

# LADI v2: Multi-label Dataset and Classifiers for Low-Altitude Disaster Imagery

Sam Scheele  
MIT Lincoln Laboratory  
244 Wood St  
Lexington, MA 02421 USA  
Samuel.Scheele@ll.mit.edu

Katie Picchione  
NASA Disasters Program  
300 Hidden Figures Way SW  
Washington, DC 20024  
Katie.Picchione@nasa.gov

Jeffrey Liu  
MIT Lincoln Laboratory  
244 Wood St  
Lexington, MA 02421 USA  
Jeffrey.Liu@ll.mit.edu

## Abstract

*ML-based computer vision models are promising tools for supporting emergency management operations following natural disasters. Imagery taken from small manned and unmanned aircraft can be available soon after a disaster and provide valuable information from multiple perspectives for situational awareness and damage assessment applications. However, emergency managers often face challenges in effectively utilizing this data due to the difficulties in finding the most relevant imagery among the tens of thousands of images that may be taken after an event. Despite this promise, there is still a lack of training data for imagery of this type from multiple perspectives and for multiple hazard types. To address this, we present the LADI v2 (Low Altitude Disaster Imagery version 2) dataset, a curated set of about 10,000 disaster images captured by the Civil Air Patrol (CAP) in response to over 50 disaster events (2015-2023) from over 30 US states and territories and annotated for multi-label classification by trained CAP volunteers. We also provide two pretrained baseline classifiers and compare their performance to state-of-the-art vision-language models in multi-label classification. The data and code are released publicly to support the development of computer vision models for emergency management research and applications.*

## 1. Introduction

Rapid and accurate assessment of post-disaster conditions is critical for effective disaster response and recovery operations. Low altitude aerial imagery, such as aerial photographs collected in small manned aircraft by the Civil Air Patrol (CAP) in the United States or images from small Unmanned Aerial Systems (sUAS), can provide valuable information about the extent and severity of damage caused by disasters. However, the large quantity of images collected during these missions can present a significant challenge for analysts tasked with identifying actionable infor-

mation in a timely manner. To address this challenge, we introduce the Low Altitude Disaster Imagery v2 (LADI v2) dataset, a multi-label image classification dataset designed to facilitate the development of computer vision models for identifying useful post-disaster aerial images. LADI v2 builds upon existing work in image-based damage assessment [12, 16, 20, 30, 31, 35] by providing a diverse, multi-hazard dataset that includes oblique and nadir aerial imagery from various locations and disaster declarations across the United States.

The technical contributions of LADI v2 are twofold. First, we present a curated dataset of post-disaster aerial imagery from multiple perspectives for multiple hazard types, with high quality, operationally-relevant labels annotated by trained CAP volunteers for multi-label classification. Second, we provide two pretrained baseline reference classifiers. We also demonstrate the utility of the dataset as a benchmark, and compare the performance of the baseline classifier to state-of-the-art vision-language models (VLM) on open-vocabulary classification as a benchmark. We outperform the open source VLM on nearly all classes and are competitive with the commercial VLM in test set, and broadly outperform both VLMs in the validation set. This demonstrates the continued need for open, domain-specific training data for specialized applications such as disaster response.

LADI v2 also offers unique characteristics that can contribute to the advancement of machine learning research. The dataset features a realistic distribution shift between the training and test sets, representing annual variation in disaster incident types, as well as changes due to new operational procedures and technologies. This characteristic makes LADI v2 a valuable benchmark for evaluating domain adaptation techniques in the context of disaster response.

The dataset, classifiers, and associated documentation are made openly available on GitHub [25] and Hugging Face [24, 26], enabling researchers and practitioners to build upon this work and adapt the models to their specific

needs. Through this contribution, we aim to streamline the process of identifying useful post-disaster aerial images, ultimately supporting more efficient and effective disaster response efforts while providing a valuable resource for the broader machine learning community.

The paper is structured as follows: Section 2 covers related work in disaster imagery datasets and necessary background on vision-language models; Section 3 provides the details of the dataset; Section 4 discusses the pretrained baseline classifiers; and Section 5 concludes with a summary and comments on limitations and future directions.

## 2. Related Work

### 2.1. Natural Disaster Imagery Datasets

There is existing work on using imagery and computer vision to facilitate post disaster damage assessment [12, 16, 20, 23, 30, 31, 35]. We highlight a few relevant examples and comment on the gap that LADI v2 addresses. The xBD dataset [12] provides 23000 labeled nadir-perspective satellite images annotated with bounding boxes for building damage for multiple locations across the globe, covering multiple hazard types. FloodNet [30] and RescueNet [31] each provide segmentation masks, classification labels, and visual question-answering captions for high resolution, low altitude aerial imagery from UAVs, but are limited to single incidents each: Hurricanes Harvey and Michael, respectively. Incidents1M [35] provides a large multi-label dataset classifying 43 incident types in 49 outdoor locations from nearly 1 million images scraped from the web. These images feature a variety of perspectives and locations, but are primarily ground-based, and only identify the type and location of the incident. CRASAR-U-DROIDS [23] provides annotated building damage polygons for orthomosaic imagery from small Unmanned Aerial Systems (sUAS). Compared to the existing literature, there is still a lack of training data to support both oblique and nadir low-altitude aerial imagery—which is increasingly common in disaster response applications with UAVs and small manned aircraft—across multiple geographies, event types, and damage and infrastructure categories.

