

# DiRecNetV2: A Transformer-Enhanced Network for Aerial Disaster Recognition

Demetris Shianios<sup>[0009–0005–8266–0727]</sup>, Panayiotis S.  
Kolios<sup>[0000–0003–3981–993X]</sup>, and Christos Kyrkou<sup>\*[0000–0002–7926–7642]</sup>

KIOS Research and Innovation Center of Excellence, University of Cyprus,  
Nicosia, Cyprus

<https://www.kios.ucy.ac.cy/>

{shianios.demetris,kolios.panayiotis,kyrkou.christos}@ucy.ac.cy

**Abstract.** The integration of Unmanned Aerial Vehicles (UAVs) with artificial intelligence (AI) models for aerial imagery processing in disaster assessment, necessitates models that demonstrate exceptional accuracy, computational efficiency, and real-time processing capabilities. Traditionally Convolutional Neural Networks (CNNs), demonstrate efficiency in local feature extraction but are limited by their potential for global context interpretation. On the other hand, Vision Transformers (ViTs) show promise for improved global context interpretation through the use of attention mechanisms, although they still remain underinvestigated in UAV-based disaster response applications. Bridging this research gap, we introduce DiRecNetV2, an improved hybrid model that utilizes convolutional and transformer layers. It merges the inductive biases of CNNs for robust feature extraction with the global context understanding of Transformers, maintaining a low computational load ideal for UAV applications. Additionally, we introduce a new, compact multi-label dataset of disasters, to set an initial benchmark for future research, exploring how models trained on single-label data perform in a multi-label test set. The study assesses lightweight CNNs and ViTs on the AIDERSv2 dataset, based on the frames per second (FPS) for efficiency and the weighted F1 scores for classification performance. DiRecNetV2 not only achieves a weighted F1 score of 0.964 on a single-label test set but also demonstrates adaptability, with a score of 0.614 on a complex multi-label test set, while functioning at 176.13 FPS on the Nvidia Orin Jetson device.

**Keywords:** Natural Disaster Recognition · Image Classification · UAV (Unmanned Aerial Vehicle) · Convolutional Neural Networks · Vision Transformers · Multi-Label Classification

## 1 Introduction

The increasing impact of natural disasters such as floods, hurricanes, earthquakes, and wildfires on global regions demonstrates the urgent need for greater

awareness and emergency preparedness. Data from Our World in Data <sup>1</sup> reveal that, on average, natural disasters claim the lives of approximately 45,000 individuals annually, accounting for approximately 0.1% of deaths worldwide. Regrettably, 2023 encountered several catastrophic events, which comprised the earthquakes in Turkey-Syria that claimed over 33,000 lives and injured thousands more [61], the extensive floods in South Africa and South America [9], and the devastating wildfires that swept across Western Canada that have burned more than 478,000 hectares of land as of the Canadian Wildland Fire Information System <sup>2</sup>. Moreover, the latest global climate science assessment indicates an increasing frequency of simultaneous multiple disasters, presenting complex challenges in disaster management [40]. For instance, typhoons can simultaneously cause structural collapses and floods due to their strong winds and rainfall. Likewise, thunderstorms can spark both fires and floods; lightning ignites fires while heavy rains lead to floods. Another complex circumstance is that of earthquakes, which might lead to the collapse of buildings followed by fires in infrastructure. For disaster mitigation and resource allocation to be effective, these multi-label instances must be accurately identified and analyzed.

Unmanned aerial vehicle (UAV) technology such as drones has advanced to the point where it can be an effective tool for minimizing the consequences of natural disasters, particularly in remote regions [31]. Through the integration of advanced AI systems and high-resolution cameras, UAVS offer rapid and secure data collection, providing images and videos crucial for mapping disasters, informing authorities promptly, and facilitating emergency response and rescue operations in real-time [29]. With the ability to be controlled remotely, they are capable of flying over large, frequently inaccessible areas and turning into intelligent drones that improve situational awareness and operational effectiveness in emergency situations. Moreover, Deep Learning can empower drones to classify disasters, facilitating rapid and precise recognition of impacted zones, evaluation of damage extent, and prioritization of response efforts.

Deep learning algorithms such as Convolutional Neural Networks (CNNs) have been pivotal in image processing and computer vision, gaining prominence since their inception in 1998 [32]. More recently, Vision Transformers (ViTs) [11], have revolutionised image recognition, offering a different approach to interpreting visual data. However, the research community is now exploring the advancements of lighter models in applications like emergency response, where quick processing is essential. These streamlined architectures strike a balance between accuracy and the requirement for quick and effective computation, making them crucial for real-time analysis.

In this work, we have created a novel hybrid model building on the advantages of CNNs and ViTs, namely the Disaster Recognition Network (*DiRecNetV2*) model. This model represents an advancement from our previous endeavors, in which we introduced DiRecNet, a benchmark dataset (AIDERSv2), and we explored visual explainability techniques like Grad-CAM [56]. The proposed Di-

<sup>1</sup> <https://ourworldindata.org/>

<sup>2</sup> <https://cwfis.cfs.nrcan.gc.ca/>

RecNetV2 is designed to address the challenges of accuracy and complexity found in current lightweight models, extending the scope of our ongoing research effort. The proposed model combines the broad contextual strengths of ViTs which help to capture long-range dependencies within an image, with the local inductive biases of CNNs which help to learn hierarchical features. The combination of simple design choices, such as depthwise convolution and a reduced number of heads and encoder blocks in the Vision Transformer, leads to a highly efficient model. This model is specifically designed for the unique needs of UAV-based disaster management. Furthermore, we thoroughly benchmark lightweight CNNs and ViTs on the AIDERSv2 dataset<sup>3</sup>. We evaluated models such as ConvNeXt Tiny [34], EfficientNet-B0[60], MnasNet [59], MobileNetV2 [52], MobileNetV3 Small [21], ShuffleNetV2 [38] and SqueezeNet [22] for the CNNs family, while for ViTs, we examine; Convit Tiny [12], GCvit XXtiny [19], MobileViT [42], MobileViTV2 [41] and Vit-Tiny [11,58].

In addition, we propose a multi-label dataset of 300 images, evenly distributed among fire-earthquake, fire-flood, and flood-earthquake disaster combinations. This dataset functions as a test set for deep learning models that are initially trained on more extensive single-label datasets. Our assessments span both single-label and multi-label classification performance, along with efficiency in terms of execution time using the FPS metric (frames processed in a second) over a Jetson Orin device. In our study, DiRecNetV2 has demonstrated outstanding performance on the AIDERSv2 dataset. Additionally, our proposed model has shown robustness in handling complex scenarios, such as identifying multiple disasters in a single image, on a specially constructed multi-label dataset. These results underscore its effectiveness in both single-label and multi-label disaster recognition tasks.

## 2 Related Work

In this section, we delve into related works on disaster classification, including studies exploring the use of lightweight CNNs [31], and the emerging application of Vision Transformers (ViTs) in disaster management scenarios. Several works examine disaster recognition for single classes, such as; earthquakes [37,8,55,25,36,66], floods [10,17,45,44,48,47,53,39], and wildfires [54,64,26,14,20,20,24]. Currently, significant research efforts in natural disaster detection have leveraged UAVs, satellites, and social media as primary sources of data.[15,1,3,67,4,6].

Multi-label learning addresses the issue of one example being simultaneously linked to multiple labels [70]. In the context of a natural disaster, a multi-label problem can be an example of a satellite image that needs to be classified for both flood and wildfire occurrences simultaneously, due to overlapping disaster events in the same geographic region. Another example is a situation where there is a fire and the infrastructure collapses as a result of a bomb explosion. While there are some works regarding multi-label text classification for disaster, [2,63,13,5],

---

<sup>3</sup> link to be provided upon publication

there are few based on computer vision for disaster recognition [57,7]. Despite the importance of identifying multiple disasters at once for effective emergency response, there is still a significant research gap concerning the application of Deep Learning models to this challenge.

