# Deep Learning Benchmarks and Datasets for Social Media Image Classification for Disaster Response

Firoj Alam[1], Ferda Ofli[1], Muhammad Imran[1], Tanvirul Alam[2], Umair Qazi[1]

[1]Qatar Computing Research Institute, HBKU, Doha, Qatar

[2]BJIT Limited, Dhaka, Bangladesh

[1]{fialam, fofli, mimran, uqazi}@hbku.edu.qa, [2]tanvirul.alam@bjitgroup.com

*Abstract*—During a disaster event, images shared on social media helps crisis managers gain situational awareness and assess incurred damages, among other response tasks. Recent advances in computer vision and deep neural networks have enabled the development of models for real-time image classification for a number of tasks, including detecting crisis incidents, filtering irrelevant images, classifying images into specific humanitarian categories, and assessing the severity of damage. Despite several efforts, past works mainly suffer from limited resources (i.e., labeled images) available to train more robust deep learning models. In this study, we propose new datasets for disaster type detection, and informativeness classification, and damage severity assessment. Moreover, we relabel existing publicly available datasets for new tasks. We identify exact- and near-duplicates to form non-overlapping data splits, and finally consolidate them to create larger datasets. In our extensive experiments, we benchmark several state-of-the-art deep learning models and achieve promising results. We release our datasets and models publicly, aiming to provide proper baselines as well as to spur further research in the crisis informatics community.

*Index Terms*—Deep learning, Disaster Image Classification, Natural disasters, Crisis computing, Social media, Benchmarking

## I. INTRODUCTION

Social media is widely used during natural or human-induced disasters as a source to quickly disseminate information and learn useful insights. People post content (i.e., through different modalities such as text, image, and video) on social media to get help and support, identify urgent needs, or share their personal feelings. Such information is useful for humanitarian organizations to plan and launch relief operations. As the volume and velocity of the content are significantly high, it is crucial to have real-time systems to automatically process social media content to facilitate rapid response.

There has been a surge of research works in this domain in the past couple of years. The focus has been to analyze the usefulness of social media data and develop computational models using different modalities to extract actionable information. Among different modalities (e.g., text and image), more focus has been given to textual content analysis compared to imagery content (see [1], [2] for a comprehensive survey). Though many past research works have demonstrated that images shared on social media during a disaster event can help humanitarian organizations in a number of ways. For example, [3] uses images shared on Twitter to assess the severity of the infrastructure damage and [4] focuses on identifying damages in infrastructure as well as environmental elements.

Current publicly available datasets for developing classification models for disaster response tasks include damage severity assessment dataset [3], CrisisMMD [5], and multi-modal damage identification dataset [4]. The annotated labels in these datasets include different classification tasks such as *(i)* disaster types, *(ii)* informativeness, *(iii)* humanitarian, and *(iv)* damage severity assessment. Upon studying these datasets, we note several limitations: *(i)* when considered independently, these datasets are fairly small in contrast to the datasets used in the computer vision community, e.g., ImageNet [6] and MS COCO [7], which entangles development of robust models for real-world applications, *(ii)* they contain exact- and near-duplicates, which often provides misleading performance scores due to the random train and test splits, *(iii)* inconsistent train/test splits have been used across different studies, which makes it difficult to compare the reported results in the literature. Another interesting aspect is that there has been significant progress in neural network architectures for image processing in the past few years; however, they have not been widely explored in the *crisis informatics*[1] for disaster response tasks.

To address such limitations, our contributions in this study are as follows:

- We developed *disaster types* and *informativeness* datasets, which are completely new for the research community.
- We relabeled existing datasets for the new tasks. The motivation of using existing datasets for new tasks is that it significantly reduces the data collection, cleaning, and annotation efforts. Upon consolidation, we report the largest datasets available to date for different tasks.
- We divided each dataset into train, dev and test splits and created a *non-overlapping test set* by eliminating exact- and near-duplicate images between the test and train sets.

---

[1]https://en.wikipedia.org/wiki/Disaster_informatics

We also unified all task-specific datasets from different sources into a single set for different tasks.

- We provide benchmark results for four tasks, on separate as well as combined datasets, using several state-of-the-art neural network architectures. These results set new baselines for the crisis informatics community for the image classification tasks. Finally, we will make the datasets with their splits publicly available.[2]

The rest of the paper is organized as follows. Section II provides a brief overview of the existing work. Section III introduces the tasks while Section IV describes the datasets prepared for this study. Section V explains the experiments and Section VI presents the results and discussion. Finally, we conclude the paper in Section VII.

## II. RELATED WORK

The studies on image processing in the crisis informatics domain are relatively fewer compared to the studies on analyzing textual content for humanitarian aid.[3] With recent successes of deep learning for image classification, research works have started to use social media images for humanitarian aid. The importance of imagery content on social media for disaster response tasks has been reported in many studies [3], [8]–[13]. For instance, the analysis of flood images has been studied in [8], in which the authors reported that the existence of images with relevant textual content is more informative. Similarly, the study by Daly and Thom [9] analyzed fire event's images, which are extracted from social media data. Their findings suggest that images with geotagged information are useful to locate the fire affected areas.

