

# Automating tephra fall building damage assessment using deep learning

Eleanor Tennant <sup>1</sup>, Susanna F. Jenkins <sup>2</sup>, Victoria Miller <sup>3</sup>, Richard Robertson <sup>4</sup>, Bihan Wen <sup>5</sup>, Sang-Ho Yun <sup>2</sup>, Benoit Taisne <sup>2</sup>

<sup>1</sup> Earth Observatory of Singapore @ NTU, Interdisciplinary Graduate Programme, Nanyang Technological University, Singapore, 639798

<sup>2</sup> Earth Observatory of Singapore and Asian School of the Environment, Nanyang Technological University, Singapore, 639798

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

<sup>4</sup> The UWI Seismic Research Centre, Saint Augustine, Trinidad, and Tobago

<sup>5</sup> School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore 639798

Correspondence to: Eleanor Tennant ([eleanorm001@e.ntu.edu.sg](mailto:leanorm001@e.ntu.edu.sg))

In the wake of a volcanic eruption, the rapid assessment of building damage is paramount for effective response and recovery planning. Uninhabited 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 ~50,000 building bounding boxes around ~2,000 individual buildings in 2,811 optical images collected during surveys conducted after the 2021 eruption of La Soufrière volcano, St Vincent and the Grenadines. We used this labelled data to train convolutional neural networks (CNNs) for: 1) Building localisation (average precision = 0.728); 2) Damage classification into two levels of granularity: No damage vs Damage (F1 score = 0.809); and Moderate damage vs Major damage, (F1 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 volcanic islands in the Caribbean where building types are similar, though 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

Deleted: optical images

Deleted: D

Deleted: all of

Deleted: shapefile

Deleted: Our

Deleted:

45 applicability worldwide. ~~To facilitate continued development and collaboration a~~l trained  
46 models and pipeline code can be downloaded from GitHub,

Deleted: A

Deleted: to facilitate collaboration and development

## 47 1 Introduction

48 Tephra fall produced by explosive volcanic eruptions can have detrimental effects on buildings,  
49 which in turn affects the ability for a community to recover and rehabilitate. These effects range  
50 from surface-level issues such as corrosion of metal roofs (e.g., Rabaul, Papua New Guinea,  
51 Blong, 2003a) or damage to non-structural components (e.g., gutters: Ambae, Vanuatu, Jenkins  
52 et al., 2024) through to complete building collapse (e.g., Pinatubo, Philippines, Spence et al.,  
53 1996).

Deleted: under review

Deleted: in press

54  
55 After, or during, an eruption, the collection of empirical data detailing the damage incurred is  
56 critical to guiding the planning and implementation of response and recovery efforts. This  
57 involves estimation of damages and losses, which are needed to determine the necessary  
58 funding for repair or reconstruction; along with an assessment of building functionality, which  
59 can inform temporary housing requirements. In addition to its use in post disaster recovery, the  
60 collection of damage data are key to the development of vulnerability models (Deligne et al.,  
61 2022), which relate hazard intensity to damage (e.g., Spence et al., 2005; Wilson et al., 2014;  
62 Williams et al., 2020), and can be used to inform resilient construction practises and/or for pre-  
63 event impact assessments.

Deleted: fragility functions

Deleted: ;

Deleted: Spence et al., 2021)

Deleted:

64  
65 Post-event building damage assessments usually consist of ground surveys, whereby the  
66 amount of damage to each building is described using a quantitative or qualitative damage state  
67 (e.g., Spence et al., 1996; Blong 2003a; Jenkins et al. 2013; Jenkins et al. 2015; Hayes et al. 2019;  
68 Meredith et al. 2022). However, tephra fall damage can extend tens or even hundreds of  
69 kilometres away from a volcano (Spence et al., 2005) meaning that comprehensive ground  
70 based damage assessments can be both time consuming and costly. Furthermore, the  
71 uncertainty that is often associated with the end of an eruption may prevent the safe completion  
72 of a ground-based damage assessment before tephra is remobilised by winds and rain. This lag  
73 between the event itself and the completion of a damage assessment, can hinder recovery  
74 efforts and compromise the accuracy of data collected for the development of forecasting  
75 models.

Deleted: ground based

76

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

Deleted: ), and

100  
101 To our knowledge, only one study attempts to automate the assessment of building damage  
102 from volcanic hazards (Wang et al., 2024). In contrast, attention has been given to more  
103 commonly occurring hazards such as earthquakes and hurricanes, with the development of  
104 both mono-temporal (post-event imagery only) and multi-temporal (images taken at different  
105 times) approaches (Table 1). Early approaches at automation with optical imagery used image  
106 processing methods, often focusing on identifying changes in pixel values between pre- and  
107 post-event imagery (e.g., Bruzzone and Fernández Prieto 2000; Ishii et al. 2002; Zhang et al.  
108 2003). Image processing methods are susceptible to user biases such as the choice of thresholds  
109 that equate to distinct levels of damage severity, or damage states, and may require  
110 recalibration when applied to a new dataset. As a result, image processing methods were  
111 succeeded by the application of traditional machine learning algorithms that use ‘handcrafted’  
112 image features. These features are observable properties that can be extracted from the image  
113 such as shape, colour, texture, and statistical properties of the image (e.g., Li et al. 2015;  
114 Anniballe et al. 2018; Lucks et al. 2019; Naito et al. 2020). The success of a given machine  
115 learning approach is dependent on the selection of the best features for the job; for example, a  
116 texture-based feature might be good for classifying buildings as damaged or not damaged due  
117 to an increased number of edges in damaged buildings but less useful for a task such as  
118 differentiating between building roof types where the difference in textures between the classes

Deleted: While efforts to automate the assessment of building damage from volcanic hazards are minimal (to our knowledge there has been one study focusing on building damage from volcanic eruptions: Wang et al., 2024), 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 (uses pre- and post-event imagery) approaches (Table 1).

Deleted: cats from dogs

130 is less significant. Deep learning, in particular the use of convolutional neural networks (CNNs),  
131 removes this need for feature selection. A CNN is a network of layers comprising filters which  
132 are small matrices of values. When an image is passed through the network, at each layer the  
133 filters are convolved with the output from the previous layer to create a new representation of  
134 the image that is progressively more abstract with depth in the network. This process reduces  
135 the image's original spatial dimensions (X and Y) while increasing the number of channels,  
136 facilitating classification. During network training the filter values (known as weights) are  
137 optimised to reduce the loss between the predicted label for the image and the true label.  
138 Through this training a CNN learns the features of the images that are useful for classification.  
139 For a detailed background on deep learning see Aggarwal, (2018).

140  
141 Thus far, deep learning models have been developed for optical image sets for hurricanes (Li et  
142 al. 2019; ~~a~~, Dung Cao and Choe 2020; Pi et al. 2020; Cheng et al. 2021; Khajwal et al. 2023);  
143 earthquakes (Nex et al. 2019; Xu et al. 2019; Duarte et al. 2020; Moradi and Shah-Hosseini  
144 2020); wildfires (Galanis et al. 2021); volcanic hazards (Wang et al., 2024); and models that  
145 have been proposed for multiple hazards (e.g., Gupta and Shah 2020; Weber and Kané 2020;  
146 Shen et al. 2021; Bouchard et al. 2022) (Table 1). However, building damage caused by different  
147 hazards looks very different (e.g., damage caused by vertical loading from volcanic tephra fall  
148 vs ground shaking from an earthquake). These observable differences mean that an optical  
149 imagery multi-hazard damage classification model that performs consistently well across the  
150 different hazards is not yet achievable. Therefore, distinct models tailored for specific hazards  
151 are required (Nex et al., 2019, Bouchard et al., 2022). It follows that models may also benefit  
152 from being regionalised, given the differences in building typologies (construction material and  
153 styles) that can also affect the observable damage (Nex et al., 2019).

154  
155 Many of the approaches for automating building damage assessment use both pre- and post-  
156 event imagery (Table 1), which makes the task more straightforward since any changes to the  
157 pre-event imagery can be considered damage. However, pre-event imagery at a high-enough  
158 resolution is not always available in post-disaster scenarios. The automated assessment of  
159 building damage from volcanic hazards using only post-event optical imagery has not yet been  
160 achieved in part due to absence of the large datasets that are needed in order to train models.  
161 The 2021 eruption of La Soufrière volcano, St Vincent and the Grenadines, provided  
162 ~~unprecedented~~ circumstances allowing for the collection of high-resolution UAV imagery

Deleted: ;

Deleted: an

Deleted: opportunity



166 enabling the development of fully automated models that can assess tephra fall building damage  
 167 from post-event data only. With their growing ubiquity and low cost, UAVs have become an  
 168 increasingly useful tool during and after volcanic eruptions (e.g., Andaru and Rau 2019; Gailler  
 169 et al. 2021; Román et al. 2022). UAVs offer a distinct advantage over satellite imagery because  
 170 they can be scheduled at any point, they do not suffer from cloud obscuring the images as they  
 171 fly at relatively low altitude, and they capture imagery from multiple perspectives, which may  
 172 lead to increased ability to capture damage information. In this study we used UAV optical  
 173 imagery collected after the 2021 eruption of La Soufrière volcano [to develop a methodology](#) for  
 174 tephra fall building damage assessment; the main contributions of our work are three-fold:

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

187  
 188  
 189 *Table 1. A non-exhaustive list of works using deep learning on optical imagery for building*  
 190 *damage assessment. Studies use different scores to evaluate performance: F1 scores are in*  
 191 *italics, mean average precision scores are underlined, accuracy scores in **bold**. For all scores, 1*  
 192 *represents a perfect model. [A detailed explanation of the scores used for evaluation is provided](#)*  
 193 *[in Section 2.3.3.](#)*

Study	Hazard	Number of damage classes	<u>Pre-disaster imagery</u>	Data type	Building localisation	Damage classification
Li et al. (2019a)	Hurricane	2	<u>No</u>	airborne		<i>0.448</i>
Weber and Kane, (2020)	Multi	4	<u>Yes</u>	satellite (xBD)	<i>0.835</i>	<i>0.697</i>

Deleted: building damage maps

Deleted: Pre and post?

Formatted Table

Deleted: P

Deleted: P & P

Dung Cao and Choe. (2020)	Hurricane	2	No	satellite	-	<b>0.972</b>
Pi et al. (2020)	Hurricane	2	No	UAV, airborne	0.745 (UAV) 0.807 (airborne)	
Cheng et al. (2021)	Hurricane	5	No	UAV	0.656	<b>0.610</b>
Galanis et al. (2021)	Wildfire	2	No	satellite		0.981
Gupta and Shah (2020)	Multi	4	Yes	satellite (xBD)	0.840	0.740
Shen et al. (2021)	Multi	4	Yes	satellite (xBD)	0.864	0.782
Bouchard et al. (2022)	Multi	2	Yes	satellite (xBD)	0.846	0.709
Khajwal et al. (2023)	Hurricane	5	No	ground airborne	-	0.650
Singh and Hoskere, (2023)	Multi	5	No	satellite		<b>0.880</b>
Wang et al (2024)	Volcanic tephra	4	Yes	satellite	0.868	0.783

Deleted: P

Deleted: P

Deleted: P

Deleted: P

Formatted: Font: Cambria, 10 pt

Deleted: P & P

Deleted: P & P

Deleted: P & P

Deleted: P

Deleted: P

Deleted: P & P

Deleted: Our work

... [1]

### 1.1 The 2020-2021 eruption of La Soufrière volcano St Vincent

La Soufrière St Vincent is an active stratovolcano standing at 1220 meters above sea level on the island of St Vincent. On 27<sup>th</sup> 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 three months until 9 April 2021, when, following two days 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 kilometres 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 centimetres in coastal communities to the north of the island (Cole et al. 2023) (Figure 1).

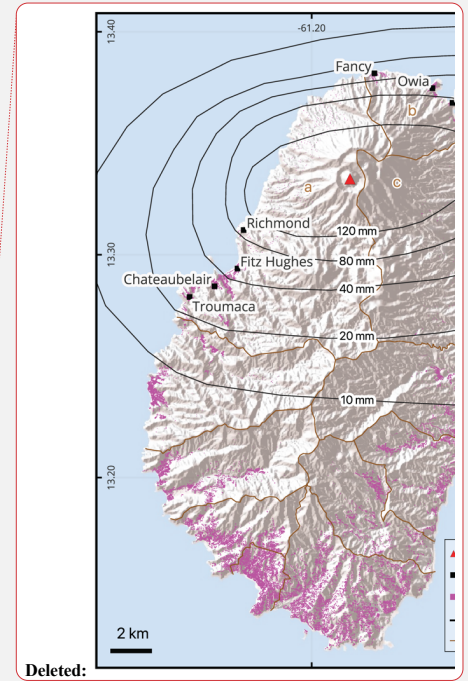
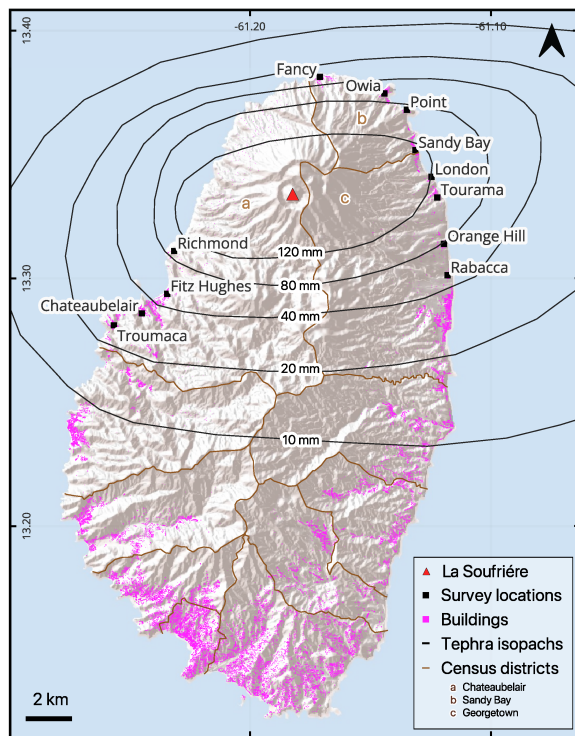
Deleted: ion

The explosive phase was anticipated, and an evacuation order was issued on 8 April 2021 for the ~16,000 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 were 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,

Deleted:

233 detached structures, with the majority in the more impacted north of the island (census  
 234 districts of Chateaubelair, Georgetown, and Sandy Bay: Figure 1) constructed using concrete  
 235 and blocks (84% in Chateaubelair, 74% in Georgetown, 50% in Sandy Bay), with sheet metal  
 236 roofs (90-92% of all buildings in these areas) (SVG population and housing census, 2012).

237



238

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

244

## 245 2 Method

246 After the 2021 eruption of La Soufrière three UAV optical imagery datasets were collected to  
 247 assess the extent of the damage. These were collected by different parties at separate times after

251 the eruption. All UAV survey locations are shown in [Figure 1](#), and representative examples of  
252 images can be found [in](#) Section S1 of the supplementary material.

