



## 1 Automating tephra fall building damage assessment using deep learning

2  
3 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  
4 Taisne <sup>2</sup>

5  
6 <sup>1</sup> Earth Observatory of Singapore @ NTU, Interdisciplinary Graduate Programme, Nanyang Technological University, Singapore, 639798

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

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

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

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

11  
12 *Correspondence to: Eleanor Tennant (eleanorm001@e.ntu.edu.sg)*

13  
14  
15 In the wake of a volcanic eruption, the rapid assessment of building damage is paramount for  
16 effective response and recovery planning. Uninhabited aerial vehicles, UAVs, offer a unique  
17 opportunity for assessing damage after a volcanic eruption, with the ability to collect on  
18 demand imagery safely and rapidly from multiple perspectives at high resolutions. In this work,  
19 we established a UAV-appropriate tephra fall building damage state framework and used it to  
20 label ~50,000 building bounding boxes around ~2,000 individual buildings in 2,811 optical  
21 images optical images collected during surveys conducted after the 2021 eruption of La  
22 Soufrière volcano, St Vincent and the Grenadines. We used this labelled data to train  
23 convolutional neural networks (CNNs) for: 1) Building localisation (average precision = 0.728);  
24 2) Damage classification into two levels of granularity: No Damage vs Damage (F1 score =  
25 0.809); and Moderate damage vs Major damage, (F1 score = 0.838) (1 is the maximum  
26 obtainable for both metrics). The trained models were incorporated into a pipeline along with  
27 all of the necessary image processing steps to generate spatial data (a shapefile with damage  
28 state attributes) for rapid tephra fall building damage mapping. Our pipeline is expected to  
29 perform well across other volcanic islands in the Caribbean where building types are similar,  
30 though would benefit from additional testing. Through cross validation, we found that the UAV  
31 look angle had a minor effect on the performance of damage classification models, while for the  
32 building localisation model, the performance was affected by both the look angle and the size  
33 of the buildings in images. These observations were used to develop a set of recommendations  
34 for data collection during future UAV tephra fall building damage surveys. This is the first  
35 attempt to automate tephra fall building damage assessment solely using post-event data. We  
36 expect that incorporating additional training data from future eruptions will further refine our  
37 model and improve its applicability worldwide. All trained models and pipeline code can be  
38 downloaded from GitHub to facilitate collaboration and development.



39 **1 Introduction**

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

46

47 After, or during, an eruption, the collection of empirical data detailing the damage incurred is  
48 critical to guiding the planning and implementation of response and recovery efforts. This  
49 involves estimation of damages and losses, which are needed to determine the necessary  
50 funding for repair or reconstruction; along with an assessment of building functionality, which  
51 can inform temporary housing requirements. In addition to its use in post disaster recovery,  
52 the collection of damage data are key to the development of fragility functions (Deligne et al.,  
53 2022), which relate hazard intensity to damage (e.g., Wilson et al., 2014; Williams et al., 2020;  
54 Spence et al., 2021) and can be used to inform resilient construction practises and/or for pre-  
55 event impact assessments.

56

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

68

69 Given the need for, but also the challenges associated with, conducting post-event building  
70 damage assessments quickly, approaches that use remotely sensed (RS) data, either optical or  
71 Synthetic Aperture Radar (SAR) imagery have been developed in volcanology (e.g., Jenkins et



72 al. 2013; Williams et al. 2020; Lerner et al. 2021; Biass et al. 2021; Meredith et al. 2022), and  
73 operationally by emergency management services (e.g., International Charter “Space and Major  
74 disasters”, Copernicus Emergency Management Service, ARIA: Advanced Rapid Imaging and  
75 Analysis system) (Yun et al., 2015)). The use of optical imagery largely consists of visual  
76 inspection, which may be influenced by image resolution and is prone to subjectivity (Novikov  
77 et al. 2018). Furthermore, visual inspection of satellite optical imagery can still be time  
78 consuming without crowd sourcing (e.g. Ghosh et al. 2011), and is constrained by satellite  
79 recurrence intervals and cloud cover. Automated SAR based methods (e.g., Yun et al., 2015) are  
80 not limited by cloud cover, but they may lack the resolution required for building level damage  
81 assessment (30 m for damage proxy maps generated from Sentinel data using the ARIA system;  
82 [https://aria-share.jpl.nasa.gov/20210409-LaSoufriere\\_volcano](https://aria-share.jpl.nasa.gov/20210409-LaSoufriere_volcano)).

83

84 While efforts to automate the assessment of building damage from volcanic hazards are  
85 minimal (to our knowledge there has been one study focusing on building damage from  
86 volcanic eruptions: Wang et al., 2024), attention has been given to more commonly occurring  
87 hazards such as earthquakes and hurricanes, with the development of both mono-temporal  
88 (post-event imagery only) and multi-temporal (uses pre- and post-event imagery) approaches  
89 (Table 1). Early approaches at automation with optical imagery used image processing  
90 methods, often focusing on identifying changes in pixel values between pre- and post-event  
91 imagery (e.g., Bruzzone and Fernández Prieto 2000; Ishii et al. 2002; Zhang et al. 2003). Image  
92 processing methods are susceptible to user biases such as the choice of thresholds that equate  
93 to distinct levels of damage severity, or damage states, and may require recalibration when  
94 applied to a new dataset. As a result, image processing methods were succeeded by the  
95 application of traditional machine learning algorithms that use ‘handcrafted’ image features.  
96 These features are observable properties that can be extracted from the image such as shape,  
97 colour, texture, and statistical properties of the image (e.g., Li et al. 2015; Anniballe et al. 2018;  
98 Lucks et al. 2019; Naito et al. 2020). The success of a given machine learning approach is  
99 dependent on the selection of the best features for the job; for example, a texture-based feature  
100 might be good for classifying buildings as damaged or not damaged due to an increased number  
101 of edges in damaged buildings but less useful for a task such as differentiating cats from dogs  
102 where the difference in textures between the classes is less significant. Deep learning, in  
103 particular the use of convolutional neural networks (CNNs), removes this need for feature  
104 selection. A CNN is a network of layers comprising filters which are small matrices of values.



105 When an image is passed through the network, at each layer the filters are convolved with the  
106 output from the previous layer to create a new representation of the image that is progressively  
107 more abstract with depth in the network. This process reduces the image's original spatial  
108 dimensions (X and Y) while increasing the number of channels, facilitating classification. During  
109 network training the filter values (known as weights) are optimised to reduce the loss between  
110 the predicted label for the image and the true label. Through this training a CNN learns the  
111 features of the images that are useful for classification. For a detailed background on deep  
112 learning see Aggarwal, (2018).

113

114 Thus far, deep learning models have been developed for optical image sets for hurricanes (Li et  
115 al. 2019; Dung Cao and Choe 2020; Pi et al. 2020; Cheng et al. 2021; Khajwal et al. 2023);  
116 earthquakes (Nex et al. 2019; Xu et al. 2019; Duarte et al. 2020; Moradi and Shah-Hosseini  
117 2020); wildfires (Galanis et al. 2021); volcanic hazards (Wang et al., 2024); and models that  
118 have been proposed for multiple hazards (e.g., Gupta and Shah 2020; Weber and Kané 2020;  
119 Shen et al. 2021; Bouchard et al. 2022) (Table 1). However, building damage caused by different  
120 hazards looks very different (e.g., damage caused by vertical loading from volcanic tephra fall  
121 vs ground shaking from an earthquake). These observable differences mean that an optical  
122 imagery multi-hazard damage classification model that performs consistently well across the  
123 different hazards is not yet achievable. Therefore, distinct models tailored for specific hazards  
124 are required (Nex et al., 2019, Bouchard et al., 2022). It follows that models may also benefit  
125 from being regionalised, given the differences in building typologies (construction material and  
126 styles) that can also affect the observable damage (Nex et al., 2019).

