$ of the ground-truth summary tweets contain *high relevance label* across the datasets. Therefore, based on these observations, we can say that the generated ground-truth summaries comprise of significantly high-relevance tweets.

**Table 6**: Table shows the percentage number of tweets of each *relevance label*, such as *high*, *medium*, and *low* for the annotated ground-truth summaries of all the three annotators for $D_1 - D_8$ datasets. (Note: we release the final relevance label based on the majority vote for each dataset.)

| Dataset | $\mathbf{A_1}$ | | | $\mathbf{A_2}$ | | | $\mathbf{A_3}$ | | |
|---|---|---|---|---|---|---|---|---|---|
| | **High** | **Medium** | **Low** | **High** | **Medium** | **Low** | **High** | **Medium** | **Low** |
| $D_1$ | 60.00% | 20.00% | 20.00% | 65.00% | 17.50% | 17.50% | 60.00% | 17.50% | 22.50% |
| $D_2$ | 57.50% | 15.00% | 27.50% | 65.00% | 12.50% | 22.50% | 60.00% | 20.00% | 20.00% |
| $D_3$ | 65.00% | 15.00% | 20.00% | 62.50% | 15.00% | 22.50% | 70.00% | 7.50% | 22.50% |
| $D_4$ | 60.00% | 12.50% | 27.50% | 60.00% | 12.50% | 27.50% | 57.50% | 12.50% | 30.00% |
| $D_5$ | 65.00% | 22.50% | 12.50% | 65.00% | 10.00% | 25.00% | 60.00% | 17.50% | 22.50% |
| $D_6$ | 60.00% | 32.50% | 7.50% | 65.00% | 20.00% | 15.00% | 65.00% | 22.50% | 12.50% |
| $D_7$ | 60.00% | 22.50% | 15.00% | 75.00% | 15.00% | 10.00% | 57.50% | 20.00% | 22.50% |
| $D_8$ | 60.00% | 20.00% | 20.00% | 65.00% | 15.00% | 20.00% | 62.50% | 15.00% | 22.50% |

**Diversity :** *Diversity* is a well-accepted metric to evaluate the summary quality in terms of novel information present in summary [9]. We follow Garg et al. [17] and Nguyen et al. [43, 44] to calculate *Diversity* of a summary, which measures the mean *Diversity* between every pair of tweets in the summary

using tweets *key-phrases*. For this purpose, we also requested a meta-annotator to provide *key-phrases* of each tweet containing sufficient justification of its importance as well as information regarding the topic of the tweet. The meta-annotator follows the guideline where he/she initially identifies all the numeric, locations, and then identifies all the keywords relevant to disaster, and finally, he/she selects a *key-phrase* as a continuous set of words from a tweet. For instance, let us consider a tweet, *Tropical storm Hagupit kills at least 21 in Philippines*. As per the guideline, the following words will be marked- *storm*, *Hagupit*, *kills*, *21*, and *Philippines*. Therefore the *key-phrase* of the tweet which includes all the important keywords identified is *storm Hagupit kills at least 21 in Philippines*. We calculate the *Aggregate Diversity Score* over all the three annotators for eight disaster datasets. Our results, as shown in Table 7, indicate that the annotated ground-truth summaries contain $0.3885 - 0.6595$ *Aggregate Diversity Score* across the disasters. The $D_2$ dataset has the highest *Aggregate Diversity Score* (0.6595), and the $D_6$ dataset has the lowest *Aggregate Diversity Score* (0.3885). Therefore, the above experiment implies that the annotated ground-truth summaries ensure good diversity in the summary tweets.

**Table 7**: Table shows the *Aggregate Diversity Score* of the ground-truth summaries generated by the three annotators for $D_1 - D_8$ disaster datasets.

| Dataset | Aggregate Diversity Score | Dataset | Aggregate Diversity Score |
|---------|---------------------------|---------|---------------------------|
| $D_1$ | 0.5404 | $D_5$ | 0.4342 |
| $D_2$ | 0.6595 | $D_6$ | 0.3885 |
| $D_3$ | 0.5175 | $D_7$ | 0.4267 |
| $D_4$ | 0.5206 | $D_8$ | 0.4202 |

**Quality Comparison of existing datasets with our dataset :**   To understand the significance of the quantitative quality of the annotated ground-truth summaries, we compare these summaries with the existing available datasets through *Coverage*, *Relevance*, and *Diversity*. For comparison, we consider four available disaster datasets provided by Dutta et al. [12] i.e., *Sandy Hook Elementary School Shooting*, *Uttrakhand Flood*, *Hagupit Typhoon*, and *Hyderabad Blast* (ground-truth summaries annotated by three annotators[12] along with a tweet set). For above-mentioned datasets, we requested our three annotators for *category label* annotations, three meta-annotators for *relevance label* annotations, and one meta-annotator for *key-phrase* annotations. Then, we use generated ground truth summaries of those datasets as published by the respective authors [12] to compute *Coverage*, *Relevance* and *Diversity*.

For *Covergae*, we show the number of categories in the annotated ground-truth summaries of three annotators and in a dataset in Table 8. Our observations indicate that the *Aggregate Coverage* of all three annotators is in the range of $51.85 - 74.09\%$ across the datasets. For *Relevance*, we show the percentage of each *relevance label* for ground-truth summaries for these datasets in Table 8. Our observation indicates that $52.78 - 69.69\%$ of the

---

[12]We refer these three annotators as $E_1$, $E_2$, and $E_3$, hereby

ground-truth summary tweets of existing datasets contain *high relevance label* across the datasets. Similarly, we show the *Aggregate Diversity Score* over all the three annotators for above four disaster datasets in Table 8. Our observations indicate that the ground-truth summary of existing datasets contains $0.4327 - 0.7034$ *Aggregate Diversity Score* across the disasters. Therefore, we conclude that the quality of the proposed annotated ground-truth summaries is in the range of the quality of ground-truth summaries of existing datasets in terms of *Coverage*, *Relevance*, and *Diversity*.

**Table 8**: We show the number of categories in a dataset and the ground-truth summaries and Aggregate Coverage (in %) of all the three annotators, the percentage number of tweets of each *relevance label*, such as *high*, *medium*, and *low* for the annotated ground-truth summaries of all the three annotators, and *Aggregate Diversity Score* of the ground-truth summaries generated by the three annotators for four existing disaster datasets, i.e., *Sandy Hook Elementary School Shooting* (SHShoot), *Uttrakhand Flood* (UFlood), *Hagupit Typhoon* (HTyphoon), and *Hyderabad Blast* (HBlast). (Note: # represents a number in this table.)

| Dataset | # of categories covered in a dataset | # of categories covered in ground-truth summary of | | | Aggregate Coverage in % | Relevance label | | | | | | | | | Aggregate Diversity Score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | $E_1$ | $E_2$ | $E_3$ | | $E_1$ | | | $E_2$ | | | $E_3$ | | | |
| | | | | | | High | Medium | Low | High | Medium | Low | High | Medium | Low | |
| SHShoot | 9 | 5 | 6 | 4 | 55.55 | 63.89 | 16.67 | 19.44 | 52.78 | 36.11 | 11.11 | 63.89 | 16.67 | 19.44 | 0.4327 |
| UFlood | 9 | 4 | 5 | 5 | 51.85 | 64.71 | 11.76 | 23.53 | 52.94 | 26.47 | 20.59 | 58.88 | 20.59 | 20.59 | 0.6319 |
| HTyphoon | 9 | 6 | 7 | 7 | 74.07 | 60.98 | 24.39 | 14.63 | 58.53 | 21.95 | 19.51 | 63.41 | 17.07 | 19.51 | 0.7034 |
| HBlast | 9 | 5 | 5 | 6 | 59.25 | 60.60 | 18.18 | 21.21 | 57.57 | 21.21 | 21.21 | 69.69 | 15.15 | 15.15 | 0.6220 |

## 5.2 Inter-annotator Agreement

In this Subsection, we check the consistency among the three annotators through the inter-annotator agreement for summary generation. We check this agreement on the two different types: 1) category assignment and 2) ground-truth summary. We discuss each of the measures next.

