

# MEDIC: A Multi-Task Learning Dataset for Disaster Image Classification

Firoj Alam,<sup>1</sup> Tanvirul Alam,<sup>2</sup> Md. Arid Hasan,<sup>3,4</sup> Abul Hasnat,<sup>5</sup>  
Muhammad Imran,<sup>1</sup> Ferda Ofli<sup>1</sup>

<sup>1</sup>Qatar Computing Research Institute, HBKU, Qatar

<sup>2</sup>Rochester Institute of Technology, USA, <sup>3</sup>Cognitive Insight Limited, Bangladesh

<sup>4</sup>Daffodil International University, Bangladesh <sup>5</sup>BLACKBIRD.AI, USA

{fialam, mimran, fofli}@hbkku.edu.qa,

tanvirul.alam@mail.rit.edu, arid.cse0325.c@diu.edu.bd, mhasnat@gmail.com

## Abstract

Recent research in disaster informatics demonstrates a practical and important use case of artificial intelligence to save human lives and sufferings during natural disasters based on social media contents (text and images). While notable progress has been made using texts, research on exploiting the images remains relatively under-explored. To advance the image-based approach, we propose MEDIC<sup>1</sup>, which is the largest social media image classification dataset for humanitarian response consisting of 71,198 images to address four different tasks in a multi-task learning setup. This is the first dataset of its kind: social media image, disaster response, and multi-task learning research. An important property of this dataset is its high potential to contribute research on *multi-task learning*, which recently receives much interest from the machine learning community and has shown remarkable results in terms of memory, inference speed, performance, and generalization capability. Therefore, the proposed dataset is an important resource for advancing image-based disaster management and multi-task machine learning research.

## 1 Introduction

Natural disasters cause significant damage (e.g., Hurricane Harvey in 2017 cost \$125 billion)<sup>2</sup> and it requires urgent assistance in time of crisis. In the last decade, various social media played important roles in humanitarian response tasks as they were widely used to disseminate information and obtain valuable insights. During disaster events, people post content (e.g., text, images, and video) on social media to ask for help (e.g., report of a person stuck on a rooftop during a flood), offer support, identify urgent needs, or share their feelings. Such information is helpful for humanitarian organizations to take immediate actions, plan and launch relief operations. Recent researches demonstrated that images shared on social media during a disaster event assist humanitarian organizations, which include assessing the severity of the infrastructure damage [56], identifying damages in infrastructure [53], identifying humanitarian information [4], detecting crisis incidents [79], and detecting disaster events with other related tasks [6]. However, the amount of research and resources to develop powerful computer vision based predictive models remains insufficient compared to the NLP based progress [30, 67, 32]. This research is motivated by these observations and aims to enrich resources to make further advancements in the computer vision based disaster management studies.

<sup>1</sup>Available at: <https://crisisnlp.qcri.org/medic/index.html>

<sup>2</sup>[https://en.wikipedia.org/wiki/List\\_of\\_disasters\\_by\\_cost](https://en.wikipedia.org/wiki/List_of_disasters_by_cost)



Figure 1: Examples of images representing all tasks. **T1:** Disaster types, **T2:** Informativeness, **T3:** Humanitarian, **T4:** Damage severity.

Several models addressing different tasks need to be deployed to track real-time disaster events and extract humanitarian and damage-related information as reported in [5, 3]. These tasks include (i) disaster types, (ii) informativeness, (iii) humanitarian, and (iv) damage severity assessment (see section 3 for a more detail). Existing works [56, 4, 53] address these tasks separately, which turns out to have higher computational complexities (e.g., computational power, training and inference time). Hence, this research aims at reducing this gap by addressing the different tasks simultaneously in multi-task learning (MTL) setup, which can also help in reducing carbon footprint [68].

Recent advances in deep convolutional neural networks (CNN) and their learning techniques provide efficient solutions for different computer vision applications. While the simple computer vision applications require applying only single task such as classification [24], semantic segmentation [50], or object detection [63], the complex computer vision applications such as autonomous vehicles, robotics, social media image streaming [5, 84] need to incorporate multiple tasks, which significantly increases the computational and memory requirements for both training and inference. MTL techniques [10, 84, 76] have emerged as the standard approach for these complex computer vision applications where a model is trained to solve multiple tasks simultaneously, which helps to improve the performance, reduce inference time and computational complexities. For example, an image posted on social media during a disaster event can contain information whether it is a flood event, shows infrastructure damage, and is severe. Such a multitude of information needs to be detected in real-time to facilitate humanitarian organizations [5, 3] where a single model solving multiple tasks can be more effective than having multiple models for multiple tasks.

Labeled public image datasets, such as ImageNet [66] and Microsoft COCO [48] made significant contributions to the advancement of today’s powerful machine learning models. Likewise, for the MTL setup, several image datasets have already been proposed, which are summarized in Table 1. These datasets include images from different domains such as indoor scenes, driving, face, handwritten digits, and animal recognition, which are already contributing to the advancement of MTL research. However, an MTL dataset for critical real-world applications which comprise humanitarian response tasks during natural disasters is yet to become available. This research proposes a novel MTL dataset for disaster image classification.

This research extends the previous work of Alam et al. [6] where the images are mostly annotated for individual tasks, and only 5,558 out of 71,198 images have labels for all four tasks mentioned above. We provide its extensive extension by annotating the images for all tasks, i.e., we annotated 155,899 more labels for these tasks in addition to existing ones.<sup>3</sup> Figure 1 shows example images with the labels for all four tasks.

Our contributions in this research can be summarized as follows: (i) we provide a social media MTL image dataset for disaster response tasks with various complexities (), which can be used as an

<sup>3</sup>For four tasks, 71,198 images results in 284,792 labels, existing annotation consisted only 128,893 labels.

evaluation benchmark for computer vision research; *(ii)* we ensured high quality annotation by making sure that at least two annotators agree on a label; *(iii)* we provide a benchmark for heterogeneous multi-task learning and baseline studies to facilitate future study; *(iv)* our experimental results can also be used as a baseline in the single-task learning setting.

The rest of the paper is organized as follows. Section 2 provides a brief overview of the existing work. Section 3 introduces the tasks and describes the dataset development process. Section 4 explains the experiments, and presents the results, and Section 5 provides a discussion. Finally, we conclude the paper in Section 6.

## 2 Related Work

This paper mainly focuses on the development of multi-task learning dataset for disaster response tasks. We first discuss the recent related work on multi-task learning and available multi-task learning datasets; finally, we discuss social media image classification literature and datasets for disaster response.

### 2.1 Multi-task Learning and Datasets

**Multi-task Learning** MTL aims to improve generalization capability by leveraging information in the training data consisting of multiple related tasks [10]. It simultaneously learns multiple tasks and has shown promising results in terms of generalization, computation, memory footprint, performance, and inference time by jointly learning through a shared representation [10, 76]. Since the seminal work by Caruana [10], research on multi-task learning has received wide attention in the last several years in NLP, computer vision, and other research areas, see related surveys in [64, 86, 76, 14, 80]. Multi-task learning brings benefits when associated tasks share complementary information. However, performance can suffer when multiple tasks have conflicting needs, and the tasks have competing priorities (i.e., one is superior to the other). This phenomenon is referred to as negative transfer. This understanding led to the question of what, when, and how to share information among tasks [73, 76]. To address these aspects, in the deep learning era, numerous architectures and optimization methods have been proposed. The architectures are categorized into hard and soft parameter sharing. Hard parameter sharing design consists of a shared network followed by task-specific heads [37, 35, 12]. In soft parameter sharing, each task has its own set of parameters, and a feature sharing mechanism to deal with cross-task task [52, 65, 20]. In the multi-task learning literature, a problem can be formulated in two different ways - homogeneous and heterogeneous [73]. While the homogeneous multi-task learning assumes that each task corresponds to a single output, the heterogeneous multi-task learning assumes each task corresponds to a unique set of output labels [10, 82]. The latter setting uses a neural network using multiple sets of outputs and losses. In this study, we aim to provide a benchmark with our heterogeneous MTL dataset using the hard parameter sharing approach.

**Datasets** Earlier studies such as [34] and [40] mostly exploited the MNIST [44] and USPS [27] datasets for MTL experiments. These datasets were originally designed for single-task classification settings. For example, the widely used MNIST dataset was originally designed for 10 digits classification, and Office-Caltech [21] was designed to categorize images in 31 classes, which are collected from different domains. However, such datasets are used with the homogeneous problem setting of multi-task learning by selecting ten target classes as ten binary classification tasks [40, 73, 81]. Numerous other widely used datasets such as MC-COCO [47] and CelebA [49] have also been used for multi-task learning in the homogeneous problem setting.