In recent years, Vision Transformers have been introduced [11], a transformative method for image classification by adapting the transformer architecture, previously used in natural language processing[62]. In the field of disaster recognition, only a few studies have been conducted using ViT architecture, focusing on specific areas such as fire detection [27], wildfire segmentation [18], flood segmentation [49], and earthquake magnitude estimation [50]. There is a dearth of research on the use of Vision Transformers for disaster classification, with a significant gap in understanding their performance in such contexts.

Furthermore, there is a growing need to incorporate AI models into drones and embedded systems for disaster detection/classification purposes [30]. For these models to process and analyse streaming and live video footage efficiently, they must be lightweight, ensuring quick inference, low memory usage, and low power consumption. As a result, researchers are investigating more compact architectures of Deep Learning models in disaster management [68,69,16,46,65,33,43,51]. Although most of these studies are CNN-based, there has been relatively little investigation into lightweight models for ViTs in the context of disaster recognition tasks.

Overall, many disaster detection approaches focus on single-disaster types, with an ongoing transition towards multi-class detection. However, these models often turn out to be complex for UAV integration, highlighting the necessity for models that align with the limitations of UAV systems. Furthermore, there is a clear gap in models that can handle multi-label disaster scenarios, presenting a vital area for further research and development. Moreover, there is a limited understanding of the effectiveness of Vision Transformers in classifying disasters from aerial images. Our study presents a hybrid CNN-ViT model that integrates lightweight architecture appropriate for embedded systems to precisely fill these gaps. Furthermore, we have developed a multi-label dataset and conducted comprehensive benchmarking to evaluate the performance of lightweight CNNs and Vision Transformers in disaster classification scenarios for both single and multi-label instances.

### 3 Multi-Label Dataset and Hybrid CNN-ViT Model

The following sections detail the methodological procedure for deriving our two main contributions, the multi-label dataset, and the hybrid CNN-ViT model.

#### 3.1 Multi-Disaster Dataset

The development of a multi-label disaster classification system is driven by the essential need to provide accurate, real-time information in response to the occurrence of simultaneous disasters. This is particularly important for facilitating the



**Fig. 1.** Sample of images in the database depicting various multi-label disaster instances.

effective and appropriate distribution of resources and preparedness in multiple disaster situations at the same time. We aim to develop a multi-label dataset that will enable us to explore the capabilities of deep learning models to simultaneously identify multiple disaster events. To curate a dataset for multi-label disaster scenarios, we scoured online platforms such as Google, Bing, and YouTube, using search queries like "flood and fire," "floods and earthquakes/collapsed buildings," and "fires and earthquakes/collapsed buildings." These cases are uncommon, making data collection challenging and time-consuming. Although such scenarios are rare, we collected 300 images, ensuring an equal distribution of 100 images for each disaster combination.

We carefully examined the chosen images in the ensuing review stage to search for any ambiguities or errors. This thorough validation procedure was essential to maintaining the integrity of the dataset, guaranteeing a low degree of bias, and reducing the possibility of mistakes that would hinder the precision and dependability of our disaster classification model. Two researchers carried out the process of collection and verification. The dataset is currently only used as a test set for models that have already been trained on larger datasets. As the dataset expands, it will become a crucial resource for training and testing models on multi-label instances, thereby improving their accuracy and performance. This work represents a first step, establishing baseline results for models trained on larger single-label datasets and then assessed on our multi-label dataset, providing a foundation for future research endeavors.

### 3.2 Single and Multi-Label Classification

Single-label classification tasks are commonly encountered in numerous Computer Vision applications, with the primary objective being to identify a single result from a range of potential outcomes. This method typically utilizes the softmax activation function (Eq.1) within the final dense layer that outputs the class probabilities. This function guarantees that the total probabilities for all classes sum up to 1, while also ensuring that they are dependent on each other, making it well-suited for single-label classification scenarios.

$$\text{softmax}(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad (1)$$

Where  $z$  is the input vector containing logits (raw scores) for each  $k$  class.

In single-label classification training, the commonly used loss function is categorical cross-entropy (Eq.2).

$$H(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i) \quad (2)$$

The categorical cross-entropy loss function compares the predicted probability  $\hat{y}$  distribution with the true distribution  $y$  making it suitable for single-label cases where each instance has only one label.

On the other hand, when a number of outcomes are possible for a given instance, multi-label classification becomes essential. In that case, the sigmoid activation function (Eq. 3) is employed that returns the independent probabilities for each class. Similar to softmax, these probabilities, range between 0 and 1, but are interdependent for each class and might not sum to 1, allowing for simultaneous classification across multiple classes.

$$\sigma(z)_j = \frac{1}{1 + e^{-z_j}} \quad (3)$$

Furthermore, the process of multi-label classification involves converting the labels into one-hot encoded vectors, since they provide a distinct and explicit representation of each label's presence or absence for a particular instance. This encoding format is the key to effectively managing multiple labels per instance, as it enables the binary cross-entropy loss function (Eq. 4) to evaluate the prediction accuracy for each individual label independently.

$$H(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \quad (4)$$

### 3.3 Lightweight Deep Learning Models

The development of lightweight deep learning models is increasingly important, particularly for drone-based applications. Given that these models are less

resource-intensive, they can operate on devices like drones that have limited processing power. Consequently, intelligence drones, outfitted with lightweight deep learning models, can conduct real-time analysis and make decisions mid-flight, a critical capability for emergency response and disaster monitoring tasks.

**Convolutional Neural Networks:** Convolutional Neural Networks (CNNs), which were first introduced in 1998 [32] and are well-known for their efficiency in pattern recognition and feature extraction from visual data, have since evolved to become an essential tool in image processing and computer vision applications. However, in an effort to preserve performance while lowering computational load, the trend towards lighter CNN architectures is being driven by the need for efficiency in embedded systems. For quick processing, these simplified models are constructed to maximize the trade-off between accuracy and computational efficiency. They employ strategies such as reduced layer depth, minimized parameters, and compact channel dimensions to simplify their architecture. Techniques like depth-wise separable convolutions are also utilized to cut computational demands by decoupling spatial and channel filtering. Moreover, recent trends in research have moved towards channel attention mechanisms, focusing on salient features to further refine the models' complexity for enhanced performance. In our study, we conducted an analysis of various lightweight CNN models such as; ConvNeXt Tiny [34], EfficientNet-B0[60], MnasNet [59], MobileNetV2 [52], MobileNetV3 Small [21], ShuffleNetV2 [38] and SqueezeNet [22]. We decided to include ConvNeXt Tiny in our evaluations even though it might not be generally considered lightweight in comparison to other architectures, however is the lightest model variant within the newly introduced ConvNeXt family.

**Vision Transformers:** Vision Transformers (ViTs) have recently emerged in computer vision applications as potential alternatives to CNNs. Initially, vision transformer models [11], inspired by their counterparts in language processing [62], were computationally heavy with more than 85 million parameters. The research community is now exploring ways to construct lightweight versions of ViTs. The goal of these efforts is to reduce the computing footprint of ViTs while maintaining their advantages, such as their adaptability and global receptive field. By developing lightweight ViTs, researchers hope to bring the advantages of transformer-based models to more resource-limited applications. To reduce the computational demands of self-attention, effective attention techniques like localized or sparse attention are integrated. By factorizing transformer blocks and adding convolutional layers, the architecture is further optimized and a hybrid model that effectively processes spatial information is created. The investigated lightweight ViTs are: Convit Tiny [12], GCvit XXtiny [19], MobileViT-(s,xx,xss) [42], MobileViTV2-(050,100) [41] and Vit-Tiny [11].

