



Automating tephra fall building damage assessment using deep learning

Eleanor Tennant¹, Susanna F. Jenkins², Victoria Miller³, Richard Robertson⁴, Bihan Wen⁵, Sang-Ho Yun², and Benoit Taisne²

¹Earth Observatory of Singapore, Interdisciplinary Graduate Programme, Nanyang Technological University, Singapore, 639798, Singapore

²Earth Observatory of Singapore, Asian School of the Environment, Nanyang Technological University, Singapore, 639798, Singapore

³GNS Science, P.O. Box 30368, 5040 Lower Hutt, Aotearoa/New Zealand

⁴The UWI Seismic Research Centre, Saint Augustine, Trinidad and Tobago

⁵School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, 639798, Singapore

Correspondence: Eleanor Tennant (eleanorm001@e.ntu.edu.sg)

Received: 7 May 2024 – Discussion started: 30 May 2024

Revised: 1 October 2024 – Accepted: 25 October 2024 – Published: 12 December 2024

Abstract. In the wake of a volcanic eruption, the rapid assessment of building damage is paramount for effective response and recovery planning. Uncrewed aerial vehicles, UAVs, offer a unique opportunity for assessing damage after a volcanic eruption, with the ability to collect on-demand imagery safely and rapidly from multiple perspectives at high resolutions. In this work, we established a UAV-appropriate tephra fall building damage state framework and used it to label $\sim 50\,000$ building bounding boxes around ~ 2000 individual buildings in 2811 optical images collected during surveys conducted after the 2021 eruption of La Soufrière volcano, St Vincent and the Grenadines. We used these labelled data to train convolutional neural networks (CNNs) for (1) building localisation (average precision equals 0.728) and (2) damage classification into two levels of granularity: no damage vs. damage (F_1 score = 0.809) and moderate damage vs. major damage (F_1 score = 0.838) (1 is the maximum obtainable for both metrics). The trained models were incorporated into a pipeline along with all the necessary image processing steps to generate spatial data (a georeferenced vector with damage state attributes) for rapid tephra fall building damage mapping. Using our pipeline, we assessed tephra fall building damage for the town of Owia, finding that 22 % of buildings that received 50–90 mm of tephra accumulation experienced at least moderate damage. The pipeline is expected to perform well across other vol-

canic islands in the Caribbean where building types are similar, though it would benefit from additional testing. Through cross-validation, we found that the UAV look angle had a minor effect on the performance of damage classification models, while for the building localisation model, the performance was affected by both the look angle and the size of the buildings in images. These observations were used to develop a set of recommendations for data collection during future UAV tephra fall building damage surveys. This is the first attempt to automate tephra fall building damage assessment solely using post-event data. We expect that incorporating additional training data from future eruptions will further refine our model and improve its applicability worldwide. To facilitate continued development and collaboration all trained models and the pipeline code can be downloaded from GitHub.

1 Introduction

Tephra fall produced by explosive volcanic eruptions can have detrimental effects on buildings, which in turn affects the ability for a community to recover and rehabilitate. These effects range from surface-level issues such as corrosion of metal roofs (e.g. Rabaul, Papua New Guinea; Blong, 2003a) or damage to non-structural components (e.g. gutters: Am-

bae, Vanuatu; Jenkins et al., 2024) through to complete building collapse (e.g. Pinatubo, Philippines; Spence et al., 1996).

After, or during, an eruption, the collection of empirical data detailing the damage incurred is critical to guide the planning and implementation of response and recovery efforts. This involves estimation of damages and losses, which are needed to determine the necessary funding for repair or reconstruction, along with an assessment of building functionality, which can inform temporary housing requirements. In addition to its use in post-disaster recovery, the collection of damage data is key to the development of vulnerability models (Deligne et al., 2022), which relate hazard intensity to damage (e.g. Spence et al., 2005; Wilson et al., 2014; Williams et al., 2020) and can be used to provide information about resilient construction practises and/or for pre-event impact assessments.

Post-event building damage assessments usually consist of ground surveys, whereby the amount of damage to each building is described using a quantitative or qualitative damage state (e.g. Spence et al., 1996; Blong, 2003a; Jenkins et al., 2013, 2015; Hayes et al., 2019; Meredith et al., 2022). However, tephra fall damage can extend tens or even hundreds of kilometres away from a volcano (Spence et al., 2005), meaning that comprehensive ground-based damage assessments can be both time-consuming and costly. Furthermore, the uncertainty that is often associated with the end of an eruption may prevent the safe completion of a ground-based damage assessment before tephra is remobilised by winds and rain. This lag between the event itself and the completion of a damage assessment can hinder recovery efforts and compromise the accuracy of data collected for the development of forecasting models.

Given the need for, but also the challenges associated with, conducting post-event building damage assessments quickly, approaches that use remotely sensed (RS) data, either optical or synthetic aperture radar (SAR) imagery, have been developed in volcanology (e.g. Jenkins et al., 2013; Williams et al., 2020; Lerner et al., 2021; Biass et al., 2021; Meredith et al., 2022) and operationally by emergency management services (e.g. International Charter “Space and Major Disasters”, Copernicus Emergency Management Service, ARIA (Advanced Rapid Imaging and Analysis) system; Yun et al., 2015). The use of optical imagery largely consists of visual inspection, which may be influenced by image resolution and is prone to subjectivity (Novikov et al., 2018). Furthermore, visual inspection of satellite optical imagery can still be time-consuming without crowd sourcing (e.g. Ghosh et al., 2011) and is constrained by satellite recurrence intervals and cloud cover. Automated SAR-based methods (e.g. Yun et al., 2015) are not limited by cloud cover, but they may lack the resolution required for building-level damage assessment (30 m for damage proxy maps generated from Sentinel data using the ARIA system; https://aria-share.jpl.nasa.gov/20210409-LaSoufriere_volcano, last access: 21 January 2024).

To our knowledge, only one study attempts to automate the assessment of building damage from volcanic hazards (Wang et al., 2024). In contrast, attention has been given to more commonly occurring hazards such as earthquakes and hurricanes, with the development of both mono-temporal (post-event imagery only) and multi-temporal (images taken at different times) approaches (Table 1). Early approaches at automation with optical imagery used image processing methods, often focusing on identifying changes in pixel values between pre- and post-event imagery (e.g. Bruzzone and Fernández Prieto, 2000; Ishii et al., 2002; Zhang et al., 2003). Image processing methods are susceptible to user biases, such as the choice of thresholds that equate to distinct levels of damage severity or damage states, and may require recalibration when applied to a new dataset. As a result, image processing methods were succeeded by the application of traditional machine learning algorithms that use “hand-crafted” image features. These features are observable properties that can be extracted from the image such as shape, colour, texture and statistical properties of the image (e.g. Li et al., 2015; Anniballe et al., 2018; Lucks et al., 2019; Naito et al., 2020). The success of a given machine learning approach is dependent on the selection of the best features for the job; for example, a texture-based feature might be good for classifying buildings as damaged or not damaged due to an increased number of edges in damaged buildings but less useful for a task such as differentiating between building roof types where the difference in textures between the classes is less significant. Deep learning, in particular the use of convolutional neural networks (CNNs), removes this need for feature selection. A CNN is a network of layers comprising filters which are small matrices of values. When an image is passed through the network, at each layer the filters are convolved with the output from the previous layer to create a new representation of the image that is progressively more abstract with depth in the network. This process reduces the image’s original spatial dimensions X and Y while increasing the number of channels, facilitating classification. During network training the filter values (known as weights) are optimised to reduce the loss between the predicted label for the image and the true label. Through this training a CNN learns the features of the images that are useful for classification. For a detailed background on deep learning, see Aggarwal (2018).

Thus far, deep learning models have been developed for optical image sets for hurricanes (Y. Li et al., 2019; Dung Cao and Choe, 2020; Pi et al., 2020; Cheng et al., 2021; Khajwal et al., 2023), earthquakes (Nex et al., 2019; Xu et al., 2019; Duarte et al., 2018; Moradi and Shah-Hosseini, 2020), wildfires (Galanis et al., 2021), volcanic hazards (Wang et al., 2024) and models that have been proposed for multiple hazards (e.g. Gupta and Shah, 2020; Weber and Kané, 2020; Shen et al., 2021; Bouchard et al., 2022) (Table 1). However, building damage caused by different hazards looks very different (e.g. damage caused by vertical loading from volcanic

Table 1. A non-exhaustive list of works using deep learning on optical imagery for building damage assessment. Studies use different scores to evaluate performance: F_1 scores are in italics, mean average precision scores are underlined, and accuracy scores are in bold. For all scores, 1 represents a perfect model. A detailed explanation of the scores used for evaluation is provided in Sect. 2.3.3.

Study	Hazard	Number of damage classes	Pre-disaster imagery	Data type	Building localisation	Damage classification
Y. Li et al. (2019)	Hurricane	2	No	Airborne		<u>0.448</u>
Weber and Kané (2020)	Multi	4	Yes	Satellite (xBD)	<i>0.835</i>	<i>0.697</i>
Dung Cao and Choe (2020)	Hurricane	2	No	Satellite	–	0.972
Pi et al. (2020)	Hurricane	2	No	UAV, airborne		<u>0.745 (UAV)</u> <u>0.807 (airborne)</u>
Cheng et al. (2021)	Hurricane	5	No	UAV	<u>0.656</u>	0.610
Galanis et al. (2021)	Wildfire	2	No	Satellite		<i>0.981</i>
Gupta and Shah (2020)	Multi	4	Yes	Satellite (xBD)	<i>0.840</i>	<i>0.740</i>
Shen et al. (2021)	Multi	4	Yes	Satellite (xBD)	<i>0.864</i>	<i>0.782</i>
Bouchard et al. (2022)	Multi	2	Yes	Satellite (xBD)	<i>0.846</i>	<i>0.709</i>
Khajwal et al. (2023)	Hurricane	5	No	Ground, airborne	–	<i>0.650</i>
Singh and Hoskere, (2023)	Multi	5	No	Satellite		0.880
Wang et al. (2024)	Volcanic tephra	4	Yes	Satellite	<i>0.868</i>	<i>0.783</i>

tephra fall vs. ground shaking from an earthquake). These observable differences mean that an optical-imagery multi-hazard damage classification model that performs consistently well across the different hazards is not yet achievable. Therefore, distinct models tailored for specific hazards are required (Nex et al., 2019; Bouchard et al., 2022). It follows that models may also benefit from being regionalised, given the differences in building typologies (construction material and styles) that can also affect the observable damage (Nex et al., 2019).

Many of the approaches for automating building damage assessment use both pre- and post-event imagery (Table 1), which makes the task more straightforward since any changes to the pre-event imagery can be considered damage. However, pre-event imagery at a high enough resolution is not always available in post-disaster scenarios. The automated assessment of building damage from volcanic haz-

ards using only post-event optical imagery has not yet been achieved in part due to the absence of the large datasets that are needed in order to train models. The 2021 eruption of La Soufrière volcano, St Vincent and the Grenadines, provided unprecedented circumstances allowing for the collection of high-resolution uncrewed aerial vehicle (UAV) imagery enabling the development of fully automated models that can assess tephra fall building damage from post-event data only. With their growing ubiquity and low cost, UAVs have become an increasingly useful tool during and after volcanic eruptions (e.g. Andaru and Rau, 2019; Gailler et al., 2021; Román et al., 2022). UAVs offer a distinct advantage over satellite imagery because they can be scheduled at any point; they do not suffer from cloud obscuring the images as they fly at relatively low altitude; and they capture imagery from multiple perspectives, which may lead to an increased ability to capture damage information. In this study we used

UAV optical imagery collected after the 2021 eruption of La Soufrière volcano to develop a methodology for tephra fall building damage assessment. The main contributions of our work are three-fold:

1. We have devised a UAV-appropriate building damage state framework, laying the foundation for future UAV tephra fall building damage surveys.
2. We have developed a deep learning pipeline that consists of all trained models and image processing steps to rapidly output spatial damage data that can facilitate prompt, post-event response and recovery and enable data collection prior to further changes by natural or human processes (tephra clean-up).
3. Imagery used in this work is diverse in terms of the flight altitude, time of acquisition after the event and UAV vantage point. We have conducted extensive testing to understand the best practises for building damage surveys and to create a series of recommendations for the collection of future UAV surveys for building damage assessment.

The 2020–2021 eruption of La Soufrière volcano in St Vincent

La Soufrière is an active stratovolcano standing at 1220 m a.s.l. (above sea level) on the island of St Vincent. On 27 December 2020 a thermal anomaly was detected inside the summit crater by the NASA Fire Information for Resource Management System (FIRMS). This was confirmed by the Soufrière Monitoring Unit to be caused by a new dome growing within the crater. Dome growth continued for 3 months until 9 April 2021, when, following 2 d of heightened seismic activity and lava effusion rate, the ongoing effusive eruption of La Soufrière entered an explosive phase (Joseph et al., 2022). Between 9–22 April, a total of 32 distinct explosions occurred, with the tallest plumes reaching heights of up to 15 km above the vent (Joseph et al., 2022). Throughout this explosive phase, tephra blanketed the island, resulting in a total deposit thickness of up to 16 cm in coastal communities to the north of the island (Cole et al., 2023) (Fig. 1).

The explosive phase was anticipated, and an evacuation order was issued on 8 April 2021 for the ~16000 residents in the northern part of the island (Joseph et al., 2022). As a result, there were no reported fatalities directly attributable to the eruption; nevertheless, the overall damage to infrastructure services and physical assets was estimated at XCD 416.07 million (equivalent to USD 153.29 million) (PDNA, 2022). Approximately 63 % of this monetary impact was borne by the housing sector. In St Vincent, residential buildings are typically single-story, detached structures, with the majority in the more impacted north of the island (census districts of Chateaubelair, Georgetown and Sandy

Bay; Fig. 1) constructed using concrete and blocks (84 % in Chateaubelair, 74 % in Georgetown, 50 % in Sandy Bay) with metal sheet roofs (90 %–92 % of all buildings in these areas) (SVG population and housing census, 2012).