The analysis of imagery content shared on social media has been recently explored using deep learning techniques for damage assessment purposes. Most of these studies categorize the severity of damage into discrete levels [3], [11], [12] whereas others quantify the damage severity as a continuous-valued index [14], [15]. Recently, [13] presented an image processing pipeline to extract meaningful information from social media images during a crisis situation, which has been developed using deep learning-based techniques. Their image processing pipeline includes collecting images, removing duplicates, filtering irrelevant images, and finally classifying them with damage severity. The study by Mouzannar et al. [4] proposed a multimodal dataset, which has been developed for training a damage detection model. Similarly, [16] explores unimodal as well as different multimodal modeling approaches based on a collection of multimodal social media posts.

Currently, publicly available datasets include damage severity assessment dataset [3], CrisisMMD [5] and damage identification multimodal dataset [4]. The former dataset is only annotated for images, whereas the latter two are annotated for both text and images. Other relevant datasets are Disaster Image Retrieval from Social Media (DIRSM) [17] and MediaEval 2018 [18]. For the image classification task, transfer learning has been a popular approach, where a pre-trained neural network is used to train a new model for a new task [19]–[22]. For this study, we follow same approach using different deep learning architectures.

Our study differs from prior works in a number of ways. We propose new datasets for different tasks, annotate existing datasets for new tasks, create non-overlapping train/dev/test splits, and finally consolidate them to create a unified, large-scale dataset for several tasks. Lastly, we use the dataset to provide benchmarks using state-of-the-art deep learning models.

## III. TASKS

For this study, we addressed four different disaster-related tasks important for humanitarian aid.

### A. Disaster type detection

When ingesting images from unfiltered social media streams, it is important to automatically detect different disaster types those images show. For instance, an image can depict a wildfire, flood, earthquake, hurricane, and other types of disasters. In the literature, disaster types have been defined in different hierarchical categories such as natural, man-made, and hybrid [23]. Natural disasters are events that result from natural phenomena (e.g., fire, flood, earthquake). Man-made disasters are events that result from human actions (e.g., terrorist attack, accidents, war, and conflicts). Hybrid disasters are events that result from human actions, which effect natural phenomena (e.g., deforestation results in soil erosion, and climate change). In this study, we focused on most frequently occurring (see in [23]) disaster event types such as *(i)* earthquake, *(ii)* fire, *(iii)* flood, *(iv)* hurricane, and *(v)* landslide. In addition, we also included two additional class labels such as *(vi)* other disaster – to cover all other disaster types (e.g., plane crash), and *(vii)* not disaster – for images that do not show any identifiable disasters. This results in a total of seven categories for the disaster type classification task. In Figure 1, we provide example images for different disaster types.

### B. Informativeness

Images posted on social media during disasters do not always contain informative (e.g., image showing damaged infrastructure due to flood, fire or any other disaster events) or useful content for humanitarian aid. It is necessary to remove any irrelevant or redundant content to facilitate crisis responders' efforts more effectively. Therefore, the purpose of this classification task is to filter irrelevant images. The class labels for this task are *(i)* informative and *(ii)* not informative.

### C. Humanitarian

An important aspect of crisis responders is to assist people based on their needs, which requires information to be classified into more fine-grained categories to take specific actions. In the literature, humanitarian categories often include *affected individuals*; *injured or dead people*; *infrastructure and*
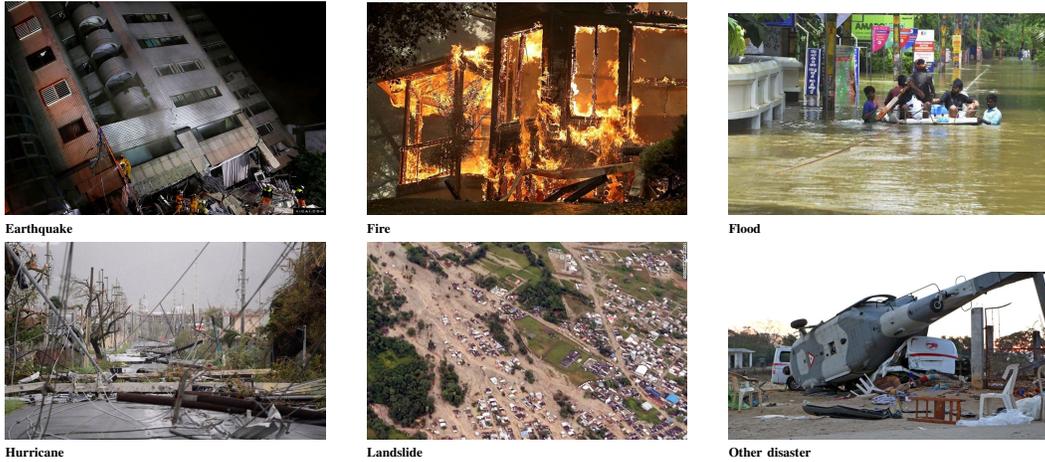
Fig. 1: Example images for different disaster types. *Not disaster* images are not shown.



Fig. 2: An image annotated as *(i)* fire event, *(ii)* informative, *(iii)* infrastructure and utility damage, and *(iv)* severe damage.

*utility damage*; *missing or found people*; *rescue, volunteering, or donation effort*; and *vehicle damage* [5]. In this study, we focus on four categories that are deemed to be the most prominent and important for crisis responders such as *(i)* affected, injured, or dead people, *(ii)* infrastructure and utility damage, *(iii)* rescue volunteering or donation effort, and *(iv)* not humanitarian.

