

NADBenchmarks - a compilation of Benchmark Datasets for Machine Learning Tasks related to Natural Disasters

Adiba Mahbub Proma, Md Saiful Islam, Stela Ciko, Raiyan Abdul Baten, Ehsan Hoque

University of Rochester, United States

Abstract

Climate change has increased the intensity, frequency, and duration of extreme weather events and natural disasters across the world. While the increased data on natural disasters improves the scope of machine learning (ML) in this field, progress is relatively slow. One bottleneck is the lack of benchmark datasets that would allow ML researchers to quantify their progress against a standard metric. The objective of this short paper is to explore the state of benchmark datasets for ML tasks related to natural disasters, categorizing them according to the disaster management cycle. We compile a list of existing benchmark datasets introduced in the past five years. We propose a web platform - NADBenchmarks - where researchers can search for benchmark datasets for natural disasters, and we develop a preliminary version of such a platform using our compiled list. This paper is intended to aid researchers in finding benchmark datasets to train their ML models on, and provide general directions for topics where they can contribute new benchmark datasets.

1 Introduction

With climate change exacerbating the intensity, frequency, and duration of extreme weather events and natural disasters across the world (Pörtner et al. 2022), rapid advancement is required in its management. Researchers are taking advantage of the huge amount of data available, applying ML for more effective solutions in managing natural disasters (Dwarakanath et al. 2021; Yu, Yang, and Li 2018). For example, the NOAA website¹ contains carefully curated climate change related data that researchers can use. However, progress is relatively slow. In fields of computer vision and natural language processing, benchmark datasets have played a crucial role in their recent rapid advancement and we believe that developing ML algorithms for natural disaster management can benefit from such rigor. Benchmark datasets are preprocessed, curated datasets for training and testing ML algorithms. These datasets provide scope for standard evaluation, allowing ML communities to quantify their progress and compare their models against each other. Existing review papers focus on the role of big data and machine learning to tackle natural disasters (Dwarakanath

et al. 2021; Yu, Yang, and Li 2018; Resch, Usländer, and Havas 2018), but to our knowledge, do not cover the scope of benchmark datasets for natural disasters, and their potential in accelerating research in this domain.

The objective of this short paper is to explore the state of benchmark datasets for ML tasks related to natural disasters, categorizing the datasets according to the disaster management cycle, which consists of four stages - mitigation, preparedness, responses and recovery (Khan et al. 2008). The **mitigation** phase deals with making long term plans to reduce the effects of a disaster; the **preparedness** phase focuses on plans for responding to a disaster; **response** activities include damage assessment and providing post-disaster coordination; and the **recovery** phase relates to recovering from the damages. For this task, we compile a list of existing benchmark datasets introduced in the past five years, find current gaps in literature, and discuss their implications. To facilitate research in this domain, we propose a web platform - NADBenchmarks - where researchers can search for benchmark datasets related to the topic. Our goal is to increase accessibility for researchers, instead of generating a leaderboard platform since leaderboards often pose the risk of simplifying progress to a singular metric (Ethayarajh and Jurafsky 2020). Our preliminary version of such a platform developed using our compiled list can be found in this link².

2 Method

Search Criteria: We searched different combinations of the keywords ‘ML datasets’, ‘natural disaster’, and ‘benchmark’ in Google Scholar, ACM digital library and Scopus. To capture the latest works, we collected additional papers from popular ML conferences CVPR, NeurIPS, ICML and ICLR; and the website climatechange.ai (Rolnick et al. 2022). We especially targeted papers that introduced datasets for a particular task. For this short review, we focused on papers from the past five years (2017-2022) since we wanted to capture the most recent works as a representation of the current state of the art for this domain.

Data Extraction and Curation: 11 criteria related to benchmark dataset characteristics are selected and data is extracted from the papers accordingly. These criteria are - dataset name, application, ML task, natural disaster topic,

Phases in Disaster Cycle	Examples of Benchmark Applications	Type of ML Tasks
Prevention - Forecasting and prediction	Drought forecasting	Multiclass ordinal classification
	Earth surface forecasting; Extreme summer prediction; seasonal cycle prediction	Video prediction
Preparedness - Early warning	Wildfire spread prediction	Image segmentation
Preparedness - Monitoring and detection	Ground deformation detection	Binary classification
	Real-time wildfire smoke detection	Binary classification
	Disaster detection	Binary classification, multiclass classification, multiclass multilabel classification
	Flood detection	Binary classification, image segmentation
Response - Damage assessment	Damage severity assessment	Multiclass ordinal classification, Semantic segmentation, multilabel classification, multitask learning
	Damage assessment of buildings	Image segmentation, multiclass ordinal classification
	Post Flood Scene Understanding	Image classification, semantic segmentation, visual question answering
Response - Post-disaster Coordination and Response	Assessing informativeness	Binary classification, multitask learning
	Categorization of humanitarian tasks	Multiclass classification, multitask learning
Recovery	Sentiment analysis	Multiclass multilabel classification

Table 1: This table summarizes some common ML applications where benchmark datasets have been introduced for natural disasters. A full list can be found in our webpage.

phase (and subphase) in disaster management cycle, timespan, geological coverage, dataset type, size, and data source. Additionally, the paper title, venue and the year published are also extracted. For space constraints, the raw data is not provided in the paper but Table 1 has summary examples and a full list is available on the webpage. Implementation details for the webpage are provided in section 5.

3 Review of Benchmark Datasets

3.1 Prevention/Mitigation

Prevention refers to long term planning on how to reduce the risk of natural disasters. This phase can be broken down into two sub-phases - Long-term Risk Assessment and Reduction, which refers to analyzing risk and taking steps to mitigate them; and Forecasting and Predicting, which focuses on methods to predict natural disasters. Two of our reviewed datasets focused on forecasting and prediction problems. EarthNet2021 was presented as a challenge for Earth surface forecasting, extreme weather prediction and seasonal cycle prediction using satellite imagery (Requena-Mesa et al. 2020). DroughtED was introduced for drought prediction, classifying drought into six categories (from no drought to exceptional) using drought and meteorological observations, and spatio-temporal data (Minixhofer et al. 2021).

3.2 Preparedness

The goal of the preparedness phase is to plan how to respond to a natural disaster, including detecting its progression, and warning citizens. So, it can be categorized into two sub-phases - Monitoring and Detection, and Early Warning. There is only one dataset on early warning - Next Day Wildfire Spread dataset - which is labelled to predict how far the wildfire would spread, given previous images (Huot et al.

2021). Most papers focus on monitoring and detection of disasters from images, but with varying complexity. Three papers introduce benchmarks for detecting types of disasters, and thus defining it as a single-label multiclass classification problem (Alam et al. 2020; Said et al. 2021; Weber et al. 2020). Recent work shows further improvement, introducing datasets for multiclass multilabel learning - Incidents1M builds on the incidents dataset by adding more labels to the images in the incidents dataset (Weber et al. 2022), and the MEDIC dataset is introduced as an extension of Alam et al. (2020)’s work (Alam et al. 2021a). The community has also started exploring the potential of video datasets - VIDI contains 4,534 video clips with 4,767 labels and provides a baseline for multilabel disaster detection (Sesver et al. 2022).

One of the most common topics in detecting specific natural disasters is flood monitoring and detection from Earth observation data such as satellites and SARs. MediaEval introduced flood-related challenges for three consecutive years (2017-2020), starting with flood detection, flooded road detection, and flood severity detection (Bischke et al. 2017, 2018, 2019). Similarly, Spacenet 8 introduced a benchmark challenge for flood mapping on road segments and buildings (Hansch et al. 2022). A benchmark for flood extent detection is introduced by Gahlot et al. (2022), where special focus is given to distinguishing flood and general water bodies. Datasets are also introduced for real-time wildfire smoke detection (Dewangan et al. 2022) and volcanic stage classification (Bountos et al. 2022).

