

# Machine Learning improves warning systems of debris flows

Małgorzata Chmiel<sup>1</sup>, Fabian Walter<sup>2</sup>, Michaela Wenner<sup>2</sup>, Zhen Zhang<sup>3</sup>, McArdell Brian W.<sup>4</sup>, and Clément Hibert<sup>5</sup>

<sup>1</sup>Laboratory of Hydraulics, Hydrology and Glaciology, ETH Zürich

<sup>2</sup>ETH Zürich

<sup>3</sup>Institute of Mountain Hazards and Environment, Chinese Academy of Sciences

<sup>4</sup>Swiss Federal Institute for Forest, Snow and Landscape Research WSL

<sup>5</sup>Institut De Physique Du Globe De Strasbourg

November 24, 2022

## Abstract

Automatic identification of debris flow signals in continuous seismic records remains a challenge. To tackle this problem we use a machine learning approach, which can be applied to continuous real-time data streams. We show that a machine learning model based on the random forest algorithm recognizes different stages of debris flow formation and propagation at the Illgraben torrent, Switzerland, with an accuracy exceeding 90%. In contrast to typical debris flow detection requiring instrumentation installed directly in the torrent, our approach provides a significant gain in warning times of tens of minutes to hours. For real-time data streams from 2020, our detector raises alarms for all 8 independently confirmed Illgraben events and gives no false alarms. We suggest that our seismic machine-learning detector is a critical step towards the next generation of debris-flow warning, which increases warning times using both simpler and cheaper instrumentation compared to existing operational systems.

# Machine Learning improves warning systems of debris flows

Małgorzata Chmiel<sup>1</sup>, Fabian Walter<sup>1</sup>, Michaela Wenner<sup>1,2</sup>, Zhen Zhang<sup>1,3,4</sup>,  
Brian W. McArdell<sup>2</sup>, Clement Hibert<sup>6</sup>

<sup>1</sup>Laboratory of Hydraulics, Hydrology and Glaciology, ETH Zürich, Zürich, Switzerland

<sup>2</sup>Swiss Federal Institute for Forest, Snow and Landscape Research, Zürich, Switzerland

<sup>3</sup>Key Laboratory of Mountain Hazards and Surface Process, Institute of Mountain Hazards and

Environment, Chinese Academy of Sciences, Chengdu, China

<sup>4</sup>University of Chinese Academy of Sciences, Beijing, China

<sup>5</sup>Institut de Physique du Globe de Strasbourg, CNRS UMR 7516, University of Strasbourg/EOST, 7 8

Strasbourg, France

## Key Points:

- A novel debris-flow detector is developed using a machine-learning model and seismic data from a Swiss torrent.
- Signals of 22 debris flows recorded by six seismic stations are used to train and test the machine-learning model.
- A detector is running on the continuous real-time 2020 data stream detecting all 13 debris flows in 3 months and raising no false alarms.

---

Corresponding author: Małgorzata Chmiel, [chmielm@ee.ethz.ch](mailto:chmielm@ee.ethz.ch)

## Abstract

Automatic identification of debris flow signals in continuous seismic records remains a challenge. To tackle this problem we use a machine learning approach, which can be applied to continuous real-time data streams. We show that a machine learning model based on the random forest algorithm recognizes different stages of debris flow formation and propagation at the Illgraben torrent, Switzerland, with an accuracy exceeding 90 %. In contrast to typical debris flow detection requiring instrumentation installed directly in the torrent, our approach provides a significant gain in warning times of tens of minutes to hours. For real-time data streams from 2020, our detector raises alarms for all 8 independently confirmed Illgraben events and gives no false alarms. We suggest that our seismic machine-learning detector is a critical step towards the next generation of debris-flow warning, which increases warning times using both simpler and cheaper instrumentation compared to existing operational systems.

## Plain Language Summary

Debris flows are fast-moving masses of mud, soil, fragmented rock, and water transporting large volume of material in mountainous areas. They pose a significant danger to human life, properties, and infrastructure. Thus, it is crucial to reliably detect debris flows early enough to send an alarm message to local communities. We propose a novel approach for debris-flow detections using recorded ground vibrations generated by 22 debris flows in Illgraben, Switzerland. We use a machine-learning algorithm that automatically learns to distinguish between debris flow generated ground vibrations and other seismic signals. This allows us to increase warning times by at least 42 min comparing to existing detection systems at Illgraben.

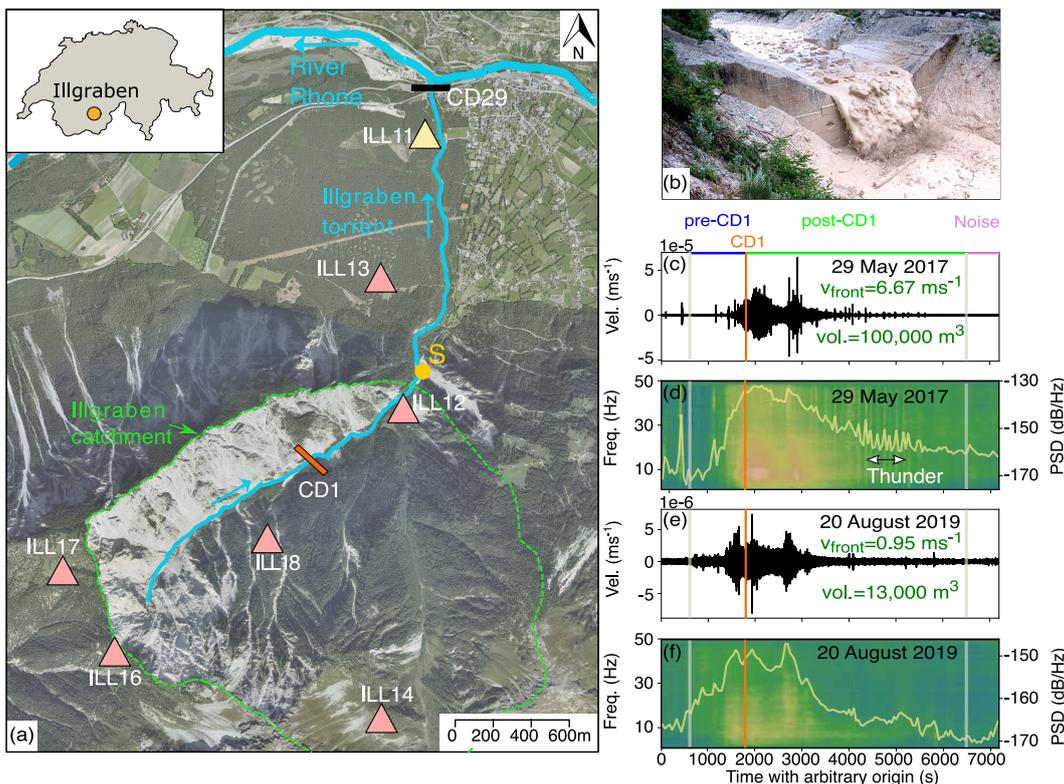
## 1 Introduction

Debris flows are mixtures of water and sediments of all sizes, which are mobilized by heavy precipitation in steep Alpine torrents. They move downstream with average velocities of several meters per seconds (Hürlimann et al., 2003) showing a flow behaviour in-between landslides and sediment transporting floods (Iverson, 1997). Debris flows have a high destructive potential, which is amplified at the flow front, where large cobbles and boulders concentrate (Iverson, 1997). The significant hazard to human life and infrastructure in Alpine regions, including Switzerland (e.g., Jakob & Hungr, 2005; Badoux et al., 2016) demands reliable warning systems to reduce risk in exposed terrain (e.g., Stähli et al., 2015).

Recently, modern seismic instrumentation has suggested new warning perspectives, because even at large distances (tens to hundreds of kilometers) seismometers can detect high-frequency ( $>1$  Hz) ground unrest induced by debris flows (for a review see Allstadt et al., 2018). This may extend warning times compared to conventional instrumentation within or near debris flow torrents, which can only be installed in accessible terrain (e.g., Arattano & Marchi, 2008; Coviello et al., 2019).

Despite recent advances in our theoretical understanding of high-frequency debris flow seismograms (Cole et al., 2009; Kean et al., 2012; Lai et al., 2018; Farin et al., 2019), seismometers installed at larger distances from torrents have yet to be implemented in operational warning systems. Identification of debris flow signals in the presence of other seismic activity remains a challenge. Since seismic debris flow signals have moderate amplitudes, simple threshold-based detection criteria cannot distinguish them from cultural noise, earthquakes and other Alpine mass movements at a permissible false detection rate (Walter et al., 2017).

Here, we introduce a machine learning approach to detect debris flows based on their seismic signature. For the Illgraben torrent, Switzerland, seismic records from an 8-station



**Figure 1.** Study site. (a) Illgraben catchment is outlined with green dashed line (source: Swisstopo). Check dam (CD) 1 and CD29 are represented with orange and black bars. Seismometer locations are indicated with triangles. Connection of east side hillslope (Sagenschleif) with Illgraben channel is marked at Point S. Inset shows Illgraben’s location in Switzerland. Station ILL11 not used in detection is marked in yellow and station ILL15 is located outside of the presented map. (b) Photo of Illgraben debris-flow passing CD29 (Source: WSL). (c) Vertical ground velocity recorded at ILL18 on 29 May 2017 (large and fast event) and the corresponding power spectral density (PSD) in (d). (e) Vertical ground velocity recorded at ILL18 on 20 August 2019 (small and slow event) and the corresponding PSD in (f). Arrival times of the debris flows at CD1 are marked with an orange line. In (d) and (f) PSDs averaged over 1-50 Hz are shown with a yellow line. Time windows between gray and orange lines divide the records in (c)-(f) into 3 signals classes.

67 network allow for debris flow detection in the upper catchment area, where instrument  
 68 deployment is not possible. Trained with data from 20 events, our detection algorithm was  
 69 subjected to real-time data streams from summer 2020 and identified all 13 debris flows with  
 70 no false alarms. Our approach adds up to an hour of warning time to the earliest possible  
 71 in-torrent detection at Illgraben.

## 72 2 Study site

73 Located in the southwestern part of Switzerland, Illgraben is one of the most active  
 74 debris-flow torrents in the European Alps (Rickenmann et al., 2001) (Figure 1a). The  
 75 catchment area extends over 10.4 km<sup>2</sup> from the summit of the Illhorn mountain [2716 m  
 76 above sea level (asl)] to the Rhone River (610 m asl). The steep slopes (~40°) of the  
 77 upper Illgraben catchment are characterized by rockfall and landslide activity (Berger et  
 78 al., 2011b). The resulting sediments accumulate downslope or in the Illgraben channel and

79 provide sliding material with volumes ranging from 500 to more than 4,000 m<sup>3</sup> (Schlunegger  
80 et al., 2009). During heavy precipitations and intense summer thunderstorms from April to  
81 October this material is regularly mobilized in form of debris flows (Badoux et al., 2009).  
82 The larger debris flow volumes (10<sup>3</sup> to 10<sup>5</sup> m<sup>3</sup>) result from cumulative sediment mobilization  
83 and entrainment along the flow path and often reach the Rhone River. Like elsewhere,  
84 Illgraben debris flows have boulder-rich fronts resulting from particle sorting phenomena  
85 (Pierson, 1986; Johnson et al., 2012) followed by turbulent slurry with a large concentration  
86 of suspended sediments of variable granulometry and water content (Costa, 1984; Iverson,  
87 1997; McCoy et al., 2010; Berger et al., 2011a, 2011b).

88 In 1961 a major landslide occurred in the upper Illgraben catchment and resulted in a  
89 debris flow destroying the bridge of the Cantonal highway along the Rhone river (Graf et  
90 al., 2007; Berger et al., 2011a). Consequently, a series of 30 Check Dams (CD; see Figure  
91 1b for lowest CD29) was placed along the lower 3.4 km of the channel in order to stabilize  
92 the current debris flow path, expand discharge capacity, and minimize erosion.

93 As debris-flows still pose a hazard to people crossing the channel and to nearby infras-  
94 tructure, an in-torrent warning system was commissioned by the governmental authorities  
95 and installed in 2007 (Badoux et al., 2009). The system consists of geophone detectors in  
96 check dams and flow depth measurements in the lower Illgraben part (Badoux et al., 2009;  
97 McArdell et al., 2007). Similar instruments and a recently re-deployed force plate form the  
98 debris flow observatory, which is operated for research purposes independently of the warn-  
99 ing system since 2000 (Hürlimann et al., 2003; McArdell et al., 2007; Berger et al., 2011a).  
100 The observatory provides estimates of debris flow depth, volume and density (Schlunegger  
101 et al., 2009).

102 Illgraben’s warning system has undergone different upgrades but is still based on the  
103 initial measurement principles. The earliest possible detection is provided by geophones  
104 installed inside check dam CD1 (Figure 1a), but this detection point is not deemed reliable  
105 as it is contingent upon solar power supply and mobile phone reception, which vary as a  
106 consequence of shadowing effects from the canyon walls. However, the CD1 arrival times  
107 can be downloaded in retrospect and are available for subsequent data analysis.

108 The present warning system in Illgraben requires instrument installation in direct con-  
109 tact with the torrent, which implies that detection is insensitive to sediment movement in  
110 the highly unstable and inaccessible upper catchment. This is a major weakness as debris  
111 flows mobilize in the upper catchment above CD1 (Schlunegger et al., 2009), where detection  
112 could increase warning times by tens of minutes.

## 113 2.1 Seismic debris-flow data

114 In past years, researchers have attempted to extend detection capabilities to the upper  
115 Illgraben catchment using seismological and acoustic measurements (Burtin et al., 2016;  
116 Walter et al., 2017; Schimmel et al., 2018; Marchetti et al., 2019; Wenner et al., 2019).  
117 In this context, since 2017, a seasonal seismometer network has been installed around the  
118 Illgraben catchment during spring/summer months (Figure 1a). The network consists of  
119 8 3-component 1 Hz seismometers recording ground velocity continuously at a sampling  
120 frequency of 100 Hz.

121 All 8 stations are equipped with a modem to stream seismic data via the 4G mobile  
122 phone network to the Swiss Seismological Service from where they are made available via a  
123 seedlink server. A python client on a work station at the Institute of Hydraulics, Hydrology  
124 and Glaciology (VAW) at ETH Zürich receives data packages every 8-9 seconds, which are  
125 concatenated to continuous records. Disregarding station ILL18 (the only station within  
126 the Illgraben canyon), realtime data return reaches 95%, except during rare instrument  
127 malfunctioning events. Mobile phone reception at Ill18 was unstable and gappy in 2018 and  
128 2019, but improved to 99%, during 2020 when the network provider was changed. Even

129 when real-time data transfer is interrupted, data are locally recorded and transmitted when  
 130 the mobile network connection is stable again.

131 Between 2017 and 2019, the seismic network recorded more than 22 debris flows. Fig-  
 132 ures 1c-f show vertical velocities seismograms and associated spectrograms of two debris  
 133 flows. Figure 1c shows the largest recorded event (vol.= 100,000 m<sup>3</sup>, velocity of the front  
 134  $v_{\text{front}}=6.67 \text{ m s}^{-1}$ ), and Figure 1e shows one of the smallest recorded events (vol.= 13,000 m<sup>3</sup>,  
 135  $v_{\text{front}}=0.95 \text{ m s}^{-1}$ ). Both debris-flow signals show emergent onsets with dominant frequen-  
 136 cies above 1 Hz reaching frequencies of 40-50 Hz. The signal emerges from the background  
 137 noise at times that depend on the distance between the debris-flow front and the recording  
 138 station (Walter et al., 2017). For the larger event (Figure 1c, d) we observe burst-like sig-  
 139 nals generated by thunder, not directly related to the debris-flow (Marchetti et al., 2019).  
 140 Seismograms generated by all events recorded on eight stations are presented in Figures  
 141 S1-S22.

### 142 3 Methods

143 We use the vertical-component seismograms of the 22 debris-flow events between 2017  
 144 and 2019 recorded on six stations (stations ILL12, ILL13, ILL14, ILL16, ILL17, ILL18  
 145 in Figure 1a) to train a machine learning model and test its detection capability. ILL15  
 146 and ILL11 were not used, because the former was deployed later in the season and the  
 147 latter is located in the Rhone Valley recording strong anthropogenic noise signals. Debris-  
 148 flow properties are shown in Table S1. The emergent nature of debris flow seismograms  
 149 precludes use of standard event detectors and instead the spectral content of the continuous  
 150 seismic signal is analyzed (Walter et al., 2017; Lai et al., 2018; Wenner et al., 2019, 2020).  
 151 We split the debris flow seismograms into 100s time windows with an overlap of 50%. This  
 152 window length is long enough to extract stable spectral characteristics and results in large  
 153 enough set of training data. The overlap is chosen to further increase the number of samples  
 154 in the data set.

#### 155 3.1 Labeled data

156 We define labels for three seismic event classes:

- 157 1. *Pre-CD1*: debris flow signals before passage of CD1
- 158 2. *Post-CD1*: debris flow signals after passage of CD1
- 159 3. *Noise*: signals not associated with debris flows

160 Dividing the debris-flow signals into two classes caters to the need to detect debris  
 161 flows in the upper catchment before CD1 passage. In the lower Illgraben part, the check  
 162 dams interrupt the flow and possibly influence particle sorting which might change debris  
 163 flow signals. For 20 events, the arrival times at CD1 are known from geophone installations  
 164 within the check dam. For 2 events geophone detections were not available and instead we  
 165 used estimates from amplitude source location (ASL), which traces the flow front location  
 166 using time-varying amplitudes of the debris flow seismograms (Walter et al., 2017). The  
 167 three different signal classes are indicated on Figures 1c-f and S23.

#### 168 3.2 Catalog compilation and processing steps

169 The construction of our debris flow detector is a supervised machine-learning classifica-  
 170 tion (Goodfellow et al., 2016), because we ask an algorithm to classify a signal of unknown  
 171 origin based on a previously trained machine learning model. Training the model requires a  
 172 labeled signal catalog with signals whose classes are known from independent observations.  
 173 We compile such a labeled data set from debris flow seismograms defined by manually picked  
 174 signal start and end times (Figure 1d, f; see TextS1 for details). Debris flow seismograms are

175 defined to be those records lying between the earliest signal start and the latest signal end  
 176 among all stations. Including all available stations, this yields 3,631 pre-CD1 time windows  
 177 and 13,046 post-CD1 time windows. We randomly choose 550 100-second long noise time  
 178 windows from 2017, 2018 and 2019 and several rainfall seismograms to compile the noise  
 179 catalog. This provides 16,614 noise time windows.

180 We use a two-iteration training and testing approach: in the first iteration we confine  
 181 ourselves to the 18 debris flows with the cleanest seismic signatures. From these we use  
 182 all 100 second time windows from 15 events with both pre-CD1 and post-CD1 labels to  
 183 train the model and test it on the seismograms from the remaining debris flows. We use  
 184 2/3 (11,076) randomly selected noise time windows for the training, and the rest (5,538)  
 185 for the testing. In the second iteration we repeat this exercise with time windows from 20  
 186 debris flows for training and 2 debris flows for testing. We furthermore inject identified false  
 187 positives (29,741 time windows) from the first iteration into the noise class. This increases  
 188 the noise class to 46,355 time windows.

### 189 *3.2.1 Detector implementation and performance*

190 Rather than using the raw seismic signals, our classification algorithm operates on 70  
 191 statistical signal features. A feature is a scalar number, which describes waveform charac-  
 192 teristics [e.g. root-mean-square amplitude (RMS), spectral content (e.g. mean and variance  
 193 of the discrete Fourier transform), and signal variations throughout the network (e.g. ratio  
 194 between maximum RMS and minimum RMS). The complete feature list is given in Table  
 195 S2 in SI and Provost et al. (2017). 59 features are extracted for each station separately.  
 196 Additional 11 network features are calculated based on all available stations.

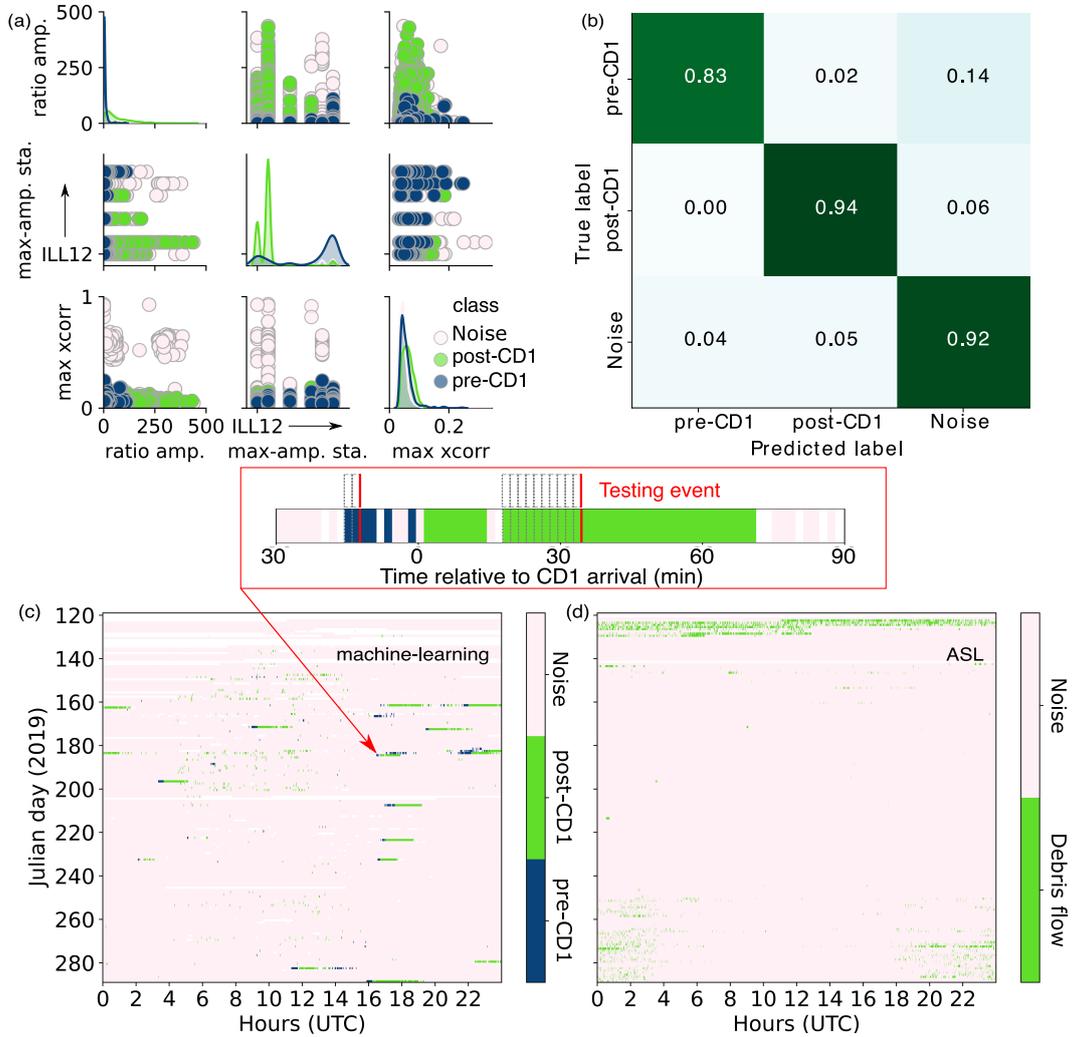
197 We use the Random Forest (RF) supervised classifier (Breiman, 2001) as the machine  
 198 learning algorithm, which comprises majority votes among an ensemble of randomized de-  
 199 cision trees. The decision trees are formed by consecutive inequality operations, which  
 200 determine if features of a signal are smaller or larger than predefined thresholds. These  
 201 thresholds, the order and the number of the inequality operations are learned during the  
 202 training phase, whereas hyperparameters (e.g. the maximum number of the inequality op-  
 203 erations and the total number of decision trees) are determined as described in Text S1.