Version 1 of the Low Altitude Disaster Imagery dataset (LADI v1) began to address these gaps by providing annotations for low altitude, multi-perspective imagery from multi-hazard and multi-geographic events [20]. It was included in a number of NIST TRECVID challenges [1–3]. However, these labels were created by untrained crowd-sourced workers, and the term “damage” was not clearly defined. Instead, LADI v2 was labeled by a team of volunteer annotators from the Civil Air Patrol who have been trained in the FEMA damage assessment process. Damage labels follow FEMA’s criteria for preliminary damage assessments (PDAs), which articulate five levels of damage:

unaffected, affected, minor damage, major damage, and destroyed. Furthermore, the label set and annotator training materials were developed in conjunction with FEMA’s Response Geospatial Office to ensure compatibility with their procedures. For these reasons, LADI v2 offers a different label set than LADI v1. Nevertheless, we still believe LADI v1 offers some value, since it features a larger number of annotated images, and thus may serve as a suitable pretraining task for classifiers trained on LADI v2.

### 2.2. Vision-Language Models

Recent advances in vision-language models represent extremely promising developments toward generalizable solutions for computer vision problems [11]. Vision-language models typically include an image encoder and text encoder for each respective data modality. These models are trained on tasks to promote alignment between the image encoding and text encoding; such tasks include contrastive image-text learning, popularized by [29], as well as input reconstruction such as masked language modeling [7] and masked image modeling [34]. These vision-language models can be used to address various computer vision tasks, including open vocabulary classification, object detection, and segmentation [11]. We discuss a few vision language models that are relevant to this paper, particularly LLaVA-NeXT [19], and GPT-4o [27].

LLaVA (Large Language and Vision Assistant) [18] and its refinements LLaVA 1.5 [17] and LLaVA-NeXT [19] build upon CLIP by incorporating a large language model decoder, such as Vicuna [5] or Mistral [14], to allow it to be trained additional tasks, such as visual question answering. Commercial offerings such as GPT-4o [27] offer similar multimodal capabilities. We benchmark our models on LADI v2 against zero-shot classification using LLaVA-NeXT and GPT-4o in Section 4.3.

## 3. LADI v2 Dataset

### 3.1. Data collection and annotation

LADI v2 images are sourced from the FEMA (Federal Emergency Management Agency) Civil Air Patrol Image Uploader Repository, publicly hosted in an Amazon Web Services s3 bucket [9], which contains aerial photographs collected in support of federally declared disasters from 2009 onward, as well as CAP training missions. Figure 1 shows a sample of images in the dataset. Image metadata includes timestamp and location information. To ensure LADI v2 contains only images collected during disasters, we compared image metadata against disaster declaration data from the OpenFEMA Disaster Declarations API [10] to identify images taken within 14 days of the start of a declared federal disaster and within an affected county’s boundaries.



Figure 1. A sample of images from our training set. LADI v2 contains images with both positive and negative examples of damage from a range of altitudes, perspectives, geographies, and lighting conditions

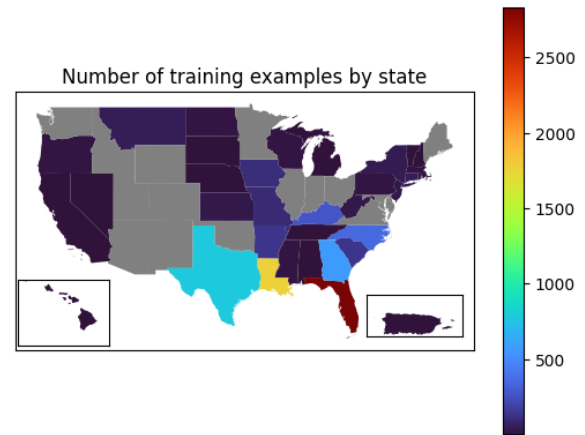
We consulted emergency management professionals at the FEMA Response Geospatial Office to develop the set of labels and labeling instructions. Labels were chosen to help emergency managers identify the most relevant images when conducting initial damage assessments, which can support disaster declarations and assistance grants. The label set is provided as the “v2” set in Table 1. The various damage levels for the buildings—affected, minor, major, and destroyed, listed in increasing order of severity—are determined based on the FEMA preliminary damage assessment criteria [8].

Images were annotated by a team of 46 Civil Air Patrol volunteers who had been previously trained in the FEMA preliminary damage assessment process [8]. Each image was shown to three volunteers, and labels were assigned by majority vote when there was disagreement between the annotators.