### 3.4 Transformer Enhanced Convolutional Network Architecture

This study integrates a custom CNN with a Vision Transformer (ViT) to present a novel architecture for UAV emergency response. To achieve high performance

with minimal computational demands and satisfy the requirements of UAV applications, this hybrid approach combines the effective feature extraction of CNN with the global power capability of ViT. To this end, we end up with the Disaster Recognition Network V2 *DiRecNetV2* model, which is an improvement over our prior work [56], now integrating transformer capabilities.

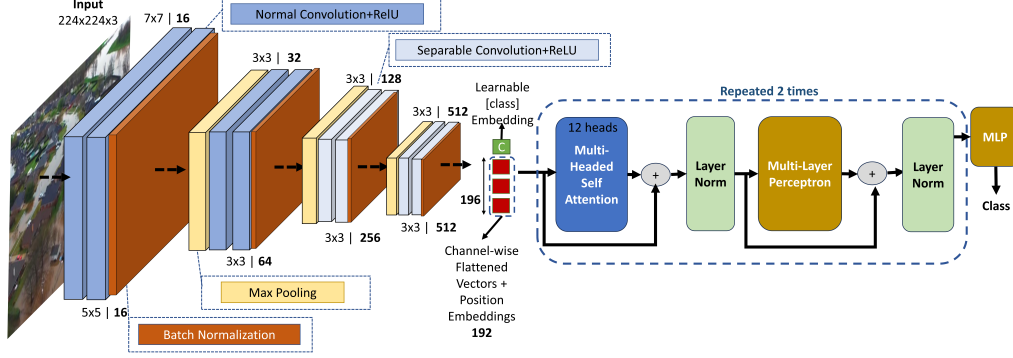
The initial model from *DiRecNet* [56], incorporates four key blocks, enabling the learning of hierarchical features while maintaining feature map resolution effectively. In the new version the output of the final block is passed through a transformer encoder block to enhance the model’s ability to capture and integrate complex patterns and relationships in the data. This block utilizes the self-attention mechanism inherent in transformers, enabling the model to consider the entire context of the input, leading to a more comprehensive understanding of the features extracted by the *DiRecNet* backbone. Table 1 presents a detailed analysis of the model blocks, including the number of parameters, and Fig. 2 illustrates the model’s architecture.

**DiRecNet Feature Extractor:** The model begins with the *DiRecNet* feature extractor with two consecutive convolutional layers, one with a  $7 \times 7$  kernel and 16 filters, and the other with a  $5 \times 5$  kernel and 16 filters, following modern network trends that apply larger kernels [35]. Batch normalization and max-pooling of stride  $2 \times 2$  follow. The subsequent block has two  $3 \times 3$  convolutional layers with 32 and 64 filters. After batch normalization, another max-pooling layer is applied. The third block utilizes two separable convolutions with 128 and 256 filters, followed by batch normalization and max-pooling. The final block incorporates two identical separable convolutions with 512 filters, followed by batch normalization and max-pooling. In contrast to the original *DiRecNet* we add a batch normalization to the last block, a modification absent in the original design. The feature map is subsequently projected to an embedding dimension of 192, as specified by the transformer hyperparameter. This results in a  $192 \times 14 \times 14$  feature map, which is then forwarded to the transformer encoder block.

**Transformer Encoder Block:** The feature map is flattened into 196 element vectors in a depth-wise fashion from  $(14 \times 14)$  patches each of 192 elements, creating a denser representation of the image’s features. A classification token ‘CLS’ is added, resulting in 197 patches, each of 192 dimensions. The patches are passed through a Transformer encoder block with a 12-way multi-head attention mechanism, offering simultaneous focus on various image areas. This approach combines local and global context for comprehensive image analysis. LayerNorm normalization is applied before multi-head attention and before passing through an MLP (Multi-Layer Perceptron) of the same dimensionality of 192. After several experiments, we settled on using two transformer encoder blocks as it resulted in slightly better accuracy. Lastly, the classification head incorporates a LayerNorm as in the standard ViT architecture to normalize the features, and Dropout for regularization. It is followed by a linear layer for class mapping



and uses a Softmax in case of single-label classification or a Sigmoid layer for multi-label classification.



**Fig. 2.** The DiRecNetV2 model architecture, showcasing how features extracted by CNN blocks are fed into the encoder blocks of ViTs of the Vision Transformer, significantly enhancing the model’s capabilities.

### 3.5 Training Process

**Datasets:** In our research, we initially divided the data, sourced from the AIDERSv2 dataset as described in [56], into training, validation, and test sets, with proportions of 80%, 10%, and 10% respectively. The distribution of these sets is detailed in Table 2. For the multi-label evaluation, it is important to highlight that the training was performed on the same dataset as used for the single-label tasks, with the only difference being in the final activation function. However, the test set was the Multi-label dataset discussed in Section 3.1. We trained the models on data with single disasters for the purpose to assess their performance on new images with more than one disaster. Future work could involve fine-tuning on the multi-label dataset; once we acquire a more extensive multi-disaster dataset. We can split it into training, validation, and test sets to investigate how model performance varies when fine-tuned on a multi-label dataset during training.

**Data Pre-Processing:** The images were scaled to  $224 \times 224 \times 3$  and standardized for DiRecNetV2; therefore, to change the distribution to have a mean of zero and a standard deviation of one. Random augmentations were applied to expand the diversity of the dataset and combat overfitting. Specifically, we applied rotation, zoom, horizontal shift, vertical shift, horizontal flip, and shear.

Layer	Input Shape	Output Shape	Param #
<b>DiRecNetV2</b>	[32, 3, 224, 224]	[32, 4]	38,016
<b>DiRecNet Feature Extractor</b>	[32, 3, 224, 224]	[32, 196, 192]	–
Conv2d	[32, 3, 224, 224]	[32, 16, 224, 224]	2,368
Conv2d	[32, 16, 224, 224]	[32, 16, 224, 224]	6,416
BatchNorm2d	[32, 16, 224, 224]	[32, 16, 224, 224]	32
MaxPool2d	[32, 16, 224, 224]	[32, 16, 112, 112]	–
Conv2d	[32, 16, 112, 112]	[32, 32, 112, 112]	4,640
Conv2d	[32, 32, 112, 112]	[32, 64, 112, 112]	18,496
BatchNorm2d	[32, 64, 112, 112]	[32, 64, 112, 112]	128
MaxPool2d	[32, 64, 112, 112]	[32, 64, 56, 56]	–
DepthwiseConv2d	[32, 64, 56, 56]	[32, 64, 56, 56]	640
PointwiseConv2d	[32, 64, 56, 56]	[32, 128, 56, 56]	8,320
DepthwiseConv2d	[32, 128, 56, 56]	[32, 128, 56, 56]	1,280
PointwiseConv2d	[32, 128, 56, 56]	[32, 256, 56, 56]	33,024
BatchNorm2d	[32, 256, 56, 56]	[32, 256, 56, 56]	512
MaxPool2d	[32, 256, 56, 56]	[32, 256, 28, 28]	–
DepthwiseConv2d	[32, 256, 28, 28]	[32, 256, 28, 28]	2,560
PointwiseConv2d	[32, 256, 28, 28]	[32, 512, 28, 28]	131,584
DepthwiseConv2d	[32, 512, 28, 28]	[32, 512, 28, 28]	5,120
PointwiseConv2d	[32, 512, 28, 28]	[32, 192, 28, 28]	98,496
BatchNorm2d	[32, 192, 28, 28]	[32, 192, 28, 28]	384
MaxPool2d	[32, 192, 28, 28]	[32, 192, 14, 14]	–
Flatten	[32, 192, 14, 14]	[32, 197, 192]	–
Dropout	[32, 197, 192]	[32, 197, 192]	–
<b>Transformer Encoder Blocks</b>	[32, 197, 192]	[32, 197, 192]	–
TransformerEncoderBlock (1)	[32, 197, 192]	[32, 197, 192]	–
MultiheadSelfAttentionBlock	[32, 197, 192]	[32, 197, 192]	148,608
MLPBlock	[32, 197, 192]	[32, 197, 192]	74,496
TransformerEncoderBlock (2)	[32, 197, 192]	[32, 197, 192]	–
MultiheadSelfAttentionBlock	[32, 197, 192]	[32, 197, 192]	148,608
MLPBlock	[32, 197, 192]	[32, 197, 192]	74,496
<b>Classifier Head</b>	[32, 197, 192]	[32, 4]	–
LayerNorm	[32, 192]	[32, 192]	384
Dropout	[32, 192]	[32, 192]	–
Linear	[32, 192]	[32, 4]	772
Softmax/Sigmoid	[32, 4]	[32, 4]	–