### D. Damage severity

Assessing the severity of the damage is important to help the affected community during disaster events. The severity of damage can be assessed based on the physical destruction to a built-structure visible in an image (e.g., destruction of bridges, roads, buildings, burned houses, and forests). Following the work reported in [3], we define the categories for this classification task as *(i)* severe damage, *(ii)* mild damage, and *(iii)* little or none.

Figure 2 shows an example image that illustrates available annotations for all four tasks.

## IV. DATA PREPARATION

### A. Datasets

For this study, we used public and in-house labeled datasets. Below, we provide the details of each dataset.

*1) Damage Assessment Dataset (DAD):* The damage assessment dataset consists of labeled imagery data with damage severity levels such as severe, mild, and little-to-no damage [3]. The images have been collected from two sources: AIDR [24] and Google. To crawl data from Google, authors used the following keywords: damage building, damage bridge, and damage road. The images from AIDR were collected from Twitter during different disaster events such as Typhoon Ruby, Nepal Earthquake, Ecuador Earthquake, and Hurricane Matthew. The dataset contains $\sim 25K$ images annotated by paid-workers as well as volunteers. In this study, we use this dataset for the informativeness and damage severity tasks. For the informativeness task, we map *mild* and *severe* images into informative class and manually sift through the *little-to-no damage* images to separate them into *informative* and *not informative* categories. For the damage severity task, we map the label *little-to-no damage* into *little or none* to align with other datasets.

*2) CrisisMMD:* This is a multimodal (i.e., text and image) and multi-task dataset, which consists of $18,082$ images collected from tweets during seven disaster events crawled by the AIDR system [5]. The data is annotated by crowd-workers using the Figure-Eight platform[4] for three different tasks: *(i)* informativeness with binary labels (i.e., informative vs. not informative), *(ii)* humanitarian with seven class labels (i.e., infrastructure and utility damage, vehicle damage, rescue, volunteering, or donation effort, injured or dead people, affected individuals, missing or found people, other relevant information and not relevant), *(iii)* damage severity assessment with three labels (i.e., severe, mild and little or no damage).

*3) AIDR Disaster Type Dataset (AIDR-DT):* For disaster type classification task, we annotated images with categories mentioned in Section III-A. We obtained tweets from 17 disaster events and 3 general collections, all of which have been collected by the AIDR system. The 17 disaster events include flood, earthquake, fire, hurricane, terrorist-attack, and armed-conflict. The tweets in general collections contains

[4]Currently acquired by https://appen.com/

keywords related to natural disasters, human-induced disasters, and security incidents. We crawled images for these collections. After collecting the images we first remove exact duplicates based on *tweet ids*. Then, we remove exact- and near-duplicates images using a duplicate filtering approach discussed in [13]. From a large number of images of these collections, we sampled ∼30K images for annotation.

The labeling of these images was performed in two steps. First, a set of images were labeled as *earthquake*, *fire*, *flood*, *hurricane*, and *none of these categories*. Then, we selected a sample of ∼2200 images, which are labeled as *none of these categories* in the previous step for annotating *not disaster* and *other disaster* categories. The rationale for choosing such a sample number of images was due to limited annotation resources.

For the landslide category, we crawled images from Google, Bing, and Flickr using keywords landslide, mudslide, "mud slides", landslip, "rock slides", rockfall, "land slide", earthslip, rockslide, and "land collapse". As images have been collected from different sources, therefore, it resulted in having duplicates. To take this into account, we applied the same duplicate filtering as before to remove exact- and near-duplicate images. Then, the remaining images were manually labeled as *landslide* and *not landslide*.

For the annotation task, we used the following definitions for the disaster types:

(i) Earthquake: images showing damaged or destroyed buildings, fractured houses, ground ruptures such as railway lines, roads, airport runways, highways, bridges, and tunnels.
(ii) Fire: images showing man-made fires or wildfires (forests, grasslands, brush, and deserts), destroyed forests, houses, or infrastructures.
(iii) Flood: images showing flooded areas, houses, roads, and other infrastructures.
(iv) Hurricane: images showing high winds, a storm surge, heavy rains, collapsed electricity polls, grids, and trees.
(v) Landslide: images showing landslide, mudslide, landslip, rockfall, rockslide, earth slip, and land collapse
(vi) Other disasters: images showing any other disaster types such as plane crash, bus, car, or train accident, explosion, war, and conflicts.
(vii) Not disaster: images showing cartoon, advertisement, or anything that cannot be easily linked to any disaster type.

In Figure 3, we report the distribution of the labeled images in different events and general collections.

*4) AIDR Informativeness Dataset (AIDR-Info):* For this dataset, we collected tweets and images using the AIDR system. We used the same duplicate filtering approach to remove duplicate images. Then, we labeled 9,936 images with two class labels, informative vs. not-informative using the definition discussed in [5].[5] In Figure 4, we report the distribution of images labeled for different events. Across

[5]If the image is useful for humanitarian aid then we label it as "informative" otherwise as "not informative".
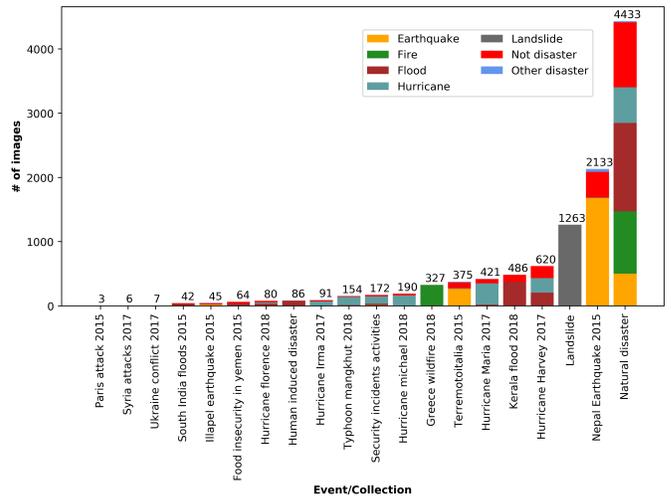


