

# Near-real-time detection of co-seismic ionospheric disturbances using machine learning

Quentin Brissaud, Elvira Astafyeva

### ▶ To cite this version:

Quentin Brissaud, Elvira Astafyeva. Near-real-time detection of co-seismic ionospheric disturbances using machine learning. Geophysical Journal International, 2022, 230, pp.2117-2130. 10.1093/gji/ggac167. insu-03748509

## HAL Id: insu-03748509 https://insu.hal.science/insu-03748509

Submitted on 30 Mar 2023

**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. *Geophys. J. Int.* (2022) **230**, 2117–2130 Advance Access publication 2022 May 2 GJI Seismology

## Near-real-time detection of co-seismic ionospheric disturbances using machine learning

Quentin Brissaud<sup>10</sup> and Elvira Astafyeva<sup>2</sup>

<sup>1</sup>NORSAR, 2007 Kjeller, Norway. E-mail: quentin@norsar.no <sup>2</sup>Institut de Physique du Globe de Paris (IPGP), Université de Paris, CNRS UMR7154, 35-39 Rue Hélène Brion, 75013 Paris, France

Accepted 2022 April 28. Received 2022 April 22; in original form 2021 December 30

#### SUMMARY

Tsunamis generated by large earthquake-induced displacements of the ocean floor can lead to tragic consequences for coastal communities. Measurements of co-seismic ionospheric disturbances (CIDs) offer a unique solution to characterize an earthquake's tsunami potential in near-real-time (NRT) since CIDs can be detected within 15 min of a seismic event. However, the detection of CIDs relies on human experts, which currently prevents the deployment of ionospheric methods in NRT. To address this critical lack of automatic procedure, we designed a machine-learning-based framework to (1) classify ionospheric waveforms into CIDs and noise, (2) pick CID arrival times and (3) associate arrivals across a satellite network in NRT. Machine-learning models (random forests) trained over an extensive ionospheric waveform data set show excellent classification and arrival-time picking performances compared to existing detection procedures, which paves the way for the NRT imaging of surface displacements from the ionosphere.

Key words: Ionosphere/atmosphere interactions; Tsunami warning; Artificial intelligence.

#### **1 INTRODUCTION**

Large seafloor displacements due to earthquakes are known to generate destructive tsunamis. Unfortunately, near-real-time (NRT) mapping of the co-seismic surface displacements to characterize the earthquake tsunami potential is still challenging for conventional methods, especially for earthquakes with magnitudes  $M_w > 8$  (Wright *et al.* 2012; Katsumata *et al.* 2013; LaBrecque *et al.* 2019). In our definition, NRT corresponds to times within 15–20 min after the earthquake onset which is crucial for early-warning application as it gives several tens of minutes for populations to evacuate before the tsunami reaches the coasts.

Recently, several research groups have demonstrated that ionospheric measurements can offer an alternative to seismo-geodetic methods to estimate the tsunami potential of earthquakes. The ionosphere is an electrically charged atmospheric layer that is concentrated between 150 and 400 km of altitude. This layer is sensitive to the vertically propagating acoustic energy excited by natural (e.g. earthquakes, tsunamis and volcanic eruptions) and man-made events (e.g. explosions, rocket launches and nuclear tests) (Heki 2006; Komjathy *et al.* 2016; Rolland *et al.* 2016; Shults *et al.* 2016; Astafyeva 2019; Astafyeva & Shults 2019). In particular, the ionospheric signature of earthquakes, known as co-seismic ionospheric disturbances (CIDs), can be detected 7–9 min after the earthquake. CID waveform characteristics are correlated to the seismic source properties which can help us constraining source parameters and might inform us about the tsunamigenic potential of an earthquake. For instance, the amplitude of the CID scales almost linearly with the magnitude of an earthquake (Astafyeva *et al.* 2013b, 2014; Cahyadi & Heki 2015; Occhipinti *et al.* 2018; Heki 2021), or—for submarine earthquakes—with the tsunami wave height or volume of water that was displaced due to an earthquake (Kamogawa *et al.* 2016; Rakoto *et al.* 2018; Manta *et al.* 2020). Additionally, CID arrival times and detection coordinates provide strong constraints on the position of the seismic source, or the origin of tsunami (Afraimovich *et al.* 2006; Heki *et al.* 2006; Astafyeva *et al.* 2009; Tsai *et al.* 2011; Lee *et al.* 2018; Bagiya *et al.* 2020; Inchin *et al.* 2021; Zedek *et al.* 2021). Moreover, Astafyeva *et al.* (2011, 2013a) and Astafyeva (2019) showed that the distribution of the first-detected CIDs matches the position of the maximum displacement on the ground and Kakinami *et al.* (2021) showed that the initial point of CID matches the maximum vertical displacement of the tsunami source.

However, despite the high potential of seismo-ionospheric assessment of natural hazards, the detection and analysis of ionospheric disturbances still rely on human experts. This manual process is problematic when processing large data volume. Only a few studies have focused on the automatization of detection procedures in the ionosphere but only at low frequencies (Efendi & Arikan 2017; Belehaki *et al.* 2020). Ravanelli *et al.* (2021) investigated the use of both Global Navigation Satellite System (GNSS) ground and ionospheric Total Electron Content (TEC) measurements for NRT tsunami genesis estimation. However, Ravanelli *et al.* (2021) did not present any detection procedure for CIDs, but only showed TEC variations in NRT scenario. In addition, their TEC processing procedure included the use of 8th order polynomial fit in order to highlight the co-seismic signature. The latter is not possible in our definition of NRT mode, that is 15–20 min after the earthquake onset time. The first NRT-compatible method detecting CID was suggested by Maletckii & Astafyeva (2021). However, their study only showed good results on 1 Hz data with CIDs showing high temporal TEC derivative. Therefore, the community needs methods allowing for rapid automatic detection and recognition of CIDs for both future NRT developments and processing of large amount of TEC data retrospectively.

The problem of earthquake waveform detection has been investigated in the seismic community since the early days of modern computers (e.g. Allen 1982). The automatization of waveform detection procedures has historically been performed in the seismic community using analytical methods such as the short-time average/longtime average (STA/LTA) filter (Allen 1982). However, the high rate of false positives generated by these analytical filters has motivated the seismic community to implement machine-leaning (ML) approaches that combine both low computational time and high accuracy (Ross et al. 2018; Mousavi et al. 2020). Even when only small labelled waveform data sets are available, ML methods provide excellent classification results (Provost et al. 2017; Wenner et al. 2021). In particular, random forests (RFs; Breiman 2001) show excellent generalization abilities and do not require an extensive hyperparameter tuning. RF is an ensemble technique that builds predictions by aggregating predictions from a set of decision trees. Aggregating results from individual decision trees built using bootstrap aggregation, which consist of randomly selecting input features to train each tree, makes RF particularly robust to new data.

To address the lack of automatic detection method, we build an RF-based architecture to classify TEC time-series, pick arrival times and associate detected arrivals. RFs are trained over an extensive CID waveform data set from 12 large-magnitude earthquakes to classify TEC waveforms between CIDs and noise and pick arrival times in NRT. Our method is, to the best of our knowledge, the first reported ML classifier and arrival-time picker of CIDs. In this paper, we first describe the generation of our waveform data set, our detection procedure and our ML models. We show classification performance results over our testing data set and against other analytical detection methods. We finally discuss the future implementation of such method for NRT applications.