### Category Assignment Inter-annotator Agreement

We compute Fleiss kappa [45] co-efficient score to calculate the inter-annotator agreement on the basis of the category assigned by an annotator given a tweet. We utilize the categories identified by all three annotators during ground-truth preparation for a disaster discussed in Section 4. Fleiss kappa provides an understanding of the similarities among the annotators. Our results, as shown in Table 9 indicate that the Fleiss kappa co-efficient score ranges between $72 - 80\%$. We also found that the score is highly substantial, within the range between $61 - 80\%$ [46]. Therefore, the kappa values across the category annotations ensure consistency across the annotators.

**Table 9**: Table shows the classification annotation agreement scores for $D_1$-$D_8$ datasets over the three annotators category annotated data.

| Dataset | Kappa score | Dataset | Kappa score |
|---------|-------------|---------|-------------|
| $D_1$ | 72.13 | $D_5$ | 72.29 |
| $D_2$ | 74.76 | $D_6$ | 79.46 |
| $D_3$ | 78.33 | $D_7$ | 77.89 |
| $D_4$ | 75.21 | $D_8$ | 78.12 |

## Ground-truth Summary Inter-annotator Agreement

To assess the consistency among the three annotators for summary creation, we measure it through summary-level Inter-annotator Agreement, where we compare their annotated summaries based on *content similarity*, *topical similarity*, and *semantic similarity*, as discussed below.

**Content similarity:** We calculate *syntactic similarity*, i.e., the distance between the tweet keywords based on their meaning [47]. To calculate the *syntactic similarity* between a pair of annotated summaries, we use Cosine Similarity [48], $CosSIM(A_i, A_j)$ between the summary generated by a pair of annotators, $A_i$ and $A_j$ as:

$$CosSIM(A_i, A_j) = \frac{|Kw(A_i) \cap Kw(A_j)|}{\sqrt{|Kw(A_i)| \ |Kw(A_j)|}} \tag{1}$$

where $Kw(A_i)$ and $Kw(A_j)$ are the keywords of $A_i$ and $A_j$ summaries, respectively. We consider only nouns, verbs and adjectives as keywords [49]. Our observations as shown in Table 10 indicate that $CosSIM(A_i, A_j)$ ranges between $72 - 96\%$ which indicates that summaries written by different annotators are syntactically similar to each other.

**Topical similarity:** To calculate *topical similarity*, $TopSIM(A_i, A_j)$ between any pair of annotators summaries, $A_i$ and $A_j$, we initially identify the top 15 topics words present in an annotator summary using Latent Dirichlet Allocation (LDA) [50]. We calculate the embedding for the topics of an annotator summary as the average of the values of the topic word vector. We consider the embedding of each topic word provided by Word2Vec [1], which was trained on 52M disaster-related messages from various types of disasters. We, finally, compute *topical similarity* as Cosine Similarity of the topic vector embedding of $A_i$ and $A_j$ as:

$$TopSIM(A_i, A_j) = \frac{\vec{T_i} \cdot \vec{T_j}}{|\vec{T_i}| \ |\vec{T_j}|} \tag{2}$$

where, $\vec{T_i}$ and $\vec{T_j}$ are the topic embedding vectors of $A_i$ and $A_j$ respectively. Our observations, as shown in Table 10 indicate that $TopSIM(A_i, A_j)$ ranges

**Table 10**: Table shows Content Similarity, Topical Similarity, Semantic Similarity (using Word2Vec), and Semantic Similarity (using BERT) % for $D_1$-$D_8$ datasets over the three annotators annotated summaries.

| Dataset | Annotator | Content Similarity | | Topical Similarity | | Semantic Similarity (using Word2Vec) | | Semantic Similarity (using BERT) | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\mathbf{A_2}$ | $\mathbf{A_3}$ | $\mathbf{A_2}$ | $\mathbf{A_3}$ | $\mathbf{A_2}$ | $\mathbf{A_3}$ | $\mathbf{A_2}$ | $\mathbf{A_3}$ |
| $D_1$ | $A_1$ | 83.85 | 81.93 | 87.06 | 87.49 | 99.12 | 98.94 | 99.58 | 99.26 |
| $D_2$ | $A_1$ | 89.27 | 91.55 | 84.63 | 90.45 | 98.93 | 98.98 | 99.64 | 99.54 |
| $D_3$ | $A_1$ | 89.64 | 90.82 | 91.36 | 85.32 | 99.03 | 98.95 | 99.41 | 99.69 |
| $D_4$ | $A_1$ | 83.17 | 78.68 | 91.08 | 86.71 | 99.38 | 98.93 | 99.43 | 99.21 |
| $D_5$ | $A_1$ | 69.93 | 71.84 | 87.61 | 84.53 | 99.44 | 99.11 | 99.40 | 99.52 |
| $D_6$ | $A_1$ | 88.34 | 86.49 | 93.55 | 93.55 | 99.73 | 99.68 | 99.72 | 99.70 |
| $D_7$ | $A_1$ | 91.14 | 95.59 | 97.66 | 97.95 | 99.71 | 99.79 | 99.87 | 99.93 |
| $D_8$ | $A_1$ | 85.12 | 91.35 | 97.41 | 98.49 | 99.45 | 99.69 | 99.55 | 99.72 |

between $85 - 98\%$, which indicates that the summaries written by different annotators are topically similar to each other.

**Semantic similarity:** In order to understand the context similarity between the annotated summaries, we calculate *semantic similarity*. *Semantic similarity* measures the distance between the semantic meanings or semantic content of a pair of keywords [51]. We compute *semantic similarity* using two different mechanisms, Word2Vec [1] and BERT [52]. We calculate $SemSIM_{W2V}(A_i, A_j)$ as the *semantic similarity* score using Word2Vec between the vector embedding of the summary generated by a pair of annotators, $A_i$ and $A_j$ as:

$$SemSIM_{W2V}(A_i, A_j) = \frac{\vec{W_i} \cdot \vec{W_j}}{|\vec{W_i}| \ |\vec{W_j}|} \qquad (3)$$

where, $\vec{W_i}$ and $\vec{W_j}$ are the embeddings of $A_i$ and $A_j$ respectively. We calculate $\vec{W_i}$ as the average of the values of the tweet embeddings. We calculate the embedding for each tweet as the average of the values of the word vector of the tweet. We consider the embedding of each word provided by Word2Vec [1], which was trained on 52M disaster-related messages from various types of disasters. We consider only nouns, verbs, and adjectives as words [49]. As shown in Table 10, our observation indicates that $SemSIM_{W2V}(A_i, A_j)$ ranges between $98 - 99\%$, which indicates that the summaries written by different annotators are highly semantic similar to each other. Similarly, we calculate $SemSIM_{BERT}(A_i, A_j)$ as the semantic similarity score using BERT between the vector embedding of the summary generated by a pair of annotators, $A_i$ and $A_j$ as:

$$SemSIM_{BERT}(A_i, A_j) = \frac{\vec{B_i} \cdot \vec{B_j}}{|\vec{B_i}| \ |\vec{B_j}|} \qquad (4)$$

where, $\vec{B}_i$ and $\vec{B}_j$ are the embedding of $A_i$ and $A_j$ respectively. We calculate $\vec{B}_i$ as the average of the values of the tweet embedding. We consider the embedding for each tweet provided by a pre-trained language model known as BERT [52]. We use an uncased version of the BERT-base model with default hyperparameters. As shown in Table 10, our observation indicates that $SemSIM_{BERT}(A_i, A_j)$ ranges above 99%, which indicates that the summaries written by different annotators are highly semantically similar to each other.