2 Method

After the 2021 eruption of La Soufrière three UAV optical-imagery datasets were collected to assess the extent of the damage. These were collected by different parties at separate times after the eruption. All UAV survey locations are shown in Fig. 1, and representative examples of images can be found in Sect. S1 of the Supplement.

2.1 Dataset description

2.1.1 Dataset 1: April–May 2021 (UWI-TV)

Collected by UWI-TV, the multimedia channel of the University of the West Indies at the request of the UWI Seismic Research Centre (SRC), this dataset consists of video footage for Chateaubelair, Fitz Hughes, Troumaca and Sandy Bay acquired with a frame rate of 30 frames per second (fps) and a resolution of 1920×1080 pixels. Flight paths were not programmed, and the vantage point varies between at nadir (directly above buildings) and very off nadir (showing the sides of buildings). Images do not contain GPS positioning or altitudes and were not manually georeferenced.

2.1.2 Dataset 2: 12–14 May 2021 (GOV)

This dataset was collected by the Government of St Vincent and the Grenadines Ministry of Transport, Works, Lands and Surveys, and Physical Planning for the purpose of assessing the eruption impact. This dataset consists of video footage for Chateaubelair, London, Richmond and Sandy Bay acquired with a frame rate of 30 fps and a resolution of 1920×1080 pixels. Buildings are imaged from an at-nadir to off-nadir vantage point at ~200 m above the ground. Buildings are lower resolution in this dataset when compared to the other two. Images contain GPS positioning and altitudes.

2.1.3 Dataset 3: August–September 2021 (SRC)

This is the most extensive dataset, collected by SRC for the purpose of assessing eruption impact. It consists of photos and videos for Belmont, Chateaubelair, Fancy, London (video only), Orange Hill (video only), Owia, Point, Rabacca (video only), Richmond, Sandy Bay and Tourama. Videos were acquired with a frame rate of 30 fps and have a resolution of 1920×1080 pixels, while photos are 4056×3040 pixels. Flight paths were programmed to follow a linear swath-like trajectory. Buildings are captured from nadir between 55–290 m above the ground. Images contain GPS positioning and altitudes.

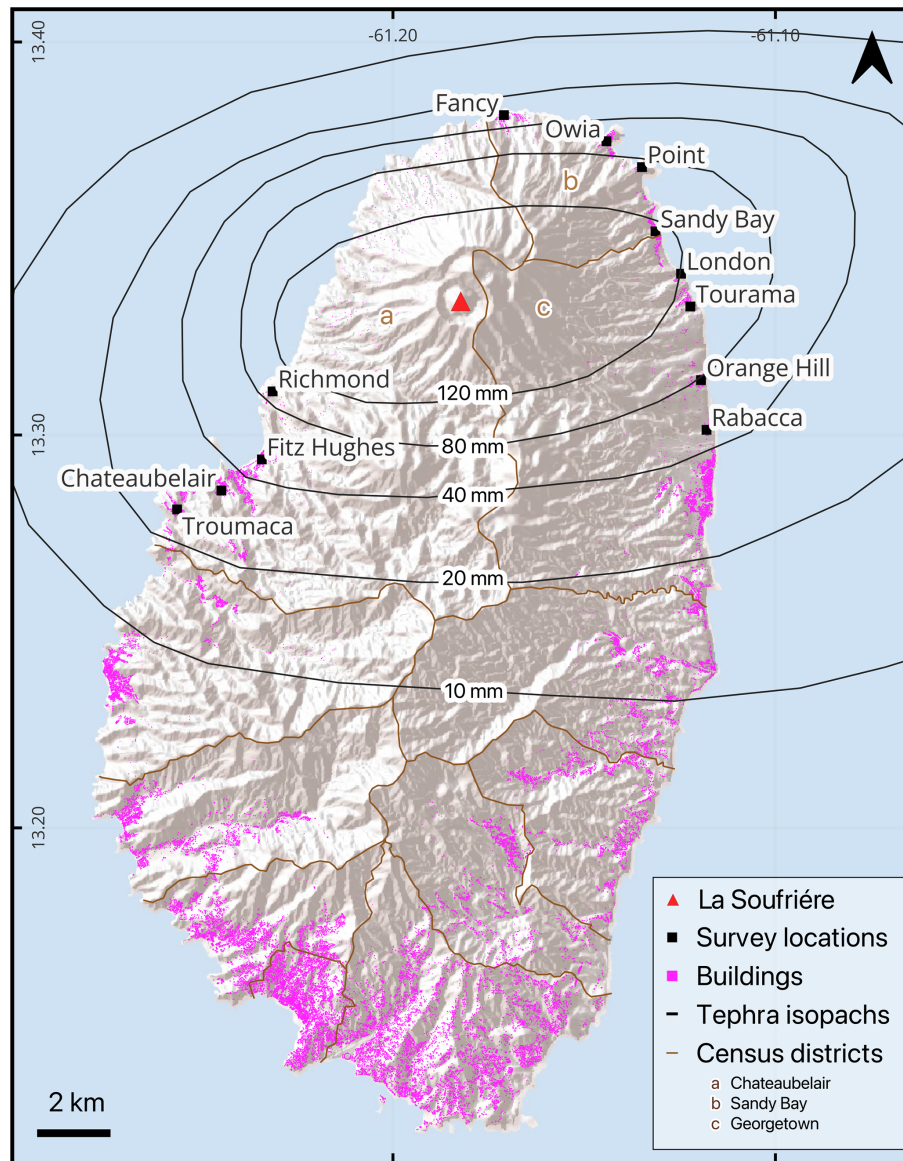


Figure 1. The island of St Vincent with UAV survey locations included in this work labelled and marked in black. Tephra isopachs (Cole et al., 2023) mark lines of constant total tephra thickness. Building footprints are marked in pink, data source: © OpenStreetMap contributors 2024. Distributed under the Open Data Commons Open Database License (ODbL) v1.0. Coordinate reference system: WGS 84 (EPSG:4326).

For all three datasets, image frames were extracted from the videos every 2 s, an interval chosen to reduce redundant homogeneous images. This resulted in a total of 7956 image frames. Due to the UAV surveying approach (i.e. hovering in one place for a while) many near-identical images were generated. To avoid potentially biasing the training towards over-represented buildings, we manually filtered out duplicate images. After filtering and the removal of images with no buildings present, the full combined dataset consisted of 2811 image frames. We labelled all images by drawing bounding boxes around each building present and storing the bounding box positions. In total 49 173 building bounding boxes were drawn around ~ 2000 individual buildings (with some

buildings being present in multiple images). Given the absence of individual building location information, this number was approximated by overlaying Open Street Map building footprints with UAV GPS tracks where available. Bounding boxes were drawn by a team of five including the lead author, and all boxes were checked by the lead author. Each box was then assigned one of three damage states, which are described below. For consistency the damage states were assigned by the lead author. All labelling, modelling and analysis were conducted using MATLAB (The MathWorks Inc., 2023).



Figure 2. Example of the three damage states used in this work: no damage to minor damage, moderate damage and major damage.

2.2 Developing and applying a building damage state framework

The first tephra fall building damage state framework was developed after the eruption of Pinatubo, Philippines, in 1991 (Spence et al., 1996) and was adapted from the macro seismic intensity scale used to evaluate seismic damage (Karnik et al., 1984). In the adapted framework damage ranges from damage state (DS) 0 (DS0) – “no damage” – through to DS5 – “complete roof collapse and severe damage to the rest of the building”. Subsequent tephra fall building damage state frameworks were modified from the work of Spence et al. (1996) with changes in the wording made to reflect the characteristics of the case study (Table 2). In the damage state descriptions, damage to three critical aspects of a building is described: the roof covering, the roof structure and the vertical structure (Blong, 2003a; Hayes et al., 2019; Jenkins et al., 2024). In our study, most images depict buildings from an at-nadir or close-to-nadir perspective, making roof damage more discernible than damage to the vertical structure. Thus, we generated a damage state framework that is based on the proportion of observable damage to the roof, as in the work of Williams et al. (2020). Our final framework, which was developed over several iterations, classifies building damage into three classes: no observable damage to minor damage, moderate damage and major damage (Table 3, Fig. 2). Damage states are deliberately generic so that the range of possible damage to the range of different building types can be captured (Blong, 2003b). Our three classes are respectively comparable to the DS0–1, DS2 and DS3–5

damage scales developed for ground surveys (Table 2). In the frameworks presented in Table 2, DS1 describes light/minor damage or superficial damage to non-structural components. In our framework we included minor damage in the no damage class since the difference between the two can be subtle and not easily discernible through remote assessment. Furthermore, buildings with minor damage are typically habitable and unlikely to require costly repairs; therefore, from a response and recovery perspective, we considered them to be better grouped with undamaged buildings. Our moderate damage class requires damage to or collapse of up to 50 % of the roof area, which closely fits with damage state 2 of Blong (2003a), Hayes et al. (2019), and Jenkins et al. (2024). The ground-based frameworks distinguish damage states 3 through 5 by increasing amounts of damage to the building walls (Table 2). However, the quantity and severity of impacted walls is not easy to differentiate in the majority of our UAV images, which show buildings from a nadir or close-to-nadir perspective. Therefore, in our framework, we grouped these states together under “major damage”.

2.3 Model development

After labelling, we split the full combined image dataset (2811 frames from the UWI-TV, GOV and SRC sets) into train, validation and test sets (Fig. 3). Given that many images lacked GPS positions, we grouped images by location to ensure independence among the sets. The partitioning was chosen to include diversity in both the image sets (UWI-TV, GOV, SRC) and the location, which affects the tephra fall thickness. We aimed for a standard data split of 80 %, 10 %

Table 2. A comparison of tephra fall building damage state frameworks available to date.

	Pinatubo, Philippines, 1991 (Spence et al., 1996)	Rabaul caldera, Papua New Guinea, 1994 (Blong, 2003a)	Calbuco, Chile, 2015 (Hayes et al., 2019)	Manaro Vuoi, Ambae island, Vanuatu, 2017–2018 (Jenkins et al., 2024)
DS0	No damage		No damage	No damage
DS1	Light roof damage: – Gutter damage. – Few tiles dislodged.	Light damage: – Damage to gutters and/or water tanks. – Cleanup required.	Minor damage to non-structural elements: – Damage to gutters. – Few tiles dislodged. – Damage to fittings, e.g. air-conditioning units and appliances. – Damage to contents. – Dents in the roof covering.	Light damage or damage to non-structural elements: – Damage to gutters. – Damage to contents. – Dents or minor slumping in roof cover.
DS2	Moderate roof damage: – Bending or excessive deflection of roof sheeting or purlins. – No damage to principal roofing supports.	Moderate damage: – Bending or excessive damage to as much as half roof sheeting and/or purlins. – Damage to roof overhangs or verandas. – Slight roof structural damage possible. – Interior requires cleaning, repainting and/or overhaul of electrical systems. – Solar heater needs replacing.	Moderate damage but vertical structure and roof supports intact: – As above. – Bending or excessive (e.g. perforation, cracking) damage (with or without collapse) to up to half of roof covering, e.g. tiles, metal sheet. – Little to no damage to principal roof supports, i.e. rafters or trusses. – Damage to roof overhangs or verandas.	Moderate damage but vertical structure and roof supports intact: – As for DS1, plus: – Bending or excessive damage (without collapse) to up to half of the roof covering. – Little or no damage to roof support trusses and rafters. – Damage to roof overhangs or verandas. – Interior requires repair.
DS3	Severe roof damage and some damage to vertical structure: – Severe damage or partial collapse of roof overhangs or verandas. – Severe deformation of main roof sheeting. – Some damage to roof supporting structure, columns, trusses.	Heavy damage: – Damage to roof structure and some damage to walls. – At least one wall damaged/misaligned. – Collapse of part of ceiling	Severe damage to the roof and supports: – As above. – Bending or excessive (e.g. perforation, cracking) damage (with or without collapse) to over half of roof covering. – Damage to any single principal roof support and some damage to walls. – Severe damage or partial collapse of roof overhangs or verandas.	Severe damage to the roof and supports: – As for DS2, plus: – Bending or excessive damage (with or without collapse) to more than half of the roof covering. – Damage to any single principal roof support and/or some damage to walls (less than half of walls affected). – Severe damage or partial collapse of roof overhangs or verandas.

Table 2. Continued.