#### **2 DATA COLLECTION**

The GNSS is widely used to sound the ionosphere. GNSS signals transmitted by satellites and captured by ground-based dualfrequency GNSS receivers enable the estimation of the differential slant TEC (sTEC), which is equal to the number of electrons along a line-of-sight (LOS) between a satellite and a receiver. The sTEC is calculated from phase and code measurements (Afraimovich et al. 2006; Hofmann-Wellenhof et al. 2008; Shults et al. 2016). The phase measurements provide precise information about the ionospheric variations and disturbances, but they are biased by an unknown phase ambiguity constant. The code measurements are noisy and less precise, but are not ambiguous, which enables to estimate the bias by averaging the code values along the arc of measurements. The sTEC is then estimated by removing the bias from the phase measurements. However, in near-real-time scenario, since the CID and other disturbances are clearly seen in phase measurements, we suggest to calculate the sTEC using solely phase measurements that



**Figure 1.** CID waveform data set. (a) Map showing the event included in the training data set. Details about each event can be found in Table A1. (b–g) vTEC waveforms against time that include a CID arrival (panels b–e, green) and that only contain noise (panels f and g, red). The CID arrival time is shown as a grey vertical line in panels (b)–(e).

can be rapidly retrieved in real-time via the Networked Transport of RTCM via Internet Protocol (NTRIP):

$$sTEC_{ph} = \frac{1}{A} * \frac{f_1^2 * f_2^2}{f_1^2 - f_2^2} * (L_1 * \lambda_1 - L_2 * \lambda_2),$$
(1)

where  $A = 40.308 \text{ m}^3 \text{ s}^{-2}$ ,  $L_1$  and  $L_2$  are phase measurements,  $\lambda_1$  and  $\lambda_2$  are wavelengths at the two Global Positioning System (GPS) frequencies:  $f_1 = 1227$ , 60 and  $f_2 = 1575$ , 42 MHz. Once the sTEC is calculated, the first data value is subtracted from all data series to remove an unknown bias. Finally, because the sTEC is affected by the elevation angle of the LOS, we convert sTEC to vertical TEC (vTEC) by using the standard 'mapping function':

vTEC = sTEC \* cos 
$$\left( \arcsin\left(\frac{R_e \cos\theta}{R_e + H_{\text{ion}}}\right) \right)$$
, (2)

where  $R_e$  is the Earth radius,  $\theta$  is the LOS elevation angle and  $H_{\rm ion}$  is the altitude of ionospheric detection. The  $H_{\rm ion}$  cannot be known because the sTEC is an integral parameter. Based on the physical principles, the  $H_{\rm ion}$  is presumed to be around the ionization maximum, that is around 250–350 km. Here we take  $H_{\rm ion} = 250$  km for all events. This choice is reasonable from the point of view of the ionospheric physics, while determining the real altitude of CID detection is out of the scope of this work. Moreover, once the system is trained, it can detect CID in TEC data series for any  $H_{\rm ion}$  value. The total electron content is measured in TEC units (TECU), with 1 TECU=10<sup>16</sup> electrons m<sup>-2</sup>.

To construct our database, we collected GNSS–TEC data with CID signatures for 12 earthquakes that occurred between 2003 and 2016 (see Fig. 1 and Table A1), including the M6.6 Chuetsu earthquake which is the smallest earthquake ever recorded by ionospheric GNSS data (Cahyadi & Heki 2015). The typical CID waveform are N-shaped and hump signatures (Fig. 1b). However, CID waveforms also depend both on the magnetic field configuration in the epicentral region and on the geometry of the GNSS sounding (Heki & Ping 2005; Astafyeva & Heki 2009; Rolland *et al.* 2013; Bagiya *et al.* 2019). Therefore, in order to correctly represent the large diversity of CID waveforms in our model, we included a variety of

different TEC signatures that could be recorded after an earthquake (examples shown in Figs 1b–e).

The GNSS data used in this study were of 1, 15 and 30 s cadences (see Table A1). Following the NRT-compatible scenario, we did not apply bandpass filter to extract or amplify CID signatures, but only worked with raw relative vTEC.

## 3 AUTOMATIC DETECTION AND ASSOCIATION MODELS

We propose a multistep RF-based procedure to automatically detect CIDs in TEC data series (see Fig. 2): (1) selection of a time window; (2) data pre-processing; (3) waveform features extraction, (4) RF-based classification of inputs features between noise and CID classes; (5) if detection probability >50 per cent at step 4, RF-based arrival time picking; (6) if three successive time windows classified as CID, confirmation of the presence of an arrival and aggregation of arrival times; and (7) if a detection is confirmed at step 6, we then associate this arrival to previously detected CIDs. Finally, we shift the time window and repeat the procedure.

#### 3.1 Pre-processing and feature extraction

To extract consistent waveform features in TEC data with different sampling times, we first downsample all waveforms down to 30 s (see Supporting Information Section S6). Consistency in sampling rate is critical as the higher-frequency spectral content can lead to substantial variations in input features. For example, energy peaks at higher frequencies, that would normally be smoothed out at lower frequencies, can drastically alter the envelope kurtosis and skewness. Additionally, TEC data may contain long-term trends (signals with periods typically greater than 30 min) due to GNSS satellite motion and other long-period TEC changes which can be considered as noise for the problem of CID detection. Therefore, we remove long-term trends by first taking the time derivative of vTEC waveforms to remove long-wavelength trends and then performing a linear de-trending. Derivatives are computed using second-order central differences in the interior points and second-order one-side (forward or backward) differences at the boundaries. Once the TEC waveforms have been pre-processed, we extract 46 features calculated from the vTEC time-series, spectra and spectrograms (see Supporting Information Section S1). These features are commonly used for signal classification tasks (e.g. Hammer et al. 2013; Hibert et al. 2014; Provost et al. 2017; Wenner et al. 2021).

#### 3.2 Building a single-station CID arrival detector

We selected an RF model (Breiman 2001) to discriminate vTEC signals between earthquakes and noise classes. Our RF model takes the features extracted from a given waveform at the previous step as inputs and outputs the probability of this waveform to be signal or noise. An input waveform is classified as CID if the detection probability predicted by the RF is over 50 per cent. RF predictions are constructed from average predictions from an ensemble of individual decision trees. Individual decision trees are built through bootstrap aggregation that consist of randomly selecting input features to train each tree. RFs have excellent generalization abilities and do not require an extensive hyperparameter tuning. We used the 'ExtraTrees' scikit implementation of the RF (Pedregosa *et al.* 2011) which introduces an additional layer of randomness when building decision trees which allow for better generalization of the

training data set (Geurts *et al.* 2006). The training procedure relies on bootstrap samples to build each tree along with out-of-bag samples to estimate the generalization score. Bootstrap aggregation is an iterative procedure where a subset of the training set is randomly selected to train the RF, called in-the-bag set, at each training step. Samples left out at each training step, that is out-of-bag samples, are used to estimate the generalization score. Bootstrapping makes decision trees less sensitive to the choice of training data set which reduces the probability of overfitting. Additionally, the error computed from out-of-bag samples provides an excellent metric for RF's classification performances.

We need to first build a data set of features to train our RF classifier. This data set building process is summarized in Fig. 3. For each station, CID wave trains are described by an arrival time and a duration. Arrival times are selected manually as the time of sudden increase in vTEC amplitudes. Wave train durations are considered uniform across satellites and stations for a given event (see Table A1). Wave train durations are used to automatically label waveforms as CIDs, that is to build our training data set. We consider that a time window contains a CID if it overlaps the true wave train, that is CID confirmed by human analyst, by at least 70 per cent which makes the RF more flexible to detect partial CID waveforms. Values picked for the wave train duration correspond to estimates of the minimum duration of the CID across the network of satellites and stations. This choice ensures that at least the arrival time and/or the time at vTEC maximum are contained in the waveforms. Similar to Ross et al. (2018), we augment our training data set by selecting four time-windows over each CID arrival by randomly shifting the beginning of the time window while still fulfilling the 70 per cent overlap condition. Noise waveforms are selected randomly over each time-series in the data set with the condition that the beginning and end time of the noise window should be at least 30 min away from any CID wave train. Before extracting features, we add artificial Gaussian noise to the waveforms in the training data set to reduce overfitting similar to Mousavi et al. (2020). We add Gaussian noise to each waveform s (for both arrival and noise classes) so that the perturbed waveform  $\overline{s}$  shows a specific signal-to-noise ratio (SNR)  $\overline{s} = s + \sqrt{\frac{\sigma^2}{\text{SNR}}}n$ , where s is the original waveform,  $\sigma^2$  is the variance of the original waveform, *n* is the added noise sampled from a normal distribution and the SNR is picked within the range SNR  $\in (1, 5)$ . Binary classification over an imbalanced training, that is different number of inputs between the two classes, may result in a classifier that is biased towards the majority class, that is the more frequent class (Brodersen et al. 2010). We therefore choose an equal number of CID and noise waveforms in the data set to ensure the balance between true positive and true negative rates. The final data set consists of 2867 CIDs and 2867 randomly picked noise waveforms.