**Table 1.** The DiRecNetV2 model’s structured layout, showcases the evolution of feature spaces across various blocks, complete with input and output dimensions utilizing a batch size of 32. In addition, we can explore, the number of parameters at each stage, culminating in a streamlined model architecture with a total of 799,380 parameters.

**Transfer Learning:** For training the different baseline networks we adopted a transfer learning approach for both lightweight Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs). This method significantly accelerated the training process by leveraging the knowledge embedded in the weights ob-

	Earthquakes	Floods	Wildfire/Fire	Normal	Total
<b>Train</b>	1927	4063	3509	3900	13399
<b>Validation</b>	239	505	439	487	1670
<b>Test</b>	239	502	436	477	1654
<b>Total</b>	2405	5070	4384	4864	<b>16723</b>

**Table 2.** Proportion of images in each class within the train, validation, and test set.

tained from training these models on the extensive ImageNet dataset [28]. The use of transfer learning offers a substantial advantage as it uses pretrained models with a rich understanding of diverse features from larger dataset, thereby enhancing learning efficiency and performance on diverse tasks. Our sole modification involved adapting the classifier head for each pretrained lightweight model. Specifically, to mitigate overfitting, we incorporated a dropout layer with a rate of 0.5, followed by a dense layer tailored to handle the number of classes in our case to four.

**Training Regime:** For our custom-designed model that was not pretrained on ImageNet, we conducted training of 300 epochs from scratch to refine its learning capabilities. In contrast, for the pretrained models, we limited the training to 40 epochs, considering their prior knowledge acquired from the larger dataset. This training strategy was uniformly applied across both single-label and multi-label classification experiments. In all cases, we utilized the Adam optimizer with a learning rate of  $1e-4$  and batch size was set to 32, while for the model selection we chose the iteration of the model where the validation accuracy was at its highest akin to early stopping.

**Configurations:** The experiments were carried out on the Linux operating system using the Tesla V100 Graphics Processor Unit, with 64GB RAM and CUDA version 10.2. We use PyTorch <sup>4</sup> 1.12.1 as the deep learning framework along with Python <sup>5</sup> version 3.8.0. Additionally, to evaluate the models’ frames per second (FPS), we deploy them on NVIDIA’s jetson Orin device.

## 4 Evaluation and Results

### 4.1 Evaluation Metrics

To evaluate the performance of the models, we investigated two key performance indicators since both accuracy and speed are crucial to detect natural disasters in real time. These are the weighted F1 score (Eq.5) and frames per second (FPS). We use Weighted F1-Score for both single-label and multi-label evaluation, because it accounts for the relative importance of each label by weighting

<sup>4</sup> <https://pytorch.org/>

<sup>5</sup> <http://www.python.org>

the F1 Score of each label according to its prevalence in the dataset. The FPS values, acquired from tests conducted on the NVIDIA Jetson Orin device.

$$\text{Weighted F1 Score} = \sum_{i=1}^N w_i \times \text{F1 Score}_i \quad (5)$$

Where,

$$w_i = \frac{\text{No. of samples in class } i}{\text{Total number of samples}}$$

We use the scoring scheme of our previous research [56] to determine the best trade-off between speed and accuracy in model performance, as defined in the Eq. 6. We can determine the model that provides the best fit for a particular application scenario by varying the parameter  $\lambda$ . In the current study, we have given the lambda a value of 0.3 when we prioritize speed and 0.7 when we prioritize accuracy. Moreover, for a broader assessment, we compare the models against a revised scoring equation from [23], as presented in Eq. 7. To match the scale of the evaluation metrics, we have set the normalization constant  $C$  at  $1e27$  in this instance.

$$\text{Score1} = \lambda \times F1_{\text{norm}} + (1 - \lambda) \times FPS_{\text{norm}} \quad (6)$$

$$\text{Score2} = \frac{2^{F1} \times FPS}{C} \quad (7)$$

However, prior to this, we use the formula in Eq. 8 to normalize the values of FPS and Weighted F1 across all models since their ranges differ. Values in  $x$  are squeezed into the range  $[a, b]$ , where  $a$  was set to 0.1 and  $b$  at 1, making the variables comparable to one another.

$$x_{\text{norm}} = (b - a) \frac{x - \min(x)}{\max(x) - \min(x)} + a \quad (8)$$

## 4.2 Computational Efficiency Evaluation

In our analysis comparing deep learning architectures for real-time disaster recognition, the introduced DiRecNetV2 model stands out as it has the smallest number of parameters at just 0.8 million. With the exception of ConvNeXt, which has 27.82 million parameters, the other examined models have parameters under 12 million. In terms of computational demand, measured in giga FLOPs, all models are below 2, except for ConvNeXt which has 4.46. When considering the size of the models, we observe a wider distribution; DiRecNetV2, at 3.20 MB, has the second-smallest footprint, closely followed by SqueezeNet’s 2.90 MB. DiRecNetV2’s low parameter count makes it particularly well-suited for embedded devices where resource efficiency is crucial.

Moreover, the examination of model’s execution reveals that DiRecNetV2’s processing speed is second with 176.13 FPS only to SqueezeNet, which leads with

Model Name	GFLOPs	No. of	Model Size	FPS
		Params (M)	(MB)	(1/s)
ConvT Tiny (2021) [12]	1.08	5.52	22.07	54.34
GCViT XXtiny (2023) [19]	1.94	11.48	45.93	31.06
MobileViT s (2021) [42]	1.42	4.94	19.76	30.71
MobileViT xs (2021)[42]	0.71	1.93	7.74	26.07
MobileViT xxs (2021) [42]	0.26	0.95	3.81	33.05
MobileViTV2 050 (2022) [41]	0.36	1.11	4.46	33.42
MobileViTV2 0100 (2021) [41]	1.41	4.39	17.56	37.65
Vit Tiny [11,58]	1.08	5.53	22.1	74.88
ConvNext Tiny (2022)[34]	4.46	27.82	111.29	89.09
EfficientNet-B0 (2019) [60]	0.41	4.01	16.05	55.74
MnasNet (2019) [59]	0.34	3.11	12.43	90.98
MobileNetV2 (2018) [52]	0.33	2.23	8.92	87.23
MobileNetV3 Small (2019) [21]	0.06	0.93	3.72	89.72
ShuffleNetV2 (2017) [38]	0.15	1.26	5.03	75.73
SqueezeNet (2016) [22]	0.26	0.73	2.90	183.08
DiRecNetV2 (Proposed)	1.09	0.80	3.20	176.13

**Table 3.** Comparison of lightweight Vision Transformers (ViTs) and lightweight CNN models, focusing on their computational complexity by giga Flops (floating-point operations per second), the number of parameters, and model size. Each model has been customized with a classifier head, including a dropout rate (0.5) between the feature extractor and the fully connected layer. This customization is aligned with the target class count, which in this case is four.