	Pinatubo, Philippines, 1991 (Spence et al., 1996)	Rabaul caldera, Papua New Guinea, 1994 (Blong, 2003a)	Calbuco, Chile, 2015 (Hayes et al., 2019)	Manaro Vuoi, Ambae island, Vanuatu, 2017–2018 (Jenkins et al., 2024)
DS4	Partial roof collapse and moderate damage to rest of building: – Collapse of sheeting but not truss. – Partial collapse of sheeting and some truss failure. – Failure of supporting structure. – Moderate damage to other parts of building resulting from roof collapse.	Severe damage: – Roof collapse and moderate to severe damage to rest of the building. – Failure of roof trusses and supporting structure. – At least half of the external walls and/or internal walls deformed or collapsed. – For two-storey buildings, collapse of external and internal walls of upper floor. – Plumbing and other services may be damaged.	Partial or total collapse of the roof and supports: – As above. – Collapse of roof covering and any single principal roof support(s). – At least half of the external walls and/or internal walls deformed or collapsed.	Partial collapse of the roof and supports: – As for DS3, plus: – Collapse to less than half of roof covering and principal roof support(s). – At least half of external and/or internal walls deformed or collapsed.
DS5	Complete roof collapse and severe damage to the rest of the building: – Collapse of roof and supporting structure over more than 50 % of roof area. – Partition walls destroyed. – External walls destabilised.	Collapse: – Collapse of roof and supporting external walls over more than 50 % of floor area of building. – Internal walls collapsed. – Damage to floor and/or foundation. – Structure is irreparable, not salvageable, beyond economic repair.	Building collapse: – As above. – Collapse of roof, principal roof supports and/or supporting external walls over > 50 % of floor area of building.	Building collapse: – As for DS4, plus: – Collapse of roof, principal roof supports and/or supporting external walls over more than half of floor area of building.

and 10 % for training, validation and testing; however given the above constraints, this produced a split of 80, 8 and 12 (considering the number of bounding boxes and not the number of images). These datasets were used to develop our approach for building damage assessment. In line with studies shown in Table 1, we chose to split the damage assessment task into two subtasks: (i) building localisation (i.e. identification of building bounding boxes within the images) and (ii) damage classification. While it is possible to develop a model that can simultaneously locate and classify buildings with different levels of damage, model training under this approach can take significantly more time and resources to converge when compared to an approach that splits the tasks (Bouchard et al., 2022). Furthermore, decoupling the two tasks allows for greater flexibility; for example, if building

locations are already known, then only the classification can be run, speeding up the remote assessment.

In machine learning, the performance of a model and its optimal hyperparameters can be highly dependent on the characteristics of the dataset used for training, and hyperparameters that work well for one dataset may not work well for another. Therefore, it is common practice to optimise hyperparameters, model architectures and training strategies to find the configuration that performs the best for a particular problem. For building localisation and damage classification we conducted a series of independent experiments using different image pre-processing approaches, CNN architectures and combinations of hyperparameters with the aim of iterating towards the best experimental set-up (model selection: Sect. 3.1.1 and 3.2.1). Each experiment consisted of three

Table 3. The damage state framework developed for our UAV optical-imagery dataset.

Damage state	Description of the damage
No damage to minor damage	– No visible damage or – up to 10 % of the roof covering missing; and/or – no roof or structural collapse; and/or – visible damage to non-structural elements, e.g. gutters or decorative elements (fascia). – Comparable to DS0–1 (Table 2).
Moderate damage	– Up to 50 % roof area damaged (evidence of bending) or collapsed; may include light damage to vertical structure (e.g. wooden slats above windows broken). – Comparable to DS2 (Table 2).
Major damage	– More than 50 % roof area damaged or collapsed; may include damage to the vertical structure including total building collapse. – Comparable to DS3–5 (Table 2).

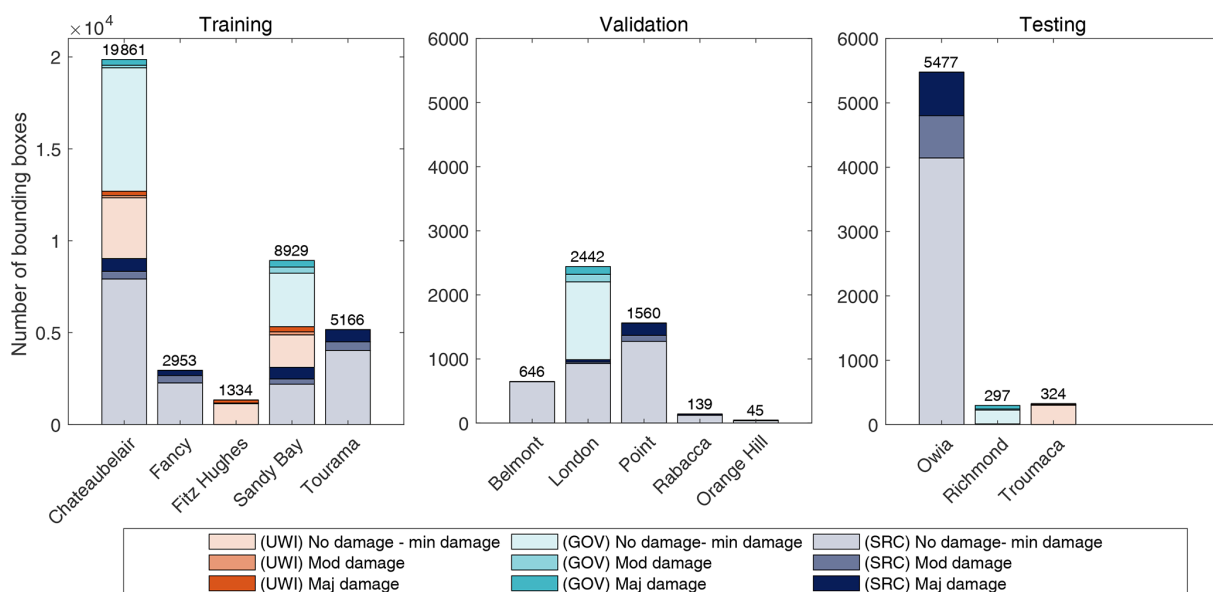


Figure 3. The number of bounding boxes of each damage state in each UAV imagery dataset (UWI-TV, GOV, SRC) for each of the locations in this study. Imagery was divided into three groups: training, validation and testing. The division of datasets between the three groups was chosen to incorporate diversity in the image sets (UWI-TV, GOV, SRC) whilst keeping images from the same location together and maintaining an approximate split of 80 % training, 10 % validation and 10 % testing.

replicates of a given combination of these aspects. Replicates were conducted since the stochastic nature of the training process can cause models to converge at slightly different points (Aggarwal, 2018). For each experiment the replicate with the highest evaluation metric was the one compared against the other experiments.

Once we identified the best-performing experimental set-up for each task, we conducted *K*-fold cross-validation on the combined training and validation sets to understand how the choice of these affects model performance (see Sect. 3.1.3 and 3.2.2).

Following model selection and cross-validation we calculated the performance of the best model identified for each task on the test set. Finally, to see if better performance could be achieved with more data available for training, we re-trained the models on the combined training and validation data before evaluating on the test data (evaluation on the test set: Sect. 3.1.3 and 3.2.3). All stages of model development, including model selection, cross-validation and final evaluation, are shown in Fig. 4, and more information about the specific experiments conducted for model selection is given in Sect. S3.

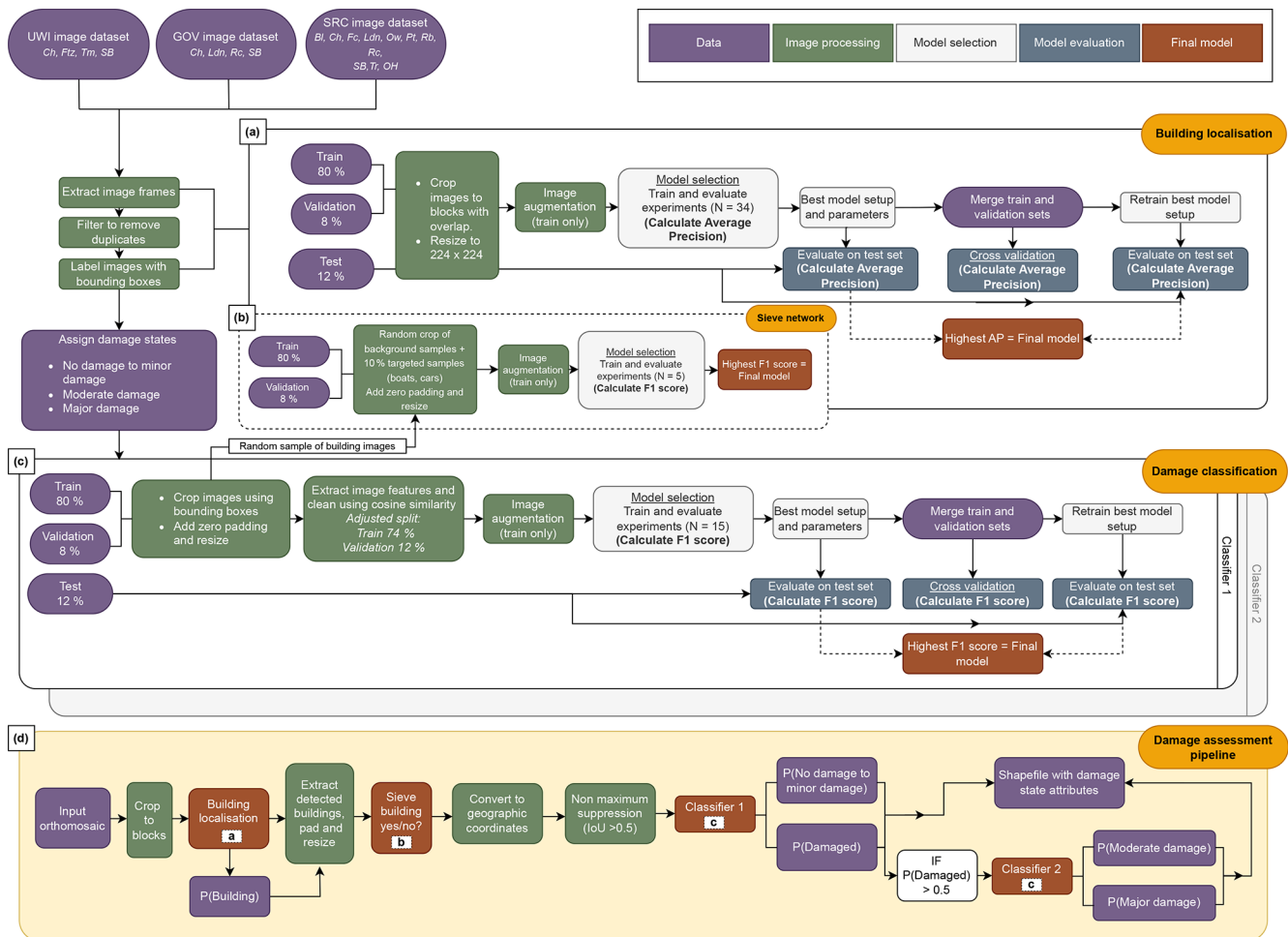


Figure 4. A schematic showing the full methodology for (a) developing a model for building localisation; (b) developing a sieve network, which acts as an add-on to the building localisation model; (c) developing a model for damage classification; and (d) the building damage assessment pipeline developed in this work. The pipeline operates on an orthomosaic image (to be generated separately) and incorporates the final trained models for building localisation and two stages of damage classification along with all the necessary processing steps to link the models. Dataset locations referred to are Bl – Belmont, Ch – Chateaubelair, Fc – Fancy, Ftz – Fitz Hughes, Ldn – London, OH – Orange Hill, Ow – Owia, Pt – Point, Rb – Rabacca, Rc – Richmond, SB – Sandy Bay, Tr – Tourama and Tm – Troumaca. Pipeline schematic generated using <https://app.diagrams.net/> (last access: 27 September 2024).

Past studies have trained deep learning algorithms on georeferenced images (i.e. each pixel has a geographical location attached) (Gupta and Shah, 2020; Shen et al., 2021; Bouchard et al., 2022) and non-georeferenced images (e.g. Y. Li et al., 2019; Pi et al., 2020; Cheng et al., 2021). In this work we labelled the non-georeferenced images and trained models on these. This was done firstly to preserve the multiple viewing angles that we have of each building, with each image counting as a different data point, and secondly due to the absence of GPS locations on a large portion of the dataset. In an operational context, spatial information must be tied to the assessed damage. Therefore, beyond the creation of distinct models for each task, we designed a comprehensive, fully automated pipeline that integrates models for building localisation and damage classification. Our pipeline

contains all the necessary processing steps to guide images through the separate models, enabling them to operate on a georeferenced orthomosaic image (to be generated separately) or on non-georeferenced images. When applied to an orthomosaic image, the output from the pipeline is a georeferenced vector dataset that can readily be plotted in a GIS to generate damage maps.

In Sect. 4 we apply the pipeline to assess building damage in Owia, St Vincent, which received 50–90 mm of tephra fall during the 2020–2021 eruption (Fig. 1). Owia was selected out of the three possible test set locations (Fig. 3) due to its large size and the existence of GPS locations that enabled the generation of a georeferenced orthomosaic image; for this we used Agisoft’s Metashape software. To compare the assessed building damage with tephra thickness, we used

the TephraFits code (Biass et al., 2019) to identify the theoretical maximum accumulation using the isopachs from Cole et al. (2023). This maximum accumulation and the isopachs were interpolated using cubic splines, and the surface was exported at a resolution of 10 m to provide a tephra thickness value for each building.

2.3.1 Building localisation

For building localisation, we used the cutting edge two-stage object detector Faster R-CNN (Ren et al., 2017). When applied to a test image containing the relevant objects, Faster R-CNN outputs the positions within the image (X , Y , width and height in pixels) of bounding boxes containing the object and a confidence score for each box. As per customary practice (Zou et al., 2019) we used a confidence of > 0.5 , meaning that only boxes with confidence greater than this are output.

For object detection, to reduce model training and inference time, full-sized images were split into image blocks. Experiments conducted as part of building localisation model selection included variations in block size and the proportion of block overlap, along with the development of separate models for images captured with different viewing angles, training for only the SRC portion of the dataset (images mostly at nadir) and the combined UWI-TV-GOV portion (images mostly off nadir). A total of 34 experiments were conducted to include all credible combinations of the varied hyperparameters and to find the best experimental set-up (Table S2).