Fig. 3: Number of disaster types labeled images from different disaster events and collections in AIDR-DT dataset.
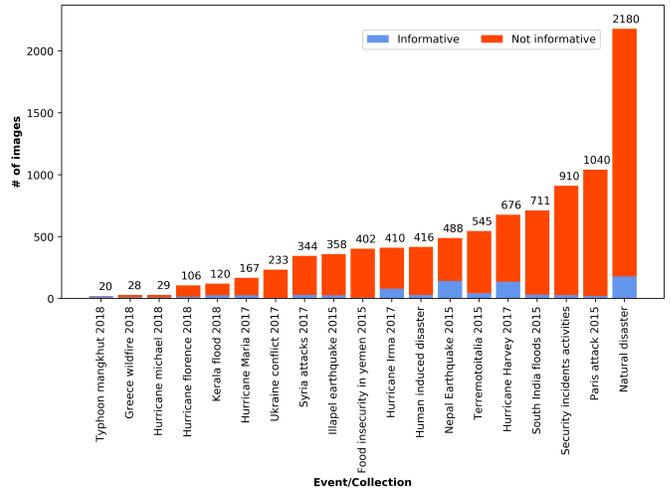


Fig. 4: Number of labeled informative *vs.* not-informative images from different disaster events and collections in AIDR-info dataset.

different collections, number of not informative images is higher than informative images.

*5) Damage Multimodal Dataset (DMD):* The multimodal damage identification dataset consists of 5,878 images collected from Instagram and Google [4]. Authors of the study crawled the images using more than 100 hashtags, which are proposed in crisis lexicon [25]. The manually labeled data consist of six damage class labels such as fires, floods, natural landscape, infrastructural, human, and non-damage. The non-damage image includes cartoons, advertisements, and images that are not relevant or useful for humanitarian tasks. For this study, we re-labeled them for all four tasks. using the same class labels discussed in the previous section. We followed the annotation instructions reported in [5] and as discussed in section III-A.

## B. Annotation

The annotation has been done by domain experts and ensured the quality of the annotation for new datasets and relabeling existing datasets for new tasks. For the disaster type labeling, the annotators followed the definition discussed in IV-A3 and for other tasks, definitions and instructions are adapted from [5].

## C. Data Split

Before consolidating the datasets we split each of them into train, dev, and test sets with 70:10:20 ratio, respectively. The purpose was threefold: *(i)* train and evaluate individual datasets on each task, *(ii)* have a close-to-equal distribution from each dataset into the final consolidated dataset, and *(iii)* provide the research community an opportunity to use the splits independently. After data split, we identify duplicate images (see in Section IV-D) across sets and move them into the training set to create a non-overlapping test set.

## D. Duplicate Image Identification

To develop a machine learning model, it is important to design non-overlapping training and test sets. A common practice is to randomly split the dataset into train and test sets. This approach often creates an overlapping train-test split with social media data. For example, exact- or near-duplicate images can be in both train and test sets. Based on the work in [13], we identified duplicate images. Since all datasets have already been manually labeled, we did not want to remove any image from any dataset. We instead attempted to create a non-overlapping train, dev, and test split. The motivation is that having exact- and near-duplicate images in the training set creates a natural augmentation in the training set.

To identify duplicate images in the test set, we first train the model using train and dev set and find the nearest images of the test set. To train the model, we first extract features using a pre-trained deep learning model.[6] Then, we use the Nearest Neighbor [27] to train the model with the training set of each respective dataset. For example, for the informative dataset of CrisisMMD, we use the training set to train the model and then use it to obtain the nearest images for each image in the test set.

Next, we manually identify duplicates by investigating each image from a given test set and the identified nearest images from the corresponding train set for four different tasks and twelve different datasets. Out of these images, we identified 5,593 exact- and near-duplicate images in different test sets. We then move the identified images to the training set to create non-overlapping test sets. It also helped us to identify an approximate threshold to automatically identify near-duplicate images. In Figure 5, we present a histogram of *Euclidean* distance measures of the exact- and near-duplicate images. It shows the number of images in different bins. With our analysis, we realized that a distance threshold of less than

---

[6]Note that the pre-trained model is trained using ResNet18 architecture [26] on the damage assessment dataset [3].
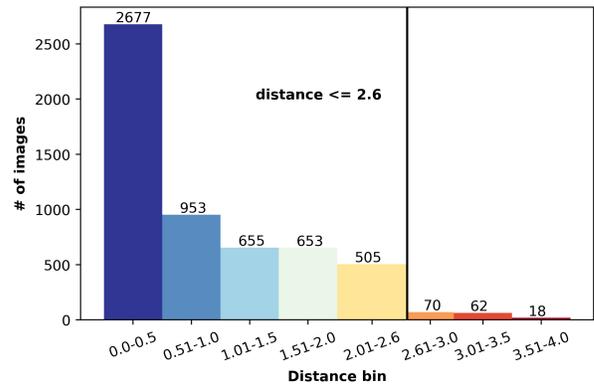


Fig. 5: Number of duplicates in different distance bins.

equal 2.6 is a reasonable choice for automatic duplicate detection. From the figure, we observe that there are also duplicate images with higher thresholds; however, that is a very small number comparatively. Also, note that choosing a higher number will lead to an increase in false positives. Using the threshold value of $d \leq 2.6$ we automatically identified duplicate images in the dev set and moved them to the train set.