204 RF has proven useful in seismological applications (e.g., Rouet-Leduc et al., 2017,  
 205 2019) and mass movements detection (e.g., Hibert et al., 2017; Maggi et al., 2017; Provost  
 206 et al., 2017). For our implementation we use Scikit-learn machine learning Python library  
 207 (Pedregosa et al., 2011).

208 In the training phase the machine learning algorithm has access to the features and  
 209 their associated labels (pre-CD1, post-CD1 and noise). Subsequently, the performance of  
 210 the machine learning model is evaluated on testing set, which were not included in the  
 211 training set. The true labels of the testing set are compared to the model predictions, which  
 212 may or may not be correct (Figure 2).

213 The RF algorithm returns the feature importance which elucidates how the model  
 214 reaches its predictions (Breiman, 2001). Figure 2a shows pairwise relations between the three  
 215 most important features. In each subplot two features are plotted against each other and the  
 216 univariate distributions of the same features are plotted on the diagonal with density plots.  
 217 The three most important features are network features: 1. ratio between the maximum  
 218 RMS and the minimum RMS in the network, 2. station number with maximum RMS, and  
 219 3. maximum coherence (normalized cross-correlation) between station pairs.

220 This shows that: (1) The machine-learning model strongly relies on the relative RMS  
 221 amplitudes throughout the network and the RMS amplitude ratio is the lowest for the pre-  
 222 CD1 class. (2) Some noise time windows are highly correlated. (3) ILL18 has the largest



**Figure 2.** Machine-learning model evaluation (second iteration training). (a) Pairwise relations of the three most important features (see text for details). Features from each class are marked in different colors. (b) Normalized confusion matrix with true and predicted labels (columns and rows). (c) Results of the ML-based detector and (d) ASL-based detector applied to the 2019 continuous data. Inset in (c) shows a zoom on the testing debris flow, which was not part of the training set. Gray dashed lines denote individual detections in time windows and red line shows the alarm raised after a fixed number of subsequent detections.

223 RMS for the pre-CD1 class, while ILL12 and ILL13 show the largest RMS for the post-CD1  
 224 class.

225 We use a confusion matrix (Figure 2b) and Receiver Operating Characteristic (ROC)  
 226 curve (Figure S25, in SI) to evaluate our model performance. The confusion matrix, also  
 227 called an error matrix (Stehman, 1997), assesses classification performance in a table layout  
 228 with true labels as columns and predicted labels as rows. For an ideal classifier all samples  
 229 locate on the diagonal where the predicted label equals the true label and the diagonal  
 230 values are normalized with 1.

231 We present within results of the second iteration, results of the first iteration are pre-  
 232 sented in Figures S24-S25. The confusion matrix on Figure 2b shows the highest misclassifi-  
 233 cation for the pre-CD1 class with 14 % of pre-CD1 time windows classified as noise. However,  
 234 we verified that  $\sim 30\%$  of these "confused" time windows are the first three time windows  
 235 of the pre-CD1 seismograms, and the normalized number of true positives increases to from  
 236 0.83 to 0.87 (pre-CD1—noise misclassification lowers to 0.11) if we remove these windows  
 237 from the testing set. Whereas these initial samples are labeled as pre-CD1 they might still  
 238 constitute noise for stations located further away from the torrent. Based on the scores on  
 239 the confusion matrix diagonal we expect that our detector identifies debris flow signals at  
 240 an accuracy near 90 %.

### 241 3.3 Detections and alarms

242 So far, we evaluated the performance of our machine learning model using the union of  
 243 predictions from all stations. For an operational real-time alarm system we define a detector,  
 244 which requires that more than half of the operational stations point towards the same class.  
 245 If such a majority does not exist the detector does not make a prediction. Consequently,  
 246 for real-time operation, we define "detection" and "alarm" as follows:

- 247 1. *Detection*: a single time window in which the majority vote over online stations  
 248 predicts the pre-CD1 or post-CD1 class.
- 249 2. *Alarm*:  $> 2$  subsequent detections for the pre-CD1 class, and  $> 10$  subsequent detec-  
 250 tions for the post-CD1 class.

251 If no majority exists among online stations, the detector freezes the current alarm status  
 252 and waits for the prediction from the next time window to update the alarm status. The  
 253 alarm definition introduces a time lag between an initial debris flow detection (200 s for  
 254 the pre-CD1 class and 16 min 40 s for the post-CD1 class, see the inset in Figure 2c for  
 255 a visual representation). This time lag is small for the pre-CD1 class which is crucial for  
 256 warning, and at the same time minimizes the number of false alarms. In future detector  
 257 improvements, more advanced logical operations will likely reduce the time lag between  
 258 initial detection and alarm, especially for the post-CD1 class.

## 259 4 Results

### 260 4.1 2019 continuous classification

261 Next, we run the detector over the 2019 archived data using 100 s time windows, this  
 262 time without overlaps. In 2019 we monitored station up- and down times, which we now  
 263 use to reproduce real-time station performance. The 2019 data contain 13 training events  
 264 and one, which was part of the testing set.

265 Figure 2c shows the detector performance over 170 days in 2020. As expected, the debris  
 266 flow detections (dark blue pixels for the pre-CD1 class, and green pixels for the post-CD1  
 267 class) are embedded in the noise windows (light pink). Debris flows consist of continuous

268 detections, but isolated detection "pixels" illustrate numerous noise time windows falsely  
 269 classified as debris flows.

270 We apply the alarm criterion to the 2019 detections and find that our debris flow  
 271 detector raises a pre-CD1 alarm for 11 events, including the testing event, and misses only  
 272 three small-volume events. For all 14 events a post-CD1 alarm was raised. 6 false positive  
 273 alarms were raised (1 post-CD1 class, and 5 pre-CD1 class). For comparison, the ASL-  
 274 based detector (e.g., Battaglia & Aki, 2003; Walter et al., 2017) catches 5 large debris  
 275 flows but raises false alarms on 64 days; for some days (e.g.: Julian days 123-126, 270-  
 276 280) it generates false alarms continuously (Figure 2d). Visual data inspection shows that  
 277 ASL detection tends to fail when amplified noise signals resulting from electronic spikes or  
 278 cultural activity are present on individual stations. The machine learning detector is less  
 279 sensitive to these spurious signals.

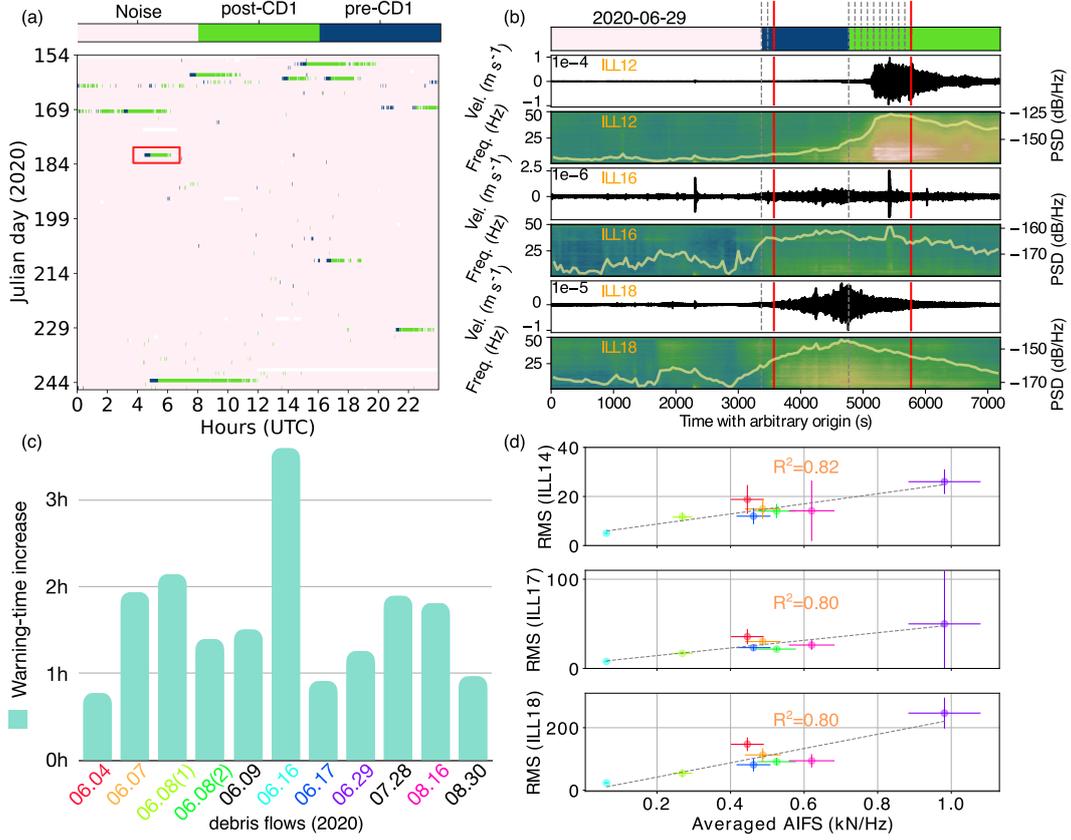
280 We stress that for this 2019 debris-flow detection comparison, the machine learning  
 281 approach is biased: The machine learning model learns 13 events in the training phase and  
 282 only one event [marked with a red arrow (Figure 2c)] is independent from the training phase.  
 283 This event happens to be missed by the ASL-detector. On the other hand, this test of the  
 284 2019 data demonstrates the drastic reduction of false alarms when moving from ASL to  
 285 machine-learning based detection.

286 Finally, we test the machine learning detector on the 2017 and 2018 data (Figure S26),  
 287 analogous to the 2019 test. The detector generates less than 3 false alarms per year and  
 288 correctly raises pre-CD1 and post-CD1 alarms for the event not included in the training  
 289 set (marked with red arrows in Figure S26). Moreover, the detector finds some previously  
 290 unknown events (Figure S27) with either pre-CD1 or post-CD1 alarms. Based on signal  
 291 strengths and characteristics, these alarms correspond to small debris flows, which did not  
 292 trigger or reach the in-torrent detection system.

## 293 4.2 2020 continuous classification

294 The final realistic and rigorous test of our machine learning detector is the real-time  
 295 classification of the 2020 data-stream. The 2020 seismic network was deployed at the end  
 296 of May 2020 and the detector has been running continuously since 2 June 2020. In the first  
 297 week of operation (3-9 June 2020) the detector correctly raised alarms for 5 debris flows  
 298 triggered by high-intensity rainfalls [cumulative rainfall over one week=52.4 mm (Swiss  
 299 Meteorological Service, Montana precipitation station)]. In total, during three months (3  
 300 June to 3 September 2020) the detector caught 13 debris flows and raised no false alarms.  
 301 Figure 3b shows an example of detections and alarms, vertical records, and spectrograms  
 302 during the initiation of the 29 June 2020 debris flow.

303 8 out of the 9 June and July 2020 alarms were independently confirmed by the WSL  
 304 observatory, although the debris-flow observation station currently undergoes maintenance  
 305 and does not provide CD1 arrival times. Nevertheless, we can compare our results with  
 306 recordings from a video camera installed at the lowermost check dam CD29 near the Rhone  
 307 River (Table S3). This comparison was not possible for events, which occurred at night  
 308 or which stopped before reaching CD29. We also estimated arrival times at a nearby seis-  
 309 mometer (ILL11) installed within a few meters from the torrent, which is not part of our  
 310 detection system. Depending on their average flow velocities, most debris flows arrived at  
 311 CD29  $\sim$  1-2 h after our pre-CD alarm times (Figure 3c, Table S3). Given typical travel  
 312 times between CD1 and CD29 of 20 minutes (Badoux et al., 2009; Walter et al., 2017), our  
 313 system therefore provides additional warning time between 20 minutes and over 1.5 hours  
 314 with respect to CD1.



**Figure 3.** Debris flow detections in 2020. (a) Results of the ML-based detector run on the continuous real-time data stream from 2020. The 29 June 2020 event is marked with a red box. (b) Detections (gray dashed lines), alarms (red lines), vertical records of seismometers and spectrograms. The top horizontal color bar shows detection type (white: noise; blue: pre-CD1; green: post-CD1). (c) Warning time gain with respect to detection at CD29. (d) Relation between signal amplitudes near CD1 and averaged apparent impact force spectra (AIFS) calculated for the lowest Illgraben stretch (Zhang et al., 2020). The horizontal error bars are taken as 10% from the averaged AIFS and the vertical error bars represent the standard deviation of RMS calculated over 10 consecutive post-CD1 detections. In Panels c and d events are indicated by the same color code.

## 5 Discussion and Conclusions

The central result of this study is that machine learning applied to real-time seismic data streams can detect debris flows in regions where conventional instrument deployment is not possible. This provides significant increases in warning times. All 8 events independently captured by the WSL debris flow observatory were detected with our approach. Several smaller debris flows, which did not reach in-torrent instrumentation but generated weak yet clear debris flow seismograms were also detected.

The performance in 2020 with no false positives or false negatives is encouraging but warrants modifications to the detector to automatically identify debris flows large enough to leave the upper catchment. This leads to the pivotal question whether our machine learning detector provides some quantitative measure of event size at the earliest alarm times, because this would allow warning against particularly destructive events. To this end we investigate if alarm-time seismic amplitudes scale with frequency-averaged apparent impact forces spectra (AIFS). The latter represent moment transfer of debris flow particles during ground impacts (Farin et al., 2019). We follow (Zhang et al., 2020) (see also Text S2) to calculate AIFS and their averages over the lowest Illgraben extent. These averages scale with boulder sizes accumulating at the flow front by the time it reaches the Rhone River near CD29 (Zhang et al., 2020) (Figure S28).

We do not find significant correlations between seismic amplitudes at the time of pre-CD1 alarms ( $R^2$  varying between stations, from 0.01 to 0.38). However, for the earliest detection time window contributing to the post-CD1 alarms, there is a clear correlation between seismic amplitudes and AIFS (Figure 3d). Not all stations correlate equally, but ILL14, ILL17, and ILL18 have an  $R^2$  of around 0.80. This shows that shortly after debris flow passage at CD1, seismic amplitudes can identify flow fronts with large boulder sizes, some 20 minutes before they arrive at CD29.

The poor correlation between seismic amplitudes during pre-CD1 alarms and AIFS raise questions about what seismogenic processes are detected at the very beginning of a debris flow. In general, initial sediment mobilization leading to debris flows may occur via pore water pressure increases or water drag forces leading to sediment failure on lateral slopes or within the torrent channel (Berti & Simoni, 2005; Godt & Coe, 2007; Gregoretto & Fontana, 2008). Our pre-CD1 detections identify time windows, when seismic amplitudes steadily increase (Figure 1d,f, 3b), rather than distinct bursts of seismic energy, which are observed in our records at other times (e.g. between 0 and 1000 s in Figure 1d). The steady increase in seismic energy argues for slow mobilization of debris flow material rather than sudden landslide failures on steep slope, which would be associated with burst-like signals.

The Illgraben site is an ideal natural laboratory to test debris flow detections, because regular event occurrence facilitates detector training. This is particularly important for machine learning algorithms relying exclusively on labeled training data. 22 training events used here can be considered a small training catalogue compared to most machine learning applications. Yet our practice to split signals into 100 second time windows increases the training data set by several orders of magnitudes to provide reliable detection. To transfer our Illgraben detector to other debris flow catchments, modifications are likely necessary to cope with fewer training events. We evaluated the accuracy of classification as a function of a number debris flow events used in a training set (Figure S29). The results show that machine-learning model trained even on a single event gives better results than a random guess, but a higher accuracy ( $> 0.7$ ) and stable predictions are obtained from 9 training events. Alternatively, it remains to be tested if the machine learning model trained at Illgraben could simply be applied in other geographic regions, i.e. if the model learned "general" characteristics of debris flow seismograms, which are independent of source-station distances and subsurface properties affecting seismic wave propagation. For machine learning algorithms applied to earthquake detection such detector transferability has already been confirmed (Ross et al., 2018).

Machine learning provides powerful tools for time series analysis and the approach presented here is only a first step to leverage this potential for natural hazard warning. Nevertheless, our relatively simple application already tackled the longstanding problem to reliably detect debris flows in an upper catchment area, which is inaccessible to existing detectors. Moreover, seismic data acquisition such as used here is a relatively cheap alternative to in-torrents instruments, which require major construction efforts. The combination of seismic monitoring and real-time data processing based on machine learning therefore provides significant advantages for Alpine mass movement detection, which have yet to be harnessed in operational warning schemes.

## Acknowledgments

Seismometer installation was funded by WSL and the Canton Valais and supported by the Swiss Military. We thank John Clinton, Roman Racine, and Stefan Wiemer; the Swiss Seismological Service and its electronic laboratory (ELAB) for technical support. The data from the Illgraben network is collected under the network code XP (<https://doi.org/10.12686/sed/networks/xp>) and all seismic data will be openly available after a 2-year embargo (in 2022) via the archives in the Swiss Seismological Service, <http://www.seismo.ethz.ch/en/research-and-teaching/products-software/waveform-data/> and the European Integrated Data Archive (EIDA), <http://www.orfeus-eu.org/data/eida/>. This work was funded by the Swiss National Science Foundation (SNSF) project Glacial Hazard Monitoring with Seismology (GlaHMSeis, grant PP00P2 157551) and Swisscom Broadcast AG. Obspy Python routines ([www.obspy.org](http://www.obspy.org)) were used to download waveforms and pre-process seismic data.

