



**HAL**  
open science

# A human-on-the-loop approach for labelling seismic recordings from landslide site via a multi-class deep-learning based classification model

Jiaxin Jiang, David Murray, Vladimir Stankovic, Lina Stankovic, Clement Hibert, Stella Pytharouli, Jean-Philippe Malet

## ► To cite this version:

Jiaxin Jiang, David Murray, Vladimir Stankovic, Lina Stankovic, Clement Hibert, et al.. A human-on-the-loop approach for labelling seismic recordings from landslide site via a multi-class deep-learning based classification model. *Science of Remote Sensing*, 2025, pp.100189. 10.1016/j.srs.2024.100189 . hal-04865297

**HAL Id: hal-04865297**

**<https://hal.science/hal-04865297v1>**

Submitted on 7 Jan 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

## Journal Pre-proof

A human-on-the-loop approach for labelling seismic recordings from landslide site via a multi-class deep-learning based classification model

Jiaxin Jiang, David Murray, Vladimir Stankovic, Lina Stankovic, Clement Hibert, Stella Pytharouli, Jean-Philippe Malet



PII: S2666-0172(24)00073-7  
DOI: <https://doi.org/10.1016/j.srs.2024.100189>  
Reference: SRS 100189

To appear in: *Science of Remote Sensing*

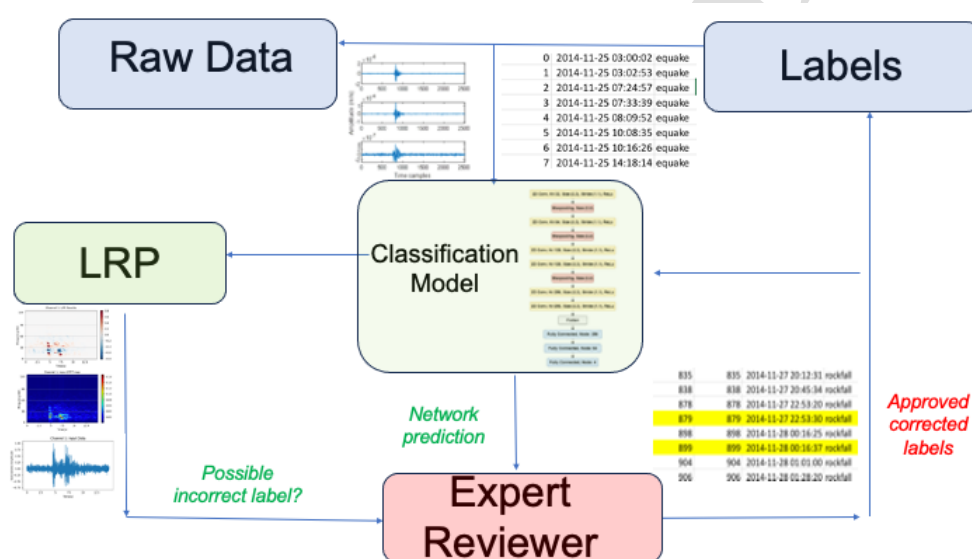
Received date: 20 August 2024  
Revised date: 19 November 2024  
Accepted date: 25 December 2024

Please cite this article as: J. Jiang, D. Murray, V. Stankovic et al., A human-on-the-loop approach for labelling seismic recordings from landslide site via a multi-class deep-learning based classification model. *Science of Remote Sensing* (2025), doi: <https://doi.org/10.1016/j.srs.2024.100189>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024 Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 1 Graphical Abstract

2 **A human-on-the-loop approach for labelling seismic recordings from**  
3 **landslide site via a multi-class deep-learning based classification**  
4 **model**5 Jiaxin Jiang, David Murray, Vladimir Stankovic, Lina Stankovic, Clement  
6 Hibert, Stella Pytharouli, Jean-Philippe Malet

7 Highlights

8 **A human-on-the-loop approach for labelling seismic recordings from**  
9 **landslide site via a multi-class deep-learning based classification**  
10 **model**

11 Jiaxin Jiang, David Murray, Vladimir Stankovic, Lina Stankovic, Clement  
12 Hibert, Stella Pytharouli, Jean-Philippe Malet

- 13 • Robust multi-class CNN-based seismic signal classifier
- 14 • LRP explainability maps for model diagnosis
- 15 • Trustworthy AI with geoscientist in the design loop

16 A human-on-the-loop approach for labelling seismic  
17 recordings from landslide site via a multi-class  
18 deep-learning based classification model

19 Jiaxin Jiang<sup>a</sup>, David Murray<sup>a</sup>, Vladimir Stankovic<sup>a</sup>, Lina Stankovic<sup>a</sup>,  
20 Clement Hibert<sup>b</sup>, Stella Pytharouli<sup>c</sup>, Jean-Philippe Malet<sup>c</sup>

<sup>a</sup>*Dept. Electronic and Electrical Engineering, University of Strathclyde, Glasgow, UK*

<sup>b</sup>*Institut Terre & Environnement de Strasbourg, University of  
Strasbourg, Strasbourg, France*

<sup>c</sup>*Dept. Civil and Environmental Engineering, University of Strathclyde, Glasgow, UK*

---

21 **Abstract**

22 With the increased frequency and intensity of landslides in recent years,  
23 there is growing research on timely detection of the underlying subsurface  
24 processes that contribute to these hazards. Recent advances in machine  
25 learning have introduced algorithms for classifying seismic events associated  
26 with landslides, such as earthquakes, rockfalls, and smaller quakes. How-  
27 ever, the opaque, “black box” nature of deep learning algorithms has raised  
28 concerns of reliability and interpretability by Earth scientists and end-users,  
29 hesitant to adopt these models. Leveraging on recent recommendations on  
30 embedding humans in the Artificial Intelligence (AI) decision making process,  
31 particularly training and validation, we propose a methodology that incor-  
32 porates data labelling, verification, and re-labelling through a multi-class  
33 convolutional neural network (CNN) supported by Explainable Artificial In-  
34 telligence (XAI) tools, specifically, Layer-wise Relevance Propagation (LRP).  
35 To ensure reproducibility, a catalogue of training events is provided as sup-  
36plementary material. Evaluation from the French Seismologic and Geodetic  
37 Network (R sif) dataset, gathered in the Alps in France, demonstrate the  
38 effectiveness of the proposed methodology, achieving a recall/sensitivity of  
39 97.3% for rockfalls and 68.4% for quakes.

40 *Keywords:* seismic signal analysis, microseismic signal classification, deep  
41 learning, explainable artificial intelligence, data annotation, model training  
42 *PACS:* 0000, 1111  
43 *2000 MSC:* 0000, 1111

---

## 44 **1. Introduction**

45 Seismic signal analysis is based on collecting, processing and performing  
46 inference on seismic signals with the goal of detecting, understanding, clas-  
47 sifying and locating seismic events, including not only earthquakes, but also  
48 rockfalls and smaller quakes or tremors that characterise landslides and their  
49 severity. The devastating effects of landslides on humans and infrastructure  
50 have been making headlines, and more recently have been often attributed  
51 to extreme weather and/or human activities. Seismometers provide accurate  
52 recordings of mechanical waves originating from various sources, but due to  
53 their high sensitivity, distinguishing between mechanical waves originating  
54 from tectonic activities and any other signals contained in the recordings  
55 (e.g., rainfall, animals, traffic, natural noise, machinery, etc.) is not an easy  
56 task. Manually identifying events based on recordings of seismometers is a  
57 time-consuming and subjective task, prone to errors and bias. Thus, manual  
58 detection has gradually been replaced by methods that automatically detect  
59 and classify seismic events. With higher availability in seismic recordings and  
60 advances in AI, seismic signal analysis has become a very much data-driven  
61 field and has spread well beyond seismology and geoscience, as it is now of  
62 interest to much broader research communities [1].

63 Deep learning has been shown to be achieve excellent detection and clas-  
64 sification performance for a range of applications where sufficient amount of  
65 labelled data is available, including automated road extraction [2], pneumo-  
66 nia diagnosis from medical imaging [3], satellite image analysis [4], [5], and  
67 car detection [6]. Due to the availability of many well-maintained datasets,  
68 the number of deep learning approaches used in seismology has also sky-

69 rocketed in recent years (see Fig. 1 in [1]) using enormous amounts of data  
70 to train the models. Consequently, recent literature is dominated by deep  
71 learning techniques applied to diverse tasks such as seismic event labelling  
72 using Residual Neural Network (ResNet) [7], magnitude estimation using  
73 a network that combines CNN and Recurrent Neural Networks (RNN) [8],  
74 event localization using CNN architectures [9], multitask learning for classifi-  
75 cation with velocity models [10] and tackling seismic inversion problems with  
76 conditional Generative Adversarial Networks (GAN) [11]. A detailed review  
77 of deep learning architectures, specifically proposed, for event classification  
78 from seismic recordings can be found in [12].

79 For example, CNN-based model ‘DeepQuake’ [13] has demonstrated ro-  
80 bust performance for high-magnitude earthquakes, though it has limitations  
81 with microseismic events, as demonstrated in [12]. In [14], RockNet, taking  
82 both 3-channel time series window and a spectrogram of the vertical channel  
83 of the window as inputs, is proposed for classifying rockfalls and earthquakes.  
84 The deep learning models achieve state-of-the-art performance in detecting  
85 and classifying seismic signals avoiding cumbersome manual feature gener-  
86 ation, selection and extraction process, with their ability to automatically  
87 learn most discriminative features from raw recordings. However, this also  
88 means that these models are limited by the used training set, and may learn  
89 specifically spurious correlations with the prediction target [15], [16]. Fur-  
90 thermore, the fact that the feature engineering task is taken away from the  
91 designer, makes deep learning models opaque, and hence often referred to as  
92 “black box”, which limits their use. Indeed, geoscientists are still reluctant  
93 to use them and rather rely on less complex interpretable methods based  
94 on hand-crafted features [17] that ensure that relevant physical features are  
95 used for detection and classification (see, e.g., Table I in [17] and Table A1  
96 in [18]).

97 Explainable artificial intelligence (XAI) [19], [20], is a research direction  
98 that provides human-interpretable explanations that can potentially enhance

99 training process, correct manual data annotation, improve models, and con-  
100 tribute towards building trust in AI-generated outputs [21], [22]. XAI tools  
101 have been extensively used in computer vision (e.g., [23]) and time-series sig-  
102 nal analysis problems (e.g., [24]); however, the work on explaining the output  
103 of deep learning models for seismic signal analysis, and using these explana-  
104 tions to improve confidence in data labelling, model training and building  
105 trust in inferred outputs, is still in its infancy.