Deleted: Figure 1

Formatted: Font: Not Italic

253

## 254 **2.1 Dataset description**

### 255 **Dataset 1: April-May 2021 (UWI-TV)**

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

Deleted: .

262

### 263 **Dataset 2: 12<sup>th</sup> – 14<sup>th</sup> May 2021 (GOV)**

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

271

### 272 **Dataset 3: August -September 2021 (SRC)**

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

280

281 For all three datasets, image frames were extracted from the videos every two seconds, an  
282 interval chosen to reduce redundant homogeneous images, this resulted in a total of 7,956  
283 image frames. Due to the UAV surveying approach (i.e., hovering in one place for a while) many

286 near-identical images were generated. To avoid potentially biasing the training towards  
 287 overrepresented buildings we manually filtered out duplicate images. After filtering, and the  
 288 removal of images with no buildings present, the full combined dataset consisted of 2,811 image  
 289 frames. We labelled all images by drawing bounding boxes around each building present and  
 290 storing the bounding box positions. In total 49,173 building bounding boxes were drawn around  
 291 ~2,000 individual buildings (with some buildings being present in multiple images). Given the  
 292 absence of individual building location information, this number was approximated by  
 293 overlaying Open Street Map building footprints with UAV GPS tracks where available. Bounding  
 294 boxes were drawn by a team of five including the lead author, and all boxes were checked by the  
 295 lead author. Each box was then assigned one of three damage states, which are described below.  
 296 For consistency the damage states were assigned by the lead author. All labelling, modelling,  
 297 and analysis were conducted using MATLAB 2023b.

## 299 2.2 Developing and applying a building damage state framework

300 The first tephra fall building damage state framework was developed after the eruption of  
 301 Pinatubo, Philippines, 1991 (Spence et al., 1996), and was adapted from the macro seismic  
 302 intensity scale used to evaluate seismic damage (Karnik et al., 1984). In the adapted framework  
 303 damage ranges from DS0 - "no damage", through to DS5 - "complete roof collapse and severe  
 304 damage to the rest of the building". Subsequent tephra fall building damage state frameworks  
 305 were modified from the work of Spence et al., (1996) with changes in the wording made to  
 306 reflect the characteristics of the case study (Table 2). In the damage state descriptions, damage  
 307 to three critical aspects of a building is described: the roof covering, the roof structure, and the  
 308 vertical structure (Blong 2003b; Hayes et al. 2019; Jenkins et al., 2024). In our study, most  
 309 images depict buildings from an at nadir or close to nadir perspective making roof damage more  
 310 discernible than damage to the vertical structure. Thus, we generated a damage state  
 311 framework that is based on the proportion of observable damage to the roof, as in the work of  
 312 Williams et al. (2020). Our final framework, which was developed over several iterations,  
 313 classifies building damage into three classes: No observable damage to minor damage,  
 314 Moderate damage, and Major damage (Table 3, Figure 2). Damage states are deliberately  
 315 generic so that the range of possible damage to the range of different building types can be  
 316 captured (Blong, 2003a). Our three classes are comparable to DS0-1, DS2, and DS3-5,  
 317 respectively, of damage scales developed for ground surveys (Table 2). In the frameworks  
 318

Deleted: detailed building inventory information

Formatted: English (UK)

Deleted: Bounding boxes were drawn by a team of five including the lead author, and all boxes were checked by the lead author. smay. Nevertheless, this was not considered an issue since deep learning models for object localisation will quickly learn to ignore background pixels (Bouchard et al., 2022).

Deleted: modeling

Deleted: to describe damage from

Deleted: developed

Deleted: For tephra fall,

Deleted: D

Deleted: d

Deleted: consisted of the following classes:

Deleted: D1: "Light roof damage", D2: "Moderate roof damage", D3: "Severe roof damage and some damage to vertical structure", D4: "Partial roof collapse and moderate damage to rest of building",

Deleted: split damage into five damage states, plus one not damaged, based on

Deleted: d

Deleted: (Spence et al., 1996;

Deleted: under review

Deleted: in press

Deleted: Ground based damage state frameworks for tephra fall have previously split damage into five damage states, plus one not damaged, based on damage to three critical aspects of a building: the roof covering, the roof structure, and the vertical structure (Spence et al., 1996; Blong 2003b; Hayes et al. 2019; Jenkins et al., under review). Remote damage assessments are often less able to resolve the detailed resolution achievable on the ground, and so a coarser resolution damage state framework is needed....

Deleted: e or

Deleted: Classes are comparable to DS0-1, DS2, and DS3-5 of damage scales developed for ground surveys respectively (Table 2) . (

Deleted: (Table 2).

Deleted: damage state

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 better grouped with undamaged buildings. Our Moderate damage class requires damage or collapse to up to 50% of the roof area, which closely fits with damage state 2 of Blong, (2003), 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'.

Deleted: very

Deleted: W

Deleted: discernable

Deleted: therefore

Deleted: (

Deleted: in press

Deleted: ¶

Deleted: ¶

In some images tarpaulins can be seen partially or fully covering roofs (~30 buildings). These were potentially placed to cover damage that occurred during the eruption, including corrosion due to prolonged presence of tephra on metal roofs or, holes generated by nails lifted out through sub-optimal cleaning approaches (VM personal communication). Alternatively, tarpaulins may have been placed as a preventative measure to help shed tephra (e.g., Ambae Vanuatu, Jenkins et al., under review 2024). Erring on the conservative side, we considered buildings with a tarpaulin to be damaged; we assessed the severity of the damage for each building based on the level of visible deformation. We assigned buildings with a tarpaulin and no visible deformation to the moderately damaged class and those with a tarpaulin and visible deformation to the major damage class. ¶

Deleted: 3

Formatted: Centered

Formatted Table

Formatted: Centered

Formatted: Centered

Formatted: Centered

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, (2003)	Calbuco, Chile, 2015 Hayes et al., (2019)	Manaro Vuoi, Ambae island, Vanuatu, 2017-2018 Jenkins et al., (2024)
DS0	<b>No damage</b>		<b>No damage</b>	<b>No damage</b>
DS1	<b>Light roof damage:</b> - Gutter damage. - Few tiles dislodged.	<b>Light damage:</b> - Damage to gutters and/or water tanks. - Cleanup required	<b>Minor damage to non-structural elements:</b> - Damage to gutters. - Few tiles dislodged. - Damage to fittings, e.g. air-conditioning units and appliances. - Damage to contents. - Dents in the roof covering.	<b>Light damage or damage to non-structural elements:</b> - Damage to gutters. - Damage to contents. - Dents or minor slumping in roof cover.
DS2	<b>Moderate roof damage:</b> - Bending or excessive deflection of roof sheeting or purlins. - No damage to principal roofing supports.	<b>Moderate damage:</b> - 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.	<b>Moderate damage but vertical structure and roof supports intact:</b> - 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.	<b>Moderate damage but vertical structure and roof supports intact:</b> - 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.

		and/or overhaul of electrical systems. - <u>Solar heater needs replacing.</u>	- Little to no damage to principal roof supports, i.e. rafters or trusses. - <u>Damage to roof overhangs or verandas.</u>	- <u>Damage to roof overhangs or verandas.</u> - <u>Interior requires repair.</u>
DS3	<b><u>Severe roof damage and some damage to vertical structure:</u></b> - <u>Severe damage or partial collapse of roof overhangs or verandas.</u> - <u>Severe deformation of main roof sheeting.</u> - <u>Some damage to roof supporting structure, columns, trusses.</u>	<b><u>Heavy damage:</u></b> - <u>Damage to roof structure and some damage to walls.</u> - <u>At least one wall damaged/misaligned.</u> - <u>Collapse of part of ceiling</u>	<b><u>Severe damage to the roof and supports:</u></b> - <u>As above.</u> - <u>Bending or excessive (e.g., perforation, cracking) damage (with or without collapse) to over half of roof covering.</u> - <u>Damage to any single principal roof supports and some damage to walls.</u> - <u>Severe damage or partial collapse of roof overhangs or verandas.</u>	<b><u>Severe damage to the roof and supports:</u></b> - <u>As for DS2, plus:</u> - <u>Bending or excessive damage (with or without collapse) to more than half of the roof covering.</u> - <u>Damage to any single principal roof supports and/or some damage to walls (less than half of walls affected).</u> - <u>Severe damage or partial collapse of roof overhangs or verandas.</u>
DS4	<b><u>Partial roof collapse and moderate damage to rest of building:</u></b> - <u>Collapse of sheeting but not truss.</u> - <u>Partial collapse of sheeting and some truss failure.</u> - <u>Failure of supporting structure.</u> - <u>Moderate damage to other parts of building resulting from roof collapse.</u>	<b><u>Severe damage:</u></b> - <u>Roof collapse and moderate to severe damage to rest of the building.</u> - <u>Failure of roof trusses and supporting structure.</u> - <u>At least half of the external walls and/or internal walls deformed or collapsed.</u> - <u>For two-storey buildings, collapse of external and internal walls of upper floor.</u> - <u>Plumbing and other services may be damaged.</u>	<b><u>Partial or total collapse of the roof and supports:</u></b> - <u>As above</u> - <u>Collapse of roof covering and any single principal roof support(s).</u> - <u>At least half of the external walls and/or internal walls deformed or collapsed.</u>	<b><u>Partial collapse of the roof and supports:</u></b> - <u>As for DS3, plus:</u> - <u>Collapse to less than half of roof covering and principal roof support(s).</u> - <u>At least half of external and/or internal walls deformed or collapsed.</u>
DS5	<b><u>Complete roof collapse and severe damage to the rest of the building:</u></b> - <u>Collapse of roof and supporting structure over more than 50 percent of roof area.</u>	<b><u>Collapse:</u></b> - <u>Collapse of roof and supporting external walls over more than 50% of floor area of building.</u> - <u>Internal walls collapsed.</u> - <u>Damage to floor and/or foundation.</u> - <u>Structure is irreparable, not</u>	<b><u>Building collapse:</u></b> - <u>As above.</u> - <u>Collapse of roof, principal roof supports and/or supporting external walls over &gt;50% of floor area of building.</u>	<b><u>Building collapse:</u></b> - <u>As for DS4, plus:</u> - <u>Collapse of roof, principal roof supports and/or supporting external walls over more than half of floor area of building.</u>



- Partition walls destroyed.
  - External walls destabilized.
- salvageable, beyond economic repair.

398

399

400

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

<b>Damage state</b>	<b>Description of the damage</b>
<b>No damage to minor damage</b>	<ul style="list-style-type: none"> <li>- No visible damage/or</li> <li>- Up to 10% of the roof covering missing; and/or</li> <li>- No roof or structural collapse; and/or.</li> <li>- <u>Visible damage to non-structural elements e.g., gutters or decorative elements (fascia).</u></li> <li>- <u>Comparable to DS0-1 (Table 2).</u></li> </ul>
<b>Moderate damage</b>	<ul style="list-style-type: none"> <li>- <u>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).</u></li> <li>- <u>Comparable to DS2 (Table 2).</u></li> </ul>
<b>Major damage</b>	<ul style="list-style-type: none"> <li>- <u>More than 50% roof area damaged or collapsed; may include damage to the vertical structure including total building collapse.</u></li> <li>- <u>Comparable to DS3-5 (Table 2).</u></li> </ul>

401

Deleted: 7

Deleted:

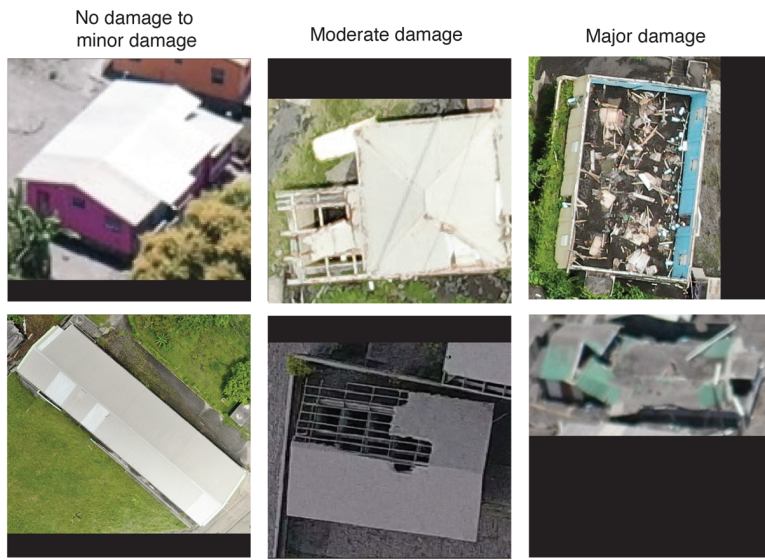
Deleted: 2.

Deleted: assessment

Formatted Table

Formatted: Centered





406  
407 *Figure 2. Example of the three damage states used in this work: No damage to minor damage,*  
408 *Moderate damage and, Major damage.*

Formatted: Font: Italic

Formatted: Left

### 410 2.3 Model development

411 After labelling, we split the full combined image dataset (2,811 frames from the UWI-TV, GOV  
412 and SRC sets) into train/validation/test sets (Figure 3). Given that many images lacked GPS  
413 positions, we grouped images by location to ensure independence among the sets. The  
414 partitioning was chosen to include diversity in both the image sets (UWI-TV/GOV/SRC) and in  
415 the location, which affects the tephra fall thickness. We aimed for a standard data split of  
416 80%/10%/10%, for train/validation/test, however given the above constraints, this produced  
417 a split of 80/8/12, (considering the number of bounding boxes and not the number of images).  
418 These datasets were used to develop our approach for building damage assessment. In line with  
419 studies shown in Table 1, we chose to split the damage assessment task into two subtasks: i)  
420 building localisation (i.e. identification of building bounding boxes within the images) and ii)  
421 damage classification. While it is possible to develop a model that can simultaneously locate  
422 and classify buildings with different levels of damage, model training under this approach can  
423 take significantly more time and resources to converge when compared to an approach that  
424 splits the tasks (Bouchard et al., 2022). Furthermore, decoupling the two tasks allows for  
425

Deleted: 2

Formatted: Font: Cambria

Deleted: a sizable proportion of the data were not geotagged, images from each location were kept together to assure the train and test sets were independent.

Deleted: thickness of

Deleted: received

Deleted: with the majority of data assigned for training

Deleted: % train,

Deleted: 8% validation, and

Deleted: % test

436 greater flexibility; for example, if building locations are already known then only the  
437 classification can be run, speeding up the remote assessment.

439 In machine learning, the performance of a model and its optimal hyperparameters can be highly  
440 dependent on the characteristics of the dataset used for training, and hyperparameters that  
441 work well for one dataset may not work well for another. Therefore, it's common practice to  
442 optimise hyperparameters, model architectures, and training strategies to find the  
443 configuration that performs the best for a particular problem. For building localisation and  
444 damage classification we conducted a series of independent experiments using different image  
445 preprocessing approaches, CNN architectures, and combinations of hyperparameters with the  
446 aim of iterating towards the best experimental setup (Model selection: Section 3.1.1; Section  
447 3.2.1). Each experiment consisted of three replicates of a given combination of these aspects.  
448 Replicates were conducted since the stochastic nature of the training process can cause models  
449 to converge at slightly different points (Aggarwal, 2018). For each experiment the replicate with  
450 the highest evaluation metric was the one compared against the other experiments.