## References

- Allstadt, K. E., Matoza, R. S., Lockhart, A. B., Moran, S. C., Caplan-Auerbach, J., Haney, M. M., ... Malone, S. D. (2018). Seismic and acoustic signatures of surficial mass movements at volcanoes. *Journal of Volcanology and Geothermal Research*, *364*, 76 - 106. Retrieved from <http://www.sciencedirect.com/science/article/pii/S0377027317306261> doi: <https://doi.org/10.1016/j.jvolgeores.2018.09.007>
- Arattano, M., & Marchi, L. (2008). Systems and sensors for debris-flow monitoring and warning. *Sensors*, *8*, 2436-2452.
- Badoux, A., Andres, N., Techel, F., & C., H. (2016). Natural hazard fatalities in Switzerland from 1946 to 2015. *Nat. Hazards Earth Syst. Sci.*, *16*, 2747-2768. doi: <https://doi.org/10.5194/nhess-16-2747-2016>
- Badoux, A., Graf, C., Rhyner, J., Kuntner, R., & McArdell, B. W. (2009). A debris-flow alarm system for the Alpine Illgraben catchment: Design and performance. *Nat. Hazards*, *49*(3), 1517-1539. doi: [doi:10.1007/s11069-008-9303-x](https://doi.org/10.1007/s11069-008-9303-x)
- Battaglia, J., & Aki, K. (2003). Location of seismic events and eruptive fissures on the Piton de la Fournaise volcano using seismic amplitudes. *J. Geophys. Res.-Solid*, *108*. doi: [1029/2002JB002193](https://doi.org/10.1029/2002JB002193)
- Berger, C., McArdell, B. W., & Schlunegger, F. (2011a). Direct measurement of channel erosion by debris flows, Illgraben, Switzerland. *Journal of Geophysical Research: Earth Surface*, *116*(F1). doi: [10.1029/2010JF001722](https://doi.org/10.1029/2010JF001722)
- Berger, C., McArdell, B. W., & Schlunegger, F. (2011b). Sediment transfer patterns at the Illgraben catchment, Switzerland: Implications for the time scales of debris flow activities. *Geomorphology*, *125*(3), 421-432. Retrieved from <http://www.sciencedirect.com/science/article/pii/S0169555X10004484> doi: <https://doi.org/10.1016/j.geomorph.2010.10.019>
- Berti, M., & Simoni, A. (2005). Experimental evidences and numerical modelling of debris flow initiated by channel runoff. *Landslides*, *2*, 171-182.
- Breiman, L. (2001). Random forests. *Mach. Learn.s*, *45*(1), 5-32. doi: [10.1023/A:1010933404324](https://doi.org/10.1023/A:1010933404324)

- 418 Burtin, A., Hovius, N., & Turowski, J. (2016). Seismic monitoring of torrential and fluvial  
419 processes. *Earth Surf. Dynam.*, *4*, 285–307. doi: [https://doi.org/10.5194/esurf-4-285-](https://doi.org/10.5194/esurf-4-285-2016)  
420 [2016](https://doi.org/10.5194/esurf-4-285-2016)
- 421 Cole, S., Cronin, S., Sherburn, S., & Manville, V. (2009). Seismic signals of snow-slurry  
422 lahars in motion: 25 September 2007, Mt Ruapehu, New Zealand. *Geophys. Res. Lett.*,  
423 *36*. doi: <https://doi.org/10.1029/2009GL038030>
- 424 Costa, J. E. (1984). Physical geomorphology of debris flows. In J. E. Costa &  
425 P. J. Fleisher (Eds.), *Developments and applications of Geomorphology* (p. 268-317).  
426 Berlin: Springer.
- 427 Coviello, V., Arattano, M., Comiti, F., Macconi, P., & Marchi, L. (2019). Seismic  
428 characterization of debris flows: insights into energy radiation and implications for  
429 warning. *Journal of Geophysical Research: Earth Surface*, *124*, 1440-1463. doi:  
430 [10.1029/2018JF004683](https://doi.org/10.1029/2018JF004683)
- 431 Farin, M., Tsai, V., Lamb, M., & Allstadt, K. (2019). A physical model of the high-  
432 frequency seismic signal generated by debris flows. *Earth Surf. Process. Landforms*,  
433 *44*, 2529–2543. doi: <https://doi.org/10.1002/esp.4677>
- 434 Godt, J., & Coe, J. (2007). Alpine debris flows triggered by a 28 July 1999 thunderstorm  
435 in the central Front Range, Colorado. *Geomorphology*, *84*, 80–97.
- 436 Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press. ([http://](http://www.deeplearningbook.org)  
437 [www.deeplearningbook.org](http://www.deeplearningbook.org))
- 438 Graf, C., Badoux, A., Dufour, F., Fritschi, B., McArdell, B., Rhyner, J., . . . Nigg, U. (2007).  
439 Alarmsystem für murgangfähige Wildbäche – Beispiel Illgraben. *Wasser Energie Luft*,  
440 *99*, 119–128.
- 441 Gregoretto, C., & Fontana, G. (2008). The triggering of debris flow due to channel-bed  
442 failure in some alpine headwater basins of the Dolomites: Analyses of critical runoff.  
443 *Hydrol. Process.*, *22*, 2248–2263.
- 444 Hibert, C., Provost, F., Malet, J.-P., Maggi, A., Stumpf, A., & Ferrazzini, V. (2017).  
445 Automatic identification of rockfalls and volcano-tectonic earthquakes at the Piton de  
446 la Fournaise volcano using a random forest algorithm. *Journal of Volcanology and*  
447 *Geothermal Research*, *340*, 130-142.
- 448 Hürlimann, M., Rickenmann, D., & Graf, C. (2003). Field and monitoring data of debris-flow  
449 events in the Swiss Alps. *Can. Geotech. J.*, *40*(1), 161-175. doi: [doi:10.1139/t02-087](https://doi.org/10.1139/t02-087)
- 450 Iverson, R. (1997). The physics of debris flows. *Rev. Geophys.*, *35*(3), 245-296. doi:  
451 [10.1029/97RG00426](https://doi.org/10.1029/97RG00426)
- 452 Jakob, M., & Hungr, O. (2005). Introduction. In M. Jakob & O. Hungr (Eds.), *Debris-*  
453 *flow hazards and related phenomena* (p. 1-8). Berlin: Springer. doi: [doi:10.1007/](https://doi.org/10.1007/3-540-27129-5_7)  
454 [3-540-27129-5\\_7](https://doi.org/10.1007/3-540-27129-5_7)
- 455 Johnson, C. G., Kokelaar, B. P., Iverson, R. M., Logan, M., LaHusen, R. G., & T.,  
456 G. J. M. N. (2012). Grain-size segregation and levee formation in geophysical  
457 mass flows. *Journal of Geophysical Research*, *117*(5A), F01032. doi: [10.1029/](https://doi.org/10.1029/2011JF002185)  
458 [2011JF002185](https://doi.org/10.1029/2011JF002185)
- 459 Kean, J., Staley, D., Leeper, R., Schmidt, K., & Gartner, J. (2012). A low-cost method to  
460 measure the timing of postfire flash floods and debris flows relative to rainfall. *Water*  
461 *Resour. Res.*, *48*. doi: <https://doi.org/10.1029/2011WR011460>
- 462 Lai, V. H., Tsai, V. C., Lamb, M. P., Ulizio, T. P., & Beer, A. R. (2018). The Seismic Sig-  
463 nature of Debris Flows: Flow Mechanics and Early Warning at Montecito, California.  
464 *Geophysical Research Letters*, *45*(11), 5528-5535. doi: [10.1029/2018GL077683](https://doi.org/10.1029/2018GL077683)
- 465 Maggi, A., Ferrazzini, V., Hibert, C., Beauducel, F., Boissier, P., & Amemoutou, A. (2017).  
466 Implementation of a multistation approach for automated event classification at piton  
467 de la fournaise volcano. *Seismological Research Letters*, *88*(3), 878-891.
- 468 Marchetti, E., Walter, F., Barfucci, G., Genco, R., Wenner, M., Ripepe, M., . . . Price,  
469 C. (2019). Infrasonic array analysis of debris flow activity and implication for early  
470 warning. *Journal of Geophysical Research: Earth Surface*, *124*(2), 567-587. doi:  
471 [10.1029/2018JF004785](https://doi.org/10.1029/2018JF004785)
- 472 McArdell, B. W., Bartelt, P., & Kowalski, J. (2007). Field observations of basal forces and

- 473 fluid pore pressure in a debris flow. *Geophys. Res. Lett.*, *34*, L07406. doi: doi:10.1029/  
474 2006GL029183
- 475 McCoy, S., Kean, J. W., Coe, J. A., Staley, D. M., Wasklewicz, T. A., & Tucker, G. E.  
476 (2010). Evolution of a natural debris flow: In situ measurements of flow dynamics,  
477 video imagery, and terrestrial laser scanning. *Geology*, *38*, 735-738. doi: doi:10.1130/  
478 G30928.1
- 479 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duch-  
480 esnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine*  
481 *Learning Research*, *12*, 2825-2830.
- 482 Pierson, T. C. (1986). Flow behavior of channelized debris flows, Mount St. Helens, Wash-  
483 ington. In A. D. Abrahams (Ed.), *Hillslope processes* (p. 269-296). Boston: Allen and  
484 Unwin.
- 485 Provost, F., Hibert, C., & Malet, J.-P. (2017). Automatic classification of endogenous  
486 landslide seismicity using the Random Forest supervised classifier. *Geophysical Research*  
487 *Letters*, *44*(1), 113. doi: 10.1002/2016GL070709
- 488 Rickenmann, D., Hürlimann, M., Graf, C., Näf, D., & Weber, D. (2001). Murgang-  
489 Beobachtungsstationen in der Schweiz. *Wasser, Energie, Luft*, *93*, 1-8.
- 490 Ross, Z. E., Meier, M.-A., Hauksson, E., & Heaton, T. H. (2018). Generalized seismic  
491 phase detection with deep learning. *Bulletin of the Seismological Society of America*,  
492 *44*(5A), 2894-2901. doi: doi.org/10.1785/0120180080
- 493 Rouet-Leduc, B., Hulbert, C., Lubbers, N., Barros, K., Humphreys, C. J., & Johnson,  
494 P. A. (2017). Machine learning predicts laboratory earthquakes. *Geophysical Research*  
495 *Letters*, *44*(18), 9276-9282. doi: 10.1002/2017GL074677
- 496 Rouet-Leduc, B., Hulbert, C., Lubbers, N., Barros, K., Humphreys, C. J., & Johnson, P. A.  
497 (2019). Continuous chatter of the cascadia subduction zone revealed by machine  
498 learning. *Nature Geosci.*, *12*, 75-79. doi: 10.1038/s41561-018-0274-6
- 499 Schimmel, A., Hübl, J., McArdell, B. W., & Walter, F. (2018). Automatic Identification of  
500 Alpine Mass Movements by a Combination of Seismic and Infrasound Sensors. *Sensors*,  
501 *18*(5), 1658.
- 502 Schlunegger, F., Badoux, A., McArdell, B. W., Gwerder, C., Schnydrig, D., Rieke-Zapp, D.,  
503 & Molnar, P. (2009). Limits of sediment transfer in an alpine debris-flow catchment,  
504 Illgraben, Switzerland. *Quat. Sci. Rev.*, *28*, 1097-1105. doi: doi:10.1016/j.quascirev  
505 .2008.10.025
- 506 Stähli, M., Sättele, M., Huggel, C., McArdell, B., Lehmann, P., Van Herwijnen, A., ...  
507 Springman, S. (2015). Monitoring and prediction in early warning systems for rapid  
508 mass movements. *Natural Hazards and Earth System Science*, *15*, 905-917.
- 509 Stehman, S. V. (1997). Selecting and interpreting measures of thematic classifica-  
510 tion accuracy. *Remote Sensing of Environment*, *62*(1), 77 - 89. Retrieved from  
511 <http://www.sciencedirect.com/science/article/pii/S0034425797000837> doi:  
512 [https://doi.org/10.1016/S0034-4257\(97\)00083-7](https://doi.org/10.1016/S0034-4257(97)00083-7)
- 513 Walter, F., Burtin, A., McArdell, B. W., Hovius, N., Weder, B., & M, T. J. (2017).  
514 Testing seismic amplitude source location for fast debris-flow detection at Illgraben,  
515 Switzerland. *Natural Hazards and Earth System Sciences*, *17*(6), 939-955. doi:  
516 <https://doi.org/10.5194/nhess-17-939-2017>
- 517 Wenner, M., Hibert, C., & Walter, F. (2020). Automatic near real-time classification of  
518 seismic signals in for slope failure detection in alpine environments. *submitted to Nat.*  
519 *Hazards Earth Syst. Sci.*
- 520 Wenner, M., Walter, F., McArdell, B., & Farinotti, D. (2019). Deciphering debris-flow  
521 seismograms at Illgraben, Switzerland. In K. J. W., C. J. A., S. P. M., & G. B. K.  
522 (Eds.), *Association of environmental and engineering geologists special publication:*  
523 *Vol. 28. debris-flow hazards mitigation: mechanics, monitoring, modeling, and as-*  
524 *essment* (p. 222-229). Colorado School of Mines: Association of Environmental and  
525 Engineering Geologists.
- 526 Zhang, Z., Walter, F., McArdell, B. W., Wenner, M., Chmiel, M., & He, S. (2020). Ex-  
527 tracting dynamics of debris flows from their seismic signature. *GRL, in submission.*

# Supporting Information for ”Machine Learning improves warning systems of debris flows”

Małgorzata Chmiel<sup>1</sup>, Fabian Walter<sup>1</sup>, Michaela Wenner<sup>1,2</sup>, Zhen Zhang<sup>1,3,4</sup>,

Brian W. McArdell<sup>2</sup>, Clement Hibert<sup>5</sup>

<sup>1</sup>Laboratory of Hydraulics, Hydrology and Glaciology, ETH Zürich, Zürich, Switzerland

<sup>2</sup>Swiss Federal Institute for Forest, Snow and Landscape Research, Zürich, Switzerland

<sup>3</sup>Key Laboratory of Mountain Hazards and Surface Process, Institute of Mountain Hazards and Environment, Chinese Academy of Sciences, Chengdu, China

<sup>4</sup>University of Chinese Academy of Sciences, Beijing, China

<sup>5</sup>Institut de Physique du Globe de Strasbourg, CNRS UMR 7516, University of Strasbourg/EOST, 7 8 Strasbourg, France

## Contents of this file

1. Text S1 to S2
2. Figures S1 to S29
3. Tables S1 and S3

## Additional Supporting Information (Files uploaded separately)

1. Captions for large Table S2.

## Introduction

This file contains supplement text, tables, and figures for the manuscript ”Machine Learning improves warning systems of debris flows”. Supplemental Text1 provides a de-

---

tailed description of methods used for the development of the machine-learning detector. Supplemental Text2 provides additional information on the apparent total impact force spectrum (AIFS) calculation. Supplemental Figures S1 to S22 show vertical-component seismograms of 22 debris flow events recorded in 2017, 2018, and 2019. Supplemental Figure S23 shows debris-flow seismograms of 21 June 2019 event, and corresponding spectrograms. Supplemental Figure S24 shows machine-learning model evaluation for the first iteration. Supplemental Figure S25 shows receiver operating characteristic (ROC) curve analysis for the first and the second iteration. Supplemental Figure S26 shows the results of the detector run (from the second iteration) over 2017 and 2018 continuous data. Supplemental Figure S27 shows seismograms of a small debris flow event newly caught by the detector. Supplemental Figure S28 shows four photos of Illgraben debris-flow events with boulder-rich fronts detected in 2020. Supplemental Figure S29 shows results of a sensitivity test (accuracy of prediction vs number of debris-flow events used in a training set). Supplemental Figure S30 shows vertical apparent total impact force spectra (AIFS) for 2020 debris flows.

Supplemental Table S1 provides detailed information on the characteristics of 22 debris flows recorded in 2017, 2018, and 2019. Supplemental Table S2 (uploaded separately) shows 70 statistical features used as input in the machine-learning model. Supplemental Table S3 provides information on the characteristics of 13 debris flow events detected in 2020.

## **Text S1. Description of methods used in the development of the machine-learning debris flow detector.**

Three classes of seismic signals are defined: pre-Check Dam 1 (CD1), post-CD1, and noise. We divide the dataset into 100 s time windows with an overlap of 50%. We find that 100s time window is long enough to contain enough of information to reliably extract signal statistics, and it is short enough to create a large catalog needed for machine learning.

We use a two-iteration training and testing approach, where in the first iteration we use a subset of 18 DF events with the cleanest seismic signature (15 for training and 3 for testing), and in the second iteration we use all 22 events (20 for training and 2 for testing) to increase the amount of information used to train the machine-learning model. In the second iteration we also inject false positives detections from the first iteration into the noise class. We use Scikit-learn (Pedregosa et al., 2011) implementation of the Random Forest classifier. When we refer to the machine-learning detector we refer to the overall system that involves pre-processing the data, extracting the features, and applying the machine-learning model to the new data.

We follow the following processing steps:

### **1. Pre-processing of the seismic data**

The data is detrended and centred, and it is not corrected for instrument response to increase the computational performance. Further, the instrumental response is flat in the 1-50 Hz frequency range. The only pre-processing we apply to data is a low-cutoff filter at 1 Hz to focus on high frequency signals.

### **2. Catalog preparation**

To extract debris-flow seismograms we manually pick the start time and the end time of the debris-flow events. The picks are based on Power Spectral Density (PSD) (Figure 1d,f in the manuscript) averaged over 1-50 Hz and root-mean squared (RMS) amplitudes calculated for each station and each event. An example of picks that are delimiting the three different classes is schematically presented on Figure S23. We choose the earliest and the latest picked times for each event as start times and end times making the picks uniform for all stations (Figure S23).

We then randomly choose 550 100-second long noise time windows from 2017, 2018 and 2019 and we manually select 41 rainfall events with a fixed duration of 30 min to compile the noise catalog. For the latter, we use rainfall data from the nearby precipitation station located in Montana, Switzerland (Swiss Meteorological Service). We choose long-duration rain events that do not overlap with the picked debris flow events (at least 45 min time difference between the debris flow events and the rainfall events). We define rain events as rainfall lasting at least 30 min with 1h moving (step=10 min) average higher than 1.4 mm. Note that in the second iteration the noise class contains also the false positive detections from the first iteration.

For the first iteration in total we obtain: 2,966 pre-CD1, 9,886 post-CD1, and 16,614 noise time windows in the catalog. For the second iteration we obtain: 3,631 pre-CD1, 13,046 post-CD1, and 46,355 noise time windows in the catalog. The number of time-windows per class in the second iteration is unbalanced, therefore we account for the non-balanced dataset in the training procedure (see details below).

### **3. Extracting statistical features of the seismic data**

For statistical feature calculation, we follow works of Hibert et al. (2017); Maggi et al. (2017); Provost, Hibert, and Malet (2017). They defined a set of significant features that proved efficient in classification of other types of mass movements, including landslides and rockfalls. In total we use 70 features based on the waveform, spectral, spectrogram, and network attributes (see Table S2 for the list of features). We slightly modify the list of features presented in Provost et al. (2017) by adding the following features: (1) waveform features: RMS, and interquartile range, (2) network features: mean Wasserstein distance, and standard deviation of the Wasserstein distance calculated between stations in the frequency band of (1-50) Hz. We use the stations separately meaning that for each station and for each 100s time window we calculate 59 individual attributes, and 11 network attributes are shared between the stations. We also tried different variations of combining the features and the stations (e.g., combining features from 4 stations or averaging features over stations), but treating the stations separately gives the best results. Finally, we eliminate time windows from the catalog when only one station was online. See Table S2 for the complete list of features, and for the formulas please refer to Provost et al. (2017).

#### 4. **Train-test dataset split**

There are two ways we might approach the splitting of the dataset: (1) by the number of time windows, and (2) by the number of events. We use the first approach to optimize the hyperparameters of the machine-learning algorithm and the second one to assess the accuracy of the classification and the performance of the machine-learning model. Hyperparameters are settings of a machine learning algorithm used to control its behavior (Goodfellow et al., 2016).

In (1) we take the whole dataset as ensemble without considering debris-flow event separation. However, we need (2) to reliably and independently from the hyperparameter optimization estimate the accuracy of the classification.

- Hyperparameter optimization

We use randomly chosen 75% of the dataset as a training set for the hyperparameter optimization. The hyperparameters are optimized by a 5-fold cross validation with a grid search, which allows for a full search over specified parameter values. The performance of a machine-learning model for different hyperparameter values is measured with accuracy, which is the proportion of sample for which the model produces the correct output.

We tune the following hyperparameters of the random forest algorithm: the number of trees in the forest, the minimum number of time windows required at node leaf, type of class balancing, the maximum depth of a tree, the minimum number of time windows required for an internal node split. The last two hyperparameters control the depth of trees. We use "balanced subsample" random forest mode to minimize a potential bias related to imbalanced classes. It automatically adjusts class weights to be inversely proportional to class frequencies in each tree.

- Optimal dataset split per debris-flow event

Next, we search for the optimal number and combination of debris-flow events used for training and testing. We first choose a subset of 8 debris-flow events with different characteristics (e.g., year of occurrence, volume, and waveforms). From these we choose 2-5 element event combinations that are used in the testing set. For each combination the testing set is independent from the training set. We use the accuracy score for each training-testing combination to assess the performance of the model. In each test we use

the optimized hyperparameters of the random forest algorithm. The accuracy score is normalized and its values are bounded between 0 and 1. The best subset reaches accuracy score of  $> 0.9$ , meaning that more than 90% of labels are predicted correctly.

The debris-flow events used in the final testing set are marked in orange in Table S1 (3 events for the first iteration), and in green (2 events for the second iteration). For the noise class, we randomly select 1/3 of the samples for the testing process. In overall we use 80 % of the dataset for training and 20 % for testing.

### 5. Building and evaluating the machine-learning model

The model is then built in a training phase in which the machine-learning algorithm has access to seismic features and the corresponding classes which are used as labels. To evaluate the model we apply it to the testing set with restricted access to the features. We then compare the predicted labels to the true labels. The accuracy score might be misleading for the imbalanced dataset, so to better assess the performance of the model we calculate a confusion matrix and Receiver Operating Characteristic (ROC) curves with cross validation (5 different splits).

The confusion matrix from the first iteration is presented in Figure S24a. It indicates that the machine learning model performs less well in classification of pre-CD1 and noise signals with a score of 0.73.

Figure S25 shows the ROC curves for the three classes. ROC curves shows the true positive rate ( $TPR=TP/(TP+FN)$ ) on the y axis and false positive rate ( $FPR=FP/(FP+TN)$ ) on the x axis, where: TP=true positive, FN=false negatives, FP=false positives, TN=true negatives.

The higher the area under the curve (AUC), the better the accuracy of the classifier. The ROC curve is computed for each class separately by pairwise comparison (one class vs. all other classes) and the dotted line shows the baseline for a random guess. In general, the ROC curves shows good classification results, strongly above the random baseline with AUC values of 0.9, and with low standard deviation values meaning that the classifier output should not be too much affected by changes in the training data (Pedregosa et al., 2011). Moreover, the ROC curves indicate the improvement of the performance of the machine-learning model in the second iteration.

Figure S24b shows pairwise relationships between the three most important features. The three most important features belong to the network attributes: 1. ratio between the maximum RMS and the minimum RMS in the network, 2. station number with maximum RMS, and 3. maximum coherence between station pairs. The three most important features are the same in the first and in the second iteration.

## **6. Applying the machine-learning model to new data**

We now apply the machine-learning model to 2019 continuous data. The results of the first iteration (Figure S24c) indicate that the model often classifies anthropogenic noise as post-CD1 class. This is interesting since it is not represented in the confusion matrix, hinting that the classification accuracy estimated over limited testing set might not fully represent the classification accuracy over the entire seismic dataset. These false positives are related to the presence of anthropogenic noise in Illgraben area because they appear only during week days, e.g., Monday-Friday:154-158, 182-186, 189-193 Julian days between ~5 a.m. UTC (7 a.m. local time) and 10 a.m-11 a.m. until 15 (17 local time). The false

positives are consequently injected into the noise class in the second iteration to teach the model how to distinguish between this anthropogenic noise and debris-flow signals.

One can ask how many debris flow events are needed to train the machine-learning model and obtain a good performance in classification. Figure S29 shows a sensitivity test: balanced accuracy score [the average of recall  $TP/(TP+TN)$  obtained on each class] as a function of number of debris-flow events used in training set with cross-validation (5 folds). The results show that using even a single debris-flow event for training gives better results than a random guess, although higher values of balanced accuracy ( $> 0.7$ ) and stable prediction are obtained from 9 debris-flow events used in the training set.

### **Text S2. Apparent total impact force spectrum (AIFS) calculation**

The vertical apparent total impact force spectrum (AIFS) of the debris flow events is calculated following Zhang et al. (2020). The peaks of seismic signals recorded at stations ILL11 and ILL2 (pink and yellow triangle on Figure 1) are used to estimate the average velocity of the debris flows. To calculate the AIFS, we assume that the velocity of the debris flow is constant during run-out and that the events originate upstream from CD1. The average debris flow velocities used in the AIFS calculations are presented in Table S3.