106 In order to pave the way towards a regulatory framework for ensuring  
107 trust in AI, the European Commission has published seven principles of  
108 Trustworthy AI [25], which include Human Agency and Oversight, Technical  
109 Robustness and Safety, Privacy and Data Governance, Transparency, Di-  
110 versity, Non-discrimination and Fairness, Societal and Environmental Well-  
111 Being and Accountability.

112 Depending on how the AI-based seismic analysis will be used, from un-  
113 derstanding the subsurface processes and mechanics to hazard and disaster  
114 management, the AI systems can be seen as minimal risk to high risk, and  
115 therefore subject to strict oversight before they can be used to ensure infras-  
116 tructure and human safety. Therefore, the following principles are important  
117 for seismic analysis. First, AI systems should empower decision makers when  
118 it comes to hazard assessment or infrastructure planning, allowing them to  
119 make informed decisions from the AI system outputs. The principle of Hu-  
120 man Agency and Oversight caters for proper oversight mechanisms that need  
121 to be ensured, which can be achieved through human-on-the-loop and human-  
122 in-command approaches. Second, the principle of technical robustness and  
123 safety, in part states that AI systems need to be accurate, reliable and re-  
124 producible to ensure unintentional harm can be minimised and prevented.  
125 Accuracy refers to the ability to correct predictions based on AI models and  
126 can be implemented via rigorous evaluation and indication of likelihood of po-  
127 tential errors. Reproducibility describes whether an AI experiment exhibits  
128 the same behaviour when repeated under the same conditions. A reliable AI



129 system is one that works properly with a range of inputs and in a range of  
130 situations. Third, the principle of privacy and data governance enables users  
131 to trust the data gathering process and that it does not contain inaccuracies,  
132 errors or mistakes, especially with respect to labelling or cataloguing by ex-  
133 pert geoscientists. Fourth, the principle of transparency states that the data  
134 and AI system should be transparent through traceability mechanisms in the  
135 form of documentation of datasets and processes that yielded in decision, in-  
136 cluding data gathering, data labelling and algorithms used. Furthermore,  
137 transparency also includes explainability, that is, AI systems and their deci-  
138 sions should be explained in a manner adapted to the stakeholder concerned.  
139 This includes XAI. Fifth, transparency also states that humans need to be  
140 aware that they are interacting with an AI system, and must be informed of  
141 the system's capabilities and limitations. Finally, the social and environmen-  
142 tal well-being principle state that the AI systems should be sustainable and  
143 environmentally friendly - this can be through taking into considering the  
144 resource usage and energy consumption for training the models. Moreover,  
145 they should consider the societal impact. Monitoring, understanding, mod-  
146 elling and predicting landslide processes due to climate change, especially  
147 rainfall, tackle United Nations (UN) Sustainable Development Goal (SDG)  
148 13 on Climate Action [26]. As explained in [27], shearing and friction be-  
149 tween the soil grains results in release of seismic energy within the landslide  
150 body. Therefore, passive seismic monitoring is a good approach to monitor  
151 and mitigate slope instabilities, as it provides high temporal resolution data  
152 in near real time that relate to the dynamics of the landslide. This means  
153 that the transition (and rapid transformation) of the landslide from slow  
154 rate sliding into a rapid slope failure may be detected and therefore mitigate  
155 associated hazards.

## 156 2. Literature review on Trustworthy AI for Seismic Signal Analysis

157 To ensure trust and expert's control of the decision process, machine  
158 learning-based seismic signal analysis has been performed either in a semi-  
159 automated manner [28] using continuous expert oversight and monitoring  
160 (human-on-the-loop), using interpretable models [17], or using non-interpretable  
161 models (such as Random Forests) but with numerous hand-crafted features  
162 [29] to ensure that the inference is made on signal characteristics identified  
163 by experts as important. In [18] a detailed study of feature importance is  
164 presented where 119 features are constructed based on seismic signal liter-  
165 ature and their importance tested using four different feature importance  
166 methods and different classifiers based on Support Vector Machine, Random  
167 Forest, and three graph signal processing based semi-supervised approaches.  
168 The features are experimentally ranked showing time-, frequency-, cepstrum  
169 and polarity features that are of highest importance in inference making per  
170 studied class. The results show that out of 119 constructed features only  
171 a subset contributed significantly to the decision. Note that this study was  
172 based on quantifying the importance of hand-crafted features in accurately  
173 classifying multiple event classes from continuous data, thus deep learning  
174 networks were not considered.

175 In [13], convolutional neural networks (CNNs) are used to classify isolated  
176 catalogued seismic events into noise, earthquake and other events. The au-  
177 thors developed a heatmap-based visualisation tool to explain model outputs  
178 via the outputs of activation functions of each filter in the convolutional lay-  
179 ers and then overlapping the result with the raw input signal. However, this  
180 study has several weaknesses when it comes to gaining trust in model out-  
181 puts. Firstly, it is not clear how explanations are formed by fusing outputs  
182 of the activation functions from different layers. Secondly, only binary clas-  
183 sification is considered, i.e., identifying relatively well-defined earthquakes  
184 from other signals. Thirdly, the approach does not exploit advanced XAI  
185 methods, and it is not used to explain any false predictions.

186 In [30], the authors proposed a Dual-Channel CNN Module where one  
187 channel contains raw time-domain waveforms, and the other channel con-  
188 tains frequency-domain information by Discrete Cosine Transform (DCT) to  
189 classify input seismic waveform into rock fracturing and noise, together with  
190 an explanation module, EUG-CAM (Elaborate Upsampling-based Gradient-  
191 weighted Class Activation Mapping). It builds upon the principles of the  
192 gradient weighted class activation mapping (GradCAM) [31], harnessing the  
193 influence of feature map values and gradients to elucidate the importance of  
194 diverse features in the last convolutional layer. Recognizing the discrepancy  
195 between feature map sizes and input data dimensions, EUG-CAM uses a  
196 strategic amalgamation of transposed convolution, unpooling, and interpo-  
197 lation, to generate feature mappings from a coarse localization map. This  
198 results in an explanation feature map that effectively encapsulates class acti-  
199 vation, learning insights, and network architecture considerations. However,  
200 the model's limitation is in classifying only two classes (rock fracturing vs.  
201 noise) and its confinement to binary classification. Furthermore, the reliance  
202 on a 1-D CNN model facilitates explanations primarily within the time do-  
203 main, possibly neglecting the benefits of frequency-domain insights garnered  
204 from the DCT. Additionally, the visualization maps cannot show the ad-  
205 verse input signal influence (negative contribution) on classification results,  
206 hampering a comprehensive and well-rounded comprehension of the model's  
207 decision-making process.

208 In [12], the authors present CNN models with six channel inputs for  
209 multi-class classification of earthquakes, quakes, rockfalls and noise and use  
210 visualization of feature maps to understand the network's internal work-  
211 ings. The authors examine feature maps at various convolutional layers and  
212 the second fully connected (FC) layer, gaining insights into feature extrac-  
213 tion. Different models, including time-series, Short-time Fourier Transform  
214 (STFT) and Continuous Wavelet Transform (CWT)-based designs, highlight  
215 the network's focus on time, frequency, and wavelet characteristics. The main

216 observation is that early layers locate event positions and extract basic fea-  
217 tures, while deeper layers refine these features into abstract representations  
218 for classification. The second FC layer’s feature distributions vary across  
219 seismic events, indicating the network’s capability to distinguish three event  
220 types from noise based on learned features. In addition, Layer-wise Rele-  
221 vance Propagation (LRP) showed promising results in identifying the most  
222 relevant features for each class, further enhancing the interpretability of the  
223 model [32].

### 224 *2.1. Contributions*

225 The goal of this paper is to provide comprehensive explanations to iden-  
226 tify key features learnt by a deep neural network for multi-class classification,  
227 demonstrate that these features are in agreement with the physical properties  
228 of seismic signals and common hand-crafted features used in the literature  
229 [17]. The generated explanations are then used to explain instances of mis-  
230 classifications and correct errors in manual labelling, jointly with a geoscient-  
231 ist, who verified the corrected labels of the classified events and the features  
232 associated with these events. This builds trust in the models confirming that  
233 the learnt feature representations agree with expert knowledge.

234 We use state-of-the-art XAI tools to explain deep learning models for  
235 detection and classification of micro-seismic signals and show how these  
236 explanations can be used to improve the designs and explain correct and  
237 wrong predictions. In particular, we use a CNN-based architecture with a  
238 frequency-domain input, for detection and classification of seismic signals  
239 into four classes: earthquake, micro-earthquake referred to as quake, rock-  
240 fall and noise. These are the same classes as used in [12] and [29]. There  
241 are three inputs to the CNN, each comprising continuous recordings from  
242 the channels of a typical three-component seismometer, usually deployed for  
243 seismic monitoring.

244 Our models are trained and tested on a publicly accessible dataset Résif  
245 [33] that has over 1000 labelled events, including earthquakes, quakes, rock-

246 falls and anthropological noise. After classification, we use Layer-wise Rel-  
247 evance Propagation (LRP) [34] to explain the decision making process. We  
248 analyse the basis of the model for event classification and communicate the  
249 reasons for misclassification of individual events. Furthermore, if the pre-  
250 dicted class is different to the expert label, and after inspection of the fil-  
251 tered signal, its STFT and LRP map, the event is sent back to the expert  
252 for re-labelling. This protocol is used to correct possible labelling mistakes  
253 in the large annotated seismic dataset.

254 In summary, our main contributions are:

- 255 1. ensuring data integrity by leveraging on a well-maintained ongoing seis-  
256 mological data portal releasing checked seismic recordings publicly, as  
257 well as cataloguing/labelling by expert geoscientists – this aspect is by  
258 nature transdisciplinary
- 259 2. traceability to enable transparency by leveraging on public datasets,  
260 where data gathering, labelling and performance with different algo-  
261 rithms are well documented
- 262 3. an additional catalogue of 829 labelled events for a period of 3 days,  
263 classified by the CNN, verified by an expert and labels corrected -  
264 provided as supplementary material
- 265 4. reproducibility by releasing the catalogue of 822 manually selected high  
266 quality training events as supplementary material
- 267 5. designing a multi-classifier robust to noisy continuous recordings for  
268 high performance but also indicating likelihood of potential errors
- 269 6. reliability of design by ensuring that the multiclassifier design works for  
270 a continuous input stream with noisy signals and low signal to noise  
271 ratio events
- 272 7. explainability for transparency by providing explanations of the multi-  
273 classifier outputs via XAI LRP maps
- 274 8. communication for transparency by clearly identifying the level of per-  
275 formance and limitations

276 9. tackling the UN SDG 13 by accurately detecting landslide related  
277 events that helps build trust in precursors to landslides such as rockfalls  
278 and quakes

279 The first three contributions are presented in Section 3, where we describe  
280 the dataset used and data pre-processing. Contributions (4)-(5) are covered  
281 in Section 4, where the proposed CNN-based architecture, the sliding-window  
282 continuous detection method, the proposed post-processing and explainabil-  
283 ity tools used are described. Section 5 demonstrates our contributions (6)-(8).  
284 The significance of this work, i.e., contribution (9) was discussed above and is  
285 demonstrated in Section 5. Finally we conclude in Section 6 with suggestions  
286 for further work.