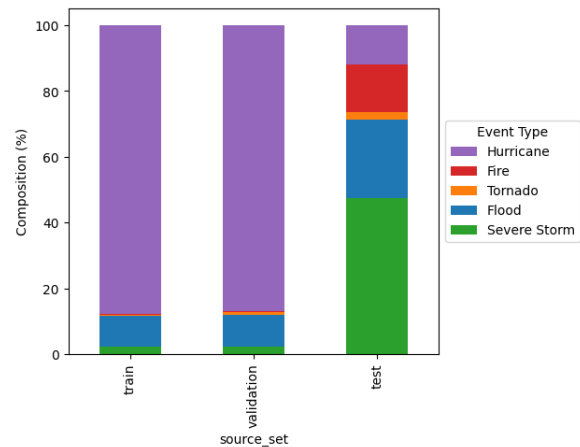
v2	v2a
bridges_any	bridges_any
bridges_damage	buildings_any
buildings_any	buildings_affected_or_greater
buildings_affected	buildings_minor_or_greater
buildings_destroyed	debris_any
buildings_major	flooding_any
buildings_minor	flooding_structures
debris_any	roads_any
flooding_any	roads_damage
flooding_structures	trees_any
roads_any	trees_damage
roads_damage	water_any
trees_any	
trees_damage	
water_any	

Table 1. The v2 and v2a set of labels.

After initial experiments with training classifiers on the “v2” label set, we found that there were very few positive examples of “bridges\_damage”. In addition, we found that performance was poor when distinguishing between various levels of building damage (“buildings\_affected”, “buildings\_minor”, “buildings\_major”, and



(a) Number of training examples per state



(b) Distribution of events in the dataset splits

Figure 2. Details of the dataset: the label sets, number of training examples by state, and event type distribution of the various splits.

“buildings\_destroyed”). FEMA staff advised that it was less important to classify building damage categories than to determine whether buildings that had sustained any level of damage were present in image. In light of this, we removed the “bridges\_damage” class and combined the building damage categories into “buildings\_affected\_or\_greater” and “buildings\_minor\_or\_greater”. This revised set of labels is called the “v2a” label set, as shown in Table 1, and is what we report our results on. The “v2a” label set

contains 12 labels, categorized into 5 non-damage-related (“bridges\_any”, “buildings\_any”, “roads\_any”, “trees\_any”, and “water\_any”) and 7 damage-related (“buildings\_affected\_or\_greater”, “buildings\_minor\_or\_greater”, “debris\_any”, “flooding\_any”, “flooding\_structures”, “roads\_damage”, and “trees\_damage”).

### 3.2. Dataset statistics and characteristics

Our dataset consists of 9,963 images, split into 8,030 train examples, 892 validation examples, and 1,041 test examples. The training and validation examples are drawn from disaster declarations between 2015-2022, and the test examples are drawn from disaster declarations in 2023. In total, the dataset draws from over 100 disaster declarations from more than 30 US states and territories, see Figure 2. Since disaster declarations can span multiple states, we estimate the number of distinct disaster events by clustering the declarations based on incident date, and we estimate that at least 50 distinct disaster events are in the dataset.

*Disaster type distribution.* Since the training and validation examples are drawn from the same set of disaster incidents, the distributions of hazard types for those two splits are quite similar, whereas the test set has comparatively more images from severe storms, floods, and fires, and many fewer images from hurricanes. The hazard type distributions of the splits are visualized in Figure 2b.

*Label distribution and co-occurrence.* Figures 3a and 3b give the label-label co-occurrence matrices and label-hazard type co-occurrence matrices for the training, validation, and test sets. As the training and validation sets are randomly drawn from the same disaster incidents, they do not have substantially different statistical distributions and are combined in these figures for the sake of brevity. We observe a substantial difference between the training/validation co-occurrence matrices and the test co-occurrence matrices. In the label-hazard type co-occurrence matrix, this is easily explained by the fact that different incident hazard types were more common in the 2023 test set. The label-label matrices require slightly more in-depth explanation.

The label-label matrices indicate that images showing damage occur less frequently in the test set. We believe this is due to a combination of factors. First, there is natural variation in the intensity and distribution of damage for different incidents. Second, CAP has recently expanded their use of the WaldoAir camera system [28]. This system takes images at regular time intervals in a grid pattern [28], whereas images taken from handheld cameras tend to focus on pre-selected targets or areas that the photographer selects, which often are areas with prominent damage. As a consequence, a much lower percentage of the WaldoAir images contain damage. The train and validation sets consist of about 50% Waldo images, while the test set consists of 65% Waldo images.

*Distribution shift.* The distribution shifts between the train/validation and test sets represent challenges in disaster applying for machine learning. The distribution of hazard types, severity, and locations change year to year based on cyclical weather patterns [15] and climate change [13]. Furthermore, changes in operational procedures and technology, such as the increased adoption of the WaldoAir system, can lead to distribution variation even among disasters resulting from the same type of hazard. This underscores one of the key domain-specific challenges of applying machine learning solutions for disaster response applications. To our knowledge, no other benchmark disaster imagery dataset explicitly addresses the shift in data distribution from year to year.

### 3.3. Comparison to Existing Datasets

Compared to LADI v1 [20], v2 offers higher-quality labels from trained annotators, and a different label set. While both versions of the dataset draw from the same repository of public domain operational FEMA CAP images [9], only about 2.4% of images from v2 also appear in v1.