Several existing datasets consisting of multiple unique output label sets were studied in the heterogeneous setting. For example, AdienceFaces [17] was designed for gender and age group classification tasks, OmniArt [74] consists of seven tasks, NYU-V2 [71] consists of three tasks, and PASCAL [18, 11] consists of 5 tasks. Very few datasets were specifically designed for multi-task learning research. Most notable ones are Taskonomy [85] and BDD100K [84]. The Taskonomy dataset consists of 4 million images of indoor scenes from 600 buildings, and each image was annotated for twenty-six visual tasks. Ground truths of this dataset were obtained programmatically, and knowledge distillation approaches. The BDD100K dataset is a diverse 100K driving video dataset consisting of

ten tasks. It was collected from Nexar,<sup>4</sup> where videos are uploaded by the drivers. In Table 1, we provide widely used datasets, which have been used for multi-task learning.

Ref.	Dataset	Source	Size	Task type	# Tasks	Tasks	# classes	Domain	Year
<b>Datasets used for multi-task learning</b>									
[76]	PASCAL [18, 11]	Flickr	12,030 (I)	Hete.	5	SS, HS, SE, and SD	-	Diverse objects	2021
[76]	NYU-V2 [71]	PC	1,449 (I)	Hete.	3	IS, SS, and SC	-	Indoor video	2021
[84]	BDD100K	Nexar	100,000 (V)	Hete.	10	ten tasks	<sup>5</sup>	Driving	2020
[73]	MNIST [44]	-	70,000 (I)	Homo.	10	10 digits cls.	10 CL.	Handwritten	2019
[73]	CIFAR10 [39]	-	60,000 (I)	Homo.	10	10 animal cls.	10 CL.	Animal	2019
[73]	UCSD-Birds [78]	-	11,788 (I)	Homo.	10	10 R/tasks	Ranking	Animal	2019
[73]	OmniGlott [42]	-	1,623 (I)	Homo.	50	50 alphabets	50 CL.	Handwritten	2019
[73]	OmniArt [74]	-	133,000 (S)	Hete.	7	7 tasks	-	Artwork	2019
[85]	Taskonomy	IC	4M (I)	Hete.	26	26 tasks	-	Indoor scenes	2018
[51]	Office-caltech [21]	-	2,533 (I)	Homo.	4	Amazon, Webcam, and DSLR, Caltech-256	10 CL/task	-	2017
[51]	Office-Home [77]	SE	15,500 (I)	Homo.	4	Artistic, clip art, product, and real-world images	65 objects	Office/Home	2017
[51]	ImageCLEF <sup>6</sup>	-	2,400 (I)	Homo.	4	Caltech-256, ImageNet	-	Diverse	2017
[81]	MNIST [44]	-	70,000 (I)	Homo.	10	Pascal and Bing	10 CL.	Handwritten	2016
[81]	AdienceFaces [17]	Flickr	16,252 (I) (G), 16,139 (I) (A)	Hete.	2	10 digits cls.	Gender: 2 Age: 8	Face	2016
<b>Disaster related datasets</b>									
[79]	Incident	Web, SM	446,684 (I) DT:17,511, Info:59,717, Hum:17,769, DS:34,896	NA	1	Incident	43	Incidents	2020
[6]	CrisisBench.	Web, SM	700,000	NA	4	DT, Info, Hum, DS	DT: 7, Info: 2, Hum:4, DS:3	Disaster	2020
[22]	xBD	Satellite	1,654 I/P	NA	-	Building damage	4	Disaster	2019
[8]	MediaEval 2018	SM	18,082	NA	1	Flood	R. and cls.: 2 CL.	Disaster	2018
[4]	CrisisMMD	SM	5878	NA	3	Info, Hum, DS	Info: 2, Hum:8, DS:3	Disaster	2018
[53]	DMD	Web	~25,000	NA	1	Damage	6	Disaster	2018
[55]	DAD	SM	T1: 6,600 (I); T2: 462 I/P	NA	1	DS	3	Disaster	2017
[9]	DIRSM	Flickr	-	NA	1	Flood	R, cls.: 2 CL	Disaster	2017
<b>Our proposed multi-task learning disaster related dataset</b>									
	MEDIC	SM	71,198 (I)	Hete.	4	DT, Info, Hum, DS	DT: 7, Info: 2, Hum:4, DS:3	Disaster	2021

Table 1: Upper part of the table presets the datasets used in multi-task learning studies in computer vision research. Middle part shows disaster related datasets, and the last row shows our proposed dataset. I: Images, V: Videos, S: Samples, SE: Search engines, SM: Social media, DT: disaster types, Info: Informativeness, Hum: Humanitarian, DS: Damage severity. CL.: number of class labels. Hete.: heterogeneous, Homo: Homogeneous. PC: Personal collection. SS: semantic segmentation, HS: human part segmentation, SE: semantic edge detection of surface normals prediction, SD: saliency detection, IS: instance segmentation, SC: scene classification, IC: Indoor scenes, cls.: classification, R/tasks: Ranking tasks, I/P: image patches.

## 2.2 Disaster Response Studies and Datasets

**Social Media Images for Disaster Response** Images posted on social media during disaster plays significant role and the importance of such content was reported in many studies [61, 15, 55, 56, 3, 5]. Recent work include categorizing the severity of damage into discrete levels [55, 56, 3] or quantifying the damage severity as a continuous-valued index [57, 46]. Such developed models were used in real-time disaster response scenarios by engaging with emergency responders [29]. Other related work include adversarial networks for data scarcity issue [45, 62]; disaster image retrieval [2]; image classification in the context of bush fire emergency [41]; flooding photo screening system [58];

<sup>4</sup><https://www.getnexar.com/>



sentiment analysis from disaster image [23]; monitoring natural disasters using satellite images [1]; and flood detection using visual features [33].

**Disaster Response Image Datasets** In crisis informatics<sup>7</sup> research the publicly available image datasets include damage severity assessment dataset (DAD) [56], multimodal dataset (CrisisMMD) [4] and damage identification multimodal dataset (DMD) [53]. The first dataset is only annotated for images, whereas the last two are annotated for both text and images. Other relevant datasets are Disaster Image Retrieval from Social Media (DIRSM) [9] and MediaEval 2018 [8]. The dataset reported in [22] was constructed for detecting damage as an anomaly using pre-and post-disaster images. It consists of 700,000 building annotations. A similar and relevant work is the Incidents dataset [79], which consists of 446,684 manually labeled images with 43 incident categories. The *Crisis Benchmark Dataset* reported in [6] is the largest social media disaster image classification dataset, which is a consolidated version of DAD, CrisisMMD, DMD, and additional labeled images. For this study, we extended the *Crisis Benchmark Dataset*. To make the dataset for multi-task learning, we additionally labeled with 155,899 more labels, which resulted in the whole dataset being aligned for such a setup.

### 3 MEDIC Dataset

The MEDIC dataset consists of four different disaster-related tasks that are important for humanitarian aid.<sup>8</sup> These tasks are defined based on prior work experience with the humanitarian response organizations such as UN-OCHA and existing literature [31, 30, 4, 5]. In this section, we first provide the details of each task and class labels and then discuss the annotation details of the dataset.

#### 3.1 Tasks

**Disaster types** During man-made and natural disasters, people post textual and visual content about the current situation, and the real-time social media monitoring system requires to detect an event when ingesting images from unfiltered social media streams. For the disaster scenario, it is important to automatically detect different disaster types from the crawled social media images. For instance, an image can depict a wildfire, flood, earthquake, hurricane, and other types of disasters. Different categories (i.e., natural, human-induced, and hybrid) and sub-categories of disaster types have been defined in the literature [69]. This research focuses on major disaster events that include (i) earthquake, (ii) fire, (iii) flood, (iv) hurricane, (v) landslide, (vi) other disaster, which covers all other types (e.g., plane, train crash), and (vii) not disaster, which includes the images that do not show any identifiable disasters.

**Informativeness** Social media contents are often noisy and contain numerous irrelevant images such as cartoons, advertisements, etc. In addition to this, the clean images that show damaged infrastructure due to flood, fire, or any other disaster events are crucial for humanitarian response tasks. Therefore, it is necessary to eliminate any irrelevant or redundant content to facilitate crisis responders' efforts more effectively. For this purpose, we define the *informativeness* task as to filter out irrelevant images, where the class labels consist (i) informative and (ii) not informative.

**Humanitarian** Fine-grained categorization of certain information significantly helps the emergency crisis responders to make an efficient actionable decision. Humanitarian categories vary depending on the type of content (text vs. image). For example, the CrisisBench dataset [7] consists of tweets labeled with 11 categories, whereas CrisisMMD [4] multimodal dataset consists of 8 categories. Such variation exists between text and images because some information can easily be presented in one modality than another modality. For example, it is possible to report *missing or found people* in text than in an image, which is also reported in [4]. This research focuses on these factors and considers the four most important categories that are useful 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.