It is important to note that creating non-overlapping datasets using duplicate identification process reduces the distribution of dev and test sets. This is reasonable given the fact that we are ensuring unbiased train/test splits.

## E. Data Consolidation:

One of the important reasons to perform data consolidation is to develop robust deep learning models with large amounts of data. For this purpose, we merge all train, dev, and test sets into the consolidated train, dev, and test sets, respectively. As combining multiple datasets can results in duplicate images in train and test set, after merging the dataset, we repeat the same duplicate identification procedure to create non-overlapping sets for different tasks.

## F. Data Statistics

Tables I, II, III, IV, and V show the label distribution of all datasets for different tasks. We report the total number of images in parenthesis for each dataset in the Tables. Some class labels are skewed in individual datasets. For example, in disaster type datasets (Table I), the distribution of "other disaster" label is low in AIDR-DT dataset, whereas the distribution of "landslide" label low in DMD dataset. For the informativeness task, low distribution is observed for the "informative" label. Moreover, for the humanitarian task, we have low distribution for "rescue volunteering or donation effort" label in DMD dataset, and for the damage severity task "mild" label in CrisisMMD and DMD datasets. However, the consolidated dataset creates a fair balance across class labels for different tasks as shown in Table V.

TABLE I: Data split for the **disaster types** task. Number in parenthesis shows total number of images.

| Dataset | Class labels | Train | Dev | Test | Total |
|---|---|---|---|---|---|
| **AIDR-DT** (11723) | Earthquake | 1910 | 201 | 376 | 2487 |
| | Fire | 990 | 105 | 214 | 1309 |
| | Flood | 2059 | 241 | 533 | 2833 |
| | Hurricane | 1188 | 142 | 279 | 1609 |
| | Landslide | 901 | 119 | 257 | 1277 |
| | Not disaster | 1507 | 198 | 415 | 2120 |
| | Other disaster | 65 | 6 | 17 | 88 |
| **DMD** (5788) | Earthquake | 130 | 17 | 35 | 182 |
| | Fire | 255 | 36 | 71 | 362 |
| | Flood | 263 | 35 | 70 | 368 |
| | Hurricane | 253 | 36 | 73 | 362 |
| | Landslide | 38 | 5 | 11 | 54 |
| | Not disaster | 2108 | 288 | 575 | 2971 |
| | Other disaster | 1057 | 145 | 287 | 1489 |

TABLE II: Data split for the **informativeness** task.

| Dataset | Class labels | Train | Dev | Test | Total |
|---|---|---|---|---|---|
| **DAD** (25820) | Informative | 15329 | 590 | 2266 | 18185 |
| | Not informative | 5950 | 426 | 1259 | 7635 |
| **CrisisMMD** (18082) | Informative | 7233 | 635 | 1507 | 9375 |
| | Not informative | 6535 | 551 | 1621 | 8707 |
| **DMD** (5878) | Informative | 2071 | 262 | 573 | 2906 |
| | Not informative | 2152 | 240 | 580 | 2972 |
| **AIDR-Info** (9936) | Informative | 627 | 66 | 172 | 865 |
| | Not informative | 6677 | 598 | 1796 | 9071 |

## V. EXPERIMENTS

### A. Experimental Settings

We employ the transfer learning approach to perform experiments, which has shown promising results for various visual recognition tasks in the literature [19]–[22]. The idea of the transfer learning approach is to use existing weights of a pre-trained model. For this study, we used several neural network architectures using the PyTorch library.[7] The architectures include ResNet18, ResNet50, ResNet101 [26], AlexNet [28], VGG16 [29], DenseNet [30], SqueezeNet [31], InceptionNet [32], MobileNet [33], and EfficientNet [34].

We use the weights of the networks trained using ImageNet [6] to initialize our model. We adapt the last layer (i.e., softmax layer) of the network according to the particular classification task at hand instead of the original 1,000-way classification. The transfer learning approach allows us to transfer the features and the parameters of the network from the broad domain (i.e., large-scale image classification) to the specific one, in our case four different classification tasks. We train the models using the Adam optimizer [35] with an initial learning rate of $10^{-5}$, which is decreased by a factor of 10 when accuracy on the dev set stops improving for 10 epochs.

We designed the binary classifier for informativeness task and multiclass classifiers for other tasks.

To measure the performance of each classifier, we use weighted average precision (P), recall (R), and F1-measure (F1). We only report F1-measure due to limited space.