To improve the performance of the building localisation model we developed a sieve network that runs as an add-on to the Faster R-CNN building detector. The sieve network reduces false positives which occur when the detector predicts a bounding box that does not have an overlapping labelled building (i.e. detects a building when there is not one). More details on its development are provided in Sect. S3.2.

2.3.2 Damage classification

We chose to divide building damage classification into two separate classifications: Classifier 1 distinguishes between “no damage to minor damage” vs. the combined classes of “moderate damage” and “major damage”, while Classifier 2 further differentiates between “moderate damage” and “major damage”. A hierarchical approach to classification has been found to be effective when the number of samples is limited or classes are unbalanced (D. Li et al., 2019; An et al., 2021). We conducted experiments separately for Classifiers 1 and 2. Experiments consisted of fine-tuning two different pre-trained CNNs to determine which was better and should be used in the final models for each classifier: ResNet50 (He et al., 2015) trained on the ImageNet dataset (Deng et al., 2009) and GoogLeNet (Szegedy et al., 2015) trained on the places365 dataset (López-Cifuentes et al., 2019). Fine-tuning

is a common approach to computer vision tasks where sufficiently large, labelled datasets are not available for the task at hand (typically hundreds of thousands of images are needed; Aggarwal, 2018). During fine-tuning, the high-level features that were learnt during the initial training on the large dataset can be leveraged for the new task. In addition to the different pre-trained CNNs used, experiments also considered different ways of balancing the number of images for each damage state class (oversampling the minority class, undersampling the majority class and no balancing). When applied to a test building image, the trained classifier outputs the highest probability class and the associated probability. A total of 15 experiments were conducted for each of the classification tasks. For each experiment three replicates were conducted, each consisting of a grid search to find the best combination of learning rate, batch size and L2 regularisation. For more information on this see Sect. S3.3.

2.3.3 Model evaluation metrics

For building localisation Faster R-CNN experiments, we evaluated performance using the average precision (AP) at an intersection over union (IoU) threshold of 0.5 and the F_1 score. AP, a common metric for evaluating object detection (Zou et al., 2019), measures how often the detector gets it right (true positives, TPs) vs. wrong (false positives, FPs, and false negatives, FNs). A TP occurs when a predicted box overlaps a labelled box by more than 50% (IoU > 0.5), an FP when there is no overlapping labelled box and an FN when the detector misses a labelled box. When the detector is run on a test image, a confidence score is output for each predicted box (0–1). Once the trained detector has been run over the full test set, the precision (TP/TP + FP) and recall (TP/TP+FN) are calculated at different confidence score thresholds, and the area underneath the resulting precision–recall curve represents the precision (P). AP depicts the trade-off between precision and recall and provides an overall measure of detection performance. AP values range between 0–1, where a higher value indicates a better performance.

For building localisation, the F_1 score was calculated at IoU and confidence thresholds of 0.5. The F_1 score is calculated as $F_1 = 2 \times (\text{precision} \times \text{recall}) / (\text{precision} + \text{recall})$. To evaluate the performance of classification models, we used the macro F_1 score, which is the unweighted mean of the F_1 scores calculated for each of the classes. Similarly to the AP, values of the F_1 score range between 0–1, where a higher value indicates a better performance.

3 Results

3.1 Building localisation

3.1.1 Model selection

The five experiments with the highest average precision are shown in Table 4, with the full list of experiments provided in Table S2. Average precisions across the 34 experiments ranged from 0.295 to 0.701 (Tables 4 and S2). We found that block size played an important role in model performance; out of the 34 experiments conducted, the top three used a block size of 550×550 pixels, which was the middle of the sizes tested (450, 550, 650). We observed that models trained on the full dataset performed better than models trained separately for the nadir (SRC) and off-nadir imagery sets (UWI-TV and GOV sets combined) (Tables 4 and S2).

All trained sieve networks achieved macro and class F_1 scores that were > 0.973 (Table S3). The sieve networks' efficacy at improving building localisation is demonstrated by comparing the results of the best detector when applied to the validation dataset pre-sieving (Table 4, Row ID 1) with the post-sieving results. Prior to sieving there were a large number of false positive detections, resulting in a precision of 0.588, and after sieving these were reduced and the precision increased to 0.695 (Table 5).

3.1.2 Cross-validation

A cross-validation was conducted for the single best-performing building localisation model (without the sieve network) to understand how the choice of training and validation data affects performance. Analysing performance variations across different testing datasets can then inform recommendations for future data collection strategies (see Sect. 6).

We found that the performance of the selected object detector varied depending upon the location (Fig. 5a) or imagery dataset (Fig. 5b) used for testing. For models tested on different locations, average precisions that are in line with the AP achieved on the full validation set (0.701) were obtained for Point and Fancy (Fig. 5a). The lowest AP values were for London (0.063) and Fitz Hughes (0.187). The standard deviation (SD) (Fig. 5) shows the variability in performance between the three replicates that were trained for each test, which arises due to the stochastic nature of the training process. For models tested on the different imagery datasets individually the AP was low, with a mean value across all datasets of < 0.2 (Fig. 5b). For all three locations (Chateaubelair, Sandy Bay, London), AP for models evaluated on the SRC dataset was lower than for the UWI-TV or GOV datasets.

3.1.3 Evaluation on the test set

Evaluation of the best detection model on the test set, which consists of completely unseen data from Owia, Richmond

and Troumaca (Fig. 3), produced an AP value that is the same as the value for the validation data (0.701) (Table 6). To understand if a better model could be achieved with more data available for training, we combined the training and validation data and used this to retrain the best experimental set-up for the detector. Evaluation of the retrained model on the test set resulted in an average precision increase from 0.701 to 0.751 for the non-sieved detector and from 0.668 to 0.728 for the sieved detector, showing that having more data available for training produced a better model (Table 6).

While the AP is higher for the retrained detector without the sieve, the addition of the sieve network creates a better balance between the precision and recall which is reflected in the higher F_1 score (Table 6). For the present application equal importance is given to (1) making correct predictions about building locations and (2) identifying as many buildings as possible. Consequently, striking the balance between precision and recall is crucial. We therefore selected the retrained detector plus sieve network as the final building localisation model and the model that is incorporated into the damage assessment pipeline (Table 6).

3.2 Damage classification

3.2.1 Model selection

The five experiments with the highest macro F_1 score are shown in Table 7, with the full lists provided in Tables S4 and S5. For Classifier 1, macro F_1 scores across all 15 experiments ranged from 0.753 to 0.836, while for Classifier 2 scores ranged from 0.776 to 0.810 (Tables 7, S4 and S5). Models trained to differentiate between the no damage to minor damage and damaged classes performed better for the no damage to minor damage class, while those trained to differentiate between moderate and major damage performed better for the major damage class (Table 7). The best-performing models for both classifiers used the ResNet50 architecture rather than GoogLeNet with an unbalanced dataset. For Classifier 1 the best model had $F_1 = 0.962$ for the no damage to minor damage class and $F_1 = 0.710$ for the damaged class, while for Classifier 2 the moderate damage class had $F_1 = 0.770$ and major damage $F_1 = 0.851$.

3.2.2 Cross-validation

A cross-validation was conducted for both of the single best-performing models for Classifiers 1 and 2 identified through model selection. As was the case for the best building localisation model, this was done to understand how the choice of training and validation datasets affected model performance and to understand how our model might perform on a new dataset.

The performance of Classifier 1 for the no damage to minor damage class is consistent across the distinct locations and datasets used for evaluation with mean F_1 scores

Table 4. Hyperparameters for the five experiments with the highest average precision conducted for building localisation, ordered by average precision. The full table consisting of all 34 experiments is provided in the supplementary material. Columns marked with “1” contain yes/no (Y/N) information. In the “training dataset” column, a signifies all of SRC, UWI-TV and GOV.

Row ID	Block size	Mixed block size ¹	Block overlap	Block resized ¹	Training dataset	Max average precision	F ₁ score
1	550	N	50 %	Y	a	0.701	0.669
2	550	N	20 %	Y	a	0.700	0.668
3	550	N	20 %	Y	a	0.700	0.642
4	650	N	50 %	Y	a	0.691	0.654
5	650	N	20 %	Y	a	0.678	0.670

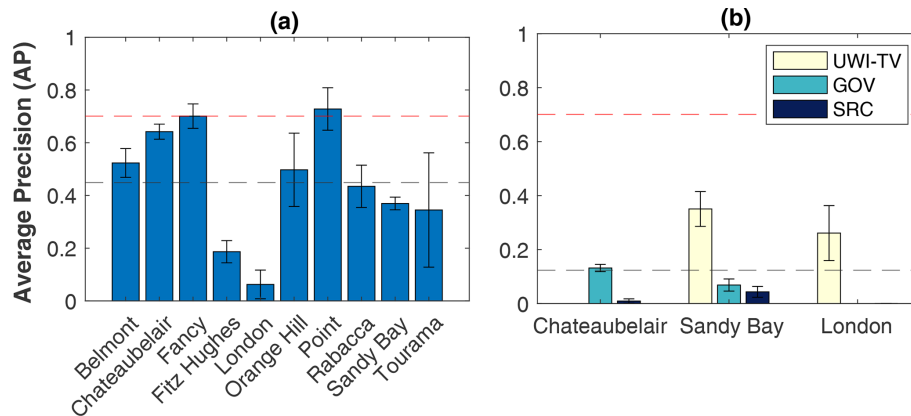


Figure 5. Cross-validation of the best experimental set-up for building localisation models which are trained to predict building box positions within the image. **(a)** The effect of changing the location used as the test set on detector average precision (AP) and **(b)** the effect of changing the imagery dataset (UWI-TV, GOV, SRC) used as the test set on AP. **(b)** For the cross-validation of the imagery dataset, models are trained on all data from that location excluding the location used for testing as indicated by the bar. For London there are data from the GOV dataset; however the number of images in the SRC dataset is insufficient for training, so no bar is shown for GOV. The AP shown is the mean value from three trained models with the same set-up, while the error bars show the standard deviation. Dashed black lines show the mean AP value across all cross-validation trained models; dashed red lines show the best AP from the experiments (0.701; Table 4).

Table 5. Comparing the performance of the best building localisation model when applied to the validation dataset before and after running the results through the sieve network.

	Precision	Recall	F ₁
Best detector pre-sieving	0.588	0.776	0.669
Best detector post-sieving	0.695	0.730	0.712

between 0.913–0.983 for locations and 0.898–0.976 for datasets (Fig. 6). For the damaged class there is more variety in the performance across the locations and datasets used for evaluation. The mean F₁ scores for the separate locations range from 0.588 (Fitz Hughes) to 0.779 (Tourama), while for the different datasets the range is 0.393 (London SRC) to 0.745 (Sandy Bay SRC).

For Classifier 2, the moderate damage class is more sensitive to the choice of location and dataset used for the evaluation than the major damage class (Fig. 6). For the differ-

ent locations the mean F₁ score ranged from 0.583–0.974. Similarly to Classifier 1, the location with the lowest mean F₁ score is Fitz Hughes, whereas the highest score was produced for Orange Hill. For the different datasets the range for the moderate damage class is between 0.522–0.746.

For the major damage class F₁ scores for the distinct locations are between 0.728–0.933, while for the different datasets the range is between 0.711–0.867.

3.2.3 Evaluation on the test set

Evaluation of the single best models for Classifier 1 and Classifier 2 on the unseen test set produced macro F₁ scores that were comparable to the scores for the validation set: 0.829 for Classifier 1 and 0.791 for Classifier 2 (Table 8). For Classifier 2, retraining the model on the combined training and testing data increased the macro F₁ score from 0.791 to 0.838, whereas for Classifier 1 retraining produced a slightly lower macro F₁ score (0.809 compared to 0.829). Nevertheless, the