## 5.3 Visualization of the Annotated Summaries

To visualize the content of the summaries provided by three annotators for a dataset, we use Word Cloud[13] based representation of the tweets in the summary [53]. For a tweet summary, the word cloud represents only the most frequent words in the summary with the font size being directly proportional to the word frequency. We generate the word cloud of the summary generated by each annotator to understand whether the most frequent words are similar across annotators. We show Word Clouds of the summaries of $A_1$, $A_2$, and $A_3$ with original tweet set for $D_1$ and $D_2$ in Figure 2. Our observations indicate that the occurrence of the most frequent words is similar across annotators and the original tweet set. From these, we can say that the summaries generated by the annotators cover all the important aspects of the given disaster event.

## 5.4 Impact on Increasing Number of Datasets on Supervised Approaches

In this Subsection, we evaluate the performance of a supervised approach when it is trained on an expanded set of datasets that includes our newly introduced datasets along with the available datasets. To show this, we consider two supervised disaster tweet summarization approaches, IKDSumm proposed by Garg et al. [32] and TSSuBERT proposed by Dusart et al. [10]. We train IKDSumm and TSSuBERT in 2 different ways: 1) train with available disaster datasets (shown in Table 1) and 2) an expanded set of datasets that includes our newly introduced dataset with the available disaster datasets. Then, we evaluate the performance of the summary generated by both ways in terms of the ROUGE-N, i.e., N=1, 2, and L, F1-score. For our experiment, we randomly select four different disaster datasets: $D_1$, $D_3$, $D_4$, and $D_6$. Our results, as shown in Table 11, indicate that IKDSumm performs better by $8.33 - 10.71\%$, $19.04 - 27.77\%$, $12 - 16.67\%$ for ROUGE-1, ROUGE-2, and ROUGE-L, respectively, when it is trained on an expanded set of datasets than the available datasets. Similarly, TSSuBERT performs better by $8 - 9.43\%$, $16 - 25\%$, and $13.33 - 16.67\%$ for ROUGE-1, ROUGE-2, and ROUGE-L, respectively, when it is trained on an expanded set of datasets than the available datasets. Therefore, based on the above experiment, we show our newly introduced datasets improve the performance of the supervised approaches.

---

[13]https://en.wikipedia.org/wiki/Tag_cloud

**Fig. 2**: Figure shows the word cloud of annotated summaries of 3 annotators and the original tweet set for $D_1$ and $D_2$ disasters, such as $A_1$ of $D_1$ in Figure 2a, $A_2$ of $D_1$ in Figure 2b, $A_3$ of $D_1$ in Figure 2c, tweet set of $D_1$ in Figure 2d, $A_1$ of $D_2$ in Figure 2e, $A_2$ of $D_2$ in Figure 2f, $A_3$ of $D_2$ in Figure 2g, and tweet set of $D_2$ in Figure 2h.

## 5.5  Case Study: Utilization of the Datasets Features in Different Tasks

In this Subsection, we study how the additional features of the proposed datasets are helpful in the various NLP tasks. These dataset features are *keyphrases*, *category labels*, and *relevance labels*. We already discussed how the

**Table 11**: Table shows F1-score of ROUGE-1, 2, and L of the summary generated by IKDSumm and TSSuBERT trained on an expanded set of datasets that includes our newly introduced datasets with the available datasets (i.e., IKDSumm-ExpData and TSSuBERT-ExpData) and trained on available disaster datasets (i.e., IKDSumm-AvlData and TSSuBERT-AvlData) for four disasters: $D_1$, $D_3$, $D_4$, and $D_6$.

| Dataset | Approach | ROUGE-1 | ROUGE-2 | ROUGE-L |
|---------|----------|---------|---------|---------|
| $D_1$ | IKDSumm-ExpData | **0.60** | **0.26** | **0.32** |
|       | IKDSumm-AvlData | 0.55 | 0.21 | 0.27 |
| $D_3$ | IKDSumm-ExpData | **0.54** | **0.21** | **0.30** |
|       | IKDSumm-AvlData | 0.49 | 0.17 | 0.25 |
| $D_4$ | IKDSumm-ExpData | **0.57** | **0.18** | **0.25** |
|       | IKDSumm-AvlData | 0.52 | 0.14 | 0.22 |
| $D_6$ | IKDSumm-ExpData | **0.56** | **0.18** | **0.26** |
|       | IKDSumm-AvlData | 0.50 | 0.13 | 0.22 |

| Dataset | Approach | ROUGE-1 | ROUGE-2 | ROUGE-L |
|---------|----------|---------|---------|---------|
| $D_1$ | TSSuBERT-ExpData | **0.57** | **0.25** | **0.30** |
|       | TSSuBERT-AvlData | 0.52 | 0.21 | 0.26 |
| $D_3$ | TSSuBERT-ExpData | **0.53** | **0.17** | **0.27** |
|       | TSSuBERT-AvlData | 0.48 | 0.13 | 0.23 |
| $D_4$ | TSSuBERT-ExpData | **0.50** | **0.14** | **0.24** |
|       | TSSuBERT-AvlData | 0.46 | 0.11 | 0.20 |
| $D_6$ | TSSuBERT-ExpData | **0.49** | **0.16** | **0.25** |
|       | TSSuBERT-AvlData | 0.45 | 0.12 | 0.21 |

*relevance labels* are helpful for summary quality evaluation in Subsection 5.1. We next discuss the utilization of *category labels* in Subsection 5.5.1 and the utilization of *key-phrases* in Subsection 5.5.2.

### 5.5.1 Utilization of Category Labels

In this Subsection, we demonstrate how the *category labels* provided with the proposed datasets can be useful in the various NLP tasks, like category identification of the tweet and determining the coverage quality of the summary. As we have already discussed how the *category labels* are helpful for summary coverage quality evaluation in Subsection 5.1, we next discuss its role for category identification.

### Identification of category of the tweets:

The classification or categorization of disaster tweets is heavily used in literature for different purposes, such as category-specific summary generation [11, 37], category-specific task assessment [54], identifying the help request in a disaster [55], planing post-disaster relief strategies [56] and so on. To evaluate the effectiveness of a tweet categorization approach, we require ground-truth *category labels*. In the literature, we found both the supervised and unsupervised state-of-the-art approaches to identifying a category of a tweet. The supervised approach is BERT-based, whereas the unsupervised approach is graph-based and ontology-based. We select one prominent approach from each of the supervised and unsupervised approaches proposed by [43] and [9, 33], respectively. For our experiment, we select 20% of the total

tweets randomly for five disaster events, such as $D_1$, $D_3$, $D_5$, and $D_7$. We then compare the identified category of the tweets by the approaches proposed by [9], [43], and [33] with the corresponding ground-truth *category labels* in the proposed dataset in terms of F1-score, as shown in Table 12. We observe that F1-scores are in the range of $0.984 - 0.997$, $0.7169 - 0.9077$, and $0.5330 - 0.7490$ for approaches used in [9], [43] and [33]. Therefore, based on the above experiment, we show the utilization of the *category label* field of the proposed datasets and can say it is beneficial in the category identification task.