### B. Datasets Comparison

To determine whether consolidated data helps achieve better performance, we train the models using training sets from

TABLE III: Data split for the **humanitarian** task.

| Dataset | Class labels | Train | Dev | Test | Total |
|---|---|---|---|---|---|
| **CrisisMMD** (11241) | Affected, injured, or dead people | 521 | 51 | 100 | 672 |
| | Infrastructure and utility damage | 3040 | 299 | 589 | 3928 |
| | Not humanitarian | 3307 | 296 | 807 | 4410 |
| | Rescue volunteering or donation effort | 1682 | 174 | 375 | 2231 |
| **DMD** (5528) | Affected, injured, or dead people | 242 | 28 | 63 | 333 |
| | Infrastructure and utility damage | 933 | 125 | 242 | 1300 |
| | Not humanitarian | 2736 | 314 | 744 | 3794 |
| | Rescue volunteering or donation effort | 74 | 9 | 18 | 101 |

TABLE IV: Data split for the **damage severity** task.

| Dataset | Class labels | Train | Dev | Test | Total |
|---|---|---|---|---|---|
| **DAD** (25820) | Little or none | 7881 | 1101 | 1566 | 10548 |
| | Mild | 2828 | 388 | 546 | 3762 |
| | Severe | 9457 | 673 | 1380 | 11510 |
| **CrisisMMD** (3198) | Little or none | 317 | 35 | 67 | 419 |
| | Mild | 547 | 56 | 125 | 728 |
| | Severe | 1629 | 144 | 278 | 2051 |
| **DMD** (5878) | Little or none | 2874 | 331 | 778 | 3983 |
| | Mild | 508 | 60 | 132 | 700 |
| | Severe | 857 | 110 | 228 | 1195 |

the individual and consolidated datasets. However, we always test the models on the consolidated test set. As our test data is same across different experiments, results are ensured to be comparable. Since we have four different tasks, which consist of fifteen different datasets, we only experimented with the ResNet18 [26] network architecture to manage the computational load.

### C. Network Architectures

Currently available neural network architectures come with different computational complexity. As one of our goals is to deploy the models in real-time applications, we exploit them to understand their performance differences. Another motivation is that current literature in crisis informatics only reports results using one or two network architectures (e.g., VGG16 in [16], InceptionNet in [4]), which we wanted to extend in this study.

## VI. RESULTS AND DISCUSSIONS

### A. Results

*1) Dataset Comparison:* In Table VI, we report classification results for different tasks and different datasets using ResNet18 network architecture. The performance of different tasks is not equally comparable as they have different levels of complexity (e.g., varying number of class labels, class imbalance, etc.). For example, the informativeness classification is a binary task, which is computationally simpler than a classification task with more labels (e.g., seven labels in disaster types). Hence, the performance is comparatively higher for informativeness. An example of a class imbalance issue can be seen in Table V with the damage severity task. The distribution of mild is comparatively small, which reflects on its and overall performance. The mild class label is also less distinctive than other class labels, and we noticed that classifiers often confuse this class label with the other two class labels. Similar findings have also been reported in [3]. For the disaster types task, the performance of the AIDR-DT

TABLE V: Data splits for the **consolidated dataset** for all tasks.

| Class labels | Train | Dev | Test | Total |
|---|---|---|---|---|
| **Disaster types (17511)** | | | | |
| Earthquake | 2058 | 207 | 404 | 2669 |
| Fire | 1270 | 121 | 280 | 1671 |
| Flood | 2336 | 266 | 599 | 3201 |
| Hurricane | 1444 | 175 | 352 | 1971 |
| Landslide | 940 | 123 | 268 | 1331 |
| Not disaster | 3666 | 435 | 990 | 5091 |
| Other disaster | 1132 | 143 | 302 | 1577 |
| **Informativeness (59717)** | | | | |
| Informative | 26486 | 1432 | 3414 | 31332 |
| Not informative | 21700 | 1622 | 5063 | 28385 |
| **Humanitarian (16769)** | | | | |
| Affected, injured, or dead people | 772 | 73 | 160 | 1005 |
| Infrastructure and utility damage | 4001 | 406 | 821 | 5228 |
| Not humanitarian | 6076 | 578 | 1550 | 8204 |
| Rescue volunteering or donation effort | 1769 | 172 | 391 | 2332 |
| **Damage severity (34896)** | | | | |
| Little or none | 11437 | 1378 | 2135 | 14950 |
| Mild | 4072 | 489 | 629 | 5190 |
| Severe | 12810 | 845 | 1101 | 14756 |

TABLE VI: Results of different classification tasks using the ResNet18 models. Trained on individual and consolidated datasets and tested on consolidated test sets. Disaster types (Disas.), Informativeness (Info.), Humanitarian (Hum.), Damage severity (Damage).

| Dataset | Disas. | Info. | Hum. | Damage |
|---|---|---|---|---|
| AIDR-DT | 0.726 | - | - | - |
| DAD | - | 0.797 | - | 0.709 |
| CrisisMMD | - | 0.790 | 0.727 | 0.374 |
| DMD | 0.591 | 0.799 | 0.636 | 0.663 |
| AIDR-Info | - | 0.725 | - | - |
| Consolidated | **0.785** | **0.851** | **0.749** | **0.736** |

model is higher compared to the DMD model. We observe that the DMD dataset is comparatively small and the model is not performing well on the consolidated dataset. This characteristic is observed in other tasks as well. As expected, overall for all tasks, the models with the consolidated datasets outperform individual datasets.

*2) Network Architectures Comparison:* In Table VII, we report results using different network architectures on consolidated datasets for different tasks, i.e., trained and tested using a consolidated dataset. Across different tasks, overall EfficientNet is performing better than other models except for humanitarian task, for which VGG16 is outperforming other models. Comparatively the second-best models are VGG16, ResNet50, ResNet101, and DenseNet (101). From the results of different tasks, we observe that InceptionNet is the worst performing model.