#### 3.3 Building an arrival-time picker

After the classification step, our detection algorithm needs to accurately select the arrival time in each window with a detection probability >50 per cent. This time picking procedure remains challenging using threshold-based conditions such as STA/LTA filters (Allen 1982). False positives will degrade the arrival time estimate when using threshold-based methods since SNR, signal duration and dispersion characteristics vary significantly between events. To overcome this problem, we build an automatic arrival-time picking procedure by using an 'ExtraTrees' RF regressor. Our RF takes a



**Figure 2.** Detection and association procedures described in Section 3: (1) selection of a time window; (2) pre-processing of the waveform; (3) extraction of waveform features from (i) time-series, (ii) spectrum and (iii) spectrogram; (4) RF classification of input waveform; (5) RF arrival time picking; (6) confirmation of an arrival if RF has classified three consecutive time windows (at times  $t^{n-2}$ ,  $t^{n-1}$ ,  $t^n$ ) as arrival; and (7) association of arrivals across different satellites and stations.



**Figure 3.** Building data sets to train our CID classifier and arrival-time picker. Each waveform in our vTEC data set contains information about the CID arrival time and wave train duration. First, four CID windows and four noise windows are extracted from each vTEC waveform. CID windows must overlap the CID wave train by at least 70 per cent while noise windows must start or end at least 1000 s, respectively, after or before the CID wave train. Each window is then pre-processed (derivative and linear detrending) to remove long-term trend. Features are extracted from the pre-processed CID and noise waveforms to build a training data set for our RF classifier with 85 per cent assigned to the training data set and 15 per cent to the validation data set. To build our arrival-time picker RF model, pre-processed CID waveforms are normalized with 85 per cent assigned to the training data set and 15 per cent to the validation data set.

normalized pre-processed waveform as input (see Fig. 2) and outputs offset in seconds from the window central time, that is a float number between -360 and 360. We trigger this arrival time picker only over windows where an arrival has been confirmed.

Similar to the RF classifier, we must build a waveform data set to train our RF arrival-time picker (see Fig. 3). We select arrival window for waveforms that overlaps the true wave train by at least 30 per cent. This overlap is significantly lower than for the detector. This choice aims at training the RF to pick arrival times over the first detection window which generally contains incomplete CID waveforms. Similar to the training of the RF classifier in Section 3.2, we augment our training data set by selecting four time-windows over each CID arrival by randomly perturbing the beginning of the time window while still fulfilling the 30 per cent overlap condition which captures the uncertainty in arrival-time picking. The final data set to train the arrival-time picker consists of 2867 CIDs.

#### 3.4 Confirming a detection on a single station

Because of the natural variability of the ionosphere, false detections can still be present after the RF classification step. These false detections generally correspond to short-time spikes in RF detection probabilities while true detections show an increase in RF detection probabilities over longer time periods. To further remove false positives, we confirm a detection if three consecutive time windows are classified as CIDs. Variations of this value between 2 and 5 have a relatively small (<1 per cent) influence on both recall and precision (see Supporting Information Section S3). Short-time decrease in detection probabilities can occur within long CID wave trains (generally caused by large earthquakes) compared to the processing time window. To reduce the number of false negatives, we determine the end time of an CID wave train when 4 consecutive time windows show a detection probability below 50 per cent.

Once a detection is confirmed, we must determine a single arrival time for the whole wave train. However, predictions in successive windows classified as CIDs and belonging to the same wave train might not have the same predicted time. Therefore, we determine the detected wave train's arrival time by computing the eighth decile of the predicted arrival times over up to 10 successive CID windows. This choice of decile removes the influence of outliers in predicted arrival times made in early detection windows. We do not include predicted arrival times beyond 10 time steps, i.e. 300 s, since these arrivals might correspond to time windows that do not include the true arrival time.

#### 3.5 Associating confirmed detections

Once arrival times are picked across multiple LOS, their spatial distribution informs us about the nature of the detected disturbance. Because large-scale disturbances (e.g. geomagnetic storms, internal gravity waves) or false positives can still pollute the detection data set after the confirmation procedure at step 5, it is critical to discriminate between CIDs and other sources. If the detected signals belong to a CID, arrival times should follow the geometry of the CID wave front, whose geometry is controlled by local sound velocities (Inchin *et al.* 2021). Therefore, the difference in CID arrival times between two detection points cannot be lower than the time it takes an acoustic wave front to propagate between these two detections at the local acoustic velocity. Furthermore, the spatial extent of the CID wave front in the ionosphere is constrained by the dimensions of the activated faults at the ground (Inchin *et al.* 2021) which is

generally below 1000 km. Arrivals detected at two LOS located at large distances from one another (i.e. >1000 km) are not likely to belong to the same CID wave front. By ignoring combinations of detections that show un-realistic travel times, we further improve the quality of our detection data set.

The association procedure is performed on a set of confirmed arrivals and consists of three steps: (1) for new detections  $d_{current}$ , give  $d_{current}$  an unused association number  $s_{current}$ , (2) for each detection  $d_{current}$  find other confirmed detections  $d_{accept}$  among LOS within an acceptable time range from the current detection  $d_{current}$ . By acceptable time range, we consider all arrivals with a time offset from the current detection  $t_{offset} < r_{max}/c_{min}$ , where  $r_{max} = 500$  km is the maximum association range between two detection points, and  $c_{min} = 0.65$  km/s is the minimum horizontal acoustic velocity.  $r_{max}$  is chosen as the maximum possible radius of a CID wave front, and  $c_{min}$  corresponds to the minimum acoustic velocity in the lower ionosphere. Finally, (3) for each detection in an acceptable time range  $d_{accept}$ , if detection has an association number  $s_{accept}$ , change  $s_{current}$  to  $s_{accept}$ .

#### 4 RESULTS

To optimize our ML models for detection and arrival-time picking, we split both data sets between 85 per cent training data and 15 per cent validation data (see Fig. 3). The classifier's validation data set is to calculate confusion matrices and measure the rate of false and true positives which is not accessible when bootstrapping samples. The performance of the classification procedure is sensitive to the window size used for training. In Fig. 4(a), we show both recall and precision metrics for both classes versus the choice of window size. Precision indicates the proportion of true detections relative to all detections (true positives plus false positives). Recall corresponds to the ratio of correct detections over all detections that should have been made (true positives plus false negatives). Because performances are also affected by the choice of overlap threshold used to build the training data set, recall and precision are averaged over four overlaps between 30 and 90 per cent. We observe that there is a clear improvement in both noise precision and arrival recall (up to  $\sim$ 94 per cent) with an increase in window size over the testing data set up to 720 s. This owes to the higher number of incomplete CID wave train for smaller windows than larger ones. For larger time windows >720 s, precision and recall values plateau as the predictive power of some input features computed over large time windows diminishes. We selected a time window of 720 s which gives excellent classification results while facilitating the arrival time picking procedure by decreasing the range of possible values compared to larger time windows.

The RF model can provide an estimate of the relative feature importance through the calculation of the Gini's impurity during training. The three best features (see Fig. 4b) include two time-series features (ratio of the envelope mean over the envelope maximum and the kurtosis of the time-series) as well as a spectral feature (energy up to the Nyquist frequency, i.e. 0.0165 Hz), which differs from other signal classification studies (e.g. Wenner *et al.* 2021). However, the calculation of feature importance can be biased when considering continuous or high-cardinality categorical variables or when inputs features are co-linear. Co-linearity is present in our input data set between spectral and time-series features (see Supporting Information Section S2) which indicates a potential bias in



**Figure 4.** Sensitivity and accuracy of the RF classification step. (a) Precision (prec.) and recall for noise and arrival classes and various window sizes averaged over multiple overlap thresholds: 30, 50, 70 and 90 per cent. The following formulae are used to compute recall and precision for arrival and noise: recall arrival = TP/(TP + FN), recall noise = TN/(TN + FP), precision arrival = TP/(TP + FP) and precision noise = TN/(TN + FN). TP, TN, FP and FN correspond to true positive, true negative, false positive and false negative. The correct detection of a CID corresponds to a TP. (b) Distribution of the three best features against each other. In the diagonal, we show univariate histograms for each feature. Best features are determined during training by calculating the Gini's impurity. W0 corresponds to the ratio of the envelope mean over the envelope maximum, W2 is the kurtosis of the time-series, and S14 is the energy up to the Nyquist frequency, that is 0.0165 Hz. (c) Confusion matrix for the detection model with window size w = 720 s and an overlap of 70 per cent. The confusion matrix is normalized over each row. (d) Arrival-class ROC curve using the detection model with window size w = 720 s. The Area Under Curve (AUC) value is shown above the panel. (e) examples of pre-processed waveforms corresponding to FP (red) and FN (green).

variable importance results. The significant overlap between distributions supports the choice of a large number of features to properly discriminate between each class. Note that this overlap between clusters is also present when using other clustering methods such as such as Principal Component Analysis and t-distributed Stochastic Neighbor Embedding (see Supporting Information Section S2), which further highlights the complexity of this classification problem.