Two large peaks in the vertical AIFS at around 2,100 m and 3,100 m visible on Figure S30 are probably caused by large CDs located at these channel locations (Zhang et al., 2020). The AIFS is lower at 2,500 m, which might be related to a denser distribution of check dams at 2,500 m. If large particles (e.g., boulders) are present in a debris flow, they gradually gather at the flow front due to the particle sorting phenomena which corresponds

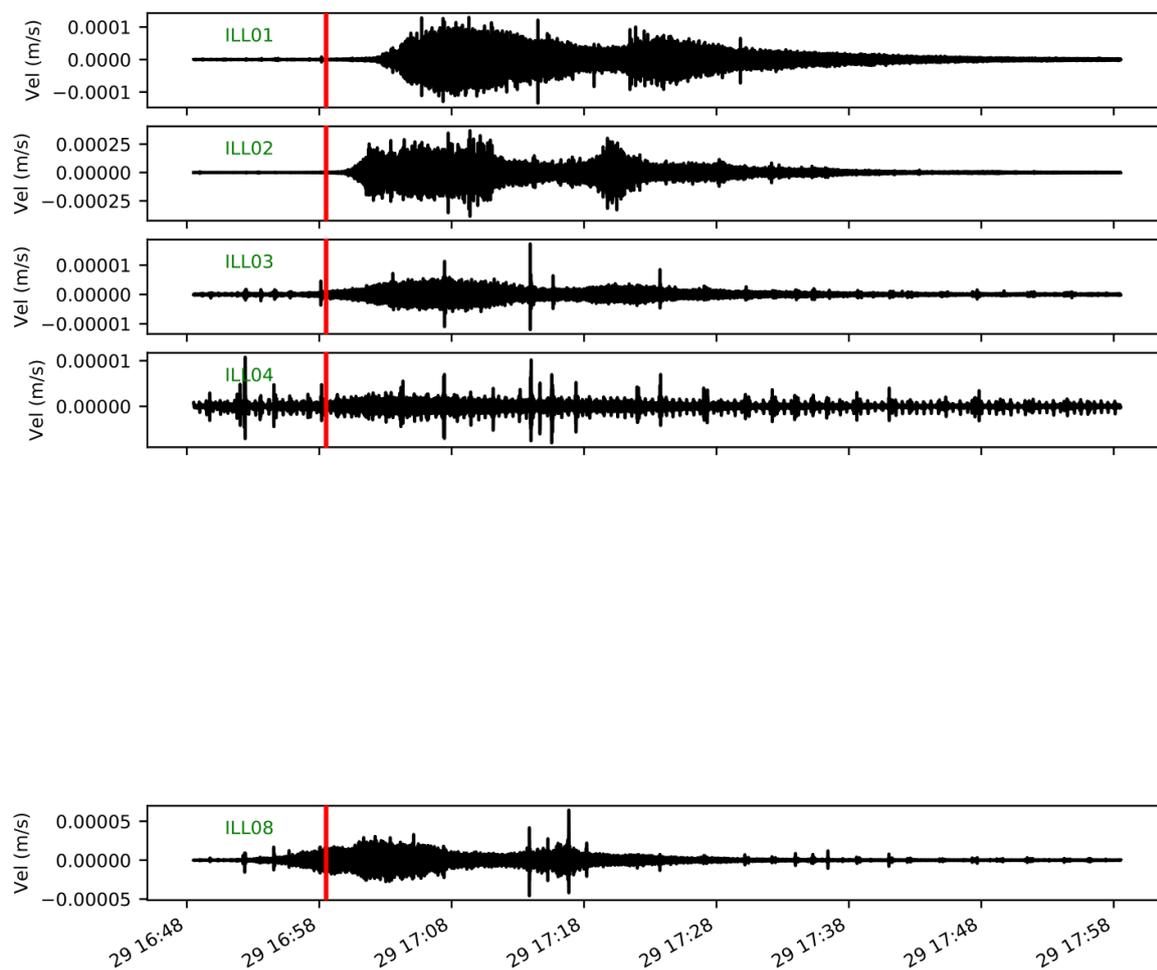
to an increase in AIFS (events on 4, 7, 8(2), 17, 29 June and 8 August 2020) . Other events seem to have lower volumes and strong deposits.

## References

- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep learning*. MIT Press. (<http://www.deeplearningbook.org>)
- Hibert, C., Provost, F., Malet, J.-P., Maggi, A., Stumpf, A., & Ferrazzini, V. (2017). Automatic identification of rockfalls and volcano-tectonic earthquakes at the Piton de la Fournaise volcano using a random forest algorithm. *Journal of Volcanology and Geothermal Research*, *340*, 130-142.
- Maggi, A., Ferrazzini, V., Hibert, C., Beauducel, F., Boissier, P., & Amemoutou, A. (2017). Implementation of a multistation approach for automated event classification at piton de la fournaise volcano. *Seismological Research Letters*, *88*(3), 878-891.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., . . . Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, *12*, 2825–2830.
- Provost, F., Hibert, C., & Malet, J.-P. (2017). Automatic classification of endogenous landslide seismicity using the Random Forest supervised classifier. *Geophysical Research Letters*, *44*(1), 113. doi: 10.1002/2016GL070709
- Schlunegger, F., Badoux, A., McArdell, B. W., Gwerder, C., Schnydrig, D., Rieke-Zapp, D., & Molnar, P. (2009). Limits of sediment transfer in an alpine debris-flow catchment, Illgraben, Switzerland. *Quat. Sci. Rev.*, *28*, 1097-1105. doi: doi:10.1016/j.quascirev.2008.10.025
- Walter, F., Burtin, A., McArdell, B. W., Hovius, N., Weder, B., & M, T. J. (2017). Testing seismic amplitude source location for fast debris-flow detection at Illgraben, Switzerland. *Natural Hazards and Earth System Sciences*, *17*(6), 1939-1955. doi: <https://doi.org/10.5194/nhess-17-1939-2017>

Zhang, Z., Walter, F., McArdell, B. W., Wenner, M., Chmiel, M., & He, S. (2020). Extracting dynamics of debris flows from their seismic signature. *GRL, in submission*.

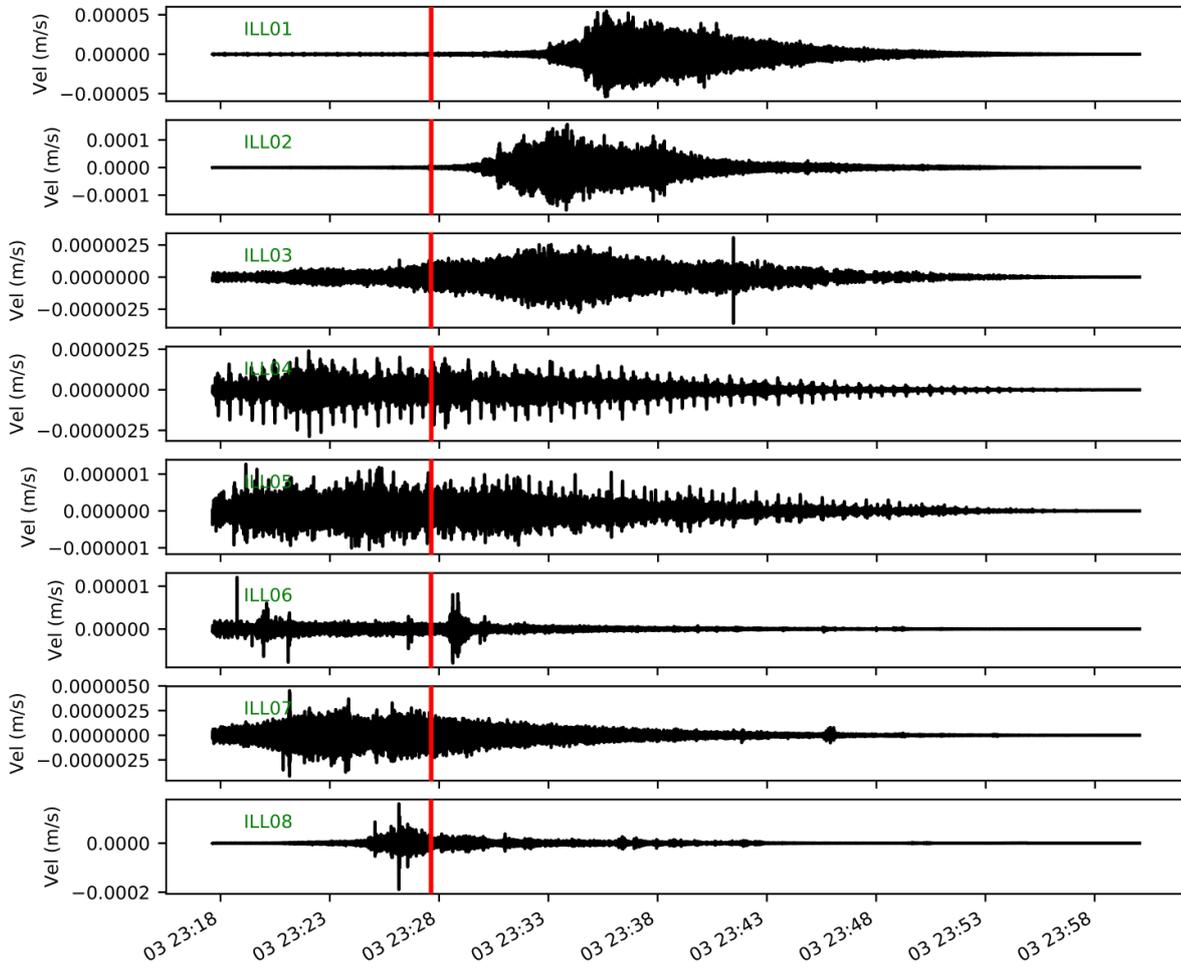
29.05.17, 16:58:31



**Figure S1.** (a) Vertical-component seismograms generated by a debris-flow event on 29 May 2017. The arrival time of the debris flow front at CD1 is marked in red.

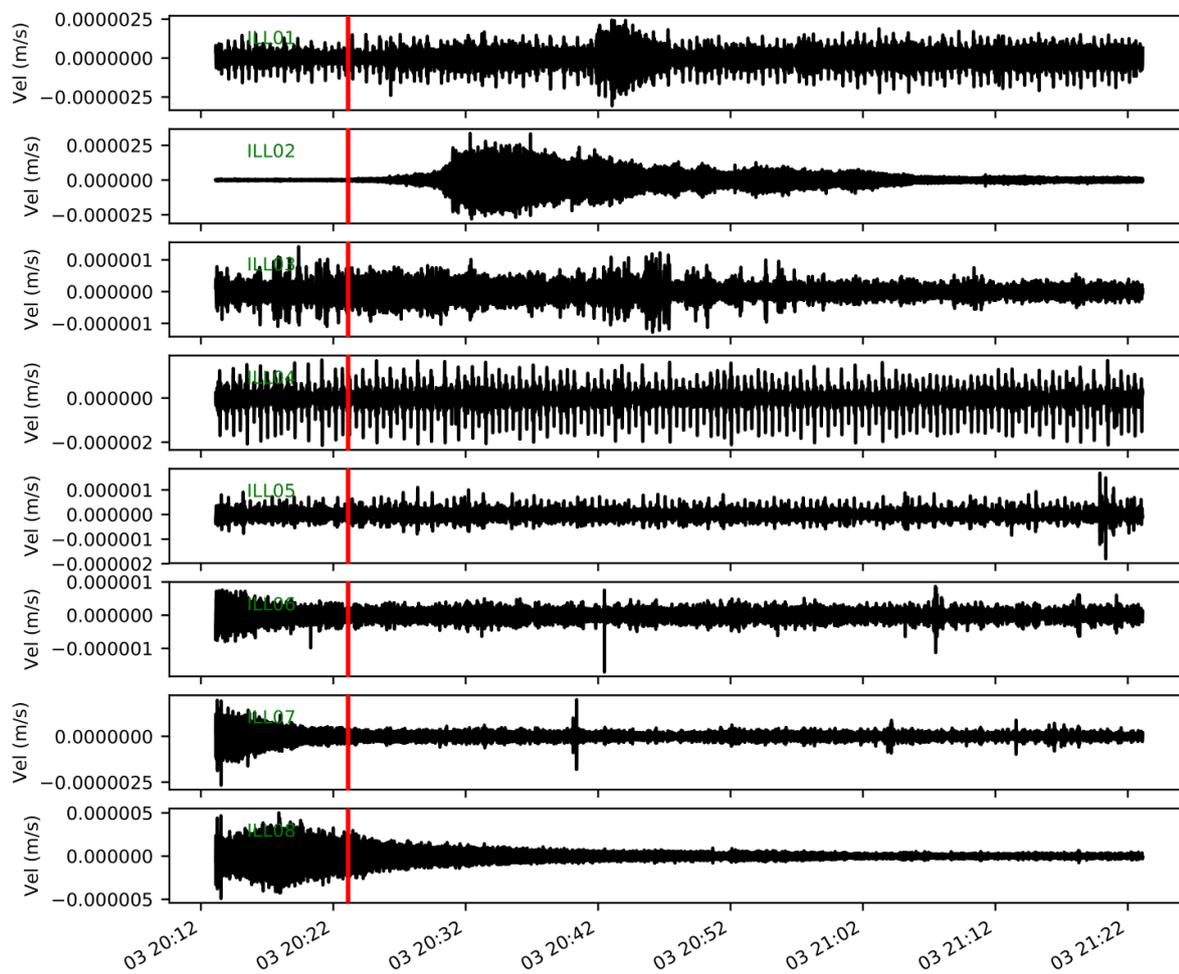
September 17, 2020, 3:01pm

03.06.17, 23:27:38



**Figure S2.** (a) Vertical-component seismograms generated by a debris-flow event on 3 June 2017. The arrival time of the debris flow front at CD1 is marked in red.

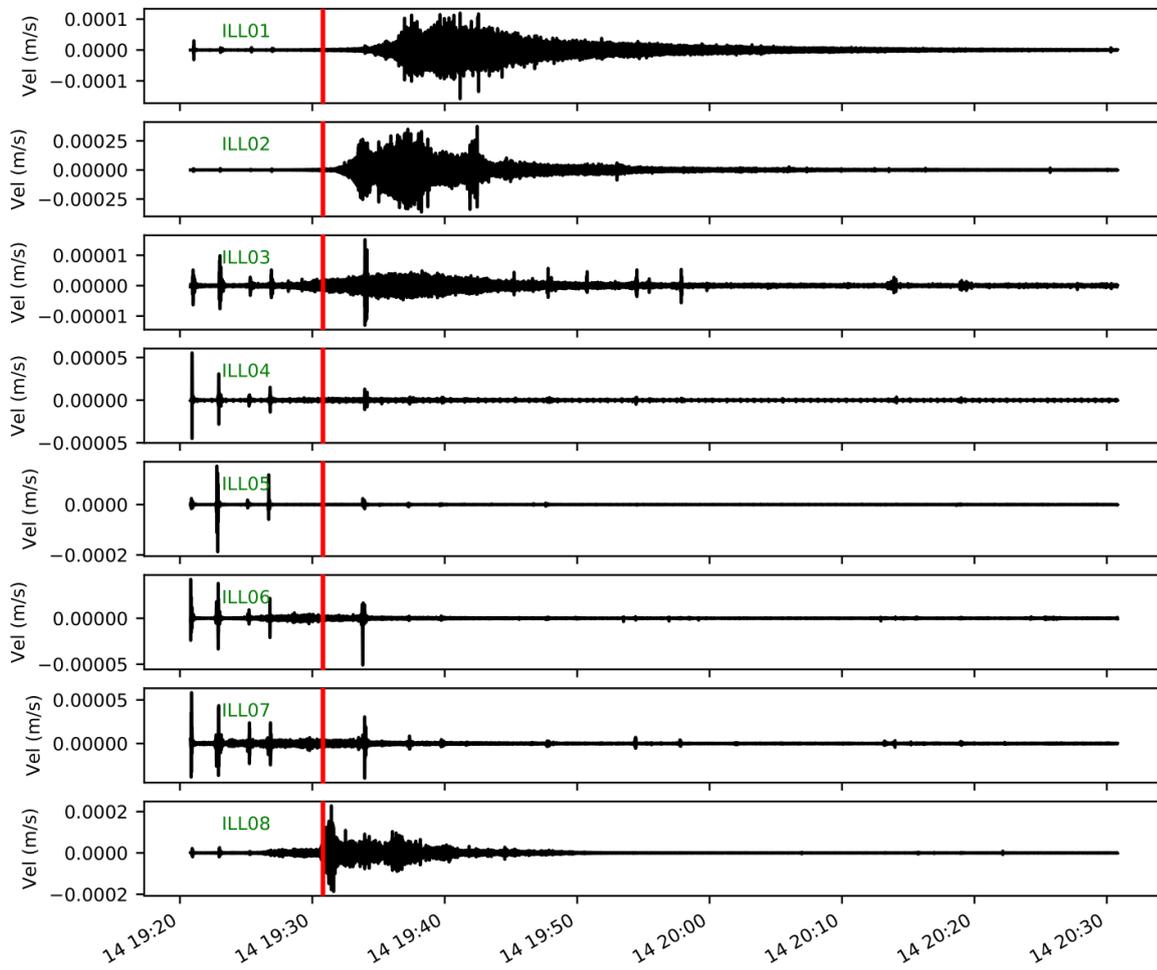
03.06.17, 20:23:07



**Figure S3.** (a) Vertical-component seismograms generated by a debris-flow event on 3 June 2017. The arrival time of the debris flow front at CD1 is marked in red.

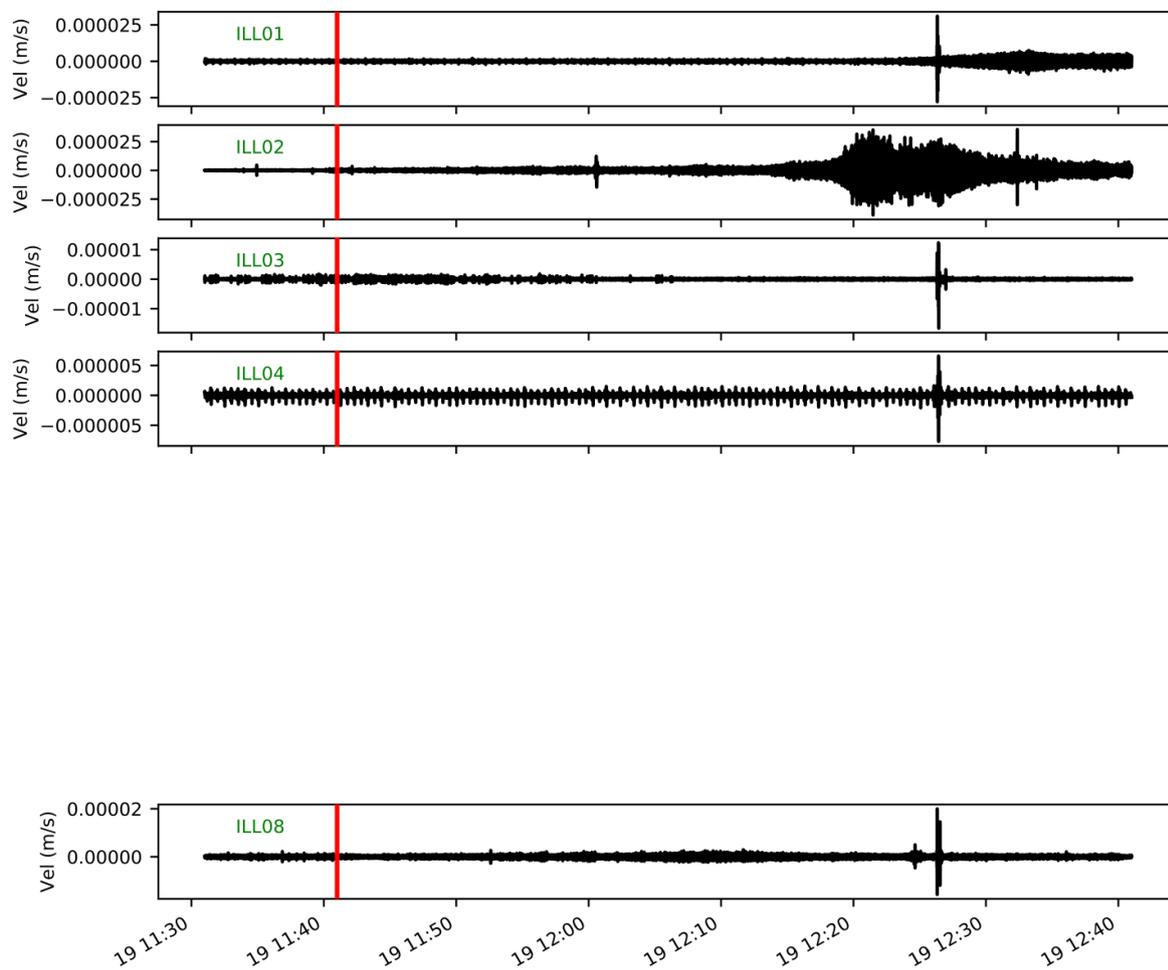
September 17, 2020, 3:01pm

14.06.17, 19:30:48



**Figure S4.** (a) Vertical-component seismograms generated by a debris-flow event on 14 June 2017. The arrival time of the debris flow front at CD1 is marked in red.