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

456 Following model selection and cross validation we calculated the performance of the best model  
457 identified for each task on the test set. Finally, to see if better performance could be achieved  
458 with more data available for training, we retrained the models on the combined training and  
459 validation data before evaluating on the test data (Evaluation on the test set: Section 3.1.3,  
460 Section 3.2.3). All stages of model development, including model selection, cross validation, and  
461 final evaluation, are shown in Figure 4 and more information about the specific experiments  
462 conducted for model selection is given in Section S3 of the supplementary material.

464 Past studies have trained deep learning algorithms on georeferenced images (i.e. each pixel has  
465 a geographical location attached) (Gupta and Shah, 2020; Shen et al., 2021; Bouchard et al.,  
466 2022) and non-georeferenced images (e.g. Li et al., 2019a; Pi et al., 2020; Cheng et al., 2021). In  
467 this work we labelled the non-georeferenced images and trained models on these. This was  
468 done firstly to preserve the multiple viewing angles that we have of each building with each

Deleted: Most previous studies have 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 (Table 1). Developing a model that can simultaneously locate and classify buildings with different levels of damage is feasible, however, model training under this approach can take a lot of time and resources to converge (Bouchard et al., 2022). Furthermore, decoupling the two tasks makes the approach more flexible in a post-disaster context, for example, if building locations are already known then only the classification can be run, speeding up the remote assessment. For these reasons, we split the building damage assessment task into two subtasks.

Deleted: : In line with the work of previous authors (Cheng et al. 2021; Bouchard et al. 2022), we split the building damage assessment task into two subtasks... [2]

Deleted: deep

Deleted: experiment with different

Deleted: z

Deleted: settings

Deleted: each task in our damage assessment approach... [3]

Deleted:

Deleted: One

Formatted: Not Highlight

Deleted: Three

Deleted: r

Deleted: ¶

Deleted: (building localisation, classification 1, ... [4]

Deleted: we combined the training and validation ... [5]

Deleted:

Deleted: To test the robustness to location, we train... [6]

Deleted: final

Deleted: (building localisation, classification 1, ... [7]

Deleted: 4

Deleted: 3

Deleted: 2

Deleted: While studies have trained deep learning ... [8]

Deleted: Past studies have

Deleted: machine

Formatted: English (UK), Not Highlight

Formatted: English (UK), Not Highlight

Deleted: here we opted to

Deleted: train models

Deleted: on

Deleted: raw

Deleted: non-rectified

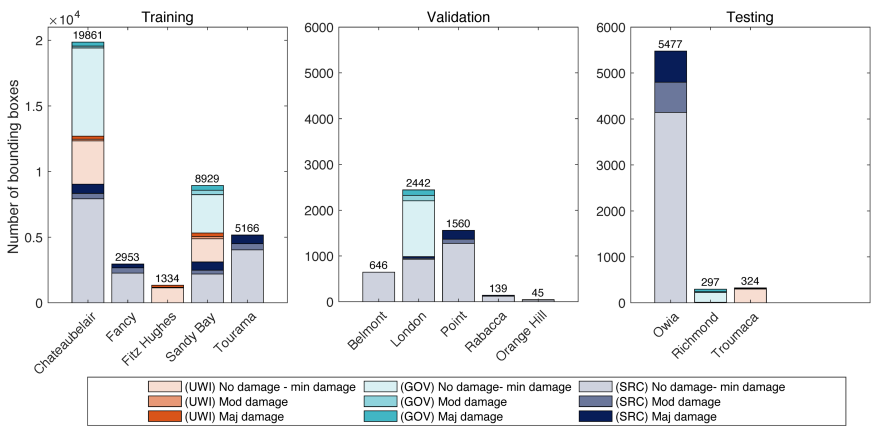
Deleted: due to the absence of GPS location on a lat... [9]

Deleted: ,

547 image counting as a different data point, and secondly, due to the absence of GPS locations on a  
 548 large portion of the dataset. In an operational context, spatial information must be tied to the  
 549 assessed damage. Therefore, beyond the creation of distinct models for each task, we designed  
 550 a comprehensive, fully automated pipeline that integrates models for building localisation and  
 551 damage classification. Our pipeline contains all the necessary processing steps to guide images  
 552 through the separate models enabling them to operate on a georeferenced orthomosaic image  
 553 (to be generated separately) or on non-georeferenced images. When applied to an orthomosaic  
 554 image the output from the pipeline is a georeferenced vector dataset that can readily be plotted  
 555 in a GIS to generate damage maps.

556  
 557 In Section 4 we apply the pipeline to assess building damage in Owia, St Vincent, which received  
 558 50-90,mm of tephra fall during the 2020-2021 eruption (Figure 1). Owia was selected out of  
 559 the three possible test set locations (Figure 3) due to its large size and the existence of GPS  
 560 locations that enabled the generation of a georeferenced orthomosaic image; for this we used  
 561 Agisoft Metashape software. To compare the assessed building damage with tephra thickness,  
 562 we used the TephraFits code (Biass et al., 2019) to identify the theoretical maximum  
 563 accumulation using the isopachs from Cole et al., (2023). This maximum accumulation and the  
 564 isopachs were interpolated using cubic splines and the surface was exported at a resolution of  
 565 10 m to provide a tephra thickness value for each building.

Deleted: which would otherwise be lost by converting and training on an orthomosaic image.  
 Deleted: I  
 Deleted: In a post-disaster context, the seamless functioning of models will benefit from a sequential workflow.  
 Deleted: B  
 Deleted: w  
 Deleted: all  
 Deleted: of  
 Deleted:  
 Deleted: allowing  
 Deleted: ).  
 Deleted: that executess all optimized final models in turn.  
 Deleted: to produce an  
 Deleted: and the required processing steps to guide images through the models (Figure 3d). The pipeline runs on an orthomosaic image and generates spatial data in shapefile format  
 Deleted: our  
 Deleted: in  
 Deleted: is the islandSt Vincent and  
 Deleted: 8  
 Deleted: ; to do this we used  
 Deleted: which enabled the generation of an orthomosaic image. We generated an orthomosaic image using Agisoft Metashape software  
 Deleted: In the following sections we provide more detail on the algorithms and architectures used for each of the tasks, and how the performance of each task was evaluated.



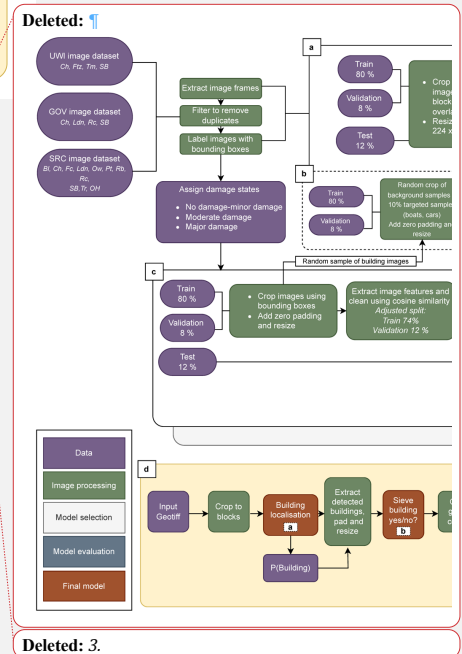
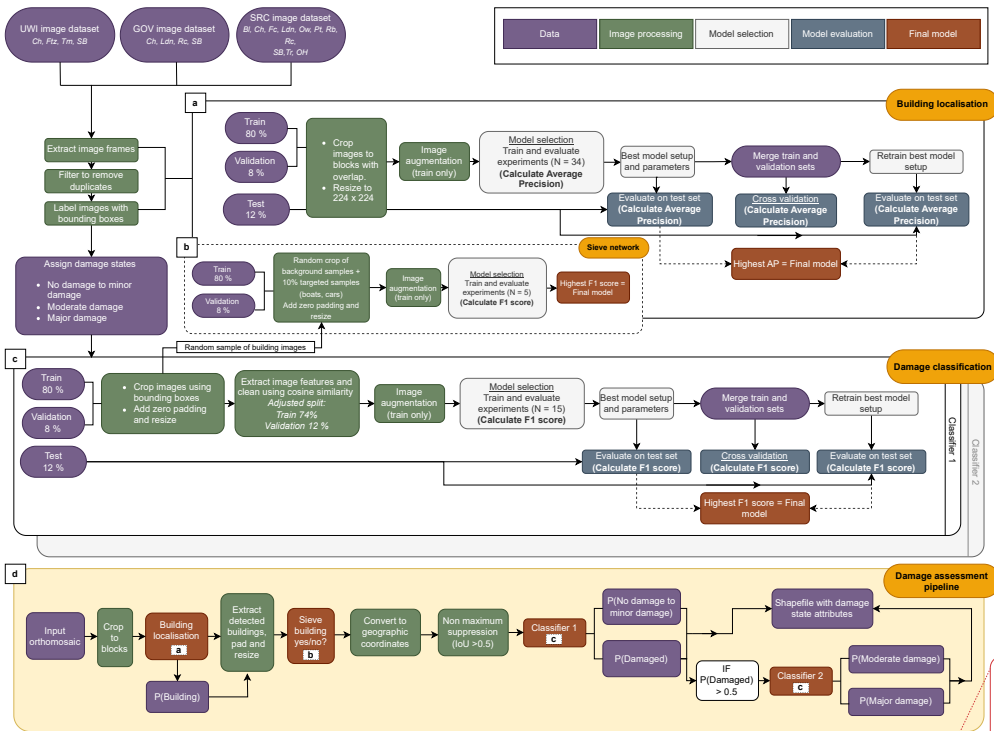
567  
 568

601

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

Deleted: 2

Deleted: 2



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

622

### 623 2.3.1 Building localisation

624

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

630

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

639 To improve the performance of the building localisation model we developed a sieve network  
640 that runs as an add on to the Faster R-CNN building detector. The sieve network reduces false  
641 positives which occur when the detector predicts a bounding box that does not have an  
642 overlapping labelled building (i.e., detects a building when there is not one). More details on its  
643 development are provided in Section 3.2 of the supplementary material.

644

### 645 2.3.2 Damage classification

646 We chose to divide building damage classification into two separate classifications. Classifier 1  
647 distinguishes between 'No damage to minor damage' versus the combined classes of 'Moderate  
648 damage' and 'Major damage', while Classifier 2 further differentiates between 'Moderate  
649 damage' and 'Major damage'. A hierarchical approach to classification has been found effective  
650 when the number of samples is limited or classes are unbalanced (Li et al. 2019b; An et al.  
651 2021). We conducted experiments separately for Classifiers 1 and 2. Experiments consisted of  
652 fine-tuning two different pretrained CNNs to determine which was better and should be used  
653 in the final models for each classifier: ResNet50 (He et al., 2015) trained on the ImageNet  
654

655

Deleted: Figure 3. 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 building damage classification and d) the building damage assessment pipeline developed in this work. The pipeline incorporates the final trained models for building localisation and two stages of building 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, Tm- Troumaca.

Deleted: conducted experiments using

Deleted: Faster R-CNN is an improvement on the Fast R-CNN algorithm proposed by Girshick, (2015). The improvement comprises an initial region proposal network (RPN) which speeds up performance. Initially In Faster R-CNN, image feature maps are extracted by passing the input image through a pretrained backbone CNN. The RPN then utilizes these features to generate proposals for potential object-containing areas, th...

Deleted: patches

Deleted: the size of these patches blocks, and the amount of overlap between patches blocks, and whether bl...

Deleted: fixed given for each experiment

Deleted: for building localisation

Deleted: see

Deleted: . More information on the values used for training can be found in the

Deleted: for details

Deleted: <#>Developing a sieve network

Deleted: small

Deleted: . Bounding boxes produced by the detector are passed to the sieve network to filter out detections...

Deleted: s

Deleted:

Deleted: Further information on the sieve networks

Deleted: can be found

Deleted: To develop the dataset used for training and evaluating the sieve network we randomly cropped...

Deleted: Building d

Formatted: Font: Times New Roman, Not Bold, English (UK)

Formatted: Normal

Deleted: .

Deleted: For building damage classification,

Deleted: w

Deleted: c

Deleted: This hierarchical approach to classification has been found effective when the number of samples...

748 dataset (Deng et al. 2009), and GoogleNet (Szegedy et al., 2015) trained on the places365  
749 dataset (López-Cifuentes et al., 2019). Fine-tuning is a common approach to computer vision  
750 tasks where sufficiently large, labelled datasets are not available for the task at hand (typically  
751 hundreds of thousands of images are needed: Aggarwal, 2015). During fine-tuning, the high-  
752 level features that were learnt during the initial training on the large dataset can be leveraged  
753 for the new task. In addition to the different pretrained CNNs used, experiments also considered  
754 different ways of balancing the number of images for each damage state class (over-sampling  
755 the minority class, under-sampling the majority class and no balancing). When applied to a test  
756 building image, the trained classifier outputs the highest probability class and the associated  
757 probability. A total of 15 experiments were conducted for each of the classification tasks. For  
758 each experiment three replicates were conducted, each consisting of a grid search to find the  
759 best combination of learning rate, batch size and L2 regularisation. For more information on  
760 this see Section 3.3 of the supplementary material.

### 762 2.3.3 Model evaluation metrics

763 For building localisation Faster R-CNN experiments, we evaluated performance using the  
764 average precision (AP) at an intersection over union (IoU) threshold of 0.5, and the F1 score.  
765 AP, a common metric for evaluating object detection (Zou et al., 2019), measures how often the  
766 detector gets it right (true positives, TP) versus wrong (false positives, FP, and false negatives,  
767 FN). A TP occurs when a predicted box overlaps a labelled box by more than 50% (IoU > 0.5), a  
768 FP when there is no overlapping labelled box, and a FN when the detector misses a labelled box.  
769 When the detector is run on a test image a confidence score is output for each predicted box (0-  
770 1). Once the trained detector has been run over the full test set, the precision ( $TP/(TP+FP)$  and  
771 recall ( $TP/(TP+FN)$ ) are calculated at different confidence score thresholds and the area  
772 underneath the resulting precision-recall curve represents the AP. AP depicts the trade-off  
773 between precision and recall and provides an overall measure of detection performance. AP  
774 values range between 0-1, where a higher value indicates a better performance.

776 For building localisation, the F1 score was calculated at IoU and confidence thresholds of 0.5.  
777 The F1 score is calculated as:  $F1 = 2x (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$ . To evaluate the  
778 performance of classification models, we used the macro-F1 score, which is the unweighted  
779 mean of the F1 scores calculated for each of the classes. Similarly to the AP, values of the F1  
780 score range between 0-1, where a higher value indicates a better performance.