---

<sup>7</sup>[https://en.wikipedia.org/wiki/Disaster\\_informatics](https://en.wikipedia.org/wiki/Disaster_informatics)

<sup>8</sup>[https://en.wikipedia.org/wiki/Humanitarian\\_aid](https://en.wikipedia.org/wiki/Humanitarian_aid)

Source	Event name	Year	# images	Source	Event name	Year	# images
Twitter	Typhoon ruby/hagupit	2014	833	Twitter	Iraq iran earthquake	2017	596
Twitter	Nepal earthquake	2015	21710	Twitter	Mexico earthquake	2017	1378
Twitter	South India floods	2015	1476	Twitter	Srilanka floods	2017	1022
Twitter	Illapel earthquake	2015	403	Twitter	Ukraine conflict	2017	240
Twitter	Food insecurity in yemen	2015	466	Twitter	Greece wildfire	2018	351
Twitter	Paris attack	2015	1043	Twitter	Hurricane florence	2018	186
Twitter	South India floods	2015	753	Twitter	Hurricane michael	2018	219
Twitter	Syria attacks	2015	350	Twitter	Kerala flood	2018	605
Twitter	Terremotoitalia	2015	919	Twitter	Typhoon mangkhut	2018	172
Twitter	Ecuador earthquake	2016	2280	Google	NA	NA	3007
Twitter	Hurricane matthew	2016	596	Twitter	Human induced disaster	NA	501
Twitter	California wildfires	2017	1585	G, B, F	NA	NA	1263
Twitter	Hurricane harvey	2017	5644	Twitter	Natural disaster	NA	6597
Twitter	Hurricane irma	2017	4973	Twitter	Security incidents activities	NA	1082
Twitter	Hurricane maria	2017	5069	G, I.	NA	NA	5879

Table 2: Data collection source, event name, year of the event and number of image annotated. G: Google, B: Bing, F: Flickr, I: Instagram.

**Damage severity** Detecting the severity of the damage is significantly important to help the affected community during disaster events. The severity of the damage can be assessed from an image based on the visual appearance of the physical destruction of a built structure (e.g., bridges, roads, buildings, burned houses, and forests). Following the work reported in [56], this research defines the following categories for the classification task (*i*) severe damage, (*ii*) mild damage, and (*iii*) little or none.

### 3.2 Datasets

#### 3.2.1 Data Curation

This research extends the labels of the Crisis Benchmark dataset reported in [6]. This Crisis Benchmark dataset has been developed by consolidating existing datasets and labeling new data for disaster type. This Crisis Benchmark dataset consists of images collected from Twitter, Google, Bing, Flickr, and Instagram. The majority of the datasets have been collected from Twitter, as shown in Table 2. The Twitter data were mainly collected during major disaster events<sup>9</sup> and using different disaster-specific keywords. The data collected from Google, Bing, Flickr, and Instagram are based on specific keywords. The dataset is diverse in terms of (*i*) number of events, (*ii*) different time frames spanning over five years, (*iii*) natural (e.g., earthquake, fire, floods) and man-made disasters (e.g., Paris attack, Syria attacks), and (*iv*) events occurred in different part of the world. The number of images in different events resulted from different factors, such as the number of tweets collected during the disaster events, the number of images crawled, filtered due to duplicates, and a random selection for the annotation. Our motivation for choosing and extending the Crisis Benchmark dataset is that it reduced the overall cost of data collection and annotation processes while also having a large dataset for multi-task learning.

#### 3.2.2 Annotation

For the manual annotation, we used Appen<sup>10</sup> crowdsourcing annotation platform. In such a platform, finding qualified workers and managing the quality of the annotation is an important issue. To ensure the quality, we used the widely used gold standard evaluation approach [13]. We designed the interface with annotation guidelines on Appen for the annotation task (see Figure 7 in Appendix). We followed the annotation guidelines from previous work [4, 6] and improved with examples for this task (see the detailed annotation guidelines with examples in Appendix A). For all tasks, we choose to annotate in a multiclass setting even though *humanitarian* and *disaster type* tasks in our context are more suitable to be framed as pure multi-label. Our decision has been influenced by several factors. The most important one was our consultation with humanitarian organizations which suggested limiting the number of classes by merging related ones and keeping only the most important information types. This is due to the information overload issue that humanitarian responders often deal with at the onset of a disaster situation if exposed to information types not important for them. Furthermore, obtaining a sufficient number of labeled instances for a large number of classes to train a pure multi-label classifier is not practical due to both annotation budget (e.g., time, cost) and

<sup>9</sup>Event names reported in Table 2 are based on Wikipedia.

<sup>10</sup><https://appen.com/>

Tasks	Fleiss ( $\kappa$ )	Krip. ( $\alpha$ )	Avg agg.	Tasks	Fleiss ( $\kappa$ )	Krip. ( $\alpha$ )	Avg agg.
Disaster types	0.46	0.46	0.70	Humanitarian	0.52	0.52	0.73
Informativeness	0.71	0.71	0.91	Damage severity	0.55	0.55	0.79

Table 3: Annotation agreement for different tasks. Fleiss Kappa ( $\kappa$ ), Krip. ( $\alpha$ ): Krippendorff’s  $\alpha$ , Avg agg.: Average observed agreement.

Class labels	Train	Dev	Test	Total	Class labels	Train	Dev	Test	Total
<b>Disaster types</b>					<b>Humanitarian</b>				
Earthquake	12,023	929	1,674	14,626	Affected, injured, or dead people	3,103	255	615	3,973
Fire	1,737	250	688	2,675	Infrastructure and utility damage	18,182	2,322	5,030	25,534
Flood	3,269	559	1,300	5,128	Not humanitarian	25,828	3,212	9,104	38,144
Hurricane	3,801	572	1,408	5,781	Rescue volunteering or donation effort	2,240	368	939	3,547
Landslide	1,046	161	327	1,534	<b>Total</b>	49,353	6,157	15,688	71,198
Not disaster	25,463	3,301	9,078	37,842	<b>Damage severity</b>				
Other disaster	2,014	385	1,213	3,612	Little or none	27,015	3,460	9,886	30,475
<b>Total</b>	49,353	6,157	15,688	71,198	Mild	4,406	812	1,708	5,218
<b>Informativeness</b>					Severe	17,932	1,885	4,094	19,817
Informative	30,547	3,699	8,603	42,849	<b>Total</b>	49,353	6,157	15,688	71,198
Not informative	18,806	2,458	7,085	28,349					
<b>Total</b>	49,353	6,157	15,688	71,198					

Table 4: Annotated dataset with data splits for different tasks.

modeling perspectives (e.g., high imbalance). For the image, which can have multiple labels, we instructed the annotators to select the label that is more important for humanitarian organizations and prominent in the image.

For the annotation, we designed a *hit* consists of 5 images. For the gold standard evaluation, we manually labeled 100 images, which are randomly assigned to the hit for the evaluation. We assigned a criterion to have at least 3 annotations per image and per task. An agreement score of 66% is used to select the final label, which ensured that at least two annotators agreed on a label. The hit was extended to more annotators if such a criterion was not met.

Since the Crisis Benchmark dataset did have task-specific labels for all images, i.e., different sets of images consisted of labels for three tasks and two tasks; therefore, we first prepared the different sets with missing labels for the annotation. For example, 25,731 images of the Crisis Benchmark dataset did not have labels for disaster types and humanitarian tasks, which we selected for the annotation tasks. In this way, we run the annotation tasks in different batches.

### 3.2.3 Crowdsourcing Results

To measure the quality of the annotation, we compute the annotation agreement using Fleiss kappa [19], Krippendorff’s alpha [38] and average observed agreement [19]. In Table 3, we present the annotation agreement for all events with different approaches mentioned above. The agreement score varies from 46% to 71% for different tasks. Note that, in the Kappa measurement, the values of ranges 0.41-0.60, 0.61-0.80, and 0.81-1 refers to moderate, substantial, and perfect agreement, respectively [43]. Based on these measurements, we conclude that our annotation agreement score leads to moderate to substantial agreement. The number of labels and subjectivity of the annotation tasks reflected the annotation agreement score. Some annotation tasks are highly subjective. For example, for the disaster-type task, hurricane or tropical cyclones often leads to heavy rain, which causes flood (e.g., an image showing a fallen tree with flood water) can be annotated as hurricane or flood. Another example is an image showing building damage and rescue effort. In such cases, the annotation task was to carefully check what is more visible in the image and select the label accordingly. Note that, the agreement score for disaster types is comparatively lower than other tasks, which is due to the high level of subjectivity in the annotation task. Annotators needed to choose one label among seven labels. The average agreement scores are comparatively higher as we made sure at least two annotators agree on a label.

After completing the annotation task, the proposed dataset added 155,899 annotated labels for four tasks in addition to the existing 128,893 labels from 71,198 images. In total, this research re-annotated 65,640 images to create the MEDIC dataset.

## 4 Experiments and Results

In Table 4, we present the dataset with task-wise data splits and distribution, which consists of 69%, 9%, and 22% for training, development, and test set respectively. We first conduct baseline experiment, followed by single task learning experiment to compare and provide a benchmark with a multitask setting.

To measure the performance of each classifier and for each task setting, we use weighted average precision (P), recall (R), and F1-score (F1), which has been widely used in the literature.

### 4.1 Baseline

For the baseline experiment we use (i) *a majority class baseline*, (ii) feature from a pre-trained model, then training and evaluation using SVM and KNN. We extracted features from the penultimate layer of the EfficientNet b1 model, which is trained using ImageNet. The majority class baseline predicts the label based on the most frequent label in the training set. This has been most commonly used in shared tasks [54]. For training SVM and KNN we used default parameters setting.