The recall for our detection model, shown in Fig. 4(c), is high for a wide range of probability thresholds indicating that the RF rarely labels true arrivals as noise. We observe in Fig. 4(d) that this value decreases rapidly for probability thresholds >50 per cent corresponding to a stricter classification. A threshold at 50 per cent is a good trade-off to balance true and false positive rates. True and false positives show strong similarities in terms of amplitude and frequency content (see Fig. 4e). However, with larger thresholds, the fall-out, that is the number of false alerts will also decrease. Changes in number of false alerts with variations in probability thresholds highlights that the threshold can be adapted to specific applications depending on the objective. For early warning applications, the number of missed alert should be low and lower thresholds could therefore be used. In contrast, when building arrival-time catalogue to invert for source parameters, precision is key and false alerts should be avoided, which necessitates larger thresholds. Additionally, results indicate that RF outperforms the other analytical methods, including STA/LTA filters, in terms of both true and false positive rates (see Appendix B).

Detection results for a waveform recorded during the 2011 Sanriku earthquake (Fig. 5a) show that both predicted (vertical grey line) and true (reported by human analyst, vertical red line in top panel) arrival times overlap, as the absolute error is low (<3 s). Note that the time used to plot detection probabilities corresponds to the end of the time window used for each classification. We observe that the duration of this wave train (~450 s) is much larger than the true wave train (~200 s), owing to the large time windows employed in our detection model. Outside of the detected wave train, detection probabilities generally remain low (<20 per cent) in accordance to the high true negative rate shown in Fig. 4(c).



**Figure 5.** Performance assessment of arrival-time picking and association steps. (a) 4-h vTEC waveform for the Sanriku event, satellite G07, station 0048 along with RF detection probabilities. The time used to plot probabilities over each window is the window end time. The true arrival is shown as a red vertical line and the RF-predicted arrival time as a dark grey vertical line. The wave train detected by the RF and heuristic models is highlighted with a grey background. (b) box plot of arrival-time picking errors (in s) versus event after 3 min since the first detection window. (c) Evolution of arrival-time picking error versus time delay since first detected window. The red curve shows the average error across all events. Red shaded background shows the  $1^{st}$  to  $3^{rd}$  quartile region computed across the events. (d–f) Tohoku's ionospheric maps with (d) hand-picked arrival times for satellites G05 and G26 along with the epicentre location (yellow star), and surface projection of the fault slip (in m) as green to yellow patches, (e) RF-based arrival-time predictions for each confirmed detection for satellites G05, G26 and G27 with an inset plot showing a newly detected CID arrival (red vertical line) for satellite G27 and station 0167 which was not reported by human analyst and (f) association classes determined from confirmed detections, along with an inset plot showing the vTEC data for satellite G26, station 0155. The vertical lines correspond to the arrival times of the two detections at the station (first is a false detection; second is a true arrival). CID coordinates were calculated at the intersection point between the LOS and the ionospheric layer using  $H_{\rm ion} = 200$  km for lower elevations, and 250 km for higher elevations. These maps are generated 15 min after the event.

In addition to the classification of individual waveform snippets, accurate arrival times are crucial for NRT applications. We assess our model's arrival-time picking accuracy by computing the error between predicted and true arrival times. Arrival-time errors for each event in our CID data set in Fig. 5(b) indicate that most arrivals (~95 per cent) are captured with an absolute error <60s, that is less than two time steps, and a large proportion of arrivals ( $\sim$ 80 per cent) are accurately reproduced with an absolute error <30 s, which is below the sampling time in each CID waveform. Some outliers are present for both Illapel and Kaikoura events. Errors for the Kaikoura earthquake owe primarily to the high noise level in the waveforms (i.e. random fluctuations of TEC background) which leads to large variations in vTEC time derivatives. For Illapel, false positives are lumped together with the true detection windows and degrade the arrival-time picking performance over 4 time steps. However, the average arrival-time picking error across the whole data set decreases significantly as the number of time steps increases, that is time since first detection (see Fig. 5c).

Confirmed detections for multiple LOS can be used to plot ionospheric maps for each event. The location of the earliest CID arrivals reported by human analysts, that is first CID arrivals (around 7 min for example after the Tohoku earthquake in Fig. 5d), should be the closest to the distribution of maximum co-seismic slip at the surface (Astafyeva et al. 2013a). In Fig. 5(d), we observe a slight shift of these first arrivals after the Tohoku earthquake to the south east of the region of maximum surface slip owing primarily to our choice of altitude of detection  $H_{ion}$  (Kakinami *et al.* 2021). In Fig. 5(d), we note that the first CID arrivals are distributed linearly from location (36°N, 144°E) to (39°N, 145°E) matching the trend of maximum co-seismic slip distribution at the surface. Comparing Tohoku's ionospheric images from human analyst picks in Fig. 5(d) and from our detection algorithm before association in Fig. 5(e), we observe that the spatial distribution of CID arrival times is accurately reproduced by our ML model. Some spurious arrivals are present in Fig. 5(e), west of the fault with early arrival times, and southeast of the fault with late arrival times. These false detections correspond to rapid changes in vTEC occurring more than 20 min before or after the true arrival and classified as earthquake signals by our model.

Our association procedure enables the discrimination between detections belonging to the same wave front and spurious arrivals. The distribution of association classes for the confirmed detections is shown in Fig. 5(f). Owing to the large time difference between spurious arrivals and the true arrivals, false detections are correctly classified in different association classes (see first vertical dark purple line in the inset plot in Fig. 5f). Note that the location of the ionospheric detection points varies from the first to the second detection at satellite G05 (inset plot in Fig. 5f) since the satellite moves with time. The time evolution of the distribution of confirmed arrivals (see Supporting Information Section S5) indicates that the entirety of the true arrivals were detected within 15 min after the event. Note that the position of ionospheric detection points is dependent on the altitude of detection  $H_{ion}$ , which could impact the association classes. However, while changing  $H_{ion}$  from 180 to 250 km for Tohoku affects the location of the ionospheric points, true CID arrivals are still correctly associated within the same class (see Supporting Information Section S7).

New detections, that is arrivals not picked by human analysts, have also been reported by our model west of the epicentre (Figs 5d and e) for the largest class corresponding the true CID (inset plot in Fig. 5e and light purple class in Fig. 5f). A low SNR pulse is visible after the predicted arrival time (red vertical line in the inset plot of Fig. 5e) at t = 9.9 min after the earthquake, which is

consistent with acoustic travel time from the source highlighted by other studies (e.g. Astafyeva *et al.* 2013a). Using our model also ensures consistency in the choice of arrival times, in contrast to human analysts who introduce a subjective uncertainty range when determining the true onset.

In order to further assess the ability of our model to detect arrivals on new unseen data, we processed waveforms recorded after the 2014 Iquique earthquake (see Table A1). In Fig. 6(a), we show the slip distribution of the Iquique earthquake along with the RF predicted arrivals times and association classes in Figs 6(b) and (c). Predicted arrival times are coherent with the region of maximum slip at the surface despite a few false detections south of the fault. This confirms the excellent detection, arrival-time picking and classification results on new data.