Deleted: the

Deleted: The AP is the most frequently used measure of an object detector's performance (Zou et al., 2019), and is calculated based on the number of times the detector gets it right (a true positive, TP) or wrong (a false positive, FP or a false negative, FN).

Deleted: A true positive occurs when the detector predicts a box that has an IoU with a labelled box of > 0.5.

Deleted: A false positive occurs when the detector predicts a bounding box that does not have an overlapping labelled box, while a false negative occurs where the detector fails to predict a box that is present in the labelled data

Deleted: The relative proportions of these are used to calculate the precision and recall, where precision is the number of things that were predicted as positive that were correct:  $\text{Precision} = TP/(TP+FP)$ , and recall is the number of things that are truly positive that were identified:  $\text{Recall} = TP/(TP+FN)$ .

Deleted: a

Deleted: (

Deleted: )

Deleted: (

Deleted: )

Deleted: which can be plotted against one another to form a curve. The AP is

Deleted:

Deleted: this

Deleted: ;

Deleted:

Deleted: it

Deleted: this

Deleted: ,

Deleted: ¶



815 **3 Results**

816 **3.1 Building localisation**

817 **3.1.1 Model selection**

818  
819 The five experiments with the highest average precision are shown in Table 4, with the full list  
820 of experiments provided in Table S2 of the supplementary material. Average precisions across  
821 the 34 experiments ranged from 0.295 to 0.701 (Table 4 and Table S2). We found that block size  
822 played an important role in model performance; out of the 34 experiments conducted, the top  
823 three used a block size of 550 x 550 pixels, which was the middle of the sizes tested (450, 550,  
824 650). We observed that models trained on the full dataset performed better than models trained  
825 separately for the nadir (SRC) and off-nadir imagery sets (UWI-TV and GOV sets combined)  
826 (Table 4 and Table S2).

827  
828 *Table 4, Hyperparameters for the five experiments with the highest average precision conducted*  
829 *for building localisation, ordered by average precision. The full table consisting of all 34*  
830 *experiments is provided in the supplementary material. Columns marked with '\*' contain Yes/No*  
831 *information. Training dataset \*\*: a= all, b= UWI-TV and GOV, c= SRC.*

Row ID	Block size	Mixed block size*	Block overlap	Block resized*	Training dataset**	Max Average Precision	F1 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

833  
834 All trained sieve networks achieved macro and class F1 scores that were > 0.973 (Table S3,  
835 Supplementary material). The sieve networks efficacy at improving building localisation is  
836 demonstrated by comparing the results of the best detector when applied to the validation  
837 dataset pre-sieving (Table 4 row ID 1) with the post-sieving results. Pre-sieving there were a  
838 large number of false positive detections, resulting in a precision of 0.588, post-sieving these  
839 were reduced and the precision increased to 0.695 (Table 5).

Deleted: top  
Deleted: (  
Deleted: )  
Deleted: conducted for building localisation  
Deleted: 3  
Deleted: 3

Deleted: 3  
Deleted: More details on the results of experiments run for building localisation model selection can be found in Section S2.1 of the supplementary material.  
Deleted: 3  
Deleted: 7  
Deleted: highest scoring (average precision)

Formatted Table  
Deleted: All training/UWI-TV&GOV/SRC  
Deleted: id  
Deleted: ?  
Deleted: ?  
Deleted: all  
Deleted: all  
Deleted: all  
Deleted: all  
Deleted: all  
Deleted: all  
Deleted: <#>Sieve Network  
Deleted: <#>Table 5  
Deleted: <#>4  
Deleted: <#>The best performing sieve network experiment achieved a macro F1 score of 0.977.  
Deleted: <#>The best detector identified through model selection (Table 43, row 1) achieved an F1 score of 0.669 (Table 76), with a precision and recall of 0.588 and 0.776, respectively, on the validation data. The lower value of precision is due to the substantial number of false positive detections. After the results of the detector were passed through the sieve network, the number of false positives was reduced, with an improved F1 score of 0.712 (Table 76).



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

	Precision	Recall	F1
Best detector pre-sieving	0.588	0.776	0.669
Best detector post-sieving	0.695	0.730	0.712

881  
882  
883

### 884 3.1.2 Cross validation

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

890 We found that the performance of the selected object detector varied, depending upon the  
 891 location (Figure 5a) or imagery dataset (Figure 5b) used for testing. For models tested on  
 892 different locations average precisions in line with the AP achieved on the full validation set  
 893 (0.701) were obtained for Point and Fancy (Figure 5a). The lowest AP values were for London  
 894 (0.063) and Fitz Hughes (0.187). The standard deviation (SD) (Figure 5) shows the variability  
 895 in performance between the three replicates that were trained for each test, which arises due  
 896 to the stochastic nature of the training process. For models tested on the different imagery  
 897 datasets individually the AP was low, with a mean value across all datasets of < 0.2 (Figure 5b).  
 898 For all three locations (Chateaubelair, Sandy Bay, London), AP for models evaluated on the SRC  
 899 dataset were lower than for the UWI-TV or GOV datasets.

900

Deleted: 4

Deleted: Experiments conducted for the sieve network, a small network that runs on the boxes produced by the object detector. Results are ordered from high to low by the Macro F1 score. ResNet50 and GoogleNet refers to the convolutional neural network architecture used in the experiment; the value after the underscore reflects the experiment ID where different IDs have different training parameters (see Section S2 of the supplementary material).

Formatted: Justified

Deleted: Experiment ID

Formatted: Font: Bold

Formatted Table

Formatted: Font: Bold

Formatted: Font: Bold

Deleted: , along with the potential for the model to generalize to a new dataset.

Deleted: s

Deleted: 4

Deleted: 4

Deleted: (Figure 4a)

Deleted: > 0.7 were

Deleted: 4

Deleted: in line with AP achieved on the full validation set (0.701)

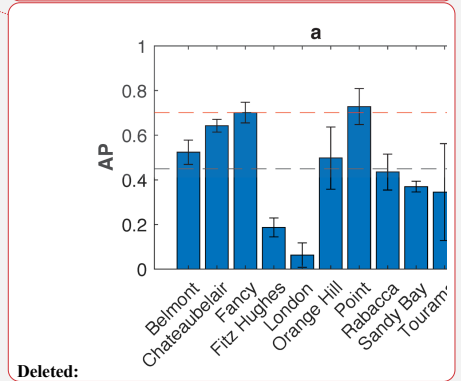
Deleted: 4

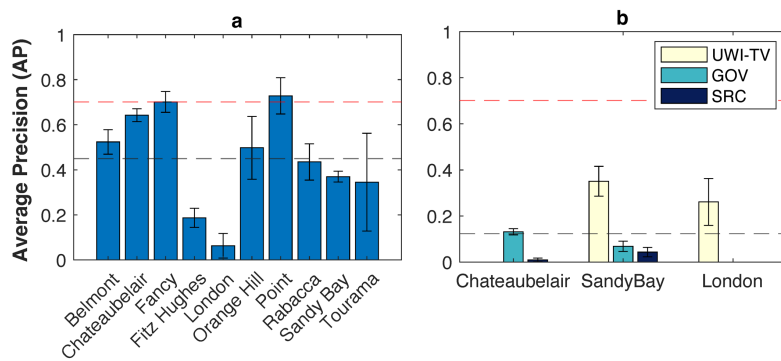
Deleted: (Figure 4b)

Deleted: 4

Deleted: higher

Formatted: Not Highlight





927  
 928 *Figure 5. Cross validation of the best experimental setup for building localisation models which*  
 929 *are trained to predict building box positions within the image. a) The effect of changing the*  
 930 *location used as the test set on detector average precision (AP) and b) the effect of changing the*  
 931 *imagery dataset (UWI-TV/GOV/SRC) used as the test set on AP. For b) cross validation of the*  
 932 *imagery dataset, models are trained on all data from that location excluding the location used for*  
 933 *testing as indicated by the bar. For London there is data from the GOV dataset, however the number*  
 934 *of images in the SRC dataset is insufficient for training, so no bar is shown for GOV. The AP shown*  
 935 *is the mean value from three trained models with the same setup while the error bars show the*  
 936 *standard deviation. Black dashed lines show the mean AP value across all cross validation trained*  
 937 *models; red dashed lines show the best AP from the experiments (0.701: Table 4).*

Deleted: 4.

Deleted:

Deleted: 3

938  
 939 **3.1.3 Evaluation on the test set**

940 Evaluation of the best detection model on the test set, which consists of completely unseen data  
 941 from Owia, Richmond and Troumaca (Figure 3) produced an AP value that is the same as the  
 942 value on the validation data (0.701) (Table 6). **To understand if a better model could be achieved**  
 943 **with more data available for training, we combined the training and validation data and used**  
 944 **this to retrain** the best experimental setup for the detector. **Evaluation of the retrained model**  
 945 **on the test set resulted in an average precision increase from 0.701 to 0.751 for the non-sieved**  
 946 **detector, and from 0.668 to 0.728 for the sieved detector, showing that having more data**  
 947 **available for training produced a better model (Table 6).**

Deleted: 2

Deleted: 4

Deleted: 3

Deleted: Retraining

Deleted:

Deleted: using the combined training and validation data caused the AP when applied to the test data to increase to 0.751 prior to sieving and 0.728 after sieving

Deleted: Comparing the precision and recall of the retrained detector and the retrained detector + sieve network shows that ...w

948  
 949 **While** the AP is higher for the retrained detector without the sieve, the addition of the sieve  
 950 network creates a better balance between the precision and recall **which is reflected in the**

Deleted: w

Deleted:

968 higher F1 score (Table 6). For the present application equal importance is given to: 1) making  
 969 correct predictions about building locations, and 2) identifying as many buildings as possible.  
 970 Consequently, striking the balance between precision and recall is crucial. We therefore selected  
 971 the retrained detector + sieve network as the final building localisation model and the model  
 972 that is incorporated into the damage assessment pipeline (Table 6).

Deleted: , reflected in a higher F1 score.

Deleted: 7

Deleted: 6

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

Formatted: Caption, Left, Keep with next

	Validation set				Test set			
	AP	P	R	F1	AP	P	R	F1
Detector (0.5 conf)	0.701	0.588	0.776	0.669	0.701	0.604	0.776	0.679
Detector + 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 + sieve					0.728	0.710	0.782	0.744

Formatted Table

978

979

## 980 3.2 Damage classification

Deleted: Building

### 981 3.2.1 Model selection

Deleted: d

982 The five experiments with the highest macro F1 score are shown in Table 7, with the full lists  
 983 provided in Tables S4 and S5 of the supplementary material. For Classifier 1, Macro F1 scores  
 984 across all 15 experiments ranged from 0.753 to 0.836, while for Classifier 2 scores ranged from  
 985 0.776 to 0.810 (Tables 7, S4, S5). Models trained to differentiate between the No damage to  
 986 minor damage and Damaged classes performed better for the No damage to minor damage  
 987 class, while those trained to differentiate between Moderate and Major damage performed  
 988 better for the Major damage class (Table 7). The best performing models for both classifiers  
 989 used the ResNet50 architecture rather than GoogleNet with an unbalanced dataset. For  
 990 Classifier 1 the best model had F1 = 0.962 for the No damage to minor damage class and F1 =  
 991 0.710 for the Damaged class. While for Classifier 2 the Moderate damage class had F1 = 0.770  
 992 and Major damage F1 = 0.851.

Deleted: top

Deleted: (

Deleted: )

Deleted: conducted for building damage classification

Deleted: 6

Deleted: 5

Deleted: 3

Deleted: 4

Deleted: and 0.776 to 0.810 for

Deleted: c

Deleted: 1 and 2 respectively

Deleted: 6

Deleted: 5

Deleted: 3

Deleted: 4

Deleted: Not Damaged

993

1015 Table 7. The top five experiments conducted for each of the building damage classifiers, ordered  
 1016 by the macro F1 score. The full list consisting of all 15 experiments for each classifier is provided  
 1017 in Tables S4 and S5 of the supplementary material.

Classifier 1					
Row ID	Architecture	Class balancing: Not Balanced/ under-sampled/ over-sampled	F1 No damage to minor damage	F1 Damaged	F1 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/ under-sampled/ over-sampled	F1 Mod damage	F1 Maj damage	F1 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

Deleted: 6

Deleted: 5

Deleted: 3

Deleted: 4

Formatted Table

Deleted: t

Deleted: D

Deleted: d

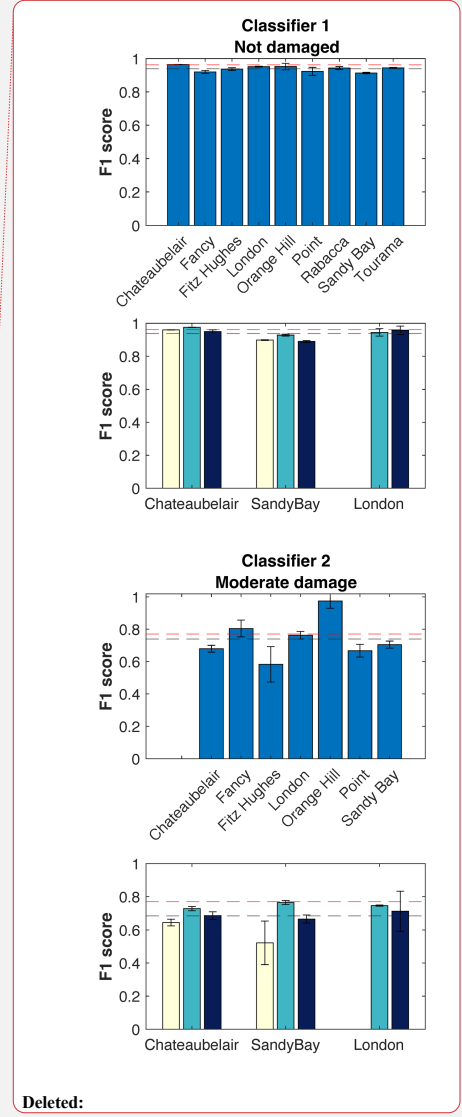
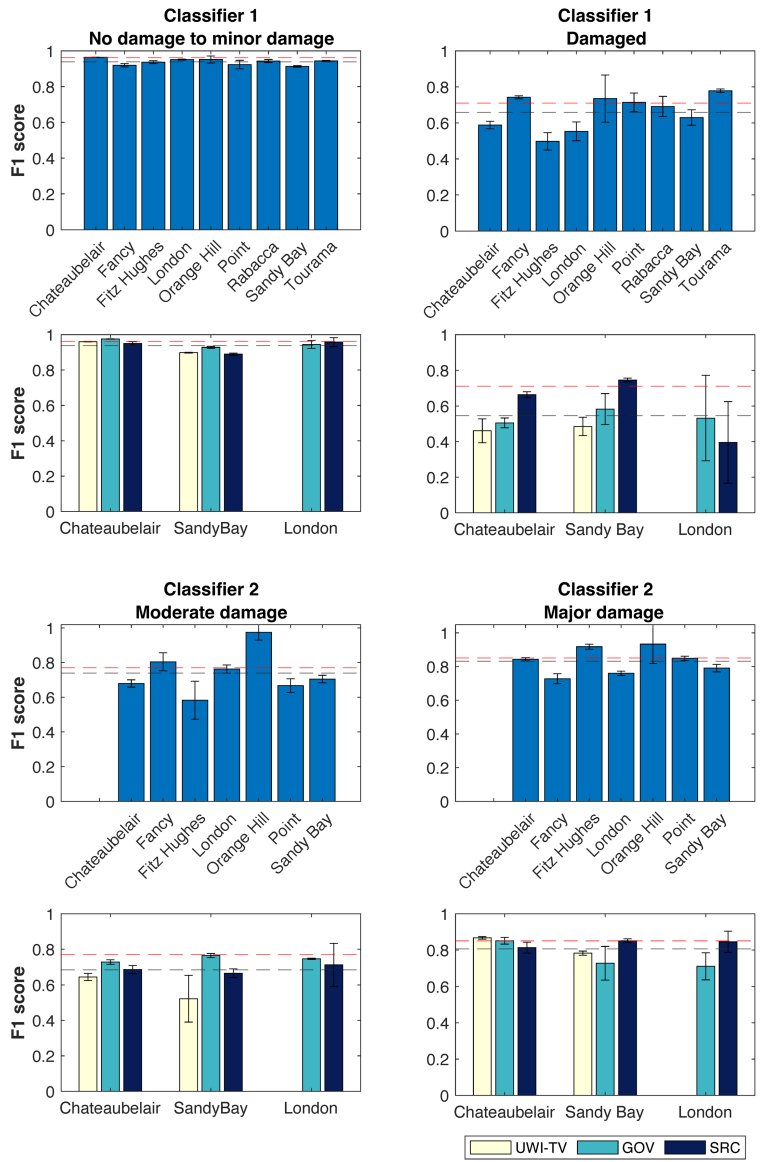
1019

1020

1021 **3.2.2 Cross validation**

1022 Cross validation was conducted for both of the single best performing models for Classifiers 1  
 1023 and 2 identified through model selection. As was the case for the best building localisation  
 1024 model, this was done to understand how the choice of training and validation datasets affected  
 1025 model performance and to understand how our model might perform on a new dataset.