### 4.2 Single-Task Learning

We used several pre-trained models for single-task learning and fine-tuned the network with the task-specific classification layer on top of the network. This approach has been popular and has been performing well for various downstream visual recognition tasks [83, 70, 60, 59]. The network architectures that we used in this study include ResNet18, ResNet50, ResNet101 [24], VGG16 [72], DenseNet [26], SqueezeNet [28], MobileNet [25], and EfficientNet [75]. We have chosen such diverse architectures to understand their relative performance and inference time. For fine-tuning, we use the weights of the networks pre-trained using ImageNet [16] to initialize our model. Our classification settings comprised binary (i.e., informativeness task) and multiclass settings (i.e., remaining three tasks). We train the models using the Adam optimizer [36] with an initial learning rate of  $10^{-3}$ , which is decreased by a factor of 10 when accuracy on the dev set stops improving for 10 epochs. The models were trained for 150 epochs. We use the model with the best accuracy on the validation set to evaluate on the test split.

### 4.3 Multi-Task Learning

In the MEDIC dataset, the tasks share similar properties; hence, we designed a simpler approach. We use the hard parameter sharing approach to reduce the computational complexity. All tasks share the same feature layers in the network, which is followed by task-specific classification layers. For optimizing the loss, we provide equal weight to each task. Assuming that the task-specific weight is  $w_i$  and task-specific loss function is  $\mathcal{L}_i$ , the optimization objective of the MTL is defined as  $\mathcal{L}_{MTL} = \sum_i w_i \cdot \mathcal{L}_i$ . During optimization (i.e., using stochastic gradient descent to minimize the objective), the network weights in the shared layers  $W_{sh}$  are updated using the following equation:

$$W_{sh} = \sum_i W_{sh} - \lambda \sum_i w_i \frac{\partial \mathcal{L}_i}{\partial W_{sh}} \quad (1)$$

We set  $w_i = 1$  in our experiments for all task-specific weights, i.e., equal weight for all tasks. We use softmax activation to get probability distribution over individual tasks and use cross-entropy as a loss function. We initialized the weight using pre-trained models mentioned above, which are trained using ImageNet.

Our implementation of multi-task learning supports all the network architectures mentioned in section 4.2. Therefore, we have run experiments using the same pre-trained models and same hyper-parameter settings for the MTL experiments.

We used the NVIDIA Tesla V100-SXM2-16 GB GPU machines consisting of 12 cores and 40GB CPU memory for all experiments.

Model	Acc	P	R	F1	Acc	P	R	F1	Acc	P	R	F1	Acc	P	R	F1
	Disaster types				Informative				Humanitarian				Damage severity			
Majority	57.9	33.5	57.9	42.4	54.8	30.1	54.8	38.8	58.0	33.7	58.0	42.6	63.0	39.7	63.0	48.7
Eff. Net Feat. + SVM	75.7	74.1	75.7	<b>73.2</b>	83.0	83.0	83.0	83.0	77.9	<b>76.1</b>	77.9	76.1	78.4	75.4	78.4	<b>75.3</b>
Eff. Net Feat. + KNN	71.0	71.7	71.0	69.9	80.4	80.3	80.4	80.3	75.3	74.8	75.3	74.6	78.3	75.1	78.3	75.1

Table 5: Baseline classification results. Eff. Net Feat.: Feature extracted from the penultimate layer of a pre-trained efficient net model.

Model	Acc	P	R	F1	Acc	P	R	F1	Acc	P	R	F1	Acc	P	R	F1
Disaster types									Humanitarian							
Single task				Multi-task				Single task				Multi-task				
ResNet18	79.1	77.8	79.1	<b>76.9</b>	78.8	77.8	78.8	76.3	79.6	78.0	79.6	<b>78.3</b>	79.7	78.0	79.7	77.8
ResNet50	79.5	78.4	79.5	77.9	80.1	79.1	80.1	<b>78.1</b>	80.5	79.1	80.5	79.3	81.0	79.4	81.0	<b>79.5</b>
ResNet101	79.8	78.3	79.8	78.4	80.6	80.1	80.6	<b>78.7</b>	80.2	78.7	80.2	<b>78.9</b>	81.0	79.7	81.0	79.8
VGG16	78.8	77.4	78.8	<b>77.2</b>	79.7	79.7	79.7	77.1	80.2	78.7	80.2	78.8	81.0	79.5	81.0	79.4
DenseNet (121)	80.3	79.6	80.3	<b>78.4</b>	80.3	79.6	80.3	<b>78.4</b>	80.3	78.8	80.3	<b>78.9</b>	80.7	79.2	80.7	79.4
SqueezeNet	76.5	74.7	76.5	<b>73.9</b>	76.2	74.4	76.2	73.3	77.9	75.9	77.9	75.9	78.2	75.8	78.2	<b>76.0</b>
MobileNet (v2)	78.7	77.4	78.7	76.8	79.2	78.3	79.2	<b>77.2</b>	80.2	77.8	80.2	78.1	80.3	78.5	80.3	<b>78.5</b>
EfficientNet (b1)	81.0	80.2	81.0	<b>79.6</b>	80.9	80.1	80.9	79.3	81.0	79.9	81.0	80.1	81.4	80.2	81.4	<b>80.4</b>
Informative									Damage severity							
Single task				Multi-task				Single task				Multi-task				
ResNet18	84.2	84.2	84.2	84.2	84.6	84.6	84.6	<b>84.6</b>	79.9	77.9	79.9	<b>78.2</b>	80.1	77.4	80.1	77.6
ResNet50	85.6	85.6	85.6	85.6	85.4	85.6	85.4	<b>85.5</b>	81.0	78.7	81.0	78.8	81.4	79.1	81.4	<b>79.5</b>
ResNet101	84.5	84.5	84.5	84.5	85.5	85.5	85.5	<b>85.5</b>	81.0	78.3	81.0	78.4	81.5	79.4	81.5	<b>79.7</b>
VGG16	84.8	85.1	84.8	84.9	85.7	85.7	85.7	85.7	80.9	78.4	80.9	78.6	81.6	79.7	81.6	79.1
DenseNet (121)	84.9	85.0	84.9	<b>84.9</b>	84.9	84.9	84.9	<b>84.9</b>	80.6	78.1	80.6	78.4	81.4	79.2	81.4	<b>79.5</b>
SqueezeNet	82.4	82.4	82.4	82.4	82.8	82.8	82.8	<b>82.8</b>	78.7	75.9	78.7	<b>76.3</b>	78.9	75.7	78.9	76.2
MobileNet (v2)	84.0	84.0	84.0	84.0	84.7	84.7	84.7	<b>84.7</b>	80.2	77.8	80.2	78.1	80.6	78.3	80.6	<b>78.8</b>
EfficientNet (b1)	84.9	85.3	84.9	85.0	86.0	86.1	86.0	<b>86.0</b>	81.3	79.4	81.3	79.9	81.8	80.1	81.8	<b>80.3</b>

Table 6: Classification results using single and multi-task settings along with different pre-trained models. Best F1 scores are highlighted.

#### 4.4 Results

In Table 5, we provide baseline results. From the majority baseline results it is clear that imbalance distribution does not play any role. Among SVM and KNN, the former is performing well in all tasks with 0.2 to 3.3% improvement.

In Table 6, we report the results for both single and multi-tasks settings using the mentioned models. Across different models, overall, EfficientNet (b1) performs better than other models. Comparing only EfficientNet (b1) models’ results for all tasks, the multi-task setting shows better than single task settings; although, the difference is minor and might not be significant. However, since we share the feature layers across the four tasks, model space requirement and inference time are reduced by a factor of four. The improved inference time is crucial for real-time disaster response systems as it reduces the operational cost that running individual models would incur.

## 5 Discussion and Future Work

The MEDIC dataset provides images from diverse events consisting of different time frames. The crowd-sourced annotation provides a reasonable annotation agreement even though they are subjective. Our experiments show that multi-task learning with neural net reduces computational complexity significantly while having comparative performance.

In Figure 2, we show the loss and accuracy plots for single and multi-task settings for EfficientNet (b1) model. We limit the plots to 40 epochs as all of the models converged by then. We notice similar convergence rates for both single and multi-task learning setups. We observe that the multi-task objective function acts as a regularizer as the training loss is consistently higher and training accuracy is lower than the single-task setting while having similar or better performance on the validation set. This suggests that the multi-task setup may benefit from models having a larger capacity.