### 287 3. Dataset

288 In this section we provide details about the data management, including  
289 collection, storage, release and labelling /cataloguing, describing the first  
290 three contributions of this paper.

#### 291 3.1. Data gathering and context

292 Our raw seismometer recordings are obtained from the publicly accessible  
293 Résif Seismological Data Portal, acquired by the French Landslide Observa-  
294 tory OMIV (Observatoire Multi-disciplinaire des Instabilités de Versants).  
295 In particular, we focus on the Super-Sauze (SZ) slow moving landslide mon-  
296 itoring array, acquired by the Super-Sauze C (SZC) station of the French  
297 Landslide Observatory on the Permanent seismological records on unstable  
298 slopes which are installed at the centre of the Super-Sauze landslide deposit  
299 in Southeast France (Latitude: 44.34787°N, Longitude: 6.67805°E). The lo-  
300 cation of the SZC station is shown on the map in Figure 1. The seismometer  
301 array consist of one central three-component sensor and three vertical one-  
302 component sensors (organized as equilateral triangle), all recording at 250Hz

303 sampling rate. In this project, we used data from the three-component sen-  
 304 sor. This choice aligns with common practices in seismic waveform classi-  
 305 fication, where a 3-channel input is standard, such as EQ-transformer [35]  
 306 and DeepQuake [13]. Additionally, it facilitates transfer learning, as many  
 307 seismometers employ three-component sensors, ensuring compatibility with  
 308 various seismic datasets and applications. Using 3 channels also reduces the  
 309 number of false positives which can occur with arrival mismatches and re-  
 310 duces the computational demand. The seismometers recorded three periods:  
 311 from 11 Oct. to 19 Nov. 2013; from 10 Nov. to 30 Nov. 2014; and from 9  
 312 June to 15 Aug. 2015.

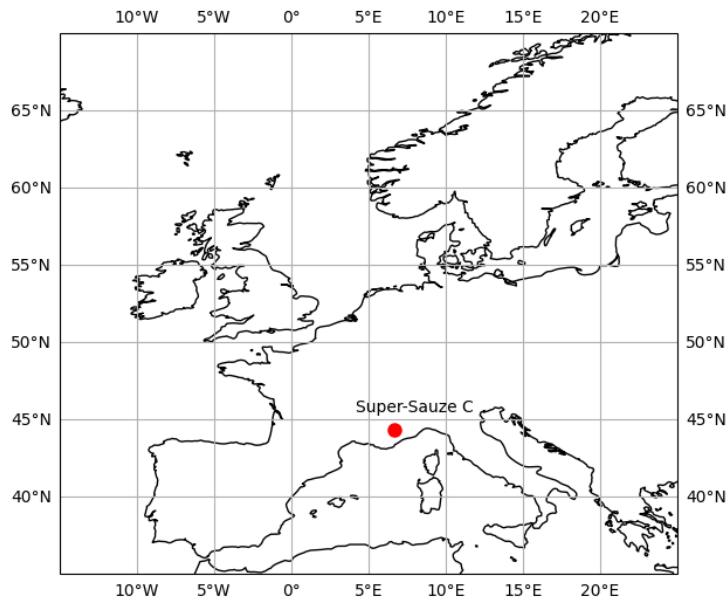


Figure 1: Map showing the location of the Super-Sauze C (SZC) station.

313 The description of the SZ slope deformation, together with the challenges  
 314 of detecting the microseismic events is well documented in [36]. Additionally,  
 315 description of how the catalogue of events was generated is documented in  
 316 [29], where events were detected by the STA/LTA algorithm applied in the  
 317 frequency domain, and manually labelled into four classes: earthquake, quake

318 (micro-earthquake events), rockfall and natural/anthropogenic (N/A) noise.  
319 All events except noise are classed as microseismic according to [27].

320 Rockfalls mainly occur at the main scarp of the landslide, where the rigid  
321 block falls from the steep slope (height > 100m). The quake is likely to  
322 be triggered by material damage, surface cracks and openings within the  
323 landslide main flow. The earthquakes class includes regional seismic events  
324 in this area and the teleseisms (global large magnitude earthquakes). N/A  
325 noise events include all anthropogenic and environmental noise, due to, e.g.,  
326 transportation, pedestrian walking, heavy rain, animals, strong wind, etc. It  
327 does not include noise in the form of instrumentation error.

### 328 *3.2. Labelling*

329 The SZ recordings over the data gathering duration described in the pre-  
330 vious subsection were labelled as described in [29], using STA/LTA in the  
331 frequency domain to pick events, and manual labelling of these events by an  
332 expert based on their amplitude, duration, spectrogram and location. The  
333 number of labels in this catalogue, which will be referred to as the origi-  
334 nal catalogue, for each class, is reported in [18] and [12], where the events  
335 were classified on continuous recordings with classifiers using manual feature  
336 generation, and deep-learning-based classifiers with automated feature ex-  
337 traction, respectively. Since detection and classification were performed on  
338 the continuous data stream, the Normalised Graph Laplacian Regularisation  
339 (normGLR)-based [18] and CNN-based [12] classifiers also reported classifi-  
340 cation of hundreds of additional non-catalogued events, with a high density  
341 of events in the period 25th to 28th Nov. 2014, which coincided with a period  
342 of high activity on the SZ slope [37].

343 As reported in [18], all four types of events are present in this 4-day  
344 time period, and in addition to the 120 events (65 rockfalls, 18 quakes, 23  
345 earthquakes and 14 noise) labelled in the original catalogue, 17 quakes, 89  
346 earthquakes and 92 rockfalls events were detected and classified by the nor-  
347 mGLR classifier whereas an additional 260 quakes, 174 earthquakes and 32



348 rockfalls were detected and classified with the CNN approach of [12]. These  
349 algorithms only missed 1 earthquake, 1 rockfall and 2 noise events that were  
350 present in the original catalogue.

351 All events detected by the normGLR classifier, the CNN classifier and  
352 an additional classifier based on Siamese networks [38] were reviewed by  
353 an expert for labelling following the methodology used to build the original  
354 catalogue, which is based on the seismic signal waveform and spectrogram  
355 features. The final outcome of the expert reviews for this 4-day period were 69  
356 quakes, 29 earthquakes and 126 rockfalls. Note that the normGLR classifier  
357 was too sensitive, overestimating the number of earthquakes[18]. The CNN-  
358 based 6-channel input multi-classifier of [12] was too sensitive for quakes and  
359 earthquakes but missed a number of rockfalls.

360 This exercise demonstrated the value of machine learning-based classifi-  
361 cation on continuous streaming recordings, since it is tedious for experts to  
362 manually review continuous data streams, as well as pick up the microseismic  
363 events, especially quakes and rockfalls, that are often “hidden” or “unclear”  
364 within ambient noise present in the recordings. These newly detected and  
365 expert-labelled events during the period 25th to 28th Nov. 2014, not present  
366 in the original catalogue, are released with this paper and are focus of this  
367 study.

#### 368 4. Methodology

369 In this section, we describe our methodology. First, building on our prior  
370 work [12], we propose an improved multi-class CNN-based classifier that uti-  
371 lizes 3-channel inputs and a modified training strategy (see Section 4.3) to  
372 enhance precision in detecting quakes and earthquakes, as well as improve  
373 recall/sensitivity rates for rockfalls. Second, we analyse the outputs of the  
374 improved multi-classifier, as part of our human-on-the-loop contribution to  
375 verify instances of labelling error, likely to occur for large volumes of continu-  
376 ous streaming seismic recordings. This is carried out via the XAI-based LRP

377 tool to visualise the features of misclassifications, which are then queried for  
378 re-evaluation by the expert.

#### 379 *4.1. Proposed CNN-based architecture*

380 An STFT-based CNN model, inspired by VGGNet [39] and adapted from  
381 [12], is used. The influence of seismometer characteristics such as sensitivity,  
382 frequency band, and axis configuration on the reliability and effectiveness of  
383 our results was explored in [12], whereby good transferability was demon-  
384 strated with recordings from different seismometers with varying sensitivity  
385 levels and sampling rates, and geographic location. Additionally, we exam-  
386 ined the performance impact of different seismometer configurations, compar-  
387 ing one-axis (single-channel) seismometers with multi-channel inputs during  
388 training. We use STFT maps as inputs for the CNN model, as these inputs  
389 were shown to provide better results on average compared to directly feeding  
390 time-series signals. Additionally, the generalisability and robustness of this  
391 architecture across different sites has been demonstrated in prior work [12].  
392 Particularly, as evidenced by the recent trend in CNN-based architectures for  
393 analysis of seismic recordings, such networks excel in extracting hierarchical  
394 and discriminative features from complex data, making them highly effective  
395 for seismic event classification. The value of binary vs multi-class networks in  
396 terms of how multi-class models are able to achieve similar performance while  
397 requiring less models to be trained and run, and hence lower overall com-  
398 plexity, was demonstrated in [12]. Multi-class CNN models offer enhanced  
399 feature extraction, adaptability to various data patterns that are often indis-  
400 tinguishable (such as local quakes and rockfalls), and improved classification  
401 performance compared to state-of-the-art DL approaches for seismic analysis,  
402 discussed in Introduction Section, that mostly focus on binary classification.

403 The architecture of the model is composed of convolutional layers, max  
404 pooling layers and FC layers, adapted to the input shapes and output cate-  
405 gories, as shown in Figure 2. Convolutional layers perform feature represen-  
406 tation and extraction, followed by max-pooling layers that downsample the

407 extracted feature into a feature map with smaller size.

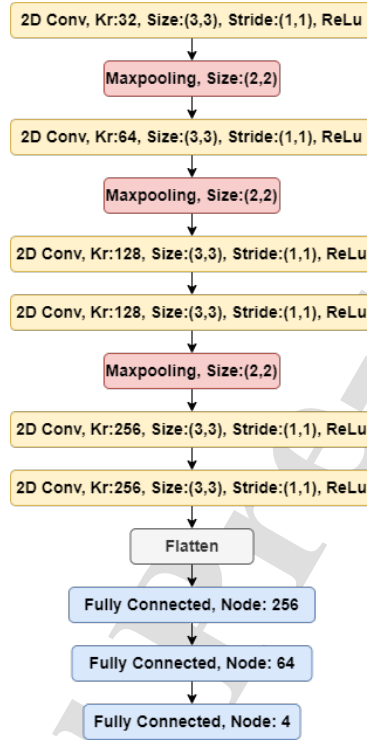


Figure 2: STFT-based CNN for seismic classification. Kr denotes the number of kernels, and ‘Flatten’ function transforms the input data into a 1D array.