Deleted: the potential for our model to generalize to



Deleted:

1034

1035 *Figure 6. Cross validation for Classifiers 1 and 2. For rows 1 and 3 the best experimental setup was*  
 1036 *retrained on all the data from locations in the combined training and validation data and*  
 1037 *evaluated on the location shown. For rows 2 and 4 the best experimental setup was retrained on*

Deleted: 5

1040 all the data from the location shown and evaluated on each dataset (UWI-TV/GOV/SRC)  
1041 separately. Each training was conducted three times, the value plotted is the mean, and the error  
1042 bars show the standard deviation. Black dashed lines show the mean F1 score across all cross  
1043 validation trained models, red dashed lines show the best F1 score for each class from the  
1044 experiments (Table 6).

Deleted: 5

1046 The performance of Classifier 1 for the No damage to minor damage class is consistent across  
1047 the distinct locations and datasets used for evaluation with mean F1 scores between 0.913-  
1048 0.983 for locations and 0.898-0.976 for datasets (Figure 6). For the Damaged class there is more  
1049 variety in the performance across the locations and datasets used for evaluation. The mean F1  
1050 scores for the separate locations range from 0.588 (Fitz Hughes) to 0.779 (Tourama) while for  
1051 the different datasets the range is 0.393 (London-SRC) to 0.745 (Sandy Bay-SRC).

Deleted: Figure 5 shows that the

Deleted: t

Deleted: d

Deleted: testing

Deleted: 5

Deleted: the choice of

1053 For Classifier 2, the Moderate damage class is more sensitive to the choice of location and  
1054 dataset used for the evaluation than the Major damage class (Figure 6). For the different  
1055 locations the mean F1 score ranged from 0.583-0.974. Similarly to Classifier 1, the location with  
1056 the lowest mean F1 score is Fitz Hughes, whereas the highest score was produced for Orange  
1057 Hill. For the different datasets the range for the Moderate damage class is between 0.522-0.746.  
1058 For the Major damage class F1 scores for the distinct locations are between 0.728-0.933 while  
1059 for the different datasets the range is between 0.711-0.867.

Deleted: validation

Deleted: 5

Moved (insertion) [1]

Deleted: .

Moved up [1]: For the Major damage class F1 scores for the distinct locations are between 0.728-0.933.

Deleted: For the Moderate damage class, the mean F1 score ranged from 0.583-0.974. Similarly to Classifier 1, Fitz Hughes produced the lowest mean F1 score, whereas the highest score was produced for Orange Hill. For the Major damage class F1 scores for the distinct locations are between 0.728-0.933. For Classifier 2 the sensitivity to the choice of dataset (UWI-TV/GOV/SRC) for the Moderate damage class is greater than for the Major damage class. For Moderate damage, the range is between 0.522-0.746, while for Major damage the range is from 0.711-0.867.

Deleted: c

Deleted: ication

Deleted: classification

Deleted: Table

Deleted: 7

Deleted: 6

Deleted: d

Deleted: The confusion matrices for both final models are plotted in Figure 6, these show class accuracy i.e., how many of the true class were correctly classified. For Classifier 1 89% of the Not damaged buildings were correctly classified, and 73% of the Damaged buildings were correctly classified. For Classifier 2 81% of the moderately damaged buildings were correctly classified, while 87% of the buildings with major damage were correctly classified. ...

### 1061 3.2.3 Evaluation on the test set

1062 Evaluation of the single best models for Classifier 1 and Classifier 2 on the unseen test set  
1063 produced Macro F1 scores that were comparable with the scores for the validation set: 0.829  
1064 for Classifier 1 and 0.791 for Classifier 2 (Table 8). For Classifier 2, retraining the model on the  
1065 combined training and testing data increased the Macro F1 score from 0.791 to 0.838. Whereas  
1066 for Classifier 1 retraining produced a slightly lower Macro F1 score (0.809 compared to 0.829).  
1067 Nevertheless, the retrained model for Classifier 1 achieved a higher recall on the Damaged class  
1068 than the non-retrained model. In an operational setting it's desirable to correctly classify as  
1069 many of the damaged buildings as possible, since in our pipeline these will be passed onto  
1070 Classifier 2, therefore we took the retrained models for both classifiers as the final models and  
1071 the models that are incorporated into the damage assessment pipeline.

1112 Table 8. Comparison of the best **damage classification models'** performance when evaluated on  
 1113 the validation and the test sets. AP is average precision, P is precision, and R is recall. \* Retrain  
 1114 models are trained on the combined training and validation sets. Results for the final models that  
 1115 are used in the damage assessment pipeline are in bold.

1116

	Validation set							Test set						
	No damage to minor damage			Damaged				No damage to minor damage			Damaged			
	P	R	F1	P	R	F1	F1 Macro	P	R	F1	P	R	F1	F1 Macro
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								<b>0.899</b>	<b>0.894</b>	<b>0.896</b>	<b>0.717</b>	<b>0.728</b>	<b>0.722</b>	<b>0.809</b>
	Mod Damage			Maj Damage				Mod Damage			Maj Damage			
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								<b>0.861</b>	<b>0.809</b>	<b>0.834</b>	<b>0.817</b>	<b>0.866</b>	<b>0.841</b>	<b>0.838</b>

1117

1118 **4 Application of the full damage assessment pipeline: Assessing tephra fall building**  
 1119 **damage in Owia**

1120

1121 In this work we have developed separate models for building localisation and two stages of  
 1122 damage classification. However, in an operational context models need to work sequentially,  
 1123 this led to the development of our damage assessment pipeline (outlined in Figure 4d). The  
 1124 pipeline operates on an orthomosaic image and outputs a georeferenced vector set, with the  
 1125 following attributes for each building that is detected: *detection* (box confidence score),  
 1126 *ClassPred\_1* (output class from Classifier 1, Damaged or No damage to minor damage),  
 1127 *ClassProb\_1* (the probability of that class), *ClassPred\_2* (output class from Classifier 2, Moderate  
 1128 damage or Major damage, this is only run if Classifier 1 outputs damage), *ClassProb\_2* (the  
 1129 probability of the class output by Classifier 2), *damageState* (the final damage state).

1130

1131 The tephra fall building damage map shown in Figure 7a was produced by overlaying the  
 1132 georeferenced vector that was output by the pipeline with the orthomosaic image in QGIS. Our  
 1133 remote damage assessment pipeline identified 442 buildings. Of these, 78% (N = 343) were  
 1134 classified as having No damage to minor damage, 9% (N = 40) as having Moderate damage and  
 1135 13% (N = 59) as having Major damage. We observed that the two upper tephra fall thickness

Deleted: 7  
 Deleted: 6  
 Deleted: '

Deleted: AP  
 Deleted: t  
 Deleted: d  
 Deleted: t  
 Deleted: d  
 Deleted: d

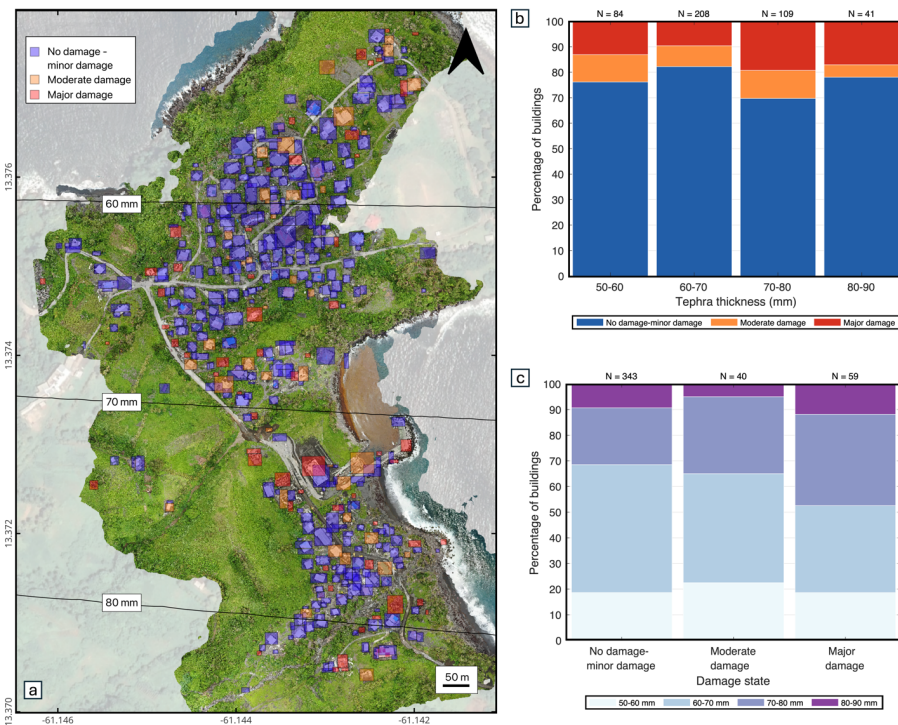
Deleted: AP  
 Deleted: Example  
 Deleted: a

Deleted: In order to  
 Deleted: 3  
 Deleted: , which can easily be generated using soft  
 Deleted: s  
 Deleted: shape file  
 Deleted: c  
 Deleted: d  
 Deleted: n  
 Deleted: t  
 Deleted: d  
 Deleted: c  
 Deleted: m  
 Deleted: m  
 Deleted: c  
 Deleted: c  
 Deleted: .  
 Deleted: Figure 7 shows an example  
 Deleted: running the pipeline on an orthomosaic i  
 Deleted: pipeline orthomosaic and output  
 Deleted: shapefile  
 Deleted: The example which consists of 417 buildi



1197 bins (70-80 mm and 80-90 mm), both had a higher proportion of buildings with Major damage  
 1198 compared to the lower thickness bins (Figure 7b, c), indicating a correlation between tephra fall  
 1199 thickness and building damage though it is not very pronounced. These findings are discussed  
 1200 in Section 5.3.

1202 The full pipeline took 1 hour to run on a standard 16GB RAM 2021 MacBook Pro, with an M1  
 1203 Pro chip. Most of the inference time was attributed to the building localisation module in the  
 1204 pipeline, which may be bypassed if building footprints are already available. When only the  
 1205 classifiers were run the time taken to run was reduced to < 5 mins.



1207 Figure 7. Application of our remote tephra fall building damage assessment pipeline to Owia,  
 1208 located in the north of St. Vincent. a) The damage map produced by overlaying the spatial data  
 1209 generated by our pipeline onto the orthomosaic image, black lines are tephra isopachs,  
 1210 interpolated from Cole et al., 2023; b) the proportion of damage states with increasing tephra

Deleted: 7  
 Deleted: the full  
 Deleted: on the orthomosaic image for  
 Deleted: .  
 Deleted: ,  
 Deleted: t



1217 *thickness: c) the proportion of tephra thickness bins with increasing damage state. Coordinate*  
1218 *reference system: WGS 84 (EPSG:4326). Satellite basemap © Google Maps 2024.*  
1219

## 1220 5 Discussion

1221

1222 In this work we have developed models for building localisation, and two levels of damage  
1223 classification for building damage resulting from tephra fall. Our final models demonstrate  
1224 strong performance for both building localisation ( $AP = 0.728$ ;  $F1 = 0.744$ ) and damage  
1225 classification (Classifier 1,  $F1 = 0.809$ , Classifier 2,  $F1 = 0.838$ ). Despite using post-event imagery  
1226 only, which makes the task more challenging than approaches using multi-temporal imagery,  
1227 our results are comparable to existing optical imagery building damage assessments developed  
1228 for various hazards that use both mono-temporal and multi-temporal images (F1 scores are  
1229 between 0.656-0.868 for building localisation and 0.650-0.981 for damage classification. Table  
1230 1).