The label set for LADI v2 includes building damage on the FEMA PDA scale, which is compatible with the damage scale used by xBD [12], RescueNet [31], and CRASAR-U-DROIDS [23]. Compared to those datasets, LADI v2 includes more distinct events (LADI v2: 50+, xBD: 19, RescueNet: 1, CRASAR-U-DROIDS: 10), total pixels: (LADI v2: 345.32e9, xBD: 23.14e9, RescueNet 53.99e9, CRASAR-U-DROIDS: 67.13e9), and among aerial datasets, area covered (LADI v2: 161.4 km<sup>2</sup>, RescueNet: 3.6 km<sup>2</sup>, CRASAR-U-DROIDS: 67.98 km<sup>2</sup>)<sup>1</sup>. LADI is also unique in its inclusion of both oblique and nadir perspective imagery. However, LADI v2’s labels only support image classification, whereas the other mentioned datasets provide segmentation polygons for building damage. Thus, LADI v2 serves to complement the capabilities of those existing datasets, and would be an ideal candidate to include in a multi-task training framework, especially due to the alignment in building damage labels.

## 4. Pretrained Classifiers

### 4.1. Architecture and training details

To support research and deployment applications, we provide two pretrained reference classifiers, LADI-v2-classifier-small-reference and LADI-v2-classifier-large-reference,<sup>2</sup> hereby referred to as the “small” and “large” classifiers for brevity. The small classifier is based on

<sup>1</sup>xBD is a satellite imagery dataset, which covers much wider area at lower resolution. It covers over 45000 km<sup>2</sup>

<sup>2</sup>We provide four classifiers in the repository. The “reference” versions, discussed in his paper, are trained only on the train set. The “main” versions are trained on all splits and intended for deployment and downstream applications.

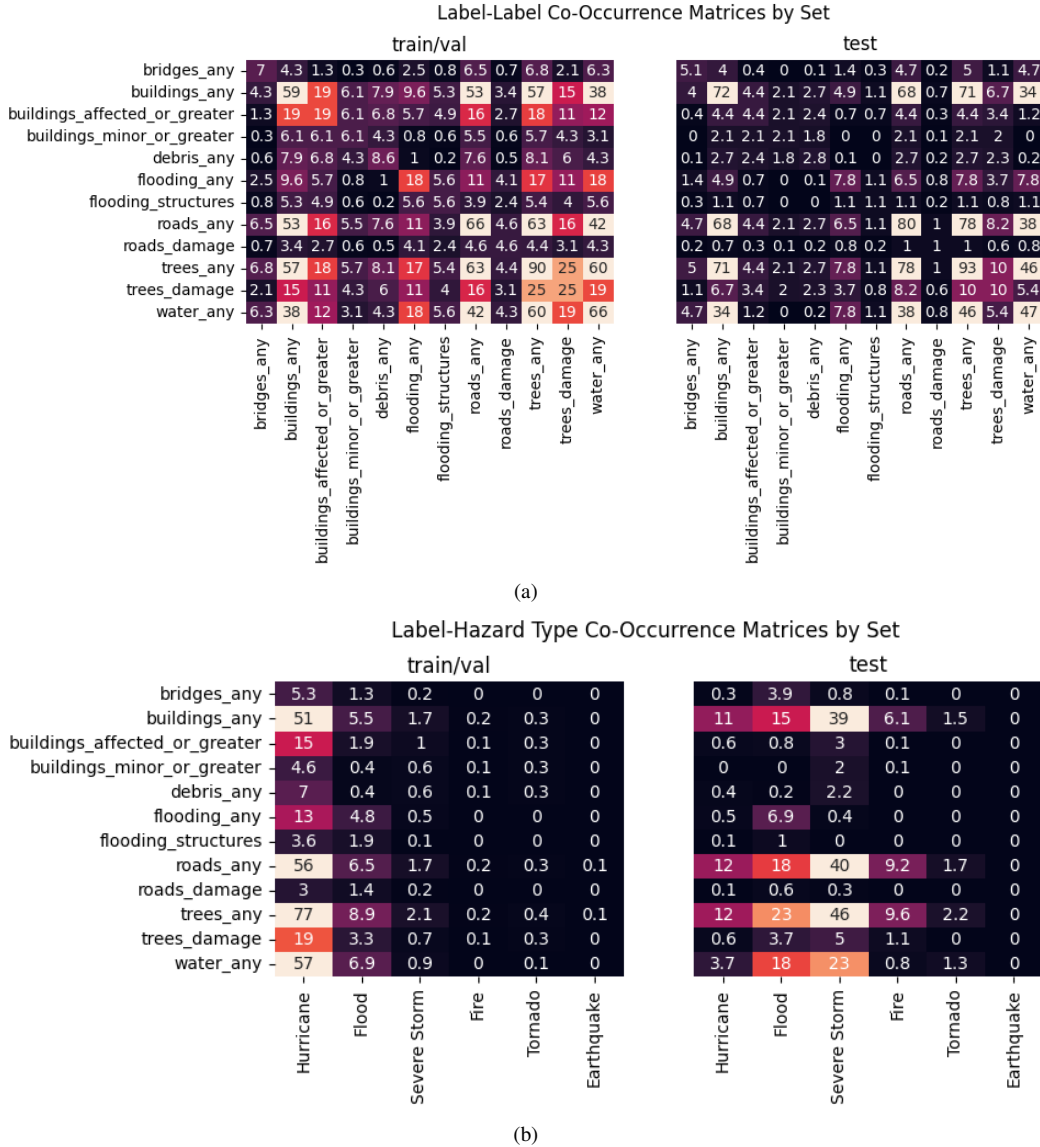


Figure 3. Co-occurrence matrices for data splits. Numbers indicate percentage of images within the given split that have the given combination of labels. Lighter colors indicate higher percentages. Note that training and validation sets are combined in these figures for brevity due to their similar distributions.

the Big-Transfer (BiT-50) architecture [4], pretrained on ImageNet-1k [6], while the large classifier is based on Swin v2 Large [21], pretrained on ImageNet-21k [32] and finetuned on ImageNet-1k [6]. Standard random augmentations of resizing, cropping, horizontal flipping, affine transformations, and color jitter were applied to the training images.