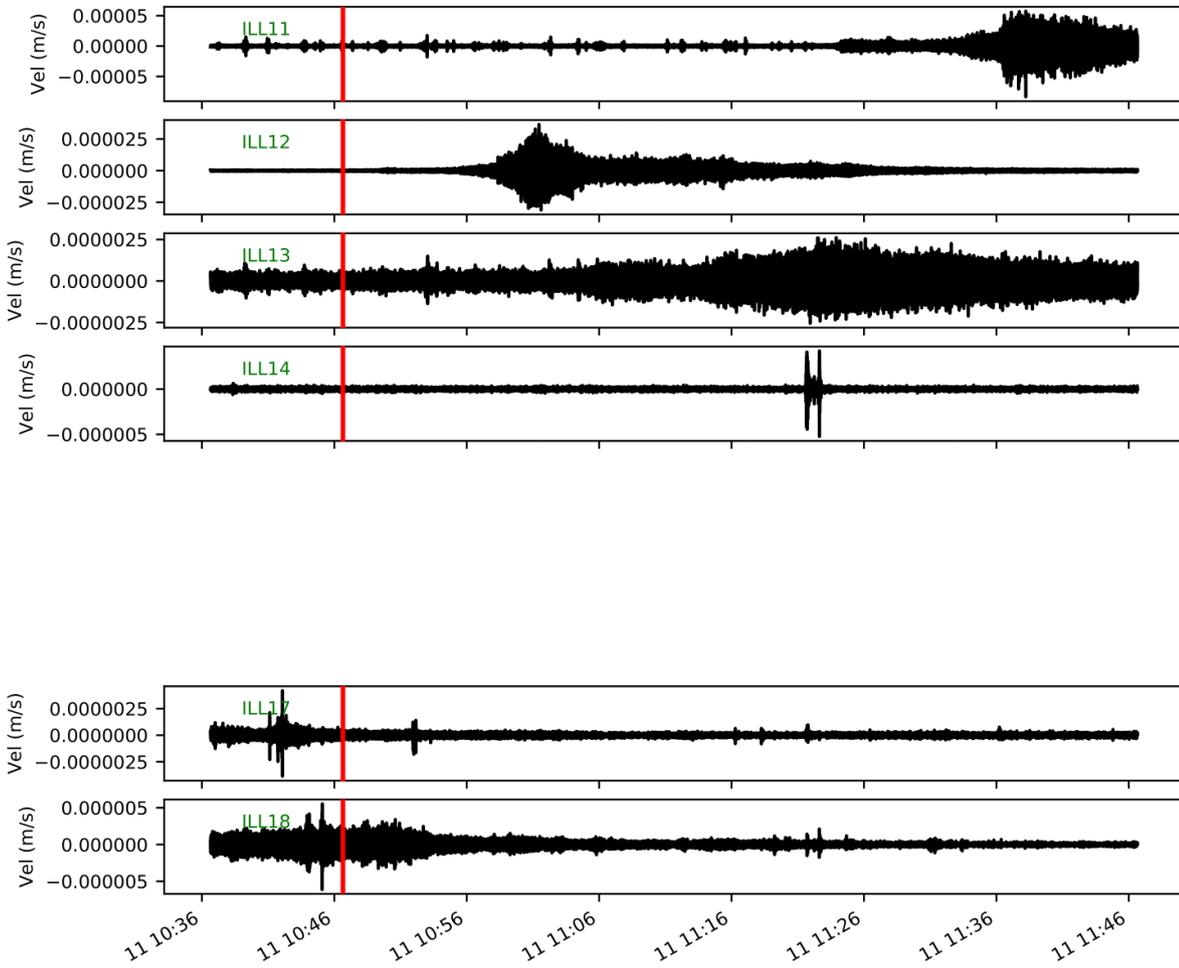
19.05.17, 11:41:00



**Figure S5.** (a) Vertical-component seismograms generated by a debris-flow event on 19 May 2017. The arrival time of the debris flow front at CD1 is marked in red.

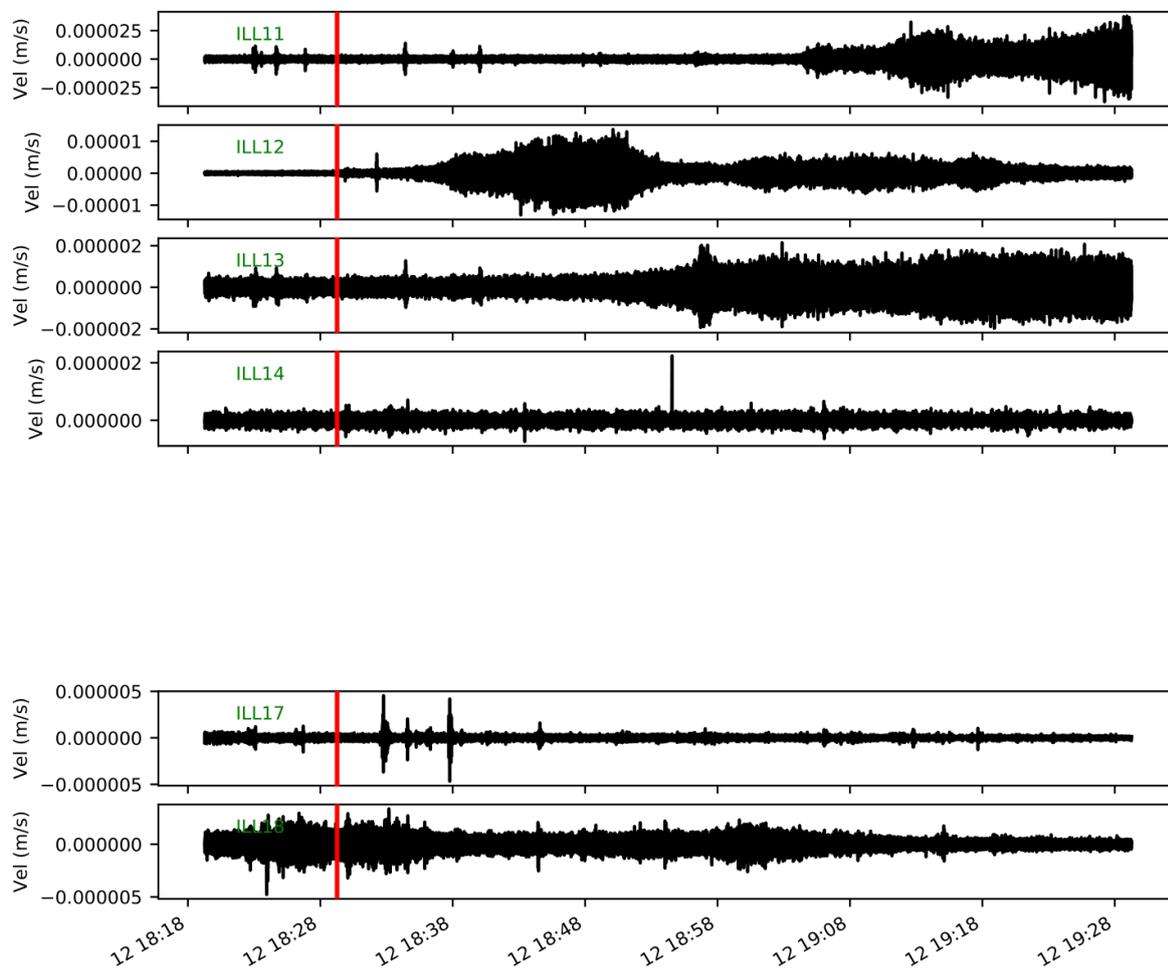
September 17, 2020, 3:01pm

11.06.18, 10:46:39



**Figure S6.** (a) Vertical-component seismograms generated by a debris-flow event on 11 June 2018. The arrival time of the debris flow front at CD1 is marked in red.

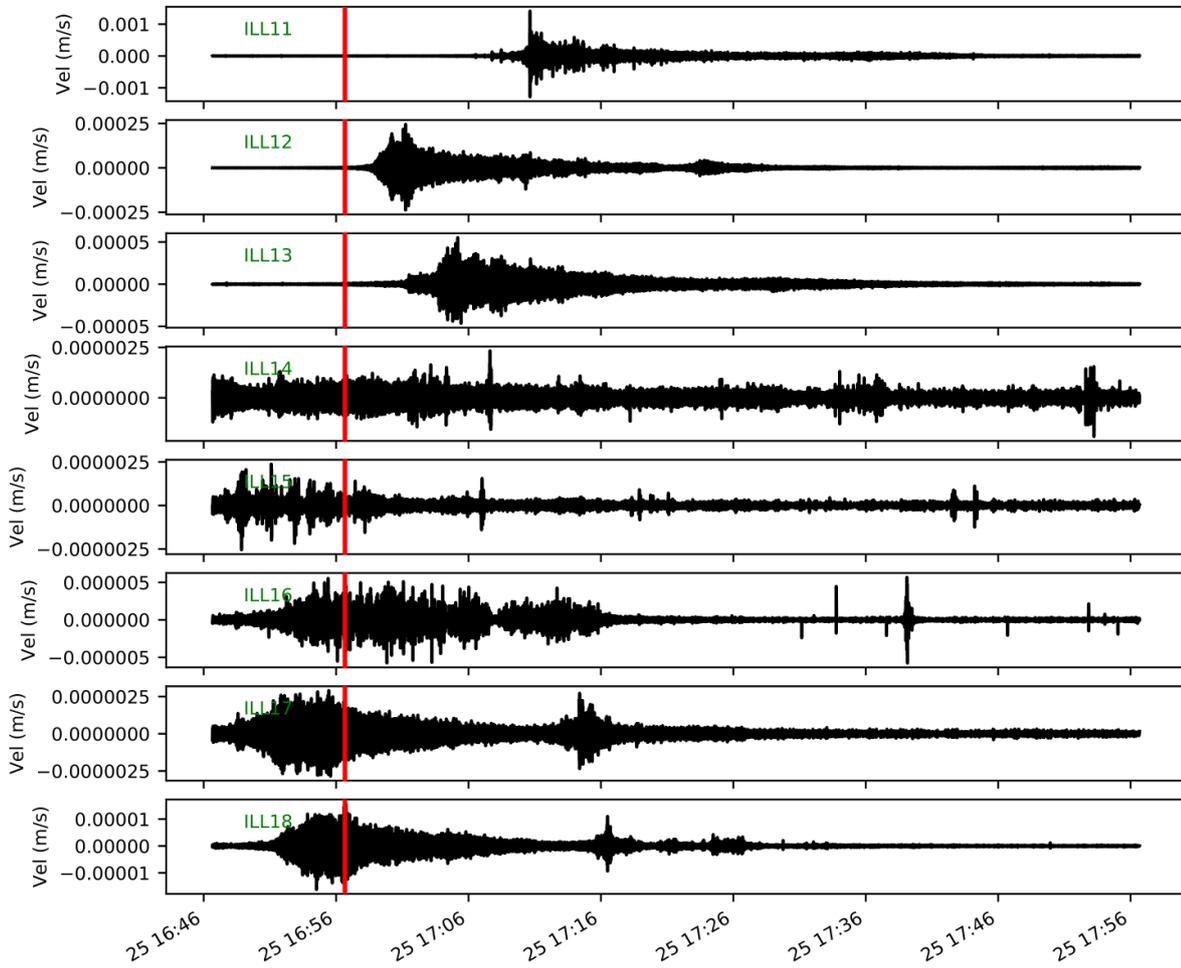
12.06.18, 18:29:16



**Figure S7.** (a) Vertical-component seismograms generated by a debris-flow event on 12 June 2018. The arrival time of the debris flow front at CD1 is marked in red.

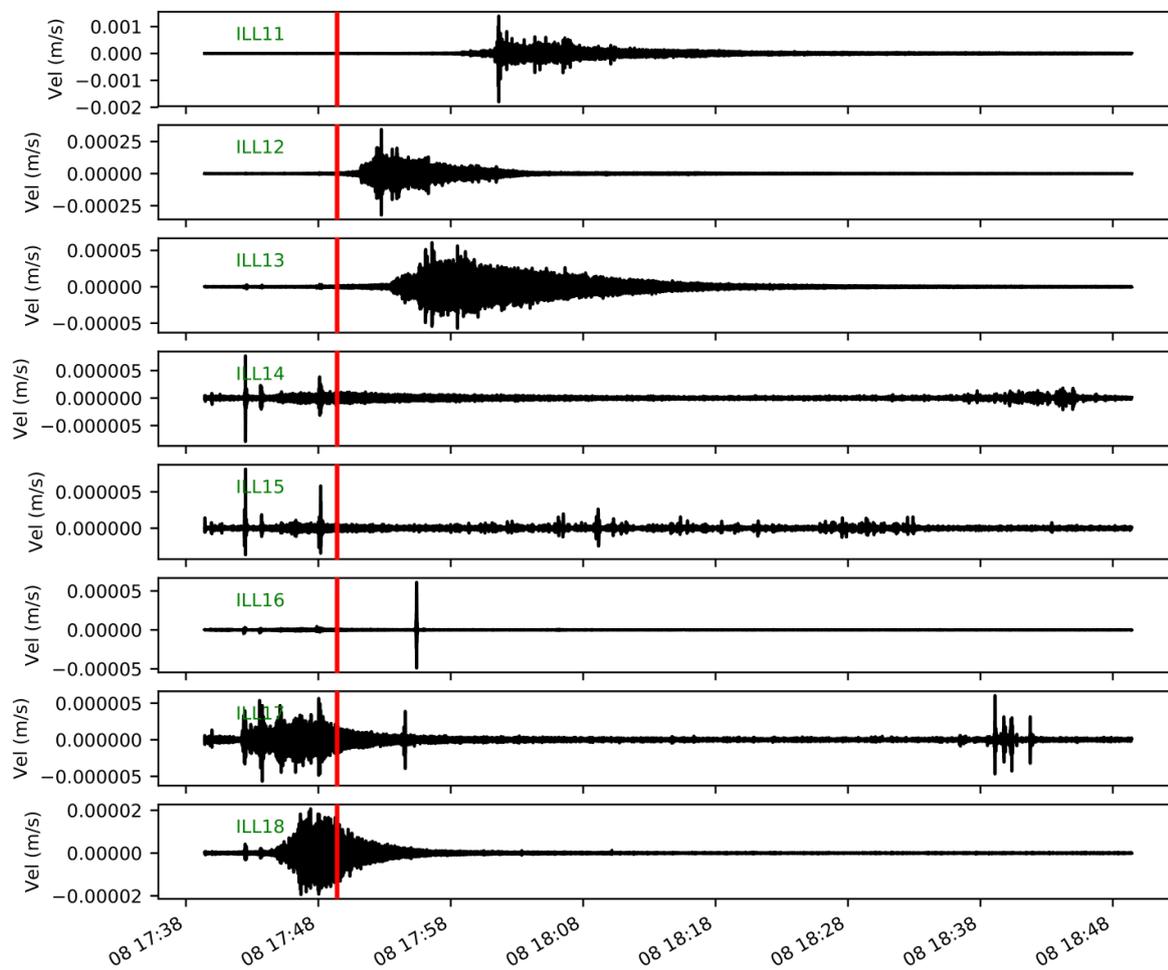
September 17, 2020, 3:01pm

25.07.18, 16:56:40



**Figure S8.** (a) Vertical-component seismograms generated by a debris-flow event on 25 July 2018. The arrival time of the debris flow front at CD1 is marked in red.

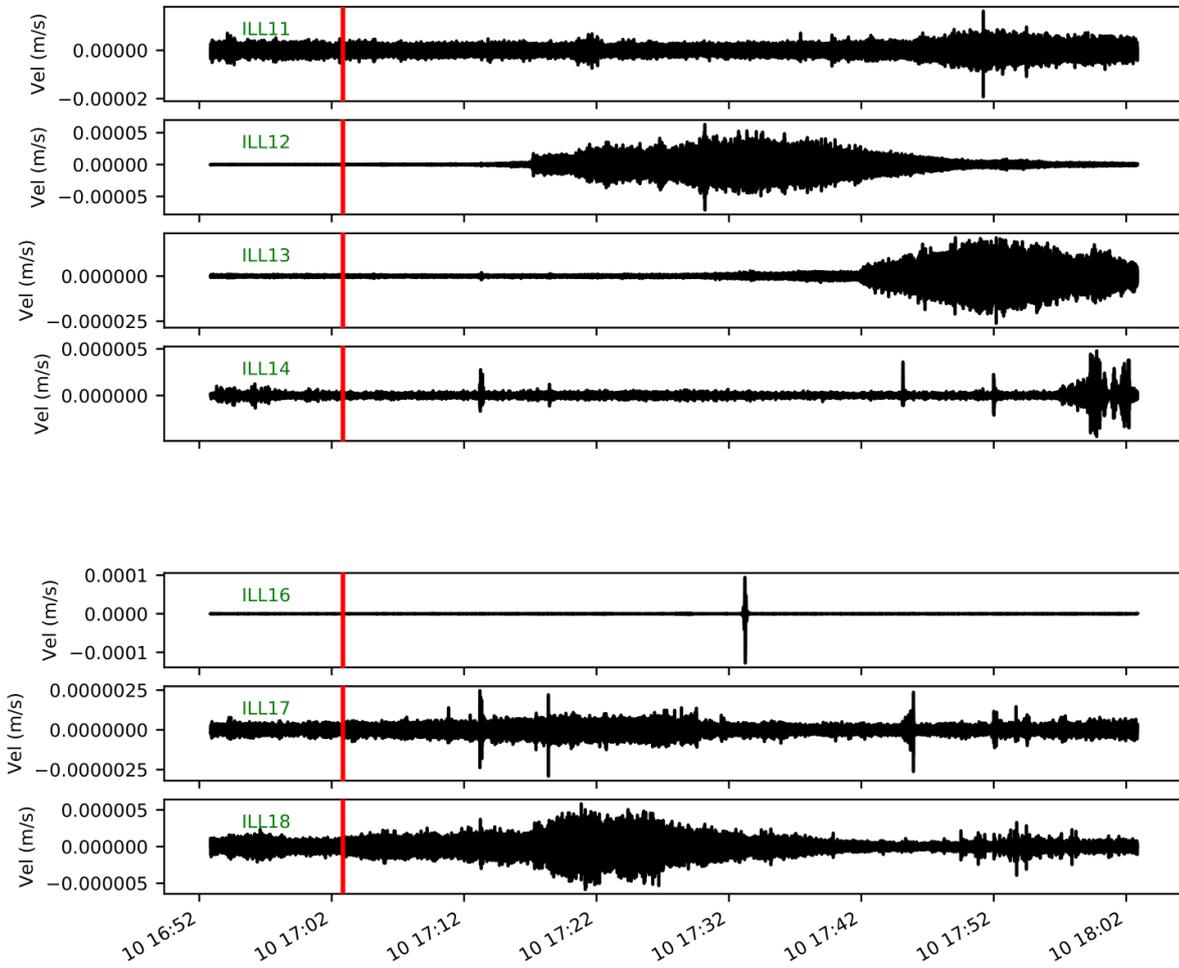
08.08.18, 17:49:25



**Figure S9.** (a) Vertical-component seismograms generated by a debris-flow event on 08 August 2018. The arrival time of the debris flow front at CD1 is marked in red.

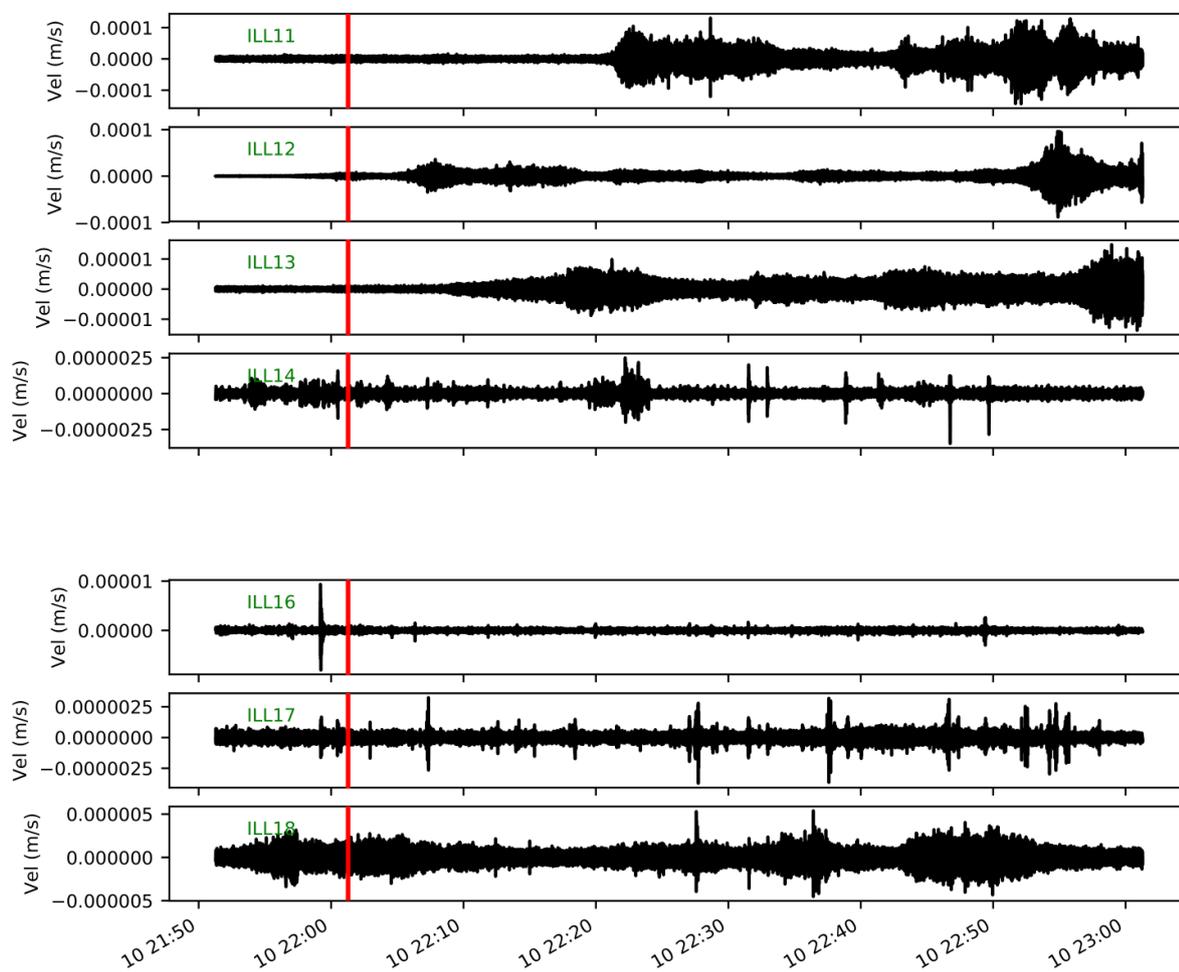
September 17, 2020, 3:01pm

10.06.19, 17:02:51



**Figure S10.** (a) Vertical-component seismograms generated by a debris-flow event on 10 June 2019. The arrival time of the debris flow front at CD1 is marked in red.