### 1232 5.1 Building localisation

1233

1234 Through running our building localisation experiments we found that the pre-processing of  
1235 images before detector training (particularly the block size) significantly influenced detector  
1236 performance. The block sizes tested were chosen as a trade-off between reducing image size  
1237 sufficiently to reduce computational cost, and retaining a large enough size such that buildings  
1238 were not dissected unnecessarily. Given that the optimum block size was the middle size of the  
1239 range tested, we are confident that this balance was achieved. Cross-validation results  
1240 demonstrated variability in average precision (AP) for models trained on different locations and  
1241 imagery datasets (UWI-TV/GOV/SRC) (Section 3.1.2; Figure 5). Deep learning models are  
1242 known to perform well when the data they are evaluated on have similar characteristics to the  
1243 data they were trained on, though have more difficulty when working with 'out of distribution'  
1244 samples (Ben-David et al., 2010). Given the relatively consistent building typology across  
1245 locations (most buildings observed are detached single storey buildings with either a gable or  
1246 hip shaped metal sheet roof; a lesser proportion have flat concrete roofs), the differences in AP  
1247 are likely due to observable variations in UAV altitude, off-nadir angles, tephra thicknesses, and  
1248 varying training sample sizes.  
1249

Deleted: AP =

Deleted: building

Deleted: that use

Deleted: both pre- and post-event

Deleted: for building localisation

Deleted: :

Deleted: ,

Deleted: damage classification:

Deleted: 3

Deleted: 4

Deleted: data come from the same distribution

Deleted: ,

Deleted: the majority of

1263 The cross-validation AP was notably lower for the London and Fitz Hughes datasets (Section  
1264 3.1.2). For the London images (from SRC and GOV datasets) this is likely caused by the smaller  
1265 apparent size of buildings in these images compared to the other locations, due to the higher  
1266 UAV altitude. Variations in object size within the training and testing data has been found to  
1267 affect the performance of deep learning models developed for building localisation, with models  
1268 often performing better for objects that are the same size as those in the training data (Nath  
1269 and Benzadan, 2020; Cheng et al., 2021; Bouchard et al., 2022). Fitz Hughes images were all  
1270 from the UWI-TV image dataset which contributed just 17% to the combined training and  
1271 validation set used for cross validation. This dataset was collected closer in time to the eruption,  
1272 therefore as a whole had more tephra on the ground than the SRC and GOV datasets, which  
1273 affects background colour. Furthermore, the UWI-TV dataset viewed buildings mostly from an  
1274 off-nadir perspective, while the other datasets were predominantly nadir images. The effect of  
1275 image background colour on localisation performance is expected to be minor. Cheng et al.  
1276 (2021) found that for the same event localisation AP dropped from 65.6 to 63.3 when their  
1277 model was tested on images containing buildings surrounded by vegetation compared to  
1278 buildings with an ocean backdrop. While Bouchard et al. (2022) suggested that models quickly  
1279 learn to ignore background pixels. On the other hand, variation in off-nadir angles is a widely  
1280 acknowledged challenge of working with UAV or aerial images (Cotrufo et al., 2018; Nex et al.,  
1281 2019; Pi et al., 2020). Under representation of the mostly off-nadir UWI-TV images in the  
1282 training data may have impacted the model's ability to recognise such instances in the test data.  
1283 During model development we experimented with different models for the different datasets  
1284 (UWI-TV, GOV, SRC), but found that models developed on the combined dataset performed  
1285 better than those developed on the separate datasets, and a combined model was the one  
1286 selected and used for cross validation. Rather than suggesting that variations in off-nadir angle  
1287 are not important, this finding likely reflects the smaller size of the individual datasets  
1288 compared to the combined datasets, meaning that less information was available to learn from.  
1289 The application of sampling approaches like those used for the damage states in the  
1290 classification model development (over or under sampling) could have been applied to balance  
1291 the data. However, the SRC dataset is much larger than either of the UWI-TV and GOV sets  
1292 (Figure 3), therefore we considered that oversampling would introduce significant bias towards  
1293 the specific examples in the under-represented dataset, whereas through under sampling we  
1294 would lose a large amount of the data that are available to learn from. Given these factors, we  
1295 did not use sampling approaches. Future work might consider the application of generative AI

Deleted: The

Deleted: and

Deleted: (UWI-TV) exhibited the lowest average precisions (Figure 4). Both London datasets featured smaller buildings than the rest of the locations, evident in Section S3 of the supplementary material, while the

Deleted: s

Deleted: a

Deleted: nd,

Deleted: The training data was, predominantly nadir images from the SRC dataset with fewer UWI-TV examples, collected further in time after the eruption had fewer UWI-TV examples which are off-nadir and, collected more closely in time after the eruption, meaning more less tephra was present in images and, had.

Deleted: Images captured with different

Deleted: differences

Deleted: s

Deleted: .

Deleted: This

Deleted: U

Deleted: u

Deleted: z

Deleted: training and testing

Deleted: images

Deleted: was

Deleted: better than

Deleted: doing this separately

Deleted: T

Deleted: is likely related

Deleted: to

Deleted: , rather than the difference in the off-nadir angles...

Deleted: ,

Deleted: h

Deleted: 2

Deleted: we did not use this approach as

1333 algorithms such as generative adversarial networks (GANs) to expand the dataset (e.g., Yi et al.  
1334 2018; Yorioka et al., 2020), although more work needs to be done to quantify the diversity in  
1335 the generated data.

1336  
1337 The variability in cross-validation results for the building localisation model likely comes from  
1338 a combination of the above factors (differences in UAV altitude, off-nadir angles, tephra  
1339 thickness, and varying training sample sizes), and suggests that there was insufficient  
1340 information in the training data for our detection models to perform well across the range of  
1341 characteristics present. This is supported by the increased performance when the best  
1342 localisation model was retrained on the combined training and validation data. However,  
1343 further investigation is required to separate the unique effect of each aspect.

Deleted: (Section 3.1.3)

Deleted: this requires

## 1345 5.2 Damage classification

1346  
1347 The final classification models achieved better performance than the final localisation model  
1348 with macro F1 scores of 0.809 and 0.838 on the test data (Table 8). Cross-validation showed  
1349 that classification models were less sensitive than the localisation model to the choice of  
1350 datasets used for training and evaluation (Section 3.2.2). We found that class wise our models  
1351 performed better on the No damage to minor damage class followed by the Major damage class.  
1352 This agrees with other multi-class studies that have found the extremities of the damage state  
1353 scheme applied easier to classify than the intermediate ones (Kerle et al., 2019, Valentijn et al.  
1354 2020).

Deleted: Building

Deleted: d

Deleted: 7

Deleted: 6

Deleted: n

Deleted: t

Deleted: d

Deleted: m

Deleted: is in agreement

Deleted: used

Deleted: to be

Formatted: Font: 12 pt, Bold, English (UK)

Formatted: Line spacing: 1.5 lines

Deleted: <#>Application of the pipeline

## 1356 5.3 Application of the full damage assessment pipeline: Assessing tephra fall building 1357 damage in Owia

1358  
1359 Application of our remote damage assessment pipeline to the town of Owia found that 22% of  
1360 buildings that received tephra accumulation in the range of 50-90 mm experienced Moderate  
1361 damage or Major damage. Within this range, the relationship between tephra thickness and  
1362 building damage was not as pronounced as in other studies (Blong, 2003b; Hayes et al., 2019;  
1363 Jenkins et al., 2024). This may be attributed to the small geographic area and therefore small  
1364 range of tephra thicknesses considered in our application when compared to other studies. In  
1365 the damage assessments of Blong, (2003b), Hayes et al., (2019) and Jenkins et al., (2024)  
1366 buildings received ~100 to 950 mm, trace to 600 mm and, trace to >220 mm respectively.

1381 Spence et al., (1996) assessed building damage over a similarly narrow range of tephra  
1382 thicknesses to this work (~150-200 mm) and found that there was considerable variation in  
1383 the level of damage despite the majority of buildings having a metal sheet roof. The spacing  
1384 between the principal roof supports (roof span) was found to be important for the amount of  
1385 damage observed, with long span buildings experiencing higher levels of damage than short  
1386 span ones (Spence et al., 1996). There are limited long span buildings in the Owia case study,  
1387 however additional characteristics such as construction style and material, building layout, age,  
1388 condition, height, and roof pitch can all affect a buildings ability to withstand tephra loading  
1389 (Spence et al., 1996; Pomonis et al., 1999; Blong, 2003b; Jenkins et al., 2014). Variation in these  
1390 characteristics across Owia could be responsible for the observed variation in building damage  
1391 over the narrow range of thicknesses considered.

1392  
1393 If we convert tephra thickness to loading, we can compare the results of our assessment with  
1394 existing relationships between tephra loading and damage for similar building types. Using a  
1395 density of 1500 kg/m<sup>2</sup> (Cole et al., 2023) suggests that a loading of at least 75-135 kg/m<sup>2</sup> was  
1396 applied to buildings for the range of thicknesses considered (50 mm-90 mm). Census data for  
1397 Owia states that 90 % of buildings have metal sheet roofs (SVG population and housing census,  
1398 2012), with the remaining 8% comprised of reinforced concrete roofs and 2% 'other material',  
1399 Given the higher resistance of the 8% of non-metal sheet roof buildings in Owia, we might  
1400 expect vulnerability models developed for metal sheet roofs to overestimate damage in the  
1401 town. Fragility functions developed for Indonesian style buildings with metal sheet roofs  
1402 (Williams et al., 2020) calculate a 48-80% probability of Owia buildings experiencing damage  
1403 exceeding Damage State 2, higher than the 22% experiencing Moderate or Major damage in our  
1404 study. Fragility curves for roof failure (Major damage) of old or poor condition metal sheet roofs  
1405 (Jenkins et al., 2014), calculate that just over 10% of buildings in Owia would experience  
1406 sufficient loading for roof collapse, comparable to the 13% observed in our study. These  
1407 comparisons highlight some of the challenges associated with using vulnerability models  
1408 developed for different locations. Moreover, they reiterate the need for the collection of both  
1409 post-event impact data and building typology information that can be used to increase the  
1410 amount of empirical data available for vulnerability model development and allow regional  
1411 vulnerability models to be developed for specific building types.

- Deleted: other
- Deleted: 10
- Deleted: ?
- Deleted: 10
- Deleted: fragility curves
- Deleted: village
- Deleted: (to change depending on what the other 10% is)

- Deleted:
- Deleted:
- Deleted: locations, and
- Deleted: Furthermore
- Deleted: It
- Deleted: building typology and

1426 Like the studies presented in Table 1, our pipeline consists of separate models for localisation  
1427 and damage classification. One of the benefits of this is that in locations where precise building  
1428 location information is available for the assessment area, the localisation step can be bypassed  
1429 and only the classifiers run. This not only enhances overall performance but also significantly  
1430 reduces computation time. Furthermore, either of the classifiers can be run independently  
1431 and/or combined with other damage assessment procedures; for example, an initial synthetic  
1432 aperture radar (SAR) based assessment (e.g., Yun et al. 2015, Jung et al., 2016), could be  
1433 followed with our Classifier 2 to provide additional granularity on the severity of the damage at  
1434 a building level rather than a pixel level.

Deleted: 0

Deleted: building

Deleted: c

#### 1436 5.4 Generalisability to other locations

1437  
1438 Our models have performed well for images collected on the island of St Vincent where building  
1439 typologies are relatively consistent. We therefore expect that our models will perform well in  
1440 other locations with similar building types, such as the other islands in the Lesser Antilles. This  
1441 hypothesis should be validated through further testing. In absence of additional UAV datasets  
1442 that include damaged buildings, testing can be done by conducting pre-event surveys to test the  
1443 performance of the building localisation model and Classifier 1 for the No damage to minor  
1444 damage class. While this is unable to assess the ability of our approach to classify damage, it  
1445 would provide some indication of performance following an event in a new location.

Deleted: c

Deleted: t

Deleted: d

1446  
1447 To develop a model that is robust to the diverse building types found across the world  
1448 necessitates assembling diverse datasets showcasing potential variations in building types and  
1449 the associated tephra fall damage. To our knowledge the UAV datasets described in this work  
1450 are the first of their kind. However, the increasing utilisation of UAVs during and after volcanic  
1451 events suggests the possibility of the emergence of more datasets in the years to come. Our  
1452 model represents a crucial initial step towards the operational implementation of this approach  
1453 globally. The compilation of global tephra fall building damage UAV datasets will facilitate the  
1454 ongoing refinement of building damage assessment approaches, including the one presented  
1455 here. In pursuit of this objective, our models stand ready for retraining as more data becomes  
1456 available. While our approach leverages images captured under a spectrum of flight conditions  
1457 (off-nadir angle, altitude, flight trajectory), our investigation has both pinpointed specific

Deleted: However

1465 conditions that are best suited for capturing building damage, which are detailed in Section 6,  
1466 and highlighted the importance of consistency in data collection.  
1467

### 1468 5.5 Improving model performance and future perspectives

1469

1470 The advantages of acquiring additional UAV datasets both before and after an event have been  
1471 outlined in Section 5.4. In addition to this, pre-event imagery can be used to construct building  
1472 inventories manually or using machine learning methods (e.g., Iannelli and Dell'Acqua, 2017;  
1473 Gonzalez et al., 2020; Meng et al., 2023). Prior to an eruption, information about how the  
1474 building typologies present will respond under certain tephra loadings (i.e., the forecasted  
1475 damage state) can be obtained through the application of fragility functions. This information  
1476 could enhance our model by serving as prior information, that is updated with outputs from our  
1477 remote damage assessment using Bayesian statistics. A similar approach has been suggested  
1478 for updating the United States Geological Survey's (USGS) Prompt Assessment of Global  
1479 Earthquakes for Response (PAGER) system (Noh et al., 2020). The framework provides a  
1480 structured way of incorporating the PAGER forecasted loss with the potentially noisy and  
1481 incomplete observations of loss in the early stages of response.  
1482

1483 Alternatively, with ample individual building inventory data available, tailored damage  
1484 classification models for specific building typologies could be developed and applied. The  
1485 rationale is that a model dedicated to a specific building type is expected to outperform a  
1486 generic multi-typology model.  
1487

1488 In this work, we established a three-class damage state framework. Existing frameworks that  
1489 were developed for ground based tephra fall damage assessment split damage into five damage  
1490 states classes and one non-damage class (Spence et al, 1996; Blong, 2003; Hayes et al., 2019;  
1491 Jenkins et al., 2024, Table 2) however in our preliminary analyses we found that: 1) in many  
1492 images we were unable to confidently apply a six-class scheme due to only being able to see one  
1493 side of the building, and 2) there were not enough examples of each damage state class to be  
1494 able to train a six-class model. With the addition of future tephra fall building damage datasets  
1495 it may be possible to apply a finer resolution damage state framework that can provide more  
1496 detail on the observable damage. However, it is unlikely that the resolution of ground-surveys  
1497 can be achieved using optical imagery, since lower damage states are still difficult to resolve

Deleted: 3

Deleted: surveys

Moved up [2]: or using machine learning methods such as the work of Meng et al., (2023).

Moved (insertion) [2]

Deleted: interrogated manually or using machine learning methods to construct building inventories that contain information such as construction materials and styles (e.g., Iannelli and Dell'Acqua, 2017; Gonzalez et al., 2020; Meng et al., 2023).

Deleted: may be particularly beneficial in constructing building inventories, which or using machine learning methods such as the work of Meng et al., (2023). include details about building typologies such as construction materials and styles. Surveys can be interrogated manually to extract building attributes

Deleted: given knowledge

Deleted: , an idea about how the buildings

Deleted: . The forecasted damage state could be subsequently refined through Bayesian updating

Deleted: based on our damage assessment models output

Deleted: based on our damage assessment models output

Deleted:

Deleted: three class

Deleted: in review

Deleted: Nevertheless, the damage states developed in our work can be equated to existing damage states generated for ground surveys such that: No damage – to minor damage = DS0-DS1, Moderate damage = DS2 and Major damage = DS3-5.