These architectures were the best performers among a number of available models. We first finetuned twenty pretrained classifiers, available on Hugging Face, on LADI v2 using default settings; the list of models evaluated and their respective performance is provided in the supplementary material. We chose the top two performing candi-

dates and performed additional optimization through hyperparameter tuning and pretraining on LADI v1 [20]. Hyperparameter tuning was done via random search on optimizer type, learning rate, and learning rate scheduler. We considered the AdamW [22] and Adafactor [33] optimizers, initial learning rates of  $2 \times 10^{-4}$ ,  $1 \times 10^{-4}$ , and  $5 \times 10^{-5}$ , and LR schedulers with exponential LR decay ( $\gamma$ : 0.5 or 0.9) or with reduce LR on plateau ( $\gamma$ : 0.5 or 0.9, patience: 5 or 10 epochs). Results for the top two architectures are shown in ablation in Table 3; we evaluate using mean Average Precision (mAP).

The final configurations for LADI-v2-classifier-small

Hyperparameter	Large	Small
Architecture	Swin v2	BiT-50
LR Scheduler	Exponential	Exponential
Optimizer	AdamW	AdamW
Initial LR	$5 * 10^{-5}$	$1 * 10^{-4}$
$\gamma$	0.9	0.9
Epochs	50	50

Table 2. Selected hyperparameters for our reference small and large models (transposed).

and LADI-v2-classifier-large are indicated in Table 3 with  $\dagger$  and  $\ddagger$  respectively. Hyperparameter tuning provides significant benefit for both architectures. Pretraining on LADI v1 provides modest benefits in the test mAP for the model based on Swin v2, but degrades performance for the model based on BiT-50. This is consistent with the observations of Beyer et al. [4], who find that pretraining BiT-50 on larger datasets does not improve, but instead degrades fine-tuning performance. We quote Beyer et al. [4]: “Not only is there limited benefit of training a large model size on a small dataset, but there is also limited (or even negative) benefit from training a small model on a larger dataset... Thus, if one uses only a ResNet50x1 [the architecture for BiT-50], one may conclude that scaling up the dataset does not bring any additional benefits.”

For the sake of brevity, all subsequent results are presented only for the “large” model.

Arch.	Tuning	Pretrain	Val_mAP	Test_mAP
BiT-50	No	No	89.6	89.0
Swin v2	No	No	88.7	86.0
$\dagger$ BiT-50	Yes	No	93.3	91.1
Swin v2	Yes	No	<b>93.8</b>	92.3
BiT-50	Yes	Yes	87.9	85.9
$\ddagger$ Swin v2	Yes	Yes	<b>93.8</b>	<b>92.6</b>

Table 3. Classifier performance ablation with hyperparameter tuning and LADI v1 pretraining. The small model configuration is indicated with  $\dagger$  and the large model with  $\ddagger$ .

## 4.2. Direct Comparison Against LADI v1

Although the label sets of LADI v1 and v2 differ, we developed a mapping that condenses both into a common set of classes to facilitate a direct comparison. This mapping enables us to compare the marginal benefit of using trained annotators over crowdsourced workers<sup>3</sup>. The mapping is presented in the supplementary material, where the condensed classes are ‘building’, ‘flooding’, ‘road’, ‘damage’,

<sup>3</sup>This analysis does not fully account for other enhancements from v1 to v2, such as the more standardized and detailed damage labels.

and ‘debris’. We evaluated both models on LADI v2’s validation and test splits, with results displayed in Table 4. In all metrics, the model trained on LADI v2 outperformed, highlighting the advantages of higher-quality labels provided by trained CAP annotators.

Split	Ver.	Prec	Rec	$F_1$	mAP
val	v1	0.799	0.808	0.800	0.869
val	v2	<b>0.877</b>	<b>0.895</b>	<b>0.886</b>	<b>0.946</b>
test	v1	0.850	0.768	0.800	0.877
test	v2	<b>0.871</b>	<b>0.850</b>	<b>0.859</b>	<b>0.917</b>

Table 4. Performance of the large classifier on the v2 test and validation sets when trained on the LADI v1 train set versus the LADI v2 train set, as measured by mean precision, recall, F1 score, and average precision. The top scores for each split are bolded.

## 4.3. Benchmarking open-vocabulary classification (LLaVA and GPT-4o)

LADI v2 also has potential to be an effective benchmark for supervised and zero-shot classification of post-disaster imagery. In particular, the validation and test sets of LADI v2 offer the ability to test model performance against a realistic distribution shift observed in practice. We demonstrate this by comparing our model’s performance to recent open-vocabulary classifiers.