408 Compared to [12], to effectively process long-duration seismic events within  
 409 continuous data streams, we increase the input window of the CNN model to  
 410 15 seconds (from 10 seconds). We also reduce convolution kernels and neural  
 411 nodes in each layer, achieving a balance between model complexity and per-  
 412 formance. Moreover, recognising the prevalence of waveforms captured by  
 413 three-component sensors, the input to the network is 3-channel input data,  
 414 in contrast to 6-channel used in [12], which significantly expands the model’s  
 415 applicability across a wider range of scenarios.

#### 416 4.2. Sliding window-based detection

417 Raw signals recorded by 3-channel (North, East and vertical direction)  
 418 seismic recorders are used. Since the classes of interest are 5-60Hz bandwidth,  
 419 we first use a band-pass filter to remove low frequency noise (denoising) as  
 420 in [12]. To allow prediction on a continuous stream of signals, a sliding win-  
 421 dow method is used to segment the continuous stream into smaller windows  
 422 as in [40], [41]. The window size and overlap are selected based on the tem-  
 423 poral resolution required for the events of interest. A window size of 3750  
 424 samples (i.e., 15 seconds) is used. The overlap between consecutive windows  
 425 is set to 93% of window size (3500 samples (14 seconds)), which corresponds  
 426 to a shift by 1 sec, allowing the CNN model to capture the temporal dynam-  
 427 ics of the events of interest. For each window, the CNN model is used to  
 428 predict the probabilities of each class being present.

429 Furthermore, since the peak amplitude of signals belonging to different  
 430 classes is large, to improve the learning ability of the models, we perform  
 431 normalization of the filtered recordings. In particular, in order to enable the  
 432 model to focus on classifying the input signals and facilitate the subsequent  
 433 explanation of the classification results, we normalise each 15-second window  
 434 by subtracting mean and dividing by the maximum of the absolute value of  
 435 each input window.

436 For the STFT map input, in order to get good time and frequency resolu-  
 437 tion, ‘Boxcar’ window with length of 128 samples with 70% overlap is used.  
 438 We perform STFT on denoised and normalized time series input window.  
 439 Thus, the input shape for the STFT-based model is  $65 \times 95 \times 3$  samples.

#### 440 4.3. Training and testing

441 The inputs to the model for both training and testing comprise STFT  
 442 maps generated from the raw recordings as discussed in the previous sub-  
 443 section. Our prior work in [12] demonstrate that CNN models tend to be  
 444 overly sensitive. To address this, we refine the sensitivity of our CNN by

445 only using the high-quality events to train the model. Specifically, we visu-  
446 ally inspected and chose events from the original catalogue to ensure that the  
447 set used for training comprised only high-quality events based on signal clar-  
448 ity and high-SNR (Signal-to-Noise Ratio) for earthquake, quake and rockfall  
449 classes. All noise events originate from the original catalogue. In addition  
450 to the manually selected events, we utilise the labelled events from the 25th  
451 November 2014 (one day) to train the model further. These additional data  
452 allows us to augment the training set with events that are not included in the  
453 high-quality subset of the original catalogue and help to improve precision  
454 and recall.

455 The list of all the high-quality events from the original catalogue as well  
456 as the events from the 25th November 2014 used for training can be found  
457 as supplementary material for the purposes of reproducibility, as the second  
458 principle of Trustworthy AI. During testing phase, we use STFT maps from  
459 26th to 28th Nov. 2014, which are not included in the training set. These  
460 labelled events are released with this paper as supplementary material.

#### 461 *4.4. Post-processing*

462 While the sliding window technique enables continuous detection, it can  
463 introduce certain challenges. One of the main issues is that it may break  
464 the continuity of the event waveform, leading to potential inconsistencies or  
465 artefacts in the classification results. This occurs because the sliding window  
466 segments are treated independently, without considering the temporal con-  
467 text or smooth transitions between adjacent windows. To address this prob-  
468 lem, post-processing techniques are often employed to refine and enhance the  
469 detection output by taking into account the temporal relationships between  
470 adjacent windows.

471 The proposed post-processing system is based on threshold filtering, me-  
472 dian filtering, and Gaussian kernel filtering of the softmax output of the  
473 CNN. In addition, a peak selection method is applied to resolve cases where  
474 two classes of events have very similar detection results. (1) Threshold fil-

475 tering: the softmax output of the CNN is filtered with a threshold value (set  
 476 to 0.5), and all values below this threshold are set to zero. This is done to  
 477 remove low-probability detections. (2) Median filtering: After the threshold  
 478 filtering step, the probability distribution may contain isolated spikes. To  
 479 remove these isolated spikes, we apply a median filter to each class sepa-  
 480 rately. In addition to removing isolated spikes, the median filter can also  
 481 merge spikes that are very close together, resulting in smoother and more  
 482 continuous probability distributions. We set the size of the median filter to  
 483 5. (3) Gaussian kernel filtering: a Gaussian kernel filter is applied to the me-  
 484 dian filtered output to smooth the probability distribution. Gaussian kernel  
 485 is defined with a sum of 1 and a length of 15. Its standard deviation is 5.  
 486 (4) Peak selection: after using Gaussian kernel filtering, we select the high-  
 487 est peak (i.e., the longest duration) as the final output. This peak selection  
 488 method allows us to choose the class of the event with the longest duration,  
 489 as it indicates a higher confidence level in the classification result.

#### 490 4.5. Explainability-informed re-labelling

491 Unlike classifiers such as RF, SVM and (norm)GLR-based classifiers that  
 492 take hand-crafted features as inputs and where feature importance was stud-  
 493 ied in detail in [18], the CNN multi-classifier is essentially a “black box”  
 494 since we do not know what features were deemed important. We therefore  
 495 utilise LRP to understand feature importance for the deep-learning CNN  
 496 multi-classifier.

497 LRP [34] is a state-of-the-art XAI method, that shows the contribution  
 498 of each sample in the input data to the classification results and can be  
 499 implemented in the pre-trained model [42]. In this paper, LRP is used to  
 500 help identify which parts of the seismic signal are most important in making  
 501 the final classification decision. This helps understanding which features of  
 502 the seismic signal are most relevant for seismic detection, and identify any  
 503 potential biases in the model. In addition, LRP can provide interpretable  
 504 and detailed explanations of the model’s decision-making process, which can

505 be useful for communicating the model’s results to human experts.

506 The LRP method starts from the output of the model, sets the output  
 507 value before activation function as relevance, and gradually back propagates  
 508 relevance, iteratively, layer by layer, to the input nodes. In the backpropa-  
 509 gation, relevance follows the conservation law, that is, a neuron’s relevance  
 510 equal to the sum of relevance as it flows out toward all other neurons. Vari-  
 511 ous propagation rules have been proposed, such as LRP- $\gamma$ , LRP- $\epsilon$  and LRP-0  
 512 rule [22]. In this paper, we used LRP- $\epsilon$  rule which is suitable for convolutional  
 513 layers and max pooling layers [43], and is defined as:

$$R_j = \sum_k \frac{a_j w_{jk}}{\epsilon + \sum_{0,j} a_j w_{jk}} R_k, \quad (1)$$

514 where  $R_j$  represents an LRP relevance score assigned to neuron  $j$ ,  $a_j$  denotes  
 515 an input activation,  $w_{jk}$  is the weight connecting neuron  $j$  to neuron  $k$  in the  
 516 layer above,  $\sum_{0,j}$  denotes that we sum over all neurons  $j$  in the lower layer  
 517 plus a bias term  $w_{0k}$  with  $a_0 = 1$ .  $\epsilon$  is a regularisation term, i.e., a small  
 518 value that prevents the denominator from being 0.

519 We generate LRP maps for all events whose CNN-based predicted class  
 520 does not correspond to the event class label as provided by the expert via  
 521 the procedure described in Subsection 3.2 (i.e., misclassification). Then, we  
 522 ask the same expert to review the recording, this time together with the  
 523 LRP feature importance map, to ensure trust in the labels. The “corrected”  
 524 labels (those that the expert agrees were originally wrongly labelled) are  
 525 then marked and released as part of the supplementary material together  
 526 with their STFT and LRP maps. The whole process is shown in Figure 3.

## 527 5. Results

528 In this section, we first demonstrate our Contribution (5 & 6), by re-  
 529 porting the performance of the proposed models on the test dataset using  
 530 standard classification performance measures as in [12]. Then, we present

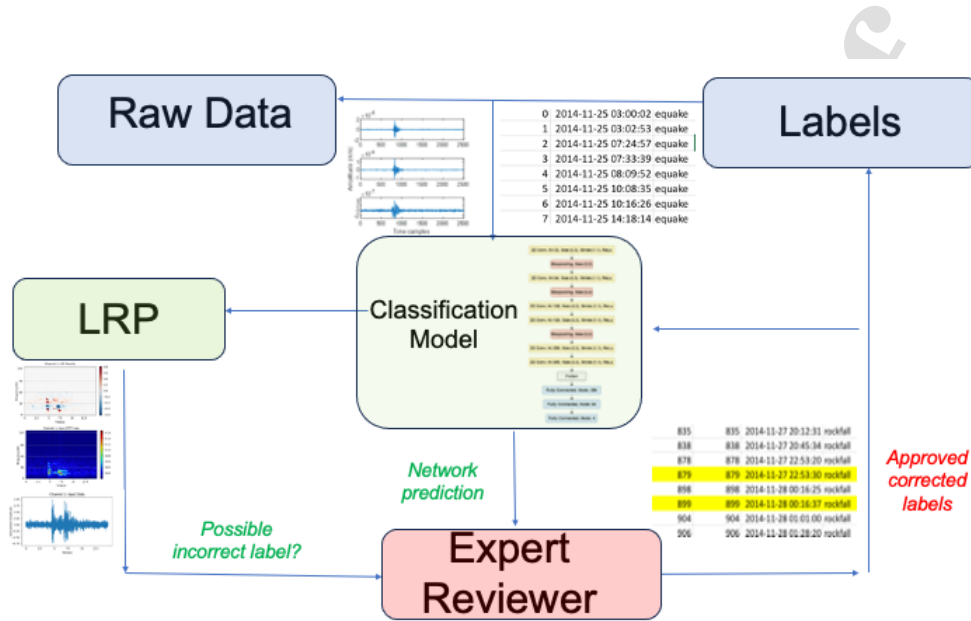


Figure 3: Flowchart of the proposed human-on-the-loop process.

531 our explainability results as per Contribution (7) and discuss main reasons  
 532 behind misclassification (Contribution (8)).