Deleted: may be applied

Deleted: is capable of providing

Deleted: confidently

1530 ~~even with very high-resolution images (Cotrufo et al., 2018). Some studies have incorporated~~  
1531 ~~3D point-cloud information into analyses (Cusicanqui et al., 2018; Vetrivel et al., 2018). While~~  
1532 ~~these approaches have shown potential, and could potentially~~ be used to provide additional  
1533 granularity to our damage states, we opted against integrating point cloud analyses into our  
1534 model, ~~due to the considerably longer processing times associated with such an approach.~~  
1535 ~~Longer processing times~~, would undermine the swift processing requirement inherent in our  
1536 methodology.

Deleted: We developed our approach using deep learning on 2D optical imagery, while s

Deleted: used

Deleted: , or combined point cloud information with deep learning on optical imagery for damage level classification (Vetrivel et al., 2018).

Deleted: the use of 3D spatial data has shown potential, and potential

Deleted: may

Deleted: . This decision was motivated by

Deleted: ,

Deleted: which

Deleted: ¶

## 1537 1538 5.6 Caveats

1539 During the assignment of building damage states, uncertainties arose, particularly concerning  
1540 the interpretation of tarpaulins and, pre-existing damage. For tarpaulins, the ambiguity arose  
1541 from whether these were either strategically placed prior to the eruption as preventative  
1542 measures to cause tephra to slide off the roof more easily; or they were placed post event to  
1543 cover damage caused by tephra fall. Additionally, in certain instances, distinguishing between a  
1544 collapsed roof and a section of the building initially lacking roofing material—possibly  
1545 functioning as a walled storage area—proved challenging. Pre-existing damage not related to  
1546 volcanic activity or buildings that were under construction at the time of image acquisition were  
1547 considered as damaged and classified accordingly. ~~The presence of buildings under~~  
1548 ~~construction at the time of image acquisition has been recognised as a challenge in studies using~~  
1549 ~~mono-temporal imagery (Nex et al., 2019; Cheng et al., 2021).~~ Pre-event imagery would have  
1550 provided clarity on ~~both of~~ these matters, however this was not available at high enough  
1551 resolution for this region.

Deleted: by several authors

1552  
1553  
1554 The majority of images used for training and evaluating our models came from the SRC dataset,  
1555 which was collected several months after the eruption. As a ~~result~~, the majority of images do  
1556 not have much tephra present. In an operational context, to expedite the recovery process, data  
1557 would ideally be collected as quickly after the eruption as it is safe to do so, therefore more  
1558 tephra would be present in the images. Given the compound effects of variations in flight angle,  
1559 image lighting, resolution and also the presence of tephra, we do not have enough information  
1560 to test the effect of tephra thickness on model performance, and caution should be taken when  
1561 using the model on data collected at different times after the eruption.

Deleted: result



1579 **6 Recommendations for UAV building damage assessment data collection**

1580  
1581 In the future we advocate for the adoption of a standardised protocol for data collection for the  
1582 purpose of UAV damage assessment. While our model was developed using a diverse dataset,  
1583 there were some disparities in performance across distinct data types. Consequently, the  
1584 standardisation of image collection serves two purposes, 1) to allow the best results to be  
1585 achieved when implementing our models, and 2) to collect data that is rich in information useful  
1586 for damage assessment with the aim of working towards the development of global datasets for  
1587 tephra fall damage. For best results we have the following recommendations:

- 1588
- 1589 • The bulk of our dataset was collected several months after the eruption of La Soufrière  
1590 however, for generating a global dataset that can be used for response and recovery,  
1591 models should ideally be trained on images collected shortly (days to weeks) after an  
1592 event.
  - 1593 • Flight paths should be pre-programmed to ensure comprehensive coverage of the area  
1594 and limit bias associated with overrepresentation of certain buildings. Ideally two flights  
1595 would be conducted with two sets of perpendicular flight lines to capture buildings from  
1596 a different perspective. GPS positioning should be enabled.
  - 1597 • A fixed altitude of 50-80 m above the ground should be maintained where possible. This  
1598 is appropriate to capture sufficient data for accurate damage classification based on the  
1599 established framework and strikes a balance between detailed information capture and  
1600 overall coverage. In mountainous areas this may not be achievable for some UAV types.  
1601 In which case a uniform height should be maintained such that the size of buildings is  
1602 consistent across image frames.
  - 1603 • We suggest a slightly off-nadir camera positioning (~5-15°), which is sufficient to  
1604 capture any bending in the roof that may not be captured from a nadir perspective.
  - 1605 • Overlap between images should be enough to generate orthoimages, 80% forward and  
1606 70% lateral overlap is sufficient.

1607  
1608 In addition to the development of optimum post-event data collection practises we advocate  
1609 for the collection of pre-event UAV datasets. Ideally, pre- and post-event imagery is collected  
1610 using the same flight paths, altitudes, and camera positioning. Pre-event datasets serve  
1611 multiple purposes:

Deleted: size



- 1613 ○ Facilitates the creation of building inventories.
- 1614 ○ Enables precise comparison of pre- and post-event imagery, reducing uncertainty
- 1615 regarding initial building conditions.
- 1616 ○ Supports the development of high-resolution change detection models
- 1617 potentially yielding more accurate results than relying solely on post-event
- 1618 imagery.
- 1619 ○ Provides an opportunity for UAV pilots to gain experience in capturing building
- 1620 datasets during 'quiet times'.

## 1621 7 Conclusions

1622 Following a large tephra fall event, building damage assessment needs to be conducted rapidly  
1623 for the purpose of response and recovery, and for the collection of data that can be used to  
1624 forecast building damage from future events. By leveraging post-event optical imagery obtained  
1625 after the 2021 eruption of La Soufrière volcano on the island of St Vincent, and convolutional  
1626 neural networks, we have developed an automated tephra fall building damage assessment  
1627 pipeline. The pipeline incorporates models for building localisation and two distinct levels of  
1628 damage classification: distinguishing between No damage to minor damage and damage, as well  
1629 as between Moderate and Major damage, which were trained and evaluated separately. When  
1630 provided with UAV optical imagery, our pipeline can rapidly generate spatial building damage  
1631 information. Our models perform well for the St Vincent datasets and are anticipated to perform  
1632 well in locations where building typologies are similar, but this requires more testing to  
1633 understand the limits of their application.

1634  
1635 Building localisation model cross validation results underscore the influence of factors such as  
1636 UAV altitude, off-nadir angles, tephra thickness, and training sample sizes on model  
1637 performance, while results show that ~~damage classification models were affected by these~~  
1638 factors to a lesser extent. We acknowledge the challenges posed by diverse datasets and by  
1639 limited data, and we propose a series of recommendations to guide the collection of future UAV  
1640 building damage datasets. In addition to the collection of post-event datasets we advocate for  
1641 the collection and incorporation of pre-event datasets, which can be used to support the  
1642 advancement of change detection models; to partially evaluate the models presented here  
1643 during quiescent times, and to develop building inventories that can be used along with fragility  
1644 functions for forecasting building damage.

Deleted: building

1647

1648 Our research marks a step forward in tephra fall building damage assessment, offering a  
1649 versatile and effective pipeline with the potential for regional applicability. As the field of UAV-  
1650 based damage assessment in volcanology continues to evolve, our work lays a foundation for  
1651 further advancements, contributing to the resilience of communities in the face of volcanic  
1652 eruptions.

1653

## 1654 **8 Author contributions**

1655

1656 Conceptualization: SFJ, RR, ET, VM. Data collection: RR and VM. Development of the  
1657 methodology: ET, SFJ, BW. Software: ET. Formal analysis: ET. Supervision: SFJ. Writing – original  
1658 draft: ET. Writing-Reviewing & Editing: ET, SFJ, VM, RR, BW, BT, SHY.

## 1659 **9 Competing interests**

1660

1661 The authors declare no competing interests.

## 1662 **10 Acknowledgements**

1663

1664 We are indebted to Monique Johnson: The UWI Seismic Research Centre, Javid Collins: UWITV,  
1665 Nikolai Lewis and Marla Mulraine: The Government of St Vincent and the Grenadines Ministry  
1666 of Transport, Works, Lands and Surveys, and Physical Planning, for sharing their UAV data and  
1667 collaborating on this work. All images and data provided in this study have been approved for  
1668 publication by the local agency responsible for monitoring geohazards in St Vincent: The UWI  
1669 Seismic Research Centre. We are very grateful to Chee Jain Hao Denny, Sim Yu Yang, Isaiah Loh  
1670 Kai En, Huang Wanxin for their assistance with data preparation, and to Vanesa Burgos, Elinor  
1671 Meredith, ~~Alberto Ardid, and Tom Wilson,~~ for interesting discussions around machine learning  
1672 and building damage assessment. We would like to thank Sébastien Biass and one anonymous  
1673 reviewer for their detailed and constructive reviews that considerably improved the  
1674 manuscript and Giovanni Macedonio for their editorial handling.

1675

## 1676 **11 Data availability**

1677

Deleted: and

1679 All trained models along with the code required to execute the damage assessment pipeline  
1680 and instructions for usage are provided at:  
1681 <https://github.com/EllyTennant/UAVdamageAssessment>  
1682

## 1683 12 Funding

1684 This research was supported by the Earth Observatory of Singapore via its funding from the  
1685 National Research Foundation Singapore and the Singapore Ministry of Education under the  
1686 Research Centres of Excellence initiative and comprises EOS contribution number 596.  
1687 Additional support was provided by the AXA Research Fund as part of the Joint Research  
1688 Initiative on Volcanic Risk in Asia.  
1689  
1690

## 13 References

- [An, G., Akiba, M., Omodaka, K., Nakazawa, T. & Yokota, H. \(2021\). Hierarchical deep learning models using transfer learning for disease detection and classification based on small number of medical images. \*Scientific Reports\*, 11\(1\). <https://doi.org/10.1038/s41598-021-83503-7>](#)
- Andaru, R. and Rau, J.Y. 2019. Lava dome changes detection at agung mountain during high level of volcanic activity using uav photogrammetry. In: *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences - ISPRS Archives*. International Society for Photogrammetry and Remote Sensing, pp. 173–179. doi: 10.5194/isprs-archives-XLII-2-W13-173-2019.
- Anniballe, R., Noto, F., Scalia, T., Bignami, C., Stramondo, S., Chini, M. and Pierdicca, N. 2018. Earthquake damage mapping: An overall assessment of ground surveys and VHR image change detection after L'Aquila 2009 earthquake. *Remote Sensing of Environment* 210, pp. 166–178. doi: 10.1016/j.rse.2018.03.004.
- Aggarwal, C. C. (2018). Neural Networks and Deep Learning. In *Neural Networks and Deep Learning*. <https://doi.org/10.1007/978-3-319-94463-0>
- [Ben-David, S., Blitzer, J., Crammer, K., Kulesza, A., Pereira, F., & Vaughan, J. W. \(2010\). A theory of learning from different domains. \*Machine Learning\*, 79\(1–2\), 151–175. <https://doi.org/10.1007/s10994-009-5152-4>](#)
- [Biass, S., Bonadonna, C., & Houghton, B. F. 2019. A step-by-step evaluation of empirical methods to quantify eruption source parameters from tephra-fall deposits. \*Journal of Applied Volcanology\*, 8\(1\). <https://doi.org/10.1186/s13617-018-0081-1>](#)
- Biass, S., Jenkins, S., Lallemand, D., Lim, T.N., Williams, G. and Yun, S.H., 2021. Remote sensing of volcanic impacts. In *Forecasting and Planning for Volcanic Hazards, Risks, and Disasters* (pp. 473-491). Elsevier.
- [Biass, S., Reyes-Hardy, M. P., Gregg, C., di Maio, L. S., Dominguez, L., Frischknecht, C., Bonadonna, C., & Perez, N. 2024. The spatiotemporal evolution of compound impacts from lava flow and tephra fallout on buildings: lessons from the 2021 Tajogaite eruption \(La Palma, Spain\). \*Bulletin of Volcanology\*, 86\(2\). <https://doi.org/10.1007/s00445-023-01700-w>](#)
- Blong, R. 2003a. *A Review of Damage Intensity Scales*. Available at: <http://www.es.mq.edu.au/NHRC/web/scales/scalesindex.htm>.
- Blong, R. 2003b. Building damage in Rabaul, Papua New Guinea, 1994. *Bulletin of Volcanology* 65(1), pp. 43–54. doi: 10.1007/s00445-002-0238-x.

Formatted: Font: Cambria, 11 pt

- Bouchard, I., Rancourt, M.È., Aloise, D. and Kalaitzis, F. 2022. On Transfer Learning for Building Damage Assessment from Satellite Imagery in Emergency Contexts. *Remote Sensing* 14(11), pp. 1–29. doi: 10.3390/rs14112532.
- Bruzzone, L. and Fernández Prieto, D. 2000. Automatic Analysis of the Difference Image for Unsupervised Change Detection. *IEEE TRANSACTIONS ON GEOSCIENCE AND REMOTE SENSING* 38(3), pp. 1171–1181.
- Cheng, C.S., Behzadan, A.H. and Noshadravan, A. 2021. Deep learning for post-hurricane aerial damage assessment of buildings. *Computer-Aided Civil and Infrastructure Engineering* 36(6), pp. 695–710. doi: 10.1111/mice.12658.
- Cole, P.D. et al. 2023. Explosive sequence of La Soufrière, St Vincent, April 2021: insights into drivers and consequences via eruptive products. Available at: <https://doi.org/10.6084/m9.figshare.c.6474317>.
- [Cotrufo, S., Sandu, C., Giulio Tonolo, F. & Boccardo, P. \(2018\). Building damage assessment scale tailored to remote sensing vertical imagery. \*European Journal of Remote Sensing\*, 51\(1\), 991–1005. <https://doi.org/10.1080/22797254.2018.1527662>](#)
- Cusicanqui, J., Kerle, N., & Nex, F. 2018. Usability of aerial video footage for 3-D scene reconstruction and structural damage assessment. *Natural Hazards and Earth System Sciences*, 18(6), 1583–1598. <https://doi.org/10.5194/nhess-18-1583-2018>
- 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., 2022. From anecdotes to quantification: advances in characterizing volcanic eruption impacts on the built environment. *Bulletin of Volcanology*, 84(1), p.7.
- Deng, J. et al., 2009. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255.
- Duarte, D., Nex, F., Kerle, N. and Vosselman, G. 2020. Satellite Image Classification of Building Damages Using Airborne and Satellite Image Samples in a Deep Learning Approach. In: *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*. Riva del Garda, Italy, pp. 4–7. Available at: <https://research.utwente.nl/en/publications/satellite-image-classification-of-building-damages-using-airborne>.
- Dung Cao, Q. and Choe, Y. 2020. *Building Damage Annotation on Post-Hurricane Satellite Imagery Based on Convolutional Neural Networks*.
- Gailler, L., Labazuy, P., Régis, E., Bontemps, M., Souriot, T., Bacques, G. and Carton, B. 2021. Validation of a new UAV magnetic prospecting tool for volcano monitoring and geohazard assessment. *Remote Sensing* 13(5), pp. 1–10. doi: 10.3390/rs13050894.
- Galanis, M., Rao, K., Yao, X., Tsai, Y.L., Ventura, J. and Fricker, G.A. 2021. DamageMap: A post-wildfire damaged buildings classifier. *International Journal of Disaster Risk Reduction* 65. doi: 10.1016/j.ijdrr.2021.102540.
- Ghosh, S. et al. 2011. Crowdsourcing for rapid damage assessment: The global earth observation catastrophe assessment network (GEO-CAN). *Earthquake Spectra* 27(SUPPL. 1). doi: 10.1193/1.3636416.
- Girshick, R. (2015). Fast R-CNN. <http://arxiv.org/abs/1504.08083>
- [Gonzalez, D., Rueda-Plata, D., Acevedo, A. B., Duque, J. C., Ramos-Pollán, R., Betancourt, A., & García, S. \(2020\). Automatic detection of building typology using deep learning methods on street level images. \*Building and Environment\*, 177. <https://doi.org/10.1016/j.buildenv.2020.106805>](#)
- Gupta, R. and Shah, M. 2020. RescueNet: Joint building segmentation and damage assessment from satellite imagery. In: *Proceedings - International Conference on Pattern Recognition*. Institute of Electrical and Electronics Engineers Inc., pp. 4405–4411. doi: 10.1109/ICPR48806.2021.9412295.
- [Hayes, J., Wilson, T. M., Deligne, N. I., Cole, L., & Hughes, M. 2017. A model to assess tephra clean-up requirements in urban environments. \*Journal of Applied Volcanology\*, 6\(1\). <https://doi.org/10.1186/s13617-016-0052-3>](#)
- Hayes, J.L. et al. 2019. Timber-framed building damage from tephra fall and lahar: 2015 Calbuco eruption, Chile. *Journal of Volcanology and Geothermal Research* 374(October 2015), pp. 142–159. Available at: <https://doi.org/10.1016/j.jvolgeores.2019.02.017>.