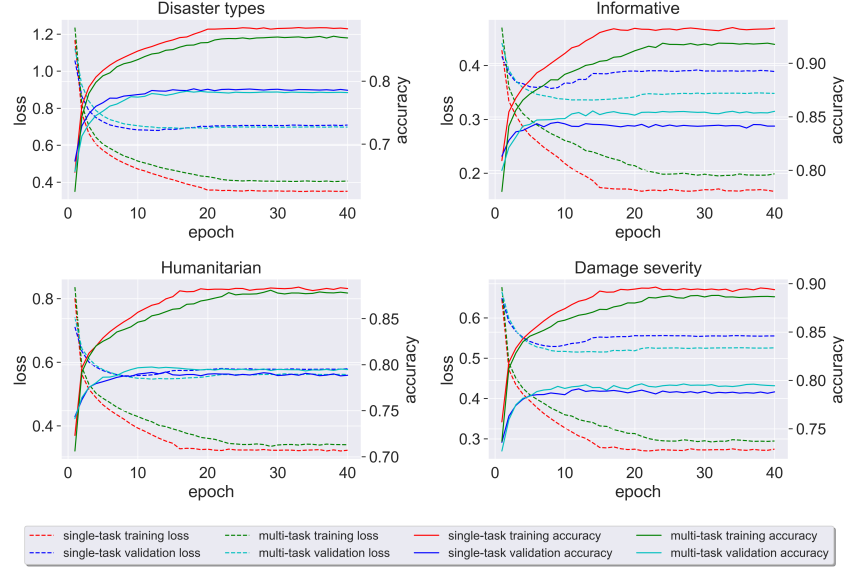


Figure 2: Training and validation loss and accuracy for EfficientNet (b1) model for single and multi-task settings.

Class distribution is an important issue that affect classifier performance. We investigated class-wise performances and confusion matrix. Our observation suggests that imbalance class distribution is not only factor for lower classification performance in certain classes. It also depends on distinguishing properties of the class label. For example, the distribution of *Fire* class label is 3.8% in the dataset but the performance is third-best among class labels. Where the distribution of *Other disaster* is 5.1%, however, the F1 is 27.0, which is the lowest performance. In appendix Section C, Table 8, we reported class-wise results.

To understand the task correlation and how they affect performance, we also run experiments with different subsets of the tasks (see Table 10 in Appendix). We obtain similar results with other task combinations. It will be an important future research avenue to explore different weighting schemes for the tasks. Regardless, our reported results can serve as a baseline for single and multi-task disaster image classification.

**Limitation** We foresee several limitations of our work. As mentioned earlier disaster types and humanitarian tasks can be annotated with multiple labels, which we annotated with the single label in this study. Even though our choice has been influenced based on the knowledge of humanitarian organizations, however, we aim to explore it further. In our experiments, we may have to explore a much larger network, which can help multi-task learning better.

**Future Work:** Our future work will include annotating images with multilabel annotation, exploring other multi-task learning methods, and investigating tasks groups and relationships. For example, it would be interesting to know why training the model with disaster types, informativeness and humanitarian tasks reduces performance as presented in Table 10. Other research avenues include multimodality (e.g., integrating text), and investigating class imbalance issues.

## 6 Conclusions

We presented a large manually annotated multi-task learning dataset, consists of 71,198 images, labeled for four tasks, which were specifically designed for multitask learning research and disaster response image classification. The dataset will not only useful to develop robust models for disaster response tasks but also will enable to evaluate the multi-task models. We provide classification results using nine different pre-trained models, which can serve as a benchmark in future work. We report that the multi-task model reduces the inference time significantly, hence, such a model can be very useful for real-time classification tasks, especially to classify social media image streams.

## References

- [1] Kashif Ahmad, Michael Riegler, Konstantin Pogorelov, Nicola Conci, Pål Halvorsen, and Francesco De Natale. Jord: a system for collecting information and monitoring natural disasters by linking social media with satellite imagery. In *Proceedings of the 15th International Workshop on Content-Based Multimedia Indexing*, pages 1–6, 2017.
- [2] Sheharyar Ahmad, Kashif Ahmad, Nasir Ahmad, and Nicola Conci. Convolutional neural networks for disaster images retrieval. In *MediaEval*, 2017.
- [3] Firoj Alam, Muhammad Imran, and Ferda Ofli. Image4Act: Online social media image processing for disaster response. In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 1–4, 2017.
- [4] Firoj Alam, Ferda Ofli, and Muhammad Imran. CrisisMMD: multimodal twitter datasets from natural disasters. In *Proceedings of the International AAAI Conference on Web and Social Media*, pages 465–473, Jun 2018.
- [5] Firoj Alam, Ferda Ofli, and Muhammad Imran. Processing social media images by combining human and machine computing during crises. *International Journal of Human Computer Interaction*, 34(4):311–327, 2018.
- [6] Firoj Alam, Ferda Ofli, Muhammad Imran, Tanvirul Alam, and Umair Qazi. Deep learning benchmarks and datasets for social media image classification for disaster response. In *IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pages 151–158, 2020.
- [7] Firoj Alam, Hassan Sajjad, Muhammad Imran, and Ferda Ofli. CrisisBench: Benchmarking crisis-related social media datasets for humanitarian information processing. In *Proceedings of the International AAAI Conference on Web and Social Media*, 2021.
- [8] Bischke Benjamin, Helber Patrick, Zhao Zhengyu, Bruijn Jens de, and Borth Damian. The multimedia satellite task at MediaEval 2018: Emergency response for flooding events. In *MediaEval*, Oct 2018.
- [9] Benjamin Bischke, Patrick Helber, Christian Schulze, Venkat Srinivasan, Andreas Dengel, and Damian Borth. The multimedia satellite task at MediaEval 2017. In *In Proceedings of the MediaEval 2017: MediaEval Benchmark Workshop*, 2017.
- [10] Rich Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [11] Xianjie Chen, Roozbeh Mottaghi, Xiaobai Liu, Sanja Fidler, Raquel Urtasun, and Alan Yuille. Detect what you can: Detecting and representing objects using holistic models and body parts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1971–1978, 2014.
- [12] Zhao Chen, Vijay Badrinarayanan, Chen-Yu Lee, and Andrew Rabinovich. Gradnorm: Gradient normalization for adaptive loss balancing in deep multitask networks. In *International Conference on Machine Learning*, pages 794–803. PMLR, 2018.
- [13] Shammur Absar Chowdhury, Marcos Calvo, Arindam Ghosh, Evgeny A Stepanov, Ali Orkan Bayer, Giuseppe Riccardi, Fernando García, and Emilio Sanchis. Selection and aggregation techniques for crowdsourced semantic annotation task. In *Sixteenth Annual Conference of the International Speech Communication Association*. ISCA, 2015.
- [14] Michael Crawshaw. Multi-task learning with deep neural networks: A survey. *arXiv preprint arXiv:2009.09796*, 2020.
- [15] Shannon Daly and J Thom. Mining and classifying image posts on social media to analyse fires. In *Proceedings of the International Conference on Information Systems for Crisis Response and Management (ISCRAM)*, pages 1–14, 2016.
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255, 2009.
- [17] Eran Eiding, Roe Enbar, and Tal Hassner. Age and gender estimation of unfiltered faces. *IEEE Transactions on Information Forensics and Security*, 9(12):2170–2179, 2014.
- [18] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.
- [19] Joseph L Fleiss, Bruce Levin, and Myunghee Cho Paik. *Statistical methods for rates and proportions*. 2013.
- [20] Yuan Gao, Jiayi Ma, Mingbo Zhao, Wei Liu, and Alan L Yuille. Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3205–3214, 2019.

- [21] Boqing Gong, Yuan Shi, Fei Sha, and Kristen Grauman. Geodesic flow kernel for unsupervised domain adaptation. In *2012 IEEE conference on computer vision and pattern recognition*, pages 2066–2073. IEEE, 2012.
- [22] Ritwik Gupta, Bryce Goodman, Nirav Patel, Ricky Hosfelt, Sandra Sajeev, Eric Heim, Jigar Doshi, Keane Lucas, Howie Choset, and Matthew Gaston. Creating xbd: A dataset for assessing building damage from satellite imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- [23] Syed Zohaib Hassan, Kashif Ahmad, Ala Al-Fuqaha, and Nicola Conci. Sentiment analysis from images of natural disasters. In *International Conference on Image Analysis and Processing*, pages 104–113. Springer, 2019.
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [25] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv:1704.04861*, 2017.
- [26] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017.
- [27] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on pattern analysis and machine intelligence*, 16(5):550–554, 1994.
- [28] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5 mb model size. *arXiv:1602.07360*, 2016.
- [29] Muhammad Imran, Firoj Alam, Umair Qazi, Steve Peterson, and Ferda Ofli. Rapid damage assessment using social media images by combining human and machine intelligence. In *17th International Conference on Information Systems for Crisis Response and Management*, pages 761–773, 2020.
- [30] Muhammad Imran, Carlos Castillo, Fernando Diaz, and Sarah Vieweg. Processing social media messages in mass emergency: A survey. *ACM Computing Surveys*, 47(4):67, 2015.
- [31] Muhammad Imran, Carlos Castillo, Ji Lucas, Patrick Meier, and Sarah Vieweg. AIDR: Artificial intelligence for disaster response. In *Proceedings of the 23rd international conference on world wide web*, pages 159–162, 2014.
- [32] Muhammad Imran, Ferda Ofli, Doina Caragea, and Antonio Torralba. Using ai and social media multi-modal content for disaster response and management: Opportunities, challenges, and future directions. *Information Processing & Management*, 57(5):102261, 2020.
- [33] Rabiul Islam Jony, Alan Woodley, and Dimitri Perrin. Flood detection in social media images using visual features and metadata. *2019 Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–8, 2019.
- [34] Zhuoliang Kang, Kristen Grauman, and Fei Sha. Learning with whom to share in multi-task feature learning. In *International Conference on Machine Learning*, 2011.
- [35] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018.
- [36] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [37] Iasonas Kokkinos. Ubernet: Training a universal convolutional neural network for low-, mid-, and high-level vision using diverse datasets and limited memory. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6129–6138, 2017.
- [38] Klaus Krippendorff. Estimating the reliability, systematic error and random error of interval data. *Educational and Psychological Measurement*, 30(1):61–70, 1970.
- [39] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Master’s thesis, Department of Computer Science, University of Toronto*, 2009.
- [40] Abhishek Kumar and Hal Daumé III. Learning task grouping and overlap in multi-task learning. In *Proceedings of the 29th International Conference on Machine Learning*, pages 1723–1730, 2012.
- [41] Ryan Lagerstrom, Yulia Arzhaeva, Piotr Szul, Oliver Obst, Robert Power, Bella Robinson, and Tomasz Bednarsz. Image classification to support emergency situation awareness. *Frontiers in Robotics and AI*, 3:54, 2016.