### 533 5.1. Analysis of classifier output

534 Our models are implemented in Keras framework. Since the activation  
 535 function of the output layer is softmax, we use categorical cross entropy as  
 536 loss function. The used optimiser is Adam with an initial learning rate of  
 537 0.0007. Adaptive learning rate adjustment is implemented, which reduces  
 538 the learning rate by a factor of 0.9 when loss improvements plateau for 5  
 539 epochs. Training is performed over 100 epochs with a batch size of 128. For  
 540 the second training session, utilizing the data from November 25, the model  
 541 is trained over a total of 50 epochs. To prevent the risk of overfitting due  
 542 to additional training, early stopping is implemented; that is, if the training  
 543 accuracy did not exhibit significant improvement within 5 consecutive epochs,  
 544 the training process is terminated early.

545 In the 3-day testing period (26th-28th Nov.), the expert labelled 46  
 546 quakes, 18 earthquakes, 74 rockfalls and 719 noise events. The confusion



547 matrix in Table 1 compares the output of the proposed CNN-based network,  
 548 with post-processing (Sec. 4.4), to the expert labels. As is common practice  
 549 for seismic signal classification on continuous data [29], the confusion matrix  
 550 also includes recall/sensitivity values in brackets. Recall is the ratio of true  
 551 positives to the sum of true positives and false negatives. In Section 3.2, it  
 552 is demonstrated that during the 4-day period from November 25th to 28th,  
 553 there are 6 additional earthquakes not labelled in the original catalogue [29].  
 554 The model discussed in [12] detected a much larger number, specifically 174  
 555 additional, earthquakes. This comparison shows the significant improvement  
 556 in the precision of earthquake classification achieved by our model. Addi-  
 557 tionally, our model achieved high recall (sensitivity) for rockfall events. As  
 558 expected, quake and noise events can be confused with the other 3 classes, due  
 559 to heterogeneity of the noise signal and very low signal amplitude of quake  
 560 signals. Next, we leverage on LRP to explain the origin of misclassifications.

Table 1: Confusion Matrix - Proposed CNN-based network with post-processing against expert labels (the numbers in brackets indicate recall/sensitivity rates).

		Model			
		Quake	Earthquake	Rockfall	Noise
Expert	Quake	<b>26</b> (56.5%)	2	9	9
	Earthquake	0	<b>15</b> (83.3%)	1	2
	Rockfall	2	0	<b>72</b> (97.2%)	0
	Noise	110	13	58	<b>538</b> (75.1%)

## 561 5.2. Explainability

562 The used package for embedding LRP into our models is iNNvestigate [44]  
 563 which supports Keras framework in Python 3. Default parameters of the  
 564 LRP- $\epsilon$  rule are used.

565 Figure 4(a) shows an example of a correctly classified earthquake event.  
 566 Positive and negative values of the LRP relevance represent positive and neg-  
 567 ative contributions to the classification results, of the corresponding STFT,

568 respectively. The distribution of LRP relevance is focused on the high fre-  
569 quencies (about 40 to 50Hz) when the P-wave is picked as well as the low  
570 frequencies (around 15 to 20Hz) of the P-wave and, after roughly 5sec, the  
571 low frequencies of the S-wave with intermediate noise shown in light blue  
572 correctly identified as not contributing (negative contribution). This exam-  
573 ple shows that the model learnt, and uses as basis for its predictions, that  
574 the P-waves of earthquake events tend to have both high and low frequen-  
575 cies (around 50Hz and 20Hz, respectively) and that high energy content of  
576 S-Waves follows in time.

577 Figure 4(b) shows an example of a correctly classified quake event. Quake  
578 events are of shorter duration than earthquakes, have lower amplitudes, and  
579 energy focused in low frequencies. LRP relevance is concentrated in the  
580 single peak (positive and negative) of the event waveform, suggesting that  
581 the normalised maximum amplitude is the key distinguishing feature. In the  
582 frequency domain, the LRP map clearly shows the importance of the peak  
583 that has energy mainly focused below 30Hz while there is also a small positive  
584 contribution between 30 to 40Hz.

585 Figure 4(c) shows an example of a correctly classified rockfall event.  
586 While the relevance score of quake events is concentrated on a single peak, rel-  
587 evance of rockfall events is concentrated on multiple peaks, which also shows  
588 an important property of rockfall events – multiple significant peaks. Look-  
589 ing at the LRP map, relevance has multiple focused points corresponding to  
590 multiple short waves – a characteristic of rockfalls. In addition, although  
591 both rockfall and quake events have a frequency band between 10 to 30Hz,  
592 LRP relevance is mostly concentrated at frequencies greater than 20Hz for  
593 rockfalls and below 20Hz for quakes.

594 Similar visualisation maps are produced for other correctly classified events.  
595 In summary, the model searches: (a) for P-wave and S-wave peaks and their  
596 corresponding frequency contributions to predict an earthquake; (b) a short  
597 wave with a single peak below 20Hz to decide quake; (c) multiple significant

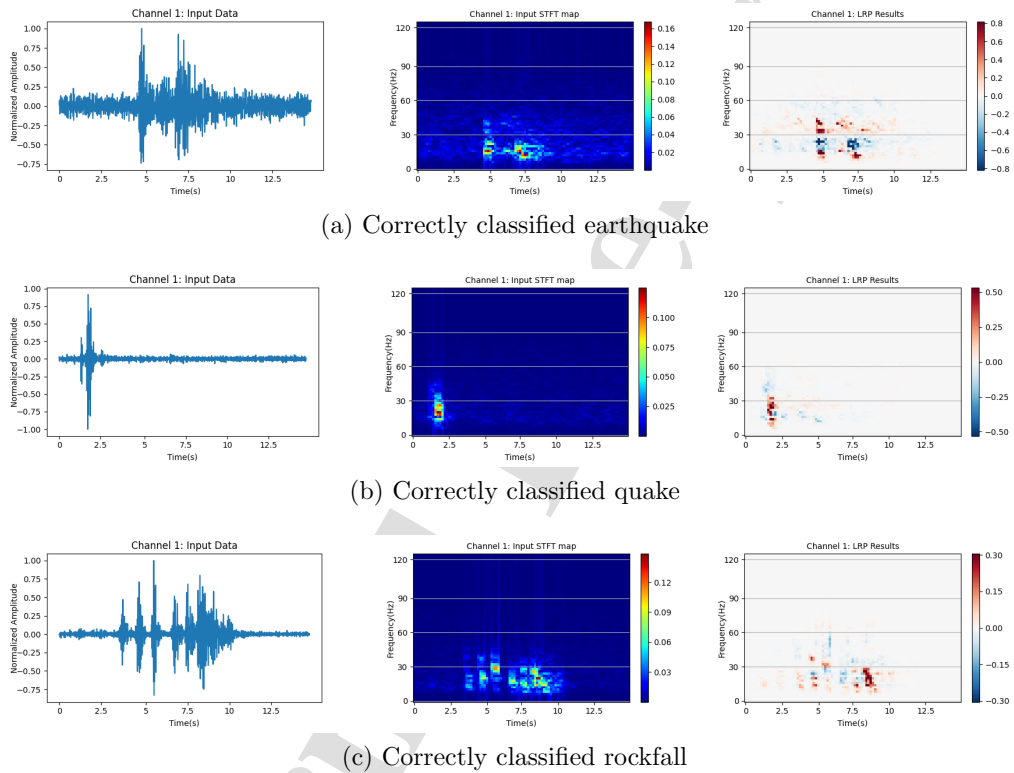


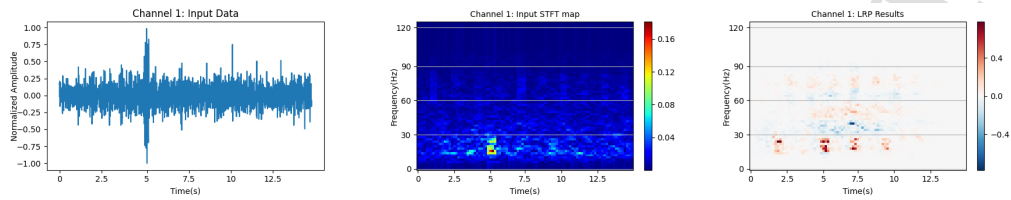
Figure 4: Correctly classified examples of earthquake, quake and rockfall: The first column shows the time-series signal, middle column the STFT, and the right column is the LRP relevance heatmap.

598 frequency components around 25Hz to decide that the target signal is rock-  
599 fall. This is in accordance to the characteristics of the three signal classes  
600 [29], [17], [12]. Next, we will analyse misclassified events to explain why they  
601 occur and how they can be avoided.

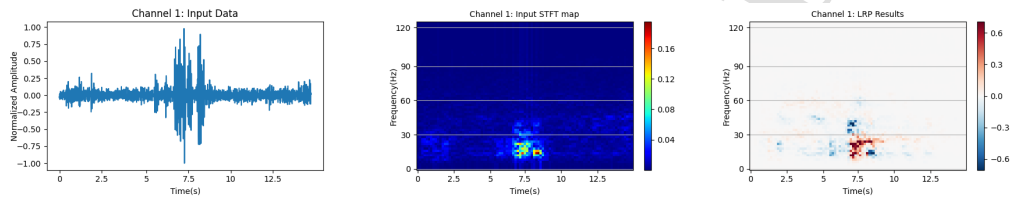
### 602 5.3. Explaining origin of misclassification

603 In this section, we show how LRP can be used for model diagnosis. The  
604 confusion matrix presented in Table 1 shows that the quake signals are some-  
605 times misclassified as rockfalls. Interestingly, however, rockfall signals are  
606 rarely misclassified as quakes (only 2 misclassified events). To investigate  
607 this further, Figure 5(a) shows an example of a quake event misclassified as  
608 rockfall. In the LRP map, relevance distribution is very scattered. That  
609 is, the LRP relevance is not focused on the quake event's peak, but instead  
610 picked up several consecutive peaks, where the positive relevance is correctly  
611 concentrated at 5 seconds. This indicates that the model correctly recog-  
612 nised a quake event's peak appearing around 5 seconds, but there was a high  
613 energy signal in nearby frequency bands, influencing the final prediction.  
614 On the other hand, there are many positive relevancies at different times  
615 that correspond to frequencies between 20Hz to 30Hz, which is akin to the  
616 learnt rockfall 'behaviour'. Thus, the main reason of misclassification be-  
617 tween quake and rockfall is that the signal-to-noise ratio of the quake event  
618 was very low, with a noise signal appearing immediately after, mimicking  
619 multiple peaks of rockfall events.