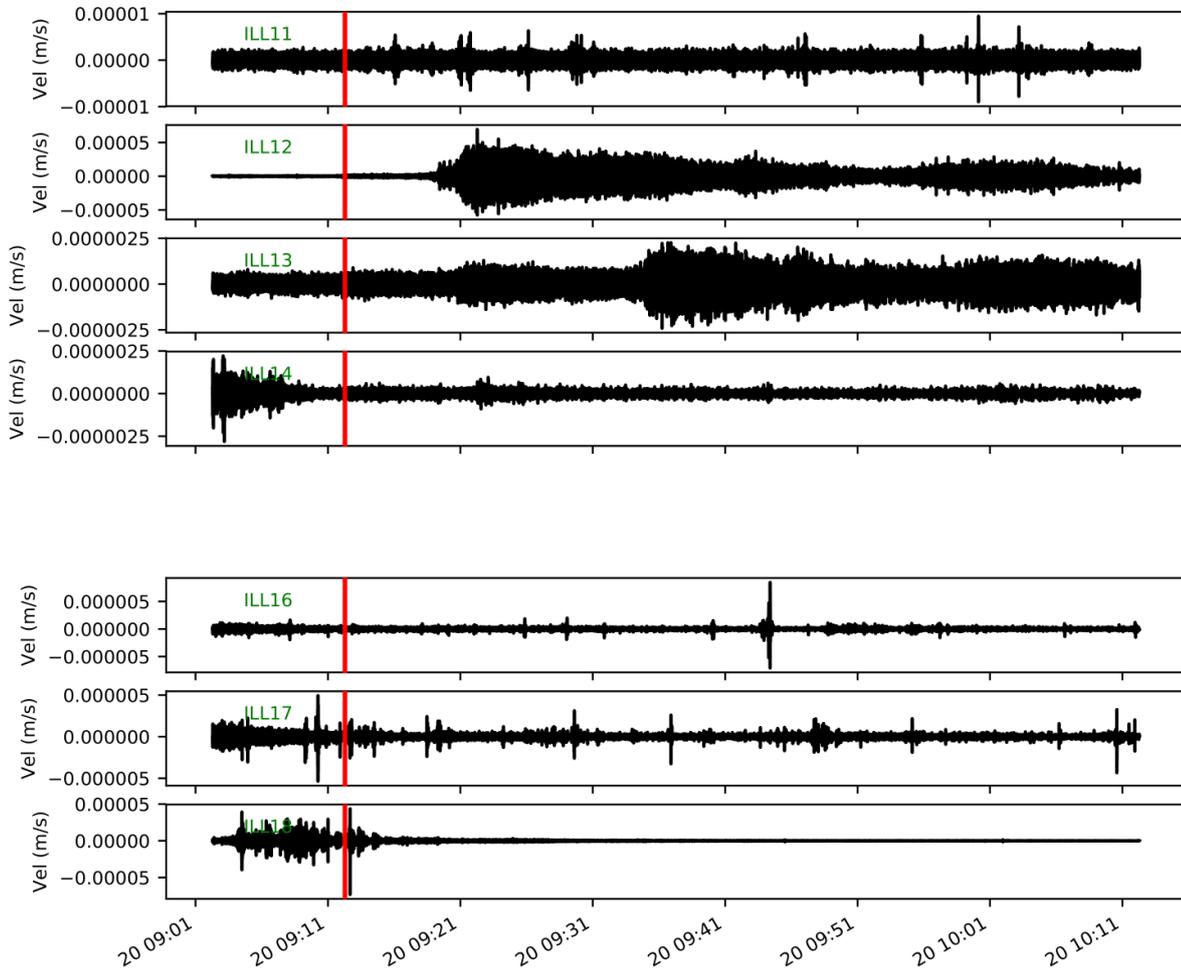
10.06.19, 22:01:17



**Figure S11.** (a) Vertical-component seismograms generated by a debris-flow event on 10 June 2019. The arrival time of the debris flow front at CD1 is marked in red.

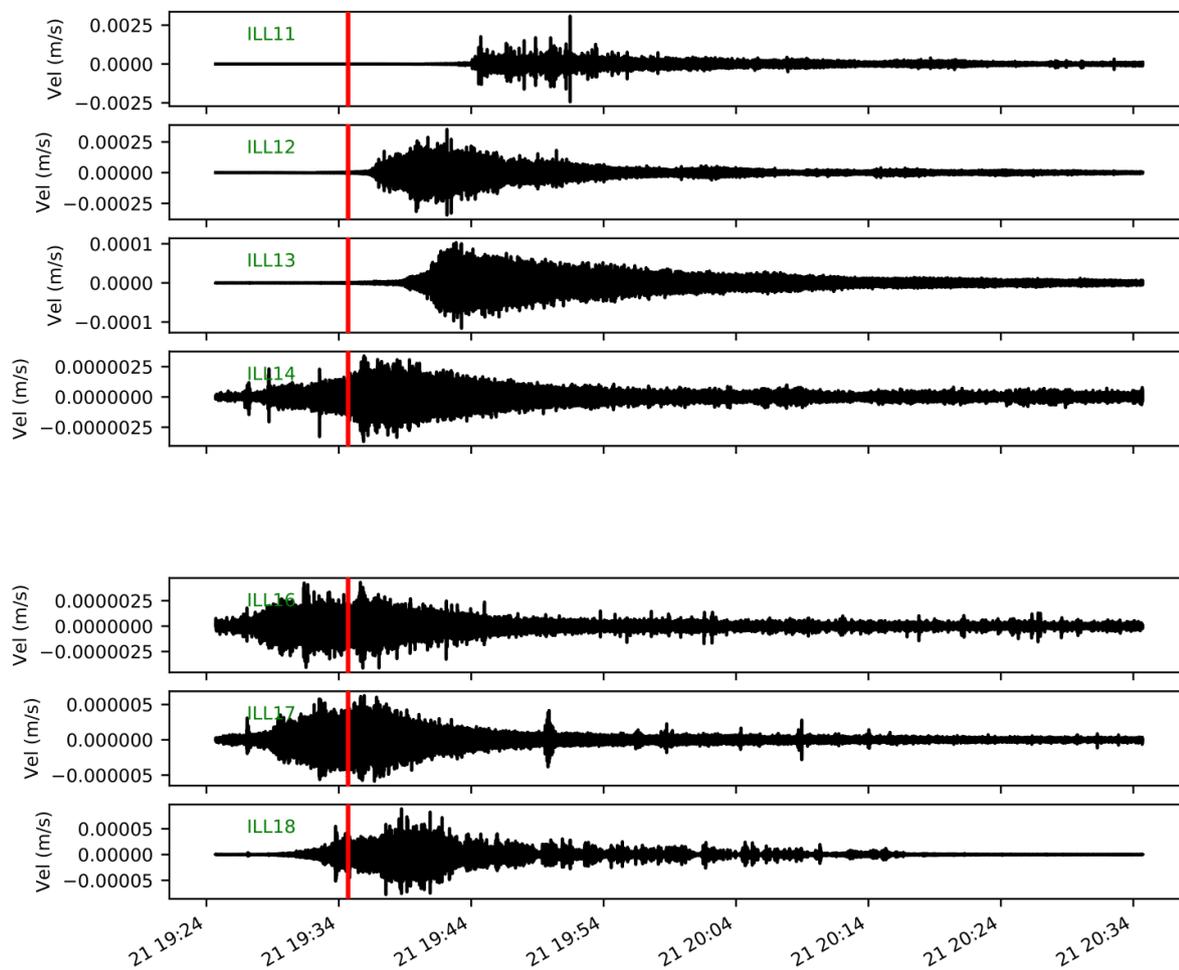
September 17, 2020, 3:01pm

20.06.19, 09:12:17



**Figure S12.** (a) Vertical-component seismograms generated by a debris-flow event on 20 June 2019. The arrival time of the debris flow front at CD1 is marked in red.

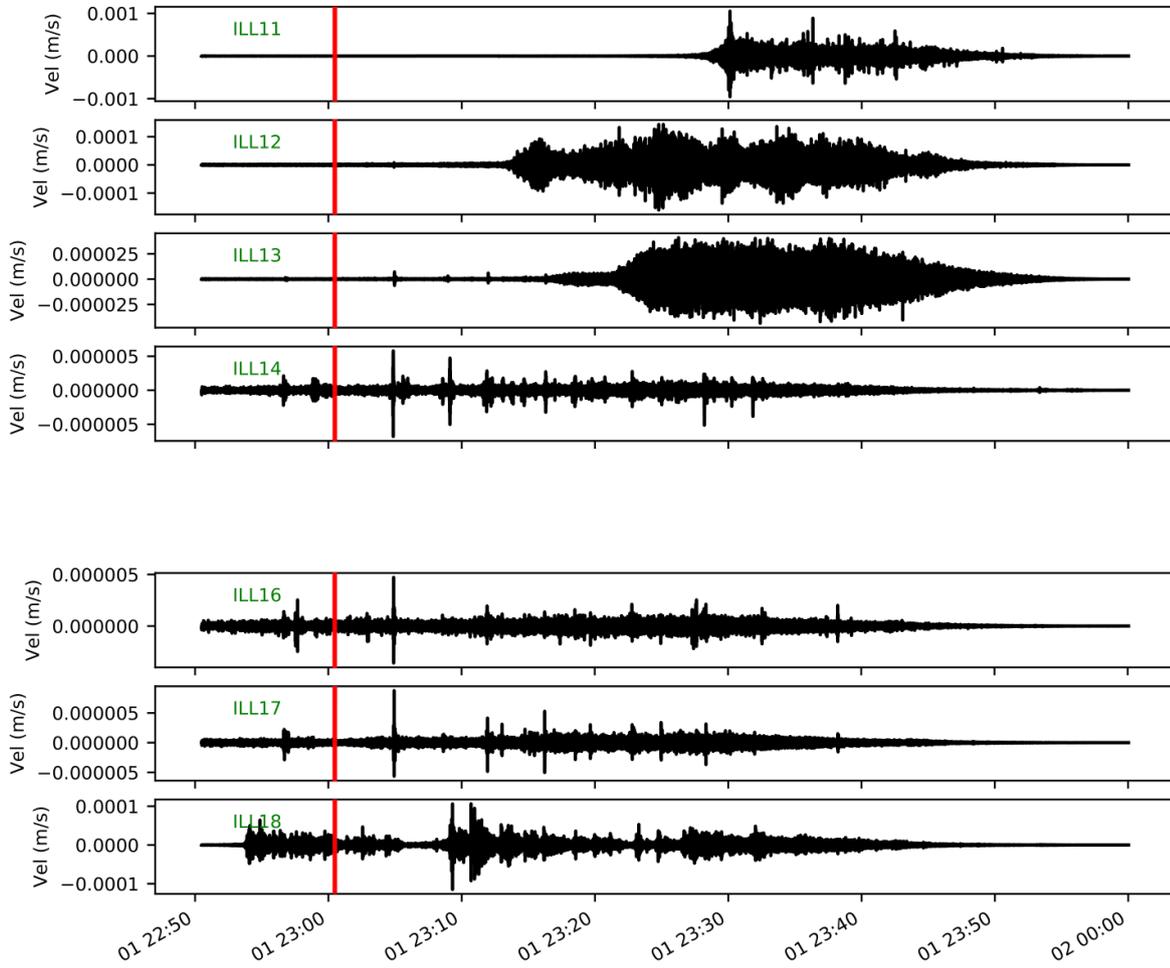
21.06.19, 19:34:42



**Figure S13.** (a) Vertical-component seismograms generated by a debris-flow event on 21 June 2019. The arrival time of the debris flow front at CD1 is marked in red.

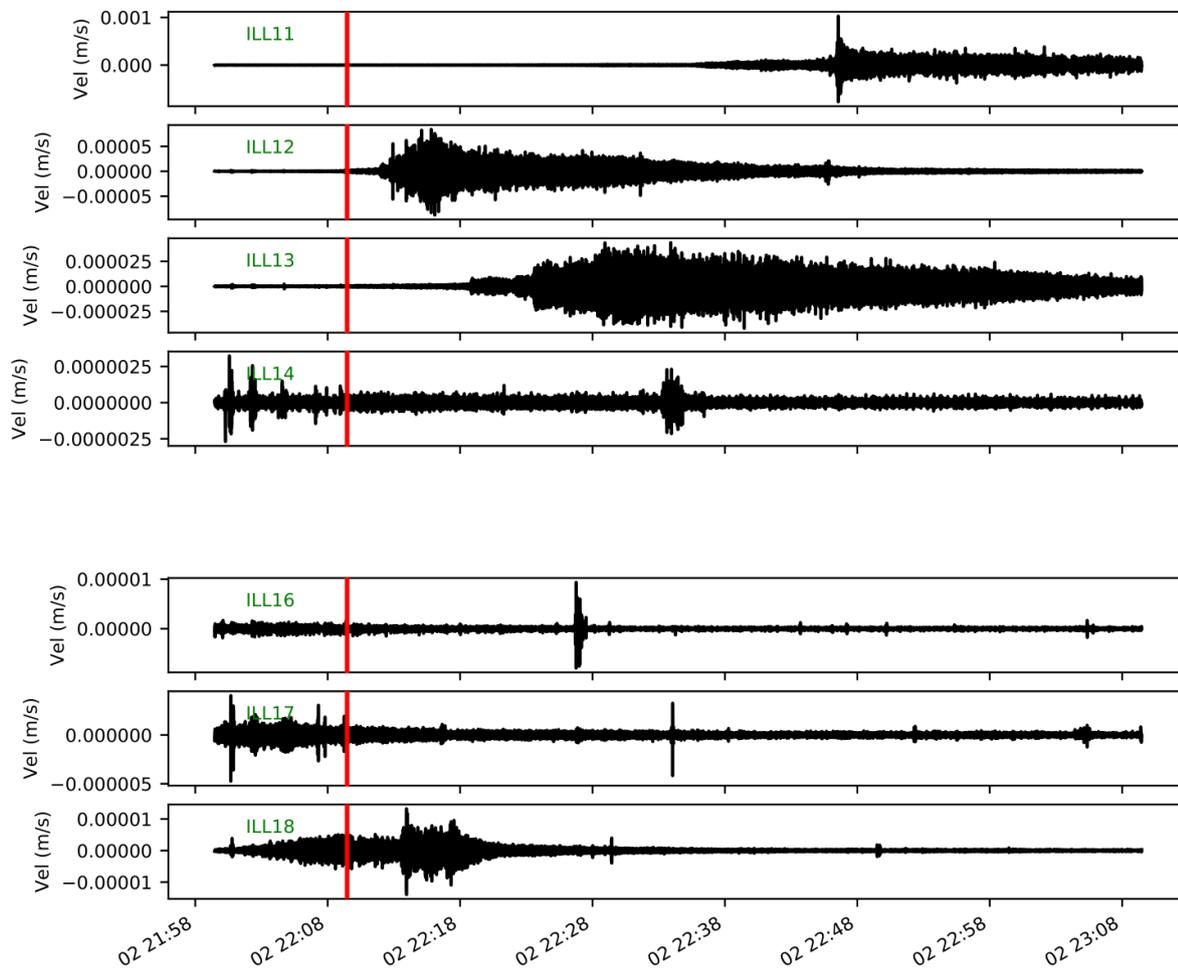
September 17, 2020, 3:01pm

01.07.19, 23:00:29



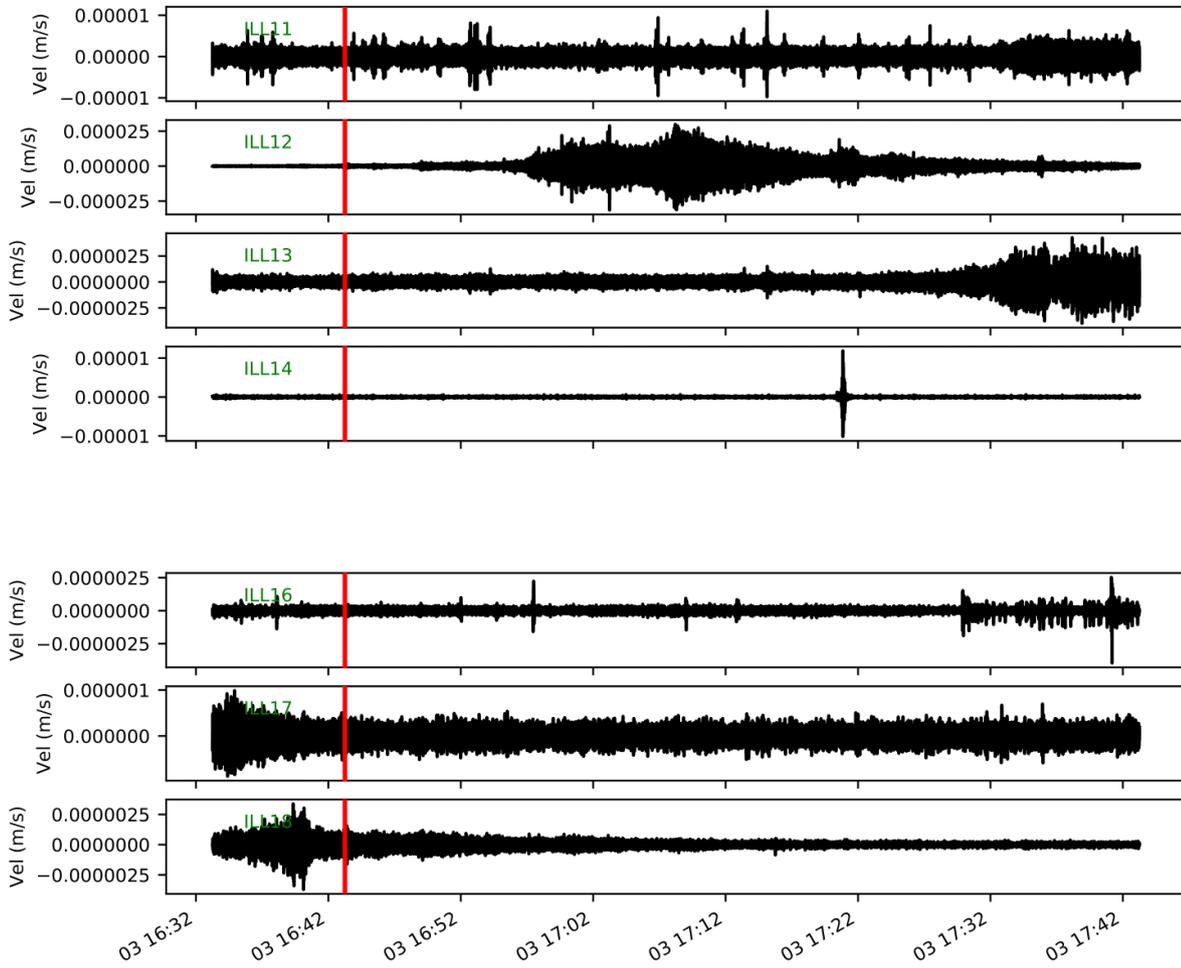
**Figure S14.** (a) Vertical-component seismograms generated by a debris-flow event on 01 July 2019. The arrival time of the debris flow front at CD1 is marked in red.

02.07.19, 22:09:28



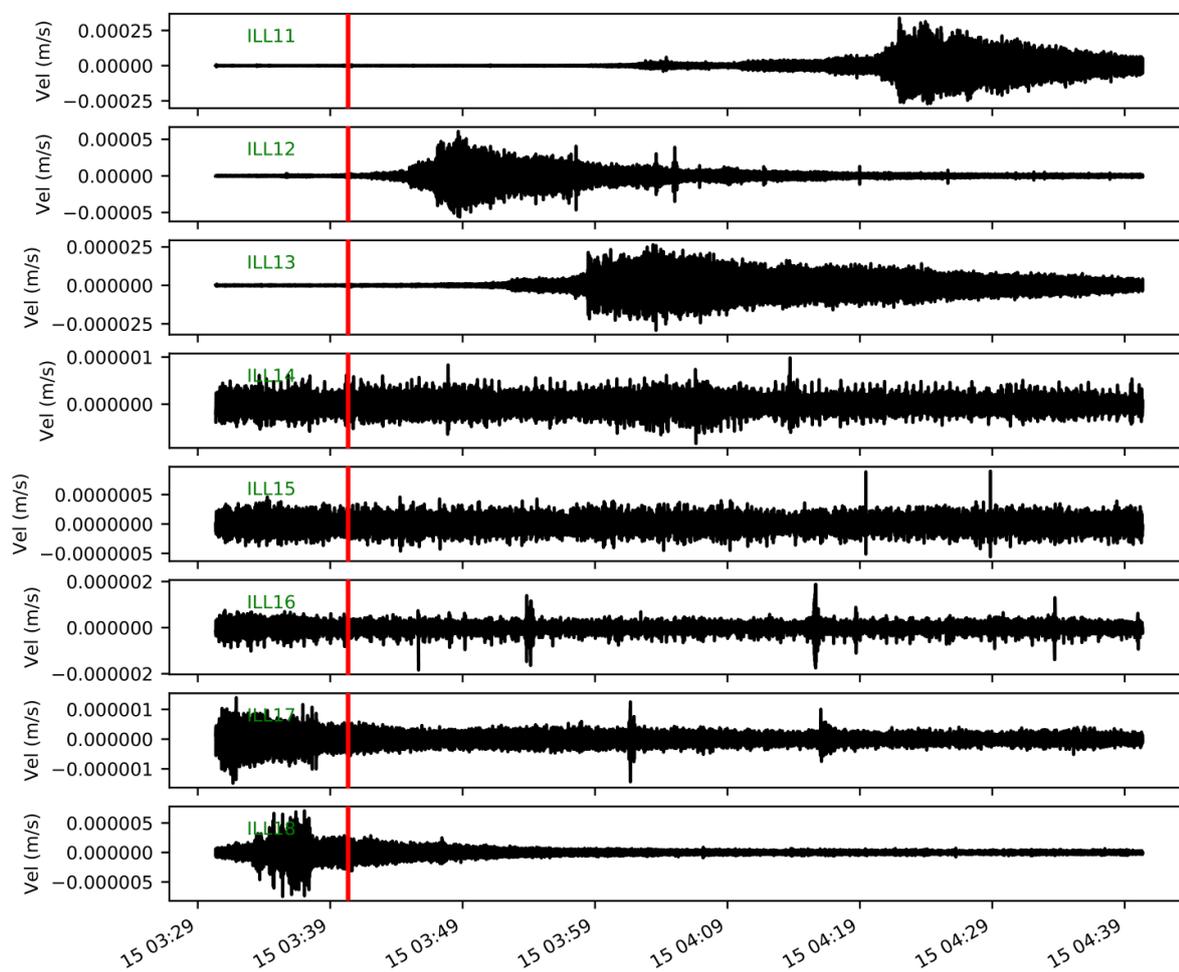
**Figure S15.** (a) Vertical-component seismograms generated by a debris-flow event on 02 July 2019. The arrival time of the debris flow front at CD1 is marked in red.