#### 5 DISCUSSION

Monitoring procedures NRT-compatible require both high accuracy and low computational time. To provide an estimate of our algorithm's computational time, we show in Fig. 7 the cost associated with detection, arrival-time picking, and association steps after the 2011 Tohoku event at station 0908 and satellite G05 (Fig. 7a) on a single CPU (Dell T5610 Intel Xeon E5-2630 v2 2.6Ghz 6 CPUs 64GB RAM on CentOS 7). The computational time for feature extraction, classification, validation and time picking for a single satellite/station pair is always below 1 s and is dominated by RF steps (Fig. 7b). This result suggests that a similar detection methodology, trained with higher sampling-rate data, could be implemented for NRT applications up to 1 Hz. Note that the time picking step is only present when a detection occurs which explains the jump in computational cost around 7 min after the earthquake.

We observe a significant increase in computational cost across the network 9 min after the earthquake in Fig. 7(c). This jump in association cost corresponds to the earthquake-induced acoustic wave reaching the ionosphere which leads to a large number of detections at each combination of satellite/station (see Fig. 7d). This association procedure is computationally expensive since it must scan through all possible neighbors of each new detection to update association classes, which scales linearly with the number of new detections. Yet, the maximum cost for one time step over the whole network is less than 6 s. It takes around 1 s to process 10 new detections, at a given time, over a network of about 100 satellites/stations. The number of associated detections reaches a plateau about 13 min after the earthquake (see Fig. 7e) which corresponds to the end of the association of all first CID arrivals.

The practical implementation of our detection/association procedure will require an efficient internet between the relevant GNSS stations to collect and extract time-series for classification in NRT. However, because the overall computational cost of one time iteration using our method is below 6 s on a single CPU using noncompiled Python codes, at least 24 s are available for data acquisition and processing with waveforms sampled at 30 s. The association step is currently the most costly ( $\sim$ 90 per cent of the total cost) but can be run in parallel to the other detection steps. Note that we also explored the feasibility of using our model to detect CIDs at a higher sampling rate by extracting input features without downsampling input data (see Supporting Information Section S6). Our RF detection model always shows detection probabilities >50 per cent using a 1 s sampling time but still predict a strong increase in detection probability around the CID arrival. This suggests that increasing the detection threshold to higher values (e.g. from 50 per cent to



Figure 6. Ionospheric maps for the 2014 Iquique earthquake. (a) Map showing the epicentre location (yellow star) and surface projection of the fault slip (in m) as green to yellow patches. (b) CID detections using our RF-based classifier and time picker, and (c) association classes determined from confirmed detections. CID coordinates were calculated at the intersection point between the LOS and the ionospheric layer using  $H_{ion} = 250$  km. These maps can be generated 15 min after the event.



**Figure 7.** Computational cost associated with detection, arrival-time picking, and association steps after the 2011 Tohoku earthquake. (a) vTEC time-series for satellite G07 and station 0048. (b) Stack plot of computational time (s) for pre-processing and feature extraction (green), RF classification (orange), RF arrival-time picking (blue) and confirmation (pink) steps. (c) Computational cost (s) at each time iteration of the association procedure. (d) Number of new detections per time iteration.

70 per cent) would enable implementation of our detection method at higher sampling-rates at the cost of a higher false positive likelihood.

Our model seems to be also able to detect vTEC variations associated with other traveling ionospheric disturbances (TIDs) such as volcanic explosions, Rayleigh waves, and tornadoes (see Supporting Information Section S8). This suggests that a data set of TID waveforms should be built to train an efficient discriminator between background noise, earthquake, and other TID phases. However, the discrimination between TEC signals from seismic origin and TIDs can easily be done by comparing the predicted arrival times at the ionospheric points to the distribution of seismic events in seismic catalogues which are available in NRT (e.g. Thompson *et al.* 2019).

#### 6 CONCLUSIONS

We introduced an automatic procedure for detection, arrival-time picking, and association of CIDs. Detection and arrival time picking steps are performed using random forests trained over a CID data set built from 12 earthquake events. These methods show excellent classification results with 96 per cent true positive rate and 96 per cent true negative rate, and arrival-time accuracy with an average error <20 s using a 120 s time delay since the first detection window. Our model also outperforms threshold-based detection methods in terms of both recall and precision. Our analytical classification procedure accurately associates all arrivals corresponding to the same wave front. Classification results also indicate that low SNR arrival that were not picked by human analysts could also captured by our RF detection model.

The performance of our automated procedure is promising for future NRT applications, including the use of CID arrival times for construction of ionospheric images of seismic sources. The first demonstration of seismo-ionospheric imagery was based on retrospective analysis of CID generated by the 2011 Tohoku earthquake (Astafyeva *et al.* 2011, 2013a). Here we show that our newly developed method can generate such images in NRT. Note that the position of ionospheric detection points is dependent on the altitude of detection  $H_{ion}$ . The latter parameter is not known precisely, but it is presumed to be around the height of ionospheric ionization maximum, that is around 250–350 km, depending on solar, geomagnetic, seasonal and diurnal conditions. Future studies should focus on development of real-time compatible methods of determining the true  $H_{ion}$  in order to obtain accurate source locations in NRT.

Acquiring labeled vTEC data from additional events which will significantly improve the generalization abilities of our RF models. Additionally, the choice of features made in this paper could be further refined to obtain better accuracy (Han & Kim 2019). Building a more accurate RF classifications could alleviate the need for a validation step presented in Section 3.4. However, RF memory costs increase exponentially with tree depth, and consequently data set size,  $\sim 2^D$ , with D the tree depth (Louppe 2014; Solé *et al.* 2014). The RF classification model is only about 70 mb but will grow considerably larger with new data. With a larger data set, image segmentation ML techniques such as standard convolutional neural networks (Ross et al. 2018, 2019), transformers (Mousavi et al. 2020) or residual networks (Mousavi et al. 2019) applied on nonengineered inputs such as spectrograms could lead to substantial improvements in accuracy and memory costs for both classification and arrival time picking steps. Finally, both detection performances and computational cost could be improved by training our ML model using higher sampling-rate ionospheric data such as 1 Hz data available for some GNSS receivers. Higher-frequencies input data might enable both the detection of smaller-magnitude events such as the Chuetsu earthquake (Cahyadi & Heki 2015).

The proposed association algorithm does not incorporate any information about the source nor the atmospheric dynamics. This procedure could be improved by assessing the consistency of arrival time differences across a network of satellites and stations using a range of possible sources, similarly to the methods used for the automated production of seismic bulletins (Draelos *et al.* 2015). In contrast to seismic media, atmospheric velocities, that is winds, are time-dependent which introduces further complexity when computing theoretical source–receiver arrival times. Fast simulations of acoustic wave propagation up to the ionosphere with realistic atmospheric specifications would greatly improve the classification between true and false arrivals and enable the localization of the largest surface displacements (Bagiya *et al.* 2019; Inchin *et al.* 2021; Zedek *et al.* 2021). Finally, to confirm the detection of an earthquake across a given network and trigger an alert for human analysts, an additional heuristic could be implemented based, for example, on the number of detections per association class.

#### ACKNOWLEDGMENTS

This work was supported by the French Space Agency (CNES, Project 'RealDetect').

#### DATA AVAILABILITY

GNSS data are available from the following web-services: Japan GNSS Earth Observation System, GEONET (http://datahouse1.gsi .go.jp/terras/terras\_english.html), GEONET Geological Hazard Information for New Zealand (https://www.geonet.org.nz), Scripps Orbit and Permanent Array Center (SOPAC, http://sopac-old.ucsd .edu/dataBrowser.shtml), National Seismological Centre, University of Chile (http://gps.csn.uchile.cl). Finite-fault data were downloaded from the US Geological Survey website (https://earthquake .usgs.gov/earthquakes). RF evaluation, validation, and associations codes are available at https://github.com/QuentinBrissaud/AIDE. Data and RF models are available at https://doi.org/10.6084/m9.fig share.19661115.