3.3 Response

Response deals with tackling the immediate aftermath of the disaster. This includes assessing and estimating the damage caused, and taking steps to aid those affected by the disas-

ter. This phase can be divided into two subphases - Damage Assessment and Post-disaster Coordination and Response.

Multiple papers have introduced benchmarks for damage classification using social media (Zhu, Liang, and Hauptmann 2021; Mouzannar, Rizk, and Awad 2018), Earth observation data (Gupta et al. 2019; Rahneemoonfar et al. 2021; Chowdhury, Murphy, and Rahneemoonfar 2022; Chen et al. 2018) or climate-simulated data (Kashinath et al. 2019) with images or in multimodal format (text and images) (Nguyen et al. 2017; Mouzannar, Rizk, and Awad 2018). xBD also introduces the Joint Damage Scale as a unified scale for damage assessment of all natural disasters through satellite imagery (Gupta et al. 2019). Some work has also been done in assessing building damage post Hurricane using satellites (Chen et al. 2018) and aerial videos (Zhu, Liang, and Hauptmann 2021); and for flood scene understanding (Rahneemoonfar et al. 2021). A more comprehensive annotation is provided by RescueNet, consisting of different damage levels for 11 different categories, including debris (Chowdhury, Murphy, and Rahneemoonfar 2022). During post-disaster coordination, the limited resources must be allocated properly for rescue and relief operations. Benchmarks created over the years deal with assessing the informativeness of data from social media, and categorizing the type of humanitarian aid required. One of the earlier benchmarks is the CrisisMMD dataset (Alam, Ofli, and Imran 2018), where tweets from seven natural disasters during 2017 were annotated for informativeness as a binary classification problem, and across eight humanitarian categories as a multiclass classification problem. HumAID builds on concepts from CrisisMMD, introducing a larger dataset for humanitarian aid classification, labelling tweets for 19 natural disasters across ten categories (Alam et al. 2021b).

Publicly available datasets were combined to increase size, and both damage severity and humanitarian aid classification were introduced for multilabel classification (Alam et al. 2020), and then extended to CrisisBench (Alam et al. 2021c). Building upon CrisisBench, the MEDIC dataset provides a benchmark for multitask learning for damage severity and humanitarian aid classification (Alam et al. 2021a).

3.4 Recovery

The last stage is the recovery phase which focuses on aiding people to get their lives back to normalcy. We did not find a paper that was directly related, but we could loosely classify image-sentiment dataset into this category (Hassan et al. 2022). Crowdworkers label images of natural disasters according to their sentiments on seeing the image, creating a benchmark for multilabel sentiment classification for natural disasters. In the future, this kind of benchmark can be useful for analyzing the long-term psychological effects on victims of natural disasters.

4 Trends and research gaps

More datasets are needed for Prevention and Recovery Phase. According to our review, most benchmark datasets have been introduced for the preparedness and the response phase, as shown in figure 1. Only two papers focused on

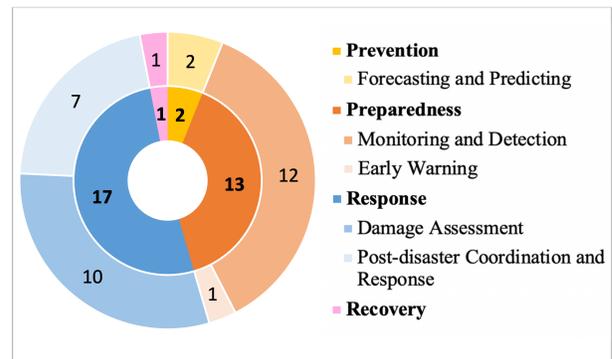


Figure 1: Pie-chart showing the number of papers found for each of the disaster management phases and subphases.

prevention, and we could loosely categorize one paper for recovery. However, the lack of benchmark datasets must not be confused with the lack of applications of ML. Multiple works have been published on risk assessment during floods (Skakun et al. 2014), earthquakes (Wilson et al. 2016; Ehrlich and Tenerelli 2013) and hurricanes (Ehrlich and Tenerelli 2013). Similarly, ML has been used for change detection in the recovery phase and also to predict the economic consequences of such outlier events (Gurrapu et al. 2021). Various data sources have been used so far, including satellites, crowdsourced data, financial records and call logs (Skakun et al. 2014; Wilson et al. 2016; Ehrlich and Tenerelli 2013; Gurrapu et al. 2021). These examples show that there is utility in generating benchmark datasets for these phases. Some potential ideas for benchmark datasets include datasets for risk assessment, evaluating or predicting recovery process, predicting displacement patterns post-disaster, and so on.

More diversification is needed in data type and data sources. Currently, there is disproportionate use of social media and Earth observation data (satellites, UAV and drones) as data sources. This can be attributed to their increased availability, accessibility and their relative inexpensiveness. Many people turn to social media for help during a crisis, thus contributing to a large pool of resources, ideal for data mining. Organizations also make it easy for researchers to use their data. Twitter, for example, has an API that allows researchers to access tweets through querying hashtags. Moreover, remote sensing makes it easy to collect images before, during and after a natural disaster without putting further lives at risk. While these are very important sources, they still have their limitations. For instance, occlusion due to clouds or smoke is a common issue for satellite imagery. This is especially exacerbated by the fact that in case of some natural disasters such as hurricanes or wildfire, clouds and smoke are inevitable. UAVs and drones can often miss specific angles, and thus miss crucial information. Moreover, although social media provides a large resource, they are often noisy and require a lot of preprocessing.

Most benchmarks reviewed in the paper are in image format. In fact, social media and Earth observation data are usually in image format and considering the rapid improvement

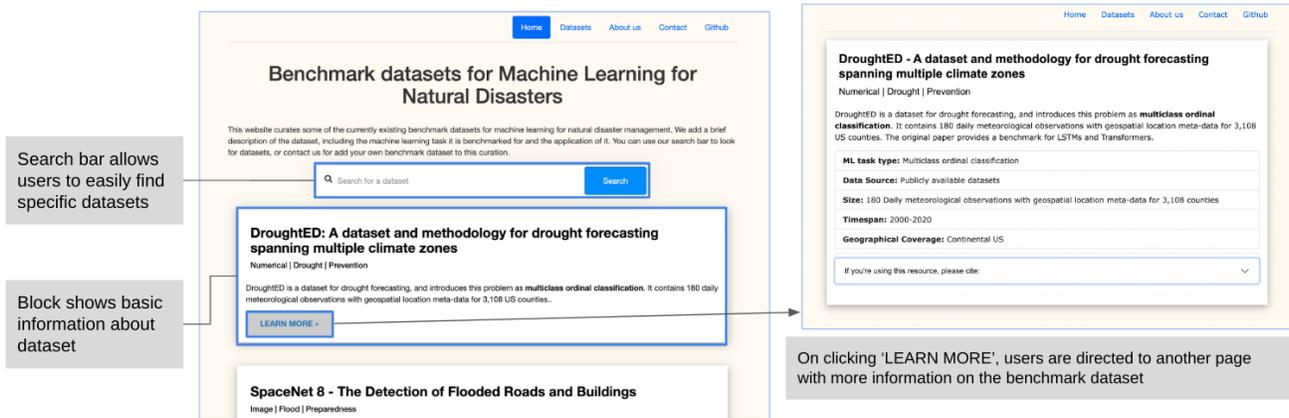


Figure 2: Our interface with key features highlighted.

in machine vision algorithms, it makes sense that most problems in this domain are presented as vision problems. However, this increases the research communities’ risk of running into limitations and missing out on the advantages of other data types and sources. Commercial data sources such as cell networks, call logs, and financial records can be good resources for damage assessment tasks and mobility prediction (Smallwood, Lefebvre, and Bengtsson 2022). Multimodal approaches that include audio and text data could improve the performance of existing algorithms.

Benchmarks for multitask learning problems. Multitask learning (MTL) has been a promising method for achieving generalizability and improving efficiency, and researchers are just starting to explore the suitability of MTL for climate models (Gonçalves, Von Zuben, and Banerjee 2015). MTL models work under the assumption that the tasks are related to one another. However, trying to build a unified approach for unrelated tasks can be detrimental to performance and there is ongoing research on what tasks learn better together (Standley et al. 2020). Some papers have started discussing the scope of building MTL benchmark datasets for this domain. So far, we were able to find only one paper on benchmark for MTL (Alam et al. 2021a), but we can expect more in the future.

5 Interface implementation

Taking inspiration from current benchmark data curation websites such as CrisisNLP (datasets for crisis management), GEM (datasets for natural language generation), and paperswithcode (datasets for vision tasks), we are building a web platform to increase data accessibility for researchers in this domain. Currently, our working prototype consists of information on the benchmark datasets reviewed in this paper. Users can scroll through our list of datasets, each of which consists of a short description explaining the ML task and the topic of the benchmark. Interested researchers can click the ‘Learn More’ button for more information, as shown in figure 2. Currently, we display the data extracted using the 11 criteria. In the future, we aim to add more features, including information on annotation methods, ML models the

benchmark dataset has been tested on, labels, data distribution and accuracy metrics.

6 Discussion and conclusion

In this paper, we briefly review existing benchmark datasets for natural disasters, categorizing them according to the disaster management cycle. We focused on the applications and ML tasks introduced by the original paper because our goal was to facilitate the task of searching for datasets for interdisciplinary researchers.

However, there are other characteristics of benchmark datasets that are also of equal importance - evaluation metrics, annotation methods, and data distribution (Bender and Friedman 2018; Olson et al. 2017). This is reflected in current vision and NLP benchmark websites - GEM aims at improving evaluation strategies for the NLG community, and paperswithcode ranks papers for vision tasks in terms of their accuracy (Gehrmann et al. 2021; Robert et al. 2018). A comprehensive analysis of these characteristics were not conducted for our review but can be included in the future. Further research can be done to determine the utility of leaderboards for this domain, and more holistic metrics to evaluate models can be proposed. For example, “data cards” have been introduced in NLP as a more holistic metric which takes bias into account (Bender and Friedman 2018). Future directions can also focus on generating such data cards for ML for environmental science-related tasks. Despite the limitations, there is scope for inclusion of such information on our webpage. Moreover, as new datasets are introduced, we plan to closely monitor and update our website regularly. We hope the community would contribute to this curation by submitting their datasets. Our goal for future work involves creating a more comprehensive list of benchmark datasets, increasing the scope of information for the datasets, and providing general analytics to inform researchers about the current state of the art.

Acknowledgement

This research was supported by funding from the Goergen Institute for Data Science.

References

- Alam, F.; Alam, T.; Hasan, M.; Hasnat, A.; Imran, M.; Ofli, F.; et al. 2021a. MEDIC: a multi-task learning dataset for disaster image classification. *arXiv preprint arXiv:2108.12828*.
- Alam, F.; Ofli, F.; and Imran, M. 2018. Crisismmd: Multimodal twitter datasets from natural disasters. In *Twelfth International AAAI Conference on Web and Social Media*.
- Alam, F.; Ofli, F.; Imran, M.; Alam, T.; and Qazi, U. 2020. 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*, 151–158. IEEE.
- Alam, F.; Qazi, U.; Imran, M.; and Ofli, F. 2021b. HumAID: Human-Annotated Disaster Incidents Data from Twitter with Deep Learning Benchmarks. In *International Conference on Web and Social Media*, 933–942.
- Alam, F.; Sajjad, H.; Imran, M.; and Ofli, F. 2021c. CrisisBench: Benchmarking Crisis-related Social Media Datasets for Humanitarian Information Processing. In *International Conference on Web and Social Media*, 923–932.
- Bender, E. M.; and Friedman, B. 2018. Data Statements for Natural Language Processing: Toward Mitigating System Bias and Enabling Better Science. *Transactions of the Association for Computational Linguistics*, 6: 587–604.
- Bischke, B.; Helber, P.; Brugman, S.; Basar, E.; Zhao, Z.; Larson, M. A.; and Pogorelov, K. 2019. The multimedia satellite task at MediaEval 2019. In *MediaEval*.
- Bischke, B.; Helber, P.; Schulze, C.; Srinivasan, V.; Dengel, A.; and Borth, D. 2017. The Multimedia Satellite Task at MediaEval 2017. In *MediaEval*.
- Bischke, B.; Helber, P.; Zhao, Z.; de Bruijn, J.; and Borth, D. 2018. The Multimedia Satellite Task at MediaEval 2018 Emergency Response for Flooding Events.
- Bountos, N. I.; Papoutsis, I.; Michail, D.; Karavias, A.; Elias, P.; and Parcharidis, I. 2022. Hephaestus: A large scale multitask dataset towards InSAR understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1453–1462.
- Chen, S. A.; Escay, A.; Haberland, C.; Schneider, T.; Staneva, V.; and Choe, Y. 2018. Benchmark dataset for automatic damaged building detection from post-hurricane remotely sensed imagery. *arXiv preprint arXiv:1812.05581*.
- Chowdhury, T.; Murphy, R.; and Rahneemoonfar, M. 2022. RescueNet: A High Resolution UAV Semantic Segmentation Benchmark Dataset for Natural Disaster Damage Assessment. *arXiv preprint arXiv:2202.12361*.
- Dewangan, A.; Pande, Y.; Braun, H.-W.; Vernon, F.; Perez, I.; Altintas, I.; Cottrell, G. W.; and Nguyen, M. H. 2022. FlgLib & SmokeyNet: Dataset and Deep Learning Model for Real-Time Wildland Fire Smoke Detection. *Remote Sensing*, 14(4): 1007.
- Dwarakanath, L.; Kamsin, A.; Rasheed, R. A.; Anandhan, A.; and Shuib, L. 2021. Automated Machine Learning Approaches for Emergency Response and Coordination via Social Media in the Aftermath of a Disaster: A Review. *IEEE Access*, 9: 68917–68931.
- Ehrlich, D.; and Tenerelli, P. 2013. Optical satellite imagery for quantifying spatio-temporal dimension of physical exposure in disaster risk assessments. *Natural Hazards*, 68(3): 1271–1289.
- Ethayarajh, K.; and Jurafsky, D. 2020. Utility is in the Eye of the User: A Critique of NLP Leaderboards. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4846–4853. Online: Association for Computational Linguistics.
- Gahlot, S.; Ramasubramanian, M.; Gurung, I.; Hansch, R.; Molthan, A.; and Maskey, M. 2022. Curating flood extent data and leveraging citizen science for benchmarking machine learning solutions. *Earth and Space Science Open Archive*, 9.
- Gehrmann, S.; Adewumi, T.; Aggarwal, K.; Ammanamanchi, P. S.; Aremu, A.; Bosselut, A.; Chandu, K. R.; Clinciu, M.-A.; Das, D.; Dhole, K.; Du, W.; Durmus, E.; Dušek, O.; Emezue, C. C.; Gangal, V.; Garbacea, C.; Hashimoto, T.; Hou, Y.; Jernite, Y.; Jhamtani, H.; Ji, Y.; Jolly, S.; Kale, M.; Kumar, D.; Ladhak, F.; Madaan, A.; Maddela, M.; Mahajan, K.; Mahamood, S.; Majumder, B. P.; Martins, P. H.; McMillan-Major, A.; Mille, S.; van Miltenburg, E.; Nadeem, M.; Narayan, S.; Nikolaev, V.; Niyongabo Rubungo, A.; Osei, S.; Parikh, A.; Perez-Beltrachini, L.; Rao, N. R.; Rاونak, V.; Rodriguez, J. D.; Santhanam, S.; Sedoc, J.; Sellam, T.; Shaikh, S.; Shimorina, A.; Sobrevilla Cabezudo, M. A.; Strobelt, H.; Subramani, N.; Xu, W.; Yang, D.; Yerukola, A.; and Zhou, J. 2021. The GEM Benchmark: Natural Language Generation, its Evaluation and Metrics. In *Proceedings of the 1st Workshop on Natural Language Generation, Evaluation, and Metrics (GEM 2021)*, 96–120. Online: Association for Computational Linguistics.
- Gonçalves, A. R.; Von Zuben, F. J.; and Banerjee, A. 2015. A multitask learning view on the earth system model ensemble. *Computing in Science & Engineering*, 17(6): 35–42.
- Gupta, R.; Hosfelt, R.; Sajeew, S.; Patel, N.; Goodman, B.; Doshi, J.; Heim, E.; Choset, H.; and Gaston, M. 2019. xbd: A dataset for assessing building damage from satellite imagery. *arXiv preprint arXiv:1911.09296*.
- Gurrapu, S.; Batarseh, F. A.; Wang, P.; Sikder, M. N. K.; Gorentala, N.; and Gopinath, M. 2021. DeepAg: Deep Learning Approach for Measuring the Effects of Outlier Events on Agricultural Production and Policy. In *2021 IEEE Symposium Series on Computational Intelligence (SSCI)*, 1–8. IEEE.
- Hänsch, R.; Arndt, J.; Lunga, D.; Gibb, M.; Pedelose, T.; Boedihardjo, A.; Petrie, D.; and Bacastow, T. M. 2022. SpaceNet 8-The Detection of Flooded Roads and Buildings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1472–1480.
- Hassan, S. Z.; Ahmad, K.; Hicks, S.; Halvorsen, P.; Al-Fuqaha, A.; Conci, N.; and Riegler, M. 2022. Visual sentiment analysis from disaster images in social media. *Sensors*, 22(10): 3628.
- Huot, F.; Hu, R. L.; Goyal, N.; Sankar, T.; Ihme, M.; and Chen, Y.-F. 2021. Next Day Wildfire Spread: A Ma-

- chine Learning Data Set to Predict Wildfire Spreading from Remote-Sensing Data. *arXiv preprint arXiv:2112.02447*.
- Kashinath, K.; Mudigonda, M.; Mahesh, A.; Chen, J.; Yang, K.; Greiner, A.; and Prabhat, M. 2019. ClimateNet: Bringing the power of Deep Learning to weather and climate sciences via open datasets and architectures. In *AGU Fall Meeting Abstracts*, volume 2019, GC33A–06.
- Khan, H.; Vasilescu, L. G.; Khan, A.; et al. 2008. Disaster Management Cycle—a theoretical approach. *Journal of Management and Marketing*, 6(1): 43–50.
- Minixhofer, C. D.; Swan, M.; McMeekin, C.; and Andreadis, P. 2021. DroughtED: A dataset and methodology for drought forecasting spanning multiple climate zones. In *ICML 2021 Workshop on Tackling Climate Change with Machine Learning*.
- Mouzannar, H.; Rizk, Y.; and Awad, M. 2018. Damage Identification in Social Media Posts using Multimodal Deep Learning. In *Information Systems for Crisis Response And Management*.
- Nguyen, D. T.; Ofli, F.; Imran, M.; and Mitra, P. 2017. Damage Assessment from Social Media Imagery Data during Disasters. In *Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017*, 569–576.
- Olson, R. S.; La Cava, W.; Orzechowski, P.; Urbanowicz, R. J.; and Moore, J. H. 2017. PMLB: a large benchmark suite for machine learning evaluation and comparison. *Bio-Data Mining*, 10(1): 1–13.
- Pörtner, H.-O.; Roberts, D. C.; Adams, H.; Adler, C.; Aldunce, P.; Ali, E.; Begum, R. A.; Betts, R.; Kerr, R. B.; Biesbroek, R.; et al. 2022. Climate change 2022: Impacts, adaptation and vulnerability. *IPCC Sixth Assessment Report*.
- Rahnemoonfar, M.; Chowdhury, T.; Sarkar, A.; Varshney, D.; Yari, M.; and Murphy, R. R. 2021. Floodnet: A high resolution aerial imagery dataset for post flood scene understanding. *IEEE Access*, 9: 89644–89654.
- Requena-Mesa, C.; Benson, V.; Denzler, J.; Runge, J.; and Reichstein, M. 2020. EarthNet2021: A novel large-scale dataset and challenge for forecasting localized climate impacts. *arXiv preprint arXiv:2012.06246*.
- Resch, B.; Usländer, F.; and Havas, C. 2018. Combining machine-learning topic models and spatiotemporal analysis of social media data for disaster footprint and damage assessment. *Cartography and Geographic Information Science*, 45(4): 362–376.
- Robert; Ross; Marcin; Elvis; Guillem; Andrew; and Thomas. 2018. Papers with code - the latest in machine learning.
- Rolnick, D.; Donti, P. L.; Kaack, L. H.; Kochanski, K.; Lacoste, A.; Sankaran, K.; Ross, A. S.; Milojevic-Dupont, N.; Jaques, N.; Waldman-Brown, A.; et al. 2022. Tackling climate change with machine learning. *ACM Computing Surveys (CSUR)*, 55(2): 1–96.
- Said, N.; Ahmad, K.; Conci, N.; and Al-Fuqaha, A. 2021. Active learning for event detection in support of disaster analysis applications. *Signal, Image and Video Processing*, 15(6): 1081–1088.
- Sesver, D.; Gençoğlu, A. E.; Yıldız, Ç. E.; Günindi, Z.; Habibi, F.; Yazıcı, Z. A.; and Ekenel, H. K. 2022. VIDl: A Video Dataset of Incidents. *arXiv preprint arXiv:2205.13277*.
- Skakun, S.; Kussul, N.; Shelestov, A.; and Kussul, O. 2014. Flood hazard and flood risk assessment using a time series of satellite images: A case study in Namibia. *Risk Analysis*, 34(8): 1521–1537.
- Smallwood, T. R.; Lefebvre, V.; and Bengtsson, L. 2022. Mobile phone usage data for disaster response. *Communications of the ACM*, 65(4): 40–41.
- Standley, T.; Zamir, A.; Chen, D.; Guibas, L.; Malik, J.; and Savarese, S. 2020. Which tasks should be learned together in multi-task learning? In *International Conference on Machine Learning*, 9120–9132. PMLR.
- Weber, E.; Marzo, N.; Papadopoulos, D. P.; Biswas, A.; Lapedriza, A.; Ofli, F.; Imran, M.; and Torralba, A. 2020. Detecting natural disasters, damage, and incidents in the wild. In *European Conference on Computer Vision*, 331–350. Springer.
- Weber, E.; Papadopoulos, D. P.; Lapedriza, A.; Ofli, F.; Imran, M.; and Torralba, A. 2022. Incidents1M: a large-scale dataset of images with natural disasters, damage, and incidents. *arXiv preprint arXiv:2201.04236*.
- Wilson, R.; zu Erbach-Schoenberg, E.; Albert, M.; Power, D.; Tudge, S.; Gonzalez, M.; Guthrie, S.; Chamberlain, H.; Brooks, C.; Hughes, C.; et al. 2016. Rapid and near real-time assessments of population displacement using mobile phone data following disasters: The 2015 Nepal earthquake. *PLoS Currents*, 8.
- Yu, M.; Yang, C.; and Li, Y. 2018. Big data in natural disaster management: a review. *Geosciences*, 8(5): 165.
- Zhu, X.; Liang, J.; and Hauptmann, A. 2021. Msnet: A multilevel instance segmentation network for natural disaster damage assessment in aerial videos. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023–2032.