127

128 Many of the approaches for automating building damage assessment use both pre- and post-  
129 event imagery (Table 1), which makes the task more straightforward since any changes to the  
130 pre-event imagery can be considered damage. However, pre-event imagery at a high-enough  
131 resolution is not always available in post-disaster scenarios. The automated assessment of  
132 building damage from volcanic hazards using only post-event optical imagery has not yet been  
133 achieved in part due to absence of the large datasets that are needed in order to train models.  
134 The 2021 eruption of La Soufrière volcano, St Vincent and the Grenadines, provided an  
135 unprecedented opportunity for the collection of high-resolution UAV imagery enabling the  
136 development of fully automated models that can assess tephra fall building damage from post-  
137 event data only. With their growing ubiquity and low cost, UAVs have become an increasingly





138 useful tool during and after volcanic eruptions (e.g., Andaru and Rau 2019; Gailler et al. 2021;  
139 Román et al. 2022). UAVs offer a distinct advantage over satellite imagery because they can be  
140 scheduled at any point, they do not suffer from cloud obscuring the images as they fly at  
141 relatively low altitude, and they capture imagery from multiple perspectives, which may lead  
142 to increased ability to capture damage information. In this study we used UAV optical imagery  
143 collected after the 2021 eruption of La Soufrière volcano for tephra fall building damage  
144 assessment; the main contributions of our work are three-fold:

145

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

157

158

159 *Table 1. A non-exhaustive list of works using deep learning on optical imagery for building*  
160 *damage assessment. Studies use different scores to evaluate performance: F1 scores are in*  
161 *italics, mean average precision scores are underlined, accuracy scores in **bold**. For all scores,*  
162 *1 represents a perfect model.*

163

Study	Hazard	Number of damage classes	Pre and post?	Data type	Building localisation	Damage classification
Li et al. (2019)	Hurricane	2	P	airborne		<i>0.448</i>
Weber and Kane, (2020)	Multi	4	P & P	satellite (xBD)	0.835	0.697
Dung Cao and Choe. (2020)	Hurricane	2	P	satellite	-	<b>0.972</b>
Pi et al. (2020)	Hurricane	2	P	UAV, airborne		<u>0.745 (UAV)</u> <u>0.807 (airborne)</u>



Cheng et al. (2021)	Hurricane	5	P	UAV	<u>0.656</u>	<b>0.610</b>
Galanis et al. (2021)	Wildfire	2	P	satellite		0.981
Gupta and Shah (2020)	Multi	4	P & P	satellite (xBD)	0.840	0.740
Shen et al. (2021)	Multi	4	P & P	satellite (xBD)	0.864	0.782
Bouchard et al. (2022)	Multi	2	P & P	satellite (xBD)	0.846	0.709
Khajwal et al. (2023)	Hurricane	5	P	ground airborne	-	0.650
Singh and Hoskere, (2023)	Multi	5	P	satellite		<b>0.880</b>
Wang et al (2024)	Volcanic tephra	4	P & P	satellite	0.868	0.783
<b>Our work</b>	<b>Volcanic tephra</b>	<b>3</b>	<b>P</b>	<b>UAV</b>	<u>0.728</u>	<i>C1 0.809, 0.812</i>
					0.744	<i>C2 0.838, 0.838</i>

164

165

### 166 1.1 The 2020-2021 eruption of La Soufrière volcano St Vincent

167 La Soufrière St Vincent is an active stratovolcano standing at 1220 meters above sea level on  
 168 the island of St Vincent. On 27<sup>th</sup> December 2020 a thermal anomaly was detected inside the  
 169 summit crater by the NASA Fire Information for Resource Management System (FIRMS). This  
 170 was confirmed by the Soufrière Monitoring Unit to be caused by a new dome growing within  
 171 the crater. Dome growth continued for three months until 9 April 2021, when, following two  
 172 days of heightened seismic activity and lava effusion rate, the ongoing effusive eruption of La  
 173 Soufrière entered an explosive phase (Joseph et al. 2022). Between 9 – 22 April, a total of 32  
 174 distinct explosions occurred, with the tallest plumes reaching heights of up to 15 kilometers  
 175 above the vent (Joseph et al. 2022). Throughout this explosive phase, tephra blanketed the  
 176 island, resulting in a total deposition thickness of up to 16 centimeters in coastal communities  
 177 to the north of the island (Cole et al. 2023) (Figure 1).

178

179 The explosive phase was anticipated, and an evacuation order was issued on 8 April 2021 for  
 180 the ~16,000 residents in the northern part of the island (Joseph et al. 2022). As a result, there  
 181 were no reported fatalities directly attributable to the eruption, nevertheless, the overall  
 182 damage to infrastructure services and physical assets were estimated at XCD 416.07 million  
 183 (equivalent to USD 153.29 million) (PDNA, 2022). Approximately 63% of this monetary impact  
 184 was borne by the housing sector. In St. Vincent, residential buildings are typically single-story,  
 185 detached structures, with the majority in the more impacted north of the island (census  
 186 districts of Chateaubelair, Georgetown, and Sandy Bay: Figure 1) constructed using concrete



187 and blocks (84% in Chateaubelair, 74% in Georgetown, 50% in Sandy Bay), with sheet metal  
188 roofs (90-92% of all buildings in these areas) (SVG population and housing census, 2012).  
189



190  
191 *Figure 1. The island of St Vincent with UAV survey locations included in this work labelled and*  
192 *marked in black. Tephra isopachs (Cole et al., 2023) mark lines of constant total tephra thickness.*  
193 *Census districts referred to in the text are: a) Chateaubelair, b) Sandy Bay and c) Georgetown.*  
194 *Building footprints are marked in pink, data source: © OpenStreetMap contributors 2024.*  
195 *Distributed under the Open Data Commons Open Database License (ODbL) v1.0. Coordinate*  
196 *reference system: WGS 84 (EPSG:4326).*

197

## 198 **2 Method**

199 After the 2021 eruption of La Soufrière three UAV optical imagery datasets were collected to  
200 assess the extent of the damage. These were collected by different parties at separate times  
201 after the eruption. All UAV survey locations are shown in Figure 1, and representative examples  
202 of images can be found Section S1 of the supplementary material.

203



204 **2.1 Dataset description**

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

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

212

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

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

221

222 **Dataset 3: August -September 2021 (SRC)**

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

230

231 For all three datasets, image frames were extracted from the videos every two seconds, an  
232 interval chosen to reduce redundant homogeneous images, this resulted in a total of 7,956  
233 image frames. Due to the UAV surveying approach (i.e., hovering in one place for a while) many  
234 near-identical images were generated. To avoid potentially biasing the training towards  
235 overrepresented buildings we manually filtered out duplicate images. After filtering, and the  
236 removal of images with no buildings present the full combined dataset consisted of 2,811 image



237 frames. We labelled all images by drawing bounding boxes around each building present and  
238 storing the bounding box positions. In total 49,173 building bounding boxes were drawn  
239 around ~2,000 individual buildings (with some buildings being present in multiple images).  
240 Given the absence of detailed building inventory information, this number was approximated  
241 by overlaying Open Street Map building footprints with UAV GPS tracks. Bounding boxes were  
242 drawn by a team of five including the lead author, and all boxes were checked by the lead  
243 author. Each box was then assigned one of three damage states, which are described below. For  
244 consistency the damage states were assigned by the lead author. All labelling, modeling, and  
245 analysis were conducted using MATLAB 2023b.

246

## 247 **2.2 Developing and applying a building damage state framework**

248

249 Ground based damage state frameworks for tephra fall have previously split damage into five  
250 damage states, plus one not damaged, based on damage to three critical aspects of a building:  
251 the roof covering, the roof structure, and the vertical structure (Spence et al., 1996; Blong  
252 2003b; Hayes et al. 2019; Jenkins et al., under review). Remote damage assessments are often  
253 less able to resolve the detailed resolution achievable on the ground, and so a coarser resolution  
254 damage state framework is needed. In our study, most images depict buildings from an at nadir  
255 or close to nadir perspective making roof damage more discernible than damage to the vertical  
256 structure. Thus, we generated a damage state framework that is based on the proportion of  
257 observable damage to the roof, as in the work of Williams et al. (2020). Our final framework,  
258 which was developed over several iterations, classifies building damage into three classes: No  
259 observable or minor damage, Moderate damage, and Major damage (Table 2). Damage states  
260 are deliberately generic so that the range of possible damage to the range of different building  
261 types can be captured (Blong, 2003a). We included minor damage in the No damage class since  
262 the difference between the two can be subtle and not easily discernible through remote  
263 assessment. Furthermore, buildings with minor damage are typically habitable and unlikely to  
264 require costly repairs; therefore from a response and recovery perspective, we considered  
265 them better grouped with undamaged buildings.