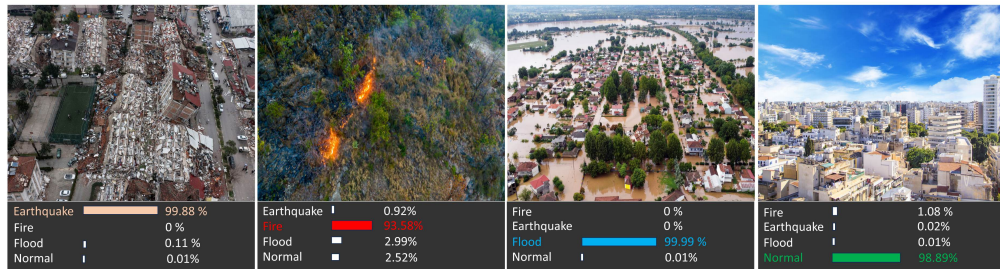
an FPS of 183.08 but at a higher parameter count and larger size. This approach is especially advantageous for high-speed drones that are tasked with wide-area surveillance missions, where there is a requirement to analyze an extensive number of frames. In such scenarios, the ability to rapidly process and interpret vast amounts of visual data is crucial. A detail analysis of the computational efficiency of the models is presented in Table 3.

### 4.3 Disaster Classification Evaluation

**Single-Label Classification Performance:** In our comparative analysis of lightweight CNNs and ViTs we observed that there was a range in performance based on the weighted F1 score. Scores for the lightweight CNNs range, from 0.845 of the SqueezeNet to the highest performance at 0.940 of ConvNeXt Tiny. Lightweight ViTs exhibited a performance spectrum, with scores spanning from 0.837 for MobileViT xs to 0.932 for GCViT XXtiny. This analysis suggests that lightweight Vision Transformers have the potential to compete with traditional CNNs in the task of disaster identification from aerial images. Notably, within both families, the only two models that achieved a score above 0.90 were the more complex ones (ConvNeXt, GCViT XXtiny), characterized by a larger number of parameters and greater model size. This pattern suggests that although efficient

Model Name	Weighted F1	Score1			Score 2
		Balance	Prioritize WF1	Prioritize FPS	
ConvNeXt Tiny	0.940	0.646	0.720	0.572	1764.609
EfficientNet-B0	0.862	0.281	0.285	0.277	4.988
MnasNet	0.897	0.502	0.513	0.490	89.600
MobileNetV2	0.893	0.479	0.490	0.467	67.401
MobileNetV3 Small	0.868	0.398	0.372	0.425	12.003
ShuffleNetV2	0.852	0.304	0.272	0.336	3.436
SqueezeNet	0.845	0.586	0.421	0.752	4.974
Convit Tiny	0.871	0.307	0.324	0.289	8.827
GCVit XXtiny	0.932	0.452	0.582	0.323	355.801
MobileVitV2 050	0.834	0.113	0.108	0.119	0.403
MobileVitV2 100	0.875	0.240	0.297	0.184	5.705
MobileViT s	0.855	0.190	0.210	0.170	1.759
MobileViT xs	0.837	0.131	0.126	0.135	0.532
MobileViT xxs	0.859	0.219	0.241	0.198	2.775
Vit Tiny	0.873	0.375	0.373	0.377	14.666
DiRecNetV2	<b>0.964</b>	<b>0.980</b>	<b>0.988</b>	<b>0.972</b>	<b>18982.892</b>

**Table 4.** Evaluating model efficiency in contexts where accuracy is crucial against situations demanding high frames per second (FPS) highlights different operational focuses on emergency management and live surveillance.



**Fig. 3.** The examples demonstrate DiRecNetV2’s proficiency in identifying diverse disaster situations. Using a subset of four test set images, these images show the model’s robust classification accuracy for earthquakes, fires, floods, and normal cases.

lightweight models can be obtained, there is still an association between higher performance and greater complexity.

Conversely, with its simplified architecture and fewer parameters, the DiRecNetV2 model achieved remarkable results, outperforming all other models with a notable weighted F1-Score of 0.964. This superior performance was further confirmed through the evaluations using the scoring formulas tailored to assess both speed and accuracy. The DiRecNetV2 model excelled in three critical scenarios: where accuracy and speed are equally prioritized, where speed is paramount, and where accuracy is of most importance as shown by the results in Table 4. When equal emphasis is placed on speed and performance, the proposed model attains a score of 0.980, followed by ConvNeXt, which scores 0.646. In scenarios where the focus is on maximizing the weighted F1 score, the proposed model again leads with a score of 0.988, followed by ConvNeXt with a score of 0.720. Moreover, when the priority shifts to frames per second (FPS), our model maintains the highest performance with a score of 0.972, with SqueezeNet trailing behind at 0.752. Furthermore, when evaluating model performance using the Score2 (Eq.7) DiRecNetV2 again surpasses competing models. Therefore, the proposed model manages to bridge the gap between high accuracy and low efficiency, delivering exceptional classification performance and execution speed. Examples of visual predictions are shown in Figure 3.

**Multi-Label Classification Performance:** During the multi-label evaluation, we applied a threshold of 0.5 to each class’s probability output, post-sigmoid activation, meaning a class was assigned if its probability was higher than this threshold. The evaluation of the different models in the multi-label classification task shows a common trend where precision is consistently high for all models over the three disaster classes. This suggests that the models perform well in terms of specificity; where the instances they label positively for a disaster type are correct, demonstrating a low false positive rate. Recall scores, on the other hand, are less consistent and typically lower in contrast, indicating that the models may not be as good at finding all pertinent examples of each class. This suggests that even though the models accurately predict outcomes, a significant number of cases or "true positives", the instances correctly classified as belonging to a specific class, are being neglected. The reason for the low recall rates observed could be that the models were not optimized for datasets with simultaneous disaster events, thus models are less likely to identify individual disasters and ignore situations where features from multiple disasters coexist in one image. In general, the results suggest that when the models were fine-tuned to transition from recognizing single disaster events to detecting dual disaster scenarios, their performance diminished due to the increased complexity and overlap of disaster characteristics.

Comparison of the models reveals that the DiRecNetV2 stands out with stronger performance across the weighted F1-score which is the harmonic mean of precision and recall. Specifically, it achieves the highest average Weighted F1 metric of 0.614 indicating that DiRecNetV2 maintains a balance between pre-

cision and recall compared to other models, making it the most reliable choice among those listed for multi-label disaster classification. An important observation is the lower recall value of 0.230 for fires compared to earthquakes and floods in the multi-label model is probably due to how the model was trained. During training, the model had been exposed primarily to single-label images where the entire image was associated with one disaster type, thus enabling it to learn distinctive features relevant to that particular disaster. Consequently, when encountering images with multiple disaster types during testing, such as an image featuring both fires and floods, the model may struggle to identify the less prominent fire-related features, leading to a reduced recall rate for the fire category. DiRecNetV2’s performance highlights the possibilities of hybrid models that combine transformer and convolutional-based architectures, not only for single-label classification tasks but also for multi-label ones. It would be interesting to discover how these models’ performance changes in subsequent work when they are fine-tuned on a multi-label dataset. Table 5 summarizes the results of multi-label classification, while visual examples are presented in Figure 4.