- [42] Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- [43] J Richard Landis and Gary G Koch. The measurement of observer agreement for categorical data. *biometrics*, pages 159–174, 1977.
- [44] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [45] Xukun Li, Doina Caragea, Cornelia Caragea, Muhammad Imran, and Ferda Ofli. Identifying disaster damage images using a domain adaptation approach. In *Proceeding of the International Conference on Information Systems for Crisis Response and Management (ISCRAM)*, pages 633–645, 2019.
- [46] Xukun Li, Doina Caragea, Huaiyu Zhang, and Muhammad Imran. Localizing and quantifying damage in social media images. In *2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pages 194–201, 2018.
- [47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.
- [48] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *European Conference on Computer Vision (ECCV)*, pages 740–755, 2014.
- [49] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE international conference on computer vision*, pages 3730–3738, 2015.
- [50] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3431–3440, 2015.
- [51] Mingsheng Long, ZHANGJIE CAO, Jianmin Wang, and Philip S Yu. Learning multiple tasks with multilinear relationship networks. *Advances in Neural Information Processing Systems*, 30:1594–1603, 2017.
- [52] Ishan Misra, Abhinav Shrivastava, Abhinav Gupta, and Martial Hebert. Cross-stitch networks for multi-task learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3994–4003, 2016.
- [53] Hussein Mouzannar, Yara Rizk, and Mariette Awad. Damage Identification in Social Media Posts using Multimodal Deep Learning. In *Proceedings of the International Conference on Information Systems for Crisis Response and Management (ISCRAM)*, pages 529–543, May 2018.
- [54] Preslav Nakov, Giovanni Da San Martino, Tamer Elsayed, Alberto Barrón-Cedeño, Rubén Míguez, Shaden Shaar, Firoj Alam, Fatima Haouari, Maram Hasanain, Watheq Mansour, Bayan Hamdan, Zien Sheikh Ali, Nikolay Babulov, Alex Nikolov, Gautam Kishore Shahi, Julia Maria Struß, Thomas Mandl, Mucahid Kutlu, and Yavuz Selim Kartal. Overview of the CLEF-2021 CheckThat! lab on detecting check-worthy claims, previously fact-checked claims, and fake news. In K. Selcuk Candan, Bogdan Ionescu, Lorraine Goeuriot, Birger Larsen, Henning Müller, Alexis Joly, Maria Maistro, Florina Piroi, Guglielmo Faggioli, and Nicola Ferro, editors, *Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Twelfth International Conference of the CLEF Association*, LNCS (12880). Springer, 2021.
- [55] Dat Tien Nguyen, Firoj Alam, Ferda Ofli, and Muhammad Imran. Automatic image filtering on social networks using deep learning and perceptual hashing during crises. In *Proc. of ISCRAM*, May 2017.
- [56] Dat Tien Nguyen, Ferda Ofli, Muhammad Imran, and Prasenjit Mitra. Damage assessment from social media imagery data during disasters. In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, pages 1–8, Aug 2017.
- [57] Karoon Rashedi Nia and Greg Mori. Building damage assessment using deep learning and ground-level image data. In *14th Conference on Computer and Robot Vision (CRV)*, pages 95–102. IEEE, 2017.
- [58] Huan Ning, Zhenlong Li, Michael E Hodgson, et al. Prototyping a social media flooding photo screening system based on deep learning. *ISPRS International Journal of Geo-Information*, 9(2):104, 2020.
- [59] Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1717–1724, 2014.
- [60] G. Ozbulak, Y. Aytaç, and H. K. Ekenel. How transferable are cnn-based features for age and gender classification? In *International Conference of the Biometrics Special Interest Group*, pages 1–6, Sept 2016.
- [61] Robin Peters and Joao Porto de Albuquerque. Investigating images as indicators for relevant social media messages in disaster management. In *Proceedings of the International Conference on Information Systems for Crisis Response and Management (ISCRAM)*, 2015.

- [62] Samira Pouyanfar, Yudong Tao, Saad Sadiq, Haiman Tian, Yuexuan Tu, Tianyi Wang, Shu-Ching Chen, and Mei-Ling Shyu. Unconstrained flood event detection using adversarial data augmentation. In *IEEE International Conference on Image Processing (ICIP)*, pages 155–159, 2019.
- [63] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [64] Sebastian Ruder. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*, 2017.
- [65] Sebastian Ruder, Joachim Bingel, Isabelle Augenstein, and Anders Søgaard. Latent multi-task architecture learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 4822–4829, 2019.
- [66] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015.
- [67] Naina Said, Kashif Ahmad, Michael Riegler, Konstantin Pogorelov, Laiq Hassan, Nasir Ahmad, and Nicola Conci. Natural disasters detection in social media and satellite imagery: a survey. *Multimedia Tools and Applications*, 78(22):31267–31302, 2019.
- [68] Roy Schwartz, Jesse Dodge, Noah A Smith, and Oren Etzioni. Green AI. *Communications of the ACM*, 63(12):54–63, 2020.
- [69] Ibrahim Mohamed Shaluf. Disaster types. *Disaster Prevention and Management: An International Journal*, 2007.
- [70] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan, and Stefan Carlsson. CNN features off-the-shelf: an astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition Workshops*, pages 806–813, 2014.
- [71] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European conference on computer vision*, pages 746–760. Springer, 2012.
- [72] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [73] Gjorgji Strezoski, Nanne van Noord, and Marcel Worring. Learning task relatedness in multi-task learning for images in context. In *Proceedings of the 2019 on International Conference on Multimedia Retrieval*, pages 78–86, 2019.
- [74] Gjorgji Strezoski and Marcel Worring. Omniart: a large-scale artistic benchmark. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(4):1–21, 2018.
- [75] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv:1905.11946*, 2019.
- [76] Simon Vandenhende, Stamatios Georgoulis, Wouter Van Gansbeke, Marc Proesmans, Dengxin Dai, and Luc Van Gool. Multi-task learning for dense prediction tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [77] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5018–5027, 2017.
- [78] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- [79] Ethan Weber, Nuria Marzo, Dim P Papadopoulos, Aritro Biswas, Agata Lapedriza, Ferda Ofli, Muhammad Imran, and Antonio Torralba. Detecting natural disasters, damage, and incidents in the wild. In *European Conference on Computer Vision*, pages 331–350. Springer, 2020.
- [80] Joseph Worsham and Jugal Kalita. Multi-task learning for natural language processing in the 2020s: where are we going? *Pattern Recognition Letters*, 136:120–126, 2020.
- [81] Yongxin Yang and Timothy Hospedales. Deep multi-task representation learning: A tensor factorisation approach. *arXiv preprint arXiv:1605.06391*, 2016.
- [82] Yongxin Yang and Timothy Hospedales. Deep multi-task representation learning: A tensor factorisation approach. In *5th International Conference on Learning Representations*, 2017.
- [83] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *Advances in Neural Information Processing Systems*, pages 3320–3328, 2014.
- [84] Fisher Yu, Haofeng Chen, Xin Wang, Wenqi Xian, Yingying Chen, Fangchen Liu, Vashisht Madhavan, and Trevor Darrell. BDD100k: A diverse driving dataset for heterogeneous multitask learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2636–2645, 2020.

- [85] Amir R Zamir, Alexander Sax, William Shen, Leonidas J Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3712–3722, 2018.
- [86] Yu Zhang and Qiang Yang. A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering*, 2021.

## 7 Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
  - (b) Did you describe the limitations of your work? **[Yes]**. See limitations in section 5.
  - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See in Section B.5.2.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
  - (b) Did you include complete proofs of all theoretical results? **[N/A]**
3. If you ran experiments (e.g. for benchmarks)...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** See in appendix B.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** See implementation details in section 3 and 4.
  - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? We will make them available in the final version.
  - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **Yes**. See such details in section 4.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **[Yes]**
  - (b) Did you mention the license of the assets? **[Yes]** See dataset details in appendix B.
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]**
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[N/A]**
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[Yes]** See in appendix B.
5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[Yes]** See in appendix B.
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[Yes]** See in appendix A.4.