We first evaluate our model compared to the recently released open 7.5 billion parameter LLaVA-NeXT model [19] on a zero-shot classification task. For each class, the LLaVA model saw an image in the test or validation set, followed by a prompt such as, “Does this image contain `class_name`? Answer with ‘yes’ or ‘no.’” For the classes involving FEMA preliminary damage assessment categories, we included a summary of the damage category criteria based on the FEMA Preliminary Damage Assessment Pocket Guide [8] in the prompt. The model outputs were converted to binary labels and used to compute the  $F_1$  score and shown in Table 5. Since the VLMs do not provide a confidence score, we are unable to calculate scores like Average Precision for the VLMs. In the test set, our method outperformed LLaVA-NeXT on all labels, including all damage-related classes, except in three categories: “bridges\_any”, “roads\_any”, and “trees\_any”, in which our model comes within 3% of LLaVA-NeXT. In the validation set, our model outperforms LLaVA-NeXT in all categories.

We also compare our model to a commercial multimodal model, GPT-4o, using the same prompt format as above. The results are shown in Table 5. On the test set, our model is competitive with GPT-4o, beating or tying its performance on 5 of the 12 classes. On the validation set, we outperform GPT-4o in all but the “water\_any” class, showing strong in-sample performance.

Compared to the open source LLaVA-NeXT, GPT-4o broadly outperforms. Since the details of GPT-4o’s training and architecture are not public, we can only speculate on the performance gap. We believe that one contributor to the performance gap may be access to sufficient openly available high-quality data for the given domain, thus demonstrating a need for high-quality, open-access labeled disaster-related aerial imagery. Though GPT-4o clearly highlights the potential of multimodal vision-language models for computer vision tasks—particularly in its demonstrated performance across the different distributions of the validation and test sets—its closed-source nature and restrictive licensing makes it difficult to build derivative works.

GPT-4o also suffers from practical challenges for operational use. It is only accessible over the internet via an API, and we found that the API frequently timed out when trying to download large high-resolution images, such as those captured by CAP. This required the images to be downloaded, resized, re-encoded, and uploaded in small batches. As such, it is not suitable for large-scale or time-critical tasks, nor for usage in offline environments. In comparison, our pretrained classifiers run much quicker than both VLMs, and can run offline, unlike GPT-4o.

We note that in many classes, we observe performance degradation between the validation and the test set, as shown in the last 3 columns of Table 5. This highlights some of the challenge encountered with predictions under a distribution shift. However, we also see that there is significant performance degradation (>10%) even in the VLMs, particularly in certain damage categories (`debris_any`, `flooding_any`, `flooding_structures` for GPT-4o, and `buildings_affected_or_greater`, `buildings_minor_or_greater`, `flooding_any`, `flooding_structures`, `roads_damage`, `trees_damage` for LLaVA), which have not seen the training data. Thus, we can infer that some of the performance degradation can be attributed to the difference in underlying difficulty of the problem, rather than overfitting to the training data.

Whereas our baseline model was trained with “standard” supervised image classification techniques, we anticipate that approaches that incorporate domain adaptation techniques, or more recent architectures, such as multimodal language models, should handle the distribution shift better than our model. The baseline model can thus be used in conjunction with the LADI v2 validation and test sets to benchmark the efficacy of such approaches.

## 5. Conclusion

*Summary of Contributions.* In this paper, we introduce the LADI v2 dataset, a curated collection of post-disaster aerial images from multiple perspectives, hazard types, and geographies across the United States. The dataset addresses

the need for high-quality, diverse training data to support the development of computer vision models for disaster response applications. To facilitate research and implementation efforts, we provide the dataset and two pretrained reference classifiers as open-source resources.

One of the key strengths of LADI v2 is the quality of its annotations, which were provided by trained Civil Air Patrol volunteers using label sets and training materials developed in collaboration with FEMA. This ensures that the labels are consistent with the standards used by emergency management professionals, enhancing its practical utility.

Furthermore, LADI v2 features a realistic distribution shift between the training, validation, and test splits, capturing the year-to-year variability in disaster events as well as changes in operational procedures and technology, such as the increased adoption of the WaldoAir system by the Civil Air Patrol. This characteristic makes the dataset a valuable benchmark for evaluating domain adaptation techniques in the context of disaster response.

The pretrained classifiers demonstrate strong performance on the LADI v2 test set, outperform state of the art open source open-vocabulary classification from LLaVA-NeXT, and are competitive with commercial offerings such as GPT-4o on the most relevant damage classes. This comparison highlights the value of open source domain-specific training data for specialized applications and underscores the potential impact of the LADI v2 dataset on the broader machine learning community.

*Limitations and potential improvements.* While LADI v2 represents a step forward in the availability of high-quality annotated disaster imagery datasets, there are some limitations to consider. First, certain hazard types and geographies are under-represented due to multiple factors, including the likelihood of those hazard types occurring, the likelihood that a federal disaster declaration is issued, and whether CAP is tasked to collect images. For example, imagery from California is comparatively underrepresented in the dataset relative to the state’s size, population, and exposure to hazards because its state-level emergency management agency is relatively well-equipped, reducing the likelihood of FEMA-supported CAP missions in the area.