In Table VII, we also report different neural network models with their number of layers, parameters, and memory consumption during inference. There can always be a trade-off between performance vs. computational complexity, i.e., number of layers, parameters, and memory consumption. In terms of memory consumption and the number of parameters, VGG16 seems expensive than others. Based on the performance and computational complexity, we can conclude that EfficientNet can be the best option to use in real-time applications. We

TABLE VII: Results using different neural network models on the consolidated dataset with four different tasks. Trained and tested using the consolidated dataset. Comparable result is shown in **bold** and best results is shown in <u>underlined</u>. #L (P) - number of layers (number of parameters in millions). Mem. (memory in MB). IncepNet (InceptionNet), MobNet (MobileNet), EffiNet (EfficientNet)

| | # L (P) | Mem. | Disas. | Info. | Hum. | Damage | Avg. |
|---|---|---|---|---|---|---|---|
| ResNet18 | 18 (11.2) | 74.6 | 0.785 | 0.851 | 0.749 | 0.736 | 0.780 |
| ResNet50 | 50 (23.5) | 233.5 | **0.808** | 0.852 | 0.762 | **0.751** | **0.793** |
| ResNet101 | 101 (42.5) | 377.6 | **0.813** | 0.852 | 0.765 | 0.737 | 0.792 |
| AlexNet | 8 (57.0) | 222.2 | 0.754 | 0.828 | 0.716 | 0.709 | 0.752 |
| VGG16 | 16 (134.3) | 673.9 | 0.798 | **0.858** | **0.773** | **0.753** | **0.796** |
| DenseNet | 121 (7.0) | 174.2 | **0.806** | **0.862** | 0.755 | 0.739 | 0.791 |
| SqueezeNet | 18 (0.7) | 48.0 | 0.755 | 0.829 | 0.719 | 0.708 | 0.753 |
| IncepNet | 42 (24.3) | 206.0 | 0.528 | 0.593 | 0.509 | 0.615 | 0.561 |
| MobNet (v2) | 20 (2.2) | 8.5 | 0.782 | 0.849 | 0.746 | 0.730 | 0.777 |
| EffiNet (b1) | 25 (7.8) | 177.8 | <u>0.816</u> | <u>0.863</u> | <u>0.765</u> | <u>0.758</u> | <u>0.801</u> |

computed throughput for EfficientNet using a batch size of 128 and it can process ∼260 images per second on an NVIDIA Tesla P100 GPU. Among different ResNet models, ResNet18 is a reasonable choice given that its computational complexity is significantly less than other ResNet models.

### B. Discussions

Achieving a better performance with deep learning models requires relatively larger datasets. To date, the developed dataset sizes for disaster response tasks are comparatively small. Hence, we address that by combining data from multiple sources and relabeled them for new tasks. The proposed datasets consists of binary and multiple class labels and addressed in binary and multiclass classification settings. However, the datasets can be turned into multi-label and multi-task settings, which we aim to address in a future study.

A significant challenge with social media data is the exact- and near-duplicate content. We address this issue, and our proposal for the community is to remove duplicates before the annotation process. Towards this direction, another important challenge is that current duplicate detection is similarity and threshold-based with deep learning feature extraction. In our analysis, we describe a procedure to determine a reasonable threshold for automatic duplicate detection. However, this requires further study, which we aim to do in the future.

Real-time event detection is an important problem from social media content. Our new disaster types dataset can help to develop models and deploy in real-time applications. We also explore several deep learning models, which vary with performance and complexities. Among them, EfficientNet appears to be a reasonable option. Note that EfficientNet has a series of network architectures (b0-b7) and for this study, we only reported results with EfficientNet (b1). We aim to further explore other architectures.

A small and low latency model is desired to deploy mobile and handheld embedded computer vision applications. The development of MobileNet [33] sheds light towards that direction. Our experimental results suggest that it is computationally simpler and provides a reasonable accuracy, only 2-3% lower than the best models for different tasks.

Comparing our results with previous state-of-the-art is not possible due to differences in data splits and the issue of duplicate images. On informativeness and humanitarian tasks, previous reported results (weighted F1) are 83.2 and 76.3, respectively, using the CrisisMMD dataset [16]. The authors in [4] reported a test accuracy of $83.98 \pm 1.72$ for six disaster types tasks using the DMD dataset with a five-fold cross-validation run. In another study, using the CrisisMMD dataset, authors report weighted-F1 of 81.22 and 86.96 for informativeness and humanitarian tasks, respectively [36]. They used a small subset of the whole CrisisMMD dataset in their study. Due to differences in data splits, these systems are hard to compare. However, we hope our datasets and splits will provide a standard ground for future studies to compare results.

## VII. CONCLUSIONS

Images shared on social media contain useful information for humanitarian organizations. There has been limited work for disaster response image classification tasks compared to text due to the limited resources to develop deep learning models. In this study, we provide new datasets for disaster type detection and informativeness classification. We also relabeled existing datasets for new tasks, and provide a consolidated dataset. We identified duplicates and created non-overlapping splits, which can ensure unbiased results. We addressed four tasks such as disaster types, informativeness, humanitarian and damage severity, that are needed for disaster response. The datasets have a unique characteristic that it can turn into multi-label, and multitask learning setups and would be useful for the deep learning community to develop new algorithms. We also aim to address this in the future. Furthermore, we used different state-of-the-art deep learning architectures to provide benchmark results on the datasets.