266

267 In some images tarpaulins can be seen partially or fully covering roofs (~30 buildings). These  
268 were potentially placed to cover damage that occurred during the eruption, including corrosion  
269 due to prolonged presence of tephra on metal roofs or holes generated by nails lifted out



270 through sub-optimal cleaning approaches (VM personal communication). Alternatively,  
271 tarpaulins may have been placed as a preventative measure to help shed tephra (e.g., Ambae  
272 Vanuatu, Jenkins et al., under review). Erring on the conservative side, we considered buildings  
273 with a tarpaulin to be damaged; we assessed the severity of the damage for each building based  
274 on the level of visible deformation. We assigned buildings with a tarpaulin and no visible  
275 deformation to the moderately damaged class and those with a tarpaulin and visible  
276 deformation to the major damage class.

277

278 *Table 2. The damage assessment framework developed for our UAV optical imagery dataset*

279

---

	- No visible damage/or
No damage to minor damage	- Up to 10% of the roof covering missing; and/or - No roof or structural collapse; and/or - Visible damage to non-structural elements e.g., gutters or decorative elements (fascia).
Moderate damage	- Up to 50% roof area damaged (evidence of bending) or collapsed; may include light damage to vertical structure (e.g. wooden slats above windows broken).
Major damage	- More than 50% roof area damaged or collapsed; may include damage to the vertical structure including total building collapse.

---

280

281

### 282 **2.3 Model development**

283

284 After labelling, we split the full combined image dataset (2,811 frames from the UWI-TV, GOV  
285 and SRC sets) into train/validation/test sets (Figure 2). Given that a sizable proportion of the  
286 data were not geotagged, images from each location were kept together to assure the train and  
287 test sets were independent. The partitioning was chosen to include diversity in both the image  
288 sets (UWI-TV/GOV/SRC) and in the location, which affects the thickness of tephra fall received.  
289 We aimed for a standard data split of 80/10/10, with the majority of data assigned for training,  
290 however given the above constraints, this produced a split of 80% train, 8% validation, and



291 12% test (considering the number of bounding boxes and not the number of images). These  
292 data were used to develop our approach for building damage assessment. In line with the work  
293 of previous authors (Cheng et al. 2021; Bouchard et al. 2022), we split the building damage  
294 assessment task into two subtasks, training and evaluating models for building localisation,  
295 which consists of identifying building bounding boxes within the images and building damage  
296 classification separately. We chose to further divide the task of building damage classification  
297 into two separate classifications based on preliminary analysis.

298

299 In deep learning, the performance of a model and its optimal hyperparameters can be highly  
300 dependent on the characteristics of the dataset used for training, and hyperparameters that  
301 work well for one dataset may not work well for another. Therefore, it's common practice to  
302 experiment with different hyperparameter settings, model architectures and training  
303 strategies to find the configuration that performs the best for a particular problem. For each  
304 task in our damage assessment approach (localisation, classification 1, classification 2) we  
305 conducted a series of experiments using different image preprocessing approaches, CNN  
306 architectures, and combinations of hyperparameters with the aim of iterating towards the best  
307 experimental setup (Model selection: Section 3.1.1; Section 3.2.1).

308

309 Once we identified the best performing experimental setup for each task (building localisation,  
310 classification 1, classification 2), we combined the training and validation datasets and  
311 conducted K-fold cross-validation using the experimental setup and optimal hyperparameters  
312 that were identified (Cross validation: Section 3.1.3, Section 3.2.2). To test the robustness to  
313 location, we trained models on 9 out of the 10 locations present in the combined training and  
314 validation sets and evaluated each model's performance on the remaining location. To test the  
315 robustness to the dataset, we trained models and evaluated the performance for each of the  
316 three locations that have data from more than one dataset (e.g., Chateaubelair-GOV,  
317 Chateaubelair-UWI-TV, Chateaubelair-SRC) separately.

318

319 Following model selection and cross validation we calculated the final performance of the best  
320 model identified for each task (building localisation, classification 1, classification 2) on the test  
321 set. Finally, to see if better performance could be achieved with more data available for training,  
322 we retrained the models on the combined training and validation data before evaluating on the  
323 test data (Evaluation on the test set: Section 3.1.4, Section 3.2.3). All stages of model





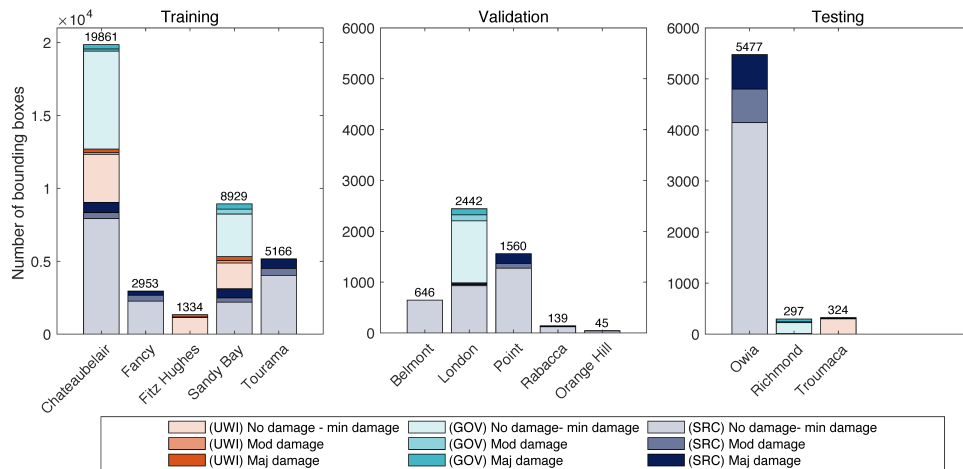
324 development, including model selection, cross validation, and final evaluation, are shown in  
 325 Figure 3 and more information about the specific experiments conducted for model selection is  
 326 given in Section S2 of the supplementary material.

327

328 In a post-disaster context, the seamless functioning of models will benefit from a sequential  
 329 workflow. Beyond the creation of distinct models for each task, we designed a comprehensive  
 330 pipeline that executes all optimized final models and the required processing steps to guide  
 331 images through the models (Figure 3d). The pipeline runs on an orthomosaic image and  
 332 generates spatial data in shapefile format that can readily be plotted in a GIS. In the following  
 333 sections we provide more detail on the algorithms and architectures used for each of the tasks,  
 334 and how the performance of each task was evaluated.