#### REFERENCES

- Afraimovich, E., Astafyeva, E. & Kiryushkin, V., 2006. Localization of the source of ionospheric disturbance generated during an earthquake, *Int. J. Geomagn. Aeronomy*, 6, GI2002, doi:10.1029/200403000092.
- Allen, R., 1982. Automatic phase pickers: their present use and future prospects, *Bull. seism. Soc. Am.*, 72(6B), S225–S242.
- Astafyeva, E., 2019. Ionospheric detection of natural hazards, *Rev. Geophys.*, **57**, 1265–1288.
- Astafyeva, E. & Heki, K., 2009. Dependence of waveform of near-field coseismic ionospheric disturbances on focal mechanisms, *Earth Planets Space*, 61, 939–943.
- Astafyeva, E., Heki, K., Afraimovich, E., Kiryushkin, V. & Shalimov, S., 2009. Two-mode long-distance propagation of coseismic ionosphere disturbances, *J. geophys. Res.*, **118**, A10307, doi:10.1029/2008JA013853.
- Astafyeva, E., Lognonné, P. & Rolland, L.M., 2011. First ionosphere images for the seismic slip on the example of the Tohoku-oki earthquake, *Geophys. Res. Lett.*, **38**, L22104 , doi:10.1029/2011GL049623.
- Astafyeva, E., Rolland, L.M., Lognonné, P., Khelfi, K. & Yahagi, T., 2013a. Parameters of seismic source as deduced from 1 Hz ionospheric GPS data: case-study of the 2011 Tohoku-oki event, *J. geophys. Res.*, **118**, 5942–5950.
- Astafyeva, E., Rolland, L.M. & Sladen, A., 2014. Strike-slip earthquakes can also be detected in the ionosphere, *Earth planet. Sci. Lett.*, 405, 180–193.
- Astafyeva, E., Shalimov, S., Olshanskaya, E. & Lognonné, P., 2013b. Ionospheric response to earthquakes of different magnitudes: larger quakes perturb the ionosphere stronger and longer, *Geophys. Res. Lett.*, 40, 1675– 1681.
- Astafyeva, E. & Shults, K., 2019. Ionospheric GNSS imagery of seismic source: possibilities, difficulties, and challenges, *J. geophys. Res.*, **124**(1), 534–543.
- Bagiya, M.S., Sunil, A., Rolland, L., Nayak, S., Ponraj, M., Thomas, D. & Ramesh, D.S., 2019. Mapping the impact of non-tectonic forcing mechanisms on gnss measured coseismic ionospheric perturbations, *Sci. Rep.*, 9(1), 1–15.
- Bagiya, M.S., Sunil, P.S., Sunil, A.S. & Ramesh, D.S., 2018. Coseismic contortion and coupled nocturnal ionospheric perturbations during 2016 Kaikoura, M<sub>w</sub> 7.8 New Zealand earthquake, J. geophys. Res., 123(2), 1477–1487.

- Bagiya, M.S., Thomas, D., Astafyeva, E., Bletery, Q., Lognonné, P. & Ramesh, D.S., 2020. The ionospheric view of the 2011 Tohoku-oki earthquake seismic source: the first 60 seconds of the rupture, *Sci. Rep.*, 10(5232), doi:10.1038/s41598-020-61749-x.
- Belehaki, A. et al., 2020. An overview of methodologies for real-time detection, characterisation and tracking of traveling ionospheric disturbances developed in the techtide project, J. Space Weather Space Clim., 10, 42, doi:10.1051/swsc/2020043.
- Bessason, B., Eiríksson, G., Thorarinsson, Ó, Thórarinsson, A. & Einarsson, S., 2007. Automatic detection of avalanches and debris flows by seismic methods, *Journal of Glaciology*, 53(182), 461–472.
- Breiman, L., 2001. Random forests, Mach. Learn., 45(1), 5-32.
- Brodersen, K.H., Ong, C.S., Stephan, K.E. & Buhmann, J.M., 2010. The balanced accuracy and its posterior distribution, in 2010 20th International Conference on Pattern Recognition, pp. 3121–3124, IEEE, Istanbul, Turkey.
- Cahyadi, M.N. & Heki, K., 2015. Coseismic ionospheric disturbance of the large strike-slip earthquakes in north Sumatra in 2012 M<sub>w</sub> dependence of the disturbance amplitudes, *Geophys. J. Int.*, 200(1), 116–129.
- Curilem, G., Vergara, J., Fuentealba, G. & Acuña, G., 2009. Classification of seismic signals at Villarrica volcano (Chile) using neural networks and genetic algorithms, 180(1), 1–8, doi:10.1016/j.jvolgeores.2008.12.002.
- Draelos, T.J., Ballard, S., Young, C.J. & Brogan, R., 2015. A new method for producing automated seismic bulletins: probabilistic event detection, association, and location, *Bull. seism. Soc. Am.*, **105**(5), 2453–2467.
- Efendi, E. & Arikan, F., 2017. A fast algorithm for automatic detection of ionospheric disturbances: DROt, *Adv. Space Res.*, **59**(12), 2923–2933.
- Geurts, P., Ernst, D. & Wehenkel, L., 2006. Extremely randomized trees, Mach. Learn., 63(1), 3–42.
- Hammer, C., Beyreuther, M. & Ohrnberger, M., 2012. A seismic-event spotting system for volcano fast-response systems, *Bulletin of the Seismological Society of America*, **102**(3), 948–960, doi.org/10.1785/0120110167.
- Hammer, C., Ohrnberger, M. & Faeh, D., 2013. Classifying seismic waveforms from scratch: a case study in the alpine environment, *Geophys. J. Int.*, **192**(1), 425–439.
- Han, S. & Kim, H., 2019. On the optimal size of candidate feature set in random forest, *Appl. Sci.*, **9**(5), 898, doi:10.3390/app9050898.
- Heki, K., 2006. Explosion energy of the 2004 eruption of the Asama Volcano, central Japan, inferred from ionospheric disturbances, *Geophys. Res. Lett.*, 33, L17101, doi:10.1029/2006GL026249.
- Heki, K., 2021. Ionospheric disturbances related to earthquakes, in *Ionospheric Dynamics and Applications*, pp. 511–526, eds Huang, C., Lu, G., Zhang, Y. & Paxton, L.J., American Geophysical Union.
- Heki, K., Otsuka, Y., Choosakul, N., Hemmakorn, N., Komolmis, T. & Maruyama, T., 2006. Detection of ruptures of andaman fault segments in the 2004 great Sumatra earthquake with coseismic ionospheric disturbances, *J. geophys. Res.*, **111**, B09313 , doi:10.1029/2005JB004202.
- Heki, K. & Ping, J., 2005. Directivity and apparent velocity of the coseismic ionospheric disturbances observed with a dense GPS array, *Earth planet. Sci. Lett.*, 236(3), 845–855.
- Hibert, C. *et al.*, 2014. Automated identification, location, and volume estimation of rockfalls at Piton de la Fournaise volcano, *J. geophys. Res.*, 119(5), 1082–1105.
- Hofmann-Wellenhof, B., Lichtenegger, H. & Wasle, E., 2008. GNSS-Global Navigation Satellite System, Springer.
- Inchin, P., Snively, J., Kaneko, Y., Z. D., M. & Komjathy, A., 2021. Inferring the evolution of a large earthquake from its acoustic impacts on the ionosphere, *AGU Adv.*, 2, doi:10.1029/e2020AV000260.
- Kakinami, Y., Saito, H., Yamamoto, T., Chen, C.-H., Yamamoto, M., Nakajima, K., Liu, J.-Y. & Watanabe, S., 2021. Onset altitudes of co-seismic ionospheric disturbances determined by multiple distributions of gnss tec after the foreshock of the 2011 Tohoku earthquake on march 9, 2011, *Earth Space Sci.*, doi:10.1029/2020EA001217.
- Kamogawa, M., Orihara, Y., Tsurudome, C., Tomida, Y., Kanaya, T., Ikeda, D., et al., 2016. A possible space-based tsunami early warning system using observations of the tsunami ionospheric hole, *Sci. Rep.*, 6, 37989, doi:10.1038/srep37989.