## Appendix

### A Data Collection

#### A.1 Data Curation and Annotation

We extended the Crisis Benchmark dataset to develop MEDIC, a multitask learning dataset for disaster response. For the annotation, we provided detailed instructions to the annotators, which they followed during the annotation tasks. Our annotation consists of four tasks in different batches, and we provided task-specific instructions along with them.

#### A.2 Annotation Instructions

The annotation task involves identifying images that are useful for humanitarian aid/response. During different disaster events (i.e., natural and human-induced or hybrid), *humanitarian aid*<sup>11</sup> involves assisting people who need help. The primary purpose of humanitarian aid is to save lives, reduce suffering, and rebuild affected communities. Among the people in need belong homeless, refugees, and victims of natural disasters, wars, and conflicts who need necessities like food, water, shelter, medical assistance, and damage-free critical infrastructure and utilities such as roads, bridges, power lines, and communication poles.

For disaster types and humanitarian tasks, it is possible that some images can be annotated with multiple labels. In such cases, the instruction is to choose a label that is critical (i.e., higher priority) for humanitarian organizations and more prominent in the image.



Figure 3: Examples of images disaster types.

##### A.2.1 Disaster types

The purpose of identifying disaster type is to understand the type of disaster events shared in an image. The annotation task involves looking into the image can carefully select one of the following disaster types based on their specific definition. There might be the case that an image shows an effect of a hurricane (destroyed house) and also flood, in such cases the task is to carefully check what is more visible and select label accordingly. Example of images demonstrating different disaster types is shown in Figure 3.

- **Earthquake:** this type of images shows damaged or destroyed buildings, fractured houses, ground ruptures such as railway lines, roads, airport runways, highways, bridges, and tunnels.
- **Fire:** image shows man-made fires or wildfires (forests, grasslands, brush, and deserts), destroyed forests, houses, or infrastructures.
- **Flood:** image shows flooded areas, houses, roads, and other infrastructures.

<sup>11</sup>[https://en.wikipedia.org/wiki/Humanitarian\\_aid](https://en.wikipedia.org/wiki/Humanitarian_aid)

- **Hurricane:** image shows high winds, a storm surge, heavy rains, collapsed electricity polls, grids, and trees.
- **Landslide:** image shows landslide, mudslide, landslip, rockfall, rockslide, earth slip, and land collapse
- **Other disasters:** image shows any other disaster types such as plane crash, bus, car, or train accident, explosion, war, and conflicts.
- **Not disaster:** image shows cartoon, advertisement, or anything that cannot be easily linked to any disaster type.



Figure 4: Example images for **informativeness**.

### A.2.2 Informativeness

The purpose of this task is to determine whether image is useful for *humanitarian aid* purposes as defined below. If the given image is useful for *humanitarian aid*, the annotation task is to select the label “Informative”, otherwise select the label “Not informative” image. Example of images demonstrating informative vs. not-informative is shown in Figure 3.

- **Informative:** if an image is useful for humanitarian aid and shows one or more of the following: cautions, advice, and warnings, injured, dead, or affected people, rescue, volunteering, or donation request or effort, damaged houses, damaged roads, damaged buildings; flooded houses, flooded streets; blocked roads, blocked bridges, blocked pathways; any built structure affected by earthquake, fire, heavy rain, strong winds, gust, etc., disaster area maps.
- **Not informative:** if the image is not useful for humanitarian aid and shows advertising, banners, logos, cartoons, and blurred.



Figure 5: Example images for **humanitarian** categories.

### A.2.3 Humanitarian Categories

Based on the *humanitarian aid* definition above, we define each **humanitarian** information category below.

- **Affected, injured or dead people:** image shows injured, dead, or affected people such as people in shelter facilities, sitting or lying outside, etc.
- **Infrastructure and utility damage:** image shows any built structure affected or damaged by the disaster. This includes damaged houses, roads, buildings; flooded houses, streets, highways; blocked roads, bridges, pathways; collapsed bridges, power lines, communication poles, etc.



- **Not humanitarian:** image is not relevant or useful for humanitarian aid and response such as non-disaster scenes, cartoons, advertisement banners, celebrities, etc.
- **Rescue, volunteering, or donation effort:** image shows any type of rescue, volunteering, or response effort such as people being transported to safe places, people being evacuated from the hazardous area, people receiving medical aid or food, donation of money, blood, or services, etc.

#### A.2.4 Damage severity

The purpose of this task is to identify the severity of damage reported in an image. It can be physical destruction to a build-structure. Our goal is to detect physical damages like broken bridges, collapsed or shattered buildings, destroyed or creaked roads. We define each damage severity category below.

1. **Severe:** Substantial destruction of an infrastructure belongs to the severe damage category. For example, a non-livable or non-usable building, a non-crossable bridge, or a non-drivable road, destroyed, burned crops, forests are all examples of severely damaged infrastructures. For example, if one or more building in the image show substantial loss of amenity or images shows a building that is not safe to use then such image should be labeled as severe damage.
2. **Mild:** Partially destroyed buildings, bridges, houses, roads belong to mild damage category. For example, if image shows a building with damage upto 50%, partial loss of amenity/roof or part of the building can has to be closed down then it should label as mild damage.
3. **Little or none:** Images that show damage-free infrastructure (except for wear and tear due to age or disrepair) belong to the little-or-no-damage category.

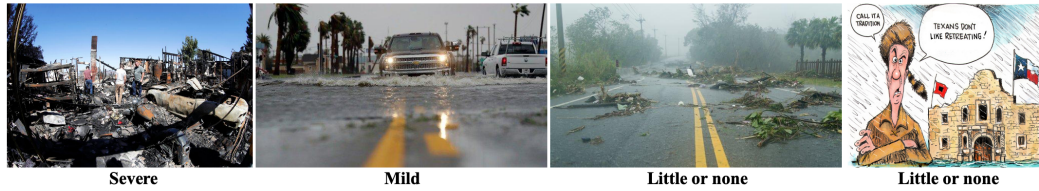


Figure 6: Example images for **damage severity**.

 <p><b>Disaster Type: (required)</b> Please select disaster type below:</p> <p><input type="radio"/> Earthquake <input type="radio"/> Fire <input type="radio"/> Flood <input type="radio"/> Hurricane <input type="radio"/> Landslide <input type="radio"/> Not disaster <input type="radio"/> Other disaster</p> <p><b>Humanitarian: (required)</b> Please select the humanitarian type below:</p> <p><input type="radio"/> Affected, injured, or dead people <input type="radio"/> Infrastructure and utility damage <input type="radio"/> Not humanitarian <input type="radio"/> Rescue volunteering or donation effort</p> <p><b>Annotation: DT, Hum</b></p>	 <p><b>Disaster Type: (required)</b> Please select disaster type below:</p> <p><input type="radio"/> Earthquake <input type="radio"/> Fire <input type="radio"/> Flood <input type="radio"/> Hurricane <input type="radio"/> Landslide <input type="radio"/> Not disaster <input type="radio"/> Other disaster</p> <p><b>Humanitarian: (required)</b> Please select the humanitarian type below:</p> <p><input type="radio"/> Affected, injured, or dead people <input type="radio"/> Infrastructure and utility damage <input type="radio"/> Not humanitarian <input type="radio"/> Rescue volunteering or donation effort</p> <p><b>Damage Severity: (required)</b> Please select damage severity below:</p> <p><input type="radio"/> Little to None <input type="radio"/> Mild <input type="radio"/> Severe</p> <p><b>Annotation: DT, Hum, DS</b></p>
--	--

Figure 7: Example of annotation interfaces on Appen crowdsourcing platform. DT: disaster type, Hum: humanitarian, DS: damage severity.

#### A.3 Annotation Interface

An example of annotation interface is shown in Figure 7. Image on the left shows annotation task is launched to annotate image for disaster type and humanitarian tasks and image on the right shows annotation task is launched for three tasks.

## A.4 Manual Annotation

In our annotation tasks through the Appen platform, more than 3000 annotators participated from more than 50 countries. For the annotation task, we estimated hourly wages and it was 6 to 8 USD per hour on average, which varied depending on the two to three labels annotation per image. We think such pay is reasonable as annotators are from various part of the world where wages varies depending on the location. In total we paid 5,159 USD for the annotation, including Appen charges.

## B The MEDIC dataset

The dataset can be downloaded from <https://crisisnlp.qcri.org/medic/index.html>.

### B.1 Data Format

The dataset format can be found in <https://crisisnlp.qcri.org/medic/index.html>.

### B.2 Terms of use, privacy and License

The MEDIC dataset is published under CC BY-NC-SA 4.0 license, which means everyone can use this dataset for non-commercial research purpose: <https://creativecommons.org/licenses/by-nc/4.0/>.

### B.3 Data maintenance

We provided data download link through <https://crisisnlp.qcri.org/medic/index.html>. We also host on dataverse<sup>12</sup> for wider access. We will maintain the data for a long period of time and make sure dataset is accessible.