## REFERENCES

[1] M. Imran, C. Castillo, F. Diaz, and S. Vieweg, "Processing social media messages in mass emergency: A survey," *ACM Computing Surveys*, vol. 47, no. 4, p. 67, 2015.

[2] N. Said, K. Ahmad, M. Riegler, K. Pogorelov, L. Hassan, N. Ahmad, and N. Conci, "Natural disasters detection in social media and satellite imagery: a survey," *Multimedia Tools and Applications*, vol. 78, no. 22, pp. 31 267–31 302, 2019.

[3] D. T. Nguyen, F. Ofli, M. Imran, and P. Mitra, "Damage assessment from social media imagery data during disasters," in *Proc. of ASONAM*, Aug 2017, pp. 1–8.

[4] H. Mouzannar, Y. Rizk, and M. Awad, "Damage Identification in Social Media Posts using Multimodal Deep Learning," in *Proc. of ISCRAM*, May 2018, pp. 529–543.

[5] F. Alam, F. Ofli, and M. Imran, "CrisisMMD: multimodal twitter datasets from natural disasters," in *Proc. of ICWSM*, Jun 2018, pp. 465–473.

[6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255.

[7] T. Lin, M. Maire, S. J. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," in *European Conference on Computer Vision (ECCV)*, 2014, pp. 740–755.

[8] R. Peters and J. Porto de Albuqerque, "Investigating images as indicators for relevant social media messages in disaster management," in *Proc. of ISCRAM*, 2015.

[9] S. Daly and J. Thom, "Mining and classifying image posts on social media to analyse fires," in *Proc. of ISCRAM*, 2016, pp. 1–14.

[10] T. Chen, D. Lu, M.-Y. Kan, and P. Cui, "Understanding and classifying image tweets," in *ACM Multimedia*, 2013, pp. 781–784.

[11] D. T. Nguyen, F. Alam, F. Ofli, and M. Imran, "Automatic image filtering on social networks using deep learning and perceptual hashing during crises," in *Proc. of ISCRAM*, May 2017.

[12] F. Alam, M. Imran, and F. Ofli, "Image4act: Online social media image processing for disaster response." in *Proc. of ASONAM*, 2017, pp. 1–4.

[13] F. Alam, F. Ofli, and M. Imran, "Processing social media images by combining human and machine computing during crises," *International Journal of Human–Computer Interaction*, vol. 34, no. 4, pp. 311–327, 2018.

[14] K. R. Nia and G. Mori, "Building damage assessment using deep learning and ground-level image data," in *14th Conference on Computer and Robot Vision (CRV)*. IEEE, 2017, pp. 95–102.

[15] X. Li, D. Caragea, H. Zhang, and M. Imran, "Localizing and quantifying damage in social media images," in *Proc. of ASONAM*, 2018, pp. 194–201.

[16] F. Ofli, F. Alam, and M. Imran, "Analysis of social media data using multimodal deep learning for disaster response," in *Proc. of ISCRAM*, May 2020.

[17] B. Bischke, P. Helber, C. Schulze, V. Srinivasan, A. Dengel, and D. Borth, "The multimedia satellite task at mediaeval 2017." in *MediaEval*, 2017.

[18] B. Benjamin, H. Patrick, Z. Zhengyu, B. J. de, and B. Damian, "The multimedia satellite task at mediaeval 2018: Emergency response for flooding events," in *MediaEval*, Oct 2018.

[19] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson, "How transferable are features in deep neural networks?" in *Advances in Neural Information Processing Systems*, 2014, pp. 3320–3328.

[20] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson, "CNN features off-the-shelf: an astounding baseline for recognition," in *Proc. of CVPR Workshops*, 2014, pp. 806–813.

[21] G. Ozbulak, Y. Aytar, and H. K. Ekenel, "How transferable are cnn-based features for age and gender classification?" in *International Conference of the Biometrics Special Interest Group*, Sept 2016, pp. 1–6.

[22] M. Oquab, L. Bottou, I. Laptev, and J. Sivic, "Learning and transferring mid-level image representations using convolutional neural networks," in *Proc. of CVPR*, 2014, pp. 1717–1724.

[23] I. M. Shaluf, "Disaster types," *Disaster Prevention and Management: An International Journal*, 2007.

[24] M. Imran, C. Castillo, J. Lucas, P. Meier, and S. Vieweg, "AIDR: Artificial intelligence for disaster response," in *Proc. of WWW*, 2014, pp. 159–162.

[25] A. Olteanu, C. Castillo, F. Diaz, and S. Vieweg, "Crisislex: A lexicon for collecting and filtering microblogged communications in crises." in *Proc. of ICWSM*, 2014.

[26] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. of CVPR*, 2016, pp. 770–778.

[27] P. Cunningham and S. J. Delany, "k-nearest neighbour classifiers," *Multiple Classifier Systems*, vol. 34, no. 8, pp. 1–17, 2007.

[28] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.

[29] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[30] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. of CVPR*, 2017, pp. 4700–4708.

[31] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5 mb model size," *arXiv:1602.07360*, 2016.

[32] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. of CVPR*, 2016, pp. 2818–2826.

[33] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv:1704.04861*, 2017.

[34] M. Tan and Q. V. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," *arXiv:1905.11946*, 2019.

[35] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Proc. of ICLR*, 2015.

[36] M. Abavisani, L. Wu, S. Hu, J. Tetreault, and A. Jaimes, "Multimodal categorization of crisis events in social media," in *Proc. of CVPR*, 2020, pp. 14 679–14 689.