- Katsumata, A., Ueno, H., Aoki, S., Yasushiro, Y. & Barrientos, S., 2013. Rapid magnitude determination from peak amplitudes at local stations, *Earth Planets Space*, 65, 843–853.
- Komjathy, A., Yang, Y., Meng, X., Vekhoglyadova, O., Mannucci, A. & Langley, R., 2016. Review and perspectives: understanding natural-hazardsgenerated ionospheric perturbations using GPS measurements and coupled modeling, *Radio Sci.*, **51**, 951–961.
- LaBrecque, J., Rundle, J., Bawden, G., Surface, E. & Area, I.F., 2019. Global navigation satellite system enhancement for tsunami early warning systems, *Global Assessment Report on Disaster Risk Reduction*.
- Lee, R., Rolland, L. & Mykesell, T., 2018. Seismo-ionospheric observations, modeling and backprojection of the 2016 Kaikoura earthquake, *Bull. seism. Soc. Am.*, **108**(3B), 1794–1806.
- Louppe, G., 2014. Understanding random forests: from theory to practice, preprint (arXiv:1407.7502).
- Maletckii, B. & Astafyeva, E., 2021. Determining spatio-temporal characteristics of coseismic travelling ionospheric disturbances (CTID) in near real-time, *Sci. Rep.*, **11**, doi:10.1038/s41598-021-99906-5.
- Manta, F., Occhipinti, G., Feng, L. & Hill, E., 2020. Rapid identification of tsunamigenic earthquakes using gnss ionospheric sounding, *Sci. Rep.*, 10, 11054, doi:10.1038/s41598-020-68097-w.
- Mousavi, S.M., Ellsworth, W.L., Zhu, W., Chuang, L.Y. & Beroza, G.C., 2020. Earthquake transformer—an attentive deep-learning model for simultaneous earthquake detection and phase picking, *Nat. Commun.*, 11(1), 1–12.
- Mousavi, S.M., Zhu, W., Sheng, Y. & Beroza, G.C., 2019. CRED: a deep residual network of convolutional and recurrent units for earthquake signal detection, *Sci. Rep.*, 9(1), 1–14.
- Occhipinti, G., Aden-Antoniow, F., Bablet, A., Molinie, J.-P. & Farges, T., 2018. Surface waves magnitude estimation from ionospheric signature of Rayleigh waves measured by Doppler sounder and OTH radar, *Sci. Rep.*, 8, 1555, doi:10.1038/s41598-018-19305-1.
- Pedregosa, F. et al., 2011. Scikit-learn: machine learning in Python, J. Mach. Learn. Res., 12, 2825–2830.
- Provost, F., Hibert, C. & Malet, J.-P., 2017. Automatic classification of endogenous landslide seismicity using the random forest supervised classifier, *Geophys. Res. Lett.*, 44(1), 113–120.
- Rakoto, V., Lognonné, P., Rolland, L. & Coisson, P., 2018. Tsunami wave height estimation from GPS-derived ionospheric data, *J. geophys. Res.*, 123, 4329–4348.
- Ravanelli, M., Occhipinti, G., Savastano, G., Komjathy, A., Shume, E.B. & Crespi, M., 2021. GNSS total variometric approach: first demonstration of a tool for real-time tsunami genesis estimation, *Sci. Rep.*, **11**(1), 1–12.
- Rolland, L., Vergnolle, M., Nocquet, J.-M., Sladen, A., Dessa, J.-X., Tavakoli, F., Nankali, H. & Cappa, F., 2013. Discriminating the tectonic and non-tectonic contributions in the ionospheric signature of the 2011, *M*<sub>w</sub>7.1, dip-slip van earthquake, eastern Turkey, *Geophys. Res. Lett.*, 40, doi:10.1002/grl.50544.
- Rolland, L.M., Occhipinti, G., Lognonné, P. & Loevenbruck, A., 2016. Ionospheric gravity waves detected offshore Hawaii after tsunami, *Geophys. Res. Lett.*, **37**, L17101, doi:110.1029/2010GL044479.
- Ross, Z.E. *et al.*, 2019. Hierarchical interlocked orthogonal faulting in the 2019 Ridgecrest earthquake sequence, *Science*, **366**(6463), 346–351.
- Ross, Z.E., Meier, M.-A. & Hauksson, E., 2018. P wave arrival picking and first-motion polarity determination with deep learning, J. geophys. Res., 123(6), 5120–5129.
- Shults, K., Astafyeva, E. & Adourian, S., 2016. Ionospheric detection and localization of volcano eruptions on the example of the april 2015 calbuco events, *J. geophys. Res.*, **121**(10), 10 303–10 315.
- Solé, X., Ramisa, A. & Torras, C., 2014. Evaluation of random forests on large-scale classification problems using a bag-of-visual-words representation, in *Artificial Intelligence Research and Development*, pp. 273–276, eds Museros, L., Pujol, O. & Agell, N., doi:10.3233/978-1-61499-452-7-273.
- Thomas, D. et al., 2018. Revelation of early detection of co-seismic ionospheric perturbations in GPS-TEC from realistic modelling approach: case study, *Sci. Rep.*, 8(1), 1–10.

- Thompson, E.M. et al., 2019. USGS near-real-time products—and their use—for the 2018 Anchorage earthquake, Seismol. Res. Lett., 91(1), 94– 113.
- Tsai, H.-F., Liu, J.-Y., Lin, C.-H. & Chen, C.-H., 2011. Tracking the epicenter and the tsunami origin with GPS ionosphere observation, *Earth Planets Space*, 63, 859–862.
- Van der Maaten, L. & Hinton, G., 2008. Visualizing data using t-SNE, Journal of machine learning research, 9, https://www.jmlr.org/papers/vo lume9/vandermaaten08a/vandermaaten08a.pdf.
- Wenner, M., Hibert, C., van Herwijnen, A., Meier, L. & Walter, F., 2021. Near-real-time automated classification of seismic signals of slope failures with continuous random forests, *Nat. Hazards Earth Syst. Sci.*, 21(1), 339–361.
- Wright, T., Houlie, N., Hildyard, M. & Iwabuchi, T., 2012. Real-time, reliable magnitudes for large earthquakes from 1 Hz GPS precise point positioning: the 2011 Tohoku-oki (Japan) earthquake, *Geophys. Res. Lett.*, 38(L12302), doi:10.1029/2012/GL051894.
- Zedek, F., Rolland, L.M., Dylan Mikesell, T., Sladen, A., Delouis, B., Twardzik, C. & Coïsson, P., 2021. Locating surface deformation induced by earthquakes using GPS, GLONASS and Galileo ionospheric sounding from a single station, *Adv. Space Res.*, 68, 3403–3416.

#### SUPPORTING INFORMATION

Supplementary data are available at GJI online.

**Figure S1**. Probability density of each input feature over our training and testing data sets. The short name of the feature for each plot is shown above the plot. The description of each feature is given in Table S1.

**Figure S2**. Spearman's correlation coefficients between each feature used for training. A description of each feature is given in Table S1.

**Figure S3**. First versus second component of (a and c) a principal component analysis (PCA) and (b and d) a T-distributed stochastic neighbour embedding (TNSE; Van der Maaten & Hinton 2008). Points are colour-coded with (a and b) the detection class and (c and d) the event name for the arrival class.

**Figure S4**. True positive rate (TPR), true negative rate (TNR), false positive rate (FPR) and false negative rate (FNR) with the choice of number of time steps for validation in the heuristic model presented in Section 3.4

**Figure S5.** Performance of RF arrival time picker. (a) Root mean square error (RMSE) versus minimum true-wave train overlap (deviation) and window size (s). The minimum true wave train overlap corresponds to the minimum fraction of the wave train that has to be included in a window to be considered for training. (b) R2 error versus minimum true wave train overlap (deviation) and window size (s). Bottom distribution of arrival-time picking errors (s) versus true time-shift from central time (s) over (c) the testing data set and (d) the training data set.

**Figure S6.** Ionospheric maps after the 2011 Tohoku earthquake generated at various times since the event. (a–c) Distribution of detected arrival times after (a) 7, (b) 11 and (c) 15 min since the event. CID coordinates were calculated at the intersection point between the LOS and the ionospheric layer using  $H_{ion} = 250$  km. The colour code corresponds to the predicted arrival time at each ionospheric point. Grey dots correspond to the location of ionospheric points where there is no detection yet but with detections after 20 mn.

**Figure S7**. Tohoku's ionospheric arrival-time maps computed 14 min after the event for (d) hand-picked arrival times along with the epicentre location (yellow star), and surface projection of the fault slip (in m) as green to yellow patches, (e) RF-based arrival-time predictions and (f) association classes determined from predicted

arrival times. CID coordinates were calculated at the intersection point between the LOS and the ionospheric layer using  $H_{\rm ion} = 180$  km.