03.07.19, 16:43:15



**Figure S16.** (a) Vertical-component seismograms generated by a debris-flow event on 03 July 2019. The arrival time of the debris flow front at CD1 is marked in red.

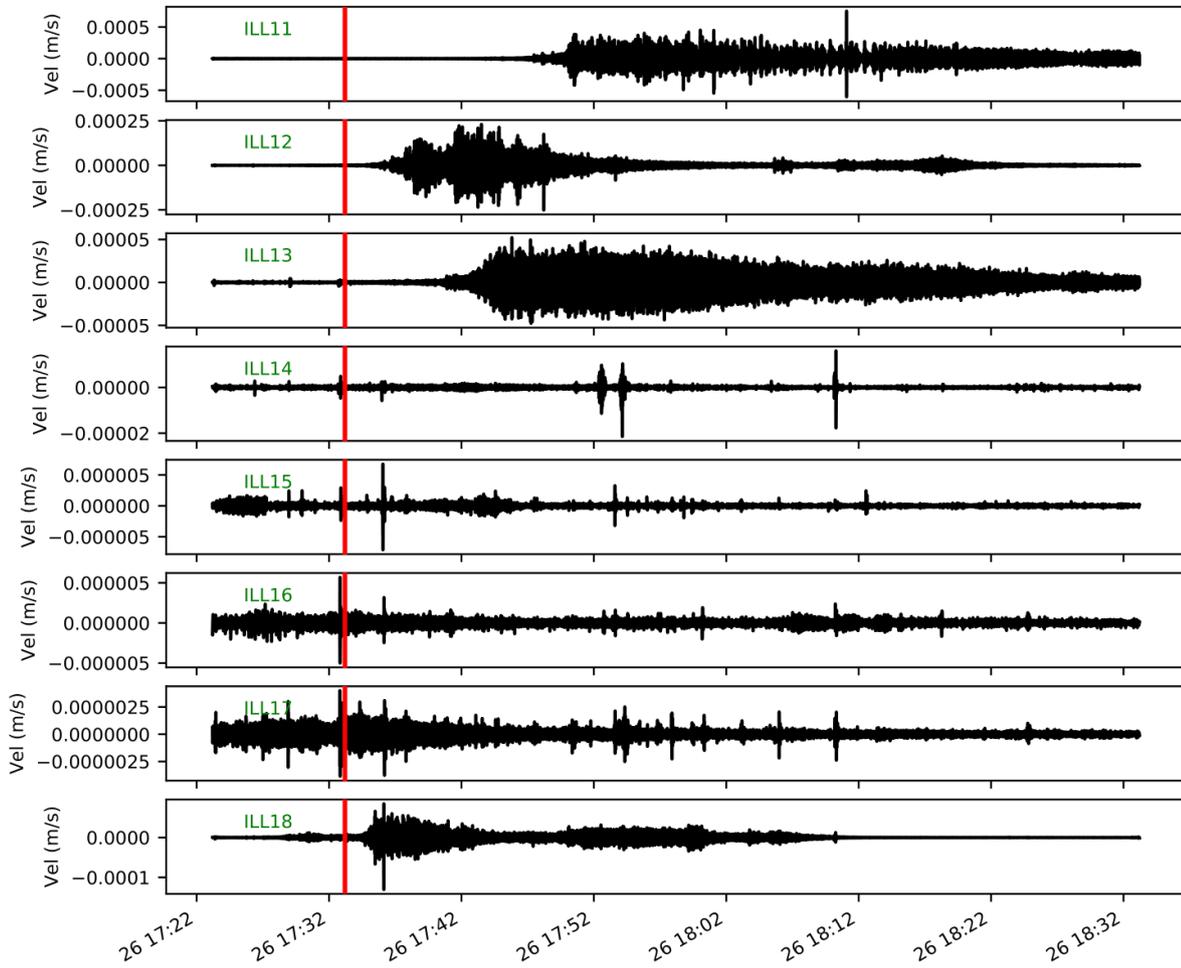
15.07.19, 03:40:21



**Figure S17.** (a) Vertical-component seismograms generated by a debris-flow event on 15 July 2019. The arrival time of the debris flow front at CD1 is marked in red.

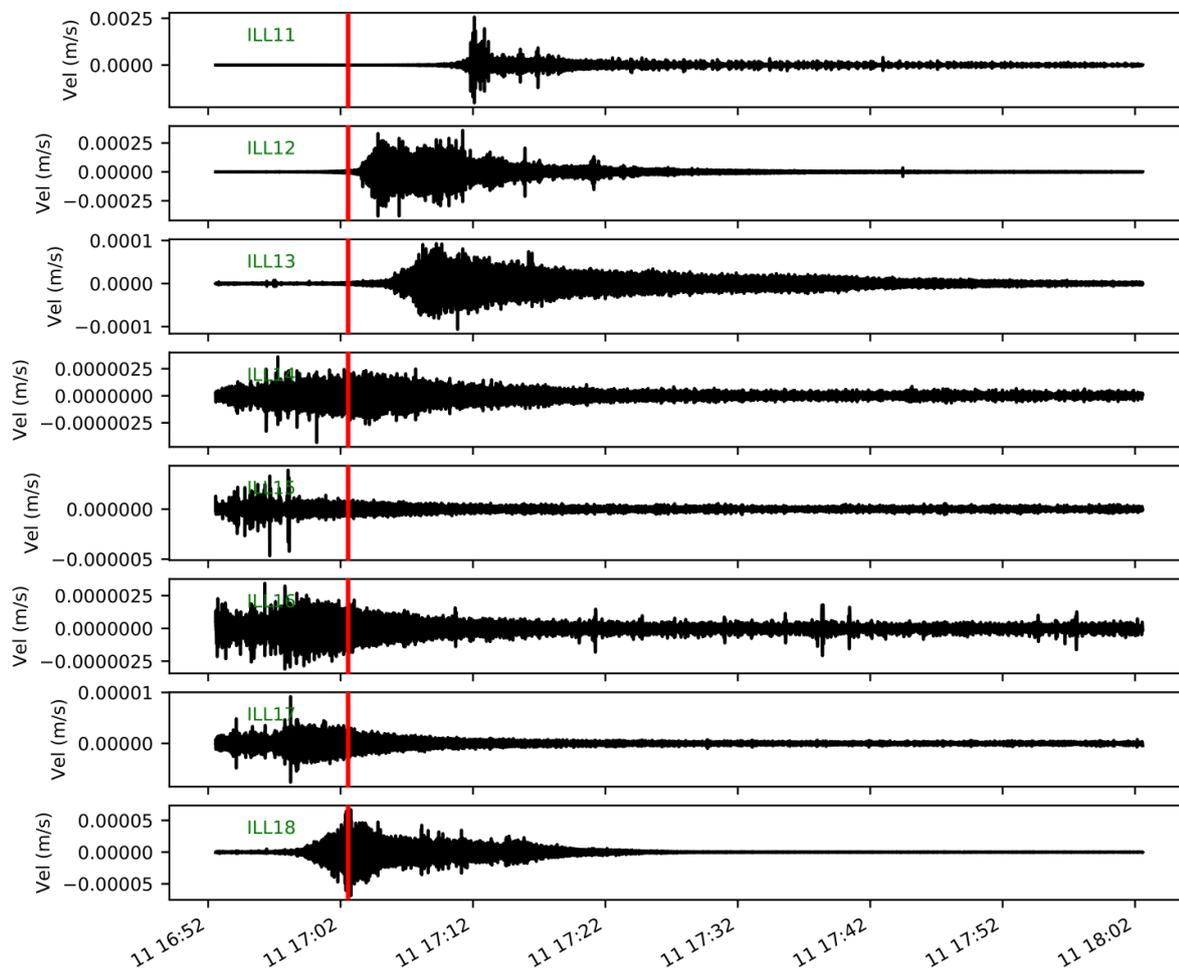
September 17, 2020, 3:01pm

26.07.19, 17:33:12



**Figure S18.** (a) Vertical-component seismograms generated by a debris-flow event on 26 July 2019. The arrival time of the debris flow front at CD1 is marked in red.

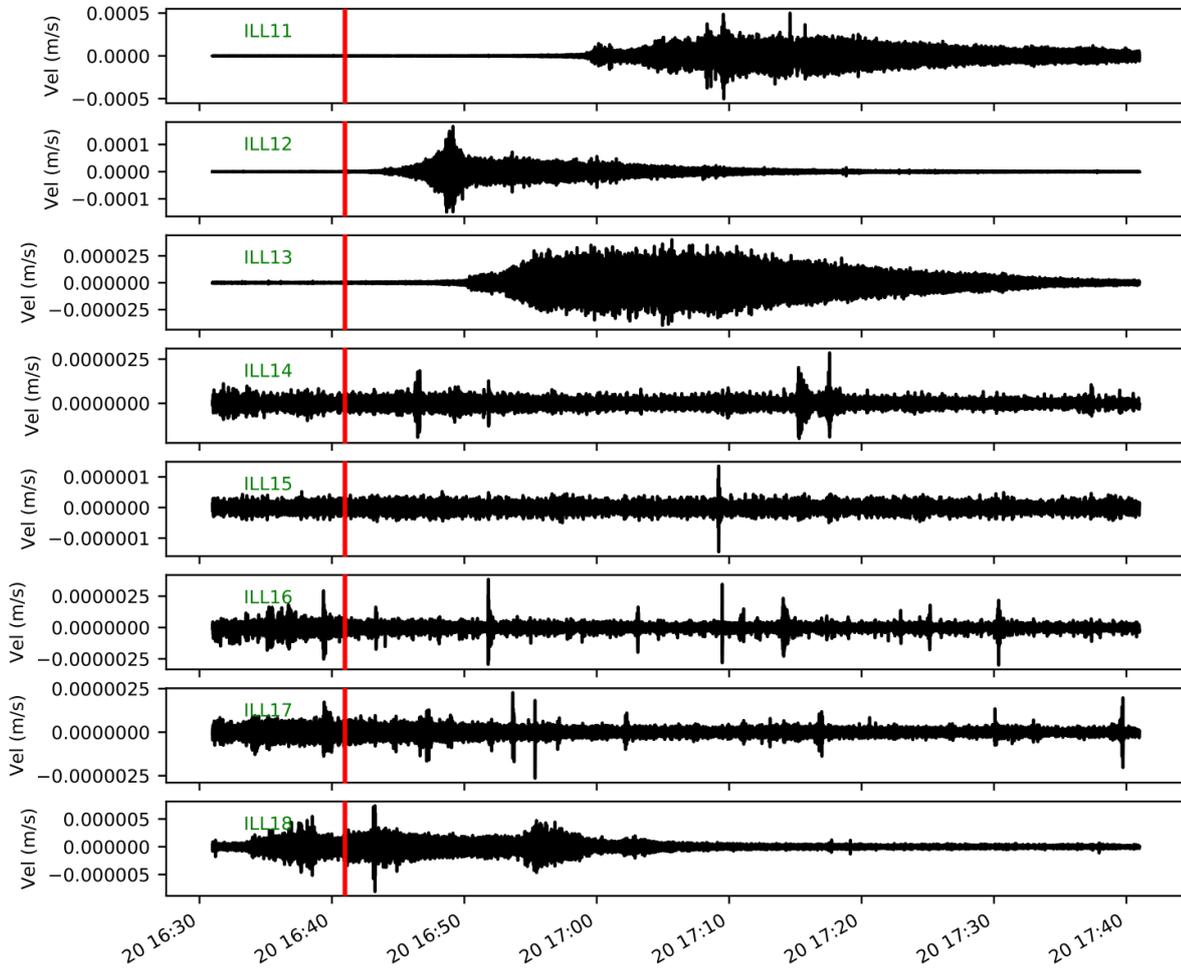
11.08.19, 17:02:34



**Figure S19.** (a) Vertical-component seismograms generated by a debris-flow event on 11 August 2019. The arrival time of the debris flow front at CD1 is marked in red.

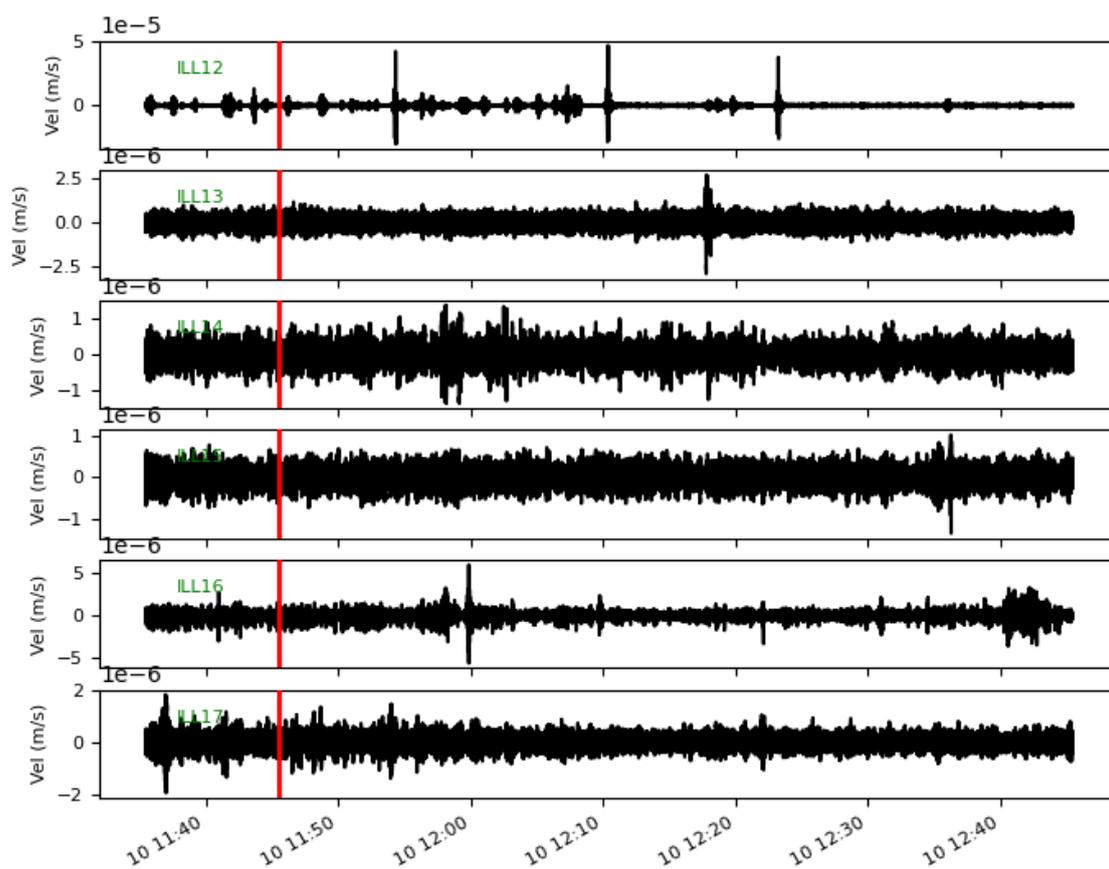
September 17, 2020, 3:01pm

20.08.19, 16:40:59



**Figure S20.** (a) Vertical-component seismograms generated by a debris-flow event on 20 August 2019. The arrival time of the debris flow front at CD1 is marked in red.

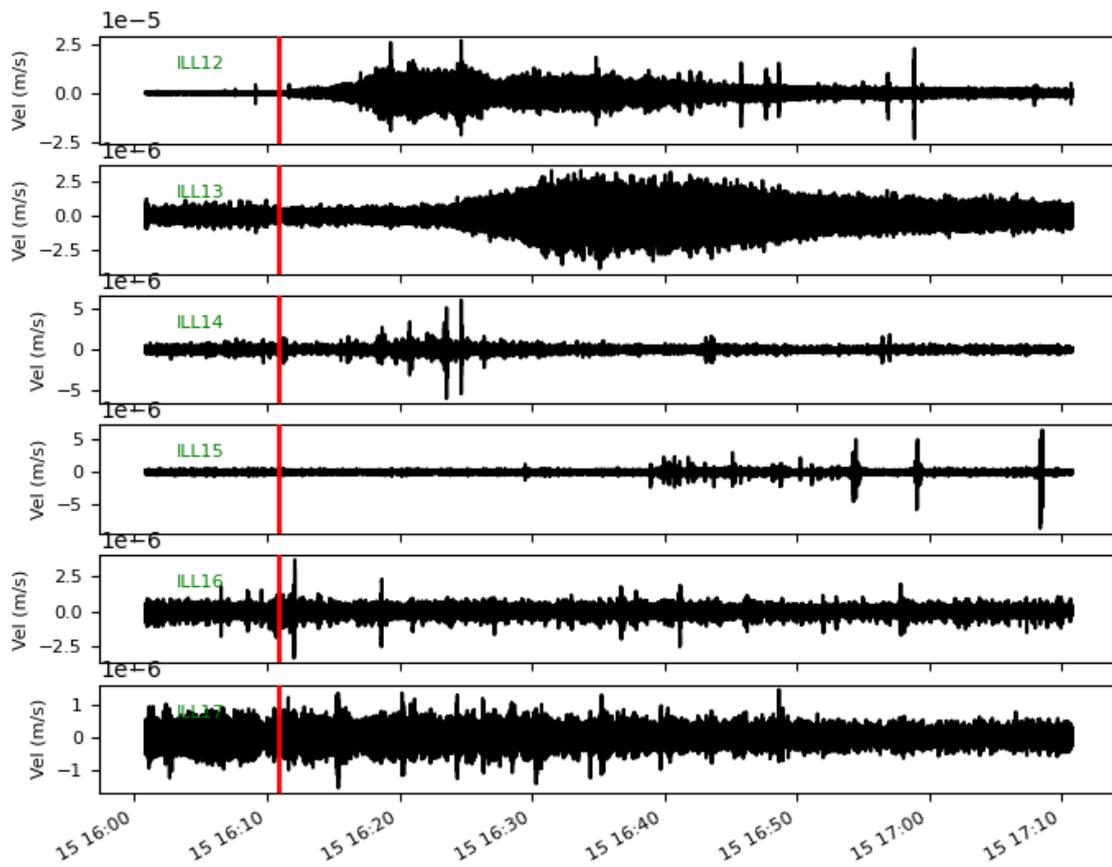
10.10.19, 11:45:28



**Figure S21.** (a) Vertical-component seismograms generated by a debris-flow event on 10 October 2019. The arrival time of the debris flow front at CD1 is marked in red.