Second, LADI v2 is specific to the United States, which means that the architecture, biomes, disaster types, and infrastructure represented in the dataset is not representative of the rest of the world. Researchers and practitioners working on disaster response applications outside of the US should augment the dataset with additional imagery from the application domain.

Finally, LADI v2 currently supports only multi-label classification tasks. Applications requiring finer-grained localization or segmentation may require additional effort to adapt the dataset or models to their specific needs. The alignment of building damage labels to existing datasets

Class	Test			Validation			Difference (Test - Validation)		
	Ours	GPT-4o	LLaVA	Ours	GPT-4o	LLaVA	Ours	GPT-4o	LLaVA
bridges_any	0.53	0.52	<b>0.56</b>	<b>0.65</b>	0.59	0.44	-0.12	-0.07	0.12
buildings_any	<b>0.94</b>	<b>0.94</b>	0.88	<b>0.96</b>	0.93	0.94	-0.02	0.01	-0.06
buildings_affected_or_greater	0.50	<b>0.66</b>	0.45	<b>0.76</b>	0.74	0.56	-0.26	-0.08	-0.11
buildings_minor_or_greater	<b>0.59</b>	0.55	0.21	<b>0.68</b>	0.50	0.38	-0.09	0.05	-0.17
debris_any	<b>0.46</b>	0.40	0.40	<b>0.66</b>	0.55	0.46	-0.20	-0.15	-0.06
flooding_any	0.53	<b>0.61</b>	0.35	<b>0.79</b>	0.73	0.50	-0.26	-0.12	-0.15
flooding_structures	<b>0.43</b>	0.39	0.11	<b>0.78</b>	0.72	0.28	-0.35	-0.33	-0.17
roads_any	0.90	<b>0.92</b>	<b>0.92</b>	<b>0.94</b>	0.90	0.92	-0.04	0.02	0.00
roads_damage	0.17	<b>0.18</b>	0.07	<b>0.44</b>	0.18	0.18	-0.27	0.00	-0.11
trees_any	0.93	<b>0.97</b>	0.95	<b>0.96</b>	<b>0.96</b>	<b>0.96</b>	-0.03	0.01	-0.01
trees_damage	0.45	<b>0.50</b>	0.26	<b>0.75</b>	0.54	0.55	-0.30	-0.04	-0.29
water_any	0.79	<b>0.84</b>	0.82	0.91	<b>0.93</b>	0.89	-0.12	-0.09	-0.07

Table 5. Comparison of  $F_1$  scores of our method, LLaVA-NeXT, and GPT-4o on the test and validation sets across classes, with differences between test and validation scores.

such as those in Gupta et al. [12], Manzini et al. [23], Rahmoonfar et al. [31] offers a possibility of combining LADI v2 with those datasets in multi-task contexts.

Nevertheless, LADI v2 offers a valuable resource for the machine learning community and has the potential to support the development of more effective and efficient disaster response tools. Future work could expand the dataset to include a wider range of geographies and disaster types, and provide additional annotation types for other vision tasks.

## References

- [1] George Awad, Asad A Butt, Keith Curtis, Jonathan Fiscus, Afzal Godil, Yooyoung Lee, Andrew Delgado, Jesse Zhang, Eliot Godard, Baptiste Chocot, Lukas Diduch, Jeffrey Liu, Alan F Smeaton, Yvette Graham, Gareth J F Jones, Wessel Kraaij, and Georges Quenot. TRECVID 2020: A comprehensive campaign for evaluating video retrieval tasks across multiple application domains. *arXiv*, 2021. 2
- [2] George Awad, Asad A Butt, Keith Curtis, Jonathan Fiscus, Afzal Godil, Yooyoung Lee, Andrew Delgado, Jesse Zhang, Eliot Godard, Baptiste Chocot, et al. Evaluating multiple video understanding and retrieval tasks at trecvid 2021. In *2021 TREC Video Retrieval Evaluation*, 2021.
- [3] George Awad, Keith Curtis, Asad Butt, Jonathan Fiscus, Afzal Godil, Yooyoung Lee, Andrew Delgado, Eliot Godard, Lukas Diduch, Jeffrey Liu, Yvette Graham, and Georges Quenot. An overview on the evaluated video retrieval tasks at TRECVID 2022. *arXiv*, 2023. 2
- [4] Lucas Beyer, Xiaohua Zhai, Amlie Royer, Larisa Markeeva, Rohan Anil, and Alexander Kolesnikov. Knowledge distillation: A good teacher is patient and consistent, 2022. 5, 6
- [5] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, 2023. Available: <https://lmsys.org/blog/2023-03-30-vicuna/>. 2
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 5
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv*, 2018. 2
- [8] FEMA. Preliminary Damage Assessment Guide, 2021. Available: [https://www.fema.gov/sites/default/files/documents/fema\\_2021-pda-guide.pdf](https://www.fema.gov/sites/default/files/documents/fema_2021-pda-guide.pdf). Accessed: May 7, 2024. 3, 6
- [9] FEMA. FEMA Civil Air Patrol Imagery S3 Bucket, 2024. Available: <http://fema-cap-imagery.s3-website-us-east-1.amazonaws.com/>. 2, 4
- [10] FEMA. OpenFEMA Disaster Declarations API, 2024. Available: <https://www.fema.gov/openfema-data-page/disaster-declarations-summaries-v2>. 2
- [11] Zhe Gan, Linjie Li, Chunyuan Li, Lijuan Wang, Zicheng Liu, and Jianfeng Gao. Vision-Language Pre-Training: Basics, Recent Advances, and Future Trends. *Foundations and Trends in Computer Graphics and Vision*, 14(34):163–352, 2022. 2
- [12] Ritwik Gupta, Richard Hosfelt, Sandra Sajeev, Nirav Patel, Bryce Goodman, Jigar Doshi, Eric Heim, Howie Choset, and Matthew Gaston. xBD: A Dataset for Assessing Building Damage from Satellite Imagery. *arXiv*, 2019. 1, 2, 4, 8
- [13] Greg Holland and Cindy L. Bruyere. Recent intense hurricane response to global climate change. *Climate Dynamics*, 42(3-4):617–627, 2014. 4
- [14] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Llio Renard Lavaud, Marie-Anne

- Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothe Lacroix, and William El Sayed. Mistral 7b, 2023. 2
- [15] R Sari Kovats, Menno J Bouma, Shakoor Hajat, Eve Worrall, and Andy Haines. El niño and health. *The Lancet*, 362 (9394):1481–1489, 2003. 4
- [16] Christos Kyrkou and Theodoris Theodoridis. 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. 1, 2
- [17] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved Baselines with Visual Instruction Tuning. *arXiv*, 2023. 2
- [18] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual Instruction Tuning. *arXiv*, 2023. 2
- [19] Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, 2024. Available: <https://llava-vl.github.io/blog/2024-01-30-llava-next/>. 2, 6
- [20] Jeffrey Liu, David Strohschein, Siddharth Samsi, and Andrew Weinert. Large scale organization and inference of an imagery dataset for public safety. In *2019 IEEE High Performance Extreme Computing Conference (HPEC)*, pages 1–6, 2019. 1, 2, 4, 5
- [21] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, Furu Wei, and Baining Guo. Swin Transformer V2: Scaling Up Capacity and Resolution. *arXiv*, 2021. 5
- [22] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. (arXiv:1711.05101), 2019. Available: <http://arxiv.org/abs/1711.05101>. arXiv:1711.05101 [cs, math]. 5
- [23] Thomas Manzini, Priyankari Perali, Raisa Karnik, and Robin Murphy. Crasar-u-droids: A large scale benchmark dataset for building alignment and damage assessment in georectified suas imagery, 2024. Available: <https://arxiv.org/abs/2407.17673>. 2, 4, 8
- [24] MIT Lincoln Lab. MIT Lincoln Lab on Hugging Face, 2024. Available: <https://huggingface.co/MITLL>. 1
- [25] MIT Lincoln Lab. LADI Overview, 2024. Available: <https://github.com/LADI-Dataset/ladi-overview>. 1
- [26] MIT Lincoln Laboratory. LADI-v2-dataset (revision 3def176), 2024. Available: <https://huggingface.co/datasets/MITLL/LADI-v2-dataset>. 1
- [27] OpenAI. Hello GPT-4o, 2024. Available: <https://openai.com/index/hello-gpt-4o/>. 2
- [28] Civil Air Patrol. CAP Mission Aircrew Waldo System Introduction, 2023. Available: [https://nesa.cap.gov/media/cms/CAP\\_WaldoAir\\_System\\_Introduction\\_S1\\_F2C29DF7C84EC.pdf](https://nesa.cap.gov/media/cms/CAP_WaldoAir_System_Introduction_S1_F2C29DF7C84EC.pdf). 4
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 2
- [30] Maryam Rahnemoonfar, Tashnim Chowdhury, Argho Sarkar, Debvrat Varshney, Masoud Yari, and Robin Robertson Murphy. Floodnet: A high resolution aerial imagery dataset for post flood scene understanding. *IEEE Access*, 9:89644–89654, 2021. 1, 2
- [31] Maryam Rahnemoonfar, Tashnim Chowdhury, and Robin Robertson Murphy. RescueNet: A High Resolution UAV Semantic Segmentation Benchmark Dataset for Natural Disaster Damage Assessment. *arXiv*, 2022. 1, 2, 4, 8
- [32] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lih Zelnik-Manor. ImageNet-21K Pretraining for the Masses. *arXiv*, 2021. 5
- [33] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. (arXiv:1804.04235), 2018. Available: <http://arxiv.org/abs/1804.04235>. arXiv:1804.04235 [cs, stat]. 5
- [34] Hao Tan and Mohit Bansal. LXMERT: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5100–5111, Hong Kong, China, 2019. Association for Computational Linguistics. Available: <https://aclanthology.org/D19-1514>. 2
- [35] Ethan Weber, Dim P. Papadopoulos, Agata Lapedriza, Ferda Ofli, Muhammad Imran, and Antonio Torralba. Incidents1M: A Large-Scale Dataset of Images With Natural Disasters, Damage, and Incidents. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4768–4781, 2023. 1, 2