**Table 12**: Table shows F1-score of the category identification approaches proposed by [9], [43] and [33], on comparing the ground-truth *category label* annotations for five disasters: $D_1$, $D_3$, $D_5$, and $D_7$.

| Dataset | Approach | F1-score |
|---------|----------|----------|
| | Approach used in [9] | **0.9950** |
| $D_1$ | Approach used in [43] | 0.8409 |
| | Approach used in [33] | 0.5330 |
| | Approach used in [9] | **0.9840** |
| $D_3$ | Approach used in [43] | 0.9077 |
| | Approach used in [33] | 0.5710 |
| | Approach used in [9] | **0.9970** |
| $D_5$ | Approach used in [43] | 0.8395 |
| | Approach used in [33] | 0.7490 |
| | Approach used in [9] | **0.9910** |
| $D_7$ | Approach used in [43] | 0.7169 |
| | Approach used in [33] | 0.5840 |

### 5.5.2 Utilization of Key-phrases

In this Subsection, we demonstrate the usefulness of *key-phrases* provided with the proposed datasets in development of robust summarization algorithm, determination of the *Diversity* in summary and the necessity of *key-phrases* highlighting in summary. We discuss the role of *key-phrases* for summary diversity quality evaluation in Subsection 5.1. We next discuss the utilization of *key-phrases* for the remaining two tasks.

### Development of a summarization algorithm:

Existing summarization approaches inherently perform two steps for creating a summary from input tweets: 1) ranking of the tweets based on the importance of tweets, and 2) selection of important tweets into the summary. For the ranking of the tweets, existing work proposed by Dusart et al. [10] integrate the tweet's word frequencies with the corresponding tweet embedding identified using the DistilBERT model to determine each tweet's importance score. With a huge training set, this approach can automatically find out important keywords and in turns can compute the importance score of tweet efficiently. However, with limited dataset this model fails to compute the importance score efficiently. Therefore, availability of *key-phrases* can aid in determination of the tweet's importance even with a sparse disaster dataset. To empirically

validate this, we consider two approaches, such as, Garg et al. [32] and Dusart et al. [10] where [32] utilizes the *key-phrases* to identify the tweet's importance and [10] does not utilizes the *key-phrases*. For our experiment, we randomly select 2% of the tweets from the $D_1$ dataset and compute the importance score of a tweet by [32] and [10] approaches, respectively. We rank these tweets along with the importance score in descending order and request a meta annotator to score which ranked list comprises of more important tweets. However, the metaannotator has no knowledge of the underlying approaches for tweet importance selection and also, which list belongs to which approach. On the basis of the meta-annotator's decision, we observe that *key-phrases* helps to identify the importance of a tweet more efficiently.

For our selection of tweets into the summary, we consider : 1) Importance of tweets from ranking and 2) *Diversity* among the already selected tweets in the summary. To ensure *Diversity* in summary tweets, we can utilize their *key-phrases*. In order to evaluate the importance of *key-phrases* in diversity calculation which effectively increases the summary quality, we follow two different variants of IKDSumm  [32] by considering *Diversity* and without *Diversity* in tweet selection into the summary. While considering *Diversity*, we iteratively select tweets with the maximum importance score and have minimum similarity (maximum diversity) with the already selected tweets into the summary. For *Diversity* calculation, we utilize *key-phrases* as shown in IKD-Summ  [32]. For the variant without considering *Diversity*, we iteratively select tweets on the basis of only the importance score into the summary. For our experiment, we consider randomly selected four different disaster datasets, i.e., $D_1$, $D_3$, $D_5$, and $D_6$ to evaluate the performance of the summary generated by both the methods with the ground-truth summary using ROUGE-N F1-score. Our observations, as shown in Table 13, indicate that the improvement in summary score of ROUGE-N F1-score ranges from $5.67 - 36.84\%$ while we with considering *key-phrases* for summary creation than without *key-phrases*. Therefore, based on the above experiments, we show that our introduced *key-phrases* filed with the proposed dataset is beneficial in developing a robust summarization approach for both the ranking as well as the selection of tweets into the summary.

**Table 13**: Table shows F1-score of ROUGE-1, 2 and L of the summary generated with considering *Diversity* (IKDSumm-withDiv) and without considering *Diversity* (IKDSumm-withoutDiv) by [10] on four disasters: $D_1$, $D_3$, $D_5$, and $D_6$.

| Dataset | Approach | ROUGE-1 | ROUGE-2 | ROUGE-L |
|---------|----------|---------|---------|---------|
| | IKDSumm-withDiv | **0.59** | **0.25** | **0.32** |
| $D_1$ | IKDSumm-withoutDiv | 0.55 | 0.21 | 0.28 |
| | IKDSumm-withDiv | **0.53** | **0.20** | **0.29** |
| $D_3$ | IKDSumm-withoutDiv | 0.50 | 0.16 | 0.25 |
| | IKDSumm-withDiv | **0.56** | **0.19** | **0.25** |
| $D_5$ | IKDSumm-withoutDiv | 0.50 | 0.12 | 0.21 |
| | IKDSumm-withDiv | **0.57** | **0.18** | **0.26** |
| $D_6$ | IKDSumm-withoutDiv | 0.49 | 0.14 | 0.22 |

**Importance of key-phrases in understanding summary:**

Rudra et al. [37] discusses that *key-phrases* makes a summary more interpretable for users and further, provide the reasoning behind the selection of a particular tweet in the summary. In order to understand the necessity of *key-phrases* in summary, we conduct an experiment where we provide two summaries 1) without *key-phrases* and 2) with *key-phrases* of each tweet to 5 meta-annotators and ask them to select the best summary out of these two summaries. For this experiment, we utilize *IKDSumm* summarization approach [32] to generate the summary and datasets as $D_1$, $D_4$, and $D_7$. Our observations, as shown in Table 14 indicate that a huge majority of meta-annotators ($D_1$ - 80%, $D_4$ - 100%, $D_7$ - 100%) found highlighting *key-phrases* provide better explanations for inclusion of a tweet into the summary compared to the other. Therefore, the identified *key-phrases* field with the proposed datasets will certainly help the research community as it provides additional important information about the summary.

**Table 14**: Table shows results of meta-annotators preference in percentage (AnnoPrefPer) of the summary generated by IKDSumm with highlighted *key-phrases* and without highlighted *key-phrases* for $D_1$, $D_4$, and $D_7$ datasets over 3 annotators.

| Dataset | Approach | AnnoPrefPer |
|---------|----------|-------------|
| $D_1$ | IKDSumm-withHighlight | 80.00 |
|  | IKDSumm-withoutHighlight | 20.00 |
| $D_4$ | IKDSumm-withHighlight | 100 |
|  | IKDSumm-withoutHighlight | NA |
| $D_7$ | IKDSumm-withHighlight | 100 |
|  | IKDSumm-withoutHighlight | NA |

## 5.6 Case Study: Evaluation of Existing Summarization Approaches

In this Subsection, we initially discuss the various existing state-of-the-art summarization approaches, followed by a comparison of the annotated ground-truth summaries for 8 disaster tweet datasets.

### 5.6.1 Existing Summarization Approaches

We categorize existing summarization approaches into *content-based, graph-based, matrix factorization-based, entropy-based, semantic similarity-based* and *deep learning-based* approaches. We select a few prominent tweet summarization approaches from each type which we discuss next.

1. *Content-based Approaches:* We discuss the existing content-based summarization approaches, which are as follows:

(a) *LUHN*: Luhn et al. [57] propose a frequency-based summarization approach that selects those tweets into a summary that has the highest frequency scoring words.

(b) *SumBasic*: Nenkova et al. [58] select those tweets into a summary that have the words with the maximum probability of occurrence.

(c) *COWTS*: Rudra et al. [16] select those tweets into a summary that has the maximum coverage of the content words (i.e., noun, main verb, and numerals).

(d) *DEPSUB*: Rudra et al. [37] propose a sub-events-based summarization approach that initially identifies the noun-verb pairs as the sub-events in each input tweet and then creates a summary by selecting the representative tweets by maximizing the information coverage of the disaster-specific keywords and the sub-events using Integer Linear Programming (ILP) based selection.

2. *Graph-based Approaches:* We discuss the existing graph-based summarization approaches, which are as follows:

(a) *Cluster Rank*: Garg et al. [59] initially identify the different clusters followed by utilizing the PageRank [60] algorithm to select tweets from each cluster in summary.

(b) *LexRank*: Erkan et al. [61] initially construct a sentence graph where nodes represent the sentences, and the edges are the content similarity between them and then determine Eigenvector [62] centrality score of each node in sentence graph. Finally, they create a summary by selecting the highest eigenvector centrality score into the summary.

(c) *EnSum*: Dutta et al. [12] propose an ensemble graph-based tweet summarization approach, EnSum, which initially selects the most important tweets by using 9 existing summarization algorithms. Then using these important tweets, they create a tweet similarity graph where the nodes represent the tweets, and the edges represent their similarity. They identify the different communities using a community detection algorithm and then create a summary by selecting the representative tweets from each community based on length, informativeness, and centrality scores.

(d) *COWEXABS*: Rudra et al. [11] propose a summarization framework initially identifying the most important tweets by maximizing the information coverage of the disaster-specific keywords in the extracted tweets. Then from these tweets, they create a graph where the nodes are disaster-specific keywords, and the edges represent the co-occurrence relationship between them. Finally, they select the tweets or tweet paths in summary, which can ensure maximum information coverage of the graph.

(e) *MEAD*: Radev et al. [63] propose a centroid-based summarization approach where they initially segregate sentences using an agglomerative clustering approach. Further, they select the tweets from each

cluster on the basis of the centrality score and diversity score in the summary.

3. *Matrix factorization-based Approaches:* We discuss the most popular matrix factorization-based summarization approaches in detail next.
    (a) *LSA*: Gong et al. [64] propose a Latent Semantic Analysis (LSA) based approach, where they select the tweets which have the highest eigenvalue after Singular Value Decomposition (SVD) of the keyword matrix created from all the input tweets.
    (b) *SumDSDR*: He et al. [65] propose a data reconstruction-based summarization approach, where they initially apply linear reconstruction and non-linear reconstruction objective functions to identify the relation between the sentences and then generate a summary by minimizing the reconstruction error.

4. *Entropy-based Approach:* Garg et al. [26] propose an entropy-based disaster tweet summarization approach, EnDSUM, that creates the summary by selecting the tweets which provide the maximum entropy and diversity in the summary.

5. *Semantic Similarity based Approach:* Garg et al. [9] propose a disaster-specific tweet summarization framework, OntoDSumm, which initially determines each tweet category using an ontology-based pseudo-relevance feedback approach followed by identification of the importance of each category which represents the number of tweets to be in summary from a category. Finally, they create a summary by selecting the representative tweets from each category based on the DMMR-based approach.

6. *Deep learning-based Approaches:* We discuss the existing deep learning-based summarization approaches, which are as follows:
    (a) *TSSuBERT:* Dusart et al. [10] propose a summarization framework, TSSuBERT, which integrates the context of the tweets using the whole vocabulary related to the event to determine the tweet importance. Finally, they iteratively select the higher important tweets in summary.
    (b) *GCNSUM:* Li et al. [27] propose a GCN-based summarization framework that initially creates a tweet-similarity graph where the nodes are tweets and edges represent their content similarity. Then, they generate tweet hidden features for each tweet by applying GCN on the tweet similarity graph. Further, they determine the importance score of each tweet by combining the tweet's hidden features and the whole event embedding. Finally, they create a summary by iteratively selecting the most important tweets in the summary.
    (c) *IKDSumm:* Garg et al. [32] propose a summarization framework, IKDSumm, which integrates the disaster-specific key-phrase identified using domain knowledge of ontology with the tweet to determine the tweet importance. Finally, they iteratively select tweets with higher importance and maximum diversity with the tweets already selected in the summary.

**Table 15**: Table shows F1-score of ROUGE-1 (R-1), ROUGE-2 (R-2), and ROUGE-L (R-L) score of the summaries generated by various baselines for $D_1$-$D_4$ datasets.

| Approach | $D_1$ | | | $D_2$ | | | $D_3$ | | | $D_4$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R-1 | R-2 | R-L | R-1 | R-2 | R-L | R-1 | R-2 | R-L | R-1 | R-2 | R-L |
| *ClusterRank* | 0.44 | 0.15 | 0.22 | 0.47 | 0.14 | 0.22 | 0.45 | 0.12 | 0.22 | 0.47 | 0.14 | 0.21 |
| *LexRank* | 0.39 | 0.09 | 0.20 | 0.40 | 0.10 | 0.18 | 0.38 | 0.08 | 0.19 | 0.33 | 0.09 | 0.19 |
| *LSA* | 0.47 | 0.17 | 0.25 | 0.47 | 0.13 | 0.22 | 0.45 | 0.14 | 0.23 | 0.47 | 0.13 | 0.21 |
| *LUHN* | 0.45 | 0.16 | 0.23 | 0.46 | 0.12 | 0.21 | 0.45 | 0.13 | 0.23 | 0.46 | 0.14 | 0.21 |
| *MEAD* | 0.39 | 0.10 | 0.22 | 0.47 | 0.11 | 0.22 | 0.46 | 0.14 | 0.24 | 0.42 | 0.07 | 0.17 |
| *SumBasic* | 0.42 | 0.13 | 0.22 | 0.46 | 0.13 | 0.21 | 0.41 | 0.09 | 0.21 | 0.43 | 0.09 | 0.20 |
| *SumDSDR* | 0.41 | 0.12 | 0.22 | 0.41 | 0.11 | 0.18 | 0.39 | 0.09 | 0.20 | 0.35 | 0.10 | 0.21 |
| *COWTS* | 0.43 | 0.14 | 0.23 | 0.46 | 0.12 | 0.22 | 0.45 | 0.13 | 0.22 | 0.44 | 0.11 | 0.21 |
| *COWEXABS* | 0.49 | 0.22 | 0.29 | 0.48 | 0.13 | 0.22 | 0.45 | 0.13 | 0.23 | 0.20 | 0.04 | 0.20 |
| *DEPSUB* | 0.52 | 0.21 | 0.23 | 0.44 | 0.12 | 0.22 | 0.44 | 0.14 | 0.23 | 0.45 | 0.11 | 0.21 |
| *EnSum* | 0.48 | 0.18 | 0.25 | 0.47 | 0.14 | 0.22 | 0.46 | 0.14 | 0.24 | 0.47 | 0.14 | 0.21 |
| *EnDSUM* | 0.55 | 0.21 | 0.27 | 0.52 | 0.17 | 0.24 | 0.52 | 0.14 | 0.26 | 0.51 | 0.16 | 0.24 |
| *OntoDSumm* | 0.57 | 0.24 | 0.30 | 0.51 | 0.17 | 0.26 | 0.52 | 0.19 | 0.27 | 0.54 | 0.17 | 0.24 |
| *TSSuBERT* | 0.55 | 0.23 | 0.28 | 0.50 | 0.15 | 0.25 | 0.51 | 0.16 | 0.26 | 0.48 | 0.12 | 0.21 |
| *GCNSUM* | 0.52 | 0.20 | 0.27 | 0.50 | 0.16 | 0.25 | 0.50 | 0.17 | 0.25 | 0.51 | 0.16 | 0.22 |
| *IKDSumm* | **0.60** | **0.26** | **0.32** | **0.55** | **0.20** | **0.27** | **0.54** | **0.21** | **0.30** | **0.57** | **0.18** | **0.25** |

**Table 16**: Table shows F1-score of ROUGE-1 (R-1), ROUGE-2 (R-2), and ROUGE-L (R-L) score of the summaries generated by various baselines for $D_5$-$D_8$ datasets.

| Approach | $D_5$ | | | $D_6$ | | | $D_7$ | | | $D_8$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | R-1 | R-2 | R-L | R-1 | R-2 | R-L | R-1 | R-2 | R-L | R-1 | R-2 | R-L |
| *ClusterRank* | 0.46 | 0.11 | 0.21 | 0.45 | 0.13 | 0.22 | 0.50 | 0.21 | 0.31 | 0.44 | 0.12 | 0.25 |
| *LexRank* | 0.41 | 0.07 | 0.14 | 0.39 | 0.08 | 0.17 | 0.44 | 0.14 | 0.25 | 0.37 | 0.08 | 0.20 |
| *LSA* | 0.47 | 0.12 | 0.22 | 0.45 | 0.10 | 0.20 | 0.47 | 0.15 | 0.26 | 0.47 | 0.14 | 0.22 |
| *LUHN* | 0.44 | 0.08 | 0.18 | 0.48 | 0.12 | 0.21 | 0.49 | 0.17 | 0.29 | 0.48 | 0.15 | 0.23 |
| *MEAD* | 0.49 | 0.13 | 0.22 | 0.44 | 0.09 | 0.20 | 0.45 | 0.14 | 0.28 | 0.42 | 0.12 | 0.22 |
| *SumBasic* | 0.42 | 0.07 | 0.16 | 0.49 | 0.12 | 0.19 | 0.55 | 0.24 | 0.31 | 0.46 | 0.12 | 0.21 |
| *SumDSDR* | 0.43 | 0.12 | 0.21 | 0.47 | 0.15 | 0.23 | 0.48 | 0.16 | 0.27 | 0.47 | 0.17 | 0.27 |
| *COWTS* | 0.43 | 0.09 | 0.19 | 0.50 | 0.16 | 0.20 | 0.46 | 0.17 | 0.26 | 0.47 | 0.14 | 0.21 |
| *COWEXABS* | 0.19 | 0.04 | 0.18 | 0.47 | 0.15 | 0.22 | 0.49 | 0.18 | 0.30 | 0.46 | 0.12 | 0.24 |
| *DEPSUB* | 0.50 | 0.12 | 0.22 | 0.43 | 0.11 | 0.19 | 0.49 | 0.17 | 0.28 | 0.44 | 0.10 | 0.20 |
| *EnSum* | 0.48 | 0.10 | 0.20 | 0.42 | 0.09 | 0.20 | 0.49 | 0.17 | 0.29 | 0.47 | 0.13 | 0.24 |
| *EnDSUM* | 0.52 | 0.13 | 0.24 | 0.47 | 0.13 | 0.21 | 0.49 | 0.21 | 0.29 | 0.48 | 0.15 | 0.25 |
| *OntoDSumm* | 0.55 | 0.17 | 0.24 | 0.53 | 0.17 | 0.25 | 0.56 | 0.25 | 0.33 | 0.49 | 0.18 | 0.28 |
| *TSSuBERT* | 0.50 | 0.12 | 0.21 | 0.47 | 0.15 | 0.23 | 0.50 | 0.18 | 0.29 | 0.40 | 0.12 | 0.22 |
| *GCNSUM* | 0.46 | 0.10 | 0.20 | 0.46 | 0.14 | 0.19 | 0.51 | 0.21 | 0.31 | 0.42 | 0.13 | 0.22 |
| *IKDSumm* | **0.56** | **0.19** | **0.25** | **0.56** | **0.18** | **0.26** | **0.57** | **0.25** | **0.34** | **0.52** | **0.20** | **0.29** |

### 5.6.2 Comparison Results and Discussions

To evaluate the performance of the generated summaries of different baselines, we compare it with the ground-truth summaries using ROUGE-N [66] scores. The ROUGE-N score is a widely recognized metric in text summarization tasks which calculates the score by comparing the number of overlapping words between the system-generated and reference (or ground-truth) summaries. We use F1-score for 3 different variants of the ROUGE-N score, i.e., N=1, 2, and L, respectively. Our observations as shown in Table 15 and 16 indicate that IKDSumm ensures best ROUGE-N F1-score on $D_1 - D_8$ datasets in comparison with different baselines followed by OntoDSumm. The reason for the high performance of IKDSumm is that it utilizes existing ontology knowledge instead of labelled training data to identify key-phrase and then utilizes these

key-phrase to determine the tweet's importance to create a summary. Further, we observe that after IKDSumm, OntoDSumm ensures the best ROUGE-N F1-score for $D_1 - D_8$, followed by TSSuBERT. The reason for the high performance of OntoDSumm is that it ensures each category's representation in a summary and handles the information diversity in summary tweets from each category. The performance of Lex Rank is the worst except for $D_4 - D_5$, and COWEXABS is the worst for $D_4 - D_5$ as both of these approaches do not ensure the category representation and information diversity in summary.

## 6 Conclusions

In this study, we introduce a compilation of annotated datasets designed for the purpose of summarizing disaster-related tweets. Annotated ground-truth summary aids in enhancing the efficiency of supervised learning methods for both training and evaluation processes. The *ADSumm* dataset comprises of tweets and corresponding ground-truth summaries for a total of eight distinct disaster events. These events encompass both man-made and natural disasters, originating from seven diverse geographical areas/country. In addition, this study also provides three supplementary components in conjunction with the datasets: *category labels*, *key-phrases*, and *relevance labels*. The *category label* is used to ensure that all categories are adequately represented in the summary. The *relevance label* is employed to assess the quality of the summary. Lastly, *key-phrases* are utilized to provide justification or explanation for the inclusion of a specific tweet in the summary. The inclusion of the recently incorporated datasets for training leads to a notable enhancement of $8 - 28\%$ in the ROUGE-N F1-score for the supervised summarization methods. This study also presents a detailed methodology for preparing the ground-truth summary. Additionally, we discuss the practical applications of the supplementary features within the datasets for various Natural Language Processing (NLP) tasks. Moreover, we provide experimental results (quantitative as well as qualitative) to ensure the quality of the annotated ground-truth summary. In addition, we evaluate the effectiveness of different state-of-the-art summarization methods on these datasets by measuring their performance using the ROUGE-N F1-score. It is anticipated that the utilization of these datasets will contribute to the advancement of more resilient and effective catastrophe tweet summarizing methods. Consequently, this will facilitate the provision of assistance by humanitarian groups and government authorities.

## References

[1] Imran, M., Mitra, P., Castillo, C.: Twitter as a lifeline: Human-annotated twitter corpora for nlp of crisis-related messages. arXiv preprint arXiv:1605.05894 (2016)

[2] Basu, M., Shandilya, A., Khosla, P., Ghosh, K., Ghosh, S.: Extracting resource needs and availabilities from microblogs for aiding post-disaster

relief operations. IEEE Transactions on Computational Social Systems **6**(3), 604–618 (2019)

[3] Priya, S., Sequeira, R., Chandra, J., Dandapat, S.K.: Where should one get news updates: Twitter or reddit. Online Social Networks and Media **9**, 17–29 (2019)

[4] Dutt, R., Basu, M., Ghosh, K., Ghosh, S.: Utilizing microblogs for assisting post-disaster relief operations via matching resource needs and availabilities. Information Processing & Management **56**(5), 1680–1697 (2019)

[5] Ghosh, S., Ghosh, K., Ganguly, D., Chakraborty, T., Jones, G.J., Moens, M.-F., Imran, M.: Exploitation of social media for emergency relief and preparedness: Recent research and trends. Information Systems Frontiers **20**(5), 901–907 (2018)

[6] Castillo, C.: Big Crisis Data: Social Media in Disasters and Time-critical Situations. Cambridge University Press, Cambridge, England (2016)

[7] Alam, F., Ofli, F., Imran, M.: Descriptive and visual summaries of disaster events using artificial intelligence techniques: case studies of hurricanes harvey, irma, and maria. Behaviour & Information Technology **39**(3), 288–318 (2020)

[8] Alam, F., Qazi, U., Imran, M., Ofli, F.: Humaid: Human-annotated disaster incidents data from twitter with deep learning benchmarks. In: Proceedings of the International AAAI Conference on Web and Social Media, vol. 15, pp. 933–942. AAAI Press, Palo Alto, California, USA (2021)

[9] Garg, P.K., Chakraborty, R., Dandapat, S.K.: Ontodsumm: Ontology-based tweet summarization for disaster events. IEEE Transactions on Computational Social Systems (2023)

[10] Dusart, A., Pinel-Sauvagnat, K., Hubert, G.: Tssubert: How to sum up multiple years of reading in a few tweets. ACM Transactions on Information Systems (2023)

[11] Rudra, K., Goyal, P., Ganguly, N., Imran, M., Mitra, P.: Summarizing situational tweets in crisis scenarios: An extractive-abstractive approach. IEEE Transactions on Computational Social Systems **6**(5), 981–993 (2019)

[12] Dutta, S., Chandra, V., Mehra, K., Das, A.K., Chakraborty, T., Ghosh, S.: Ensemble algorithms for microblog summarization. IEEE Intelligent Systems **33**(3), 4–14 (2018)

[13] Rudra, K., Ganguly, N., Goyal, P., Ghosh, S.: Extracting and summarizing situational information from the twitter social media during disasters. ACM Transactions on the Web (TWEB) **12**(3), 1–35 (2018)

[14] Olteanu, A., Vieweg, S., Castillo, C.: What to expect when the unexpected happens: Social media communications across crises. In: Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing, pp. 994–1009. ACM, New York, NY, USA (2015)

[15] Alam, F., Ofli, F., Imran, M.: Crisismmd: Multimodal twitter datasets from natural disasters. In: Proceedings of the 12th International AAAI Conference on Web and Social Media (ICWSM). AAAI Press, Palo Alto, California, USA (2018)

[16] Rudra, K., Ghosh, S., Ganguly, N., Goyal, P., Ghosh, S.: Extracting situational information from microblogs during disaster events: A classification-summarization approach. In: Proceedings of the 24th ACM International on Conference on Information and Knowledge Management. CIKM '15, pp. 583–592. ACM, New York, NY, USA (2015). https://doi.org/10.1145/2806416.2806485. https://doi.org/10.1145/2806416.2806485

[17] Garg, P.K., Chakraborty, R., Dandapat, S.K.: PORTRAIT: a hybrid aPproach tO cReate extractive ground-TRuth summAry for dIsaster evenT (2023)

[18] Nazari, N., Mahdavi, M.: A survey on automatic text summarization. Journal of AI and Data Mining **7**(1), 121–135 (2019)

[19] Jain, D., Borah, M.D., Biswas, A.: Bayesian optimization based score fusion of linguistic approaches for improving legal document summarization. Knowledge-Based Systems **264**, 110336 (2023)

[20] Duan, Y., Jatowt, A.: Across-time comparative summarization of news articles. In: Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, pp. 735–743. ACM, New York, NY, USA (2019)

[21] Ansah, J., Liu, L., Kang, W., Kwashie, S., Li, J., Li, J.: A graph is worth a thousand words: Telling event stories using timeline summarization graphs. In: The World Wide Web Conference, pp. 2565–2571. ACM, New York, NY, USA (2019)

[22] Chakraborty, R., Bhavsar, M., Dandapat, S.K., Chandra, J.: Tweet summarization of news articles: An objective ordering-based perspective. IEEE Transactions on Computational Social Systems **6**(4), 761–777 (2019)

[23] Chakraborty, R., Bhavsar, M., Dandapat, S., Chandra, J.: A network based stratification approach for summarizing relevant comment tweets of news articles. In: International Conference on Web Information Systems Engineering, pp. 33–48. Springer, New York, NY, USA (2017)

[24] Roy, S., Mishra, S., Matam, R.: Classification and summarization for informative tweets. In: 2020 IEEE International Students' Conference on Electrical, Electronics and Computer Science (SCEECS), pp. 1–4. IEEE, New York, NY, USA (2020)

[25] Dutta, S., Chandra, V., Mehra, K., Ghatak, S., Das, A.K., Ghosh, S.: Summarizing microblogs during emergency events: A comparison of extractive summarization algorithms. In: Emerging Technologies in Data Mining and Information Security, pp. 859–872. Springer, New York, NY, USA (2019)

[26] Garg, P.K., Chakraborty, R., Dandapat, S.K.: Endsum: Entropy and diversity based disaster tweet summarization. In: Proceedings of Text2Story - Fifth Workshop on Narrative Extraction From Texts Held in Conjunction with the 44th European Conference on Information Retrieval (ECIR 2022), Stavanger, Norway, April 10, 2022, vol. 3117, pp. 91–96 (2022)

[27] Li, Q., Zhang, Q.: Twitter event summarization by exploiting semantic terms and graph network. In: Proceedings of the The Thirty-Third Annual Conference on Innovative Applications of Artificial Intelligence (IAAI-21), vol. 35, pp. 15347–15354. AAAI Press, Palo Alto, California, USA (2021)

[28] Dutta, S., Ghatak, S., Roy, M., Ghosh, S., Das, A.K.: A graph based clustering technique for tweet summarization. In: 2015 4th International Conference on Reliability, Infocom Technologies and Optimization (ICRITO)(trends and Future Directions), pp. 1–6. IEEE, New York, NY, USA (2015)

[29] De Maio, C., Fenza, G., Gallo, M., Loia, V., Parente, M.: Time-aware adaptive tweets ranking through deep learning. Future Generation Computer Systems **93**, 924–932 (2019)

[30] Rudra, K., Banerjee, S., Ganguly, N., Goyal, P., Imran, M., Mitra, P.: Summarizing situational tweets in crisis scenario. In: Proceedings of the 27th ACM Conference on Hypertext and Social Media. HT '16, pp. 137–147. ACM, New York, NY, USA (2016). https://doi.org/10.1145/2914586.2914600. https://doi.org/10.1145/2914586.2914600