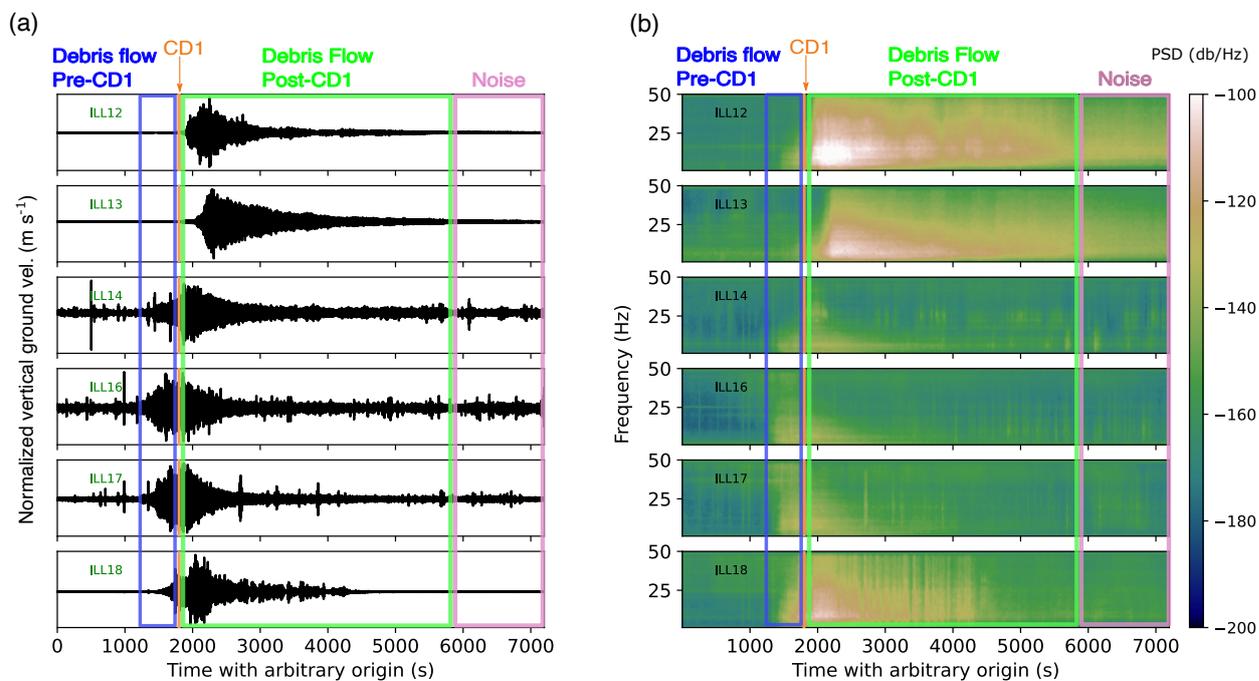
September 17, 2020, 3:01pm

15.10.19, 16:10:50

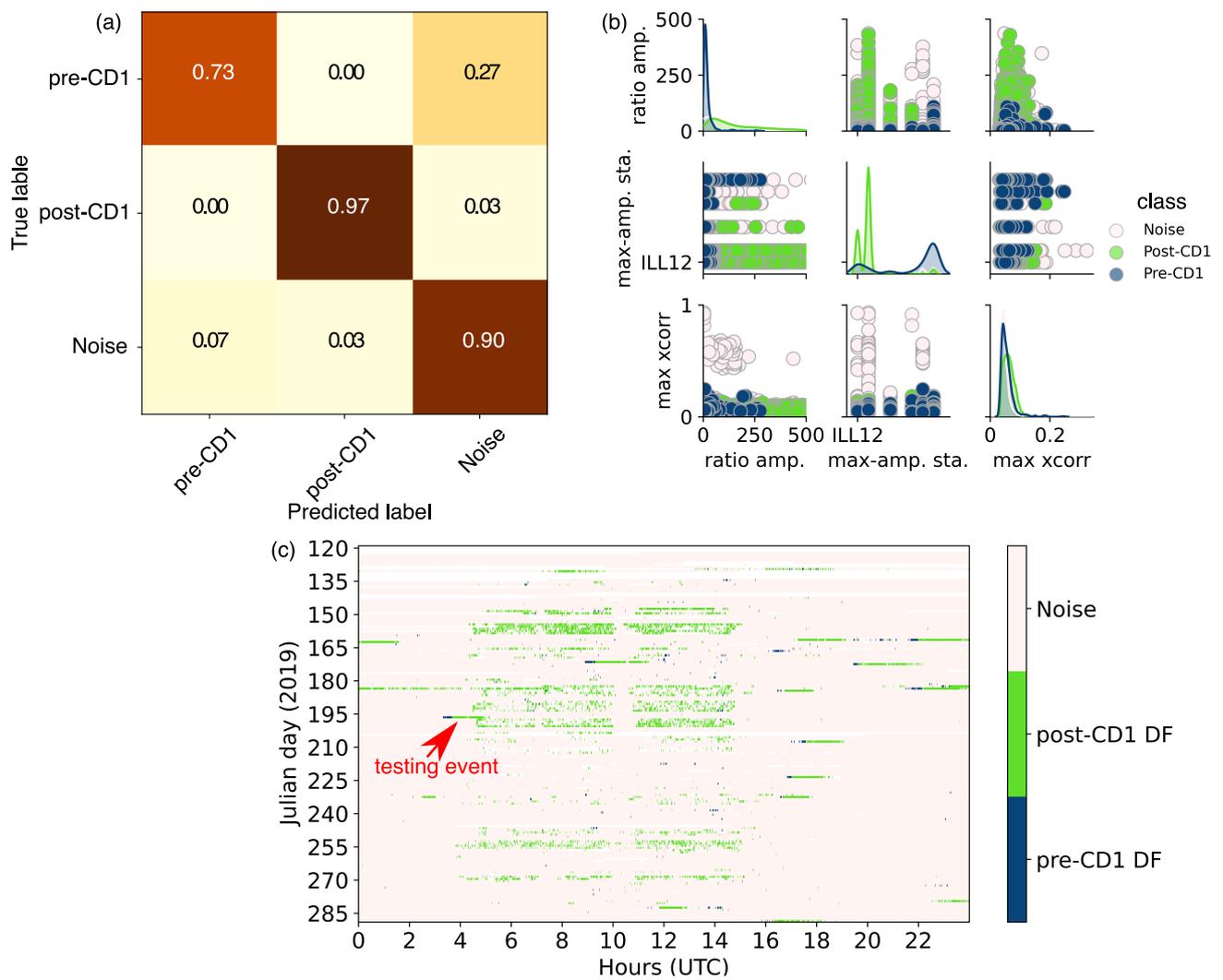


**Figure S22.** (a) Vertical-component seismograms generated by a debris-flow event on 15 October 2019. The arrival time of the debris flow front at CD1 is marked in red.

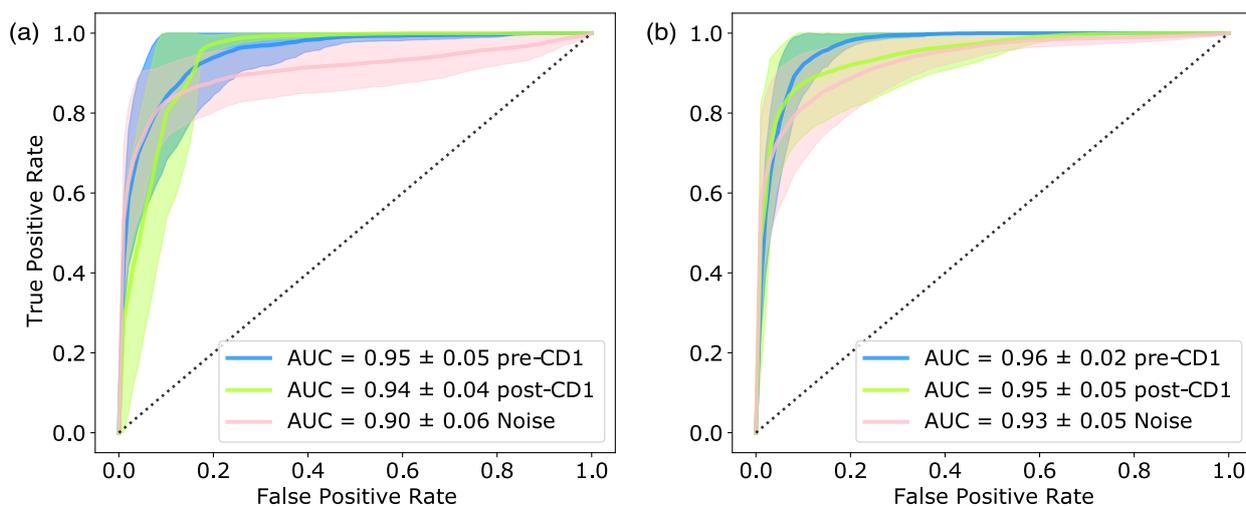
September 17, 2020, 3:01pm



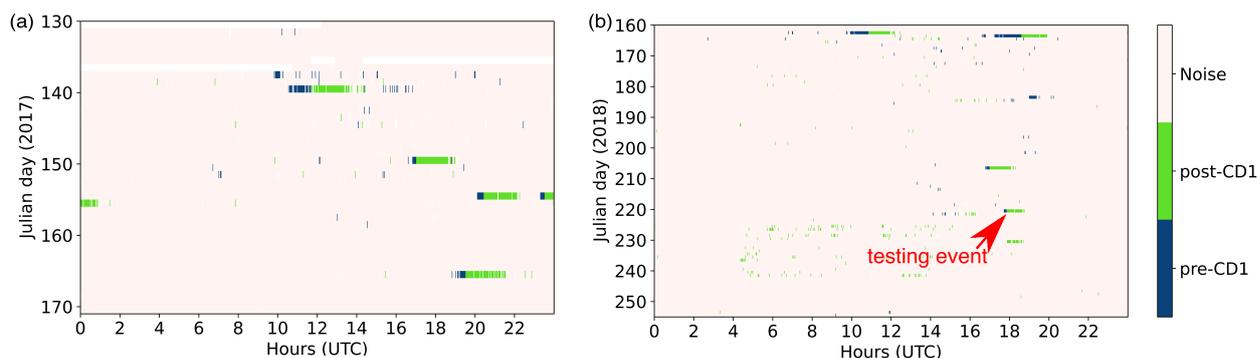
**Figure S23.** (a) Debris-flow seismograms generated by 21 June 2019 event with vol. = 83,000 m<sup>3</sup> recorded over six stations in the network. Corresponding spectrograms are showed in panel b. Three classes of seismic events used in machine-learning detector are schematically represented with different colours: pre-CD1 (blue), post-CD1 (green) and noise (pink). The arrival time of the debris flow front at CD1 is marked in orange.



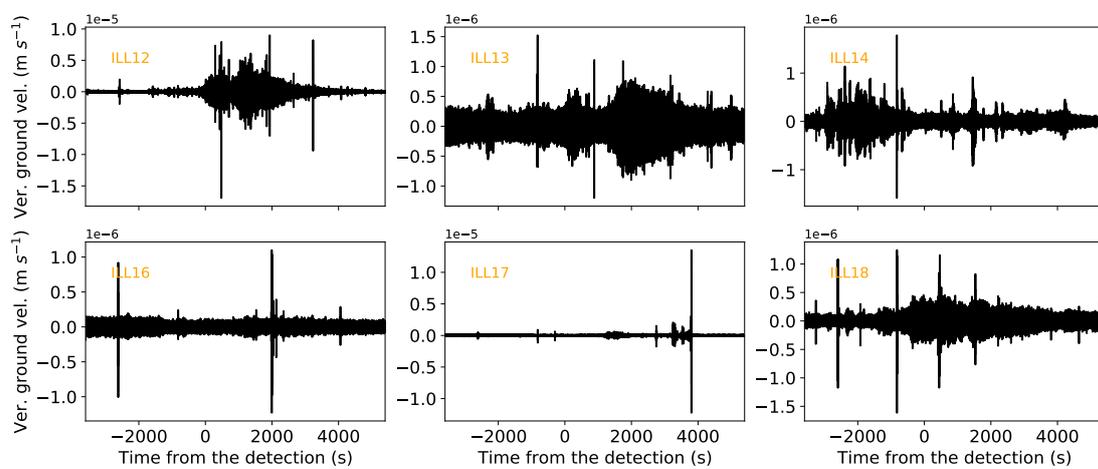
**Figure S24.** Machine-learning model evaluation for the first iteration. (a) Normalized confusion matrix with true labels as columns and predicted labels as rows. (b) Pairwise relationships of the three most important features. In each subplot, two features are plotted against each other (the same features are plotted on diagonal, which show univariate distribution of features). Features from each class are marked in different colors. (c) Results of the machine-learning detector executed on 2019 continuous data.



**Figure S25.** Mean ROC curves for the first (a), and the second iteration (b), calculated using 5-fold cross-validation. The mean ROC curves are marked in solid lines with shaded standard deviations. The true positive rate (TPR) is presented on the y axis and false positive rate (FPR) on the x axis. The area below the curve (AUC) measures model accuracy.



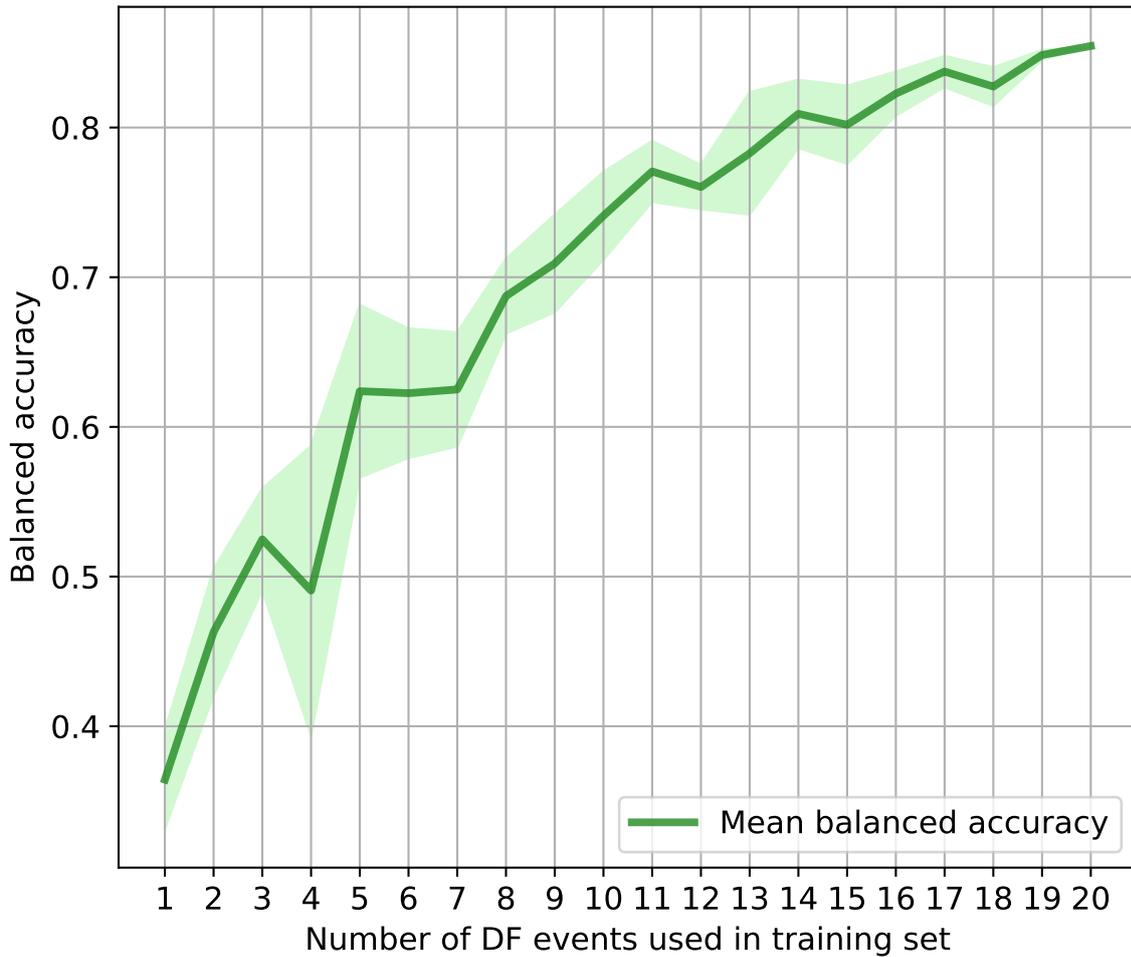
**Figure S26.** Results of the machine-learning detector run (from the second iteration) over 2017 (a), and 2018 (b) continuous data.



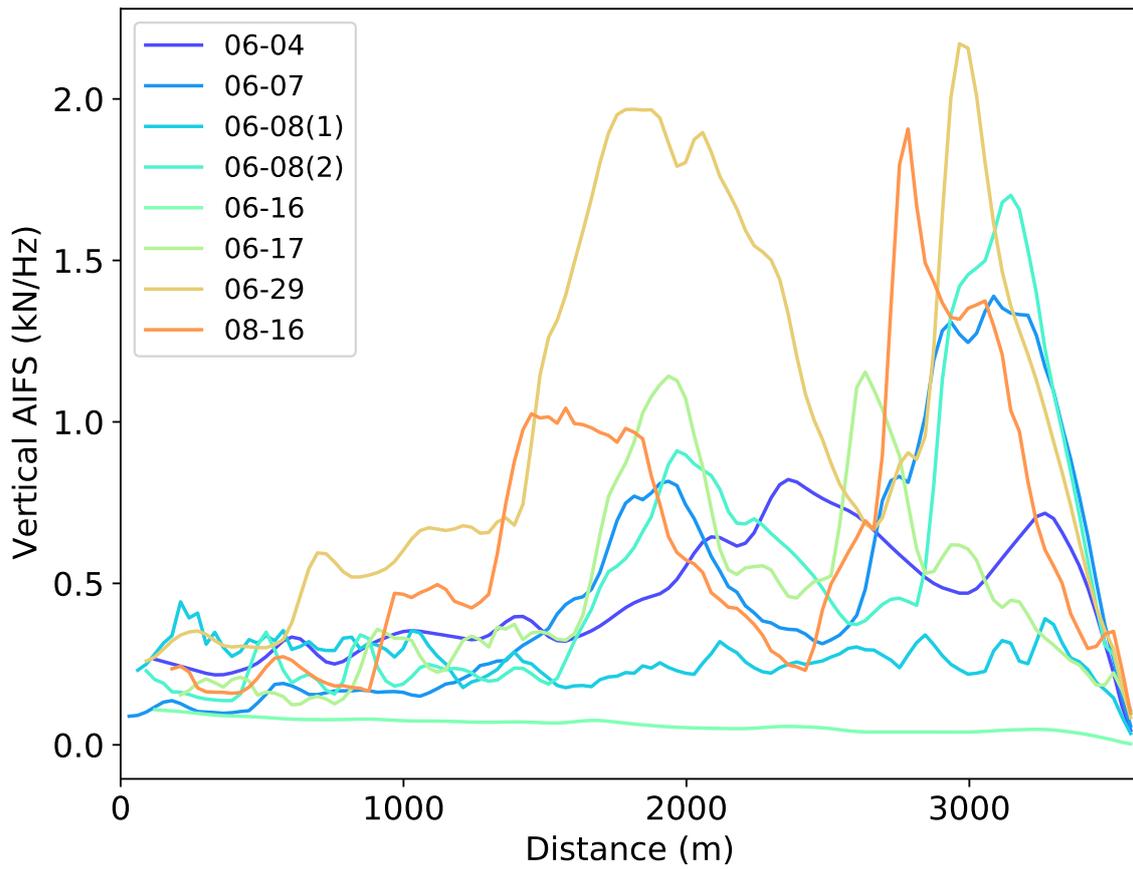
**Figure S27.** Waveforms generated by a small debris flow event found by the ML-based DF detector, and detected on 18 August 2018, 17:53:20.



**Figure S28.** Photos of Illgraben debris-flow events with boulder-rich fronts passing through CD29 detected by the machine-learning detector on 4, 8(2), 17, and 29 June 2020 (Source: WSL).



**Figure S29.** Sensitivity test: balanced accuracy score as a function of  $n$  number of debris-flow (DF) events used in training set with cross-validation (5 folds). For each test a subset of  $n$  events is chosen from 20 events to train the ML-model and two events marked in green in TableS1 are used as testing set. The mean balanced accuracy is marked in solid green line with shaded standard deviation. The values are normalized between 0 and 1. Even a model trained on a single event gives better classification results than a random guess (for a 3-class classification problem balanced accuracy of a random guess converges towards  $1/3$ ). Higher values of balanced accuracy ( $> 0.7$ ) and stable prediction are obtained from  $n=9$  events used in the training set.



**Figure S30.** Vertical apparent total impact force spectra (AIFS) for 2020 debris flows.

N°	Date	Arrival time CD1 (UTC)	Vol.(m <sup>3</sup> )	Vel.(m s <sup>-1</sup> )	$h_{99}$ (m)
0	2017-05-19	11:41:00	n.a.	n.a.	n.a.
1	2017-05-29	16:58:31	100,000	6.67	4.80
2	2017-06-03 (1)	20:23:07	n.a.	n.a.	n.a.
3	2017-06-03 (2)	23:27:38	25,000	5.10	3.30
4	2017-06-14	19:30:48	n.a.	7.10	3.40
5	2018-06-11	10:46:39	35,000	7.00	3.50
6	2018-06-12	18:29:16	n.a.	n.a.	n.a.
7	2018-07-25	16:56:40	<50,000	4.69	2.00
8	2018-08-08	17:49:25	<100,000	6.70	n.a.
9	2019-06-10 (1)	17:02:51	3,300	0.90	0.64
10	2019-06-10 (2)	22:01:17	6,600	2.38	0.59
11	2019-06-20	09:12:17	n.a.	n.a.	n.a.
12	2019-06-21	19:34:42	83,000	5.60	2.45
13	2019-07-01	23:00:29	78,000	3.80	1.62
14	2019-07-02	22:09:28	39,000	2.50	0.71
15	2019-07-03	16:43:15	n.a.	n.a.	n.a.
16	2019-07-15	03:40:21	16,000	5.00	0.68
17	2019-07-26	17:33:12	64,000	6.97	1.21
18	2019-08-11	17:02:34	53,000	5.56	n.a.
19	2019-08-20	16:40:59	13,000	0.95	0.89
20	2019-10-09	11:45:28	n.a.	n.a.	n.a.
21	2019-10-15	16:10:50	n.a.	n.a.	n.a.

**TableS1.** Characteristics of 22 debris flow events recorded in 2017, 2018, and 2019. Volume is the integrated sum of discharge over the entire debris-flow wave. Flow velocity is calculated from the travel time between in-channel sensors as described in (Schlunegger et al., 2009). Flow depth  $h_{99}$  is the depth where 99% of the depth values are smaller. n.a. denotes values that were not estimated. Volume and flow depth are estimated at the instrumented wall, CD29. The arrival times at CD1 come mostly from the measurements of a geophone installed at CD1, although the arrival times of events 0 and 2 were estimated based on the ASL results (Walter et al., 2017). Events used in the first iteration in the testing set are marked in orange, and events used in the testing set in the second iteration are marked in green.

N°	Date	Alarm	ILL11 arrivals	Warning time increase (min:s)	Peak ampl. (n° counts)	DF Vel. (m s <sup>-1</sup> )
0	2020-06-04	14:55:08	15:41:51	43:43	250,849	3.2
1	2020-06-07	07:33:28	09:29:32	136:04	117,228	0.7
2	2020-06-08 (1)	13:43:28	15:51:59	128:31	28,383	0.6
3	2020-06-08 (2)	16:35:08	17:59:10	84:02	132,500	0.8
4	2020-06-09	23:35:08	00:55:43(+1d)	90:35	84,089	2.1
5	2020-06-16	20:20:08	23:56:18	216:10	2,486	3.2
6	2020-06-17	03:11:48	04:06:58	55:10	7,468	0.7
7	2020-06-29	04:33:29	05:49:13	75:44	159,970	1.5
8	2020-07-22	15:41:50	n.a.	n.a.	n.a.	n.a.
9	2020-07-28	16:05:10	17:58:55	113:45	2,461	n.a.
10	2020-08-16	21:15:08	23:04:11	109:03	96,098	0.6.
11	2020-08-30	04:54:24	05:52:23	57:59	23,455	n.a.
12	2020-09-01	04:54:24	n.a.	n.a.	n.a.	n.a.

**TableS3.** Characteristics of 13 alarms in 2020. n.a. denotes values that were not estimated.

22 July 2020 and 01 September 2020 debris flows stopped before ILL11, and 30 August 2020 event had multiple surges which makes difficult reliable debris-flow velocity estimation and AIFS calculation.

### Captions for large Table S2.

**TableS2.** 70 statistical features including waveform, spectral, spectrogram, and network attributes used as input in the machine-learning model. The table is uploaded separately as an excel file.