335

336

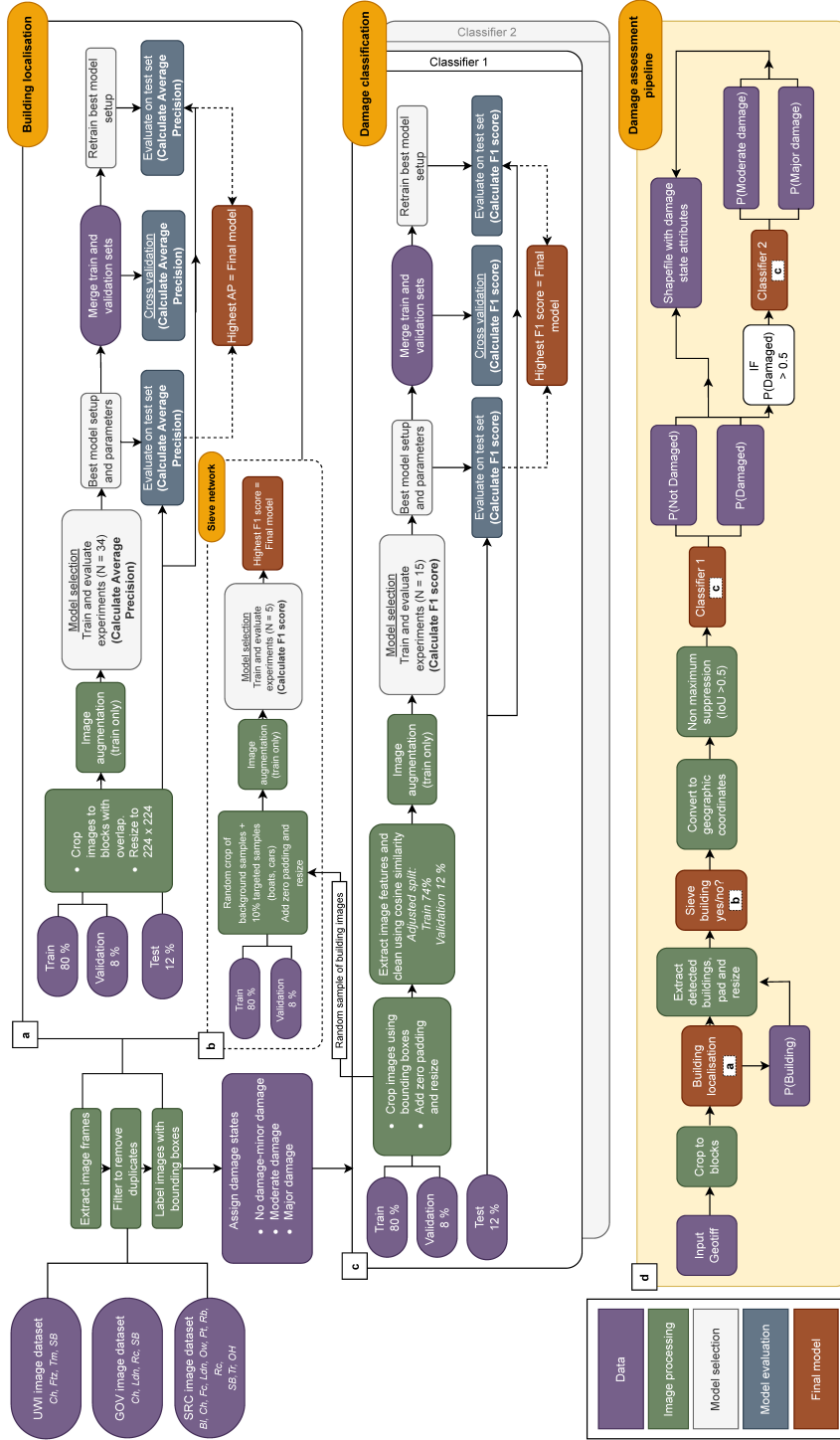


337

338

339

340 *Figure 2. The number of bounding boxes of each damage state in each UAV imagery dataset (UWI-*  
 341 *TV, GOV, SRC) for each of the locations in this study. Imagery was divided into three groups:*  
 342 *training, validation, and testing. The division of datasets between the three groups was chosen to*  
 343 *incorporate diversity in the image sets (UWI-TV/GOV/SRC), whilst keeping images from the same*  
 344 *location together and maintaining an approximate split of 80% training/10% validation/10%*  
 345 *testing.*





348

349 *Figure 3. A schematic showing the full methodology for a) developing a model for building*  
350 *localisation, b) developing a sieve network, which acts as an add on to the building localisation*  
351 *model, c) developing a model for building damage classification and d) the building damage*  
352 *assessment pipeline developed in this work. The pipeline incorporates the final trained models for*  
353 *building localisation and two stages of building damage classification along with all the necessary*  
354 *processing steps to link the models. Dataset locations referred to are: Bl – Belmont, Ch –*  
355 *Chateaubelair, Fc – Fancy, Ftz – Fitz Hughes, Ldn – London, OH – Orange Hill, Ow – Owia, Pt –*  
356 *Point, Rb – Rabacca, Rc – Richmond, SB – Sandy Bay, Tr – Tourama, Tm- Troumaca.*

### 357 **2.3.1 Building localisation**

358

359 For building localisation, we conducted experiments using the cutting edge two-stage object  
360 detector Faster R-CNN (Ren et al. 2017). Faster R-CNN is an improvement on the Fast R-CNN  
361 algorithm proposed by Girshick, (2015). The improvement comprises an initial region proposal  
362 network (RPN) which speeds up performance. Initially, image feature maps are extracted by  
363 passing the input image through a pretrained backbone CNN. The RPN then utilizes these  
364 features to generate proposals for potential object-containing areas, this is achieved by tiling a  
365 set of anchor boxes of assorted sizes across the extracted feature maps. The resulting region  
366 proposals are subsequently processed by the Fast R-CNN module, which includes a classifier  
367 that is used to determine the probability of the proposal containing an object, and a regressor  
368 that is used for adjusting the proposal box positions. When applied to a test image containing  
369 the relevant objects, Faster R-CNN outputs the positions within the image (X, Y, width, and  
370 height in pixels) of bounding boxes containing the object, and a confidence score for each box.  
371 As per customary practice (Zou et al. 2019) we used a confidence of  $> 0.5$  meaning that only  
372 boxes with confidence greater than this are output.

373

374 For object detection, to reduce model training and inference time, full sized images were split  
375 into image patches. Experiments conducted as part of building localisation model selection  
376 included variations in the size of these patches and the amount of overlap between patches. We  
377 also experimented with the development of separate models for images captured with different  
378 viewing angles, training for only the SRC portion of the dataset (images mostly at nadir) and  
379 the combined UWI-TV-GOV portion (images mostly off-nadir). A total of 34 experiments were  
380 conducted to find the best experimental setup for building localisation.



381 **2.3.2 Developing a sieve network**

382  
383 To improve the performance of the building localisation model we developed a small sieve  
384 network that runs as an add on to the Faster R-CNN building detector. Bounding boxes  
385 produced by the detector are passed to the sieve network to filter out detections that are false  
386 positives (i.e., detect a building when there is not one). To develop the dataset used for training  
387 and evaluating the sieve network we randomly cropped background samples from full sized  
388 images in the training and validation sets. Samples were cropped from each of the datasets, and  
389 samples containing buildings were removed until 100 no-building samples were achieved for  
390 each dataset. These samples were supplemented with an additional 10% targeted image  
391 samples on the observation that trained detectors were mistakenly detecting cars and boats.  
392 For the building dataset we stochastically sampled the equivalent number (n=990 train, 660  
393 validation) from the building images. Experiments for the sieve network were conducted using  
394 two different CNN architectures (ResNet50 and GoogleNet). A total of five experiments were  
395 conducted.

396

397 **2.3.3 Building damage classification**

398  
399 For building damage classification, we conducted experiments separately for classifiers 1 and  
400 2. Experiments consisted of fine-tuning two different pretrained CNNs to determine which was  
401 better and should be used in the final models for each classifier: ResNet50 (He et al., 2015)  
402 trained on the ImageNet dataset (Deng et al. 2009), and GoogleNet (Szegedy et al., 2015) trained  
403 on the places365 dataset (López-Cifuentes et al., 2019). Fine-tuning is a common approach to  
404 computer vision tasks where sufficiently large, labelled datasets are not available for the task  
405 at hand (typically hundreds of thousands of images are needed: Aggarwal, 2015). During fine-  
406 tuning, the high-level features that were learnt during the initial training on the large dataset  
407 can be leveraged for the new task. In addition to the different pretrained CNNs used,  
408 experiments also considered different ways of balancing the number of images for each damage  
409 state class (over-sampling the minority class, under-sampling the majority class and no  
410 balancing). When applied to a test building image, the trained classifier outputs the highest  
411 probability class and the associated probability. A total of 15 experiments were conducted for  
412 each of the classification tasks.

413



#### 414 **2.3.4 Model evaluation metrics**

415 For building localisation Faster R-CNN experiments, we evaluated performance using the  
416 average precision (AP) at an intersection over union (IoU) threshold of 0.5, and the F1 score.  
417 The AP is the most frequently used measure of an object detector's performance (Zou et al.,  
418 2019), and is calculated based on the number of times the detector gets it right (a true positive,  
419 TP) or wrong (a false positive, FP or a false negative, FN). A true positive occurs when the  
420 detector predicts a box that has an IoU with a labelled box of  $> 0.5$ . A false positive occurs when  
421 the detector predicts a bounding box that does not have an overlapping labelled box, while a  
422 false negative occurs where the detector fails to predict a box that is present in the labelled  
423 data. The relative proportions of these are used to calculate the precision and recall, where  
424 precision is the number of things that were predicted as positive that were correct: Precision =  
425  $TP/(TP+FP)$ , and recall is the number of things that are truly positive that were identified:  
426 Recall =  $TP/(TP+FN)$ . When a detector is run on a test image a confidence score is output for  
427 each predicted box (0-1). Once the trained detector has been run over the full test set, the  
428 precision and recall are calculated at different confidence score thresholds which can be plotted  
429 against one another to form a curve. The AP is the area underneath this precision-recall curve;  
430 it depicts the tradeoff between precision and recall and provides an overall measure of  
431 detection performance. AP values range between 0-1, where a higher value indicates a better  
432 performance.

433

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

439

### 440 **3 Results**

#### 441 **3.1 Building localisation**

##### 442 **3.1.1 Model selection**

443

444 The top five experiments (highest average precision) conducted for building localisation are  
445 shown in Table 3, with the full list of experiments provided in Table S2 of the supplementary  
446 material. Average precisions across the 34 experiments ranged from 0.295 to 0.701 (Table 3



447 and Table S2). We found that block size played an important role in model performance; out of  
448 the 34 experiments conducted, the top three used a block size of 550 x 550 pixels, which was  
449 the middle of the sizes tested (450, 550, 650). We observed that models trained on the full  
450 dataset performed better than models trained separately for the nadir (SRC) and off-nadir  
451 imagery sets (UWI-TV and GOV sets combined) (Table 3 and Table S2). More details on the  
452 results of experiments run for building localisation model selection can be found in Section S2.1  
453 of the supplementary material.

454

455 *Table 3. The five highest scoring (average precision) experiments conducted for building*  
456 *localisation, ordered by average precision. The full table consisting of all 34 experiments is*  
457 *provided in the supplementary material.*

458

Row id	Block size	Mixed block size?	Block overlap	Block resized?	Pretrained on best classifier?	Remove boxes < 32 x 32?	All training/ UWI-TV&GOV/ SRC	Max Average Precision	F1 score
1	550	N	50%	Y	N	Y	all	0.701	0.669
2	550	N	20%	Y	N	Y	all	0.700	0.668
3	550	N	20%	Y	Y	Y	all	0.700	0.642
4	650	N	50%	Y	N	Y	all	0.691	0.654
5	650	N	20%	Y	N	Y	all	0.678	0.670

459

### 460 3.1.2 Sieve Network

461 All trained sieve networks achieved macro and class F1 scores that were > 0.973 (Table 4). The  
462 best performing sieve network experiment achieved a macro F1 score of 0.977. The best  
463 detector identified through model selection (Table 3, row 1) achieved an F1 score of 0.669  
464 (Table 6), with a precision and recall of 0.588 and 0.776, respectively, on the validation data.  
465 The lower value of precision is due to the substantial number of false positive detections. After  
466 the results of the detector were passed through the sieve network, the number of false positives  
467 was reduced, with an improved F1 score of 0.712 (Table 6).

468

469 *Table 4. Experiments conducted for the sieve network, a small network that runs on the boxes*  
470 *produced by the object detector. Results are ordered from high to low by the Macro F1 score.*  
471 *ResNet50 and GoogleNet refers to the convolutional neural network architecture used in the*



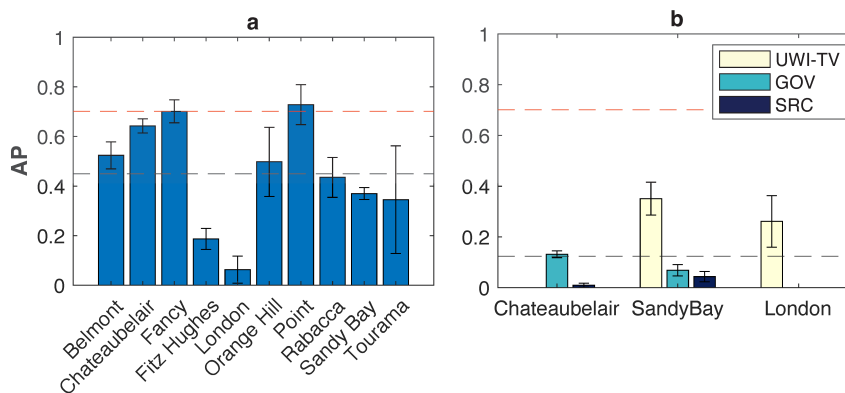
472 *experiment; the value after the underscore reflects the experiment ID where different IDs have*  
473 *different training parameters (see Section S2 of the supplementary material).*

Experiment ID	F1 building	F1 background	F1 macro
ResNet50_4	0.978	0.977	0.977
ResNet50_3	0.977	0.976	0.977
ResNet50_1	0.975	0.975	0.975
ResNet50_2	0.976	0.974	0.975
GooglNet_1	0.973	0.973	0.973

474  
475  
476

### 477 3.1.3 Cross validation

478 Cross validation was conducted for the single best performing building localisation model  
479 (without the sieve network) to understand how the choice of training and validation data  
480 affects performance, along with the potential for the model to generalize to a new dataset. We  
481 found that the performance of the selected object detector varies, depending upon the location  
482 (Figure 4a) or imagery dataset (Figure 4b) used for testing. For models tested on different  
483 locations (Figure 4a) average precisions > 0.7 were obtained for Point and Fancy in line with  
484 AP achieved on the full validation set (0.701). The lowest AP values were for London (0.063)  
485 and Fitz Hughes (0.187). The standard deviation (SD) (Figure 4) shows the variability in  
486 performance between the three replicates that were trained for each test, which arises due to  
487 the stochastic nature of the training process. For models tested on the different imagery  
488 datasets individually the AP was low (Figure 4b), with a mean value across all datasets of < 0.2.  
489 For all three locations (Chateaubelair, Sandy Bay, London), AP for models evaluated on the SRC  
490 dataset were higher than for the UWI-TV or GOV datasets.



491





492

493 *Figure 4. Cross validation of the best experimental setup for building localisation models which*  
494 *are trained to predict building box positions within the image. a) The effect of changing the*  
495 *location used as the test set on detector average precision (AP) and b) the effect of changing the*  
496 *imagery dataset (UWI-TV/GOV/SRC) used as the test set on AP. For b) cross validation of the*  
497 *imagery dataset, models are trained on all data from that location excluding the location used for*  
498 *testing as indicated by the bar. For London there is data from the GOV dataset, however the*  
499 *number of images in the SRC dataset is insufficient for training, so no bar is shown for GOV. The*  
500 *AP shown is the mean value from three trained models with the same setup while the error bars*  
501 *show the standard deviation. Black dashed lines show the mean AP value across all cross*  
502 *validation trained models, red dashed lines show the best AP from the experiments (0.701: Table*  
503 *3).*

504

#### 505 **3.1.4 Evaluation on the test set**

506 Evaluation of the best detection model on the test set, which consists of completely unseen data  
507 from Owia, Richmond and Troumaca (Figure 2) produced an AP value that is the same as the  
508 value on the validation data (0.701) (Table 3). Retraining the best experimental setup for the  
509 detector using the combined training and validation data caused the AP when applied to the  
510 test data to increase to 0.751 prior to sieving and 0.728 after sieving. Comparing the precision  
511 and recall of the retrained detector and the retrained detector + sieve network shows that while  
512 the AP is higher for the retrained detector without the sieve, the addition of the sieve network  
513 creates a better balance between the precision and recall, reflected in a higher F1 score. We  
514 therefore selected the retrained detector + sieve network as the final building localisation  
515 model (Table 6).

516

### 517 **3.2 Building damage classification**

#### 518 **3.2.1 Model selection**

519 The top five experiments (highest macro F1 score) conducted for building damage classification  
520 are shown in Table 5, with the full lists provided in Tables S3 and S4 of the supplementary  
521 material. Macro F1 scores ranged from 0.753 to 0.836 and 0.776 to 0.810 for classifier 1 and 2  
522 respectively (Tables 5, S3, S4). The best performing models for both classifiers used the  
523 ResNet50 architecture rather than GoogleNet with an unbalanced dataset. For Classifier 1 the



524 best model had  $F1 = 0.962$  for the Not Damaged class and  $F1 = 0.710$  for the Damaged class.  
525 While for Classifier 2 the Moderate damage class had  $F1 = 0.770$  and Major damage  $F1 = 0.851$ .

526

527 *Table 5. The top five experiments conducted for each of the building damage classifiers, ordered*  
528 *by the macro F1 score. The full list consisting of all 15 experiments for each classifier is provided*  
529 *in Tables S3 and S4 of the supplementary material.*

530

Classifier 1						
Row ID	Architecture	Class balancing: Not Balanced/ under-sampled/ over-sampled	Dropout	F1 Not Damaged	F1 Damaged	F1 Macro
1	Resnet50	not	0.4	0.962	0.710	0.836
2	Resnet50	not	0	0.960	0.696	0.828
3	Resnet50	not	0.6	0.957	0.699	0.828
4	Resnet50	not	0.2	0.962	0.692	0.827
5	Resnet50	under	0	0.951	0.646	0.799

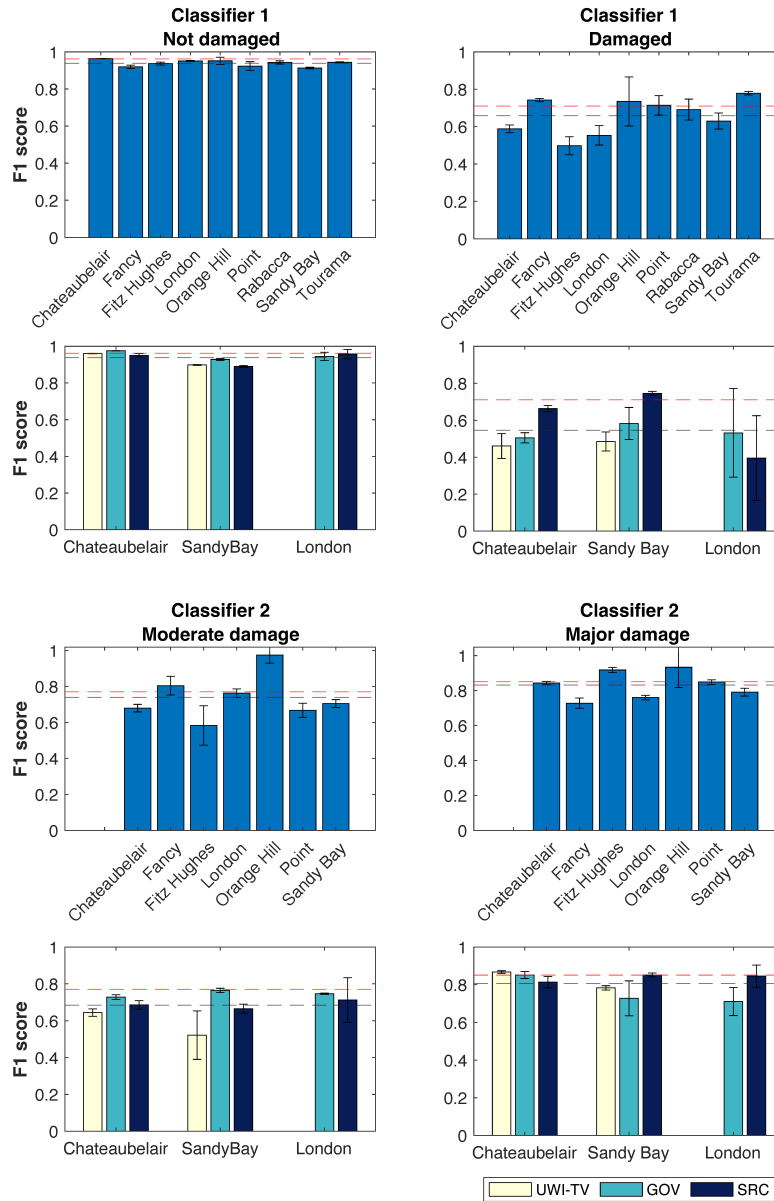
Classifier 2						
Row ID	Architecture	Class balancing: Not Balanced/ under-sampled/ over-sampled	Dropout	F1 Mod damage	F1 Maj damage	F1 Macro
1	Resnet50	not	0	0.770	0.851	0.810
2	GoogleNet	over	0	0.737	0.848	0.793
3	Resnet50	over	0	0.749	0.835	0.792
4	Resnet50	not	0.4	0.749	0.835	0.792
5	Resnet50	under	0.6	0.735	0.845	0.790

531

532

### 533 3.2.2 Cross validation

534 Cross validation was conducted for both of the single best performing models for Classifiers 1  
535 and 2 identified through model selection. As was the case for the best building localisation  
536 model, this was done to understand how the choice of training and validation datasets affected  
537 model performance and to understand the potential for our model to generalize to a new  
538 dataset.



539

540 *Figure 5. Cross validation for Classifiers 1 and 2. For rows 1 and 3 the best experimental setup*  
 541 *was retrained on all the data from locations in the combined training and validation data and*  
 542 *evaluated on the location shown. For rows 2 and 4 the best experimental setup was retrained on*  
 543 *all the data from the location shown and evaluated on each dataset (UWI-TV/GOV/SRC)*  
 544 *separately. Each training was conducted three times, the value plotted is the mean, and the error*  
 545 *bars show the standard deviation. Black dashed lines show the mean F1 score across all cross*



546 *validation trained models, red dashed lines show the best F1 score for each class from the*  
547 *experiments (Table 5).*

548

549 Figure 5 shows that the performance of Classifier 1 for the Not damaged class is consistent  
550 across the distinct locations and datasets used for testing with mean F1 scores between 0.913-  
551 0.983 for locations and 0.898-0.976 for datasets. For the Damaged class there is more variety  
552 in the performance the choice of location and dataset used for evaluation. The mean F1 scores  
553 for the separate locations range from 0.588 (Fitz Hughes) to 0.779 (Tourama) while for the  
554 different datasets the range is 0.393 (London-SRC) to 0.745 (Sandy Bay-SRC).

555

556 For Classifier 2, the Moderate damage class is more sensitive to the choice of location used for  
557 the validation than the Major damage class (Figure 5). For the Moderate damage class, the mean  
558 F1 score ranged from 0.583-0.974. Similarly to Classifier 1, Fitz Hughes produced the lowest  
559 mean F1 score, whereas the highest score was produced for Orange Hill. For the Major damage  
560 class F1 scores for the distinct locations are between 0.728-0.933. For Classifier 2 the sensitivity  
561 to the choice of dataset (UWI-TV/GOV/SRC) for the Moderate damage class is greater than for  
562 the Major damage class. For Moderate damage, the range is between 0.522-0.746, while for  
563 Major damage the range is from 0.711-0.867.

564

### 565 **3.2.3 Evaluation on the test set**

566 Evaluation of the single best models for classification 1 and classification 2 on the unseen test  
567 set produced Macro F1 scores that were comparable with the scores for the validation set:  
568 0.829 for Classifier 1 and 0.791 for Classifier 2 (Table 6). For Classifier 2, retraining the model  
569 on the combined training and testing data increased the Macro F1 score from 0.791 to 0.838.  
570 Whereas for Classifier 1 retraining produced a slightly lower Macro F1 score (0.809 compared  
571 to 0.829). Nevertheless, the retrained model for Classifier 1 achieved a higher recall on the  
572 damaged class than the non-retrained model. In an operational setting it's desirable to correctly  
573 classify as many of the damaged buildings as possible, since in our pipeline these will be passed  
574 onto Classifier 2, therefore we took the retrained models for both classifiers as the final models.  
575 The confusion matrices for both final models are plotted in Figure 6, these show class accuracy  
576 i.e., how many of the true class were correctly classified. For Classifier 1 89% of the Not  
577 damaged buildings were correctly classified, and 73% of the Damaged buildings were correctly



578 classified. For Classifier 2 81% of the moderately damaged buildings were correctly classified,  
 579 while 87% of the buildings with major damage were correctly classified.

580

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

585

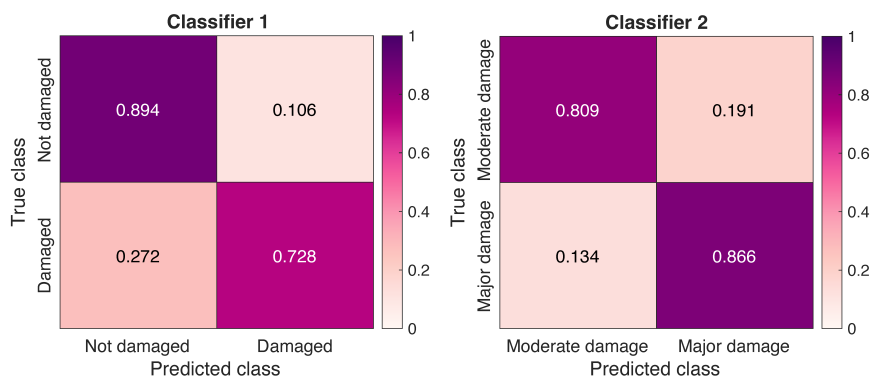
	Validation set				Test set									
	AP	P	R	F1	AP	P	R	F1						
Detector (0.5 conf)	0.701	0.588	0.776	0.669	0.701	0.604	0.776	0.679						
Detector + Sieve (0.5 conf)	0.681	0.695	0.730	0.712	0.668	0.606	0.757	0.673						
Detector retrain					0.751	0.642	0.816	0.719						
Detector retrain +sieve					<b>0.728</b>	<b>0.710</b>	<b>0.782</b>	<b>0.744</b>						
	Not damaged			Damaged			F1 Macro	Not damaged			Damaged			F1 Macro
	P	R	F1	P	R	F1		P	R	F1	P	R	F1	
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			F1 Macro	Mod Damage			Maj Damage			F1 Macro
	P	R	F1	P	R	F1		P	R	F1	P	R	F1	
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>

586

587

588

589



590

591

592 *Figure 6. Confusion matrices for the final models for Classifiers 1 and 2 (Classifier 1 retrain*  
593 *and Classifier 2 retrain; Table 6) evaluated on the test dataset. Confusion matrices show the*  
594 *proportions of each class that are classified correctly. Horizontal values sum to 100%.*

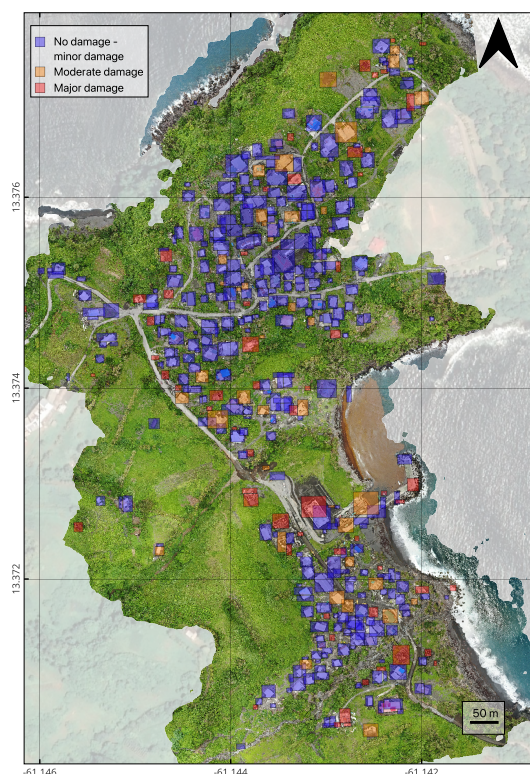
595

#### 596 **4 Example application of the full damage assessment pipeline**

597

598 In this work we have developed separate models for building localisation and two stages of  
599 damage classification. However, in an operational context models need to work sequentially,  
600 this led to the development of our damage assessment pipeline (outlined in Figure 3d). The  
601 pipeline operates on an orthomosaic image, which can easily be generated using software such  
602 as Agisoft Metashape, OpenDroneMap or Pix4D. The pipeline outputs a shape file, with the  
603 following *attributes* for each building that is detected: *detection* (box confidence score),  
604 *ClassPred\_1* (output class from classifier 1, damaged or not damaged), *ClassProb\_1* (the  
605 probability of that class), *ClassPred\_2* (output class from classifier 2, moderate damage or major  
606 damage, this is only run if classifier 1 outputs damage), *ClassProb\_2* (the probability of the class  
607 output by classifier 2). Figure 7 shows an example tephra fall building damage map produced  
608 by running the pipeline on an orthomosaic image generated using Agisoft Metashape software  
609 and plotting both the orthomosaic and output shapefile in QGIS. The example which consists of  
610 417 buildings took 1 hour to run on a standard 16GB RAM 2021 MacBook Pro, with an M1 Pro  
611 chip. Most of the inference time was attributed to the building localisation module in the  
612 pipeline, which may be bypassed if building footprints are already available. When only the  
613 classifiers were run the time taken to run was reduced to < 5 mins.

614



615

616 *Figure 7. Application of the full tephra fall building damage assessment pipeline on the*  
617 *orthomosaic image for Owia. Coordinate reference system: WGS 84 (EPSG:4326). Satellite*  
618 *basemap © Google Maps 2024.*

619

## 620 **5 Discussion**

621

622 In this work we have developed models for building localisation, and two levels of damage  
623 classification for building damage resulting from tephra fall. Our final models demonstrate  
624 strong performance for both building localisation (AP = 0.728; F1 = 0.744) and building damage  
625 classification (Classifier 1, F1 = 0.809, Classifier 2, F1 = 0.838). Despite using post-event  
626 imagery only, which makes the task more challenging than approaches that use both pre- and  
627 post-event imagery, our results are comparable to existing optical imagery building damage  
628 assessments developed for various hazards that use both mono-temporal and multitemporal  
629 images (Table 1).

630





631 **5.1 Building localisation**

632  
633 Through running our building localisation experiments we found that the pre-processing of  
634 images before detector training (particularly the block size) significantly influenced detector  
635 performance. Cross-validation results demonstrated variability in average precision (AP) for  
636 models trained on different locations and imagery datasets (UWI-TV/GOV/SRC) (Section 3.1.3;  
637 Figure 4). Deep learning models are known to perform well when data come from the same  
638 distribution, though have more difficulty when working with out of distribution samples. Given  
639 the relatively consistent building typology across locations (the majority of buildings observed  
640 are detached single storey buildings with either a gable or hip shaped metal sheet roof; a lesser  
641 proportion have flat concrete roofs), the differences in AP are likely due to observable  
642 variations in UAV altitude, off-nadir angles, tephra thicknesses, and varying training sample  
643 sizes.

644  
645 The London images (from SRC and GOV datasets) and Fitz Hughes images (UWI-TV) exhibited  
646 the lowest average precisions (Figure 4). Both London datasets featured smaller buildings than  
647 the rest of the locations, evident in Section S3 of the supplementary material, while the UWI-TV  
648 images had more tephra on the ground, which affects background colour and, viewed buildings  
649 from an off-nadir perspective. The training data, predominantly nadir images from the SRC  
650 dataset, had fewer UWI-TV examples which are off-nadir and, collected more closely in time  
651 after the eruption, meaning more tephra was present in images. This under representation in  
652 the training data may have impacted the model's ability to recognize such instances in the test  
653 data. The application of sampling approaches like those used for the damage states in the  
654 classification model development (over or under sampling) could have been applied to balance  
655 the data, however the SRC dataset is much larger than either of the UWI-TV and GOV sets  
656 (Figure 2), therefore we did not use this approach as we considered that oversampling would  
657 introduce significant bias towards the specific examples in the under-represented dataset,  
658 whereas through under sampling we would lose a large amount of the data that are available  
659 to learn from. Future work might consider the application of generative AI algorithms such as  
660 generative adversarial networks (GANs) to expand the dataset (e.g., Yi et al. 2018; Yorioka et  
661 al., 2020), although more work needs to be done to quantify the diversity in the generated data.

662  
663 The variability in cross-validation results for the building localisation model (Section 3.1.3)  
664 likely comes from a combination of the above factors (differences in UAV altitude, off-nadir



665 angles, tephra thickness, and varying training sample sizes), and suggests that there was  
666 insufficient information in the training data for our detection models to perform well across the  
667 range of characteristics present. However, this requires further investigation to separate the  
668 unique effect of each aspect.

669

## 670 **5.2 Building damage classification**

671

672 The final classification models achieved better performance than the final localisation model  
673 with macro F1 scores of 0.809 and 0.838 on the test data (Table 6). Cross-validation showed  
674 that classification models were less sensitive than the localisation model to the choice of  
675 datasets used for training and evaluation (Section 3.2.2). We found that class wise our models  
676 performed better on the not damaged class followed by the major damage class. This is in  
677 agreement with other multi-class studies that have found the extremities of the damage state  
678 scheme used to be easier to classify than the intermediate ones (Kerle et al. 2019).

679

## 680 **5.3 Application of the pipeline**

681

682 Our pipeline consists of separate models for localisation and building damage classification.  
683 One of the benefits of this is that in locations where precise building location information is  
684 available for the assessment area, the localisation step can be bypassed and only the classifiers  
685 run. This not only enhances overall performance but also significantly reduces computation  
686 time. Furthermore, either of the classifiers can be run independently and/or combined with  
687 other damage assessment procedures; for example, an initial synthetic aperture radar (SAR)  
688 based assessment (e.g., Yun et al. 2015, Jung et al., 2016), could be followed with our classifier  
689 2 to provide additional granularity on the severity of the damage at a building level rather than  
690 a pixel level.

691

## 692 **5.4 Generalisability to other locations**

693

694 Our models have performed well for images collected on the island of St Vincent where building  
695 typologies are relatively consistent. We therefore expect that our models will perform well in  
696 other locations with similar building types, such as the other islands in the Lesser Antilles. This  
697 hypothesis should be validated through further testing. In absence of additional UAV datasets  
698 that include damaged buildings, testing can be done by conducting pre-event surveys to test the



699 performance of the building localisation model and classifier 1 for the Not damaged class. While  
700 this is unable to assess the ability of our approach to classify damage, it would provide *some*  
701 indication of performance following an event in a new location.