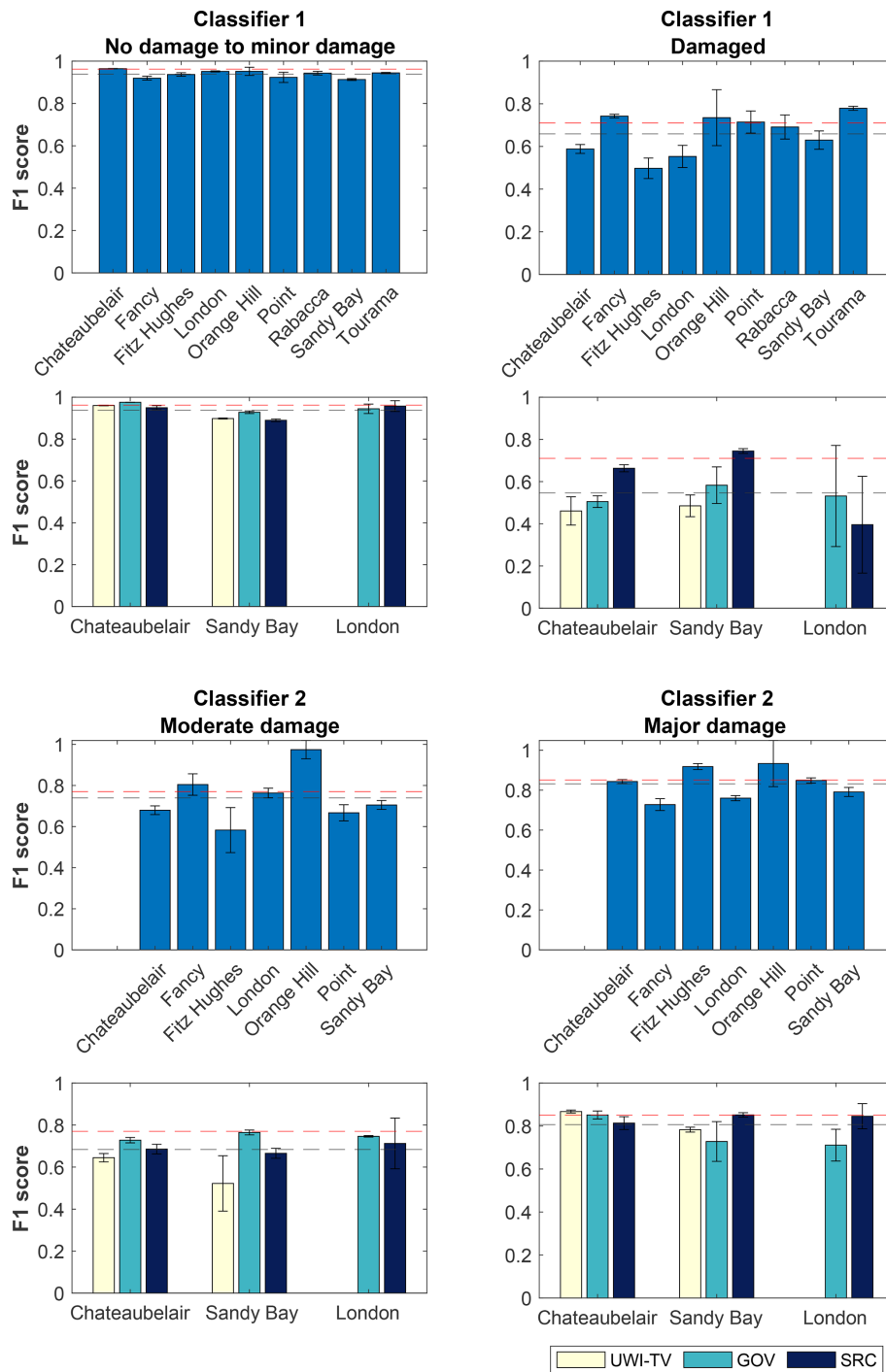


Figure 6. Cross-validation for Classifiers 1 and 2. For rows 1 and 3 the best experimental set-up was retrained on all the data from locations in the combined training and validation data and evaluated on the location shown. For rows 2 and 4 the best experimental set-up was retrained on all the data from the location shown and evaluated on each dataset (UWI-TV, GOV, SRC) separately. Each training was conducted three times, the value plotted is the mean, and the error bars show the standard deviation. Dashed black lines show the mean F_1 score across all cross-validation trained models; dashed red lines show the best F_1 score for each class from the experiments (Table 6).

Table 6. Comparison of the best building localisation models' performance when evaluated on the validation and the test sets. AP is average precision, P is precision, and R is recall. Retrain models are trained on the combined training and validation sets. Results for the final model that is used in the damage assessment pipeline are in bold.

	Validation set				Test set			
	AP	P	R	F_1	AP	P	R	F_1
Detector (0.5 conf)	0.701	0.588	0.776	0.669	0.701	0.604	0.776	0.679
Detector plus sieve (0.5 conf)	0.681	0.695	0.730	0.712	0.668	0.606	0.757	0.673
Detector retrain					0.751	0.642	0.816	0.719
Detector retrain plus sieve					0.728	0.710	0.782	0.744

Table 7. The top five experiments conducted for each of the building damage classifiers, ordered by the macro F_1 score. The full list consisting of all 15 experiments for each classifier is provided in Tables S4 and S5.

Classifier 1					
Row ID	Architecture	Class balancing: not balanced/ undersampled/ oversampled	F_1 no damage to minor damage	F_1 damaged	F_1 macro
1	Resnet50	Not	0.962	0.710	0.836
2	Resnet50	Not	0.960	0.696	0.828
3	Resnet50	Not	0.957	0.699	0.828
4	Resnet50	Not	0.962	0.692	0.827
5	Resnet50	Under	0.951	0.646	0.799
Classifier 2					
Row ID	Architecture	Class balancing: not balanced/ undersampled/ oversampled	F_1 moderate damage	F_1 major damage	F_1 macro
1	Resnet50	Not	0.770	0.851	0.810
2	GoogLeNet	Over	0.737	0.848	0.793
3	Resnet50	Over	0.749	0.835	0.792
4	Resnet50	Not	0.749	0.835	0.792
5	Resnet50	Under	0.735	0.845	0.790

retrained model for Classifier 1 achieved a higher recall on the damaged class than the non-retrained model. In an operational setting it is desirable to correctly classify as many of the damaged buildings as possible, since in our pipeline these will be passed onto Classifier 2; therefore we took the retrained models for both classifiers as the final models and the models that are incorporated into the damage assessment pipeline.

4 Application of the full damage assessment pipeline: assessing tephra fall building damage in Owia

In this work we have developed separate models for building localisation and two stages of damage classification. However, in an operational context models need to work sequentially; this led to the development of our damage assessment pipeline (outlined in Fig. 4d). The pipeline operates on an orthomosaic image and outputs a georeferenced vector set, with the following attributes (in italics) for each building that is detected: *detection* (box confidence score), *ClassPred_1* (output class from Classifier 1, damaged or no damage to minor damage), *ClassProb_1* (the probability of that class),

Table 8. Comparison of the best damage classification models' performance when evaluated on the validation and the test sets. P is precision, and R is recall. Retrain models are trained on the combined training and validation sets. Results for the final models that are used in the damage assessment pipeline are in bold.

	Validation set							Test set						
	No damage to minor damage			Damaged			F_1 macro	No damage to minor damage			Damaged			F_1 macro
	P	R	F_1	P	R	F_1		P	R	F_1	P	R	F_1	
Classifier 1	0.950	0.976	0.962	0.793	0.643	0.710	0.836	0.891	0.940	0.915	0.809	0.689	0.744	0.829
Classifier 1 retrain								0.899	0.894	0.896	0.717	0.728	0.722	0.809
	Moderate damage			Major damage			F_1	Moderate damage			Major damage			F_1
	P	R	F_1	P	R	F_1		P	R	F_1	P	R	F_1	
Classifier 2	0.769	0.660	0.770	0.852	0.825	0.851	0.810	0.903	0.663	0.765	0.730	0.927	0.817	0.791
Classifier 2 retrain								0.861	0.809	0.834	0.817	0.866	0.841	0.838

ClassPred_2 (output class from Classifier 2, moderate damage or major damage; this is only run if Classifier 1 outputs damage), *ClassProb_2* (the probability of the class output by Classifier 2) and *damageState* (the final damage state).

The tephra fall building damage map shown in Fig. 7a was produced by overlaying the georeferenced vector that was output by the pipeline with the orthomosaic image in QGIS. Our remote damage assessment pipeline identified 442 buildings. Of these, 78 % ($N = 343$) were classified as having no damage to minor damage, 9 % ($N = 40$) as having moderate damage and 13 % ($N = 59$) as having major damage. We observed that the two upper tephra fall thickness bins (70–80 and 80–90 mm) both had a higher proportion of buildings with major damage compared to the lower thickness bins (Fig. 7b and c), indicating a correlation between tephra fall thickness and building damage, though it is not very pronounced. These findings are discussed in Sect. 5.3.

The full pipeline took 1 h to run on a standard 16 GB RAM 2021 MacBook Pro with an M1 Pro chip. Most of the inference time was attributed to the building localisation module in the pipeline, which may be bypassed if building footprints are already available. When only the classifiers were run, the time taken to run was reduced to < 5 min.

5 Discussion

In this work we have developed models for building localisation and two levels of damage classification for building damage resulting from tephra fall. Our final models demonstrate strong performance for both building localisation (AP = 0.728; $F_1 = 0.744$) and damage classification (Classifier 1, $F_1 = 0.809$; Classifier 2, $F_1 = 0.838$). Despite using post-event imagery only, which makes the task more challenging than approaches using multi-temporal imagery, our results are comparable to existing optical imagery building damage assessments developed for various hazards that use both mono-temporal and multi-temporal images

(F_1 scores are between 0.656–0.868 for building localisation and 0.650–0.981 for damage classification; Table 1).

5.1 Building localisation

Through running our building localisation experiments we found that the pre-processing of images before detector training (particularly the block size) significantly influenced detector performance. The block sizes tested were chosen as a trade-off between reducing image size sufficiently to reduce computational cost and retaining a large enough size such that buildings were not dissected unnecessarily. Given that the optimum block size was the middle size of the range tested, we are confident that this balance was achieved. Cross-validation results demonstrated variability in average precision (AP) for models trained on different locations and imagery datasets (UWI-TV, GOV, SRC) (Sect. 3.1.2; Fig. 5). Deep learning models are known to perform well when the data they are evaluated on have similar characteristics to the data they were trained on, although they have more difficulty when working with “out-of-distribution” samples (Ben-David et al., 2010). Given the relatively consistent building typology across locations (most buildings observed are detached single storey buildings with either a gable or hip-shaped metal sheet roof; a lesser proportion have flat concrete roofs), the differences in AP are likely due to observable variations in UAV altitude, off-nadir angles, tephra thicknesses and varying training sample sizes.

The cross-validation AP was notably lower for the London and Fitz Hughes datasets (Sect. 3.1.2). For the London images (from SRC and GOV datasets) this is likely caused by the smaller apparent size of buildings in these images compared to the other locations due to the higher UAV altitude. Variations in object size within the training and testing data have been found to affect the performance of deep learning models developed for building localisation, with models often performing better for objects that are the same size as those in the training data (Nath and Behzadan, 2020; Cheng et al., 2021; Bouchard et al., 2022). Fitz Hughes im-

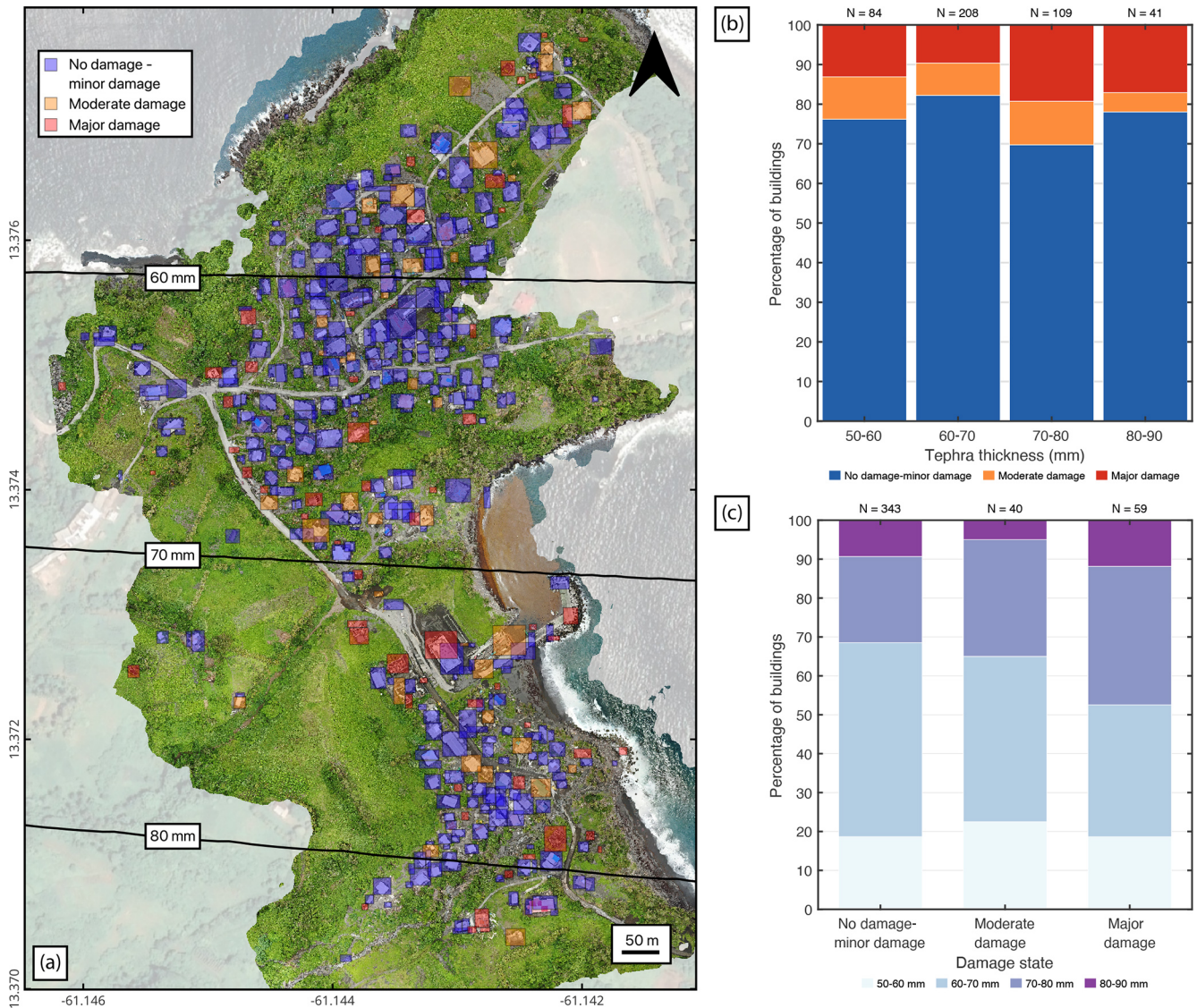


Figure 7. Application of our remote tephra fall building damage assessment pipeline to Owia, located in the north of St Vincent. **(a)** The damage map produced by overlaying the spatial data generated by our pipeline onto the orthomosaic image; black lines are tephra isopachs interpolated from Cole et al. (2023). **(b)** The proportion of damage states with increasing tephra thickness. **(c)** The proportion of tephra thickness bins with increasing damage state. Coordinate reference system: WGS 84 (EPSG:4326). Satellite basemap © Google Maps 2024.

ages were all from the UWI-TV image dataset which contributed just 17 % to the combined training and validation set used for cross-validation. This dataset was collected closer in time to the eruption; therefore as a whole had more tephra on the ground than the SRC and GOV datasets, which affects background colour. Furthermore, the UWI-TV dataset viewed buildings mostly from an off-nadir perspective, while the other datasets were predominantly nadir images. The effect of image background colour on localisation performance is expected to be minor. Cheng et al. (2021) found that for the same event localisation AP dropped from 65.6 to 63.3 when their model was tested on images containing buildings surrounded by vegetation compared to buildings with

an ocean backdrop, while Bouchard et al. (2022) suggested that models quickly learn to ignore background pixels. On the other hand, variation in off-nadir angles is a widely acknowledged challenge of working with UAV or aerial images (Cotrufo et al., 2018; Nex et al., 2019; Pi et al., 2020). Under-representation of the mostly off-nadir UWI-TV images in the training data may have impacted the model’s ability to recognise such instances in the test data. During model development we experimented with different models for the different datasets (UWI-TV, GOV, SRC) but found that models developed on the combined dataset performed better than those developed on the separate datasets, and a combined model was the one selected and used for cross-validation. Rather than

suggesting that variations in off-nadir angle are not important, this finding likely reflects the smaller size of the individual datasets compared to the combined datasets, meaning that less information was available to learn from. The application of sampling approaches like those used for the damage states in the classification model development (over- or undersampling) could have been applied to balance the data. However, the SRC dataset is much larger than either of the UWI-TV and GOV sets (Fig. 3). Therefore we considered that oversampling would introduce significant bias towards the specific examples in the under-represented dataset, whereas through undersampling we would lose a large amount of the data that are available to learn from. Given these factors, we did not use sampling approaches. Future work might consider the application of generative AI algorithms such as generative adversarial networks (GANs) to expand the dataset (e.g. Yi et al., 2018; Yorioka et al., 2020), although more work needs to be done to quantify the diversity in the generated data.