- He, K., Zhang, X., Ren, S., & Sun, J. 2015. Deep Residual Learning for Image Recognition. <http://arxiv.org/abs/1512.03385>
- Iannelli, G., & Dell'Acqua, F. (2017). Extensive Exposure Mapping in Urban Areas through Deep Analysis of Street-Level Pictures for Floor Count Determination. *Urban Science*, 1(2), 16. <https://doi.org/10.3390/urbansci1020016>
- Ishii, M., Goto, T., Sugiyama, T., Saji, H. and Abe, K. 2002. Detection of Earthquake Damaged Areas from Aerial Photographs by Using Color and Edge Information. pp. 23–25.
- Jenkins, S., & Spence, R. 2009. Vulnerability curves for buildings and agriculture The MIAVITA project is financed by the European Commission under the 7th Framework Programme for Research and Technological Development, Area "Environment", Activity 6.1 "Climate Change, Pollution and Risks."
- Jenkins, S.F., McSporry, A., Wilson, T.M., Stewart, C.S., Leonard, G.A., Cevuar, S., Garaebiti, E., In preparation. Tephra fall impacts to buildings: The 2017-2018 Manaro Voui eruption, Vanuatu. *Journal of Volcanology and Geothermal Research*
- Jenkins, S., Komorowski, J.C., Baxter, P.J., Spence, R., Picquout, A., Lavigne, F. and Surono. 2013. The Merapi 2010 eruption: An interdisciplinary impact assessment methodology for studying pyroclastic density current dynamics. *Journal of Volcanology and Geothermal Research* 261, pp. 316–329. Available at: <http://dx.doi.org/10.1016/j.jvolgeores.2013.02.012>.
- Jenkins, S. F., Spence, R. J. S., Fonseca, J. F. B. D., Solidum, R. U., & Wilson, T. M. 2014. Volcanic risk assessment: Quantifying physical vulnerability in the built environment. *Journal of Volcanology and Geothermal Research*, 276, pp 105–120. <https://doi.org/10.1016/j.jvolgeores.2014.03.002>
- Jenkins, S.F., Phillips, J.C., Price, R., Feloy, K., Baxter, P.J., Hadmoko, D.S. and de Bézizal, E. 2015. Developing building-damage scales for lahars: application to Merapi volcano, Indonesia. *Bulletin of Volcanology* 77(9). doi: 10.1007/s00445-015-0961-8.
- Johnson, J.M. and Khoshgoftaar, T.M. 2019. Survey on deep learning with class imbalance. *Journal of Big Data* 6(1). doi: 10.1186/s40537-019-0192-5.
- Joseph, E.P. et al. 2022. Responding to eruptive transitions during the 2020–2021 eruption of La Soufrière volcano, St. Vincent. *Nature Communications* 13(1). doi: 10.1038/s41467-022-31901-4.
- Jung, J., Kim, D. J., Lavalle, M., & Yun, S. H. (2016). Coherent Change Detection Using InSAR Temporal Decorrelation Model: A Case Study for Volcanic Ash Detection. *IEEE Transactions on Geoscience and Remote Sensing*, 54(10), 5765–5775. <https://doi.org/10.1109/TGRS.2016.2572166>
- Karnik, V. Schenkov, Z. Schenk, V. 1984. Vulnerability and the MSK scale. *Engineering Geology*, 20 (1984) pp161-168
- Kerle, N., Nex, F., Gerke, M., Duarte, D. and Vetrivel, A. 2019. UAV-based structural damage mapping: A review. *ISPRS International Journal of Geo-Information* 9(1), pp. 1–23. doi: 10.3390/ijgi9010014.
- Khajwal, A.B., Cheng, C.S. and Noshadravan, A. 2023. Post-disaster damage classification based on deep multi-view image fusion. *Computer-Aided Civil and Infrastructure Engineering* 38(4), pp. 528–544. doi: 10.1111/mice.12890.
- Lerner, G.A. et al. 2021. The hazards of unconfined pyroclastic density currents : a new synthesis and classification according to their deposits , dynamics , and thermal and impact This manuscript is a non-peer reviewed preprint submitted to *Journal of Volcanology and Geothermal* . pp. 1–48.
- López-Cifuentes, A., Escudero-Viñolo, M., Bescós, J., & García-Martín, Á. (2019). Semantic-Aware Scene Recognition. <https://doi.org/10.1016/j.patcog.2020.107256>
- Li, S., Tang, H., He, S., Shu, Y., Mao, T., Li, J. and Xu, Z. 2015. Unsupervised Detection of Earthquake-Triggered Roof-Holes from UAV Images Using Joint Color and Shape Features. *IEEE Geoscience and Remote Sensing Letters* 12(9), pp. 1823–1827. doi: 10.1109/LGRS.2015.2429894.
- Li, Y., Hu, W., Dong, H. and Zhang, X. 2019a. Building damage detection from post-event aerial imagery using single shot multibox detector. *Applied Sciences (Switzerland)* 9(6). doi: 10.3390/app9061128.

- [Li, D., Cong, A., & Guo, S. 2019b. Sewer damage detection from imbalanced CCTV inspection data using deep convolutional neural networks with hierarchical classification. \*Automation in Construction\*, 101, 199–208. <https://doi.org/10.1016/j.autcon.2019.01.017>](#)
- Lin, T.-Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C. L., & Dollár, P. 2014. Microsoft COCO: Common Objects in Context. <http://arxiv.org/abs/1405.0312>
- Lucks, L., Bulatov, D., Thönnessen, U. and Böge, M. 2019. 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, pp. 211–220. doi: 10.5220/0007253802110220.
- [Noh, H. Y., Jaiswal, K. S., Engler, D., & Wald, D. J. \(2020\). An efficient Bayesian framework for updating PAGER loss estimates. \*Earthquake Spectra\*, 36\(4\), 1719–1742. <https://doi.org/10.1177/8755293020944177>](#)
- Meng, S., Soleimani-Babakamali, M. H., & Taciroglu, E. 2023. Automatic Roof Type Classification Through Machine Learning for Regional Wind Risk Assessment. <http://arxiv.org/abs/2305.17315>
- Meredith, E.S., Jenkins, S.F., Hayes, J.L., Deligne, N.I., Lallemand, D., Patrick, M. and Neal, C. 2022. Damage assessment for the 2018 lower East Rift Zone lava flows of Kilauea volcano, Hawai'i. *Bulletin of Volcanology* 84(7). doi: 10.1007/s00445-022-01568-2.
- Moradi, M. and Shah-Hosseini, R. 2020. Earthquake Damage Assessment Based on Deep Learning Method Using VHR Images. *Environmental Sciences Proceedings* 5(1), p. 16. doi: 10.3390/iecg2020-08545.
- Naito, S. et al. 2020. Building-damage detection method based on machine learning utilizing aerial photographs of the Kumamoto earthquake. *Earthquake Spectra* 36(3), pp. 1166–1187. doi: 10.1177/8755293019901309.
- Nex, F., Duarte, D., Steenbeek, A. and Kerle, N. 2019. Towards real-time building damage mapping with low-cost UAV solutions. *Remote Sensing* 11(3), pp. 1–14. doi: 10.3390/rs11030287.
- Novikov, G., Trekin, A., Potapov, G., Ignatiev, V. and Burnaev, E. 2018. Satellite imagery analysis for operational damage assessment in emergency situations. In: *Lecture Notes in Business Information Processing*. Springer Verlag, pp. 347–358. doi: 10.1007/978-3-319-93931-5\_25.
- Post Disaster Needs Assessment (PDNA). 2022. St Vincent and the Grenadines
- [Pomonis, A. A., Spence, R., & Baxter, P. 1999. Risk assessment of residential buildings for an eruption of Furnas Volcano, Sao Miguel, the Azores. \*Journal of Volcanology and Geothermal Research\*, 92, pp 107-131.](#)
- Pi, Y., Nath, N.D. and Behzadan, A.H. 2020. Convolutional neural networks for object detection in aerial imagery for disaster response and recovery. *Advanced Engineering Informatics* 43. doi: 10.1016/j.aei.2019.101009.
- Ren, S., He, K., Girshick, R. and Sun, J. 2017. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39(6), pp. 1137–1149. doi: 10.1109/TPAMI.2016.2577031.
- [Román, A., Tovar-Sánchez, A., Roque-Atienza, D., Huertas, I.E., Caballero, I., Fraile-Nuez, E. and Navarro, G. 2022. Unmanned aerial vehicles \(UAVs\) as a tool for hazard assessment: The 2021 eruption of Cumbre Vieja volcano, La Palma Island \(Spain\). \*Science of the Total Environment\* 843. doi: 10.1016/j.scitotenv.2022.157092.](#)
- Shen, Y. et al. 2021. BDANet: Multiscale Convolutional Neural Network with Cross-directional Attention for Building Damage Assessment from Satellite Images. Available at: <http://arxiv.org/abs/2105.07364>.
- Singh, D. K., & Hoskere, V. 2023. Post Disaster Damage Assessment Using Ultra-High-Resolution Aerial Imagery with Semi-Supervised Transformers. *Sensors*, 23(19). <https://doi.org/10.3390/s23198235>
- Spence, R.J.S., Pomonis, A., Baxter, P.J., Coburn, A.W., White, M., Dayrit, M., and Field Epidemiology Training Program Team. 1996. 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

Deleted: ¶

- Spence, R., Martínez-Cuevas, S. and Baker, H. 2021. Fragility estimation for global building classes using analysis of the Cambridge earthquake damage database (CEQID). *Bulletin of Earthquake Engineering* 19(14), pp. 5897–5916. doi: 10.1007/s10518-021-01178-x.
- Spence, R.J.S., Kelman, I., Baxter, P.J., Zuccaro, G. and Petrazzuoli, S. 2005. *Natural Hazards and Earth System Sciences Residential building and occupant vulnerability to tephra fall*. St Vincent and the Grenadines population and housing census, 2012
- Szegedy, C., Vanhoucke, V., Ioffe, S., & Shlens, J. 2015. Rethinking the Inception Architecture for Computer Vision.
- Valentijn, T., Margutti, J., van den Homberg, M., & Laaksonen, J. (2020). Multi-hazard and spatial transferability of a CNN for automated building damage assessment. *Remote Sensing*, 12(17), 1–29. <https://doi.org/10.3390/rs12172839>
- Vetrivel, A., Gerke, M., Kerle, N., Nex, F., & Vosselman, G. 2018. 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 Journal of Photogrammetry and Remote Sensing*, 140, 45–59. <https://doi.org/10.1016/j.isprsjprs.2017.03.001>
- Wang, Z., Zhang, F., Wu, C., & Xia, J. (2024). Rapid mapping of volcanic eruption building damage: A model based on prior knowledge and few-shot fine-tuning. *International Journal of Applied Earth Observation and Geoinformation*, 126. <https://doi.org/10.1016/j.jag.2023.103622>
- Weber, E. and Kané, H. 2020. Building Disaster Damage Assessment in Satellite Imagery with Multi-Temporal Fusion. Available at: <http://arxiv.org/abs/2004.05525>.
- Williams, G.T., Jenkins, S.F., Biass, S., Wibowo, H.E. and Harijoko, A. 2020. Remotely assessing tephra fall building damage and vulnerability: Kelud Volcano, Indonesia. *Journal of Applied Volcanology* 9(1), pp. 1–18. doi: 10.1186/s13617-020-00100-5.
- Wilson, G., Wilson, T.M., Deligne, N.I. and Cole, J.W. 2014. Volcanic hazard impacts to critical infrastructure: A review. *Journal of Volcanology and Geothermal Research* 286, pp. 148–182. Available at: <http://dx.doi.org/10.1016/j.jvolgeores.2014.08.030>.
- Xu, J.Z., Lu, W., Li, Z., Khaitan, P. and Zaytseva, V. 2019. Building Damage Detection in Satellite Imagery Using Convolutional Neural Networks. (NeurIPS). Available at: <http://arxiv.org/abs/1910.06444>.
- Yi, W., Sun, Y., & He, S. 2018. Data Augmentation Using Conditional GANs for Facial Emotion Recognition. Progress In Electromagnetics Research Symposium. Japan. 1-4 August.
- Yorioka, D., Kang, H., Iwamura, K. 2020. Data Augmentation For Deep Learning Using Generative Adversarial Networks. IEEE 9th Global Conference on Consumer Electronics (GCCE)
- Yun, S.H. et al. 2015. Rapid damage mapping for the 2015 Mw 7.8 Gorkha Earthquake Using synthetic aperture radar data from COSMO-SkyMed and ALOS-2 satellites. *Seismological Research Letters* 86(6), pp. 1549–1556. doi: 10.1785/0220150152.
- Zhang, J.F., Xie, L.L. and Tao, X.X. 2003. Change Detection of Earthquake-damaged Buildings on Remote Sensing Image and its Application in Seismic Disaster Assessment. In: *International Geoscience and Remote Sensing Symposium (IGARSS)*. pp. 2436–2438. doi: 10.1109/igarss.2003.1294467.
- Zou, Z., Shi, Z., Guo, Y. and Ye, J. 2019. Object Detection in 20 Years: A Survey. Available at: <http://arxiv.org/abs/1905.05055>.

Formatted: Font: Cambria, 11 pt

Deleted: ¶