702

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

715

## 716 **5.5 Improving model performance and future perspectives**

717

718 The advantages of acquiring additional UAV datasets both before and after an event have been  
719 outlined in Section 5.3. In addition to this, pre-event surveys may be particularly beneficial in  
720 constructing building inventories, which include details about building typologies such as  
721 construction materials and styles. Surveys can be interrogated manually to extract building  
722 attributes or using machine learning methods such as the work of Meng et al., (2023). Prior to  
723 an eruption, given knowledge about the building typologies, an idea about how the buildings  
724 will respond under certain tephra loadings (i.e., the forecasted damage state) can be obtained  
725 through the application of fragility functions. This information could enhance our model by  
726 serving as prior information. The forecasted damage state could be subsequently refined  
727 through Bayesian updating based on our damage assessment models output.

728

729 Alternatively, with ample individual building inventory data available, tailored damage  
730 classification models for specific building typologies could be developed and applied. The



731 rationale is that a model dedicated to a specific building type is expected to outperform a  
732 generic multi-typology model.

733

734 In this work, we established a three class damage state framework. Existing frameworks that  
735 were developed for ground based tephra fall damage assessment split damage into five damage  
736 states classes and one non-damage class (Spence et al, 1996; Blong, 2003; Hayes et al., 2019;  
737 Jenkins et al., in review) however in our preliminary analyses we found that: 1) in many images  
738 we were unable to confidently apply a six-class scheme due to only being able to see one side  
739 of the building, and 2) there were not enough examples of each damage state class to be able to  
740 train a six-class model. Nevertheless, the damage states developed in our work can be equated  
741 to existing damage states generated for ground surveys such that: No damage – minor damage  
742 = DS0-DS1, Moderate damage = DS2 and Major damage = DS3-5. With the addition of future  
743 tephra fall building damage datasets a finer resolution damage state framework may be applied  
744 that is capable of providing more detail on the observable damage. We developed our approach  
745 using deep learning on 2D optical imagery, while some studies have used 3D point-cloud  
746 information (Cusicanqui et al., 2018), or combined point cloud information with deep learning  
747 on optical imagery for damage level classification (Vetrivel et al., 2018). While the use of 3D  
748 spatial data has shown potential, and may be used to provide additional granularity to our  
749 damage states we opted against integrating point cloud analyses into our model. This decision  
750 was motivated by the considerably longer processing times associated with such an approach,  
751 which would undermine the swift processing requirement inherent in our methodology.

752

753

## 754 **5.6 Caveats**

755

756 During the assignment of building damage states, uncertainties arose, particularly concerning  
757 the interpretation of tarpaulins and, pre-existing damage. For tarpaulins, the ambiguity arose  
758 from whether these were either strategically placed prior to the eruption as preventative  
759 measures to cause tephra to slide off the roof more easily; or they were placed post event to  
760 cover damage caused by tephra fall. Additionally, in certain instances, distinguishing between  
761 a collapsed roof and a section of the building initially lacking roofing material—possibly  
762 functioning as a walled storage area—proved challenging. Pre-existing damage not related to  
763 volcanic activity or buildings that were under construction at the time of image acquisition



764 were considered as damaged and classified accordingly. Pre-event imagery would have  
765 provided clarity on these matters, however this was not available at high enough resolution for  
766 this region.

767

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

776

777

## 778 **6 Recommendations for UAV building damage assessment data collection**

779

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

787

- 788 • The bulk of our dataset was collected several months after the eruption of La Soufrière  
789 however, for generating a global dataset that can be used for response and recovery,  
790 models should ideally be trained on images collected shortly (days to weeks) after an  
791 event.
- 792 • Flight paths should be pre-programmed to ensure comprehensive coverage of the area  
793 and limit bias associated with overrepresentation of certain buildings. Ideally two flights  
794 would be conducted with two sets of perpendicular flight lines to capture buildings from  
795 a different perspective. GPS positioning should be enabled.



- 796       • A fixed altitude of 50-80 m above the ground should be maintained where possible. This  
797       is appropriate to capture sufficient data for accurate damage classification based on the  
798       established framework and strikes a balance between detailed information capture and  
799       overall coverage. In mountainous areas this may not be achievable for some UAV types.  
800       In which case a uniform height should be maintained such that buildings size is  
801       consistent across image frames.
- 802       • We suggest a slightly off-nadir camera positioning (~5-15°), which is sufficient to  
803       capture any bending in the roof that may not be captured from a nadir perspective.
- 804       • Overlap between images should be enough to generate orthoimages, 80% forward and  
805       70% lateral overlap is sufficient.

806

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

- 811           ○ Facilitates the creation of building inventories.
- 812           ○ Enables precise comparison of pre- and post-event imagery, reducing uncertainty  
813           regarding initial building conditions.
- 814           ○ Supports the development of high-resolution change detection models  
815           potentially yielding more accurate results than relying solely on post-event  
816           imagery.
- 817           ○ Provides an opportunity for UAV pilots to gain experience in capturing building  
818           datasets during 'quiet times'.

## 819   7   Conclusions

820

821       Following a large tephra fall event, building damage assessment needs to be conducted rapidly  
822       for the purpose of response and recovery, and for the collection of data that can be used to  
823       forecast building damage from future events. By leveraging post-event optical imagery  
824       obtained after the 2021 eruption of La Soufrière volcano on the island of St Vincent, and  
825       convolutional neural networks, we have developed an automated tephra fall building damage  
826       assessment pipeline. The pipeline incorporates models for building localisation and two  
827       distinct levels of damage classification: distinguishing between no damage and damage, as well  
828       as between moderate and major damage, which were trained and evaluated separately. When



829 provided with UAV optical imagery, our pipeline can rapidly generate spatial building damage  
830 information. Our models perform well for the St Vincent datasets and are anticipated to  
831 perform well in locations where building typologies are similar, but this requires more testing  
832 to understand the limits of their application.

833

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

844

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

850

## 851 **8 Author contributions**

852

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

## 856 **9 Competing interests**

857

858 The authors declare no competing interests.

## 859 **10 Acknowledgements**

860





861 We are indebted to Monique Johnson: The UWI Seismic Research Centre, Javid Collins: UWITV,  
862 Nikolai Lewis and Marla Mulraine: The Government of St Vincent and the Grenadines Ministry  
863 of Transport, Works, Lands and Surveys, and Physical Planning, for sharing their UAV data and  
864 collaborating on this work. All images and data provided in this study have been approved for  
865 publication by the local agency responsible for monitoring geohazards in St Vincent: The UWI  
866 Seismic Research Centre. We would like to thank Chee Jain Hao Denny, Sim Yu Yang, Isaiah Loh  
867 Kai En, Huang Wanxin for their assistance with data preparation. We are very grateful to  
868 Sébastien Biass, Vanesa Burgos, Elinor Meredith, and Alberto Ardid, for interesting discussions  
869 around machine learning and building damage assessment.

870

## 871 **11 Data availability**

872

873 All trained models along with the code required to execute the damage assessment pipeline  
874 and instructions for usage are provided at:

875 <https://github.com/EllyTennant/UAVdamageAssessment>

876

## 877 **12 Funding**

878

879 This research was supported by the Earth Observatory of Singapore via its funding from the  
880 National Research Foundation Singapore and the Singapore Ministry of Education under the  
881 Research Centres of Excellence initiative and comprises EOS contribution number 596.  
882 Additional support was provided by the AXA Research Fund as part of the Joint Research  
883 Initiative on Volcanic Risk in Asia.

884

## **13 References**

- 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>



- 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.
- 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.
- Bouchard, I., Rancourt, M.É., 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>.
- 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>
- 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.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>



- 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.F., McSporran, A., Wilson, T.M., Stewart, C.S., Leonard, G.A., Cevuard, 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., Phillips, J.C., Price, R., Feloy, K., Baxter, P.J., Hadmoko, D.S. and de Bélizal, 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., Lavalley, 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>
- 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. 2019. Building damage detection from post-event aerial imagery using single shot multibox detector. *Applied Sciences (Switzerland)* 9(6). doi: 10.3390/app9061128.
- 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.
- 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
- 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
- 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.
- 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>.