ModelName	Earthquakes			Fires			Flood			Average
	Precision	Recall	Weighted F1	Precision	Recall	Weighted F1	Precision	Recall	Weighted F1	
ConvNext Tiny	1.000	0.485	0.653	1.000	0.695	0.820	1.000	0.180	0.305	0.593
EfficientNet-B0	1.000	0.590	0.742	1.000	0.100	0.182	1.000	0.355	0.524	0.483
MnasNet	0.982	0.560	0.713	1.000	0.195	0.326	1.000	0.260	0.413	0.484
MobileNetV2	1.000	0.440	0.611	1.000	0.270	0.425	1.000	0.195	0.326	0.454
MobileNetV3	1.000	0.585	0.738	1.000	0.285	0.444	1.000	0.285	0.444	0.542
ShuffleNetV2	1.000	0.285	0.444	1.000	0.050	0.095	1.000	0.110	0.198	0.246
SqueezeNet	1.000	0.550	0.710	1.000	0.165	0.283	1.000	0.290	0.450	0.481
Convit Tiny	1.000	0.350	0.519	1.000	0.665	<b>0.799</b>	1.000	0.145	0.253	0.524
GCVit XXtiny	1.000	0.325	0.591	1.000	0.650	0.788	0.960	0.120	0.213	0.497
MobileViT s	0.920	0.630	0.748	0.984	0.305	0.466	0.977	0.430	0.597	0.604
MobileViT xs	0.780	0.640	0.703	1.000	0.285	0.444	0.925	0.430	0.587	0.578
MobileViT xxs	0.862	0.685	0.763	1.000	0.185	0.312	0.940	0.625	<b>0.751</b>	0.609
MobileViT V2 050	0.850	0.735	0.788	1.000	0.305	0.467	1.000	0.270	0.425	0.560
MobileViT V2 0100	0.916	0.655	0.764	1.000	0.255	0.406	0.976	0.200	0.332	0.501
Vit Tiny	1.000	0.460	0.630	1.000	0.640	0.780	0.872	0.205	0.332	0.581
DiRecNetV2	1.000	0.605	<b>0.754</b>	1.000	0.280	0.438	0.926	0.500	0.649	<b>0.614</b>

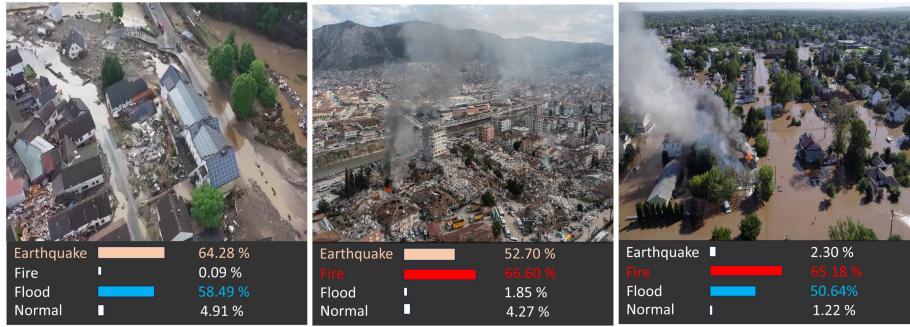
**Table 5.** Model performance comparison for multi-label classification.

#### 4.4 Accuracy-Efficiency Spectrum Gap

The experimental results highlight that architectures of greater complexity achieve greater classification performance, whereas architectures of lesser complexity exhibit higher efficiency and faster execution for both the CNNs and ViTs deep learning frameworks. Our findings serve as a benchmark and demonstrate the gap researchers need to address to design models achieving high classification performance while maintaining FPS within the same order of magnitude. With this objective in mind, our proposed DiRecNetV2 model is designed to bridge this gap.

DiRecNetV2’s architecture combines the strength capabilities of standard convolutions, the efficiency of depth-wise separable convolutions, and the global





**Fig. 4.** Examples of images from the multi-label dataset showcase the predictions of DiRecNetV2, trained for multi-label scenarios. The predictions illustrate the model’s accurate identification of dual instances, with probabilities exceeding 50% for two classes within the same image, underscoring its proficiency in handling complex multi-label classifications.

contextual understanding afforded by the attention mechanism. Unlike the typical vision transformer approach that splits the image into patches, our methodology employs CNN’s feature extraction vectors, which are then fed into the encoder block of the vision transformer. The results indicate that such a hybrid model is not only suitable but potentially superior for the task of disaster recognition from aerial images, where the integration of both local and global features is paramount. The DiRecNetV2 performance demonstrates the effective combination of high classification accuracy, rapid processing (high fps), and low complexity, setting a new benchmark for efficient and accurate image classification tasks. Moreover, architectures like DiRecNetV2 show promise for multi-label scenarios, where images may contain multiple events, as our model also outperforms competitors in multi-label test datasets.

## 5 Conclusion and Future Work

Advances in machine learning and computer vision anticipate a new era of technological empowerment for humanitarian relief by providing tools that improve the efficacy and speed of life-saving steps in times of disaster. At the same time, the integration of unmanned aerial vehicles (UAVs), like drones, with cutting-edge deep learning algorithms marks a significant step forward in enhancing disaster relief efforts. In this research, we propose a hybrid model that combines the benefits of Convolutional Neural Networks (CNNs) and Vision Transformers (ViTs) to develop an accurate, resource-efficient framework suitable for UAV-based disaster detection applications. We introduce the DiRecNetV2, that demonstrates remarkable results by fusing the global contextual awareness of ViTs with the feature extraction power of CNNs. This model not only achieves the highest accuracy among the evaluated models but also ranks second in terms

of processing speed. Implementing evaluation criteria that take into account both processing speed and accuracy using two distinct scoring systems, we proved that our model is particularly suitable for embedded systems. This result highlights the effectiveness and efficiency of such hybrid architectures in real-world applications, emphasizing their potential for future advances in a variety of computer vision tasks. By integrating CNNs and ViTs into a lightweight framework, the hybrid model offers promising solutions for resource-constrained environments, such as edge computing and mobile applications.

Additionally, we present benchmark results for both the AIDERSv2 dataset and our newly introduced multi-label dataset, which comprises 300 images featuring overlapping disasters. These benchmarks encompass performance evaluations for lightweight CNNs and ViTs. To illustrate how these models perform when trained on single-label datasets and then assessed on the proposed multi-label dataset, we provide baseline performances for these models. An important aspect of our study is the provision of benchmark results for ViTs, which remains underexplored in the field of multi-disaster recognition within existing literature. The promising performance of DiRecNetV2 in multi-label tasks underscores the effectiveness of this hybrid approach in handling such scenarios, suggesting its potential applicability in various domains such as environmental monitoring, and industrial safety. For instance, in agricultural monitoring, multi-label identification could help detect various crop diseases and nutrient deficiencies from drone or satellite imagery. Similarly, in urban planning, the model’s ability to identify multiple urban features like buildings, roads, and green spaces from satellite or aerial imagery streamlines city development and infrastructure management processes.

In future studies, we intend to grow the size of the multi-label dataset and fine-tune the examined lightweight model, evaluating the effects of training on this enriched multi-label dataset as opposed to a single-label one on performance. We also intend to explore the explainability side of these algorithms, investigating how CNNs and ViTs differ in feature recognition and comprehending how hybrid models discriminate between different image features in classification tasks. This study will primarily focus on the attention mechanisms in ViTs and hybrid models, providing insights into how they prioritize distinct image regions or features throughout the decision-making process. We expect that our proposed DiRecNetV2 model on the AIDERSv2 and multi-label datasets, combined with the benchmark results of lightweight CNNs and ViTs, will provide a solid foundation for future research. This endeavor aims to encourage the advancement of novel approaches for the application of disaster response, ultimately contributing to the communities impacted by such situations.

## Acknowledgements

This work is supported by the European Union Civil Protection Call for proposals UCPM-2022-KN grant agreement No 101101704 (COLLARIS Network). The work is partially supported by the European Union’s Horizon 2020 research and

innovation program under grant agreement No 739551 (KIOS CoE - TEAMING) and from the Republic of Cyprus through the Deputy Ministry of Research, Innovation and Digital Policy.

Christos Kyrkou would like to acknowledge the support of NVIDIA with the donation of GPU platform.