[31] Devlin, J., Chang, M.-W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018)

[32] Garg, P.K., Chakraborty, R., Gupta, S., Dandapat, S.K.: Ikdsumm: Incorporating key-phrases into bert for extractive disaster tweet summarization. Computer Speech & Language **87**, 101649 (2024). https://doi.org/10.1016/j.csl.2024.101649

[33] Dutta, S., Das, A.K., Bhattacharya, A., Dutta, G., Parikh, K.K., Das, A., Ganguly, D.: Community detection based tweet summarization. In: Emerging Technologies in Data Mining and Information Security, pp. 797–808. Springer, New York, NY, USA (2019)

[34] Fortunato, S.: Community detection in graphs. Physics reports **486**(3-5), 75–174 (2010)

[35] Gazit, H.: An optimal randomized parallel algorithm for finding connected components in a graph. SIAM Journal on Computing **20**(6), 1046–1067 (1991)

[36] Kim, T.-Y., Kim, J., Lee, J., Lee, J.-H.: A tweet summarization method based on a keyword graph. In: Proceedings of the 8th International Conference on Ubiquitous Information Management and Communication, pp. 1–8. ACM, New York, NY, USA (2014)

[37] Rudra, K., Goyal, P., Ganguly, N., Mitra, P., Imran, M.: Identifying sub-events and summarizing disaster-related information from microblogs. In: The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval, pp. 265–274. ACM, New York, NY, USA (2018)

[38] Imran, M., Castillo, C., Lucas, J., Meier, P., Vieweg, S.: Aidr: Artificial intelligence for disaster response. In: Proceedings of the 23rd International Conference on World Wide Web, pp. 159–162. ACM, New York, NY, USA (2014)

[39] Arachie, C., andSam Anzaroot, M.G., Groves, W., Zhang, K., Jaimes, A.: Unsupervised detection of sub-events in large scale disasters. In: The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational-Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pp. 354–361. AAAI Press, Palo Alto, California, USA (2020). https://aaai.org/ojs/index.php/AAAI/article/view/5370

[40] Poddar, S., Samad, A.M., Mukherjee, R., Ganguly, N., Ghosh, S.: Caves: A dataset to facilitate explainable classification and summarization of concerns towards covid vaccines. In: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 3154–3164. Association for Computing Machinery,

New York, NY, USA (2022). https://doi.org/10.1145/3477495.3531745. https://doi.org/10.1145/3477495.3531745

[41] Imran, M., Castillo, C.: Towards a data-driven approach to identify crisis-related topics in social media streams. In: Proceedings of the 24th International Conference on World Wide Web, pp. 1205–1210. ACM, New York, NY, USA (2015)

[42] Gaur, M., Shekarpour, S., Gyrard, A., Sheth, A.: empathi: An ontology for emergency managing and planning about hazard crisis. In: 2019 IEEE 13th International Conference on Semantic Computing (ICSC), pp. 396–403. IEEE, New York, NY, USA (2019)

[43] Nguyen, T.H., Rudra, K.: Towards an interpretable approach to classify and summarize crisis events from microblogs. In: Proceedings of the ACM Web Conference 2022, pp. 3641–3650 (2022)

[44] Nguyen, T.H., Rudra, K.: Rationale aware contrastive learning based approach to classify and summarize crisis-related microblogs. In: Proceedings of the 31st ACM International Conference on Information & Knowledge Management, pp. 1552–1562 (2022)

[45] Fleiss, J.L., Levin, B., Paik, M.C.: Statistical Methods for Rates and Proportions. john wiley & sons, New Jersey, USA (2013)

[46] Landis, J.R., Koch, G.G.: The measurement of observer agreement for categorical data. biometrics **33**, 159–174 (1977)

[47] Little, C., Mclean, D., Crockett, K., Edmonds, B.: A semantic and syntactic similarity measure for political tweets. IEEE Access **8**, 154095–154113 (2020)

[48] Nguyen, H.V., Bai, L.: Cosine similarity metric learning for face verification. In: Asian Conference on Computer Vision, pp. 709–720. Springer, New York, NY, USA (2010)

[49] Khan, M.A.H., Bollegala, D., Liu, G., Sezaki, K.: Multi-tweet summarization of real-time events. In: 2013 International Conference on Social Computing, pp. 128–133. IEEE, New York, NY, USA (2013)

[50] Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. the Journal of machine Learning research **3**, 993–1022 (2003)

[51] Miller, G.A., Charles, W.G.: Contextual correlates of semantic similarity. Language and cognitive processes **6**(1), 1–28 (1991)

[52] Liu, Y., Lapata, M.: Text summarization with pretrained encoders. In:

Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 3730–3740. Association for Computational Linguistics, Stroudsburg, USA (2019)

[53] Castella, Q., Sutton, C.: Word storms: Multiples of word clouds for visual comparison of documents. In: Proceedings of the 23rd International Conference on World Wide Web, pp. 665–676. ACM, New York, NY, USA (2014)

[54] Priya, S., Bhanu, M., Dandapat, S.K., Ghosh, K., Chandra, J.: Taqe: tweet retrieval-based infrastructure damage assessment during disasters. IEEE transactions on computational social systems **7**(2), 389–403 (2020)

[55] Ullah, I., Khan, S., Imran, M., Lee, Y.-K.: Rweetminer: Automatic identification and categorization of help requests on twitter during disasters. Expert Systems with Applications **176**, 114787 (2021)

[56] Rudra, K., Sharma, A., Ganguly, N., Imran, M.: Classifying information from microblogs during epidemics. In: Proceedings of the 2017 International Conference on Digital Health, pp. 104–108 (2017)

[57] Luhn, H.P.: The automatic creation of literature abstracts. IBM Journal of research and development **2**(2), 159–165 (1958)

[58] Nenkova, A., Vanderwende, L.: The impact of frequency on summarization. Microsoft Research, Redmond, Washington, Tech. Rep. MSR-TR-2005 **101** (2005)

[59] Garg, N., Favre, B., Reidhammer, K., Hakkani Tür, D.: Clusterrank: a graph based method for meeting summarization. Technical report, Idiap (2009)

[60] Page, L., Brin, S., Motwani, R., Winograd, T.: The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab (1999)

[61] Erkan, G., Radev, D.R.: Lexrank: Graph-based lexical centrality as salience in text summarization. Journal of artificial intelligence research **22**, 457–479 (2004)

[62] Borgatti, S.P.: Centrality and network flow. Social networks **27**(1), 55–71 (2005)

[63] Radev, D.R., Allison, T., Blair-Goldensohn, S., Blitzer, J., Celebi, A., Dimitrov, S., Drabek, E., Hakim, A., Lam, W., Liu, D., *et al.*: Mead-a

platform for multidocument multilingual text summarization. In: Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04). European Language Resources Association (ELRA), Lisbon, Portugal (2004)

[64] Gong, Y., Liu, X.: Generic text summarization using relevance measure and latent semantic analysis. In: Proceedings of the 24th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 19–25. ACM, New York, NY, USA (2001)

[65] He, Z., Chen, C., Bu, J., Wang, C., Zhang, L., Cai, D., He, X.: Document summarization based on data reconstruction. In: Twenty-sixth AAAI Conference on Artificial Intelligence, pp. 620–626. AAAI Press, Palo Alto, California, USA (2012)

[66] Lin, C.-Y.: Rouge: A package for automatic evaluation of summaries. In: Text Summarization Branches Out, pp. 74–81. Association for Computational Linguistics, Stroudsburg, USA (2004)