The variability in cross-validation results for the building localisation model likely comes from a combination of the above factors (differences in UAV altitude, off-nadir angles, tephra thickness and varying training sample sizes) and suggests that there was insufficient information in the training data for our detection models to perform well across the range of characteristics present. This is supported by the increased performance when the best localisation model was retrained on the combined training and validation data. However, further investigation is required to separate the unique effect of each aspect.

5.2 Damage classification

The final classification models achieved better performance than the final localisation model with macro F_1 scores of 0.809 and 0.838 on the test data (Table 8). Cross-validation showed that classification models were less sensitive than the localisation model to the choice of datasets used for training and evaluation (Sect. 3.2.2). We found that class-wise our models performed better on the no damage to minor damage class followed by the major damage class. This agrees with other multi-class studies that have found the extremities of the damage state scheme applied easier to classify than the intermediate ones (Kerle et al., 2019; Valentijn et al., 2020).

5.3 Application of the full damage assessment pipeline: assessing tephra fall building damage in Owia

The application of our remote damage assessment pipeline to the town of Owia found that 22 % of buildings that received tephra accumulation in the range of 50–90 mm experienced moderate damage or major damage. Within this range, the relationship between tephra thickness and building damage was not as pronounced as in other studies (Blong, 2003a; Hayes et al., 2019; Jenkins et al., 2024). This may

be attributed to the small geographic area and therefore small range of tephra thicknesses considered in our application when compared to other studies. In the damage assessments of Blong (2003a), Hayes et al. (2019), and Jenkins et al. (2024), buildings received ~ 100 to 950 mm, trace to 600 mm and trace to > 220 mm, respectively. Spence et al. (1996) assessed building damage over a similarly narrow range of tephra thicknesses to this work (~ 150 –200 mm) and found that there was considerable variation in the level of damage despite the majority of buildings having a metal sheet roof. The spacing between the principal roof supports (roof span) was found to be important for the amount of damage observed, with long-span buildings experiencing higher levels of damage than short-span ones (Spence et al., 1996). There are limited long-span buildings in the Owia case study; however additional characteristics such as construction style and material, building layout, age, condition, height, and roof pitch can all affect a building's ability to withstand tephra loading (Spence et al., 1996; Pomonis et al., 1999; Blong, 2003a; Jenkins et al., 2014). Variation in these characteristics across Owia could be responsible for the observed variation in building damage over the narrow range of thicknesses considered.

If we convert tephra thickness to loading, we can compare the results of our assessment with existing relationships between tephra loading and damage for similar building types. Using a density of 1500 kg m^{-2} (Cole et al., 2023) suggests that a loading of at least 75 – 135 kg m^{-2} was applied to buildings for the range of thicknesses considered (50–90 mm). Census data for Owia states that 90 % of buildings have metal sheet roofs (SVG population and housing census, 2012), with the remaining 8 % comprised of reinforced concrete roofs and 2 % “other material”. Given the higher resistance of the 8 % of buildings with non-metal-sheet roofs in Owia, we might expect vulnerability models developed for metal sheet roofs to overestimate damage in the town. Fragility functions developed for Indonesian style buildings with metal sheet roofs (Williams et al., 2020) calculate a 48 %–80 % probability of Owia buildings experiencing damage exceeding DS2, higher than the 22 % experiencing moderate or major damage in our study. Fragility curves for roof failure (major damage) of old or poor-condition metal sheet roofs (Jenkins et al., 2014) calculate that just over 10 % of buildings in Owia would experience sufficient loading for roof collapse, comparable to the 13 % observed in our study. These comparisons highlight some of the challenges associated with using vulnerability models developed for different locations. Moreover, they reiterate the need for the collection of both post-event impact data and building typology information that can be used to increase the amount of empirical data available for vulnerability model development and allow regional vulnerability models to be developed for specific building types.

Like the studies presented in Table 1, our pipeline consists of separate models for localisation and damage classification. One of the benefits of this is that in locations where precise

building location information is available for the assessment area, the localisation step can be bypassed and only the classifiers run. This not only enhances overall performance but also significantly reduces computation time. Furthermore, either of the classifiers can be run independently and/or combined with other damage assessment procedures; for example, an initial SAR-based assessment (e.g. Yun et al., 2015; Jung et al., 2016) could be followed with our Classifier 2 to provide additional granularity to the severity of the damage at a building level rather than a pixel level.

5.4 Generalisability to other locations

Our models have performed well for images collected on the island of St Vincent where building typologies are relatively consistent. We therefore expect that our models will perform well in other locations with similar building types, such as the other islands in the Lesser Antilles. This hypothesis should be validated through further testing. In the absence of additional UAV datasets that include damaged buildings, testing can be done by conducting pre-event surveys to test the performance of the building localisation model and Classifier 1 for the no damage to minor damage class. While this is unable to assess the ability of our approach to classify damage, it would provide *some* indication of performance following an event in a new location.

To develop a model that is robust to the diverse building types found across the world necessitates assembling diverse datasets showcasing potential variations in building types and the associated tephra fall damage. To our knowledge the UAV datasets described in this work are the first of their kind. However, the increasing utilisation of UAVs during and after volcanic events suggests the possibility of the emergence of more datasets in the years to come. Our model represents a crucial initial step towards the operational implementation of this approach globally. The compilation of global UAV tephra fall building damage datasets will facilitate the ongoing refinement of building damage assessment approaches, including the one presented here. In pursuit of this objective, our models stand ready for retraining as more data become available. While our approach leverages images captured under a spectrum of flight conditions (off-nadir angle, altitude, flight trajectory), our investigation has both pinpointed specific conditions that are best suited for capturing building damage, which are detailed in Sect. 6, and highlighted the importance of consistency in data collection.

5.5 Improving model performance and future perspectives

The advantages of acquiring additional UAV datasets both before and after an event have been outlined in Sect. 5.4. In addition to this, pre-event imagery can be used to construct building inventories manually or using machine learning methods (e.g. Iannelli and Dell'Acqua, 2017; Gonzalez

et al., 2020; Meng et al., 2023). Prior to an eruption, information about how the building typologies present will respond under certain tephra loadings (i.e. the forecasted damage state) can be obtained through the application of fragility functions. This information could enhance our model by serving as prior information that is updated with outputs from our remote damage assessment using Bayesian statistics. A similar approach has been suggested for updating the United States Geological Survey's (USGS) Prompt Assessment of Global Earthquakes for Response (PAGER) system (Noh et al., 2020). The framework provides a structured way of incorporating the PAGER-forecasted loss with the potentially noisy and incomplete observations of loss in the early stages of response.

Alternatively, with ample individual building inventory data available, tailored damage classification models for specific building typologies could be developed and applied. The rationale is that a model dedicated to a specific building type is expected to outperform a generic multi-typology model.

In this work, we established a three-class damage state framework. Existing frameworks that were developed for ground-based tephra fall damage assessment split damage into five damage state classes and one non-damage class (Spence et al., 1996; Blong, 2003a; Hayes et al., 2019; Jenkins et al., 2024); however in our preliminary analyses we found that (1) in many images we were unable to confidently apply a six-class scheme due to only being able to see one side of the building and (2) there were not enough examples of each damage state class to be able to train a six-class model. With the addition of future tephra fall building damage datasets it may be possible to apply a finer-resolution damage state framework that can provide more detail on the observable damage. However, it is unlikely that the resolution of ground surveys can be achieved using optical imagery, since lower damage states are still difficult to resolve even with very high-resolution images (Cotrufo et al., 2018). Some studies have incorporated 3D point-cloud information into analyses (Cusicanqui et al., 2018; Vetrivel et al., 2018). While these approaches have shown potential and could potentially be used to provide additional granularity to our damage states, we opted against integrating point cloud analyses into our model due to the considerably longer processing times associated with such an approach. Longer processing times would undermine the swift processing requirement inherent in our methodology.

5.6 Caveats

During the assignment of building damage states, uncertainties arose, particularly concerning the interpretation of tarpaulins and pre-existing damage. For tarpaulins, the ambiguity arose from whether these were either strategically placed prior to the eruption as preventative measures to cause tephra to slide off the roof more easily or placed post-event to cover damage caused by tephra fall. Additionally, in cer-

tain instances, distinguishing between a collapsed roof and a section of the building initially lacking roofing material – possibly functioning as a walled storage area – proved challenging. Pre-existing damage not related to volcanic activity or buildings that were under construction at the time of image acquisition were considered as damaged and classified accordingly. The presence of buildings under construction at the time of image acquisition has been recognised as a challenge in studies using mono-temporal imagery (Nex et al., 2019; Cheng et al., 2021). Pre-event imagery would have provided clarity on both of these matters; however this was not available at high enough resolution for this region.

The majority of images used for training and evaluating our models came from the SRC dataset, which was collected several months after the eruption. As a result, the majority of images do not have much tephra present. In an operational context, to expedite the recovery process, data would ideally be collected as quickly after the eruption as is safe to do so; therefore more tephra would be present in the images. Given the compound effects of variations in flight angle, image lighting, resolution and also the presence of tephra, we do not have enough information to test the effect of tephra thickness on model performance, and caution should be taken when using the model on data collected at different times after the eruption.

6 Recommendations for UAV building damage assessment data collection

In the future we advocate for the adoption of a standardised protocol for data collection for the purpose of UAV damage assessment. While our model was developed using a diverse dataset, there were some disparities in performance across distinct data types. Consequently, the standardisation of image collection serves two purposes: (1) to allow the best results to be achieved when implementing our models and (2) to collect data that are rich in information useful for damage assessment with the aim of working towards the development of global datasets for tephra fall damage. For best results we have the following recommendations:

- The bulk of our dataset was collected several months after the eruption of La Soufrière; however, for generating a global dataset that can be used for response and recovery, models should ideally be trained on images collected shortly (days to weeks) after an event.
- Flight paths should be pre-programmed to ensure comprehensive coverage of the area and limit bias associated with over-representation of certain buildings. Ideally two flights would be conducted with two sets of perpendicular flight lines to capture buildings from a different perspective. GPS positioning should be enabled.
- A fixed altitude of 50–80 m above the ground should be maintained where possible. This is appropriate to cap-

ture sufficient data for accurate damage classification based on the established framework and strikes a balance between detailed information capture and overall coverage. In mountainous areas this may not be achievable for some UAV types, in which case a uniform height should be maintained such that the size of buildings is consistent across image frames.

- We suggest a slightly off-nadir camera positioning ($\sim 5\text{--}15^\circ$), which is sufficient to capture any bending in the roof that may not be captured from a nadir perspective.
- Overlap between images should be enough to generate orthoimages; 80 % forward and 70 % lateral overlap is sufficient.

In addition to the development of optimum post-event data collection practises we advocate for the collection of pre-event UAV datasets. Ideally, pre- and post-event imagery is collected using the same flight paths, altitudes and camera positioning. Pre-event datasets serve multiple purposes:

- They facilitate the creation of building inventories.
- They enable precise comparison of pre- and post-event imagery, reducing uncertainty regarding initial building conditions.
- They support the development of high-resolution change-detection models, potentially yielding more accurate results than relying solely on post-event imagery.
- They provide an opportunity for UAV pilots to gain experience in capturing building datasets during “quiet times”.

7 Conclusions

Following a large tephra fall event, building damage assessment needs to be conducted rapidly for the purpose of response and recovery and for the collection of data that can be used to forecast building damage from future events. By leveraging post-event optical imagery obtained after the 2021 eruption of La Soufrière volcano on the island of St Vincent, as well as convolutional neural networks, we have developed an automated tephra fall building damage assessment pipeline. The pipeline incorporates models for building localisation and two distinct levels of damage classification: distinguishing between no damage to minor damage and damage, as well as between moderate and major damage, which were trained and evaluated separately. When provided with UAV optical imagery, our pipeline can rapidly generate spatial building damage information. Our models perform well for the St Vincent datasets and are anticipated to perform well in locations where building typologies are similar, but this requires more testing to understand the limits of their application.