## Declarations

- *Funding*: See Acknowledgements Section
- *Conflict of interest/Competing interests*: Not Applicable
- *Ethics approval and consent to participate*: Not Applicable
- *Consent for publication*: Not Applicable
- *Data availability*: Immediately with paper publication
- *Materials availability*: Not Applicable
- *Code availability*: Immediately with paper publication
- *Author contribution*: Demetris Shianios made significant contributions to this project through implementation and writing. Christos Kyrkou provided guidance, mentoring, and idea formulation. Panayiotis Kolios provided supervision throughout.

## References

1. Agrawal, T., Meleet, M., et al.: Classification of natural disaster using satellite & drone images with cnn using transfer learning. In: 2021 International Conference on Innovative Computing, Intelligent Communication and Smart Electrical Systems (ICES). pp. 1–5. IEEE (2021)
2. Aipe, A., Mukuntha, N., Ekbal, A., Kurohashi, S.: Deep learning approach towards multi-label classification of crisis related tweets. In: Proceedings of the 15th ISCRAM Conference (2018)
3. Alam, F., Alam, T., Hasan, M., Hasnat, A., Imran, M., Ofli, F., et al.: Medic: A multi-task learning dataset for disaster image classification. arXiv preprint arXiv:2108.12828 (2021)
4. Alam, F., Ofli, F., Imran, M., Alam, T., Qazi, U.: Deep learning benchmarks and datasets for social media image classification for disaster response. In: 2020 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM). pp. 151–158. IEEE (2020)
5. Anggraeni, S.R., Ranggianto, N.A., Ghozali, I., Fatichah, C., Purwitasari, D.: Deep learning approaches for multi-label incidents classification from twitter textual information. Journal of Information Systems Engineering & Business Intelligence **8**(1) (2022)
6. Bhadra, P., Balabantaray, A., Pasayat, A.K.: Mfemanet: an effective disaster image classification approach for practical risk assessment. Machine Vision and Applications **34**(5), 76 (2023)
7. Cao, Q., Liu, Y., Wang, G., He, Y., Wang, K., Liao, S.S., Pu, L.: Building a deep learning model for multi-label classification of natural disasters. In: 2023 IEEE 3rd International Conference on Information Technology, Big Data and Artificial Intelligence (ICIBA). vol. 3, pp. 505–509. IEEE (2023)

8. Chen, F., Yu, B.: Earthquake-induced building damage mapping based on multi-task deep learning framework. *IEEE Access* **7**, 181396–181404 (2019)
9. Chen, J., Shi, X., Gu, L., Wu, G., Su, T., Wang, H.M., Kim, J.S., Zhang, L., Xiong, L.: Impacts of climate warming on global floods and their implication to current flood defense standards. *Journal of Hydrology* **618**, 129236 (2023)
10. Doshi, J., Basu, S., Pang, G.: From satellite imagery to disaster insights. *arXiv preprint arXiv:1812.07033* (2018)
11. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)
12. d’Ascoli, S., Touvron, H., Leavitt, M.L., Morcos, A.S., Biroli, G., Sagun, L.: Convit: Improving vision transformers with soft convolutional inductive biases. In: *International Conference on Machine Learning*. pp. 2286–2296. PMLR (2021)
13. Elangovan, A., Sasikala, S.: A multi-label classification of disaster-related tweets with enhanced word embedding ensemble convolutional neural network model. *Informatica* **46**(7) (2022)
14. Frizzi, S., Bouchouicha, M., Ginoux, J.M., Moreau, E., Sayadi, M.: Convolutional neural network for smoke and fire semantic segmentation. *IET Image Processing* **15**(3), 634–647 (2021)
15. Gadhavi, V.B., Degadwala, S., Vyas, D.: Transfer learning approach for recognizing natural disasters video. In: *2022 Second International Conference on Artificial Intelligence and Smart Energy (ICAIS)*. pp. 793–798. IEEE (2022)
16. Ge, X., Zhao, Q., Wang, B., Chen, M.: Lightweight landslide detection network for emergency scenarios. *Remote Sensing* **15**(4), 1085 (2023)
17. Gebrehiwot, A., Hashemi-Beni, L., Thompson, G., Kordjamshidi, P., Langan, T.E.: Deep convolutional neural network for flood extent mapping using unmanned aerial vehicles data. *Sensors* **19**(7), 1486 (2019)
18. Ghali, R., Akhloufi, M.A., Jmal, M., Soudene Mseddi, W., Attia, R.: Wildfire segmentation using deep vision transformers. *Remote Sensing* **13**(17), 3527 (2021)
19. Hatamizadeh, A., Yin, H., Heinrich, G., Kautz, J., Molchanov, P.: Global context vision transformers. In: *International Conference on Machine Learning*. pp. 12633–12646. PMLR (2023)
20. Hossain, F.A., Zhang, Y., Yuan, C., Su, C.Y.: Wildfire flame and smoke detection using static image features and artificial neural network. In: *2019 1st international conference on industrial artificial intelligence (iai)*. pp. 1–6. IEEE (2019)
21. Howard, A., Sandler, M., Chu, G., Chen, L.C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., et al.: Searching for mobilenetv3. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 1314–1324 (2019)
22. Iandola, F.N., Han, S., Moskewicz, M.W., Ashraf, K., Dally, W.J., Keutzer, K.: Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size. *arXiv preprint arXiv:1602.07360* (2016)
23. Ignatov, A., Malivenko, G., Timofte, R.: Fast and accurate quantized camera scene detection on smartphones, mobile ai 2021 challenge: Report. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 2558–2568 (2021)
24. Jadon, A., Omama, M., Varshney, A., Ansari, M.S., Sharma, R.: Firenet: a specialized lightweight fire & smoke detection model for real-time iot applications. *arXiv preprint arXiv:1905.11922* (2019)

25. Ji, M., Liu, L., Zhang, R., F Buchroithner, M.: Discrimination of earthquake-induced building destruction from space using a pretrained cnn model. *Applied Sciences* **10**(2), 602 (2020)
26. Jiao, Z., Zhang, Y., Mu, L., Xin, J., Jiao, S., Liu, H., Liu, D.: A yolov3-based learning strategy for real-time uav-based forest fire detection. In: 2020 Chinese Control And Decision Conference (CCDC). pp. 4963–4967. IEEE (2020)
27. Khudayberdiev, O., Zhang, J., Elkhail, A., Balde, L.: Fire detection approach based on vision transformer. In: International Conference on Adaptive and Intelligent Systems. pp. 41–53. Springer (2022)
28. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. *Communications of the ACM* **60**(6), 84–90 (2017)
29. Kyrkou, C., Kolios, P., Theocharides, T., Polycarpou, M.: Machine learning for emergency management: A survey and future outlook. *Proceedings of the IEEE* **111**(1), 19–41 (2023). <https://doi.org/10.1109/JPROC.2022.3223186>
30. Kyrkou, C., Theocharides, T.: Deep-learning-based aerial image classification for emergency response applications using unmanned aerial vehicles. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 517–525 (2019). <https://doi.org/10.1109/CVPRW.2019.00077>
31. Kyrkou, C., Theocharides, T.: Emergencynet: Efficient aerial image classification for drone-based emergency monitoring using atrous convolutional feature fusion. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* **13**, 1687–1699 (2020). <https://doi.org/10.1109/JSTARS.2020.2969809>
32. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* **86**(11), 2278–2324 (1998)
33. Lee, G.Y., Dam, T., Ferdaus, M.M., Poenar, D.P., Duong, V.N.: Watt-effnet: A lightweight and accurate model for classifying aerial disaster images. *IEEE Geoscience and Remote Sensing Letters* (2023)
34. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 11976–11986 (2022)
35. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. pp. 11976–11986 (June 2022)
36. Ma, H., Liu, Y., Ren, Y., Wang, D., Yu, L., Yu, J.: Improved cnn classification method for groups of buildings damaged by earthquake, based on high resolution remote sensing images. *Remote Sensing* **12**(2), 260 (2020)
37. Ma, H., Liu, Y., Ren, Y., Yu, J.: Detection of collapsed buildings in post-earthquake remote sensing images based on the improved yolov3. *Remote Sensing* **12**(1), 44 (2019)
38. Ma, N., Zhang, X., Zheng, H.T., Sun, J.: Shufflenet v2: Practical guidelines for efficient cnn architecture design. In: *Proceedings of the European conference on computer vision (ECCV)*. pp. 116–131 (2018)
39. Mao, J., Harris, K., Chang, N.R., Pennell, C., Ren, Y.: Train and deploy an image classifier for disaster response. In: 2020 IEEE High Performance Extreme Computing Conference (HPEC). pp. 1–5. IEEE (2020)
40. Masson-Delmotte, V., Zhai, P., Pirani, A., Connors, S.L., Péan, C., Berger, S., Caud, N., Chen, Y., Goldfarb, L., Gomis, M., et al.: Climate change 2021: the physical science basis. Contribution of working group I to the sixth assessment report of the intergovernmental panel on climate change p. 2 (2021)
41. Mehta, S., Rastegari, M.: Separable self-attention for mobile vision transformers. arxiv 2022. arXiv preprint arXiv:2206.02680