### B.4 Benchmark code

The benchmark code is available at: <https://github.com/firojalam/medic/>.

### B.5 Ethics Statement

#### B.5.1 Dataset Collection

The dataset contains images from multiple sources such as Twitter, Google, Bing, Flickr, and Instagram. Twitter developer terms and conditions suggests that one can release 50K tweet objects<sup>13</sup> and here we only provide images not whole JSON objects. The total number of images from Twitter is less than 50,000. Hence, by releasing the data by maintaining such terms and conditions. From Google, Bing, Yahoo and Instagram images are publicly available. In addition, we also maintain licenses and cite prior work based upon we built our work.

#### B.5.2 Potential Negative Societal Impacts

The dataset consists of images collected from social media and different search engines. We have given our best efforts to eliminate any adult content during data preparation and annotation. Hence, we believe that the presence of such content in the dataset might be very unlikely. Our annotation does not contain any identifiable information such as age, gender, or race. However, the images in the dataset have many faces and one might apply facial recognition to identify someone. Intervention with human moderation would be required in order to ensure this does not lead to any misuse. We also would like to highlight that the models' prediction should be used carefully as the purpose of the models' prediction is to facilitate its user, not to make any direct decision. Model designers also need to be careful for any adversarial attack that can lead to creation and spread of any mis/disinformation.

---

<sup>12</sup><https://dataverse.org/>

<sup>13</sup><https://developer.twitter.com/en/developer-terms/agreement-and-policy>

### B.5.3 Biases

The datasets are not representative of a geolocation, user gender, age, race, so should not be used in analyses requiring a representative sample. Instead, the datasets are more suitable to be combined with existing datasets and used for training supervised machine learning models.

We also would like to highlight that some of the annotations are subjective, and we have clearly indicated in the text which of these are. Thus, it is inevitable that there would be biases in our dataset. Note that, we have very clear annotation instructions with examples in order to reduce such biases.

### B.5.4 Intended Use

The dataset can enable an analysis of image content for disaster response, which could be of interest to crisis responders humanitarian response organizations, and policymakers. There are only very few datasets available for multitask learning research. This dataset can significantly help towards this direction. Having a single model for multiple tasks can also foster Green AI.

## C Additional Experimental Details

We have done extensive analysis to understand whether multitask learning setup reduces computational time. In Table 7, we provide such findings for all the models we used in our experiments. From the results, it is clear that multitask learning setup can significantly reduce the computation time both in terms of training and inference.

Given that class distribution can play a significant role in classifier performance, therefore, we wanted to see whether low prevalent classes have any significant impact. In table 8, we report task-wise classification results for both single and multi-task settings in which the model is trained using EfficientNet model. It appears that low prevalent classes have lower performance. However, this is not always the case. For example, the distribution of *Fire* class label is 3.8% in the dataset but the performance is third-best among class labels. Where the distribution of *Other disaster* is 5.1%, however, the F1 is 27.0, which is the lowest performance. With our analysis, we found that this *Other disaster* confused with *Not disaster*.

In Table 9, we report classification results of EfficientNet (b1) multitask learning model, which shows that disaster type class label predictions and their prediction with other tasks. The results suggests that the higher distribution of *Not disaster* does not effect the classification performance much.

In Table 10, we show results obtained using combination of different subset of tasks. We observe that the results remain consistent with other combinations of tasks as well.

## D Multitask Datasets for Disaster Response

In Table 11, we present the datasets containing aligned labels for multitask learning setup. The last row represents, MEDIC, the dataset we propose in this study, which have labels for all tasks and labels for 71,198 images.



Model	Single task					Multitask
	DT	Info	Hum	DS	Sum	
Training time on train set with 49353 images						
ResNet18	16:48:36	18:25:40	15:50:21	16:37:27	2 days, 19:42:04	15:55:33
ResNet50	15:40:12	15:47:36	15:43:45	15:47:25	2 days, 14:58:58	15:45:12
ResNet101	1 day, 15:57:48	23:49:53	1 day, 17:38:58	1 day, 1:36:44	5 days, 11:03:23	1 day, 17:24:21
VGG16	15:31:48	1 day, 10:47:12	1 day, 10:44:30	1 day, 10:39:48	4 days, 23:43:18	2 days, 10:08:56
DenseNet (121)	1 day, 1:52:00	17:08:10	17:03:30	1 day, 2:05:41	3 days, 14:09:21	17:50:09
SqueezeNet	15:16:46	15:50:48	15:14:34	15:39:38	2 days, 14:01:46	15:22:03
MobileNet (v2)	15:53:40	15:23:22	15:01:28	15:41:26	2 days, 13:59:56	16:01:59
EfficientNet (b1)	23:21:11	17:10:36	17:08:05	23:41:41	3 days, 9:21:33	23:49:12
Inference time on test set with 15688 images						
ResNet18	0:04:39	0:02:17	0:01:59	0:02:03	0:10:58	0:01:56
ResNet50	0:02:06	0:02:01	0:01:54	0:01:58	0:07:59	0:01:54
ResNet101	0:01:55	0:01:55	0:02:01	0:02:33	0:08:24	0:02:00
VGG16	0:04:45	0:01:58	0:01:56	0:01:57	0:10:36	0:02:18
DenseNet (121)	0:01:59	0:01:58	0:01:53	0:01:56	0:07:46	0:01:57
SqueezeNet	0:01:55	0:02:10	0:05:15	0:02:08	0:11:28	0:01:55
MobileNet (v2)	0:02:08	0:05:38	0:01:54	0:01:52	0:11:32	0:01:59
EfficientNet (b1)	0:01:53	0:02:00	0:01:59	0:01:58	0:07:50	0:01:55

Table 7: Training and inference time in single vs. multitask settings with a batch size of 32. Time is in day, hour:minute:second format.

Class label	P	R	F1	P	R	F1
	Single-task			Multi-task		
Disaster types						
Earthquake	69.6	79.6	74.3	68.5	79.9	73.8
Fire	76.1	83.9	79.8	73.5	85.3	79.0
Flood	78.3	80.8	79.5	79.1	79.3	79.2
Hurricane	65.3	63.9	64.6	66.4	63.8	65.1
Landslide	59.1	77.4	67.0	60.8	74.0	66.8
Not disaster	88.0	92.0	89.9	87.7	92.4	90.0
Other disaster	63.7	19.8	30.2	65.4	17.0	27.0
Informative						
Informative	89.2	82.4	85.7	88.6	85.6	87.0
Not-informative	80.5	87.9	84.0	83.2	86.6	84.8
Humanitarian						
Affected, injured, or dead people	44.5	24.2	31.4	48.0	29.1	36.2
Infrastructure and utility damage	78.3	75.3	76.8	77.5	84.2	80.7
Not humanitarian	82.5	89.7	85.9	86.3	88.8	87.5
Rescue volunteering or donation effort	49.0	32.5	39.1	55.5	29.1	38.2
Damage severity						
Little or none	89.0	91.5	90.2	89.8	91.6	90.7
Mild	42.0	22.4	29.2	45.5	21.4	29.1
Severe	71.9	81.1	76.2	71.1	83.3	76.7

Table 8: Class-wise results for both single and multi-task settings using EfficientNet (b1) model.

Disaster Type	Informativeness		Humanitarian				Damage Severity		
Class Name	Informative	Not Informative	Affected, injured, or dead people	Infrastructure and utility damage	Not humanitarian	Rescue volunteering or donation effort	Little or none	Mild	Severe
Earthquake (1953)	1952	1	98	1783	24	48	17	129	1807
Fire (799)	787	12	2	668	109	20	114	23	662
Flood (1304)	1293	11	51	1090	116	47	190	109	1005
Hurricane (1354)	1350	4	1	1081	269	3	411	226	717
Landslide (398)	396	2	0	340	58	0	70	27	301
Other Disaster (315)	310	5	84	205	15	11	88	69	158
Not Disaster (9565)	2224	7341	137	293	8772	363	9196	222	147

Table 9: Results with EfficientNet (b1) multitask learning model demonstrating disaster type labels with other three tasks.

Tasks	DT	Info	Hum	DS	Tasks	DT	Info	Hum	DS
DT-Info-Hum-DS	79.3	86.0	80.4	80.3	DT-DS	78.6			79.4
DT-Info-Hum	76.3	85.0	78.0		Info-Hum-DS		85.6	80.2	79.8
DT-Info-DS	79.2	86.0		79.9	Info-Hum		85.3	78.8	
DT-Info	79.1	85.7			Info-DS		85.6		79.6
DT-Hum	79.5		80.5		Hum-DS			80.1	79.9

Table 10: Results (F1) with different combination of tasks using EfficientNet (b1). DT: Disaster type, Info: Informativeness, Hum: Humanitarian, DS: Damage severity.

	Disaster types	Informativeness	Humanitarian	Damage severity	Total
CrisisMMD [4]		✓	✓	✓	3,533
CrisisBench [6]	✓	✓	✓	✓	5,558
MEDIC	✓	✓	✓	✓	71,198

Table 11: Multitask learning datasets for disaster image classification tasks.