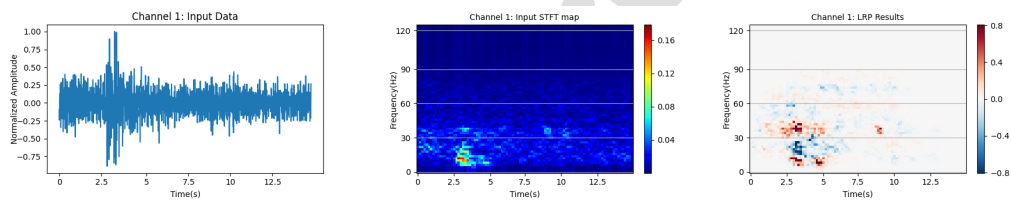
620 In Figure 5(b), we show an instance in which a rockfall event is mis-  
621 classified as a quake. The rockfall event displays multiple peaks; however,  
622 these peaks, aside from the principal one, are of low magnitude and the event  
623 has a very short time span. Analysis of the LRP representation illustrates  
624 a concentration of positive effects (depicted in red) at the primary peak of  
625 the event. Conversely, numerous negative contributions (depicted in blue)  
626 are observed at the secondary peaks, suggesting that the presence of these  
627 multiple peaks is not taken into account due to their limited magnitudes;



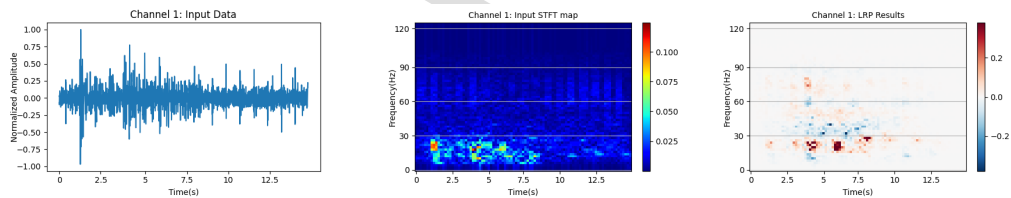
(a) Quake misclassified as rockfall



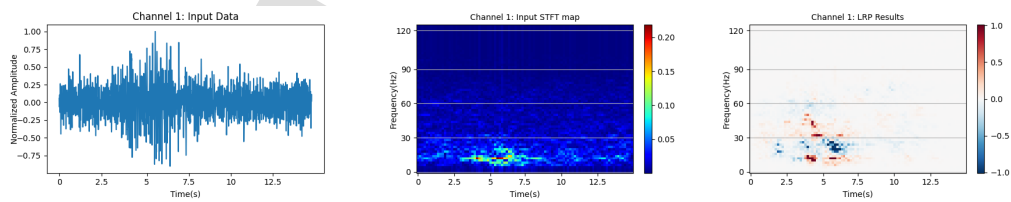
(b) Rockfall misclassified as quake



(c) Quake misclassified as earthquake



(d) Earthquake misclassified as rockfall



(e) Noise misclassified as earthquake

Figure 5: Misclassified examples.

628 hence, the model finally classifies this event as a quake.

629 In Figure 5(c), we present an instance of a quake misclassified as an  
630 earthquake. This misclassification is evident in the LRP map, where both  
631 high-frequency and low-frequency components simultaneously exhibit posi-  
632 tive contributions around the 3-second period. Thus, the model interprets  
633 this segment as a P-wave. Furthermore, at approximately 5 seconds into the  
634 waveform, a positive contribution appears in the low-frequency range. Al-  
635 though the primary peak of this event occurs around 3 seconds, the spectro-  
636 gram reveals that the low-frequency component persists for an extended dura-  
637 tion. Moreover, the event is influenced by higher-frequency noise (exceeding  
638 30Hz), and this high-frequency noise coincides with the primary waveform  
639 peak around the 3 seconds. Consequently, this led the model to mistakenly  
640 identify it as a P-wave, with the prolonged low-frequency component be-  
641 ing mistakenly identify as a S-wave. These observations align with seismic  
642 features of earthquakes, thereby causing the model's misclassification as an  
643 earthquake event.

644 In Figure 5(d), we encounter an instance where an earthquake is mis-  
645 takenly classified as a rockfall. The LRP map highlights multiple spectral  
646 peaks, which is a feature of rockfall events. However, this event may have  
647 resulted from an earthquake occurring amidst background noise, exhibiting  
648 a distinctive multi-peak pattern. Thus, despite the presence of a P-wave  
649 at approximately 1 second and an S-wave at roughly 4 seconds, complex  
650 background noise caused misclassification.

651 In Figure 5(e), the misclassification of noise as an earthquake is shown.  
652 The noise signal exhibits prominent peaks around 4 seconds and 5.5 sec-  
653 onds. Examination of the LRP map reveals the model's recognition of low-  
654 frequency and high-frequency components (15-20Hz) around the 4-second  
655 mark, along with low-frequency signals at 5.5 seconds (15Hz). This aligns  
656 with the characteristic features of P-waves and S-waves in earthquake sig-  
657 nals, resulting in the model's misclassification as an earthquake. The result

658 might have been different if time-series signals were inputted to the network  
 659 instead of the STFT maps as can be seen from the left time-series plot that  
 660 shows high level of noise throughout the signal.

661 We can see from these examples that most misclassifications are due to  
 662 high level of background noise. The next example highlights another origin  
 663 of error related to the filtering process. Figure 6 displays an unfiltered earth-  
 664 quake waveform with a frequency below 3 Hz, characteristic of low-frequency  
 665 earthquakes that are rarely associated with active landslides [45]. Since our  
 666 focus is on detecting local seismic events related to landslides, we apply a  
 667 bandpass filter in the 5-60 Hz range (see Sec. 4.2), which excludes these low-  
 668 frequency earthquakes. Consequently, this filter removed the low-frequency  
 669 event's waveform, leaving only background noise as input to the CNN. As il-  
 670 lustrated in Figure 7, the LRP map indicates that the model failed to extract  
 671 meaningful features from the filtered input, resulting in the earthquake being  
 672 misclassified as noise. This misclassification can be attributed to the rarity  
 673 and uniqueness of low-frequency earthquakes on landslides, as our filter in-  
 674 advertently eliminated their distinctive waveforms, confounding the CNN's  
 675 classification process.

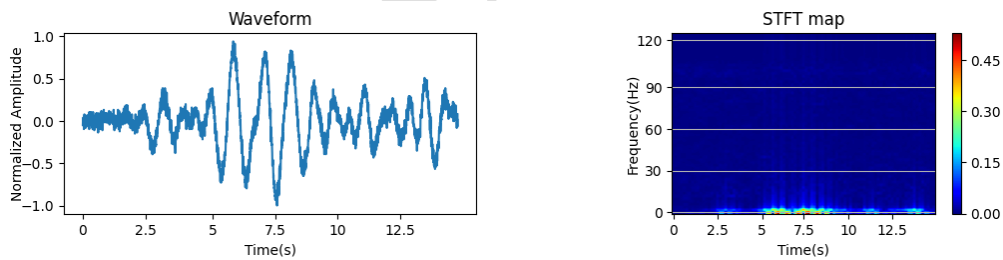


Figure 6: Waveform (left) and STFT map (right) of the unfiltered low-frequency earthquake.

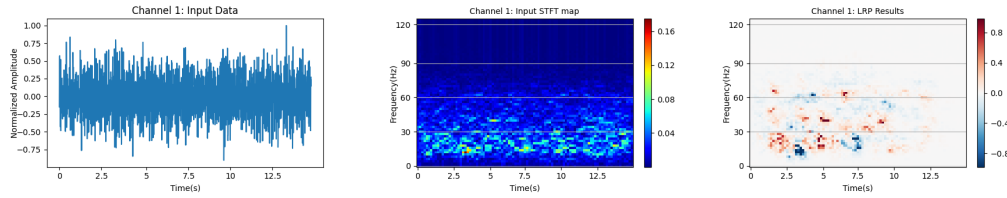


Figure 7: Waveform (left), STFT map (middle) and the LRP map (right) of the filtered low-frequency earthquake.

#### 676 5.4. Re-labelling results

677 Figure 8 shows three examples of misclassifications, which could be due  
 678 to human error during expert labelling. The example shown in Figure 8(a), is  
 679 an event classified by the model as noise, though the domain experts labelled  
 680 it as a quake. In the STFT representation of the signal, no obvious peak  
 681 corresponding to the event was discernible. Moreover, the LRP map exhibits  
 682 a disordered distribution of relevance. Collectively, these findings lead to the  
 683 argument that the event in question is more likely to be anthropogenic noise  
 684 rather than a quake. Figure 8(b) illustrates a similar situation where the  
 685 event is mistakenly labelled as an earthquake. There are no clear P-waves  
 686 at both low and high frequencies, and there are no S-waves with high en-  
 687 ergy content. For this earthquake event, we also examined the unfiltered raw  
 688 signal, and it still did not exhibit any earthquake waveform characteristics.  
 689 Figure 8(c) shows an example that was classified as a rockfall by the CNN  
 690 model, while the expert labelled it as a seismic quake. It can be concluded  
 691 from the LRP map that the model focused on multiple peaks in the event,  
 692 with a frequency distribution centred around 30Hz, characteristics that align  
 693 with typical rockfall patterns. In contrast, quakes tend to exhibit a single  
 694 dominant peak, a feature that was notably absent in the input STFT map,  
 695 where multiple peaks were discernible. Consequently, based on these dis-  
 696 tinctive patterns and spectral features, it becomes evident that the event in  
 697 question is more accurately classified as a rockfall.