Model building localisation cross-validation results underscore the influence of factors such as UAV altitude, off-nadir angles, tephra thickness and training sample sizes on model performance, while results show that damage classification models were affected by these factors to a lesser extent. We acknowledge the challenges posed by diverse datasets and by limited data, and we propose a series of recommendations to guide the collection of future UAV building damage datasets. In addition to the collection of post-event datasets we advocate for the collection and incorporation of pre-event datasets, which can be used to support the advancement of change-detection models, to partially evaluate the models presented here during quiescent times and to develop building inventories that can be used along with fragility functions for forecasting building damage.

Our research marks a step forward in tephra fall building damage assessment, offering a versatile and effective pipeline with the potential for regional applicability. As the field of UAV-based damage assessment in volcanology continues to evolve, our work lays a foundation for further advancements, contributing to the resilience of communities in the face of volcanic eruptions.

Data availability. All trained models along with the code required to execute the damage assessment pipeline and instructions for usage are provided at <https://doi.org/10.5281/zenodo.14375616> (Tennant et al., 2024).

Supplement. The supplement related to this article is available online at: <https://doi.org/10.5194/nhess-24-4585-2024-supplement>.

Author contributions. Conceptualisation: SFJ, RR, ET, VM. Data collection: RR and VM. Development of the methodology: ET, SFJ, BW. Software: ET. Formal analysis: ET. Supervision: SFJ. Writing – original draft: ET. Writing – reviewing and editing: ET, SFJ, VM, RR, BW, BT, SHY.

Competing interests. The contact author has declared that none of the authors has any competing interests.

Disclaimer. Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims made in the text, published maps, institutional affiliations, or any other geographical representation in this paper. While Copernicus Publications makes every effort to include appropriate place names, the final responsibility lies with the authors.

Acknowledgements. We are indebted to Monique Johnson of the UWI Seismic Research Centre; Javid Collins of UWI-TV; and Nikolai Lewis and Marla Mulraine of the Government of St Vin-

cent and the Grenadines Ministry of Transport, Works, Lands and Surveys, and Physical Planning for sharing their UAV data and collaborating on this work. All images and data provided in this study have been approved for publication by the local agency responsible for monitoring geohazards in St Vincent: the UWI Seismic Research Centre. We are very grateful to Chee Jain Hao Denny, Sim Yu Yang, Isaiah Loh Kai En and Huang Wanxin for their assistance with data preparation and to Vanesa Burgos, Elinor Meredith, Alberto Ardid and Tom Wilson for interesting discussions around machine learning and building damage assessment. We would like to thank Sébastien Biass and one anonymous reviewer for their detailed and constructive reviews that considerably improved the manuscript, as well as Giovanni Macedonio for their editorial handling.

Financial support. This research has been supported by the Earth Observatory of Singapore via its funding from the National Research Foundation Singapore and the Singapore Ministry of Education under the Research Centres of Excellence initiative and comprises EOS contribution number 596. Additional support was provided by the AXA Research Fund through their Joint Research Initiative.

Review statement. This paper was edited by Giovanni Macedonio and reviewed by Sébastien Biass and one anonymous referee.

References

- Aggarwal, C. C. (Eds.): Neural networks and deep learning, Springer, eBook ISBN 978-3-319-94463-0, 2018.
- An, G., Akiba, M., Omodaka, K., Nakazawa, T., and Yokota, H.: Hierarchical deep learning models using transfer learning for disease detection and classification based on small number of medical images, *Sci. Rep.*, 11, 4250, <https://doi.org/10.1038/s41598-021-83503-7>, 2021.
- Andaru, R. and Rau, J. Y.: Lava dome changes detection at Agung Mountain during high level of volcanic activity using UAV photogrammetry, *Int. Arch. Photogramm. Remote Sens. Spatial Inf. Sci.*, XLII-2/W13, 173–179, <https://doi.org/10.5194/isprs-archives-XLII-2-W13-173-2019>, 2019.
- Anniballe, R., Noto, F., Scalia, T., Bignami, C., Stramondo, S., Chini, M., and Pierdicca, N.: Earthquake damage mapping: An overall assessment of ground surveys and VHR image change detection after L'Aquila 2009 earthquake, *Remote Sens. Environ.*, 210, 166–178, <https://doi.org/10.1016/j.rse.2018.03.004>, 2018.
- Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., and Vaughan, J. W.: A theory of learning from different domains, *Mach. Learn.*, 79, 151–175, <https://doi.org/10.1007/s10994-009-5152-4>, 2010.
- Biass, S., Bonadonna, C., and Houghton, B. F.: A step-by-step evaluation of empirical methods to quantify eruption source parameters from tephra-fall deposits, *J. Appl. Volcanol.*, 8, 1–16, <https://doi.org/10.1186/s13617-018-0081-1>, 2019.
- Biass, S., Jenkins, S., Lallemand, D., Lim, T. N., Williams, G., and Yun, S. H.: Remote sensing of volcanic impacts, in: *Forecasting and Planning for Volcanic Hazards, Risks, and Disasters*, edited

- by: Papale, P., Elsevier, 473–491, <https://doi.org/10.1016/B978-0-12-818082-2.00012-3>, 2021.
- Blong, R.: Building damage in Rabaul, Papua New Guinea, 1994, *Bull. Volcanol.*, 65, 43–54, <https://doi.org/10.1007/s00445-002-0238-x>, 2003a.
- Blong, R.: A Review of Damage Intensity Scales, *Nat. Hazards*, 29, 57–76, <https://doi.org/10.1023/A:1022960414329>, 2003b.
- Bouchard, I., Rancourt, M. È., Aloise, D., and Kalaitzis, F.: On Transfer Learning for Building Damage Assessment from Satellite Imagery in Emergency Contexts, *Remote Sens.*, 14, 1–29, <https://doi.org/10.3390/rs14112532>, 2022.
- Bruzzone, L. and Fernández Prieto, D.: Automatic Analysis of the Difference Image for Unsupervised Change Detection, *IEEE T. Geosci. Remote*, 38, 1171–1181, <https://doi.org/10.1109/36.843009>, 2000.
- Cheng, C. S., Behzadan, A. H., and Noshadravan, A.: Deep learning for post-hurricane aerial damage assessment of buildings, *Comput.-Aided Civ. Inf.*, 36, 695–710, <https://doi.org/10.1111/mice.12658>, 2021.
- Cole, P. D., Barclay, J., Robertson, R. E. A., Mitchell, S., Davies, B. V., Constantinescu, R., Sparks, R. S. J., Aspinall, W., and Stinton, A.: Explosive sequence of La Soufrière, St Vincent, April 2021: insights into drivers and consequences via eruptive products, *Geol. Soc. Spec. Publ.*, 539, 81–106, <https://doi.org/10.6084/m9.figshare.c.6474317>, 2023.
- Cotrufo, S., Sandu, C., Tonolo, F. G., and Boccardo, P.: Building damage assessment scale tailored to remote sensing vertical imagery, *Eur. J. Remote Sens.*, 51, 991–1005, <https://doi.org/10.1080/22797254.2018.1527662>, 2018.
- Cusicanqui, J., Kerle, N., and Nex, F.: Usability of aerial video footage for 3-D scene reconstruction and structural damage assessment, *Nat. Hazards Earth Syst. Sci.*, 18, 1583–1598, <https://doi.org/10.5194/nhess-18-1583-2018>, 2018.
- Deligne, N. I., Jenkins, S. F., Meredith, E. S., Williams, G. T., Leonard, G. S., Stewart, C., Wilson, T. M., Biass, S., Blake, D. M., Blong, R. J., and Bonadonna, C.: From anecdotes to quantification: advances in characterizing volcanic eruption impacts on the built environment, *Bull. Volcanol.*, 84, 7, <https://doi.org/10.1007/s00445-021-01506-8>, 2022.
- Deng, J., Dong, R., Socher, L., Li, L.-J., Li, K., and Fei-Fei, L.: Imagenet: A large-scale hierarchical image database, in: *Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition*, 20–25 June 2009, Miami, USA, 248–255, <https://doi.org/10.1109/CVPR.2009.5206848>, 2009.
- Duarte, D., Nex, F., Kerle, N., and Vosselman, G.: Satellite Image Classification of Building Damages Using Airborne and Satellite Image Samples in a Deep Learning Approach, *ISPRS Ann. Photogramm. Remote Sens. Spatial Inf. Sci.*, IV-2, 89–96, <https://doi.org/10.5194/isprs-annals-IV-2-89-2018>, 2018.
- Cao, Q. D. and Choe, Y.: Building Damage Annotation on Post-Hurricane Satellite Imagery Based on Convolutional Neural Networks, *Nat. Hazards*, 103, 3357–3376, <https://doi.org/10.1007/s11069-020-04133-2>, 2020.
- Gailler, L., Labazuy, P., Régis, E., Bontemps, M., Souriot, T., Bacques, G., and Carton, B.: Validation of a new UAV magnetic prospecting tool for volcano monitoring and geohazard assessment, *Remote Sens.*, 13, 1–10, <https://doi.org/10.3390/rs13050894>, 2021.
- Galanis, M., Rao, K., Yao, X., Tsai, Y. L., Ventura, J., and Fricker, G. A.: DamageMap: A post-wildfire damaged buildings classifier, *Int. J. Disast. Risk Reduct.*, 65, 102540, <https://doi.org/10.1016/j.ijdr.2021.102540>, 2021.
- Ghosh, S., Huyck, C. K., Greene, M., Gill, S. P., Bevington, J., Svekla, W., DesRoches, R., and Eguchi, R. T.: Crowdsourcing for rapid damage assessment: The global earth observation catastrophe assessment network (GEO-CAN), *Earthq. Spectra*, 27, 179–198, <https://doi.org/10.1193/1.3636416>, 2011.
- Gonzalez, D., Rueda-Plata, D., Acevedo, A. B., Duque, J. C., Ramos-Pollán, R., Betancourt, A., and García, S.: Automatic detection of building typology using deep learning methods on street level images, *Build. Environ.*, 177, 106805, <https://doi.org/10.1016/j.buildenv.2020.106805>, 2020.
- Gupta, R. and Shah, M.: RescueNet: Joint building segmentation and damage assessment from satellite imagery, in: *Proceedings of the International Conference on Pattern Recognition*, Institute of Electrical and Electronics Engineers Inc., 10–15 January 2020, Milan, Italy, 4405–4411, <https://doi.org/10.1109/ICPR48806.2021.9412295>, 2020.
- Hayes, J., Wilson, T. M., Deligne, N. I., Cole, J., and Hughes, M.: A model to assess tephra clean-up requirements in urban environments, *J. Appl. Volcanol.*, 6, 1–9, <https://doi.org/10.1186/s13617-016-0052-3>, 2017.
- Hayes, J. L.; Calderón B, R.; Deligne, N. I.; Jenkins, S. F.; Leonard, G. S.; McSparran, A. M.; Williams, G. T.; Wilson, T. M. Timber-Framed Building Damage from Tephra Fall and Lahar: 2015 Calbuco Eruption, Chile. *J. Volcanol. Geoth. Res.*, 374, 142–159, <https://doi.org/10.1016/j.jvolgeores.2019.02.017>, 2019.
- He, K., Zhang, X., Ren, S., and Sun, J.: Deep Residual Learning for Image Recognition, *arXiv [preprint]*, arXiv:1512.03385, <https://doi.org/10.48550/arXiv.1512.03385>, 2015.
- Iannelli, G. and Dell’Acqua, F.: Extensive Exposure Mapping in Urban Areas through Deep Analysis of Street-Level Pictures for Floor Count Determination, *Urban Sci.*, 1, 16, <https://doi.org/10.3390/urbansci1020016>, 2017.
- Ishii, M., Goto, T., Sugiyama, T., Saji, H., and Abe, K.: Detection of earthquake damaged areas from aerial photographs by using colour and edge information, in: *Proceedings of the 5th Asian Conference on Computer Vision*, 23–25 January 2002, Melbourne, Australia, 27–32, http://aprs.dictaconference.org/accv2002/accv2002_proceedings/Ishii27.pdf (last access: 11 December 2024), 2002.
- Jenkins, S., Komorowski, J. C., Baxter, P. J., Spence, R., Picquout, A., Lavigne, F., and Surono: The Merapi 2010 eruption: An interdisciplinary impact assessment methodology for studying pyroclastic density current dynamics, *J. Volcanol. Geoth. Res.*, 261, 316–329, <https://doi.org/10.1016/j.jvolgeores.2013.02.012>, 2013.
- Jenkins, S. F., Spence, R. J. S., Fonseca, J. F. B. D., Solidum, R. U., and Wilson, T. M.: Volcanic risk assessment: Quantifying physical vulnerability in the built environment, *J. Volcanol. Geoth. Res.*, 276, 105–120, <https://doi.org/10.1016/j.jvolgeores.2014.03.002>, 2014.
- Jenkins, S. F., Phillips, J. C., Price, R., Feloy, K., Baxter, P. J., Hadmoko, D. S., and de Bèlizal, E.: Developing building-damage scales for lahars: Application to Merapi volcano, Indonesia, *Bull. Volcanol.*, 77, 75, <https://doi.org/10.1007/s00445-015-0961-8>, 2015.