42. Mehta, S., Rastegari, M.: Mobilevit: light-weight, general-purpose, and mobile-friendly vision transformer. arXiv preprint arXiv:2110.02178 (2021)
43. Mo, P., Li, D., Liu, M., Jia, J., Chen, X.: A lightweight and partitioned cnn algorithm for multi-landslide detection in remote sensing images. *Applied Sciences* **13**(15), 8583 (2023)
44. Munawar, H.S., Hammad, A., Ullah, F., Ali, T.H.: After the flood: A novel application of image processing and machine learning for post-flood disaster management. In: *Proceedings of the 2nd International Conference on Sustainable Development in Civil Engineering (ICSDC 2019)*, Jamshoro, Pakistan. pp. 5–7 (2019)
45. Munawar, H.S., Ullah, F., Qayyum, S., Khan, S.I., Mojtahedi, M.: Uavs in disaster management: Application of integrated aerial imagery and convolutional neural network for flood detection. *Sustainability* **13**(14), 7547 (2021)
46. Munsif, M., Afridi, H., Ullah, M., Khan, S.D., Cheikh, F.A., Sajjad, M.: A lightweight convolution neural network for automatic disasters recognition. In: *2022 10th European Workshop on Visual Information Processing (EUVIP)*. pp. 1–6. IEEE (2022)
47. Pally, R., Samadi, S.: Application of image processing and convolutional neural networks for flood image classification and semantic segmentation. *Environmental Modelling & Software* **148**, 105285 (2022)
48. Rahnemounfar, M., Chowdhury, T., Sarkar, A., Varshney, D., Yari, M., Murphy, R.R.: Floodnet: A high resolution aerial imagery dataset for post flood scene understanding. *IEEE Access* **9**, 89644–89654 (2021)
49. Roy, R., Kulkarni, S.S., Soni, V., Chittora, A., et al.: Transformer-based flood scene segmentation for developing countries. arXiv preprint arXiv:2210.04218 (2022)
50. Saad, O.M., Chen, Y., Savvaadis, A., Fomel, S., Chen, Y.: Real-time earthquake detection and magnitude estimation using vision transformer. *Journal of Geophysical Research: Solid Earth* **127**(5), e2021JB023657 (2022)
51. Saini, N., Chattopadhyay, C., Das, D.: E2alertnet: An explainable, efficient, and lightweight model for emergency alert from aerial imagery. *Remote Sensing Applications: Society and Environment* **29**, 100896 (2023)
52. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., Chen, L.C.: Mobilenetv2: Inverted residuals and linear bottlenecks. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 4510–4520 (2018)
53. Sarp, S., Kuzlu, M., Cetin, M., Sazara, C., Guler, O.: Detecting floodwater on roadways from image data using mask-r-cnn. In: *2020 International Conference on INnovations in Intelligent SysTems and Applications (INISTA)*. pp. 1–6. IEEE (2020)
54. Shamsoshoara, A., Afghah, F., Razi, A., Zheng, L., Fulé, P.Z., Blasch, E.: Aerial imagery pile burn detection using deep learning: the flame dataset. *Computer Networks* **193**, 108001 (2021)
55. Shi, L., Zhang, F., Xia, J., Xie, J., Zhang, Z., Du, Z., Liu, R.: Identifying damaged buildings in aerial images using the object detection method. *Remote Sensing* **13**(21), 4213 (2021)
56. Shianios, D., Kyrkou, C., Kolios, P.S.: A benchmark and investigation of deep-learning-based techniques for detecting natural disasters in aerial images. In: *International Conference on Computer Analysis of Images and Patterns*. pp. 244–254. Springer (2023)
57. Singh, S., Ghosh, S., Maity, A., Bag, B.C., Koley, C., Maity, H.K.: Disasternet: a multi-label disaster aftermath image classification model. In: *ICT Systems and Sustainability: Proceedings of ICT4SD 2021, Volume 1*. pp. 481–490. Springer (2022)

58. Steiner, A., Kolesnikov, A., Zhai, X., Wightman, R., Uszkoreit, J., Beyer, L.: How to train your vit? data, augmentation, and regularization in vision transformers. arXiv preprint arXiv:2106.10270 (2021)
59. Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., Le, Q.V.: Mnasnet: Platform-aware neural architecture search for mobile. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2820–2828 (2019)
60. Tan, M., Le, Q.: Efficientnet: Rethinking model scaling for convolutional neural networks. In: International conference on machine learning. pp. 6105–6114. PMLR (2019)
61. Uwishema, O.: Addressing the effects of the earthquakes on türkiye’s health-care system. *The Lancet* **401**(10378), 727 (2023)
62. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
63. Xie, S., Hou, C., Yu, H., Zhang, Z., Luo, X., Zhu, N.: Multi-label disaster text classification via supervised contrastive learning for social media data. *Computers and Electrical Engineering* **104**, 108401 (2022)
64. Xiong, C., Yu, A., Rong, L., Huang, J., Wang, B., Liu, H.: Fire detection system based on unmanned aerial vehicle. In: 2021 IEEE International Conference on Emergency Science and Information Technology (ICESIT). pp. 302–306. IEEE (2021)
65. Yang, N.T.S., Tham, M.L., Chua, S.Y., Lee, Y.L., Owada, Y., Poomrittigul, S.: Efficient device-edge inference for disaster classification. In: 2022 Thirteenth International Conference on Ubiquitous and Future Networks (ICUFN). pp. 314–319. IEEE (2022)
66. Yang, W., Zhang, X., Luo, P.: Transferability of convolutional neural network models for identifying damaged buildings due to earthquake. *Remote Sensing* **13**(3), 504 (2021)
67. Yuan, J., Ding, X., Liu, F., Cai, X.: Disaster classification net: A disaster classification algorithm on remote sensing imagery. *Frontiers in Environmental Science* **10**, 2690 (2023)
68. Yuan, J., Ma, X., Han, G., Li, S., Gong, W.: Research on lightweight disaster classification based on high-resolution remote sensing images. *Remote Sensing* **14**(11), 2577 (2022)
69. Yuan, J., Ma, X., Zhang, Z., Xu, Q., Han, G., Li, S., Gong, W., Liu, F., Cai, X.: Effic-net: lightweight fully convolutional neural networks in remote sensing disaster images. *Geo-spatial Information Science* pp. 1–12 (2023)
70. Zhang, M.L., Zhou, Z.H.: A review on multi-label learning algorithms. *IEEE transactions on knowledge and data engineering* **26**(8), 1819–1837 (2013)