698 Here we list all corrections made to the expert catalogue, following above



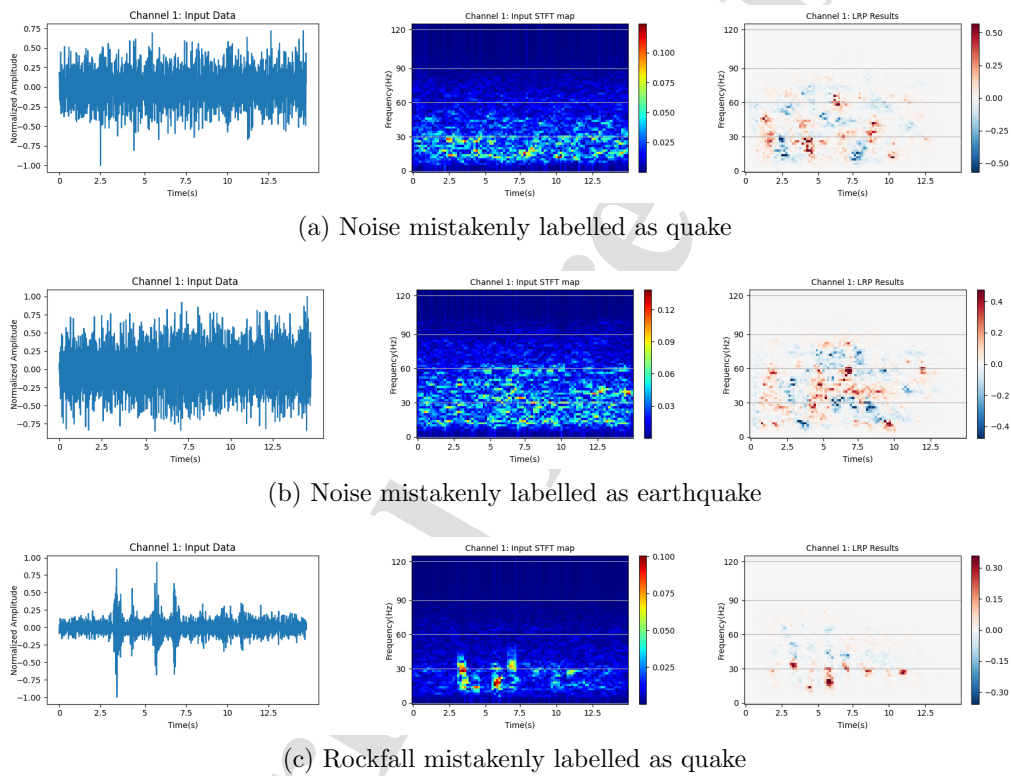


Figure 8: Three examples of events with labels corrected.

699 explainability and queries. Specifically, 7 quakes were relabelled as noise  
 700 as per example Figure 8(a), 1 earthquake was relabelled as noise (shown in  
 701 Figure 8(b)), and 1 quake as rockfall (Figure 8(c)). In addition, some noise  
 702 events were labelled by the expert though these events occurred very close to  
 703 earthquake, quake and rockfall events, which potentially caused confusion.  
 704 Hence, we removed all noise events that occurred in close proximity (within  
 705 30s) to the earthquake, quake and rockfall events - this way 38 noise events  
 706 were removed.

707 Thus, after this relabelling there are 38 quakes, 17 earthquakes, 75 rock-  
 708 falls and 689 anthropogenic noise events in total. The verified catalogue of  
 709 events is provided as supplementary material to this paper, as a contribu-  
 710 tion to address the second and third principles of Trustworthy AI, related to  
 711 reproducibility and data access. Specifically, the 260 verified events on the  
 712 25th Nov. 2015 are listed in the Training events supplementary material,  
 713 identified by the date. The 819 verified events on 26th to 28th Nov. 2014  
 714 are listed in the Additional 3-day catalogue supplementary material. In or-  
 715 der for other researchers to enable benchmarking, Table 2 and Table 3 show  
 716 the confusion matrix and classification performance after the re-labelling, re-  
 717 spectively. Although the F1-score for quake events is low, we have a high  
 718 recall but precision is low because of 8 instances of false positives for rockfall.  
 719 There are relatively few instances of quake and earthquake, which explains  
 720 why the F1-score is not the best indicator of performance and the confusion  
 721 matrix provides a more explainable and trustworthy measure of performance.

Table 2: The confusion matrix after label correction. The numbers in the brackets show the recall values.

		Model			
		Quake	Earthquake	Rockfall	Noise
Expert	Quake	<b>26</b> (68.4%)	2	8	2
	Earthquake	0	<b>15</b> (88.2%)	1	1
	Rockfall	2	0	<b>73</b> (97.3%)	0
	Noise	95	11	37	<b>546</b> (79.2%)

Table 3: The classification performance after label correction.

	Precision	Recall	F1-score
Quake	0.21	0.68	0.32
Earthquake	0.54	0.88	0.67
Rockfall	0.61	0.97	0.75
Noise	0.99	0.79	0.88

## 722 6. Conclusions and Future Work

723 The paper discusses the significance of the 7 principles of Trustworthy  
 724 AI, including human oversight, technical robustness, data governance and  
 725 transparency to the challenging problem of micro-seismic signal analysis. To  
 726 this effect, we propose a human-on-the-loop microseismic multi-class classi-  
 727 fication method together with LRP to shed light on feature importance in  
 728 order to in turn verify any possible human labelling error.

729 We demonstrate that the generated LRP maps assist human experts in  
 730 manual event classification. LRP clearly identifies properties of the signals  
 731 extracted by the network when making decisions. Based on this, we con-  
 732 cluded, for example, that the main reason why quake events are often mis-  
 733 classified as rockfall is due to appearance of a noise signal at multiple higher  
 734 frequencies that mimics rockfalls. Due to human error, experts may occa-  
 735 sionally mislabel events in the catalogue due to the similarity of event char-  
 736 acteristics, complexity of seismic data and large volume of data that needs to  
 737 be processed. However, the availability of LRP maps as a visual aid can offer  
 738 a valuable tool to verify and refine the expert’s classifications. This collabo-  
 739 rative synergy between automated and manual classification can enhance the  
 740 accuracy of microseismic catalogues, contributing to a better understanding  
 741 of geological processes.

742 Besides assisting with event labelling, another application of the LRP  
 743 maps is improving the model’s performance. Indeed, by observing the in-  
 744 sights gained through XAI tools, we discern specific features of input events  
 745 that are prone to misclassification by the CNN, which is instrumental in en-

746 hancing the robustness and generalisability of the model that can be achieved  
747 by adding more events in the training set that closely resemble the challenging  
748 input patterns identified through XAI. For example, when we discover that  
749 certain event features consistently lead to misclassifications, we collect and  
750 add more events with similar attributes into the training dataset. This tar-  
751 geted data augmentation approach has the potential to improve the model's  
752 ability to distinguish between challenging seismic events, thereby increasing  
753 model's robustness and classification performance.

754 Since LRP assigns relevance scores to highlight the most influential fea-  
755 tures for each classification, it is important to determine if these relevance  
756 patterns remain stable across various geographic areas and seismometer char-  
757 acteristics, such as sensitivity, sampling rate, and axis configurations. This  
758 evaluation will help ascertain the reliability of LRP explanations across di-  
759 verse equipment types and environments. In future work, we plan to test our  
760 system in various geographic regions and with different seismometer config-  
761 urations to assess the consistency and robustness of LRP interpretability,  
762 enhancing the broader applicability and trustworthiness of our approach.

763 Given the potential variability in expert interpretations, it is important  
764 to explore how different experts' insights may affect labeling. Future studies  
765 could employ a multi-expert assessment framework that incorporates confi-  
766 dence levels, based on the methodologies proposed by [46], to better under-  
767 stand this variability and further enhance the reliability of the classification  
768 process.

769 Since classification of quakes remains challenging, the current model could  
770 be adapted to classify a broader range of events, including low frequency  
771 events and types of anthropogenic noise, by expanding the training set and  
772 retraining the model, with LRP providing the explanations. To maximize  
773 accuracy and trust in AI-driven seismic signal analysis, integrating human  
774 expertise with AI models is important. Developing interactive explainability  
775 tools that facilitate iterative feedback from geoscientists could lead to con-

776 tinuous improvements in model performance and foster greater confidence in  
777 AI-generated outputs.

### 778 **Acknowledgement**

779 This work was supported in part by the EPSRC Prosperity Partnership  
780 research and innovation programme under grant agreement EP/S005560/1,  
781 and in part by the EPSRC New Horizons research programme EP/X01777X/1.  
782 The contextual data interpretation and labelling work by experts on the SZ  
783 dataset was supported by RSE Saltire International Collaboration Awards.

### 784 **CRedit authorship contribution statement**

785 **Jiaxin Jiang:** Methodology, Software, Validation, Formal Analysis, In-  
786 vestigation, Writing – Original Draft, Visualization. **David Murray:** Data  
787 curation, Validation, Formal Analysis, Writing- Reviewing and Editing. **Vladimir**  
788 **Stankovic:** Conceptualization, Writing- Reviewing and Editing, Supervi-  
789 sion, Project administration, Funding acquisition. **Lina Stankovic:** Con-  
790 ceptualization, Writing- Reviewing and Editing, Supervision, Project admin-  
791 istration, Funding acquisition. **Clement Hibert:** Investigation, Validation,  
792 Data curation. **Stella Pytharouli:** Funding acquisition, Project adminis-  
793 tration. **Jean-Philippe Malet:** Resources, Data curation.

### 794 **Declaration of Competing Interest**

795 The authors declare that they have no known competing financial inter-  
796 ests or personal relationships that could have appeared to influence the work  
797 reported in this paper.

### 798 **Data availability**

799 The code is provided in a GitHub repository at [https://github.com/](https://github.com/kanata2020/Explainable-seismic-classification)  
800 [kanata2020/Explainable-seismic-classification](https://github.com/kanata2020/Explainable-seismic-classification). This includes the mod-

801 els and data used for classification, as well as algorithms for explainable  
802 visualization.

### 803 References

- 804 [1] S. M. Mousavi, G. C. Beroza, Deep-learning seismology, *Science*  
805 377 (6607) (2022) eabm4470.
- 806 [2] Z. Bayramoğlu, M. Uzar, Performance analysis of rule-based classifica-  
807 tion and deep learning method for automatic road extraction, *Internation-  
808 al Journal of Engineering and Geosciences* 8 (1) (2023) 83–97.
- 809 [3] O. D. Gülgün, H. Erol, Classification performance comparisons of deep  
810 learning models in pneumonia diagnosis using chest x-ray images, *Turk-  
811 ish Journal of Engineering* 4 (3) (2020) 129–141.
- 812 [4] B. Sariturk, B. Bayram, Z. Duran, D. Z. Seker, Feature extraction from  
813 satellite images using segnet and fully convolutional networks (fcn), *Inter-  
814 national Journal of Engineering and Geosciences* 5 (3) (2020) 138–143.
- 815 [5] M. E. Dos, Determination of city change in satellite images with deep  
816 learning structures, *Advanced Remote Sensing* 2 (1) (2022) 16–22.
- 817 [6] Y. Kaya, H. İ. Şenol, A. Y. Yiğit, M. Yakar, Car detection from very  
818 high-resolution uav images using deep learning algorithms, *Photogram-  
819 metric Engineering & Remote Sensing* 89 (2) (2023) 117–123.
- 820 [7] D. Yi, S. Yiran, C. Ismet, L. Xun, S. Guangyao, Classification of clus-  
821 tered microseismic events in a coal mine using machine learning, *Journal  
822 of Rock Mechanics and Geotechnical Engineering* 13 (6) (2021) 1256–  
823 1273.
- 824 [8] M. Shakeel, K. Itoyama, K. Nishida, K. Nakadai, Etc: Earthquake  
825 magnitudes classification on seismic signals via convolutional recurrent

- 826 networks, in: 2021 IEEE/SICE International Symposium on System  
827 Integration (SII), IEEE, 2021, pp. 388–393.
- 828 [9] T. Perol, M. Gharbi, M. Denolle, Convolutional neural network for  
829 earthquake detection and location, *Science Advances* 4 (2) (2018)  
830 e1700578.
- 831 [10] J. Li, M. Ye, L. Stankovic, V. Stankovic, S. Pytharouli, Domain knowl-  
832 edge informed multitask learning for landslide-induced seismic classifi-  
833 cation, *IEEE Geoscience and Remote Sensing Letters* 20 (2023) 1–5.
- 834 [11] A. Parasyris, L. Stankovic, V. Stankovic, Synthetic data generation for  
835 deep learning-based inversion for velocity model building, *Remote Sens-  
836 ing* 15 (11) (2023) 2901.
- 837 [12] J. Jiang, V. Stankovic, L. Stankovic, E. Parastatidis, S. Pytharouli,  
838 Microseismic event classification with time-, frequency-, and wavelet-  
839 domain convolutional neural networks, *IEEE Transactions on Geo-  
840 science and Remote Sensing* 61 (2023) 1–14.
- 841 [13] L. Trani, G. A. Pagani, J. P. P. Zanetti, C. Chapeland, L. Evers, Deep-  
842 quake—an application of cnn for seismo-acoustic event classification in  
843 the netherlands, *Computers & Geosciences* 159 (2022) 104980.
- 844 [14] W.-Y. Liao, E.-J. Lee, C.-C. Wang, P. Chen, F. Provost, C. Hiber, J.-  
845 P. Malet, C.-R. Chu, G.-W. Lin, RockNet: Rockfall and earthquake  
846 detection and association via multitask learning and transfer learning,  
847 *Authorea* (2022).
- 848 [15] C. Sonesson, S. Gerster, M. Delorenzi, Batch effect confounding leads to  
849 strong bias in performance estimates obtained by cross-validation, *PloS  
850 one* 9 (6) (2014) e100335.
- 851 [16] M. Hägele, P. Seegerer, S. Lapuschkin, M. Bockmayr, W. Samek,  
852 F. Klauschen, K.-R. Müller, A. Binder, Resolving challenges in deep

- 853 learning-based analyses of histopathological images using explanation  
854 methods, *Scientific reports* 10 (1) (2020) 1–12.
- 855 [17] J. Li, L. Stankovic, S. Pytharouli, V. Stankovic, Automated platform  
856 for microseismic signal analysis: Denoising, detection, and classification  
857 in slope stability studies, *IEEE Transactions on Geoscience and Remote*  
858 *Sensing* 59 (9) (2020) 7996–8006.
- 859 [18] J. Li, L. Stankovic, V. Stankovic, S. Pytharouli, C. Yang, Q. Shi, Graph-  
860 based feature weight optimisation and classification of continuous seis-  
861 mic sensor array recordings, *Sensors* 23 (1) (2023).
- 862 [19] D. Bau, J.-Y. Zhu, H. Strobel, A. Lapedriza, B. Zhou, A. Torralba,  
863 Understanding the role of individual units in a deep neural network,  
864 *Proceedings of the National Academy of Sciences* 117 (48) (2020) 30071–  
865 30078.
- 866 [20] A. Holzinger, From machine learning to explainable ai, in: 2018 world  
867 symposium on digital intelligence for systems and machines (DISA),  
868 IEEE, IEEE, 2018, pp. 55–66.
- 869 [21] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, K.-R. Müller, Eval-  
870 uating the visualization of what a deep neural network has learned, *IEEE*  
871 *transactions on neural networks and learning systems* 28 (11) (2016)  
872 2660–2673.
- 873 [22] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, K.-R. Müller, Layer-  
874 Wise Relevance Propagation: An Overview, Springer-Verlag, Berlin,  
875 Heidelberg, 2022, p. 193–209.
- 876 [23] S. Lapuschkin, A. Binder, G. Montavon, K.-R. Müller, W. Samek,  
877 Analyzing classifiers: Fisher vectors and deep neural networks, in:  
878 2016 IEEE Conference on Computer Vision and Pattern Recognition  
879 (CVPR), IEEE, 2016, pp. 2912–2920.



- 880 [24] D. Murray, L. Stankovic, V. Stankovic, Transparent ai: Explainability  
881 of deep learning based load disaggregation, in: Proceedings of the 8th  
882 ACM International Conference on Systems for Energy-Efficient Build-  
883 ings, Cities, and Transportation, BuildSys '21, Association for Comput-  
884 ing Machinery, New York, NY, USA, 2021, p. 268–271.
- 885 [25] E. Commission, C. Directorate-General for Communications Networks,  
886 Technology, Ethics guidelines for trustworthy AI, Publications Office,  
887 2019.
- 888 [26] U. N. S. D. G. 13, accessed: 2023. [online]. Available: [https://sdgs.  
889 un.org/goals](https://sdgs.un.org/goals) (2023).
- 890 [27] N. Vouillamoz, S. Rothmund, M. Joswig, Characterizing the complexity  
891 of microseismic signals at slow-moving clay-rich debris slides: the super-  
892 sauze (southeastern france) and pechgraben (upper austria) case studies,  
893 Earth Surface Dynamics 6 (2) (2018) 525–550.
- 894 [28] A. Renouard, A. Maggi, M. Grunberg, C. Doubre, C. Hibert, Toward  
895 false event detection and quarry blast versus earthquake discrimination  
896 in an operational setting using semiautomated machine learning, Seis-  
897 mological Society of America 92 (6) (2021) 3725–3742.
- 898 [29] F. Provost, C. Hibert, J.-P. Malet, Automatic classification of endoge-  
899 nous landslide seismicity using the random forest supervised classifier,  
900 Geophy. Research Let. 44 (1) (2017) 113–120.
- 901 [30] X. Bi, C. Zhang, Y. He, X. Zhao, Y. Sun, Y. Ma, Explainable time-  
902 frequency convolutional neural network for microseismic waveform clas-  
903 sification, Information Sciences 546 (2021) 883–896.
- 904 [31] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra,  
905 Grad-CAM: Visual explanations from deep networks via gradient-based

- 906 localization, *International Journal of Computer Vision* 128 (2) (2019)  
907 336–359.
- 908 [32] J. Jiang, V. Stankovic, L. Stankovic, D. Murray, S. Pytharouli, Explain-  
909 able ai for transparent seismic signal classification, in: *IGARSS 2024-*  
910 *2024 IEEE International Geoscience and Remote Sensing Symposium*,  
911 IEEE, 2024, pp. 8801–8805.
- 912 [33] FrenchLandslideObservatorySeismologicalDatacenter/RESIF, Observa-  
913 toire multi-disciplinaire des instabilites de versants (omiv), accessed:  
914 2021. [online]. Available: <https://seismology.resif.fr/>. DOI:  
915 10.15778/RESIF.MT.
- 916 [34] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller,  
917 W. Samek, On pixel-wise explanations for non-linear classifier decisions  
918 by layer-wise relevance propagation, *PLOS ONE* 10 (7) (2015) 1–46.
- 919 [35] S. M. Mousavi, W. L. Ellsworth, W. Zhu, L. Y. Chuang, G. C. Beroza,  
920 Earthquake transformer—an attentive deep-learning model for simulta-  
921 neous earthquake detection and phase picking, *Nature Communications*  
922 11 (1) (2020) 1–12.
- 923 [36] F. Provost, J.-P. Malet, C. Hibert, A. Helmstetter, M. Radiguet, D. Ami-  
924 trano, N. Langet, E. Larose, C. Abancó, M. Hürlimann, T. Lebourg,  
925 C. Levy, G. Le Roy, P. Ulrich, M. Vidal, B. Vial, Towards a standard  
926 typology of endogenous landslide seismic sources, *Earth Surface Dynam-*  
927 *ics* 6 (4) (2018) 1059–1088.
- 928 [37] F. Provost, J.-P. Malet, J. Gance, A. Helmstetter, C. Doubre, Automatic  
929 approach for increasing the location accuracy of slow-moving landslide  
930 endogenous seismicity: the APOLoc method, *Geophysical Journal In-*  
931 *ternational* 215 (2) (2018) 1455–1473.

- 932 [38] D. Murray, L. Stankovic, S. Pytharouli, V. Stankovic, Semi-supervised  
933 seismic event detection using siamese networks, 2023, eGU General As-  
934 sembly 2023, April 2023.
- 935 [39] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-  
936 scale image recognition (2015). [arXiv:1409.1556](https://arxiv.org/abs/1409.1556).
- 937 [40] O. M. Saad, Y. Chen, Earthquake detection and p-wave arrival time  
938 picking using capsule neural network, *IEEE Transactions on Geoscience  
939 and Remote Sensing* 59 (7) (2020) 6234–6243.
- 940 [41] O. M. Saad, Y. Chen, Capsphase: Capsule neural network for seismic  
941 phase classification and picking, *IEEE Transactions on Geoscience and  
942 Remote Sensing* 60 (2021) 1–11.
- 943 [42] A. Chan, M. Schneider, M. Körner, Xai for early crop classification, in:  
944 *IGARSS 2023-2023 IEEE International Geoscience and Remote Sensing  
945 Symposium*, IEEE, 2023, pp. 2657–2660.
- 946 [43] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, K.-R. Müller, Ex-  
947 plaining nonlinear classification decisions with deep taylor decomposi-  
948 tion, *Pattern recognition* 65 (2017) 211–222.
- 949 [44] M. Alber, S. Lapuschkin, P. Seegerer, M. Hägele, K. T. Schütt, G. Mon-  
950 tavon, W. Samek, K.-R. Müller, S. Dähne, P.-J. Kindermans, inves-  
951 tigate neural networks!, *Journal of Machine Learning Research* 20 (93)  
952 (2019) 1–8.
- 953 [45] K. Masuda, S. Ide, K. Ohta, T. Matsuzawa, Bridging the gap between  
954 low-frequency and very-low-frequency earthquakes, *Earth, Planets and  
955 Space* 72 (1) (2020) 1–9.
- 956 [46] T. Sobot, V. Stankovic, L. Stankovic, Human in the loop active learning  
957 for time-series electrical measurement data, *Engineering Applications of  
958 Artificial Intelligence* 133 (2024) 108589.

## Highlights

**A human-on-the-loop approach for labelling landslide-induced seismic recordings via a multi-class deep-learning based classification model**

Jiaxin Jiang, David Murray, Vladimir Stankovic, Lina Stankovic, Clement Hibert, Stella Pytharouli, Jean-Philippe Malet

- Robust multi-class CNN-based seismic signal classifier
- LRP explainability maps for model diagnosis
- Trustworthy AI with geoscientist in the design loop

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

---

Vladimir Stankovic reports financial support was provided by University of Strathclyde. If there are other authors, they declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

---