- Jenkins, S. F., McSporran, A., Wilson, T. M., Stewart, C., Leonard, G., Cevuard, S., and Garaebiti, E.: Tephra fall impacts to buildings: The 2017–2018 Manaro Voui eruption, Vanuatu, *Front. Earth Sci.*, 12, 1–19, <https://doi.org/10.3389/feart.2024.1392098>, 2024.
- Joseph, E. P., Camejo-Harry, M., Christopher, T., Contreras-Arratia, R., Edwards, S., Graham, O., Johnson, M., Juman, A., Latchman, J. L., Lynch, L., Miller, V. L., Papadopoulos, I., Pascal, K., Robertson, R., Ryan, G. A., Stinton, A., Grandin, R., Hamling, I., Jo, M.-J., Barclay, J., Cole, P., Davies, B. V., and Sparks, R. S. J.: Responding to eruptive transitions during the 2020–2021 eruption of La Soufrière volcano, St. Vincent, *Nat. Commun.*, 13, 4129, <https://doi.org/10.1038/s41467-022-31901-4>, 2022.
- Jung, J., Kim, D. J., Lavalle, M., and Yun, S. H.: Coherent change detection using InSAR temporal decorrelation model: A case study for volcanic ash detection, *IEEE T. Geosci. Remote*, 54, 5765–5775, <https://doi.org/10.1109/TGRS.2016.2572166>, 2016.
- Karnik, V., Schenkov, Z., and Schenk, V.: Vulnerability and the MSK scale, *Eng. Geol.*, 20, 161–168, 1984.
- Kerle, N., Nex, F., Gerke, M., Duarte, D., and Vetrivel, A.: UAV-based structural damage mapping: A review, *ISPRS Int. J. Geo-Inf.*, 9, 1–23, <https://doi.org/10.3390/ijgi9010014>, 2019.
- Khajwal, A. B., Cheng, C. S., and Noshadravan, A.: Post-disaster damage classification based on deep multi-view image fusion, *Comp.-Aided Civ. Inf.*, 38, 528–544, <https://doi.org/10.1111/mice.12890>, 2023.
- Lerner, G. A., Jenkins, S. F., Charbonnier, S. J., Komorowski, J., and Baxter, P. J.: The hazards of unconfined pyroclastic density currents: A new synthesis and classification according to their deposits, dynamics, and thermal and impact characteristics, *J. Volcanol. Geoth. Res.*, 421, 107429, <https://doi.org/10.1016/j.jvolgeores.2021.107429>, 2021.
- López-Cifuentes, A., Escudero-Viñolo, M., Bescós, J., and García-Martín, Á.: Semantic-aware scene recognition, *Pattern Recognit.*, 102, 107256, <https://doi.org/10.1016/j.patcog.2020.107256>, 2019.
- Li, S., Tang, H., He, S., Shu, Y., Mao, T., Li, J., and Xu, Z.: Unsupervised detection of earthquake-triggered roof-holes from UAV images using joint colour and shape features, *IEEE Geosci. Remote Sens.*, 12, 1823–1827, <https://doi.org/10.1109/LGRS.2015.2429894>, 2015.
- Li, Y., Hu, W., Dong, H., and Zhang, X.: Building damage detection from post-event aerial imagery using single shot multibox detector, *Appl. Sci.*, 9, 1128, <https://doi.org/10.3390/app9061128>, 2019.
- Li, D., Cong, A., and Guo, S.: Sewer damage detection from imbalanced CCTV inspection data using deep convolutional neural networks with hierarchical classification, *Automat. Constr.*, 101, 199–208, <https://doi.org/10.1016/j.autcon.2019.01.017>, 2019.
- Lucks, L., Bulatov, D., Thönnessen, U., and Böge, M.: Superpixel-wise assessment of building damage from aerial images, in: VISIGRAPP 2019 – Proceedings of the 14th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications, SciTePress, 25–27 February 2019, Prague, Czech Republic, 211–220, <https://doi.org/10.5220/0007253802110220>, 2019.
- Meng, S., Soleimani-Babakamali, M. H., and Taciroglu, E.: Automatic Roof Type Classification Through Machine Learning for Regional Wind Risk Assessment, *arXiv [preprint]*, arXiv:2305.17315, <https://doi.org/10.48550/arXiv.2305.17315>, 2023.
- Meredith, E. S., Jenkins, S. F., Hayes, J. L., Deligne, N. I., Lallemand, D., Patrick, M., and Neal, C.: Damage assessment for the 2018 lower East Rift Zone lava flows of Kīlauea volcano, Hawaii, *Bull. Volcanol.*, 84, 65, <https://doi.org/10.1007/s00445-022-01568-2>, 2022.
- Moradi, M. and Shah-Hosseini, R.: Earthquake Damage Assessment Based on Deep Learning Method Using VHR Images, *Environ. Sci. Proc.*, 5, 16, <https://doi.org/10.3390/iecg2020-08545>, 2020.
- Naito, S., Tomozawa, H., Mori, Y., Nagata, T., Monma, N., Hakamura, H., Fujiwara, H., and Shoji, G.: Building-damage detection method based on machine learning utilizing aerial photographs of the Kumamoto earthquake, *Earthq. Spectra*, 36, 1166–1187, <https://doi.org/10.1177/8755293019901309>, 2020.
- Nath, N. D. and Behzadan, A. H.: Deep Convolutional Networks for Construction Object Detection Under Different Visual Conditions, *Front. Built Environ.*, 6, 97, 1–22, <https://doi.org/10.3389/fbuil.2020.00097>, 2020.
- Nex, F., Duarte, D., Steenbeek, A., and Kerle, N.: Towards real-time building damage mapping with low-cost UAV solutions, *Remote Sens.*, 11, 287, <https://doi.org/10.3390/rs11030287>, 2019.
- Noh, H. Y., Jaiswal, K. S., Engler, D., and Wald, D. J.: An efficient Bayesian framework for updating PAGER loss estimates, *Earthq. Spectra*, 36, 1719–1742, <https://doi.org/10.1177/8755293020944177>, 2020.
- Novikov, G., Trekin, A., Potapov, G., Ignatiev, V., and Burnaev, E.: Satellite imagery analysis for operational damage assessment in emergency situations, in: *Lecture Notes in Business Information Processing*, Springer Verlag, 347–358, https://doi.org/10.1007/978-3-319-93931-5_25, 2018.
- PDNA – Post Disaster Needs Assessment: La Soufrière Volcanic Eruption, UNDP – United Nations Development Programme, <https://www.undp.org/barbados/publications/post-disaster-needs-assessment-pdna-st-vincent-and-grenadines> (last access: 3 October 2022), 2022.
- Pi, Y., Nath, N. D., and Behzadan, A. H.: Convolutional neural networks for object detection in aerial imagery for disaster response and recovery, *Adv. Eng. Inform.*, 43, 101009, <https://doi.org/10.1016/j.aei.2019.101009>, 2020.
- Pomonis, A. A., Spence, R., and Baxter, P.: Risk assessment of residential buildings for an eruption of Furnas Volcano, Sao Miguel, the Azores, *J. Volcanol. Geoth. Res.*, 92, 107–131, 1999.
- Ren, S., He, K., Girshick, R., and Sun, J.: Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks, *IEEE Trans. Pattern Anal. Mach. Intel.*, 39, 1137–1149, <https://doi.org/10.1109/TPAMI.2016.2577031>, 2017.
- Román, A., Tovar-Sánchez, A., Roque-Atienza, D., Huertas, I. E., Caballero, I., Fraile-Nuez, E., and Navarro, G.: Unmanned aerial vehicles (UAVs) as a tool for hazard assessment: The 2021 eruption of Cumbre Vieja volcano, La Palma Island (Spain), *Sci. Total Environ.*, 843, 157092, <https://doi.org/10.1016/j.scitotenv.2022.157092>, 2022.
- Shen, Y., Zhu, S., Yang, T., Chen, C., Pan, D., Chen, J., Xiao, L., and Du, Q.: BDANet: Multiscale Convolutional Neural Network With Cross-Directional Attention for Building Damage Assessment From Satellite Images, *IEEE T. Geosci. Remote*, 60, 1–14, <https://doi.org/10.1109/TGRS.2021.3080580>, 2022.

- Singh, D. K. and Hoskere, V.: Post Disaster Damage Assessment Using Ultra-High-Resolution Aerial Imagery with Semi-Supervised Transformers, *Sensors*, 23, 8235, <https://doi.org/10.3390/s23198235>, 2023.
- Spence, R. J. S., Pomonis, A., Baxter, P. J., Coburn, A. W., White, M., Dayrit, M., and Field Epidemiology Training Program Team: Building Damage Caused by the Mount Pinatubo Eruption of 15 June 1991, in: *Fire and Mud: Eruptions and Lahars of Mount Pinatubo, Philippines*, edited by: Newhall, C. G., and Punongbayan, R. S., University of Washington Press, London, UK, 1055–1061, ISBN 9780295975856, 1996.
- Spence, R. J. S., Kelman, I., Baxter, P. J., Zuccaro, G., and Petrazzuoli, S.: Residential building and occupant vulnerability to tephra fall, *Nat. Hazards Earth Syst. Sci.*, 5, 477–494, <https://doi.org/10.5194/nhess-5-477-2005>, 2005.
- SVG population and housing census: SVG population and housing census 2012, <https://stats.gov.vc/data/databases/> (last access: 8 December 2023), 2012.
- Szegedy, C., Vanhoucke, V., Ioffe, S., and Shlens, J.: Rethinking the Inception Architecture for Computer Vision, arXiv [preprint], arXiv:1512.00567, <https://doi.org/10.48550/arXiv.1512.00567>, 2015.
- Tennant, E., Jenkins, S. F., Miller, V., Robertson, R., Wen, B., Yun, S.-H., and Taisne, B.: UAVdamageAssessment_v1.0, Zenodo [code], <https://doi.org/10.5281/zenodo.14375616>, 2024.
- The MathWorks Inc.: MATLAB version: R2023b, Natick, Massachusetts, USA, <https://www.mathworks.com/products/matlab.html> (last access: 18 June 2024), 2023.
- Valentijn, T., Margutti, J., van den Homberg, M., and Laaksonen, J.: Multi-hazard and spatial transferability of a CNN for automated building damage assessment, *Remote Sens.*, 12, 1–29, <https://doi.org/10.3390/rs12172839>, 2020.
- Vetrivel, A., Gerke, M., Kerle, N., Nex, F., and Vosselman, G.: Disaster damage detection through synergistic use of deep learning and 3D point cloud features derived from very high resolution oblique aerial images, and multiple-kernel-learning, *ISPRS J. Photogramm. Remote Sens.*, 140, 45–59, <https://doi.org/10.1016/j.isprsjprs.2017.03.001>, 2018.
- Wang, Z., Zhang, F., Wu, C., and Xia, J.: Rapid mapping of volcanic eruption building damage: A model based on prior knowledge and few-shot fine-tuning, *Int. J. Appl. Earth Obs. Geoinf.*, 126, 103622, <https://doi.org/10.1016/j.jag.2023.103622>, 2024.
- Weber, E. and Kané, H.: Building Disaster Damage Assessment in Satellite Imagery with Multi-Temporal Fusion, arXiv [preprint], doi:10.48550/arXiv.2004.05525, 2020.
- Williams, G. T., Jenkins, S. F., Biass, S., Wibowo, H. E., and Harijoko, A.: Remotely assessing tephra fall building damage and vulnerability: Kelud Volcano, Indonesia, *J. Appl. Volcanol.*, 9, 1–18, <https://doi.org/10.1186/s13617-020-00100-5>, 2020.
- Wilson, G., Wilson, T. M., Deligne, N. I., and Cole, J. W.: Volcanic hazard impacts to critical infrastructure: A review, *J. Volcanol. Geoth. Res.*, 286, 148–182, <https://doi.org/10.1016/j.jvolgeores.2014.08.030>, 2014.
- Xu, J. Z., Lu, W., Li, Z., Khaitan, P., and Zaytseva, V.: Building Damage Detection in Satellite Imagery Using Convolutional Neural Networks, arXiv [preprint], <https://doi.org/10.48550/arXiv.1910.06444>, 2019.
- Yi, W., Sun, Y., and He, S.: Data Augmentation Using Conditional GANs for Facial Emotion Recognition, in: *Proceedings of Progress in Electromagnetics Research Symposium*, 1–4 August 2018, Toyama, Japan, 710–714, <https://doi.org/10.23919/PIERS.2018.8598226>, 2018.
- Yorioka, D., Kang, H., and Iwamura, K.: Data Augmentation for Deep Learning Using Generative Adversarial Networks, in: *IEEE 9th Global Conference on Consumer Electronics (GCCE)*, 13–16 October 2020, Kobe, Japan, 516–518, <https://doi.org/10.1109/GCCE50665.2020.9291963>, 2020.
- Yun, S.-H., Hudnut, K., Owen, S., Webb, F., Simons, M., Sacco, P., Gurrrola, E., Manipon, G., Liang, C., Fielding, E., Milillo, P., Hua, H., and Coletta, A.: Rapid damage mapping for the 2015 M_w 7.8 Gorkha Earthquake using synthetic aperture radar data from COSMO-SkyMed and ALOS-2 satellites, *Seismol. Res. Lett.*, 86, 1549–1556, <https://doi.org/10.1785/0220150152>, 2015.
- Zhang, J. F., Xie, L. L., and Tao, X. X.: Change Detection of Earthquake-damaged Buildings on Remote Sensing Image and its Application in Seismic Disaster Assessment, in: *Proceedings of IEEE International Geoscience and Remote Sensing Symposium*, 21–25 July 2003, Toulouse, France, 2436–2438, <https://doi.org/10.1109/igarss.2003.1294467>, 2003.
- Zou, Z., Shi, Z., Guo, Y., and Ye, J.: Object Detection in 20 Years: A Survey, *Proc. IEEE*, 111, 257–276, <https://doi.org/10.1109/JPROC.2023.3238524>, 2019.