**Figure S8**. Performance assessment of RF detection and arrivaltime picking at a higher sampling rate of 1 s. 2-h vTEC waveform for the Sanriku event, satellite G07, station 0048 along with detection probabilities predicted by our RF detection model (bottom). The true arrival is shown as a red vertical line.

**Figure S9**. vTEC waveform for the Calbuco eruption, satellite G03, station antc along with detection probabilities predicted by our detection procedure (see Section 3) using a window size w = 720 s. Volcano-associated ionospheric perturbations are present between 21.3 and 22.5UT. The RF-predicted arrival time as a dark grey vertical line. The detected wave train using the RF is highlighted with a grey background.

**Figure S10**. vTEC waveform from seismic Rayleigh waves recorded after the 1994 earthquake in Kuril Islands (Astafyeva *et al.* 2009), satellite G06, station tskb along with detection probabilities predicted by our detection procedure using a window size w = 720 s. Rayleigh-wave-associated ionospheric perturbations are present between 13.6UT and 13.8UT. The RF-predicted arrival time as a dark grey vertical line. The detected wave train using the RF is highlighted with a grey background.

**Figure S11**. vTEC waveforms high-passed over  $1.5e^{-4}$  Hz extracted (a) on 2013 May 19, the day before the 2013 Moore EF5 tornado, and (b) on 2013 May 20, the day of the tornado, from stations hees and lesv in the United States and satellites G28, G08, and R11 ordered by distance between the ionospheric detection points and the city of Moore, Oklahoma, U.S. (numbers shown on the right of the plot). Background colours behind the waveforms represent the detection probability computed by the ML detection model.

**Table S1.** List of attributes. Nyf = 0.0165 Hz is the Nyquist frequency. These attributes are commonly used in signal-classification studies. We refer the reader to the following references for more details: Bessason *et al.* (2007), Curilem *et al.* (2009), Hammer *et al.* (2012), Hibert *et al.* (2014), Provost *et al.* (2017) and Wenner *et al.* (2021).

Please note: Oxford University Press is not responsible for the content or functionality of any supporting materials supplied by the authors. Any queries (other than missing material) should be directed to the corresponding author for the paper.

#### APPENDIX A: LIST OF EVENTS

The list of events compiled in our CID data set is described in Table A1.

#### APPENDIX B: COMPARISON OF RF-BASED METHOD TO ANALYTICAL DETECTORS

To further assess the RF classification performance, we compare the results to two analytical detection methods: (1) an STA/LTA detection method and (2) a derivative-based threshold method. The STA/LTA method requires to set four parameters: the STA and LTA time windows and two thresholds to activate and deactivate the detection trigger. The STA window represents the average duration of expected earthquake signals while the LTA window captures the average TEC noise amplitude. The STA/LTA method employed here uses a 60 s STA window and a 400 s LTA window. A detection is triggered if the STA/LTA threshold reaches 2.5 while the end of

Table A1. List of events included in the data set. Events are sorted by magnitude.

Event			Date	Time	Min. signal		
Reference	Mag.	Lat.; Lon.	(DD/MM/YY)	(UTC)	duration (s)	Sat.	Samp
Tohoku Astafyeva <i>et al</i> .	9.1 . (2011, 201	38.3; 142.37 (3a)	11/03/2011	05:46:23	800	G26 G05	1s, 30s
Sumatra 1 Astafyeva <i>et al.</i>	8.6 . (2014)	2.35; 92.8	11/04/2012	08:38:37	300	G32	15s
Tokachi Heki & Ping (2	8.3 005)	41.78; 143.90	25/09/2003	19:50:06	440	G13 G24	30s
<b>Illapel</b> bagiya2019map	8.3 oping	-31.57; -71.61	16/09/2015	22:54:32	600	G25,G12 G24	15s, 30s
Sumatra 2 Astafyeva <i>et al.</i>	8.2 . (2014)	0.90; 92.31	11/04/2012	10:43:09	300	G32	15s
<b>Iquique</b> Bagiya <i>et al.</i> (2	8.2 019)	-19.61; -70.77	01/04/2014	23:46:47	700	G01,G20 G23	15s, 30s
Macquarie Astafyeva <i>et al.</i>	8.1 . (2014)	-49.91; 161.25	23/12/2004	14:59:03	550	G05	30s
<b>Fiordland</b> Astafyeva <i>et al.</i>	7.8 (2013b)	-45.75; 166.58	15/07/2009	09:22:29	300	G20	30s
Kaikoura Bagiya <i>et al.</i> (2	7.8 018)	42.757; 173.077	13/11/2016	11:02:56	550	G20 G29	1s, 30s
<b>Sanriku</b> Thomas <i>et al.</i> (	7.3 2018); Asta	38.44; 142.84 fyeva & Shults (2019)	09/03/2011	02:45:20	200 G08	G07, G10	1s, 30s
Kii Heki & Ping (2	7.2 005)	33.1; 136.6	05/09/2004	10:07:07	425	G15	30s
<b>Chuetsu</b> Cahyadi & Hek	6.6 ti ( <mark>2015</mark> )	37.54; 138.45	16/07/2007	01:12:22	300	G26	30s

a wave train is chosen where the threshold goes below 0.5. This trigger value of 2.5, lower than employed at seismic stations, is used to make sure we capture each arrival, that is to increase the true positive rate. Parameters are chosen empirically and could be improved with a thorough investigation of the STA/LTA accuracy over the whole data set. However, fine tuning the hyperparameters increases the likelihood of overfitting a specific data set. This shows the advantage of using an ML-based approach that relies on an efficient optimization procedure enabling us to reach high accuracy without strong overfitting.

The analytical method used for comparison, referred to as 'AN', is based on the analysis of TEC rate-of-change. Maletckii & Astafyeva (2021) noted that, in a majority of cases, the CIDs are characterized by a rapid and high increase of TEC. To capture the CID arrival, we therefore suggest to analyse the rate of TEC change between the two consecutive epochs, between every two and every three epochs:

$$\partial vTEC_1 = |vTEC_i - vTEC_{i+1}|,$$
 (B1)

$$\partial v TEC_2 = |v TEC_i - v TEC_{i+2}|,$$
(B2)

$$\partial v TEC_3 = |v TEC_i - v TEC_{i+3}|,$$
 (B3)

$$\partial vTEC_4 = |vTEC_i - vTEC_{i+4}|,$$
 (B4)

where the subscript *i* corresponds to the time step  $t_i$ . The vTEC at epoch *i* is considered as the CID arrival if each slope  $\partial vTEC_1$ ,

 $\partial v TEC_2$  and  $\partial v TEC_3$  (and  $\partial v TEC_4$  for 1s data) are greater than the thresholds shown in Table B1. These threshold values were determined analytically over multiple events. Detections are confirmed if 12 consecutive time steps fulfil the threshold conditions described in Table B1.

To assess the performance of each method, we determine the false and true negative and positive rates over the waveforms included in the testing data set. To provide meaningful results, we scan entire waveforms (from 1-hr to 2-hr duration) instead of a few windows as done for RF training. Including entire waveforms means that more noise windows will be included than CID windows, which is an excellent test to assess the performance of each method in more realistic conditions (where CIDs are rare). We consider that a wave train, that is a time window characterized by an arrival time and a duration, classified as CID by any method is a true positive if it overlaps the true arrival by at least 70 per cent.

Our RF-based detection method outperforms AN and STA in terms of true positive and negative rates (see Fig. A1). We observe a lower true negative rate than determined during the RF validation step (see Fig. 4c). This owes to the presence of much larger number of noise windows in the data set. The STA/LTA filter also performs well to detect true arrivals. However, this high true positive rate comes at the cost of a low false positive rate, that is a large number of false alerts. The analytical method using only local time derivatives shows a large number of false negatives owing to presence of noise in the data.



Table B1. Slope parameters for different sampling rates used by the analytical detector AN.

**Figure A1.** Confusion matrices calculated over the RF testing data set consisting of 1–2-hr long waveforms for (a) the RF classification model, (b) the analytical time-derivative based model and (c) the STA/LTA filter. Confusion matrices show from top to bottom and left to right, the TPR, FPR, FNR and TNR, such that: