

Deep-learning-based phase picking for volcano seismicity

Yiyuan Zhong¹ and Yen Joe Tan²

¹The Chinese University of Hong Kong

²Earth and Environmental Sciences Programme, Faculty of Science, The Chinese University of Hong Kong

April 16, 2024

Abstract

The application of deep-learning-based seismic phase pickers for earthquake monitoring has surged in recent years. However, the efficacy of these models when applied to monitoring volcano seismicity has yet to be evaluated. Here, we first compile a dataset of seismic waveforms from various volcanoes globally. We then show that the performances of two widely used deep-learning pickers deteriorate systematically as the earthquakes' frequency content decreases. Therefore, the performances are especially poor for long-period earthquakes often associated with fluid/magma movement. Subsequently, we train new models which perform significantly better, including when tested on volcanic earthquake waveforms from northern California where no training data are used and tectonic low-frequency earthquakes along the Nankai Trough. Our model/workflow can be applied to improve monitoring of volcano seismicity globally while our compiled dataset can be used to benchmark future methods for characterizing volcano seismicity, especially long-period earthquakes which are difficult to monitor.

Deep-learning-based phase picking for volcano seismicity

Yiyuan Zhong¹, Yen Joe Tan¹

¹Earth and Environmental Sciences Programme, Faculty of Science, The Chinese University of Hong

Kong, Hong Kong S.A.R., China

Key Points:

- We compile the first data set of seismic waveforms from various volcanic regions globally.
- We show that existing deep-learning phase pickers' performances deteriorate with decreasing volcanic earthquake frequency content.
- Our retrained models perform better and are more generalizable for monitoring volcano seismicity, especially long-period earthquakes.

Corresponding author: Yen Joe Tan, yjtan@cuhk.edu.hk

Abstract

The application of deep-learning-based seismic phase pickers for earthquake monitoring has surged in recent years. However, the efficacy of these models when applied to monitoring volcano seismicity has yet to be evaluated. Here, we first compile a dataset of seismic waveforms from various volcanoes globally. We then show that the performances of two widely used deep-learning pickers deteriorate systematically as the earthquakes' frequency content decreases. Therefore, the performances are especially poor for long-period earthquakes often associated with fluid/magma movement. Subsequently, we train new models which perform significantly better, including when tested on volcanic earthquake waveforms from northern California where no training data are used and tectonic low-frequency earthquakes along the Nankai Trough. Our model/workflow can be applied to improve monitoring of volcano seismicity globally while our compiled dataset can be used to benchmark future methods for characterizing volcano seismicity, especially long-period earthquakes which are difficult to monitor.

Plain Language Summary

Earthquake activity at volcanic regions is often monitored to indicate volcanic activity. Identifying the time when the energy radiated from an earthquake source arrives at a seismometer is essential for locating the earthquake, which can be difficult for volcanic earthquakes because of high noise levels, high event rates, and obscured onsets. Previous studies have demonstrated that deep learning, a type of artificial intelligence, can excel in picking the arrival times of regular earthquakes. However, the efficacy of these models when applied to monitoring volcanic earthquakes has yet to be evaluated. Here, we first compile a dataset of earthquakes from various volcanoes globally. We then show that existing deep-learning-based models do not work well for these events, especially those with predominantly low-frequency energy. We then train two new models which

perform better than existing models for volcanic earthquakes. Our model/workflow can be applied to improve monitoring of volcanic earthquakes globally.

1 Introduction

Detecting and identifying onsets of seismic phases is fundamental to locating seismicity. Manual inspection by experienced analysts is viewed as the gold standard but is extremely laborious and time-consuming. This makes it difficult to handle the ever-increasing volumes of seismic data and periods with extremely high seismicity rate such as during volcanic unrests. On the other hand, early automatic methods, such as the short-term average over long-term average method (STA/LTA) (Allen, 1978), suffer from low accuracy and require a number of parameters to be tuned carefully. Over the past two decades, the matched-filter technique has been shown to be an effective method (Gibbons & Ringdal, 2006; Chamberlain et al., 2017) to search for repeating or near-repeating earthquakes based on waveform similarity. However, this method is only capable of detecting earthquakes in the vicinity of known template events. In recent years, deep-learning-based pickers (e.g. Ross et al., 2018; Zhu & Beroza, 2019; Mousavi et al., 2020; Soto & Schurr, 2021) have been gaining increasing attention due to their picking accuracy being comparable to human analysts (Chai et al., 2020) and high efficiency. Their application has surged in recent years, including for delineating seismicity in fault zones, subduction zones, oceanic transform faults, and volcanoes (e.g. Tan et al., 2021; Jiang et al., 2022; Chen et al., 2022; Gong et al., 2023; Liu et al., 2023; Wilding et al., 2023; Garza-Girón et al., 2023). However, it can be difficult to predict deep-learning models' performance for out-of-distribution data that are not well represented by training data (Wenzel et al., 2022; Teney et al., 2022).

Seismicity which often correlate with magmatic/volcanic processes and sometimes represent eruption precursors (White & McCausland, 2019; Acocella et al., 2023) is an important monitoring observable at volcanoes. Two types of earthquakes are commonly

64 observed in volcanic regions: volcano-tectonic earthquakes (VTs) and long-period earth-
65 quakes (LPs), which are classified mainly based on their waveform frequency content but
66 may imply different source processes (e.g. Chouet & Matoza, 2013; Saccorotti & Lok-
67 mer, 2021; Matoza & Roman, 2022, and references therein). VTs share common spec-
68 tral characteristics with regular tectonic earthquakes and have impulsive onsets. They
69 mostly originate from shear fractures in the solid part of an edifice or the underlying crust,
70 hence only indirectly indicate magmatic activity. In comparison, most conceptual source
71 models of LPs involve fluids, e.g. resonating fluid-filled cracks (Chouet & Matoza, 2013),
72 thermal stresses in cooling magmas (Aso & Tsai, 2014), pressurization of exsolved volatiles
73 from stalled magmas (Wech et al., 2020), and rapidly growing bubble in ascending mag-
74 mas (Melnik et al., 2020). Therefore, LPs are often interpreted as a more direct evidence
75 of fluid movement (e.g. Song et al., 2023). However, compared to VTs, LPs are more
76 difficult to detect because they are depleted of high frequency content and have emer-
77 gent phase onsets (Pitt et al., 2002; Shapiro et al., 2017).

78 Some recent studies have applied existing deep-learning phase pickers, which were
79 trained using regular tectonic earthquake waveforms, to monitor volcano seismicity (Mittal
80 et al., 2022; Bannister et al., 2022; Suarez et al., 2023; Li et al., 2023; Garza-Girón et
81 al., 2023; Wilding et al., 2023). However, there is currently no large-scale, systematic eval-
82 uation of the efficacy of these existing models for volcano monitoring. For instance, their
83 performances for volcanic earthquakes may be impaired by different waveform charac-
84 teristics, emergent onsets of long-period events, and high/different background noise in
85 volcanic regions (Lapins et al., 2021). While there have been a few models trained with
86 seismic data near volcanoes (Lapins et al., 2021; Kim et al., 2023; Armstrong et al., 2023),
87 limited data distribution (individual volcano) make these models less generalizable to
88 other volcanic regions. In addition, none of these studies explicitly included long-period
89 earthquakes in their analyses (Lapins et al., 2021; Kim et al., 2023; Armstrong et al., 2023).

90 In this study, we first compile a data set of seismic waveforms from various volcanic
91 regions. We then show that the performances of two widely used deep-learning pickers,
92 PhaseNet (Zhu & Beroza, 2019) and EQTransformer (Mousavi et al., 2020), deteriorate
93 when applied off-the-shelf to volcanic seismic data, especially for long-period earthquakes.
94 We then train new models that achieve significantly better performances for monitor-
95 ing volcano seismicity.

96 **2 Dataset of seismic waveforms from volcanic regions**

97 We assemble a data set of 156,272 LP waveforms (34,980 events), 156,498 VT wave-
98 forms (38,115 events), and 20,000 noise waveforms recorded by seismic stations deployed
99 around 34 volcanoes in Alaska (Power et al., 2019), 6 volcanoes in Hawaii (Hawaiian Vol-
100 cano Observatory/USGS, 1956), 8 volcanoes in northern California (NCEDC, 2014) and
101 88 volcanoes in Japan (National Research Institute for Earth Science and Disaster Re-
102 siliance, 2019). The geographical distribution of the events is shown in Figure 1. See Ta-
103 ble S1 in the supporting information for more details about data set splitting, Figure S1
104 for the distribution of recording stations, Figure S2 for the distribution of volcanoes and
105 Figures S3-S14 for other properties of the data. All the event waveforms have both man-
106 ually picked P and S phase arrivals. Most waveforms contain 3 components (77%) (Fig-
107 ure S3) and are from earthquakes located within 50 km of an active volcano (95%) (Fig-
108 ure S4). Since there are far more available VTs than LPs, we only include a similar num-
109 ber of VT waveforms as the number of available LP waveforms. We remove data with
110 large spikes and errors (e.g. events with S pick prior to P pick). For waveforms from Japan,
111 we download event waveforms whose length may vary for different events and different
112 stations. For waveforms from the US, we download event waveforms starting from 60s
113 before the P pick and ending 60s after the S pick. Hence waveforms in our data set have
114 different lengths, which will be trimmed in the subsequent processing stages. Compared
115 with previous datasets, e.g. STEAD (Mousavi et al., 2019) and INSTANCE (Michelin
116 et al., 2021), our data set has a wider distribution of frequency index (Figures S7-S10)

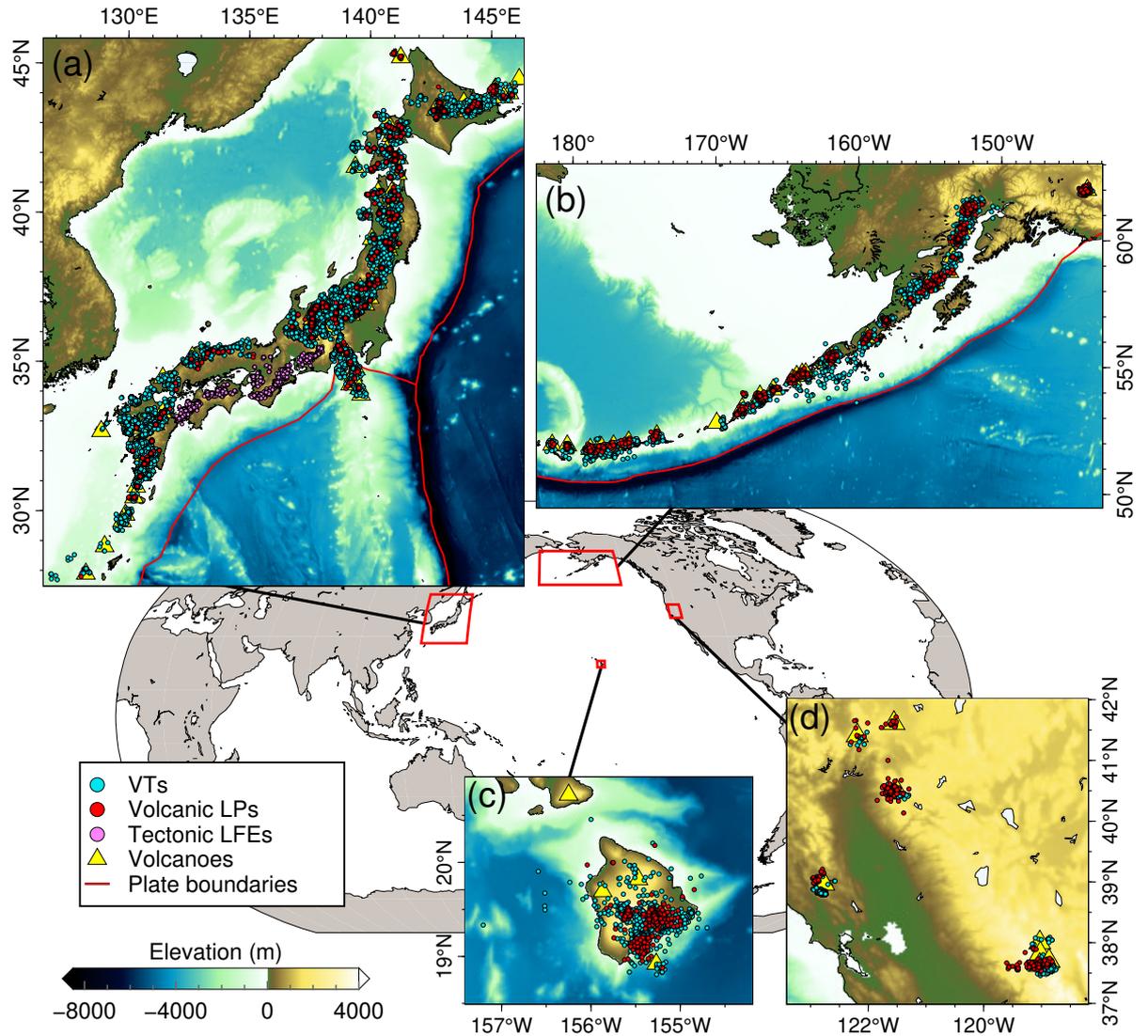


Figure 1. Geographical distribution of the earthquakes used in this study. The seismic data of volcano-tectonic earthquakes (cyan circles) and volcanic long-period earthquakes (red circles) from Japan (a), Alaska (b) and Hawaii (c) are split into a training set, a validation set and a test set, while the data from northern California (d) and the tectonic low-frequency earthquakes (LFEs) (purple circles) from Japan are only used for testing. Yellow triangles mark active volcanoes with seismic events used in this study.

117 which is a measure of the dominant frequency content of an earthquake (Buurman & West,
118 2010) (Text S1), suggesting it includes a greater variety of seismic events. To the best
119 of our knowledge, this is the first data set of seismic waveforms compiled from various
120 volcanic regions globally for machine learning.

121 **3 Evaluation of existing deep-learning phase pickers**

122 We use 15,078 LP waveforms and 15,057 VT waveforms from Alaska, Hawaii and
123 Japan to evaluate two most widely used models: PhaseNet (Zhu & Beroza, 2019) and
124 EQTransformer (Mousavi et al., 2020), which are the best performing architectures in
125 a recent benchmark study (Münchmeyer et al., 2022). PhaseNet is a U-net with 1D con-
126 volutional layers originally trained on earthquakes from northern California. EQTrans-
127 former is a stack of convolutional layers, long short-term memory (LSTM) units, and self-
128 attentive layers originally trained on the global data set STEAD (Mousavi et al., 2019).
129 We divide the testing waveforms into subsets according to frequency index values to eval-
130 uate how the model performance varies with the dominant frequency content. We ran-
131 domly extract 30s windows around the manual picks of the testing waveforms. For each
132 waveform, the same window is used to test different models. Since EQTransformer op-
133 erates on a 60s window, we will only focus on the 30s target window of the output (Münchmeyer
134 et al., 2022). We use precision, recall and F1-score to evaluate the results. Precision is
135 the fraction of output picks that are actually correct. Recall is the fraction of manual
136 picks that are correctly identified by the model. F1 score is the harmonic mean of pre-
137 cision and recall (Text S2). Considering that the original EQTransformer and PhaseNet
138 were trained under the TensorFlow framework (Abadi et al., 2015) that is different from
139 the platform we use (pyTorch) and that they were not trained on the same data set, we
140 also include the variants of EQTransformer and PhaseNet trained on the INSTANCE
141 data set (Michellini et al., 2021) for comparison, which were trained by Münchmeyer et
142 al. (2022) and available in the SeisBench package (Woollam et al., 2022). The model out-
143 put is time series of “probability” of P and S. To get predicted picks from the probabil-

144 ity time series output by the models, we first extract segments of probability curves above
145 a given threshold and the peak positions of these extracted segments are considered as
146 pick times. The model-specific threshold is tuned (Figure S15) on the validation set (Ta-
147 ble S1).

148 The recalls, precisions and F1 scores of the original models decrease systematically
149 with decreasing frequency index (Figure 2). For example, the F1 score of PhaseNet de-
150 creases from ~ 0.9 to ~ 0.5 for P picking and from ~ 0.85 to ~ 0.25 for S picking as the
151 frequency index decreases from ~ 0.5 to ~ 1.7 . Compared with precision, the recall ex-
152 hibits a greater deterioration, which can be as low as 0.4 for P picking and 0.2 for S pick-
153 ing, indicating that most LPs in the test set have been overlooked. We observe a sim-
154 ilar trend for the models trained on INSTANCE (Münchmeyer et al., 2022). This is un-
155 likely to be related to changes in signal-to-noise ratio since we do not observe significant
156 systematic changes in signal-to-noise ratio with frequency index (Figure S17). Our re-
157 sults suggest that these existing models will likely underreport LPs compared to VTs
158 when directly applied to monitoring volcano seismicity (Bannister et al., 2022; Mittal
159 et al., 2022; Wilding et al., 2023; Garza-Girón et al., 2023; Suarez et al., 2023; Li et al.,
160 2023), which is not ideal since LPs often indicate fluid/magma movements (Chouet &
161 Matoza, 2013; Matoza & Roman, 2022). Therefore, we decided it would be valuable to
162 train a new phase picker specifically for volcano seismicity.

163 **4 Training deep-learning phase pickers for volcano seismicity**

164 Among our data set, 151,431 LP waveforms, 151,657 VT waveforms and 20,000 noise
165 waveforms from Alaska, Hawaii and Japan corresponding to 70,352 events are grouped
166 into a training set (83.64%), a validation set (5.49%) and a test set (10.87%) (Table S1).
167 Here, the earthquake waveforms in the test set are the same as those presented in the
168 previous section. An extra test set comprising 4,841 waveforms from 1,094 LP events and
169 4,841 waveforms from 1,649 VT events near 8 volcanoes in northern California is used

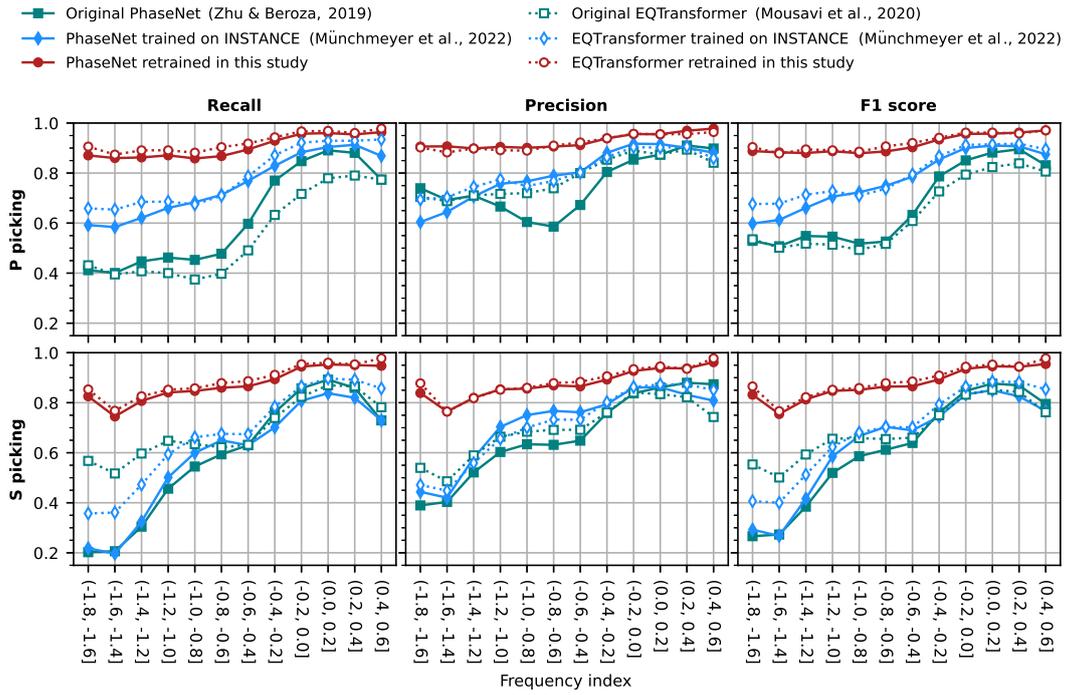


Figure 2. Performances of various models on subsets of testing waveforms with different frequency index values. The F1 scores here is slightly higher than those in Figure 3a because noise waveforms, to which frequency index is not applicable, are not included in this test.

170 to test how our model generalizes to a region where no training data have been used. In
171 addition, 6,224 waveforms of 2,356 tectonic low-frequency earthquakes (LFEs) along the
172 Nankai trough in Japan are used as another test set to investigate whether our model
173 works for tectonic LFEs associated with shear slip on the subduction zone plate inter-
174 face (Obara & Kato, 2016).

175 We use our data set to train two new models based on the PhaseNet and EQTrans-
176 former architectures implemented in the SeisBench package (Woollam et al., 2022). All
177 the waveforms are resampled to 100 Hz. We normalize each component of a waveform
178 by removing the mean and dividing it by the maximum value. We perform data augmen-
179 tation by randomly modifying the waveforms at each step of training. The modifications
180 include randomly shifting waveforms, adding gaps to waveforms, adding Gaussian noise
181 and superimposing a training example on the shifted and rescaled version of another train-
182 ing example. Each type of augmentation is performed with a given probability. Normal-
183 ization is performed before and after data augmentation. The labels for phase arrivals
184 are Gaussian functions with peaks aligning with manual picks. At each step of training,
185 a batch of waveform examples are randomly selected, normalized, randomly augmented,
186 labelled, and input into the Adam optimization algorithm (Kingma & Ba, 2015) to ad-
187 just the model weights.

188 The validation set is used to tune hyperparameters. We try various learning rates
189 0.0001/0.0005/0.001 and batch sizes 512/1024 to obtain a series of models. Each model
190 is trained for 400 epochs. Loss function on the validation set is monitored for each epoch
191 and the model snapshot at the epoch with the lowest validation loss is used as the final
192 model. For each model, we test different decision thresholds and choose the one with the
193 highest F1-score as the optimal threshold. Then we evaluate each model on the valida-
194 tion set and choose the one with the highest F1-score (Tables S2-3). The preferred learn-
195 ing rate and batch size for PhaseNet are 0.0005 and 512, respectively. They are 0.001
196 and 1024 for EQTransformer, respectively. We also compare random initialization and

197 initialization from the network weights pre-trained on the INSTANCE data set (Melnik
 198 et al., 2020; Münchmeyer et al., 2022), and we choose the one with the highest F1-score
 199 on the validation set (Table S4).

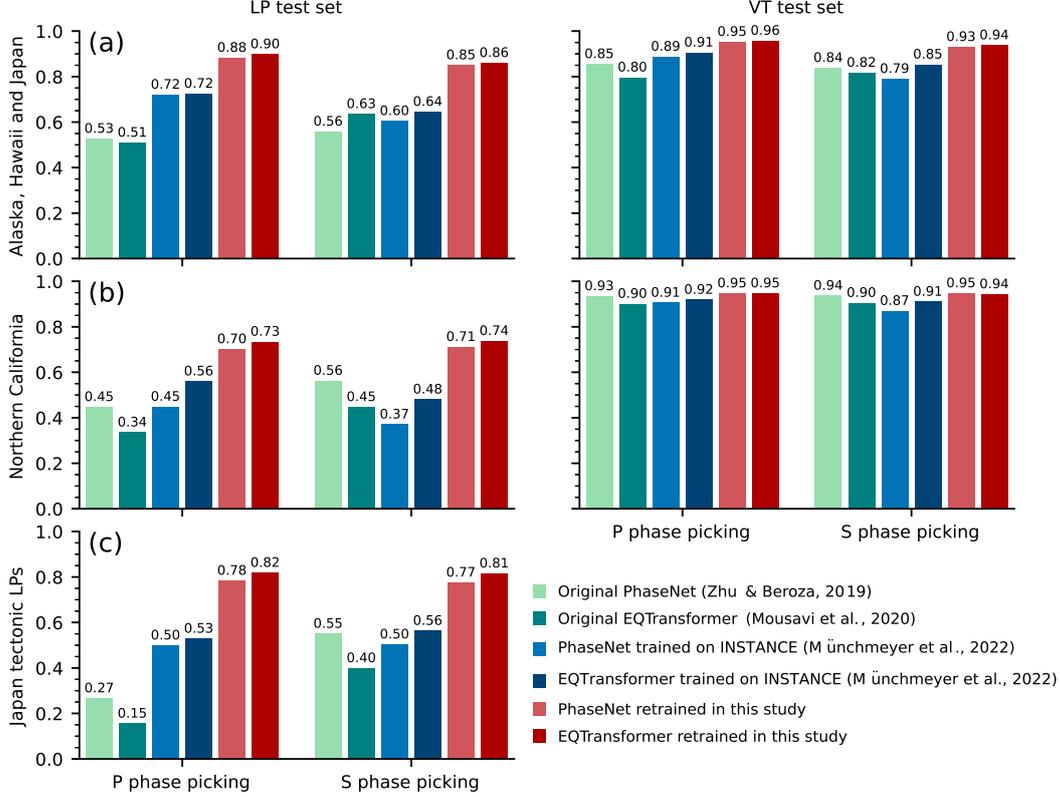


Figure 3. F1 scores of different models evaluated on the testing waveforms from (a) the same regions as the training data, (b) northern California from where no training data are used and (c) tectonic LFs in Japan. The precision and recall are given in Figures S24-S25 in the supplement.

200 We first test our models on subsets with different frequency index values as described
 201 in the previous section. Our models trained for volcano seismicity show significant per-
 202 formance improvement for waveforms with low frequency index values compared to ex-
 203 isting models, with F1 scores for P and S picking of ~ 0.9 and ~ 0.8 , respectively (Fig-
 204 ure 2). There is also a slight improvement for waveforms with high frequency index. The
 205 overall performances of various models on the whole test set are shown in Figure 3a, where
 206 our models show the best performances for both LPs and VTs for both P and S pick-

207 ing. For the LPs, the EQTransformer-based network trained in this study achieves an F1
208 score of 0.9 for P picking and 0.86 for S picking, which are 0.39 (P picking) and 0.23 (S
209 picking) higher than those of the original EQTransformer. The performance improve-
210 ment is smaller for the VTs: the retrained EQTransformer achieves F1 scores 0.16 and
211 0.12 higher than the original EQTransformer model for P and S picking respectively. The
212 EQTransformer trained on INSTANCE has similar performance to the original EQTrans-
213 former except for P picking on the LPs, for which the F1 score of the INSTANCE-based
214 EQTransformer is ~ 0.2 higher than that of the original EQTransformer but ~ 0.2 lower
215 than that of our retrained EQTransformer. A similar amount of improvement is obtained
216 by the PhaseNet-based network trained on our data set. Furthermore, our models give
217 lower picking residuals as indicated by the narrower histograms of residuals (Figure S19-
218 S20). The retrained EQTransformer shows only a marginally higher F1 score than the
219 retrained PhaseNet, suggesting that the data set plays a more important role than the
220 network architecture in differences in model performances.

221 Subsequently, we use the test set from northern California to investigate how our
222 models generalize to regions where no training data are used (Figure 3b). All the mod-
223 els show great performance for VTs, with F1 scores for P picking larger than 0.9 and F1
224 scores for S picking larger than 0.87, and our models achieve the highest F1 scores (0.95).
225 Notably, the existing pickers perform poorly for LPs, with F1 score ranging from 0.34
226 to 0.56. Although all the models experience some performance degradation for LPs com-
227 pared with the previous test, our retrained models still perform significantly better than
228 the existing models, with F1 scores ranging from 0.70 to 0.74. The performance varia-
229 tion with frequency index for this test set (Figure S18) also suggests that our models have
230 better generalization abilities when applied to a new region. The poorer performances
231 for LPs could be partly explained by the LP waveforms in this test set having lower signal-
232 to-noise ratios than VT waveforms (Figures S6 and S18).

233 Finally, we investigate whether our models also work for tectonic LFEs since both
234 tectonic LFEs and volcanic LPs appear to have similar frequency content, though they
235 are often inferred to reflect different source processes (Aso et al., 2013). Our training set
236 does not explicitly include any tectonic LFE. Here we test the models on LFEs along the
237 Nankai trough from Japan. The result is shown in Figure 3c. Our retrained models out-
238 perform the original models and the INSTANCE-based models by a large margin for both
239 P and S picking, with F1 scores of ~ 0.8 . We further confirmed that our models also work
240 for regular tectonic earthquakes, since they achieve F1 scores of 0.89 and 0.75 for P and
241 S picking respectively when tested on the INSTANCE data set (Michelini et al., 2021),
242 which is slightly better than the original EQTransformer and PhaseNet but unsurpris-
243 ingly inferior to the models trained on the INSTANCE data set (Figure S29).

244 5 Discussion

245 5.1 Comparison with existing methods

246 Deep-learning-based pickers have higher accuracy and require less parameters to
247 manually tune than traditional pickers, e.g. STA/LTA (Allen, 1978) and the Baer-Kradolfer
248 picker (Baer & Kradolfer, 1987), as demonstrated in previous studies (e.g. Zhu & Beroza,
249 2019; Mousavi et al., 2020; Münchmeyer et al., 2022). Also, deep-learning-based pick-
250 ers have greater flexibility than template matching as they are not limited by the avail-
251 ability of suitable template events. Compared with previous deep-learning models aimed
252 at tectonic earthquakes, our models can better pick volcano seismicity and thus can help
253 to improve volcano monitoring. Our compiled waveform dataset can also be used to bench-
254 mark future methods for monitoring volcanic earthquakes.

255 Our study is different from a few recent studies that have also trained models on
256 volcanic earthquakes (Lapins et al., 2021; Kim et al., 2023; Armstrong et al., 2023) in
257 two aspects. First, the previous studies focused exclusively on one volcano and thus it
258 is unclear how well these models can generalize to other volcanoes, while we use data around

259 136 active volcanoes from different regions. Second, LPs were not considered in the pre-
260 vious studies despite being an important form of volcano seismicity, while we included
261 LP earthquakes for training. We subsequently demonstrated that our models perform
262 well for both LPs and VTs, and can be generalized to other volcanoes. However, since
263 these studies adopted different data formats, input/output formats, machine-learning frame-
264 works and not all of these models are available, it would be hard to make direct com-
265 parisons.

266 Finally, our study is different from recent studies which focused on tectonic LFEs
267 (Thomas et al., 2021; Lin et al., 2023; Münchmeyer et al., 2023) in terms of training data
268 and targets. These studies focused on tectonic LFEs which are a manifestation of creep
269 or slow fault slips (Behr & Bürgmann, 2021), while our target is to pick volcano seismic-
270 ity including both VTs and LPs. The capability of our models to pick tectonic LFEs is
271 a side benefit and demonstrates that (1) our models are generalizable to other tectonic
272 environments and (2) tectonic LFEs and volcanic LPs have relatively similar waveform
273 characteristics.

274 **5.2 Different ways of performance evaluation**

275 The presented evaluation results for different models depend on the metrics used
276 and how they are calculated, which may vary in different studies. Therefore, it might
277 not be appropriate to directly compare the values reported in different papers. For in-
278 stance, some studies calculate true positive (TP), false positive (FP), true negative (TN)
279 and false negative (FN) based on waveform traces so that any of the four outcomes TP/FP/TN/FN
280 is assigned to each testing waveform (e.g. Zhu & Beroza, 2019; Mousavi et al., 2020).
281 In this case, a waveform is considered as a true positive as long as there is a predicted
282 pick sufficiently close to the manual pick even if there may also be some falsely predicted
283 picks for the same waveform. Hence, false predictions may be underreported. In contrast,
284 the definition of positive and negative in this paper is based on sampling points, where
285 any of TP/FP/TN/FN is assigned to each sampling point of a waveform rather than the

286 whole waveform (Text S2). The different definitions of FP and FN lead to different val-
287 ues of recall and precision. We have also calculated the model performances using the
288 definition of positive/negative based on waveform traces (Zhu & Beroza, 2019; Mousavi
289 et al., 2020), and the results (Figure S26-S27) show similar trends as those presented in
290 the previous section (Figure 2-3) except that the absolute values are slightly higher.

291 Alternatively, Münchmeyer et al. (2022) decomposed the evaluation into 3 tasks:
292 event detection, phase identification and onset time picking. This evaluation workflow
293 avoids the ambiguity in the definition of positive/negative for phase picking. However,
294 it uses the maximum probability value within the tested window as the prediction re-
295 sult, which may be inconsistent with the practical application of a deep-learning picker
296 where a trigger algorithm is used to retrieve picks from an output probability curve. Nev-
297 ertheless, our models also show better performances than existing models when evalu-
298 ated on the 3 tasks following Münchmeyer et al. (2022)’s workflow (Figure S21-S23 and
299 Table S5-S6), although existing models also perform well on the task of event detection
300 which is easier than phase picking. Therefore, our models show consistently better per-
301 formances than existing models regardless of the method of performance evaluation.

302 **6 Conclusion**

303 In this study, we first compile a dataset of seismic waveforms from various volcanic
304 regions globally, which has a wider distribution of frequency index than previous datasets
305 of tectonic earthquakes. We then show that existing deep-learning-based phase pickers
306 do not generalize well for volcanic earthquakes, with their performances deteriorating
307 as the earthquakes’ frequency content decreases, hence direct applications for monitor-
308 ing volcano seismicity is suboptimal with biases. Finally, we train and test new models
309 using our data set. The test results show that our models can better pick P and S phases
310 of VTs and LPs, and can be generalized to other regions not included in our training data

311 set, including for tectonic LFEs. Therefore, our results can benefit future efforts to im-
312 prove monitoring of volcano seismicity.

313 **Open Research Section**

314 Our models have been uploaded for peer review, with the archiving at Zenodo cur-
315 rently underway. All seismic data used in this study are publicly available. The seismic
316 waveforms and catalogs in Japan are from the Japan Meteorological Agency ([http://](http://www.jma.go.jp)
317 www.jma.go.jp) and the National Research Institute for Earth Science and Disaster Re-
318 siliance (<https://www.hinet.bosai.go.jp>) (National Research Institute for Earth Sci-
319 ence and Disaster Resilience, 2019). The seismic data and catalogs for Hawaii and Alaska
320 are from USGS (Hawaiian Volcano Observatory/USGS, 1956; Alaska Volcano Observa-
321 tory/USGS, 1988) and Incorporated Research Institutions for Seismology Data Manage-
322 ment center (IRIS-DMC, <https://ds.iris.edu/ds/nodes/dmc>). The seismic data and
323 catalogs for northern California are from the Northern California Earthquake Data Cen-
324 ter (NCEDC, 2014) (<https://ncedc.org>). We use the plate boundaries by Bird (2003)
325 in Figure 1. The volcano locations are from the Japan Meteorological Agency ([https://](https://www.data.jma.go.jp/vois/data/tokyo/STOCK/souran.eng/menu.htm)
326 www.data.jma.go.jp/vois/data/tokyo/STOCK/souran.eng/menu.htm), Geological Sur-
327 vey of Japan (https://gbank.gsj.jp/volcano/Quat.Vol/index_e.html), Alaska Vol-
328 cano Observatory (<https://www.avo.alaska.edu/volcano/>), Hawaiian Volcano Ob-
329 servatory (<https://www.usgs.gov/observatories/hvo>) and California Volcano Ob-
330 servatory (www.usgs.gov/observatories/calvo). We use ObsPy (Krischer et al., 2015)
331 and HinetPy (Tian et al., 2022) to facilitate waveform downloading. We use the network
332 architectures implemented in the SeisBench package (Woollam et al., 2022). We train
333 the networks under the PyTorch framework (Paszke et al., 2019) using the pytorch-lightning
334 package (Falcon & The PyTorch Lightning team, 2019).

335 **Acknowledgments**

336 This work is supported by the Direct Grant for Research (Grant 4053512) from the Chi-
337 nese University of Hong Kong, Hong Kong RGC General Research Fund (Grant 14300422),
338 and the Croucher Tak Wah Mak Innovation Award.

339 **References**

- 340 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X.
341 (2015, November). *TensorFlow, Large-scale machine learning on heterogeneous*
342 *systems*. doi: 10.5281/zenodo.4724125
- 343 Acocella, V., Ripepe, M., Rivalta, E., Peltier, A., Galetto, F., & Joseph, E. (2023).
344 Towards scientific forecasting of magmatic eruptions. *Nature Reviews Earth &*
345 *Environment*, 1–18.
- 346 Alaska Volcano Observatory/USGS. (1988). *Alaska volcano observatory*. Interna-
347 tional Federation of Digital Seismograph Networks. Retrieved from [https://](https://www.fdsn.org/networks/detail/AV/)
348 www.fdsn.org/networks/detail/AV/ doi: 10.7914/SN/AV
- 349 Allen, R. V. (1978). Automatic earthquake recognition and timing from single
350 traces. *Bulletin of the Seismological Society of America*, 68(5), 1521-1532.
- 351 Armstrong, A. D., Claerhout, Z., Baker, B., & Koper, K. D. (2023). A deep-learning
352 phase picker with calibrated bayesian-derived uncertainties for earthquakes
353 in the yellowstone volcanic region. *Bulletin of the Seismological Society of*
354 *America*, 113(6), 2323–2344.
- 355 Aso, N., Ohta, K., & Ide, S. (2013). Tectonic, volcanic, and semi-volcanic deep low-
356 frequency earthquakes in western japan. *Tectonophysics*, 600, 27–40.
- 357 Aso, N., & Tsai, V. C. (2014). Cooling magma model for deep volcanic long-period
358 earthquakes. *Journal of Geophysical Research: Solid Earth*, 119(11), 8442-
359 8456. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/abs/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2014JB011180)
360 [10.1002/2014JB011180](https://doi.org/10.1002/2014JB011180) doi: <https://doi.org/10.1002/2014JB011180>
- 361 Baer, M., & Kradolfer, U. (1987). An automatic phase picker for local and tele-

- 362 seismic events. *Bulletin of the Seismological Society of America*, 77(4), 1437-
363 1445.
- 364 Bannister, S., Bertrand, E. A., Heimann, S., Bourguignon, S., Asher, C., Shanks, J.,
365 & Harvison, A. (2022). Imaging sub-caldera structure with local seismicity,
366 okataina volcanic centre, taupo volcanic zone, using double-difference seismic
367 tomography. *Journal of Volcanology and Geothermal Research*, 431, 107653.
368 doi: <https://doi.org/10.1016/j.jvolgeores.2022.107653>
- 369 Behr, W. M., & Bürgmann, R. (2021). What's down there? The structures, ma-
370 terials and environment of deep-seated slow slip and tremor. *Philosophical*
371 *Transactions of the Royal Society A: Mathematical, Physical and Engineering*
372 *Sciences*, 379(2193), 20200218. doi: 10.1098/rsta.2020.0218
- 373 Bird, P. (2003). An updated digital model of plate boundaries. *Geochemistry, Geo-*
374 *physics, Geosystems*, 4(3). doi: <https://doi.org/10.1029/2001GC000252>
- 375 Buurman, H., & West, M. E. (2010). Seismic precursors to volcanic explosions dur-
376 ing the 2006 eruption of Augustine Volcano. In J. A. Power, M. L. Coombs, &
377 J. T. Freymueller (Eds.), *The 2006 eruption of Augustine Volcano, Alaska* (pp.
378 41–57). U.S. Geological Survey.
- 379 Chai, C., Maceira, M., Santos-Villalobos, H. J., Venkatakrishnan, S. V., Schoenball,
380 M., Zhu, W., ... Team, E. C. (2020). Using a deep neural network and trans-
381 fer learning to bridge scales for seismic phase picking. *Geophysical Research*
382 *Letters*, 47(16), e2020GL088651.
- 383 Chamberlain, C. J., Hopp, C. J., Boese, C. M., Warren-Smith, E., Chambers, D.,
384 Chu, S. X., ... Townend, J. (2017). EQcorrscan: Repeating and Near-
385 Repeating Earthquake Detection and Analysis in Python. *Seismological*
386 *Research Letters*, 89(1), 173-181.
- 387 Chen, H., Yang, H., Zhu, G., Xu, M., Lin, J., & You, Q. (2022). Deep outer-rise
388 faults in the southern mariana subduction zone indicated by a machine-
389 learning-based high-resolution earthquake catalog. *Geophysical Research*

390 *Letters*, 49(12), e2022GL097779.

391 Chouet, B. A., & Matoza, R. S. (2013). A multi-decadal view of seismic methods for
 392 detecting precursors of magma movement and eruption. *Journal of Volcanology*
 393 *and Geothermal Research*, 252, 108-175.

394 Falcon, W., & The PyTorch Lightning team. (2019, March). *PyTorch Lightning*. Re-
 395 trieved from <https://github.com/Lightning-AI/lightning> doi: 10.5281/
 396 zenodo.3828935

397 Garza-Girón, R., Brodsky, E. E., Spica, Z. J., Haney, M. M., & Webley, P. W.
 398 (2023). A specific earthquake processing workflow for studying long-lived,
 399 explosive volcanic eruptions with application to the 2008 okmok volcano,
 400 alaska, eruption. *Journal of Geophysical Research: Solid Earth*, 128(5),
 401 e2022JB025882.

402 Gibbons, S. J., & Ringdal, F. (2006). The detection of low magnitude seismic events
 403 using array-based waveform correlation. *Geophysical Journal International*,
 404 165(1), 149-166.

405 Gong, J., Fan, W., & Parnell-Turner, R. (2023). Machine learning-based new
 406 earthquake catalog illuminates on-fault and off-fault seismicity patterns at
 407 the discovery transform fault, east pacific rise. *Geochemistry, Geophysics,*
 408 *Geosystems*, 24(9), e2023GC011043.

409 Hawaiian Volcano Observatory/USGS. (1956). *Hawaiian volcano observatory net-*
 410 *work*. International Federation of Digital Seismograph Networks. Retrieved
 411 from <https://www.fdsn.org/networks/detail/HV/> doi: 10.7914/SN/HV

412 Jiang, C., Zhang, P., White, M. C. A., Pickle, R., & Miller, M. S. (2022). A Detailed
 413 Earthquake Catalog for Banda Arc–Australian Plate Collision Zone Using
 414 Machine-Learning Phase Picker and an Automated Workflow. *The Seismic*
 415 *Record*, 2(1), 1-10.

416 Kim, A., Nakamura, Y., Yukutake, Y., Uematsu, H., & Abe, Y. (2023). Develop-
 417 ment of a high-performance seismic phase picker using deep learning in the

- 418 hakone volcanic area. *Earth, Planets and Space*, 75(1), 1–15.
- 419 Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. In
420 Y. Bengio & Y. LeCun (Eds.), *3rd international conference on learning rep-*
421 *resentations, ICLR 2015, san diego, ca, usa, may 7-9, 2015, conference track*
422 *proceedings*.
- 423 Krischer, L., Megies, T., Barsch, R., Beyreuther, M., Lecocq, T., Caudron, C., &
424 Wassermann, J. (2015). ObsPy: a bridge for seismology into the scientific
425 Python ecosystem. *Computational Science & Discovery*, 8(1), 014003.
- 426 Lapins, S., Goitom, B., Kendall, J.-M., Werner, M. J., Cashman, K. V., & Ham-
427 mond, J. O. S. (2021). A little data goes a long way: Automating seismic
428 phase arrival picking at nabro volcano with transfer learning. *Journal of Geo-*
429 *physical Research: Solid Earth*, 126(7), e2021JB021910.
- 430 Li, J., Tian, Y., Zhao, D., Yan, D., Li, Z., & Li, H. (2023). Magmatic system and
431 seismicity of the Arxan volcanic group in Northeast China. *Geophysical Re-*
432 *search Letters*, 50(6), e2022GL101105.
- 433 Lin, J.-T., Thomas, A., Bachelot, L., Toomey, D., Searcy, J., & Melgar, D. (2023).
434 Detection of Hidden Low-Frequency Earthquakes in Southern Vancouver Is-
435 land with Deep Learning.
- 436 Liu, M., Li, L., Zhang, M., Lei, X., Nedimović, M. R., Plourde, A. P., ... Li, H.
437 (2023). Complexity of initiation and evolution of the 2013 yunlong earth-
438 quake swarm. *Earth and Planetary Science Letters*, 612, 118168. Re-
439 trieved from [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0012821X23001814)
440 S0012821X23001814 doi: <https://doi.org/10.1016/j.epsl.2023.118168>
- 441 Matoza, R. S., & Roman, D. C. (2022). One hundred years of advances in volcano
442 seismology and acoustics. *Bulletin of Volcanology*, 84(9), 86.
- 443 Melnik, O., Lyakhovsky, V., Shapiro, N. M., Galina, N., & Bergal-Kuvikas, O.
444 (2020). Deep long period volcanic earthquakes generated by degassing of
445 volatile-rich basaltic magmas. *Nature communications*, 11(1), 3918.

- 446 Michelini, A., Cianetti, S., Gaviano, S., Giunchi, C., Jozinović, D., & Lauciani, V.
447 (2021). Instance – the italian seismic dataset for machine learning. *Earth*
448 *System Science Data*, *13*(12), 5509–5544.
- 449 Mittal, T., Jordan, J. S., Retailleau, L., Beauducel, F., & Peltier, A. (2022). May-
450 otte 2018 eruption likely sourced from a magmatic mush. *Earth and Planetary*
451 *Science Letters*, *590*, 117566. doi: <https://doi.org/10.1016/j.epsl.2022.117566>
- 452 Mousavi, S. M., Ellsworth, W. L., Zhu, W., Chuang, L. Y., & Beroza, G. C. (2020).
453 Earthquake transformer—an attentive deep-learning model for simultaneous
454 earthquake detection and phase picking. *Nature communications*, *11*(1), 3952.
- 455 Mousavi, S. M., Sheng, Y., Zhu, W., & Beroza, G. C. (2019). Stanford earthquake
456 dataset (stead): A global data set of seismic signals for ai. *IEEE Access*, *7*,
457 179464–179476.
- 458 Münchmeyer, J., Giffard-Roisin, S., Malfante, M., Frank, W., Poli, P., Marsan, D.,
459 & Socquet, A. (2023). *Deep learning detects uncataloged low-frequency earth-*
460 *quakes across regions.*
- 461 Münchmeyer, J., Woollam, J., Rietbrock, A., Tilmann, F., Lange, D., Bornstein, T.,
462 ... others (2022). Which picker fits my data? a quantitative evaluation of
463 deep learning based seismic pickers. *Journal of Geophysical Research: Solid*
464 *Earth*, *127*(1), e2021JB023499.
- 465 National Research Institute for Earth Science and Disaster Resilience. (2019). *NIED*
466 *Hi-net, National Research Institute for Earth Science and Disaster Resilience.*
467 <https://www.hinet.bosai.go.jp>. doi: 10.17598/NIED.0003
- 468 NCEDC. (2014). *Northern california earthquake data center.* <https://ncedc.org/>.
469 doi: 10.7932/NCEDC
- 470 Obara, K., & Kato, A. (2016). Connecting slow earthquakes to huge earthquakes.
471 *Science*, *353*(6296), 253-257. doi: 10.1126/science.aaf1512
- 472 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... Chin-
473 tala, S. (2019). PyTorch: An Imperative Style, High-Performance Deep

- 474 Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché
475 Buc, E. Fox, & R. Garnett (Eds.), *Advances in Neural Information Process-*
476 *ing Systems 32* (pp. 8024–8035). Curran Associates, Inc. Retrieved from
477 [http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
478 [-high-performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
- 479 Pitt, A. M., Hill, D. P., Walter, S. W., & Johnson, M. J. S. (2002, 03). Midcrustal,
480 Long-period Earthquakes beneath Northern California Volcanic Areas. *Seismo-*
481 *logical Research Letters*, 73(2), 144-152. doi: 10.1785/gssrl.73.2.144
- 482 Power, J. A., Friberg, P. A., Haney, M. M., Parker, T., Stihler, S. D., & Dixon,
483 J. P. (2019). *A unified catalog of earthquake hypocenters and magnitudes at*
484 *volcanoes in alaska—1989 to 2018* (Tech. Rep.). US Geological Survey.
- 485 Ross, Z. E., Meier, M., Hauksson, E., & Heaton, T. H. (2018). Generalized Seis-
486 mic Phase Detection with Deep Learning. *Bulletin of the Seismological Society*
487 *of America*, 108(5A), 2894-2901.
- 488 Saccorotti, G., & Lokmer, I. (2021). Chapter 2 - a review of seismic methods for
489 monitoring and understanding active volcanoes. In P. Papale (Ed.), *Forecasting*
490 *and planning for volcanic hazards, risks, and disasters* (Vol. 2, p. 25-73). El-
491 sevier. Retrieved from [https://www.sciencedirect.com/science/article/](https://www.sciencedirect.com/science/article/pii/B9780128180822000020)
492 [pii/B9780128180822000020](https://www.sciencedirect.com/science/article/pii/B9780128180822000020) doi: <https://doi.org/10.1016/B978-0-12-818082-2>
493 [.00002-0](https://doi.org/10.1016/B978-0-12-818082-2)
- 494 Shapiro, N. M., Droznin, D., Droznina, S. Y., Senyukov, S., Gusev, A., & Gordeev,
495 E. (2017). Deep and shallow long-period volcanic seismicity linked by fluid-
496 pressure transfer. *Nature Geoscience*, 10(6), 442–445.
- 497 Song, Z., Tan, Y. J., & Roman, D. C. (2023). Deep long-period earthquakes at
498 akutan volcano from 2005 to 2017 better track magma influxes compared
499 to volcano-tectonic earthquakes. *Geophysical Research Letters*, 50(10),
500 e2022GL101987.
- 501 Soto, H., & Schurr, B. (2021). DeepPhasePick: a method for detecting and picking

- 502 seismic phases from local earthquakes based on highly optimized convolu-
503 tional and recurrent deep neural networks. *Geophysical Journal International*,
504 *227*(2), 1268-1294.
- 505 Suarez, E., Domínguez-Cerdeña, I., Villaseñor, A., Aparicio, S. S.-M., del Fresno,
506 C., & García-Cañada, L. (2023). Unveiling the pre-eruptive seismic series
507 of the la palma 2021 eruption: Insights through a fully automated analy-
508 sis. *Journal of Volcanology and Geothermal Research*, *444*, 107946. doi:
509 <https://doi.org/10.1016/j.jvolgeores.2023.107946>
- 510 Tan, Y. J., Waldhauser, F., Ellsworth, W. L., Zhang, M., Zhu, W., Michele, M.,
511 ... Segou, M. (2021). Machine-Learning-Based High-Resolution Earthquake
512 Catalog Reveals How Complex Fault Structures Were Activated during the
513 2016–2017 Central Italy Sequence. *The Seismic Record*, *1*(1), 11-19.
- 514 Teney, D., Lin, Y., Oh, S. J., & Abbasnejad, E. (2022). Id and ood performance
515 are sometimes inversely correlated on real-world datasets. In *Neural informa-*
516 *tion processing systems*.
- 517 Thomas, A. M., Inbal, A., Searcy, J., Shelly, D. R., & Bürgmann, R. (2021). Iden-
518 tification of low-frequency earthquakes on the san andreas fault with deep
519 learning. *Geophysical Research Letters*, *48*(13), e2021GL093157.
- 520 Tian, D., Kriegerowski, M., & Sawaki, Y. (2022). *seisman/hinetpy: 0.7.1*. Zen-
521 odo. Retrieved from <https://doi.org/10.5281/zenodo.6810553> doi: 10
522 .5281/zenodo.6810553
- 523 Wech, A. G., Thelen, W. A., & Thomas, A. M. (2020). Deep long-period earth-
524 quakes generated by second boiling beneath mauna kea volcano. *Science*,
525 *368*(6492), 775-779. Retrieved from [https://www.science.org/doi/abs/](https://www.science.org/doi/abs/10.1126/science.aba4798)
526 [10.1126/science.aba4798](https://www.science.org/doi/abs/10.1126/science.aba4798) doi: 10.1126/science.aba4798
- 527 Wenzel, F., Dittadi, A., Gehler, P. V., Simon-Gabriel, C.-J., Horn, M., Zietlow, D.,
528 ... Locatello, F. (2022). Assaying out-of-distribution generalization in transfer
529 learning. In *Neural information processing systems*.

- 530 White, R. A., & McCausland, W. A. (2019). A process-based model of pre-eruption
531 seismicity patterns and its use for eruption forecasting at dormant stratovolca-
532 noes. *Journal of Volcanology and Geothermal Research*, *382*, 267–297.
- 533 Wilding, J. D., Zhu, W., Ross, Z. E., & Jackson, J. M. (2023). The magmatic web
534 beneath hawai'i. *Science*, *379*(6631), 462-468.
- 535 Woollam, J., Münchmeyer, J., Tilmann, F., Rietbrock, A., Lange, D., Bornstein, T.,
536 ... Soto, H. (2022). SeisBench—A Toolbox for Machine Learning in Seismology.
537 *Seismological Research Letters*, *93*(3), 1695-1709.
- 538 Zhu, W., & Beroza, G. C. (2019). Phasenet: a deep-neural-network-based seismic
539 arrival-time picking method. *Geophysical Journal International*, *216*(1), 261–
540 273.

Deep-learning-based phase picking for volcano seismicity

Yiyuan Zhong¹, Yen Joe Tan¹

¹Earth and Environmental Sciences Programme, Faculty of Science, The Chinese University of Hong

Kong, Hong Kong S.A.R., China

Key Points:

- We compile the first data set of seismic waveforms from various volcanic regions globally.
- We show that existing deep-learning phase pickers' performances deteriorate with decreasing volcanic earthquake frequency content.
- Our retrained models perform better and are more generalizable for monitoring volcano seismicity, especially long-period earthquakes.

Corresponding author: Yen Joe Tan, yjtan@cuhk.edu.hk

Abstract

The application of deep-learning-based seismic phase pickers for earthquake monitoring has surged in recent years. However, the efficacy of these models when applied to monitoring volcano seismicity has yet to be evaluated. Here, we first compile a dataset of seismic waveforms from various volcanoes globally. We then show that the performances of two widely used deep-learning pickers deteriorate systematically as the earthquakes' frequency content decreases. Therefore, the performances are especially poor for long-period earthquakes often associated with fluid/magma movement. Subsequently, we train new models which perform significantly better, including when tested on volcanic earthquake waveforms from northern California where no training data are used and tectonic low-frequency earthquakes along the Nankai Trough. Our model/workflow can be applied to improve monitoring of volcano seismicity globally while our compiled dataset can be used to benchmark future methods for characterizing volcano seismicity, especially long-period earthquakes which are difficult to monitor.

Plain Language Summary

Earthquake activity at volcanic regions is often monitored to indicate volcanic activity. Identifying the time when the energy radiated from an earthquake source arrives at a seismometer is essential for locating the earthquake, which can be difficult for volcanic earthquakes because of high noise levels, high event rates, and obscured onsets. Previous studies have demonstrated that deep learning, a type of artificial intelligence, can excel in picking the arrival times of regular earthquakes. However, the efficacy of these models when applied to monitoring volcanic earthquakes has yet to be evaluated. Here, we first compile a dataset of earthquakes from various volcanoes globally. We then show that existing deep-learning-based models do not work well for these events, especially those with predominantly low-frequency energy. We then train two new models which

perform better than existing models for volcanic earthquakes. Our model/workflow can be applied to improve monitoring of volcanic earthquakes globally.

1 Introduction

Detecting and identifying onsets of seismic phases is fundamental to locating seismicity. Manual inspection by experienced analysts is viewed as the gold standard but is extremely laborious and time-consuming. This makes it difficult to handle the ever-increasing volumes of seismic data and periods with extremely high seismicity rate such as during volcanic unrests. On the other hand, early automatic methods, such as the short-term average over long-term average method (STA/LTA) (Allen, 1978), suffer from low accuracy and require a number of parameters to be tuned carefully. Over the past two decades, the matched-filter technique has been shown to be an effective method (Gibbons & Ringdal, 2006; Chamberlain et al., 2017) to search for repeating or near-repeating earthquakes based on waveform similarity. However, this method is only capable of detecting earthquakes in the vicinity of known template events. In recent years, deep-learning-based pickers (e.g. Ross et al., 2018; Zhu & Beroza, 2019; Mousavi et al., 2020; Soto & Schurr, 2021) have been gaining increasing attention due to their picking accuracy being comparable to human analysts (Chai et al., 2020) and high efficiency. Their application has surged in recent years, including for delineating seismicity in fault zones, subduction zones, oceanic transform faults, and volcanoes (e.g. Tan et al., 2021; Jiang et al., 2022; Chen et al., 2022; Gong et al., 2023; Liu et al., 2023; Wilding et al., 2023; Garza-Girón et al., 2023). However, it can be difficult to predict deep-learning models' performance for out-of-distribution data that are not well represented by training data (Wenzel et al., 2022; Teney et al., 2022).

Seismicity which often correlate with magmatic/volcanic processes and sometimes represent eruption precursors (White & McCausland, 2019; Acocella et al., 2023) is an important monitoring observable at volcanoes. Two types of earthquakes are commonly

64 observed in volcanic regions: volcano-tectonic earthquakes (VTs) and long-period earth-
65 quakes (LPs), which are classified mainly based on their waveform frequency content but
66 may imply different source processes (e.g. Chouet & Matoza, 2013; Saccorotti & Lok-
67 mer, 2021; Matoza & Roman, 2022, and references therein). VTs share common spec-
68 tral characteristics with regular tectonic earthquakes and have impulsive onsets. They
69 mostly originate from shear fractures in the solid part of an edifice or the underlying crust,
70 hence only indirectly indicate magmatic activity. In comparison, most conceptual source
71 models of LPs involve fluids, e.g. resonating fluid-filled cracks (Chouet & Matoza, 2013),
72 thermal stresses in cooling magmas (Aso & Tsai, 2014), pressurization of exsolved volatiles
73 from stalled magmas (Wech et al., 2020), and rapidly growing bubble in ascending mag-
74 mas (Melnik et al., 2020). Therefore, LPs are often interpreted as a more direct evidence
75 of fluid movement (e.g. Song et al., 2023). However, compared to VTs, LPs are more
76 difficult to detect because they are depleted of high frequency content and have emer-
77 gent phase onsets (Pitt et al., 2002; Shapiro et al., 2017).

78 Some recent studies have applied existing deep-learning phase pickers, which were
79 trained using regular tectonic earthquake waveforms, to monitor volcano seismicity (Mittal
80 et al., 2022; Bannister et al., 2022; Suarez et al., 2023; Li et al., 2023; Garza-Girón et
81 al., 2023; Wilding et al., 2023). However, there is currently no large-scale, systematic eval-
82 uation of the efficacy of these existing models for volcano monitoring. For instance, their
83 performances for volcanic earthquakes may be impaired by different waveform charac-
84 teristics, emergent onsets of long-period events, and high/different background noise in
85 volcanic regions (Lapins et al., 2021). While there have been a few models trained with
86 seismic data near volcanoes (Lapins et al., 2021; Kim et al., 2023; Armstrong et al., 2023),
87 limited data distribution (individual volcano) make these models less generalizable to
88 other volcanic regions. In addition, none of these studies explicitly included long-period
89 earthquakes in their analyses (Lapins et al., 2021; Kim et al., 2023; Armstrong et al., 2023).

90 In this study, we first compile a data set of seismic waveforms from various volcanic
91 regions. We then show that the performances of two widely used deep-learning pickers,
92 PhaseNet (Zhu & Beroza, 2019) and EQTransformer (Mousavi et al., 2020), deteriorate
93 when applied off-the-shelf to volcanic seismic data, especially for long-period earthquakes.
94 We then train new models that achieve significantly better performances for monitor-
95 ing volcano seismicity.

96 **2 Dataset of seismic waveforms from volcanic regions**

97 We assemble a data set of 156,272 LP waveforms (34,980 events), 156,498 VT wave-
98 forms (38,115 events), and 20,000 noise waveforms recorded by seismic stations deployed
99 around 34 volcanoes in Alaska (Power et al., 2019), 6 volcanoes in Hawaii (Hawaiian Vol-
100 cano Observatory/USGS, 1956), 8 volcanoes in northern California (NCEDC, 2014) and
101 88 volcanoes in Japan (National Research Institute for Earth Science and Disaster Re-
102 siliance, 2019). The geographical distribution of the events is shown in Figure 1. See Ta-
103 ble S1 in the supporting information for more details about data set splitting, Figure S1
104 for the distribution of recording stations, Figure S2 for the distribution of volcanoes and
105 Figures S3-S14 for other properties of the data. All the event waveforms have both man-
106 ually picked P and S phase arrivals. Most waveforms contain 3 components (77%) (Fig-
107 ure S3) and are from earthquakes located within 50 km of an active volcano (95%) (Fig-
108 ure S4). Since there are far more available VTs than LPs, we only include a similar num-
109 ber of VT waveforms as the number of available LP waveforms. We remove data with
110 large spikes and errors (e.g. events with S pick prior to P pick). For waveforms from Japan,
111 we download event waveforms whose length may vary for different events and different
112 stations. For waveforms from the US, we download event waveforms starting from 60s
113 before the P pick and ending 60s after the S pick. Hence waveforms in our data set have
114 different lengths, which will be trimmed in the subsequent processing stages. Compared
115 with previous datasets, e.g. STEAD (Mousavi et al., 2019) and INSTANCE (Michelin
116 et al., 2021), our data set has a wider distribution of frequency index (Figures S7-S10)

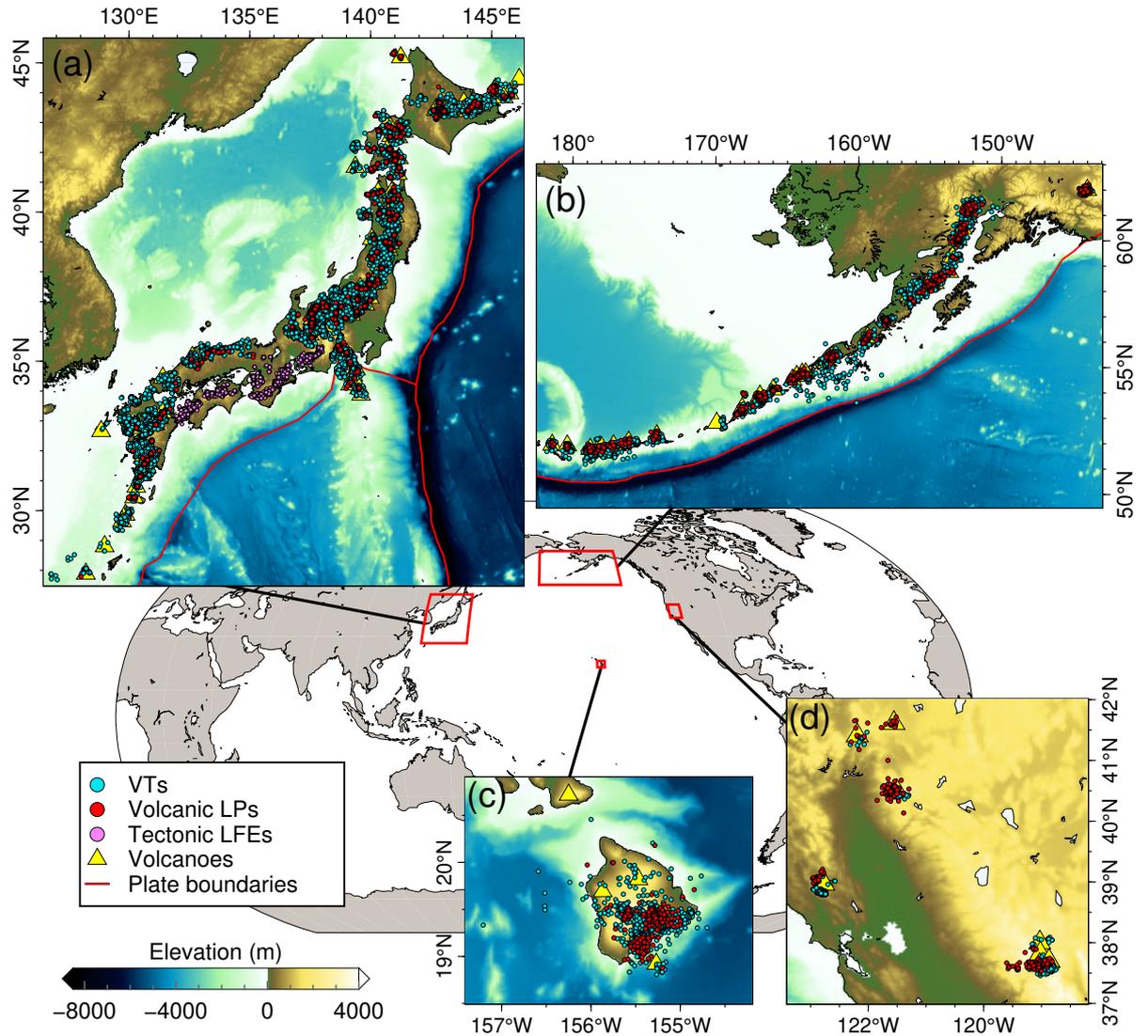


Figure 1. Geographical distribution of the earthquakes used in this study. The seismic data of volcano-tectonic earthquakes (cyan circles) and volcanic long-period earthquakes (red circles) from Japan (a), Alaska (b) and Hawaii (c) are split into a training set, a validation set and a test set, while the data from northern California (d) and the tectonic low-frequency earthquakes (LFEs) (purple circles) from Japan are only used for testing. Yellow triangles mark active volcanoes with seismic events used in this study.

117 which is a measure of the dominant frequency content of an earthquake (Buurman & West,
118 2010) (Text S1), suggesting it includes a greater variety of seismic events. To the best
119 of our knowledge, this is the first data set of seismic waveforms compiled from various
120 volcanic regions globally for machine learning.

121 **3 Evaluation of existing deep-learning phase pickers**

122 We use 15,078 LP waveforms and 15,057 VT waveforms from Alaska, Hawaii and
123 Japan to evaluate two most widely used models: PhaseNet (Zhu & Beroza, 2019) and
124 EQTransformer (Mousavi et al., 2020), which are the best performing architectures in
125 a recent benchmark study (Münchmeyer et al., 2022). PhaseNet is a U-net with 1D con-
126 volutional layers originally trained on earthquakes from northern California. EQTrans-
127 former is a stack of convolutional layers, long short-term memory (LSTM) units, and self-
128 attentive layers originally trained on the global data set STEAD (Mousavi et al., 2019).
129 We divide the testing waveforms into subsets according to frequency index values to eval-
130 uate how the model performance varies with the dominant frequency content. We ran-
131 domly extract 30s windows around the manual picks of the testing waveforms. For each
132 waveform, the same window is used to test different models. Since EQTransformer op-
133 erates on a 60s window, we will only focus on the 30s target window of the output (Münchmeyer
134 et al., 2022). We use precision, recall and F1-score to evaluate the results. Precision is
135 the fraction of output picks that are actually correct. Recall is the fraction of manual
136 picks that are correctly identified by the model. F1 score is the harmonic mean of pre-
137 cision and recall (Text S2). Considering that the original EQTransformer and PhaseNet
138 were trained under the TensorFlow framework (Abadi et al., 2015) that is different from
139 the platform we use (pyTorch) and that they were not trained on the same data set, we
140 also include the variants of EQTransformer and PhaseNet trained on the INSTANCE
141 data set (Michellini et al., 2021) for comparison, which were trained by Münchmeyer et
142 al. (2022) and available in the SeisBench package (Woollam et al., 2022). The model out-
143 put is time series of “probability” of P and S. To get predicted picks from the probabil-

144 ity time series output by the models, we first extract segments of probability curves above
145 a given threshold and the peak positions of these extracted segments are considered as
146 pick times. The model-specific threshold is tuned (Figure S15) on the validation set (Ta-
147 ble S1).

148 The recalls, precisions and F1 scores of the original models decrease systematically
149 with decreasing frequency index (Figure 2). For example, the F1 score of PhaseNet de-
150 creases from ~ 0.9 to ~ 0.5 for P picking and from ~ 0.85 to ~ 0.25 for S picking as the
151 frequency index decreases from ~ 0.5 to ~ 1.7 . Compared with precision, the recall ex-
152 hibits a greater deterioration, which can be as low as 0.4 for P picking and 0.2 for S pick-
153 ing, indicating that most LPs in the test set have been overlooked. We observe a sim-
154 ilar trend for the models trained on INSTANCE (Münchmeyer et al., 2022). This is un-
155 likely to be related to changes in signal-to-noise ratio since we do not observe significant
156 systematic changes in signal-to-noise ratio with frequency index (Figure S17). Our re-
157 sults suggest that these existing models will likely underreport LPs compared to VTs
158 when directly applied to monitoring volcano seismicity (Bannister et al., 2022; Mittal
159 et al., 2022; Wilding et al., 2023; Garza-Girón et al., 2023; Suarez et al., 2023; Li et al.,
160 2023), which is not ideal since LPs often indicate fluid/magma movements (Chouet &
161 Matoza, 2013; Matoza & Roman, 2022). Therefore, we decided it would be valuable to
162 train a new phase picker specifically for volcano seismicity.

163 **4 Training deep-learning phase pickers for volcano seismicity**

164 Among our data set, 151,431 LP waveforms, 151,657 VT waveforms and 20,000 noise
165 waveforms from Alaska, Hawaii and Japan corresponding to 70,352 events are grouped
166 into a training set (83.64%), a validation set (5.49%) and a test set (10.87%) (Table S1).
167 Here, the earthquake waveforms in the test set are the same as those presented in the
168 previous section. An extra test set comprising 4,841 waveforms from 1,094 LP events and
169 4,841 waveforms from 1,649 VT events near 8 volcanoes in northern California is used

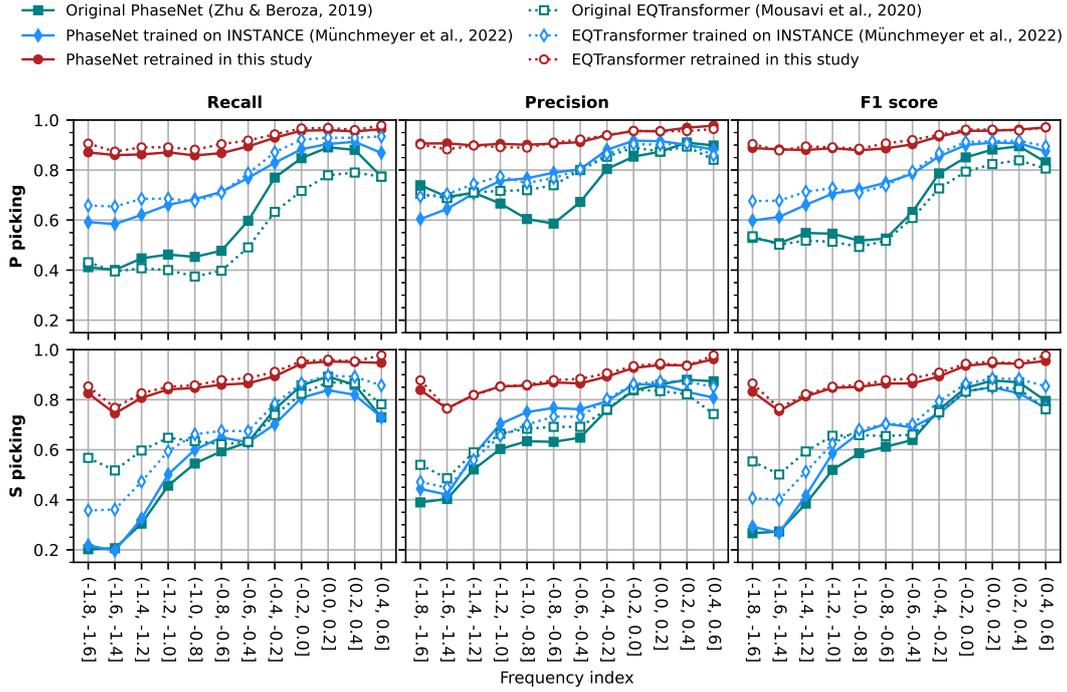


Figure 2. Performances of various models on subsets of testing waveforms with different frequency index values. The F1 scores here is slightly higher than those in Figure 3a because noise waveforms, to which frequency index is not applicable, are not included in this test.

170 to test how our model generalizes to a region where no training data have been used. In
171 addition, 6,224 waveforms of 2,356 tectonic low-frequency earthquakes (LFEs) along the
172 Nankai trough in Japan are used as another test set to investigate whether our model
173 works for tectonic LFEs associated with shear slip on the subduction zone plate inter-
174 face (Obara & Kato, 2016).

175 We use our data set to train two new models based on the PhaseNet and EQTrans-
176 former architectures implemented in the SeisBench package (Woollam et al., 2022). All
177 the waveforms are resampled to 100 Hz. We normalize each component of a waveform
178 by removing the mean and dividing it by the maximum value. We perform data augmen-
179 tation by randomly modifying the waveforms at each step of training. The modifications
180 include randomly shifting waveforms, adding gaps to waveforms, adding Gaussian noise
181 and superimposing a training example on the shifted and rescaled version of another train-
182 ing example. Each type of augmentation is performed with a given probability. Normal-
183 ization is performed before and after data augmentation. The labels for phase arrivals
184 are Gaussian functions with peaks aligning with manual picks. At each step of training,
185 a batch of waveform examples are randomly selected, normalized, randomly augmented,
186 labelled, and input into the Adam optimization algorithm (Kingma & Ba, 2015) to ad-
187 just the model weights.

188 The validation set is used to tune hyperparameters. We try various learning rates
189 0.0001/0.0005/0.001 and batch sizes 512/1024 to obtain a series of models. Each model
190 is trained for 400 epochs. Loss function on the validation set is monitored for each epoch
191 and the model snapshot at the epoch with the lowest validation loss is used as the final
192 model. For each model, we test different decision thresholds and choose the one with the
193 highest F1-score as the optimal threshold. Then we evaluate each model on the valida-
194 tion set and choose the one with the highest F1-score (Tables S2-3). The preferred learn-
195 ing rate and batch size for PhaseNet are 0.0005 and 512, respectively. They are 0.001
196 and 1024 for EQTransformer, respectively. We also compare random initialization and

197 initialization from the network weights pre-trained on the INSTANCE data set (Melnik
 198 et al., 2020; Münchmeyer et al., 2022), and we choose the one with the highest F1-score
 199 on the validation set (Table S4).

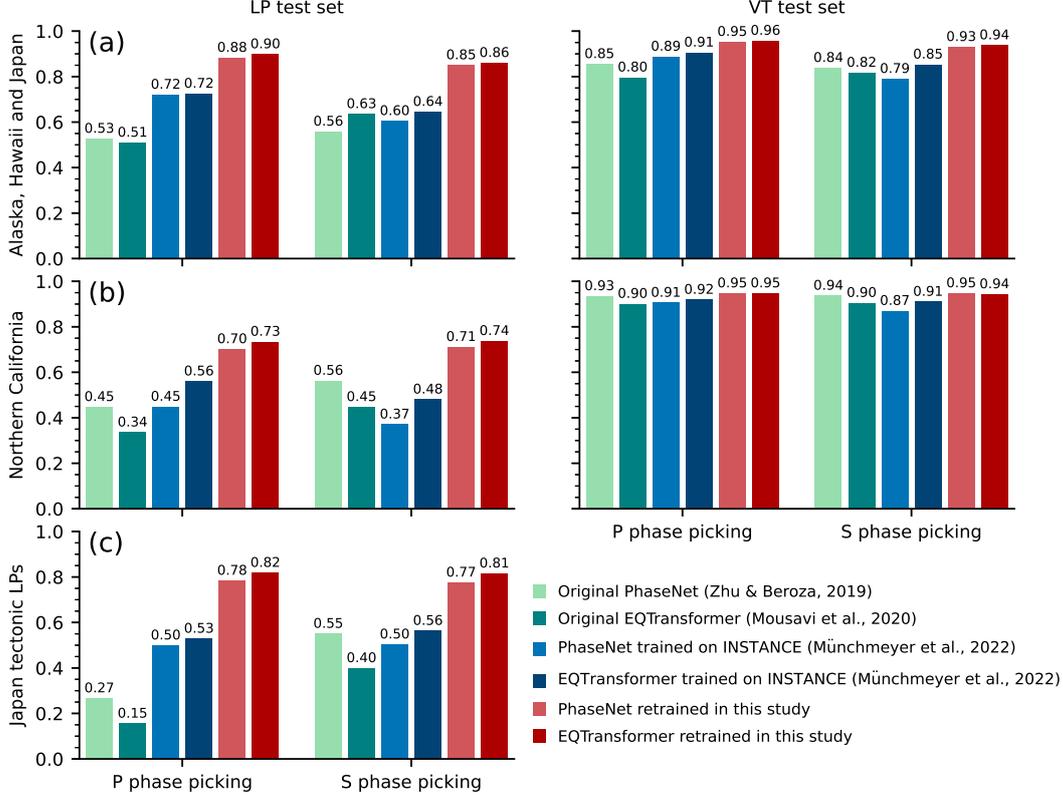


Figure 3. F1 scores of different models evaluated on the testing waveforms from (a) the same regions as the training data, (b) northern California from where no training data are used and (c) tectonic LFEs in Japan. The precision and recall are given in Figures S24-S25 in the supplement.

200 We first test our models on subsets with different frequency index values as described
 201 in the previous section. Our models trained for volcano seismicity show significant per-
 202 formance improvement for waveforms with low frequency index values compared to ex-
 203 isting models, with F1 scores for P and S picking of ~ 0.9 and ~ 0.8 , respectively (Fig-
 204 ure 2). There is also a slight improvement for waveforms with high frequency index. The
 205 overall performances of various models on the whole test set are shown in Figure 3a, where
 206 our models show the best performances for both LPs and VTs for both P and S pick-

207 ing. For the LPs, the EQTransformer-based network trained in this study achieves an F1
208 score of 0.9 for P picking and 0.86 for S picking, which are 0.39 (P picking) and 0.23 (S
209 picking) higher than those of the original EQTransformer. The performance improve-
210 ment is smaller for the VTs: the retrained EQTransformer achieves F1 scores 0.16 and
211 0.12 higher than the original EQTransformer model for P and S picking respectively. The
212 EQTransformer trained on INSTANCE has similar performance to the original EQTrans-
213 former except for P picking on the LPs, for which the F1 score of the INSTANCE-based
214 EQTransformer is ~ 0.2 higher than that of the original EQTransformer but ~ 0.2 lower
215 than that of our retrained EQTransformer. A similar amount of improvement is obtained
216 by the PhaseNet-based network trained on our data set. Furthermore, our models give
217 lower picking residuals as indicated by the narrower histograms of residuals (Figure S19-
218 S20). The retrained EQTransformer shows only a marginally higher F1 score than the
219 retrained PhaseNet, suggesting that the data set plays a more important role than the
220 network architecture in differences in model performances.

221 Subsequently, we use the test set from northern California to investigate how our
222 models generalize to regions where no training data are used (Figure 3b). All the mod-
223 els show great performance for VTs, with F1 scores for P picking larger than 0.9 and F1
224 scores for S picking larger than 0.87, and our models achieve the highest F1 scores (0.95).
225 Notably, the existing pickers perform poorly for LPs, with F1 score ranging from 0.34
226 to 0.56. Although all the models experience some performance degradation for LPs com-
227 pared with the previous test, our retrained models still perform significantly better than
228 the existing models, with F1 scores ranging from 0.70 to 0.74. The performance varia-
229 tion with frequency index for this test set (Figure S18) also suggests that our models have
230 better generalization abilities when applied to a new region. The poorer performances
231 for LPs could be partly explained by the LP waveforms in this test set having lower signal-
232 to-noise ratios than VT waveforms (Figures S6 and S18).

233 Finally, we investigate whether our models also work for tectonic LFEs since both
234 tectonic LFEs and volcanic LPs appear to have similar frequency content, though they
235 are often inferred to reflect different source processes (Aso et al., 2013). Our training set
236 does not explicitly include any tectonic LFE. Here we test the models on LFEs along the
237 Nankai trough from Japan. The result is shown in Figure 3c. Our retrained models out-
238 perform the original models and the INSTANCE-based models by a large margin for both
239 P and S picking, with F1 scores of ~ 0.8 . We further confirmed that our models also work
240 for regular tectonic earthquakes, since they achieve F1 scores of 0.89 and 0.75 for P and
241 S picking respectively when tested on the INSTANCE data set (Michelini et al., 2021),
242 which is slightly better than the original EQTransformer and PhaseNet but unsurpris-
243 ingly inferior to the models trained on the INSTANCE data set (Figure S29).

244 **5 Discussion**

245 **5.1 Comparison with existing methods**

246 Deep-learning-based pickers have higher accuracy and require less parameters to
247 manually tune than traditional pickers, e.g. STA/LTA (Allen, 1978) and the Baer-Kradolfer
248 picker (Baer & Kradolfer, 1987), as demonstrated in previous studies (e.g. Zhu & Beroza,
249 2019; Mousavi et al., 2020; Münchmeyer et al., 2022). Also, deep-learning-based pick-
250 ers have greater flexibility than template matching as they are not limited by the avail-
251 ability of suitable template events. Compared with previous deep-learning models aimed
252 at tectonic earthquakes, our models can better pick volcano seismicity and thus can help
253 to improve volcano monitoring. Our compiled waveform dataset can also be used to bench-
254 mark future methods for monitoring volcanic earthquakes.

255 Our study is different from a few recent studies that have also trained models on
256 volcanic earthquakes (Lapins et al., 2021; Kim et al., 2023; Armstrong et al., 2023) in
257 two aspects. First, the previous studies focused exclusively on one volcano and thus it
258 is unclear how well these models can generalize to other volcanoes, while we use data around

259 136 active volcanoes from different regions. Second, LPs were not considered in the pre-
260 vious studies despite being an important form of volcano seismicity, while we included
261 LP earthquakes for training. We subsequently demonstrated that our models perform
262 well for both LPs and VTs, and can be generalized to other volcanoes. However, since
263 these studies adopted different data formats, input/output formats, machine-learning frame-
264 works and not all of these models are available, it would be hard to make direct com-
265 parisons.

266 Finally, our study is different from recent studies which focused on tectonic LFEs
267 (Thomas et al., 2021; Lin et al., 2023; Münchmeyer et al., 2023) in terms of training data
268 and targets. These studies focused on tectonic LFEs which are a manifestation of creep
269 or slow fault slips (Behr & Bürgmann, 2021), while our target is to pick volcano seismic-
270 ity including both VTs and LPs. The capability of our models to pick tectonic LFEs is
271 a side benefit and demonstrates that (1) our models are generalizable to other tectonic
272 environments and (2) tectonic LFEs and volcanic LPs have relatively similar waveform
273 characteristics.

274 **5.2 Different ways of performance evaluation**

275 The presented evaluation results for different models depend on the metrics used
276 and how they are calculated, which may vary in different studies. Therefore, it might
277 not be appropriate to directly compare the values reported in different papers. For in-
278 stance, some studies calculate true positive (TP), false positive (FP), true negative (TN)
279 and false negative (FN) based on waveform traces so that any of the four outcomes TP/FP/TN/FN
280 is assigned to each testing waveform (e.g. Zhu & Beroza, 2019; Mousavi et al., 2020).
281 In this case, a waveform is considered as a true positive as long as there is a predicted
282 pick sufficiently close to the manual pick even if there may also be some falsely predicted
283 picks for the same waveform. Hence, false predictions may be underreported. In contrast,
284 the definition of positive and negative in this paper is based on sampling points, where
285 any of TP/FP/TN/FN is assigned to each sampling point of a waveform rather than the

286 whole waveform (Text S2). The different definitions of FP and FN lead to different val-
287 ues of recall and precision. We have also calculated the model performances using the
288 definition of positive/negative based on waveform traces (Zhu & Beroza, 2019; Mousavi
289 et al., 2020), and the results (Figure S26-S27) show similar trends as those presented in
290 the previous section (Figure 2-3) except that the absolute values are slightly higher.

291 Alternatively, Münchmeyer et al. (2022) decomposed the evaluation into 3 tasks:
292 event detection, phase identification and onset time picking. This evaluation workflow
293 avoids the ambiguity in the definition of positive/negative for phase picking. However,
294 it uses the maximum probability value within the tested window as the prediction re-
295 sult, which may be inconsistent with the practical application of a deep-learning picker
296 where a trigger algorithm is used to retrieve picks from an output probability curve. Nev-
297 ertheless, our models also show better performances than existing models when evalu-
298 ated on the 3 tasks following Münchmeyer et al. (2022)’s workflow (Figure S21-S23 and
299 Table S5-S6), although existing models also perform well on the task of event detection
300 which is easier than phase picking. Therefore, our models show consistently better per-
301 formances than existing models regardless of the method of performance evaluation.

302 **6 Conclusion**

303 In this study, we first compile a dataset of seismic waveforms from various volcanic
304 regions globally, which has a wider distribution of frequency index than previous datasets
305 of tectonic earthquakes. We then show that existing deep-learning-based phase pickers
306 do not generalize well for volcanic earthquakes, with their performances deteriorating
307 as the earthquakes’ frequency content decreases, hence direct applications for monitor-
308 ing volcano seismicity is suboptimal with biases. Finally, we train and test new models
309 using our data set. The test results show that our models can better pick P and S phases
310 of VTs and LPs, and can be generalized to other regions not included in our training data

311 set, including for tectonic LFEs. Therefore, our results can benefit future efforts to im-
312 prove monitoring of volcano seismicity.

313 **Open Research Section**

314 Our models have been uploaded for peer review, with the archiving at Zenodo cur-
315 rently underway. All seismic data used in this study are publicly available. The seismic
316 waveforms and catalogs in Japan are from the Japan Meteorological Agency ([http://](http://www.jma.go.jp)
317 www.jma.go.jp) and the National Research Institute for Earth Science and Disaster Re-
318 siliance (<https://www.hinet.bosai.go.jp>) (National Research Institute for Earth Sci-
319 ence and Disaster Resilience, 2019). The seismic data and catalogs for Hawaii and Alaska
320 are from USGS (Hawaiian Volcano Observatory/USGS, 1956; Alaska Volcano Observa-
321 tory/USGS, 1988) and Incorporated Research Institutions for Seismology Data Manage-
322 ment center (IRIS-DMC, <https://ds.iris.edu/ds/nodes/dmc>). The seismic data and
323 catalogs for northern California are from the Northern California Earthquake Data Cen-
324 ter (NCEDC, 2014) (<https://ncedc.org>). We use the plate boundaries by Bird (2003)
325 in Figure 1. The volcano locations are from the Japan Meteorological Agency ([https://](https://www.data.jma.go.jp/vois/data/tokyo/STOCK/souran.eng/menu.htm)
326 www.data.jma.go.jp/vois/data/tokyo/STOCK/souran.eng/menu.htm), Geological Sur-
327 vey of Japan (https://gbank.gsj.jp/volcano/Quat.Vol/index_e.html), Alaska Vol-
328 cano Observatory (<https://www.avo.alaska.edu/volcano/>), Hawaiian Volcano Ob-
329 servatory (<https://www.usgs.gov/observatories/hvo>) and California Volcano Ob-
330 servatory (www.usgs.gov/observatories/calvo). We use ObsPy (Krischer et al., 2015)
331 and HinetPy (Tian et al., 2022) to facilitate waveform downloading. We use the network
332 architectures implemented in the SeisBench package (Woollam et al., 2022). We train
333 the networks under the PyTorch framework (Paszke et al., 2019) using the pytorch-lightning
334 package (Falcon & The PyTorch Lightning team, 2019).

335 **Acknowledgments**

336 This work is supported by the Direct Grant for Research (Grant 4053512) from the Chi-
337 nese University of Hong Kong, Hong Kong RGC General Research Fund (Grant 14300422),
338 and the Croucher Tak Wah Mak Innovation Award.

339 **References**

- 340 Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X.
341 (2015, November). *TensorFlow, Large-scale machine learning on heterogeneous*
342 *systems*. doi: 10.5281/zenodo.4724125
- 343 Acocella, V., Ripepe, M., Rivalta, E., Peltier, A., Galetto, F., & Joseph, E. (2023).
344 Towards scientific forecasting of magmatic eruptions. *Nature Reviews Earth &*
345 *Environment*, 1–18.
- 346 Alaska Volcano Observatory/USGS. (1988). *Alaska volcano observatory*. Interna-
347 tional Federation of Digital Seismograph Networks. Retrieved from [https://](https://www.fdsn.org/networks/detail/AV/)
348 www.fdsn.org/networks/detail/AV/ doi: 10.7914/SN/AV
- 349 Allen, R. V. (1978). Automatic earthquake recognition and timing from single
350 traces. *Bulletin of the Seismological Society of America*, 68(5), 1521-1532.
- 351 Armstrong, A. D., Claerhout, Z., Baker, B., & Koper, K. D. (2023). A deep-learning
352 phase picker with calibrated bayesian-derived uncertainties for earthquakes
353 in the yellowstone volcanic region. *Bulletin of the Seismological Society of*
354 *America*, 113(6), 2323–2344.
- 355 Aso, N., Ohta, K., & Ide, S. (2013). Tectonic, volcanic, and semi-volcanic deep low-
356 frequency earthquakes in western japan. *Tectonophysics*, 600, 27–40.
- 357 Aso, N., & Tsai, V. C. (2014). Cooling magma model for deep volcanic long-period
358 earthquakes. *Journal of Geophysical Research: Solid Earth*, 119(11), 8442-
359 8456. Retrieved from [https://agupubs.onlinelibrary.wiley.com/doi/abs/](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1002/2014JB011180)
360 [10.1002/2014JB011180](https://doi.org/10.1002/2014JB011180) doi: <https://doi.org/10.1002/2014JB011180>
- 361 Baer, M., & Kradolfer, U. (1987). An automatic phase picker for local and tele-

- 362 seismic events. *Bulletin of the Seismological Society of America*, *77*(4), 1437-
363 1445.
- 364 Bannister, S., Bertrand, E. A., Heimann, S., Bourguignon, S., Asher, C., Shanks, J.,
365 & Harvison, A. (2022). Imaging sub-caldera structure with local seismicity,
366 okataina volcanic centre, taupo volcanic zone, using double-difference seismic
367 tomography. *Journal of Volcanology and Geothermal Research*, *431*, 107653.
368 doi: <https://doi.org/10.1016/j.jvolgeores.2022.107653>
- 369 Behr, W. M., & Bürgmann, R. (2021). What's down there? The structures, ma-
370 terials and environment of deep-seated slow slip and tremor. *Philosophical*
371 *Transactions of the Royal Society A: Mathematical, Physical and Engineering*
372 *Sciences*, *379*(2193), 20200218. doi: 10.1098/rsta.2020.0218
- 373 Bird, P. (2003). An updated digital model of plate boundaries. *Geochemistry, Geo-*
374 *physics, Geosystems*, *4*(3). doi: <https://doi.org/10.1029/2001GC000252>
- 375 Buurman, H., & West, M. E. (2010). Seismic precursors to volcanic explosions dur-
376 ing the 2006 eruption of Augustine Volcano. In J. A. Power, M. L. Coombs, &
377 J. T. Freymueller (Eds.), *The 2006 eruption of Augustine Volcano, Alaska* (pp.
378 41–57). U.S. Geological Survey.
- 379 Chai, C., Maceira, M., Santos-Villalobos, H. J., Venkatakrisnan, S. V., Schoenball,
380 M., Zhu, W., ... Team, E. C. (2020). Using a deep neural network and trans-
381 fer learning to bridge scales for seismic phase picking. *Geophysical Research*
382 *Letters*, *47*(16), e2020GL088651.
- 383 Chamberlain, C. J., Hopp, C. J., Boese, C. M., Warren-Smith, E., Chambers, D.,
384 Chu, S. X., ... Townend, J. (2017). EQcorrscan: Repeating and Near-
385 Repeating Earthquake Detection and Analysis in Python. *Seismological*
386 *Research Letters*, *89*(1), 173-181.
- 387 Chen, H., Yang, H., Zhu, G., Xu, M., Lin, J., & You, Q. (2022). Deep outer-rise
388 faults in the southern mariana subduction zone indicated by a machine-
389 learning-based high-resolution earthquake catalog. *Geophysical Research*

390 *Letters*, 49(12), e2022GL097779.

391 Chouet, B. A., & Matoza, R. S. (2013). A multi-decadal view of seismic methods for
392 detecting precursors of magma movement and eruption. *Journal of Volcanology*
393 *and Geothermal Research*, 252, 108-175.

394 Falcon, W., & The PyTorch Lightning team. (2019, March). *PyTorch Lightning*. Re-
395 trieved from <https://github.com/Lightning-AI/lightning> doi: 10.5281/
396 zenodo.3828935

397 Garza-Girón, R., Brodsky, E. E., Spica, Z. J., Haney, M. M., & Webley, P. W.
398 (2023). A specific earthquake processing workflow for studying long-lived,
399 explosive volcanic eruptions with application to the 2008 okmok volcano,
400 alaska, eruption. *Journal of Geophysical Research: Solid Earth*, 128(5),
401 e2022JB025882.

402 Gibbons, S. J., & Ringdal, F. (2006). The detection of low magnitude seismic events
403 using array-based waveform correlation. *Geophysical Journal International*,
404 165(1), 149-166.

405 Gong, J., Fan, W., & Parnell-Turner, R. (2023). Machine learning-based new
406 earthquake catalog illuminates on-fault and off-fault seismicity patterns at
407 the discovery transform fault, east pacific rise. *Geochemistry, Geophysics,*
408 *Geosystems*, 24(9), e2023GC011043.

409 Hawaiian Volcano Observatory/USGS. (1956). *Hawaiian volcano observatory net-*
410 *work*. International Federation of Digital Seismograph Networks. Retrieved
411 from <https://www.fdsn.org/networks/detail/HV/> doi: 10.7914/SN/HV

412 Jiang, C., Zhang, P., White, M. C. A., Pickle, R., & Miller, M. S. (2022). A Detailed
413 Earthquake Catalog for Banda Arc–Australian Plate Collision Zone Using
414 Machine-Learning Phase Picker and an Automated Workflow. *The Seismic*
415 *Record*, 2(1), 1-10.

416 Kim, A., Nakamura, Y., Yukutake, Y., Uematsu, H., & Abe, Y. (2023). Develop-
417 ment of a high-performance seismic phase picker using deep learning in the

- 418 hakone volcanic area. *Earth, Planets and Space*, *75*(1), 1–15.
- 419 Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. In
420 Y. Bengio & Y. LeCun (Eds.), *3rd international conference on learning rep-*
421 *resentations, ICLR 2015, san diego, ca, usa, may 7-9, 2015, conference track*
422 *proceedings*.
- 423 Krischer, L., Megies, T., Barsch, R., Beyreuther, M., Lecocq, T., Caudron, C., &
424 Wassermann, J. (2015). ObsPy: a bridge for seismology into the scientific
425 Python ecosystem. *Computational Science & Discovery*, *8*(1), 014003.
- 426 Lapins, S., Goitom, B., Kendall, J.-M., Werner, M. J., Cashman, K. V., & Ham-
427 mond, J. O. S. (2021). A little data goes a long way: Automating seismic
428 phase arrival picking at nabro volcano with transfer learning. *Journal of Geo-*
429 *physical Research: Solid Earth*, *126*(7), e2021JB021910.
- 430 Li, J., Tian, Y., Zhao, D., Yan, D., Li, Z., & Li, H. (2023). Magmatic system and
431 seismicity of the Arxan volcanic group in Northeast China. *Geophysical Re-*
432 *search Letters*, *50*(6), e2022GL101105.
- 433 Lin, J.-T., Thomas, A., Bachelot, L., Toomey, D., Searcy, J., & Melgar, D. (2023).
434 Detection of Hidden Low-Frequency Earthquakes in Southern Vancouver Is-
435 land with Deep Learning.
- 436 Liu, M., Li, L., Zhang, M., Lei, X., Nedimović, M. R., Plourde, A. P., ... Li, H.
437 (2023). Complexity of initiation and evolution of the 2013 yunlong earth-
438 quake swarm. *Earth and Planetary Science Letters*, *612*, 118168. Re-
439 trieved from [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0012821X23001814)
440 S0012821X23001814 doi: <https://doi.org/10.1016/j.epsl.2023.118168>
- 441 Matoza, R. S., & Roman, D. C. (2022). One hundred years of advances in volcano
442 seismology and acoustics. *Bulletin of Volcanology*, *84*(9), 86.
- 443 Melnik, O., Lyakhovsky, V., Shapiro, N. M., Galina, N., & Bergal-Kuvikas, O.
444 (2020). Deep long period volcanic earthquakes generated by degassing of
445 volatile-rich basaltic magmas. *Nature communications*, *11*(1), 3918.

- 446 Michelini, A., Cianetti, S., Gaviano, S., Giunchi, C., Jozinović, D., & Lauciani, V.
447 (2021). Instance – the italian seismic dataset for machine learning. *Earth*
448 *System Science Data*, *13*(12), 5509–5544.
- 449 Mittal, T., Jordan, J. S., Retailleau, L., Beauducel, F., & Peltier, A. (2022). May-
450 otte 2018 eruption likely sourced from a magmatic mush. *Earth and Planetary*
451 *Science Letters*, *590*, 117566. doi: <https://doi.org/10.1016/j.epsl.2022.117566>
- 452 Mousavi, S. M., Ellsworth, W. L., Zhu, W., Chuang, L. Y., & Beroza, G. C. (2020).
453 Earthquake transformer—an attentive deep-learning model for simultaneous
454 earthquake detection and phase picking. *Nature communications*, *11*(1), 3952.
- 455 Mousavi, S. M., Sheng, Y., Zhu, W., & Beroza, G. C. (2019). Stanford earthquake
456 dataset (stead): A global data set of seismic signals for ai. *IEEE Access*, *7*,
457 179464–179476.
- 458 Münchmeyer, J., Giffard-Roisin, S., Malfante, M., Frank, W., Poli, P., Marsan, D.,
459 & Socquet, A. (2023). *Deep learning detects uncataloged low-frequency earth-*
460 *quakes across regions.*
- 461 Münchmeyer, J., Woollam, J., Rietbrock, A., Tilmann, F., Lange, D., Bornstein, T.,
462 ... others (2022). Which picker fits my data? a quantitative evaluation of
463 deep learning based seismic pickers. *Journal of Geophysical Research: Solid*
464 *Earth*, *127*(1), e2021JB023499.
- 465 National Research Institute for Earth Science and Disaster Resilience. (2019). *NIED*
466 *Hi-net, National Research Institute for Earth Science and Disaster Resilience.*
467 <https://www.hinet.bosai.go.jp>. doi: 10.17598/NIED.0003
- 468 NCEDC. (2014). *Northern california earthquake data center.* <https://ncedc.org/>.
469 doi: 10.7932/NCEDC
- 470 Obara, K., & Kato, A. (2016). Connecting slow earthquakes to huge earthquakes.
471 *Science*, *353*(6296), 253-257. doi: 10.1126/science.aaf1512
- 472 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., ... Chin-
473 tala, S. (2019). PyTorch: An Imperative Style, High-Performance Deep

- 474 Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché
475 Buc, E. Fox, & R. Garnett (Eds.), *Advances in Neural Information Process-*
476 *ing Systems 32* (pp. 8024–8035). Curran Associates, Inc. Retrieved from
477 [http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
478 [-high-performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)
- 479 Pitt, A. M., Hill, D. P., Walter, S. W., & Johnson, M. J. S. (2002, 03). Midcrustal,
480 Long-period Earthquakes beneath Northern California Volcanic Areas. *Seismo-*
481 *logical Research Letters*, *73*(2), 144-152. doi: 10.1785/gssrl.73.2.144
- 482 Power, J. A., Friberg, P. A., Haney, M. M., Parker, T., Stihler, S. D., & Dixon,
483 J. P. (2019). *A unified catalog of earthquake hypocenters and magnitudes at*
484 *volcanoes in alaska—1989 to 2018* (Tech. Rep.). US Geological Survey.
- 485 Ross, Z. E., Meier, M., Hauksson, E., & Heaton, T. H. (2018). Generalized Seis-
486 mic Phase Detection with Deep Learning. *Bulletin of the Seismological Society*
487 *of America*, *108*(5A), 2894-2901.
- 488 Saccorotti, G., & Lokmer, I. (2021). Chapter 2 - a review of seismic methods for
489 monitoring and understanding active volcanoes. In P. Papale (Ed.), *Forecasting*
490 *and planning for volcanic hazards, risks, and disasters* (Vol. 2, p. 25-73). El-
491 sevier. Retrieved from [https://www.sciencedirect.com/science/article/](https://www.sciencedirect.com/science/article/pii/B9780128180822000020)
492 [pii/B9780128180822000020](https://www.sciencedirect.com/science/article/pii/B9780128180822000020) doi: <https://doi.org/10.1016/B978-0-12-818082-2>
493 [.00002-0](https://doi.org/10.1016/B978-0-12-818082-2)
- 494 Shapiro, N. M., Droznin, D., Droznina, S. Y., Senyukov, S., Gusev, A., & Gordeev,
495 E. (2017). Deep and shallow long-period volcanic seismicity linked by fluid-
496 pressure transfer. *Nature Geoscience*, *10*(6), 442–445.
- 497 Song, Z., Tan, Y. J., & Roman, D. C. (2023). Deep long-period earthquakes at
498 akutan volcano from 2005 to 2017 better track magma influxes compared
499 to volcano-tectonic earthquakes. *Geophysical Research Letters*, *50*(10),
500 e2022GL101987.
- 501 Soto, H., & Schurr, B. (2021). DeepPhasePick: a method for detecting and picking

- 502 seismic phases from local earthquakes based on highly optimized convolu-
503 tional and recurrent deep neural networks. *Geophysical Journal International*,
504 *227*(2), 1268-1294.
- 505 Suarez, E., Domínguez-Cerdeña, I., Villaseñor, A., Aparicio, S. S.-M., del Fresno,
506 C., & García-Cañada, L. (2023). Unveiling the pre-eruptive seismic series
507 of the la palma 2021 eruption: Insights through a fully automated analy-
508 sis. *Journal of Volcanology and Geothermal Research*, *444*, 107946. doi:
509 <https://doi.org/10.1016/j.jvolgeores.2023.107946>
- 510 Tan, Y. J., Waldhauser, F., Ellsworth, W. L., Zhang, M., Zhu, W., Michele, M.,
511 ... Segou, M. (2021). Machine-Learning-Based High-Resolution Earthquake
512 Catalog Reveals How Complex Fault Structures Were Activated during the
513 2016–2017 Central Italy Sequence. *The Seismic Record*, *1*(1), 11-19.
- 514 Teney, D., Lin, Y., Oh, S. J., & Abbasnejad, E. (2022). Id and ood performance
515 are sometimes inversely correlated on real-world datasets. In *Neural informa-*
516 *tion processing systems*.
- 517 Thomas, A. M., Inbal, A., Searcy, J., Shelly, D. R., & Bürgmann, R. (2021). Iden-
518 tification of low-frequency earthquakes on the san andreas fault with deep
519 learning. *Geophysical Research Letters*, *48*(13), e2021GL093157.
- 520 Tian, D., Kriegerowski, M., & Sawaki, Y. (2022). *seisman/hinetpy: 0.7.1*. Zen-
521 odo. Retrieved from <https://doi.org/10.5281/zenodo.6810553> doi: 10
522 .5281/zenodo.6810553
- 523 Wech, A. G., Thelen, W. A., & Thomas, A. M. (2020). Deep long-period earth-
524 quakes generated by second boiling beneath mauna kea volcano. *Science*,
525 *368*(6492), 775-779. Retrieved from [https://www.science.org/doi/abs/
526 10.1126/science.aba4798](https://www.science.org/doi/abs/10.1126/science.aba4798) doi: 10.1126/science.aba4798
- 527 Wenzel, F., Dittadi, A., Gehler, P. V., Simon-Gabriel, C.-J., Horn, M., Zietlow, D.,
528 ... Locatello, F. (2022). Assaying out-of-distribution generalization in transfer
529 learning. In *Neural information processing systems*.

- 530 White, R. A., & McCausland, W. A. (2019). A process-based model of pre-eruption
531 seismicity patterns and its use for eruption forecasting at dormant stratovolca-
532 noes. *Journal of Volcanology and Geothermal Research*, *382*, 267–297.
- 533 Wilding, J. D., Zhu, W., Ross, Z. E., & Jackson, J. M. (2023). The magmatic web
534 beneath hawai'i. *Science*, *379*(6631), 462-468.
- 535 Woollam, J., Münchmeyer, J., Tilmann, F., Rietbrock, A., Lange, D., Bornstein, T.,
536 ... Soto, H. (2022). SeisBench—A Toolbox for Machine Learning in Seismology.
537 *Seismological Research Letters*, *93*(3), 1695-1709.
- 538 Zhu, W., & Beroza, G. C. (2019). Phasenet: a deep-neural-network-based seismic
539 arrival-time picking method. *Geophysical Journal International*, *216*(1), 261–
540 273.

Supporting Information for “Deep-learning-based phase picking for volcano seismicity”

Yiyuan Zhong¹, Yen Joe Tan¹

¹Earth and Environmental Sciences Programme, Faculty of Science, The Chinese University of Hong Kong, Hong Kong SAR, China

Contents of this files

1. Text S1 to S2
2. Tables S1 to S6
3. Figures S1 to S29

Introduction

Figures S1-S14 and Table S1 show the properties of the data used in this study. Tables S2-S4 and Figures S15-S16 show the process of hyperparameter tuning. Figure S17 shows the distribution of signal to noise ratio for different subsets of the test data from Alaska, Japan and Hawaii. Figure S18 shows the model performance versus frequency index for the testing waveforms from northern California. Figure S19 and S20 show the histograms of picking residuals for VTs and LPs, respectively. Tables S6-S5 and Figures S21-S23 show the evaluation results for the 3 tasks defined in (Münchmeyer et al., 2022). Figures S24 and S25 show the recalls and precisions of different models, respectively. Figures S26

and S27 show the F1 scores calculated using the definition of positive and negative based on waveform traces (Zhu & Beroza, 2019; Mousavi et al., 2020). Figures S28 and S29 show the performances of different models on the INSTANCE data set (Michelini et al., 2021).

Text S1. Frequency index

Frequency index (FI) is a metric used to quantify the dominant frequency content of an earthquake from seismic waveforms (Buurman & West, 2010),

$$FI = \log_{10} \frac{\bar{A}_{upper}}{\bar{A}_{lower}}, \quad (1)$$

where \bar{A}_{lower} and \bar{A}_{upper} are the mean spectral amplitudes in a predefined high-frequency band and a low-frequency band, respectively. Following Song, Tan, and Roman (2023), we choose 1-5 Hz and 10-15 Hz as the low and high frequency bands, respectively. Time windows starting 1s prior to and ending 6s after P arrivals are extracted to calculate FIs. If there are multiple components at a station, the average of FI values of available components is used as the FI value for this station. The frequency index of a seismic event is defined as the average of FIs at all stations that have recorded this event (Matoza et al., 2014).

Text S2. Performance metrics

Since phase picking is not a binary classification task, we need to redefine positive and negative to calculate precision, recall and F1-score. There are discrepancies in performance reporting among different researchers. Some studies consider a waveform trace as a true positive as long as there is a predicted pick sufficiently close to the labeled pick on this waveform (Zhu & Beroza, 2019; Mousavi et al., 2020). However, false predictions may be

underestimated when the model predicts incorrect picks at the same time, leading to a higher reported precision.

Here, we base the definition of positives and negatives on sampling points (points sampled from a continuous analog signal) instead of entire waveform traces. The model output is time series of “probability” of P and S. To get predicted picks from the probability time output by the models, we first extract segments of probability curves above a given decision threshold and the peak positions of these extracted segments are considered as predicted pick times. If a predicted pick occurs within a threshold around a true pick, it is counted as a true positive prediction (TP). Following (Mousavi et al., 2020), the threshold is chosen as 0.5s. Note that this threshold for distinguishing true picks from false picks is different from the probability threshold which is used to extract picks from a “probability” time series output by the model. In the case where there are multiple predicted picks near the true pick, they are counted as only one true positive. If there are no predictions within 0.5s around a true pick, it is counted as a false negative (FN). If there are no true picks around a predicted pick, it is counted as a false positive (FP). Precision is the fraction of predicted picks that are actually correct, calculated as $TP/(TP + FP)$. Recall is the fraction of testing manual picks that have been correctly identified by the model, calculated as $TP/(TP + FN)$. F1 score is calculated as $2 \times (\text{Precision} + \text{Recall})/(\text{Precision} + \text{Recall})$, which is the harmonic mean of the precision and recall. Those samples that are not labeled as true phase arrivals and also not picked by the model are considered as true negatives (TN). For example, considering a 30s waveform with a sampling rate of 100Hz which contains 3001 samples, if there is one manual

P pick and the model gives 10 predicted picks one of which is close to the manual P pick, there are 1 TP, 9 FPs and 2991 TNs. Considering that true negatives are not involved in precision and recall and they heavily outnumber TP, FP and FN, we do not count the number of true negatives when calculating precision, recall and F1 score.

Münchmeyer et al. (2022) evaluated the performance of a model in terms of 3 tasks: (1) event detection, (2) phase identification and (3) onset time picking. For event detection, they used 1 minus noise probability as the score for detection. If the peak detection score for a waveform is above a given threshold, the waveform is considered as a positive detection. ROC (receiver operating characteristic curve) and its AUC (area under the curve) value are used to evaluate the detection performance. In phase identification, they used the ratio of the maximum value of P probability to the maximum value of S probability as the decision score. The Matthews correlation coefficient was used to evaluate the phase identification. The fraction of outliers, root mean square error and mean absolute error were used to evaluate onset time picking. To generate a proper testing set for phase identification and onset time picking, they randomly selected a 10s window around P or S arrivals for each testing waveform, and make sure only one phase is located in the selected window. However, this way of evaluation use the maximum probability value within the tested window as the prediction result, which is different from practical applications of a deep-learning picker where a trigger algorithm is used to retrieve picks from an output probability curve.

References

Buurman, H., & West, M. E. (2010). Seismic precursors to volcanic explosions during the

- 2006 eruption of Augustine Volcano. In J. A. Power, M. L. Coombs, & J. T. Freymueller (Eds.), *The 2006 eruption of Augustine Volcano, Alaska* (pp. 41–57). U.S. Geological Survey.
- Matoza, R. S., Shearer, P. M., & Okubo, P. G. (2014). High-precision relocation of long-period events beneath the summit region of kīlauea volcano, hawai'i, from 1986 to 2009. *Geophysical Research Letters*, *41*(10), 3413-3421.
- Michellini, A., Cianetti, S., Gaviano, S., Giunchi, C., Jozinović, D., & Lauciani, V. (2021). Instance – the italian seismic dataset for machine learning. *Earth System Science Data*, *13*(12), 5509–5544.
- 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. *Nature communications*, *11*(1), 3952.
- Mousavi, S. M., Sheng, Y., Zhu, W., & Beroza, G. C. (2019). Stanford earthquake dataset (stead): A global data set of seismic signals for ai. *IEEE Access*, *7*, 179464–179476.
- Münchmeyer, J., Woollam, J., Rietbrock, A., Tilmann, F., Lange, D., Bornstein, T., ... others (2022). Which picker fits my data? a quantitative evaluation of deep learning based seismic pickers. *Journal of Geophysical Research: Solid Earth*, *127*(1), e2021JB023499.
- Song, Z., Tan, Y. J., & Roman, D. C. (2023). Deep long-period earthquakes at akutan volcano from 2005 to 2017 better track magma influxes compared to volcano-tectonic earthquakes. *Geophysical Research Letters*, *50*(10), e2022GL101987.
- Zhu, W., & Beroza, G. C. (2019). Phasenet: a deep-neural-network-based seismic arrival-

time picking method. *Geophysical Journal International*, 216(1), 261–273.

Table S1. The number of waveform traces in our dataset, including volcano-tectonic earthquakes (VTs), long-period earthquakes (LPs) and noise. The number of corresponding events is given in brackets. Since different waveforms may originate from the same source, the sum of events in the training, validation and testing sets does not necessarily equal the total events. Note that splitting waveforms from the same event to different data sets does not result in data leakage, because waveforms recorded at different stations have been influenced by different path effects and different background noise, thus representing unique examples. The 4,841 LP waveforms and 4,841 VT waveforms in northern California as well as the 6,224 LFE waveforms in Japan that are used as extra test sets are not shown in this table.

	Total traces	Training set	Validation set	Test set
Whole dataset	323,088 (70,352)	270,224 (68,996)	17,744 (13,346)	35,120 (23,700)
Earthquake	303,088 (70,352)	257,763 (68,996)	15,190 (13,346)	30,135 (23,700)
Noise	20,000 (0)	12,461 (0)	2,554 (0)	4,985 (0)
LP earthquakes	151,431 (33,886)	128,802 (33,364)	7,551 (6,609)	15,078 (11,798)
VT earthquakes	151,657 (36,466)	128,961 (35,632)	7,639 (6,737)	15,057 (11,902)
Alaska LPs	51,942 (15,701)	44,263 (15,511)	2,544 (2,370)	5,135 (4,497)
Alaska VTs	50,899 (15,519)	43,198 (15,151)	2,598 (2,377)	5,103 (4,354)
Hawaii LPs	16,906 (2,351)	14,404 (2,323)	811 (666)	1,691 (1,132)
Hawaii VTs	16,814 (2,766)	14,346 (2,702)	806 (653)	1,662 (1,119)
Japan LPs	82,583 (15,834)	70,135 (15,530)	4,196 (3,573)	8,252 (6,169)
Japan VTs	83,944 (18,181)	71,417 (17,779)	4,235 (3,707)	8,292 (6,429)

Table S2. Performance metrics on the validation set for 12 PhaseNet networks trained with different hyperparameters: learning rate, batch size and σ_{label} which is the standard deviation of the Gaussian function used for labeling training data. The networks were randomly initialized. Each model is trained up to 400 epochs, and the epoch at which the loss on the validation set is the lowest is saved as the final result. MAE is the mean absolute error of picks. The picking residuals outside the interval $(-1, 1)s$ are considered as outliers and not involved in the calculation of MAE. The row with the highest F1 score is highlighted in bold face. For each network, we have tried various decision thresholds and choose the one with the highest F1-score as the optimal threshold. Figure S16 presents the threshold tuning for preferred models.

Network	Hyperparameters			Decision threshold		F1 score		MAE (s)	
	Batch size	Learning rate	σ_{label}	P picking	S picking	P picking	S picking	P picking	S picking
PhaseNet	1024	0.0010	20	0.31	0.34	0.9169	0.8842	0.0767	0.1148
PhaseNet	1024	0.0010	10	0.29	0.23	0.9110	0.8762	0.0750	0.1186
PhaseNet	1024	0.0005	20	0.32	0.31	0.9158	0.8844	0.0779	0.1162
PhaseNet	1024	0.0005	10	0.29	0.25	0.9124	0.8787	0.0762	0.1162
PhaseNet	1024	0.0001	20	0.32	0.31	0.9090	0.8773	0.0810	0.1182
PhaseNet	1024	0.0001	10	0.30	0.25	0.9005	0.8643	0.0766	0.1200
PhaseNet	512	0.0010	20	0.31	0.31	0.9157	0.8843	0.0778	0.1173
PhaseNet	512	0.0010	10	0.29	0.24	0.9115	0.8756	0.0745	0.1187
PhaseNet	512	0.0005	20	0.39	0.34	0.9181	0.8866	0.0755	0.1146
PhaseNet	512	0.0005	10	0.28	0.24	0.9134	0.8782	0.0758	0.1184
PhaseNet	512	0.0001	20	0.37	0.34	0.9106	0.8805	0.0788	0.1171
PhaseNet	512	0.0001	10	0.27	0.24	0.9001	0.8644	0.0809	0.1230

Table S3. Performance metrics on the validation set for 12 EQTransformer networks trained with different hyperparameters: learning rate, batch size and σ_{label} which is the standard deviation of the Gaussian function used for labeling training data. The networks were randomly initialized. Each model is trained up to 400 epochs, and the epoch at which the loss on the validation set is the lowest is saved as the final result. MAE is the mean absolute error of picks. The picking residuals outside the interval $(-1, 1)s$ are considered as outliers and not involved in the calculation of MAE. The row with the highest F1 score is highlighted in bold face. For each network, we have tried various decision thresholds and choose the one with the highest F1-score as the optimal threshold. Figure S16 presents the threshold tuning for preferred models.

Network	Hyperparameters			Decision threshold		F1 score		MAE (s)	
	Batch size	Learning rate	σ_{label}	P picking	S picking	P picking	S picking	P picking	S picking
EQTransformer	1024	0.0010	20	0.22	0.25	0.9245	0.8919	0.0877	0.1242
EQTransformer	1024	0.0010	10	0.15	0.16	0.9212	0.8878	0.0856	0.1182
EQTransformer	1024	0.0005	20	0.23	0.24	0.9216	0.8905	0.0911	0.1271
EQTransformer	1024	0.0005	10	0.16	0.15	0.9176	0.8861	0.0860	0.1241
EQTransformer	1024	0.0001	20	0.23	0.27	0.9149	0.8842	0.0956	0.1307
EQTransformer	1024	0.0001	10	0.15	0.16	0.9148	0.8814	0.0924	0.1283
EQTransformer	512	0.0010	20	0.19	0.27	0.9232	0.8887	0.0887	0.1238
EQTransformer	512	0.0010	10	0.17	0.13	0.9213	0.8869	0.0855	0.1230
EQTransformer	512	0.0005	20	0.22	0.23	0.9216	0.8916	0.0895	0.1243
EQTransformer	512	0.0005	10	0.13	0.15	0.9191	0.8855	0.0868	0.1215
EQTransformer	512	0.0001	20	0.20	0.18	0.9166	0.8817	0.0957	0.1350
EQTransformer	512	0.0001	10	0.15	0.14	0.9133	0.8798	0.0900	0.1264

Table S4. Performance metrics of the models trained with random initial weights and those first initialized with pre-trained weights. The performance is evaluated on the validation set. For pre-training, we use the network weights pre-trained on INSTANCE dataset (Münchmeyer et al., 2022) as the starting point before training. The hyperparameters batch size, learning rate and σ_{label} are the same as the preferred ones highlighted in bold face in Table S2-S3.

Network	Initialized with weights pre-trained on	Decision threshold		F1 score		MAE (s)	
		P picking	S picking	P picking	S picking	P picking	S picking
EQTransformer	None	0.22	0.25	0.9245	0.8919	0.0877	0.1242
EQTransformer	INSTANCE	0.22	0.22	0.9250	0.8916	0.0876	0.1256
PhaseNet	None	0.39	0.34	0.9181	0.8866	0.0755	0.1146
PhaseNet	INSTANCE	0.39	0.34	0.9175	0.8833	0.0750	0.1165

Table S5. AUC scores for the event detection task defined by (Münchmeyer et al., 2022, section 2.1.1), which is much simpler than picking. If the peak of the output probability curve for a test example is larger than the threshold, it is considered as a positive prediction.

Model	LP test set	VT test set
EQTransformer retrained in this study	0.9993	0.9994
PhaseNet retrained in this study	0.9992	0.9994
Original EQTransformer(Mousavi et al., 2020)	0.9776	0.9839
Original PhaseNet (Zhu & Beroza, 2019)	0.9932	0.9937
EQTransformer trained on INSTANCE (Münchmeyer et al., 2022)	0.9934	0.9907
PhaseNet trained on INSTANCE (Münchmeyer et al., 2022)	0.9695	0.9784

Table S6. Matthews correlation coefficients for the phase discrimination task (Münchmeyer et al., 2022).

Model	LP test set	VT test set
EQTransformer retrained in this study	0.9621	0.9787
PhaseNet retrained in this study	0.9570	0.9764
Original EQTransformer (Mousavi et al., 2020)	0.7899	0.9086
Original PhaseNet (Zhu & Beroza, 2019)	0.7333	0.9354
EQTransformer trained on INSTANCE (Münchmeyer et al., 2022)	0.7717	0.9422
PhaseNet trained on INSTANCE (Münchmeyer et al., 2022)	0.8330	0.9463

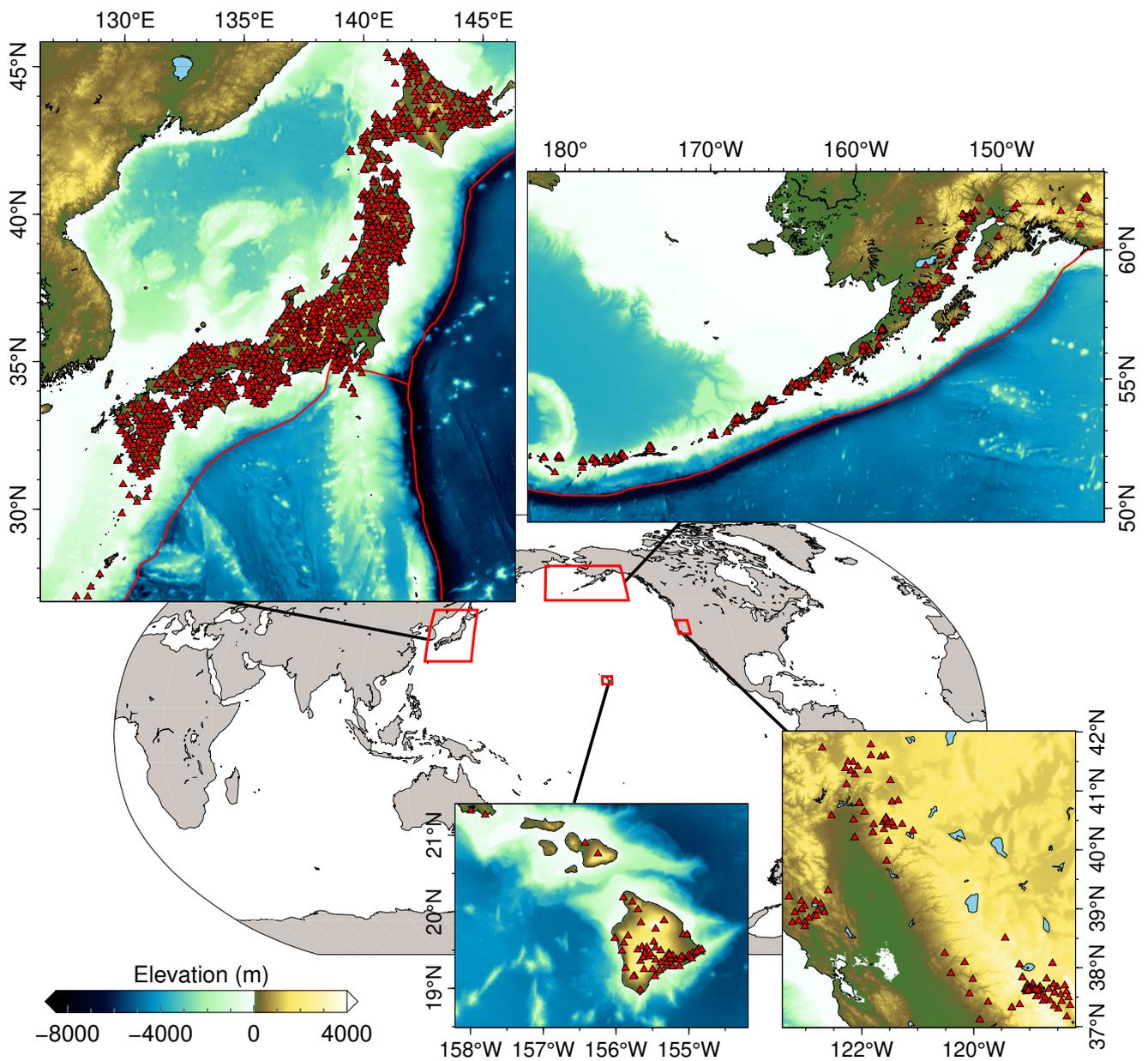


Figure S1. The geographical distribution of seismic stations (red triangles) with waveforms included in our data set, including the data set in Table S1, the northern California test set and the test set of Japan tectonic LPs.

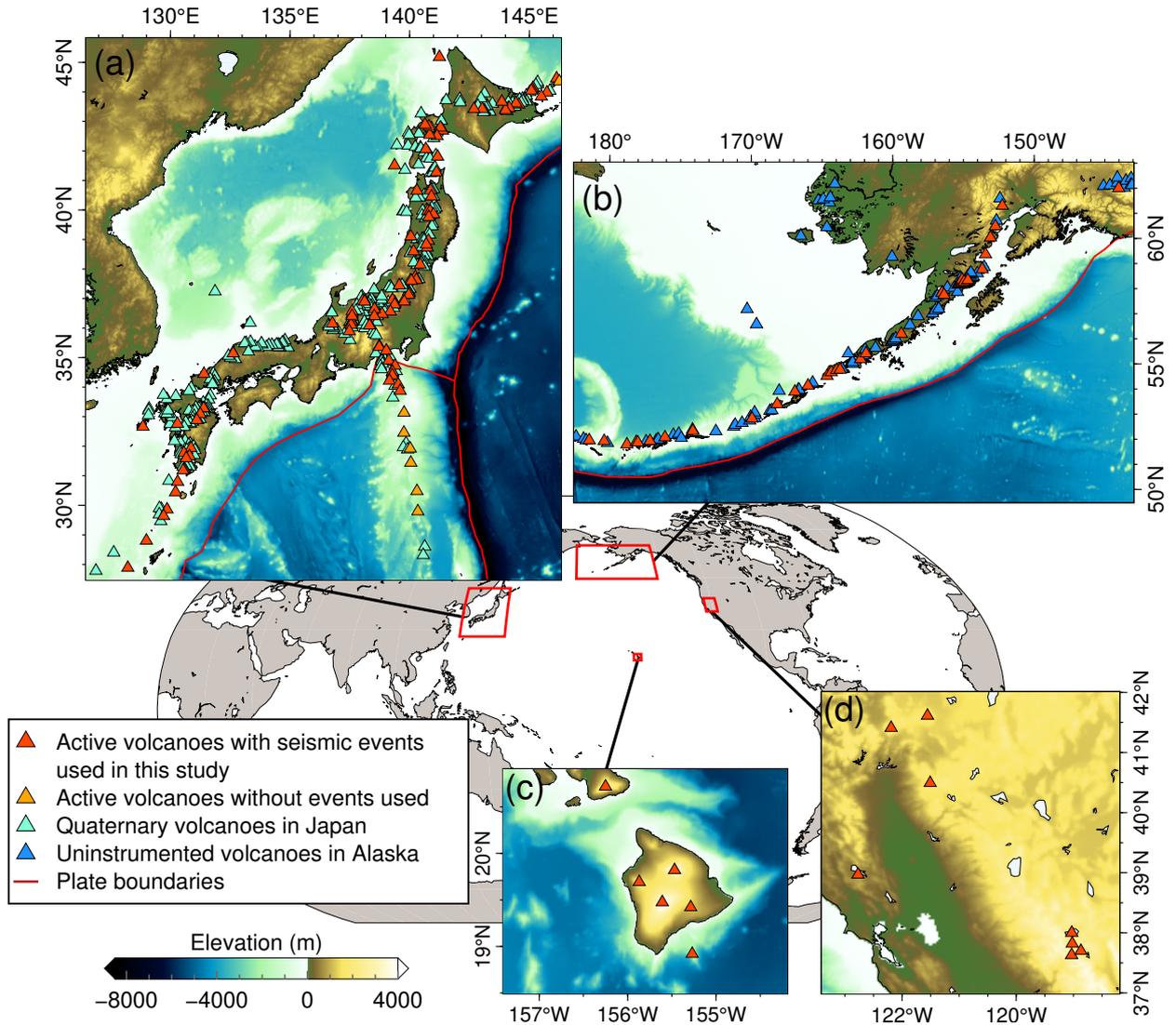


Figure S2. The geographical distribution of the active volcanoes with seismic events included in our data set (red triangles). The active volcanoes in Japan without seismic events use (orange triangles), quaternary volcanoes in Japan (green triangles) and uninstrumented volcanoes in Alaska (blue triangles) are also shown.

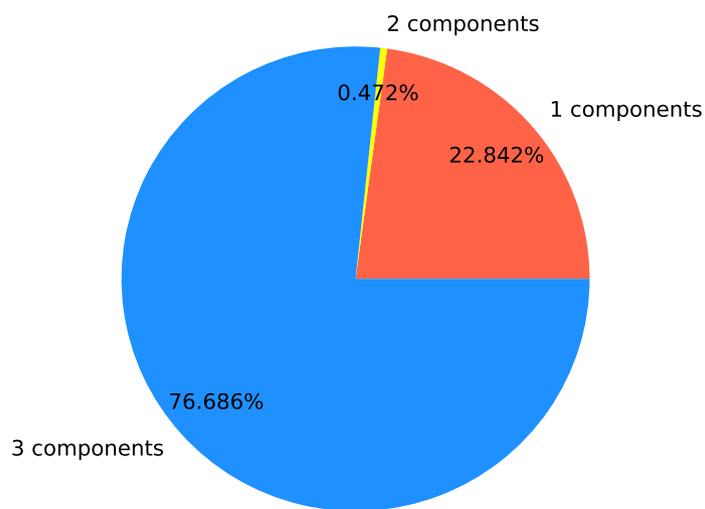


Figure S3. The proportion of seismograms with different numbers of components in Table S1. For one-component or two-component records, we fill in zeros for the remaining components before training.

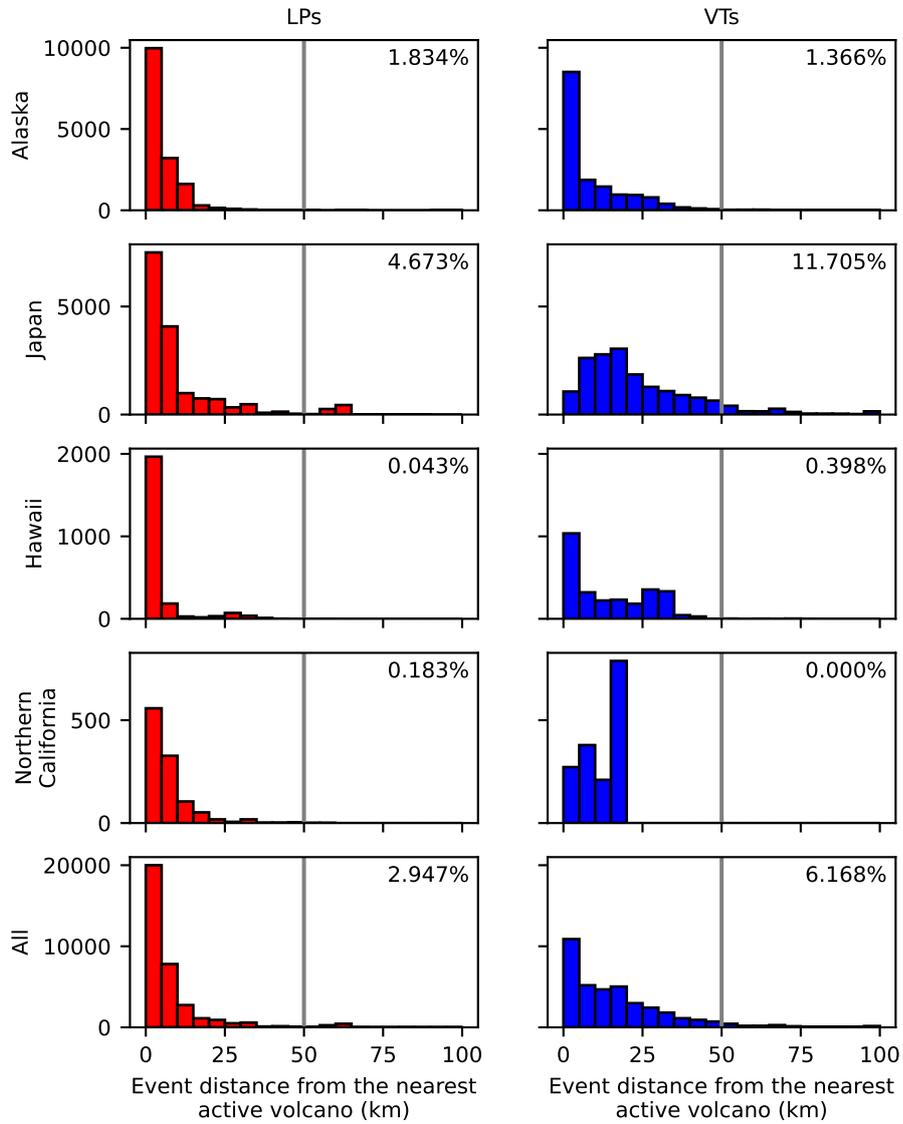


Figure S4. Histogram of the distances of events in our data set to the nearest active volcano. The numbers in the top right corner indicate the fractions of events that are more than 50 km away from the nearest volcano.

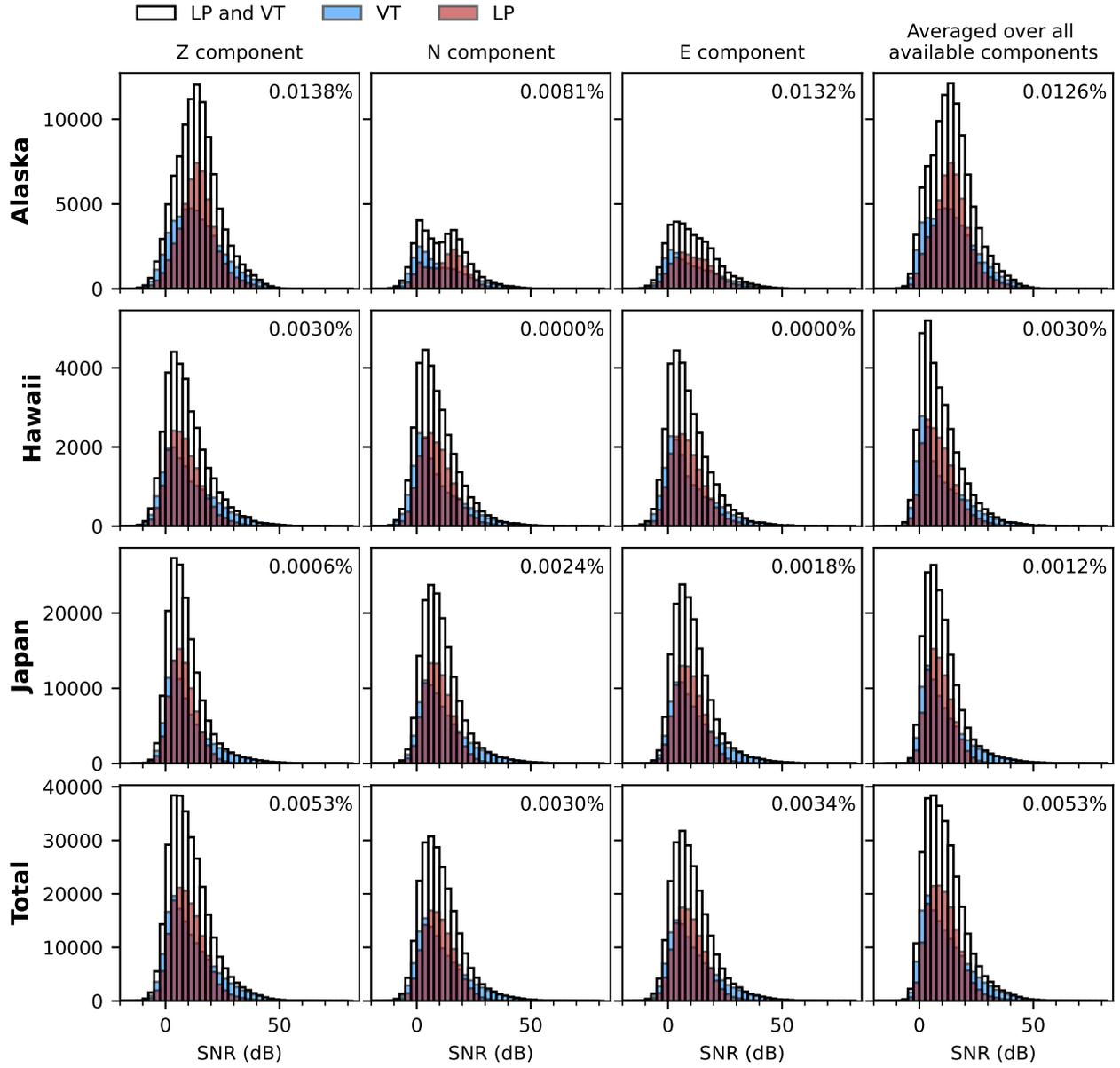


Figure S5. The distribution of signal-to-noise ratios (SNR) of the waveform traces in Table S1. Each column shows SNRs of waveforms from different regions for the same component, while each row represents SNRs of waveforms from the same region but for different components. The numbers in the top right corner indicate the fraction of samples outside the range of the x -axis. SNR is calculated as $\text{SNR} = 20 \log_{10} \frac{|S|_{95}}{|N|_{95}}$, where $|S|_{95}$ is the 95 percentile of absolute amplitudes in a 5s window right after the S arrival and $|N|_{95}$ is that in a 5s window before the P arrival.

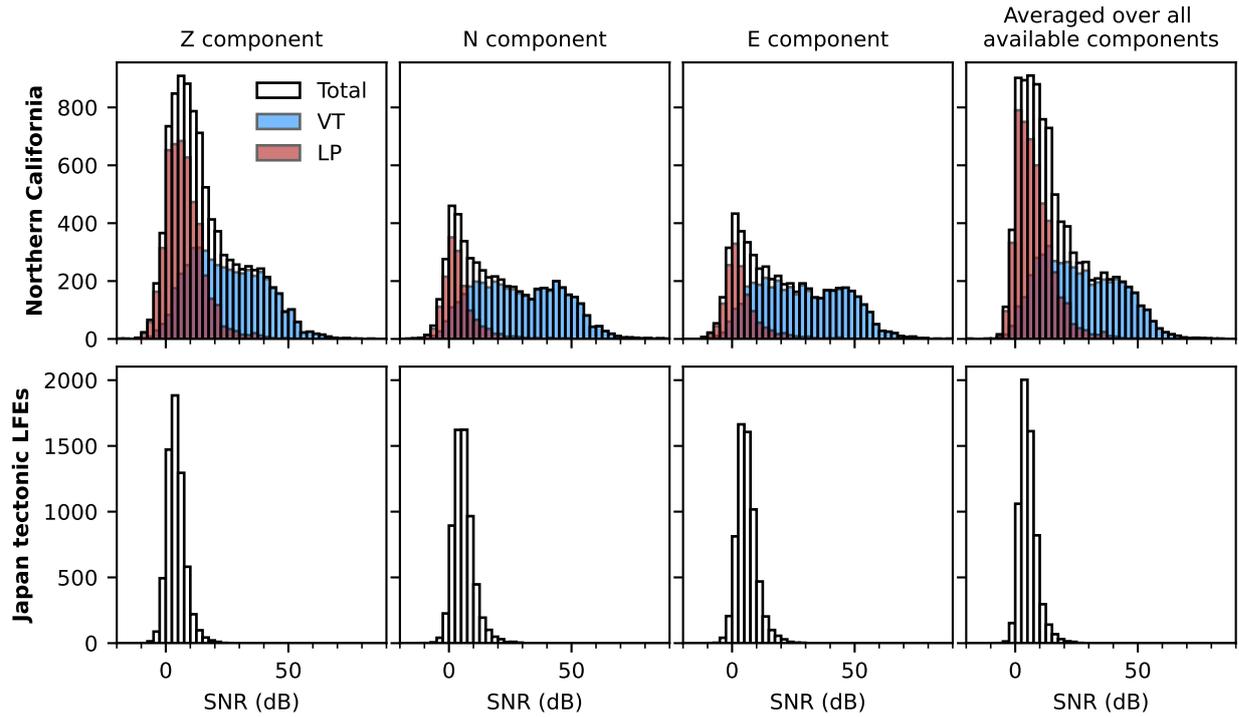


Figure S6. The distribution of signal-to-noise ratios (SNR) of the waveform traces in the northern California test set and the test set of Japan tectonic LFEs.

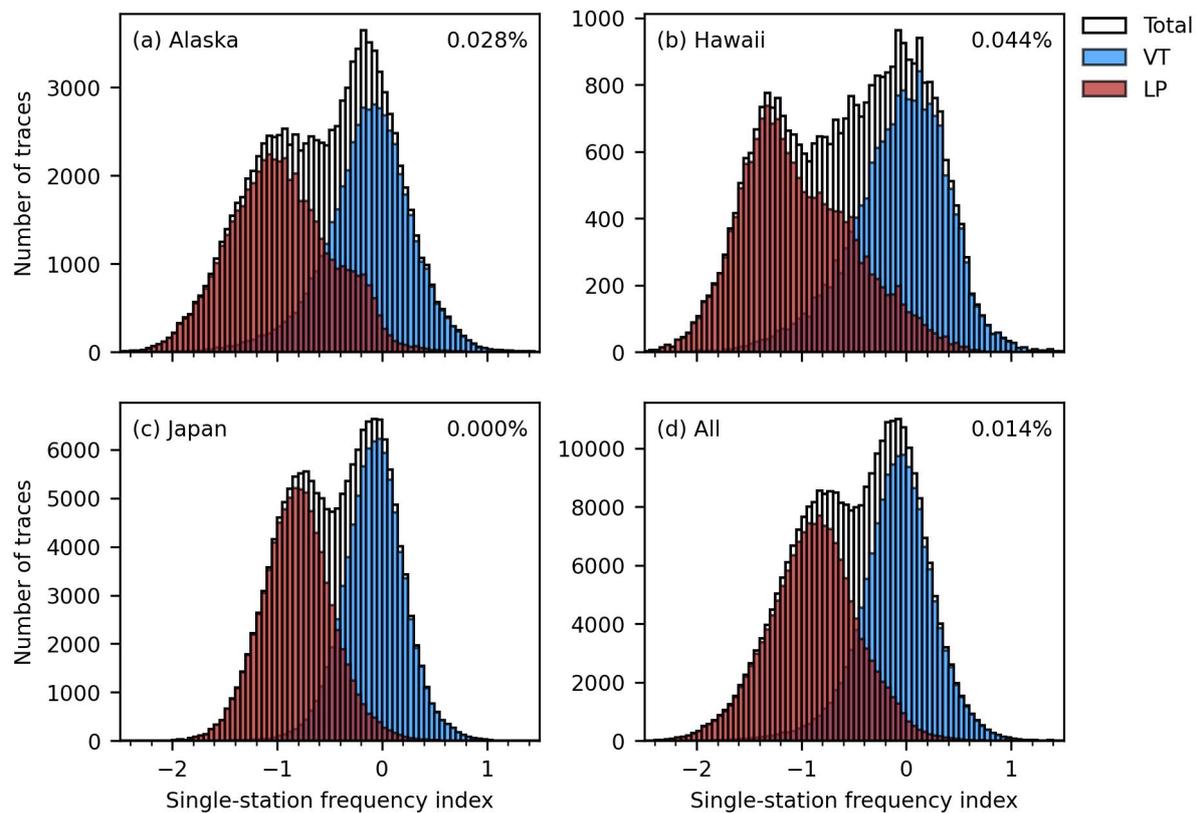


Figure S7. The distribution of single-station frequency index (FI) values of the earthquake waveforms in Table S1. The numbers in the top right corner indicate the fraction of samples outside the range of the plotted x -axis.

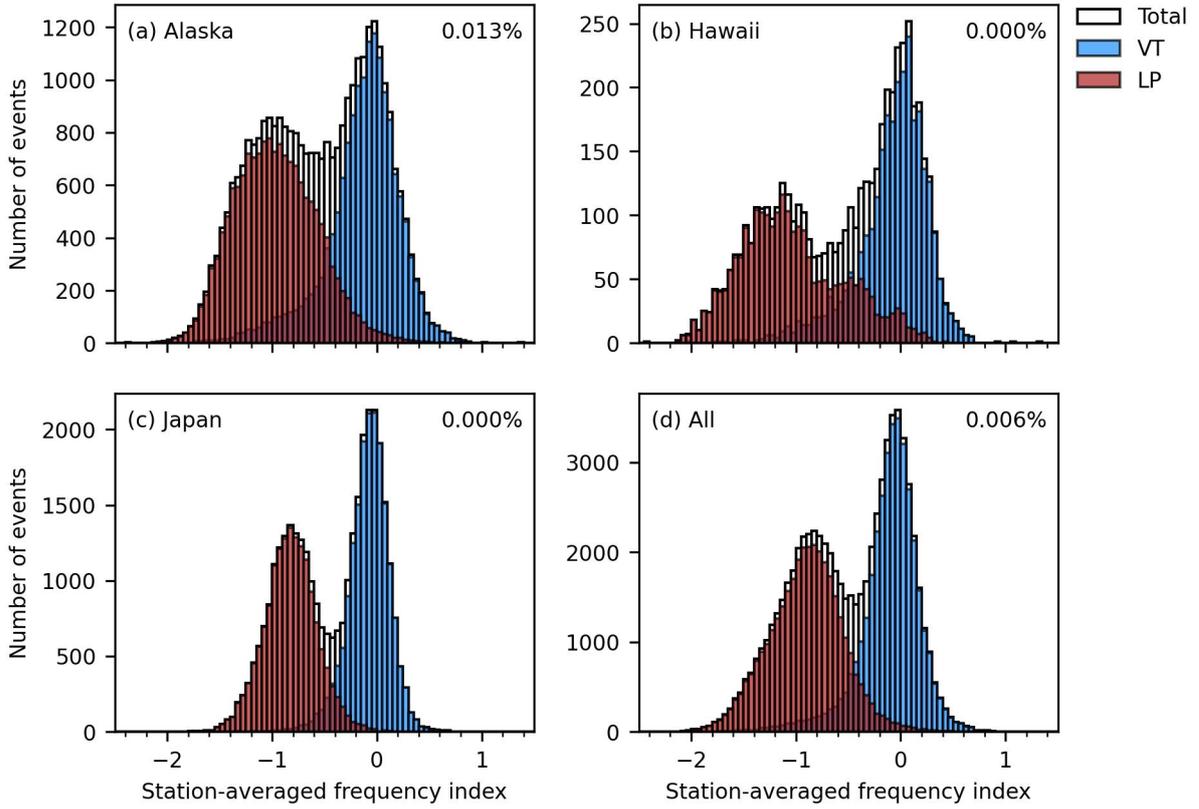


Figure S8. The distribution of event-based frequency index (FI) values of the earthquakes in Table S1, which are calculated by averaging FI values over all recording stations. The numbers in the top right corner indicate the fraction of samples outside the range of the plotted x -axis.

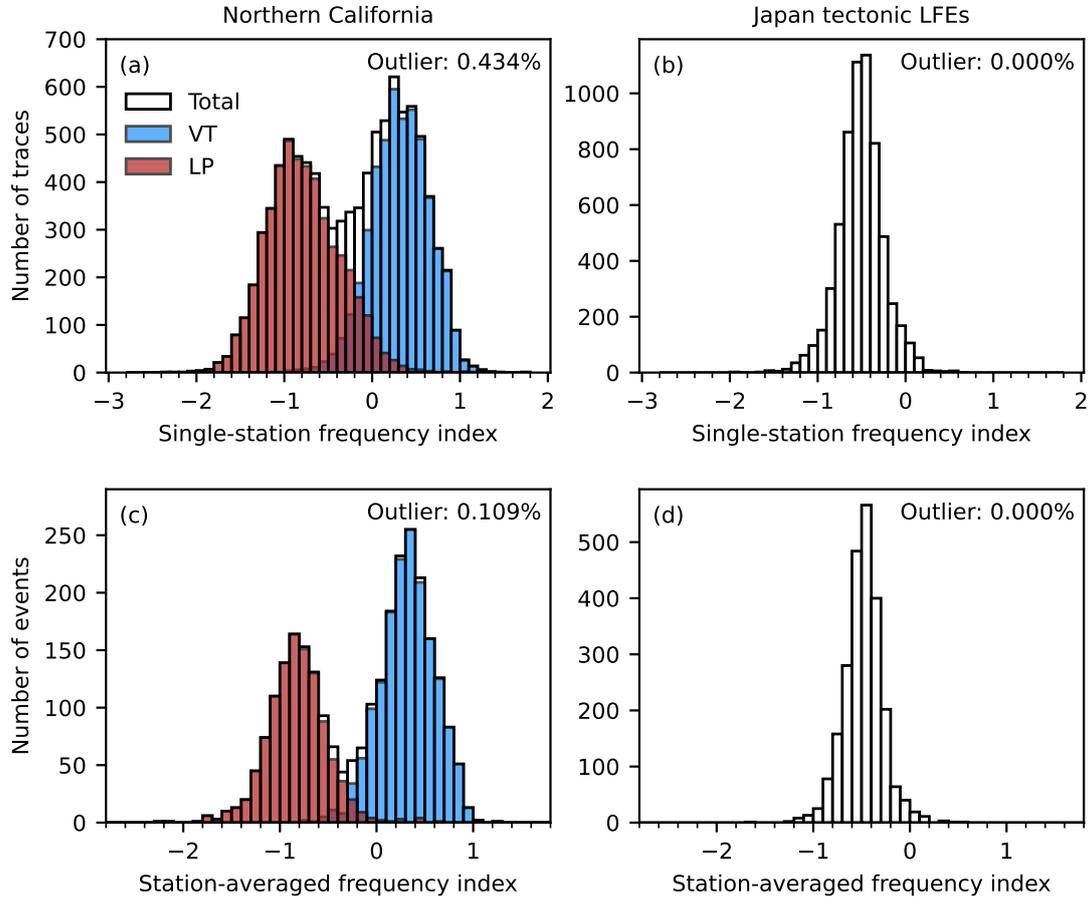


Figure S9. The distribution of frequency index (FI) values for the northern California test set (a, c) and the test set of Japan tectonic LFEs (b, d). The top row (a, b) and the bottom row (c, d) show the single-station FI and the event-based FI, respectively. The numbers in the top right corner indicate the fraction of samples outside the range of the plotted x -axis.

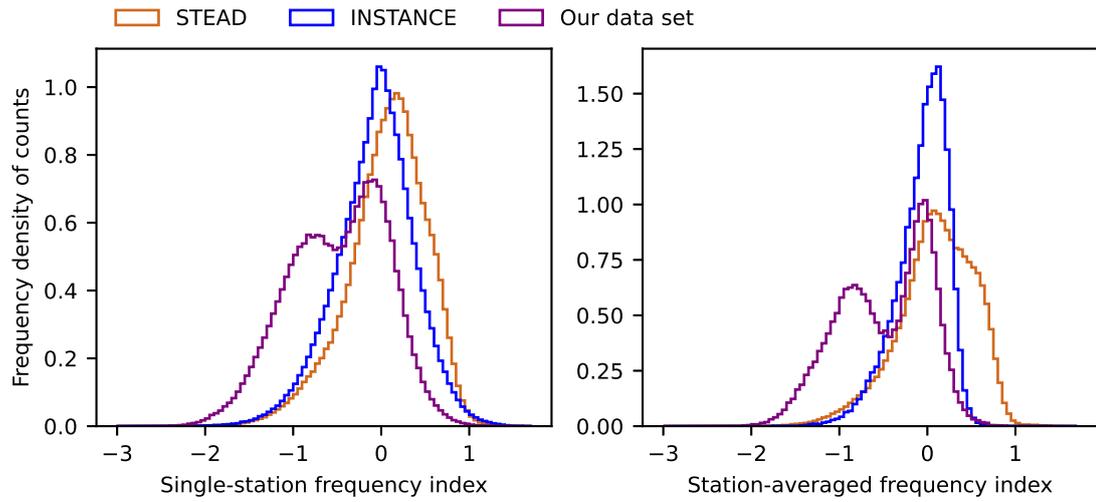


Figure S10. Comparison of frequency index distributions of INSTANCE (Michelini et al., 2021), STEAD (Mousavi et al., 2019) and our data set. The y axis, frequency density, is defined

as $\frac{\text{Counts in a bin}/\text{Total counts}}{\text{Bin width}}$.

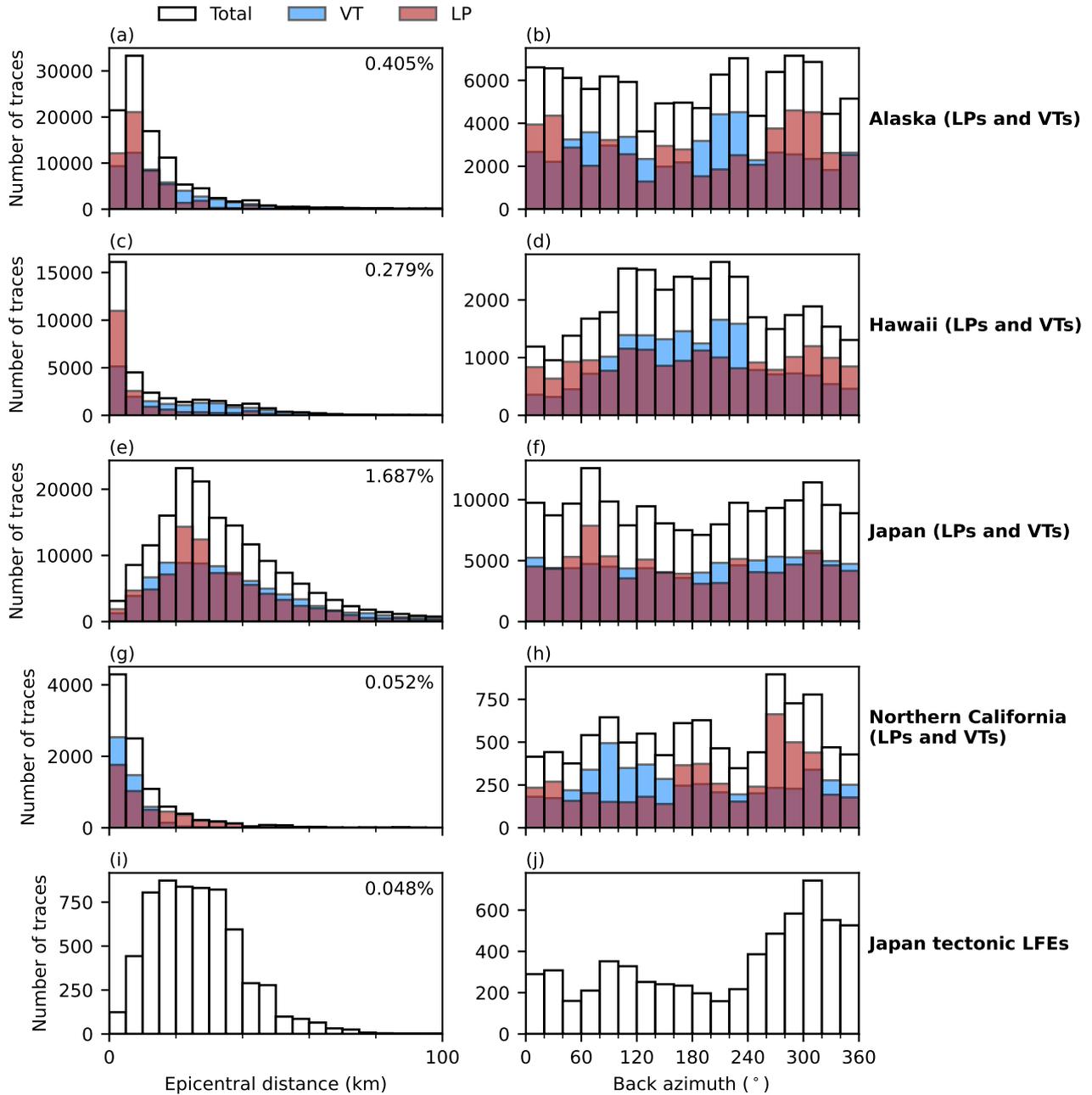


Figure S11. The distribution of epicentral distances (the first column) and back azimuths (the second column) of LP and VT waveforms from Alaska (a, b), Hawaii (c, d), Japan (e, f), northern California (g, h) and tectonic LFE waveforms from Japan (i, j). We adopt the logarithmic scale in the first column to make the number of traces with large epicentral distances visible. The fraction of traces recorded at an epicentral distance greater than 100 km is shown in the top right corner of each panel in the first column.

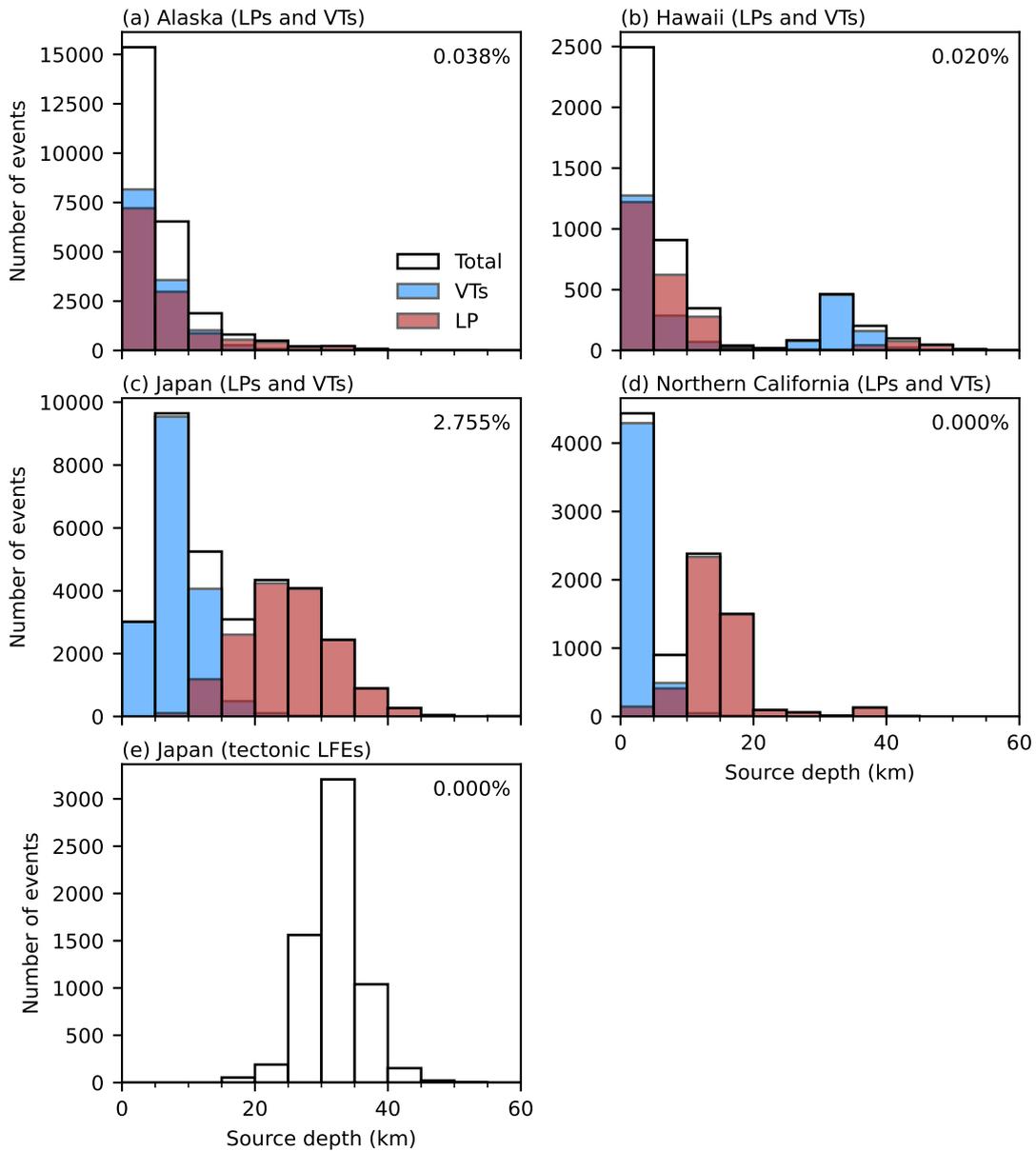


Figure S12. The distribution of source depths of the VTs and LPs in our data set from Alaska (a), Hawaii (b), Japan (c) northern California (d) as well as tectonic LFEs near the Nankai trough from Japan (e). The fraction of events deeper than 60 km is shown in the top right corner of each panel.

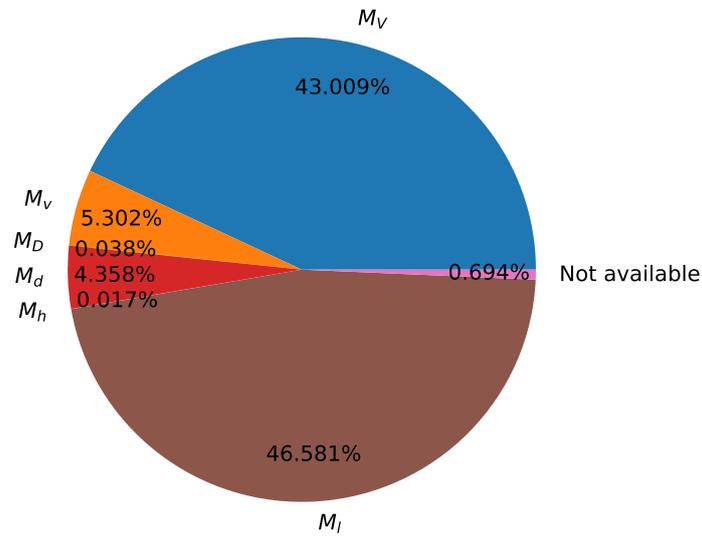


Figure S13. The proportion of different magnitude types. M_l is the local magnitude. M_d is the duration magnitude. M_h is nonstandard magnitudes used by USGS (<https://www.usgs.gov/programs/earthquake-hazards/magnitude-types>). M_V , M_v and M_D are magnitudes used by JMA (Japan Meteorological Agency), where M_V is the velocity magnitude, M_v is similar to M_V but for only 2 or 3 stations, M_D is the displacement magnitude (https://www.data.jma.go.jp/svd/eqev/data/bulletin/catalog/notes_e.html).

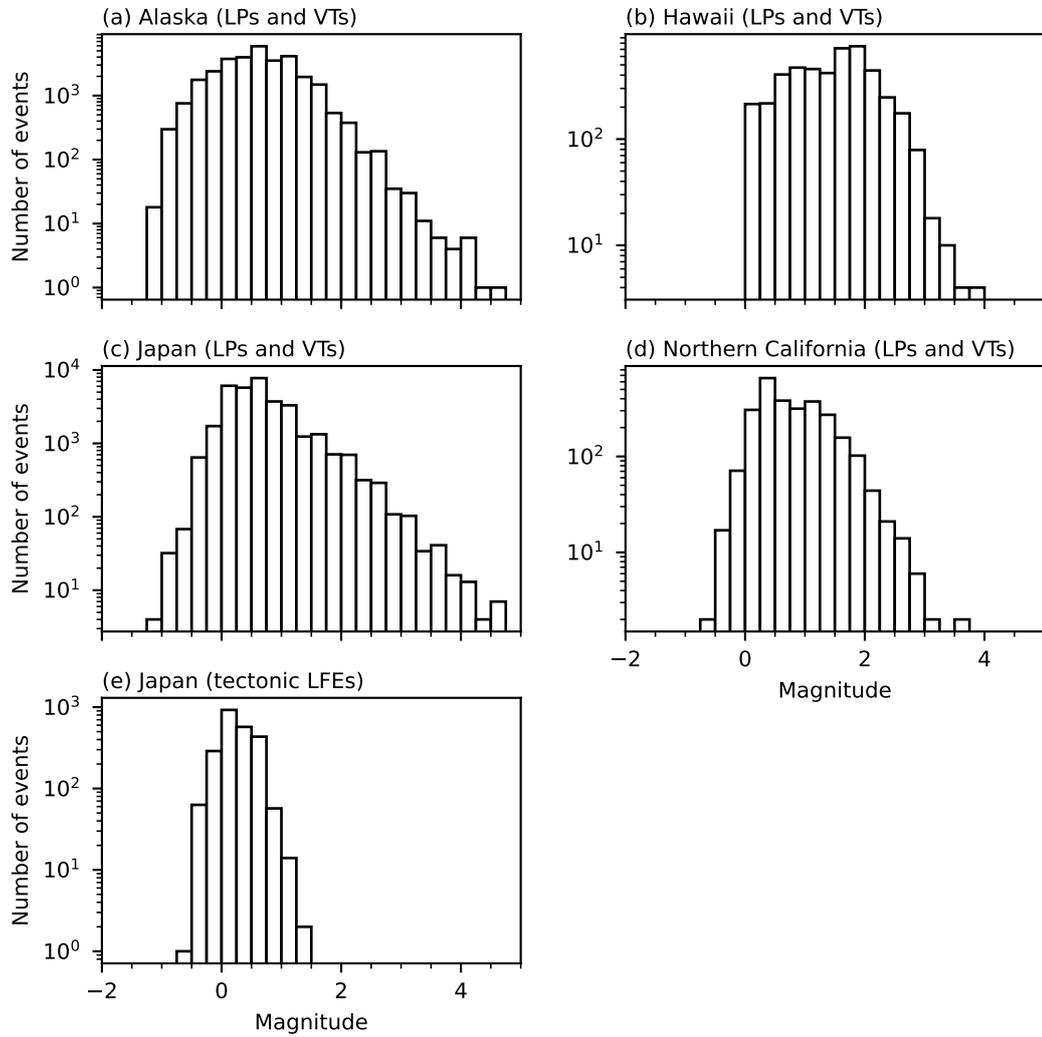


Figure S14. Histogram of magnitudes of the VTs and LPs in our data set from Alaska (a), Hawaii (b), Japan (c) northern California (d) as well as tectonic LFEs near the Nankai trough from Japan (e).

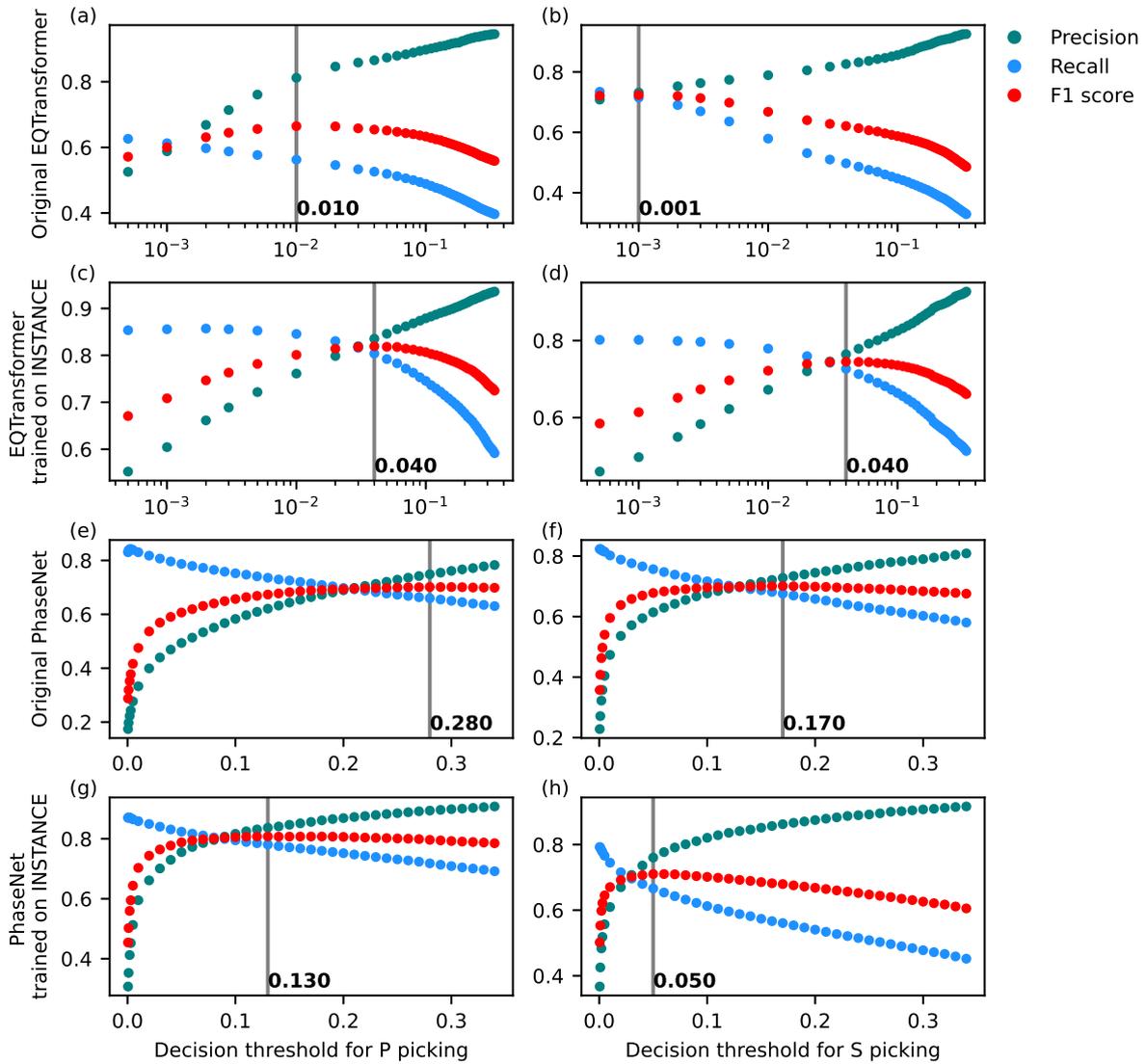


Figure S15. Threshold tuning for the original PhaseNet (Zhu & Beroza, 2019), EQTransformer (Mousavi et al., 2020) and their variants trained on the INSTANCE data set (Münchmeyer et al., 2022) The performance is evaluated on the validation set in Table S1. The left and right columns show the performance metrics for P picking and S picking, respectively. The first (a, b) and second (c, d) rows show the performance metrics of the original EQTransformer network trained on STEAD (Mousavi et al., 2019, 2020) and the EQTransformer network trained on INSTANCE (Michelini et al., 2021; Münchmeyer et al., 2022), respectively. The third (e, f) and fourth (g, h) rows show the performance metrics of the original PhaseNet network trained on California earthquakes (Zhu & Beroza, 2019) and the PhaseNet network trained on INSTANCE (Michelini et al., 2021; Münchmeyer et al., 2022), respectively. The gray lines and the numbers show the optimal thresholds found at the highest F1 scores.

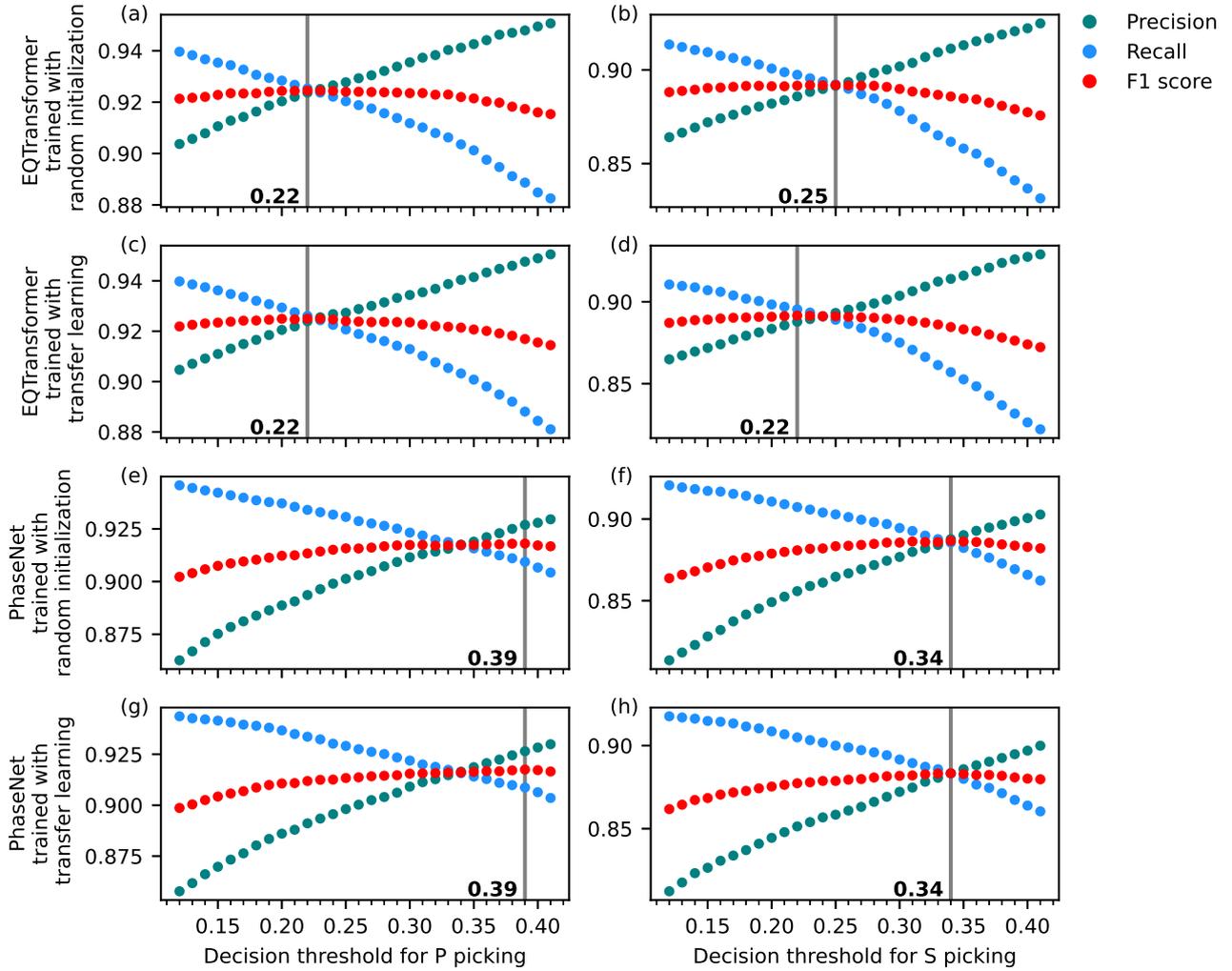


Figure S16. Threshold tuning for EQTransformer and PhaseNet networks trained for volcano seismicity in this study (Table S4). The performance is evaluated on the validation set (development set). The left and right columns show the precision, recall and F1 score for P picking and S picking, respectively. The first (a, b) and second (c, d) rows show the performance metrics of the EQTransformer networks trained with randomly initialized weights and initial weights pre-trained on INSTANCE, respectively. The third (e, f) and fourth (g, h) rows show the performance metrics of the PhaseNet networks trained with randomly initialized weights and initial weights pre-trained on INSTANCE, respectively. The gray lines and the numbers show the optimal thresholds found at the highest F1 scores.

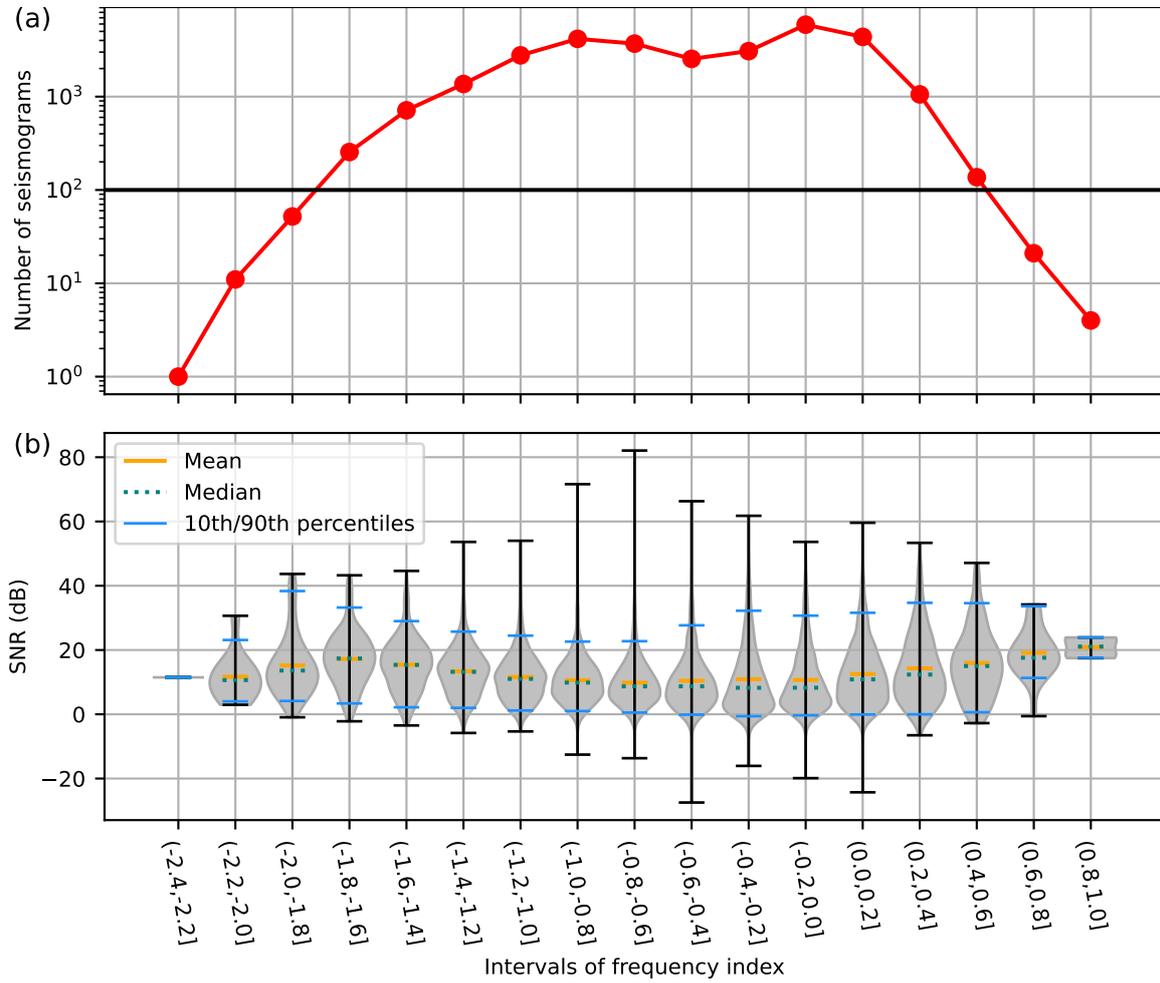
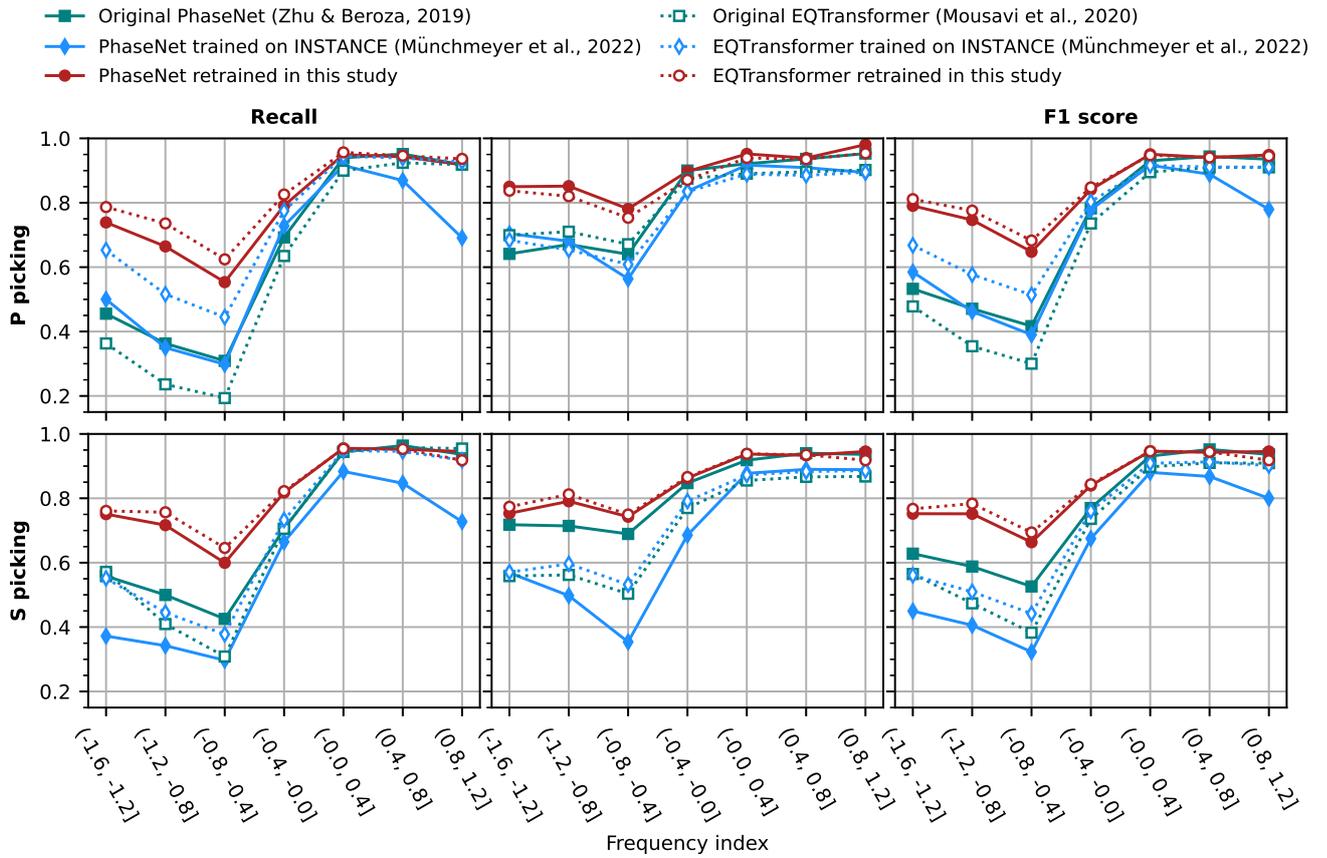
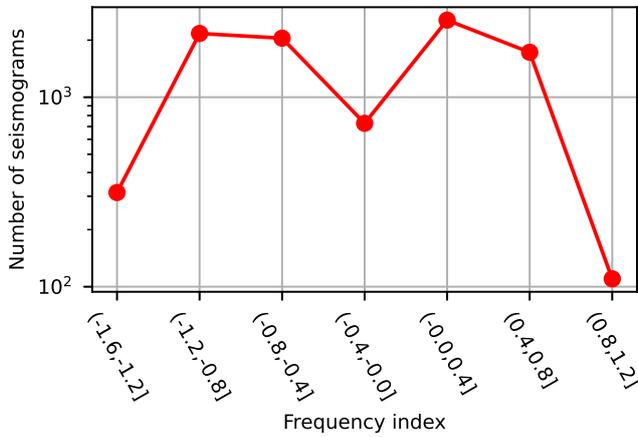


Figure S17. (a) The numbers of waveform traces in the test set for different frequency index bins. We only use the subsets with more than 100 traces for testing, i.e. those above the horizontal black line. (b) The distribution of signal to noise ratio for each subset. The vertical lines show the SNR ranges. The gray area is the estimated probability density for the SNR distribution.

(a) Model performance



(b) Number of seismograms in each bin



(c) SNR distribution

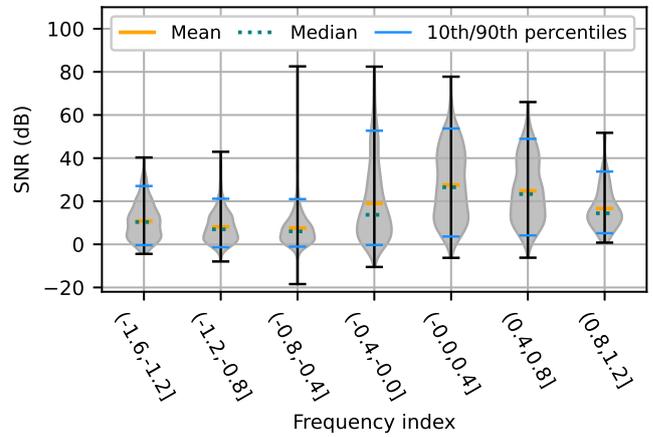


Figure S18. (a) Model performance on subsets of the testing waveforms from northern California. (b) Number of waveforms in each subset. (c) SNR distribution in each subset.

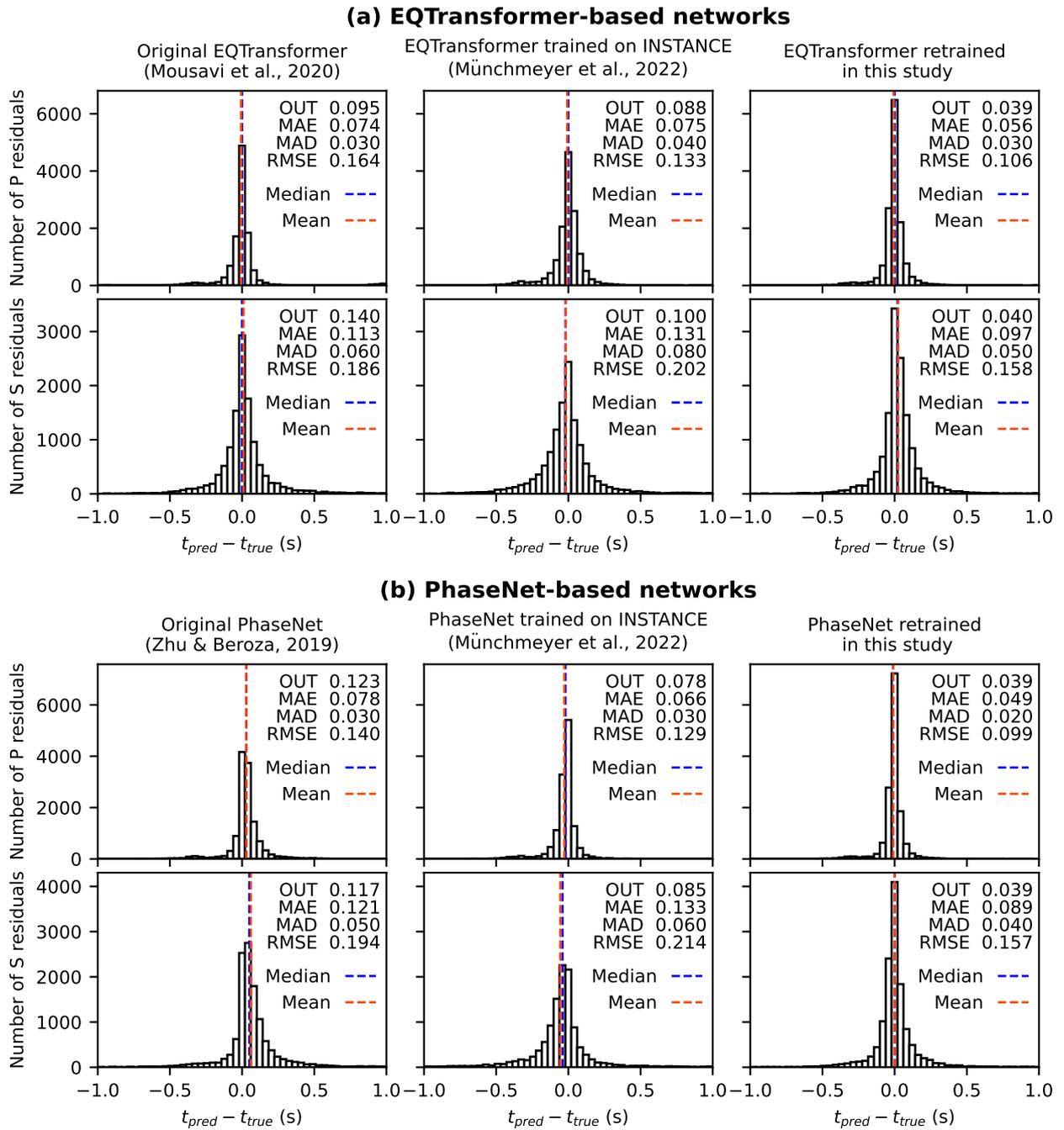


Figure S19. Histogram of residuals between the manual picks in the VT test set and the picks predicted by the EQTransformer-based networks (a) and the PhaseNet-based networks (b). The numbers in the upper right corner show the fraction of residual outside the $(-1, 1)$ s interval (OUT), the mean absolute error (MAE), the median absolute deviation (MAD) and the root mean square error (RMSE). The MAE, RMSE and MAD are calculated only for the residuals within $(-1, 1)$ s to avoid strong influence of outliers.

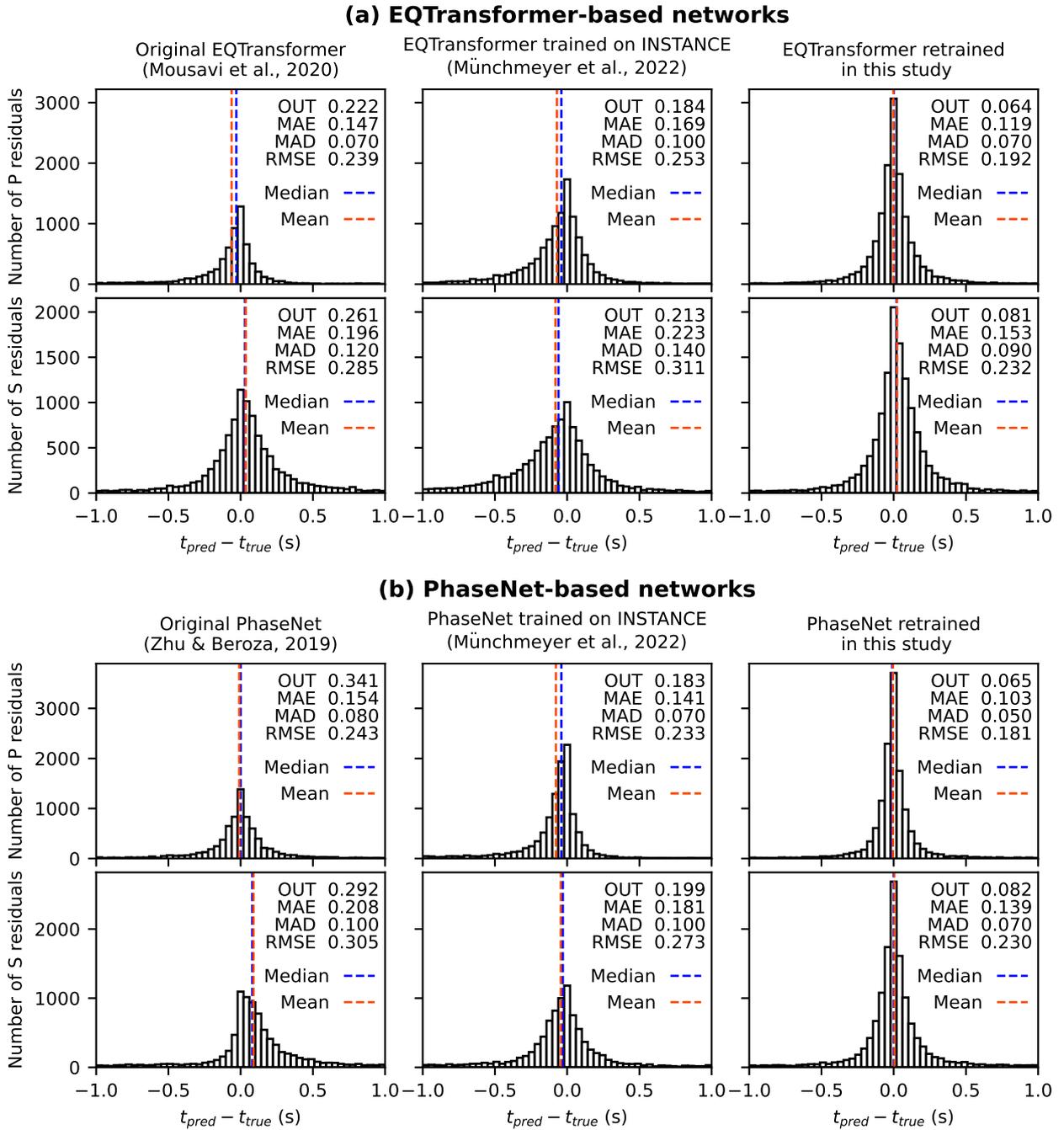


Figure S20. Similar to Figure S19 but for the LP test set.

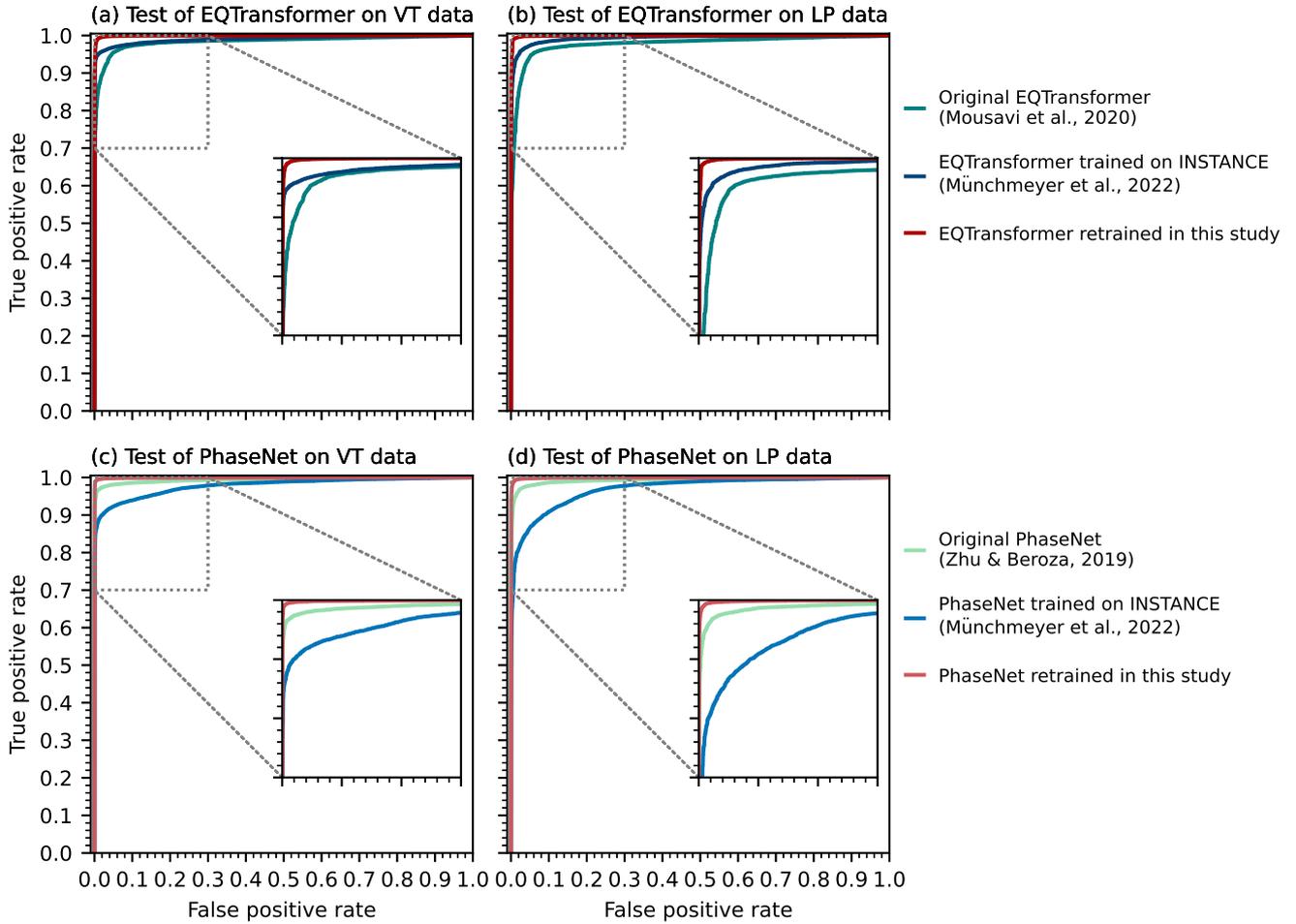


Figure S21. Receiver operating characteristics (ROC) for event detection. The first row shows the ROC curves for the EQTransformer-based networks while the second row is for the PhaseNet-based networks. If the output probability curve for a test example is larger than the decision threshold, it is considered as a positive prediction. The test data are from Alaska, Hawaii and Japan. The LP test set and VT test set (Table S1) are evaluated separately.

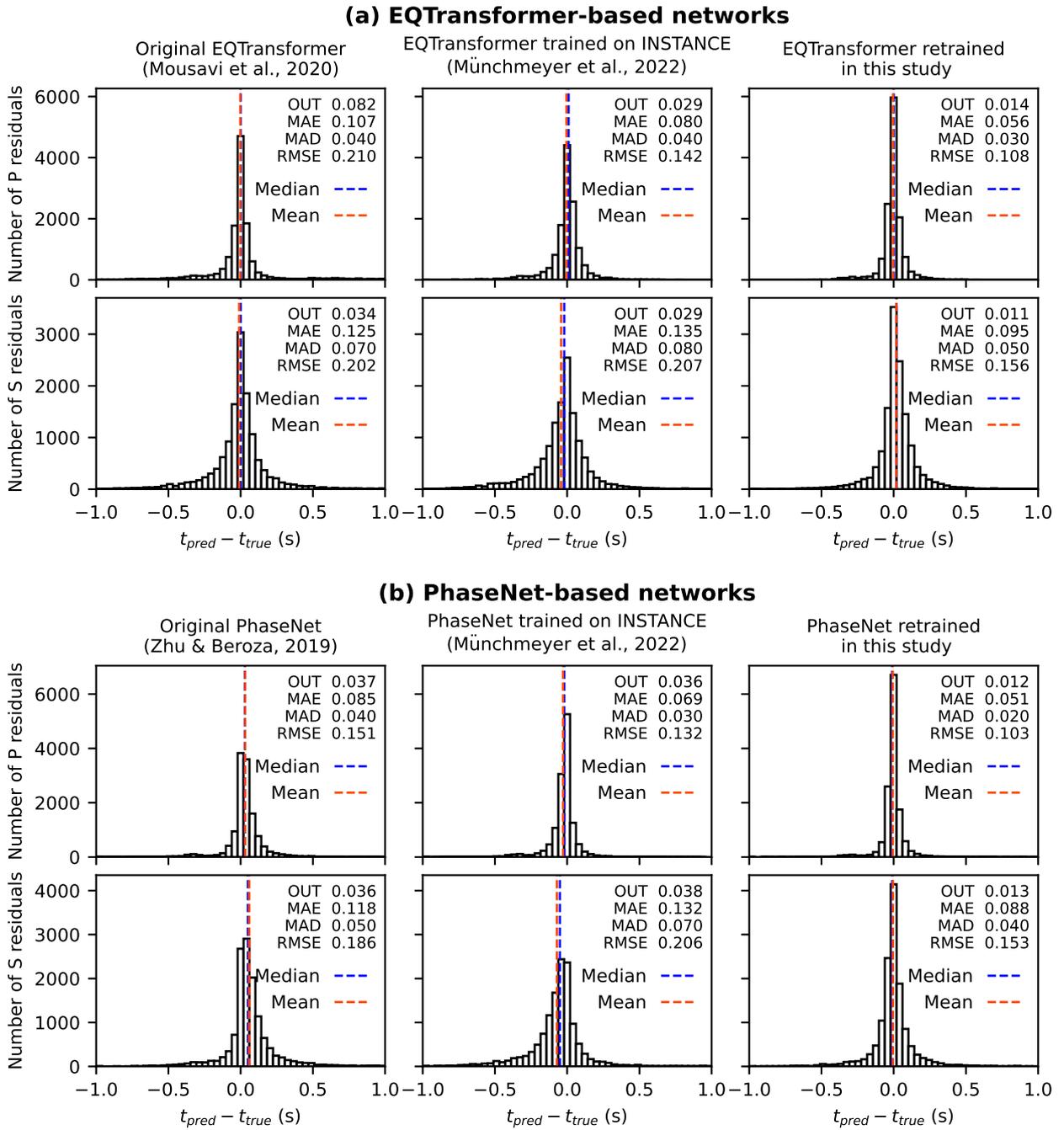


Figure S22. Residuals of phase picks for the VT test set calculated using Münchmeyer et al. (2022)’s evaluation workflow. The difference from Figure S19 is due to the different ways in pre-processing and post-processing. In Münchmeyer et al. (2022)’s workflow, 10s windows containing only P or only S are randomly generated, where the peak position of the output phase probability is taken as the model pick. See (Münchmeyer et al., 2022, Data and Method) for more details. In our evaluation workflow a 30s window containing the P manual pick is randomly generated, which may or may not contains the S pick. We run a trigger algorithm on the output probability curves and find the peaks between trigger on and off times, which may produce more than one model pick for a waveform trace even though there is only one ground truth.

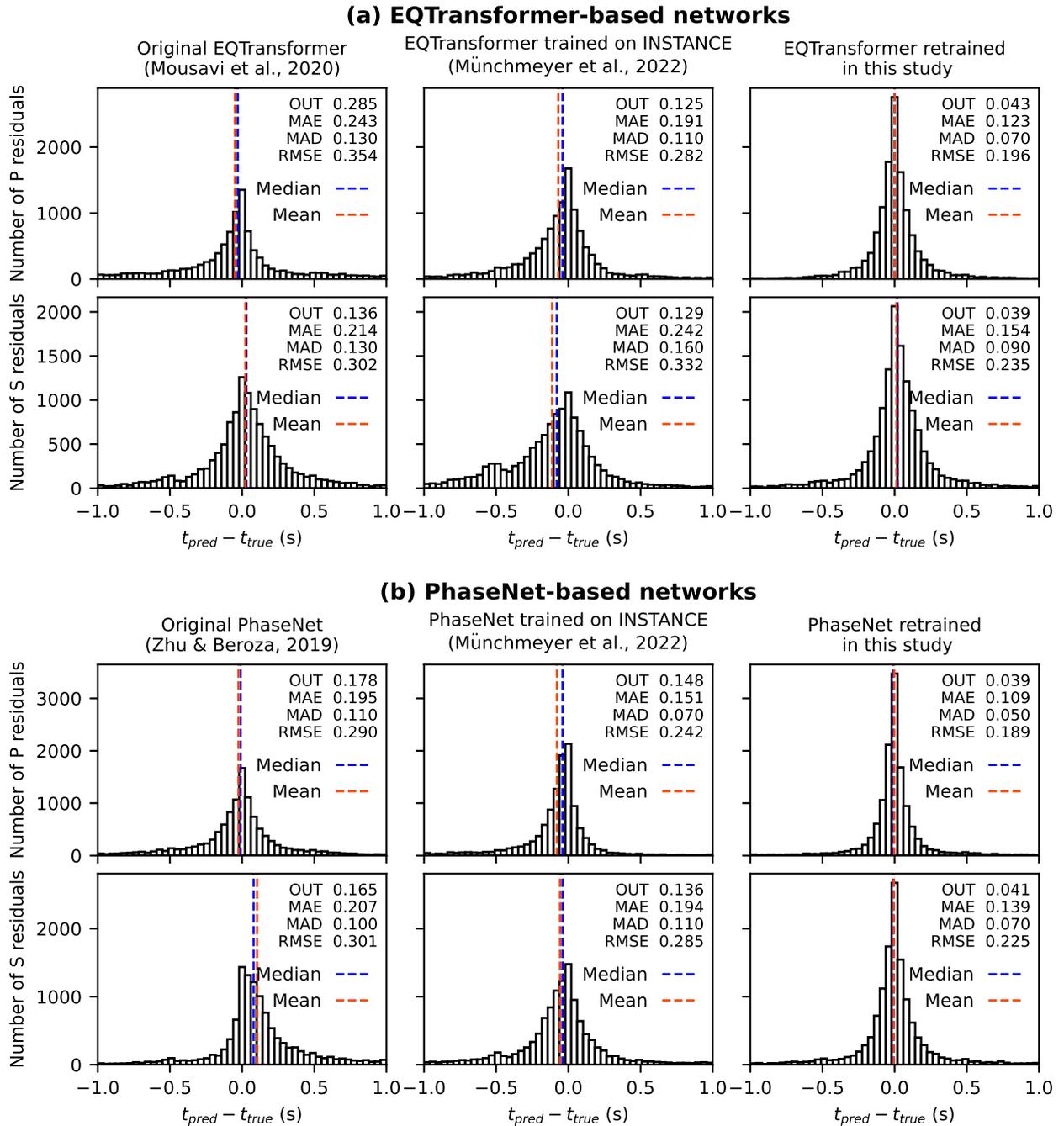


Figure S23. Similar to Figure S22 but for the LP test set. The difference from Figure S20 is due to the different ways of pre-processing and post-processing as explained in the caption of Figure S22.

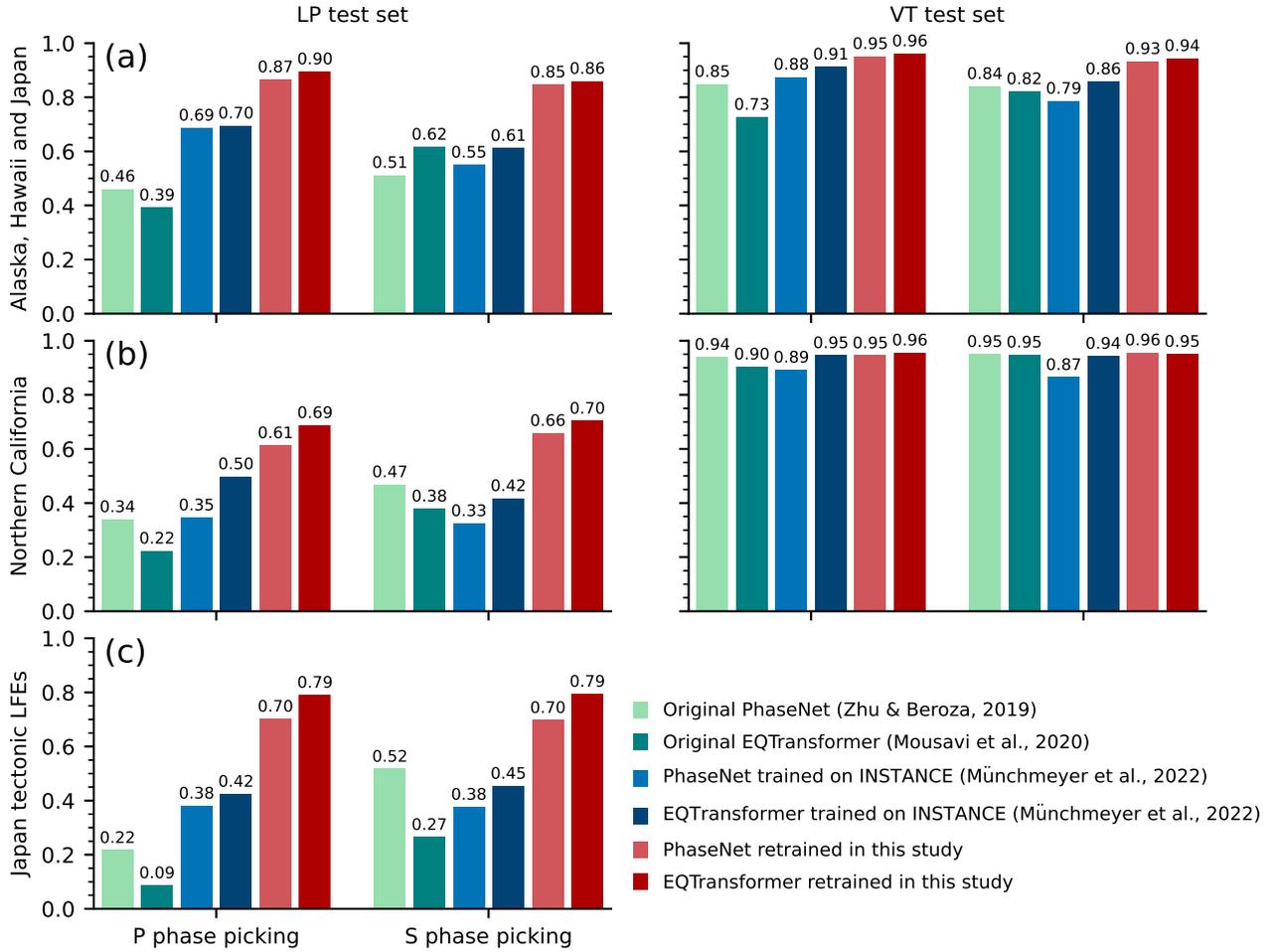


Figure S24. Recalls of different models evaluated on the test waveforms from (a) the same regions as the training data, (b) northern California from where no training data are used and (c) tectonic LP earthquakes in Japan which are generally considered different from volcanic long-period earthquakes in terms of source processes.

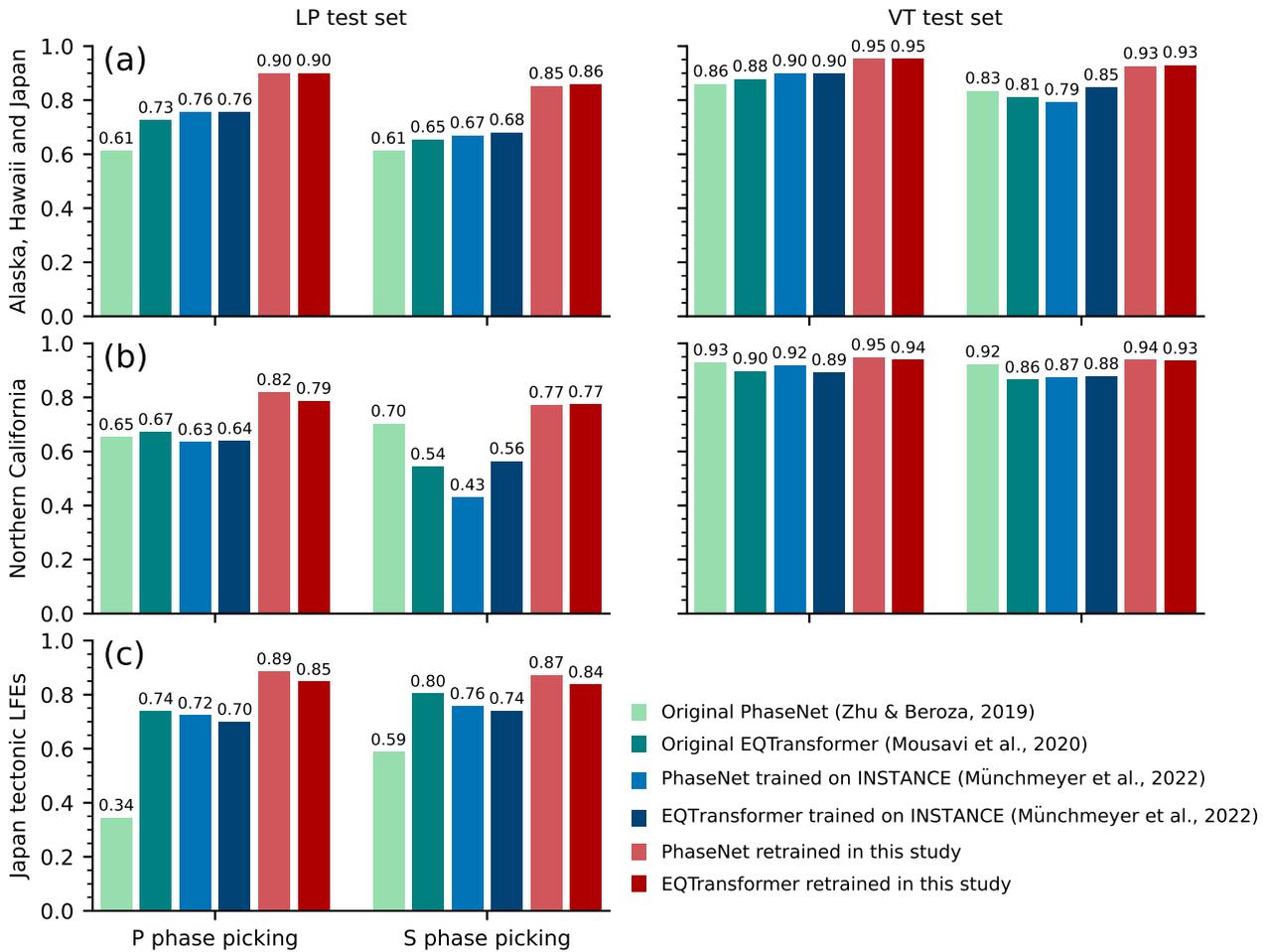


Figure S25. Precisions of different models evaluated on the test waveforms from (a) the same regions as the training data, (b) northern California from where no training data are used and (c) tectonic LP earthquakes in Japan which are generally considered different from volcanic long-period earthquakes in terms of source processes.

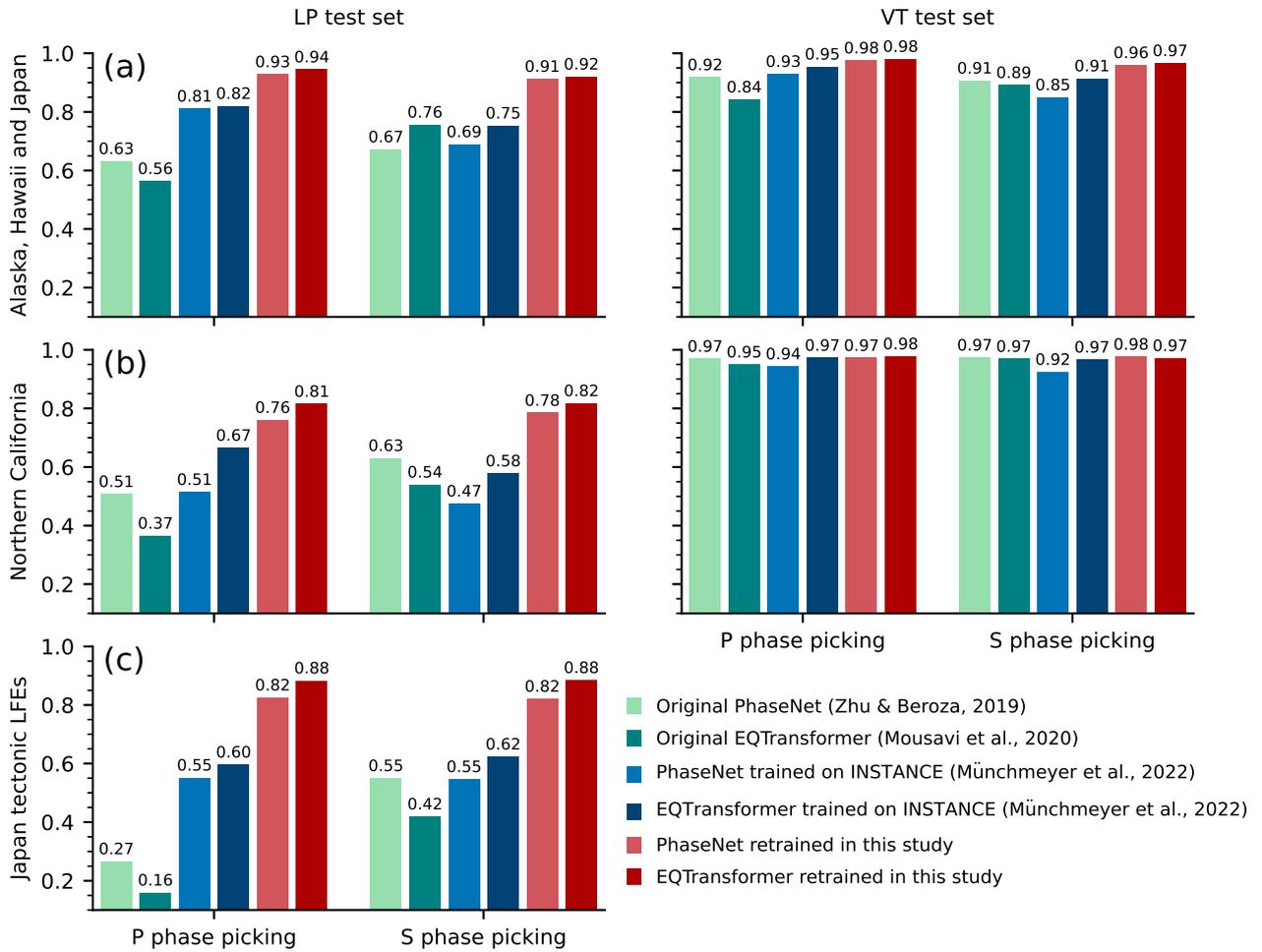


Figure S26. F1 scores of different models calculated using the definition of FP/TN/FN/TN based on waveform traces rather than sampling points. Each row shows the performance for test data from different regions: (a) the same regions as the training data, (b) northern California from where no training data are used, (c) tectonic LP earthquakes in Japan which are generally considered different from volcanic long-period earthquakes in terms of source processes. The precision and recall are given in Figure S30-S31 in the supplement.

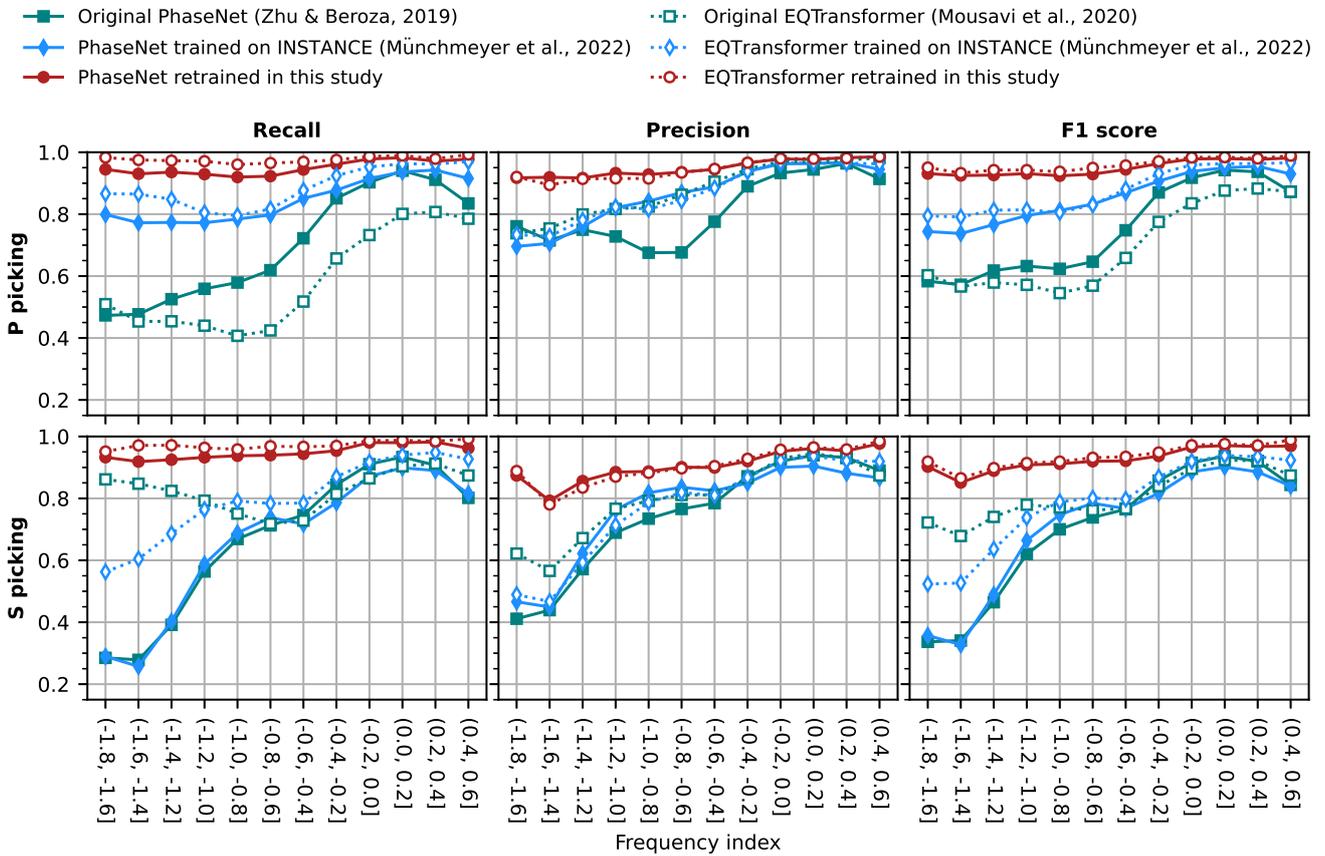


Figure S27. Model performance on subsets of testing waveforms with different frequency index values. Different from Figure 3 in the main paper, the performance in this figure is calculated using the definition of FP/TN/FN/TN based on waveform traces rather than sampling points.

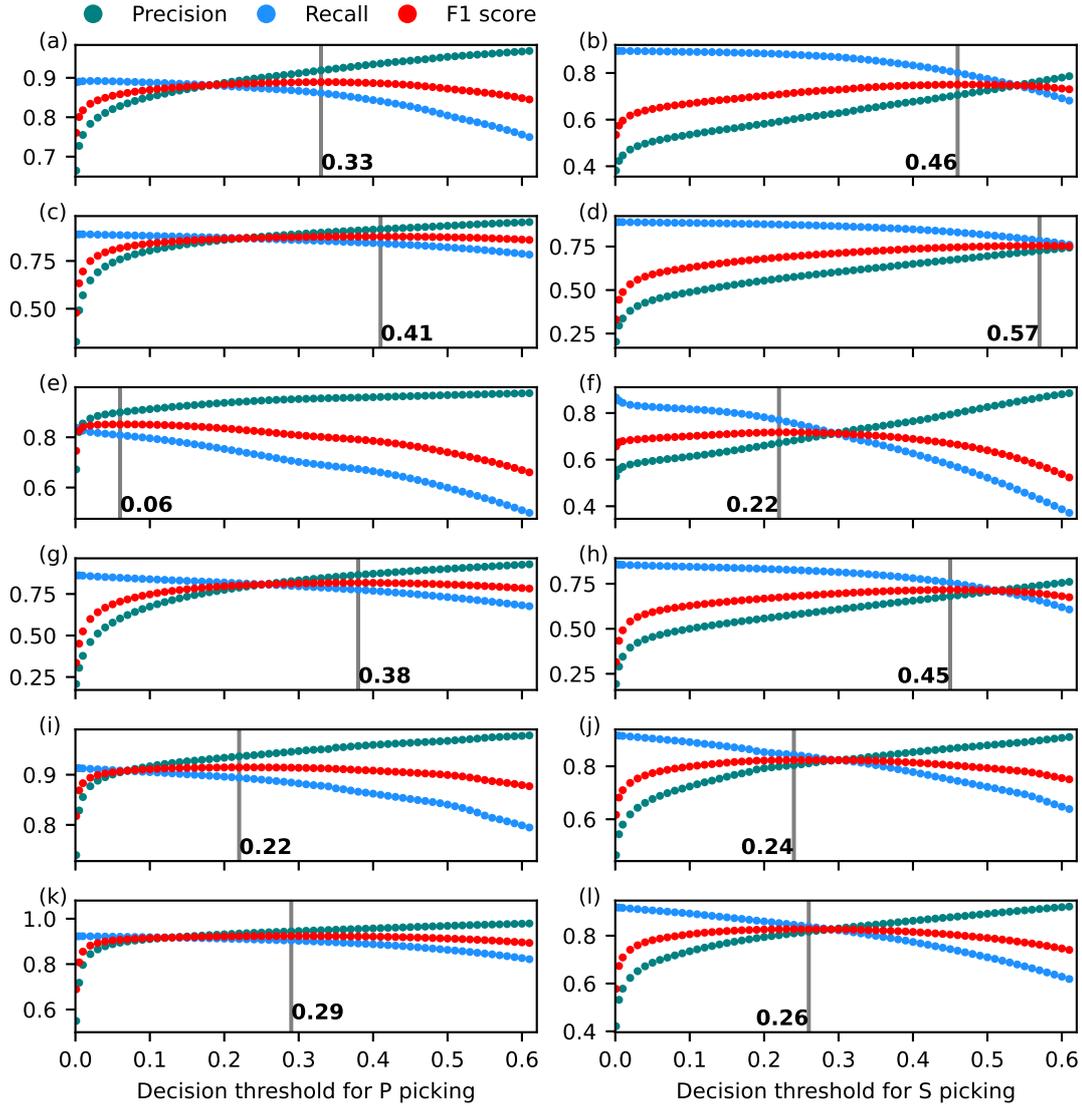


Figure S28. Model performance on the validation set of the INSTANCE data set for the EQTransformer retrained in this study (a-b), PhaseNet retrained in this study (c-d), original EQTransformer (e-f), original PhaseNet (g-h), EQTransformer trained on INSTANCE (i-j), PhaseNet trained on INSTANCE (k-l). The optimal decision thresholds (vertical gray lines) are selected to maximize the F1 score on the validation set.

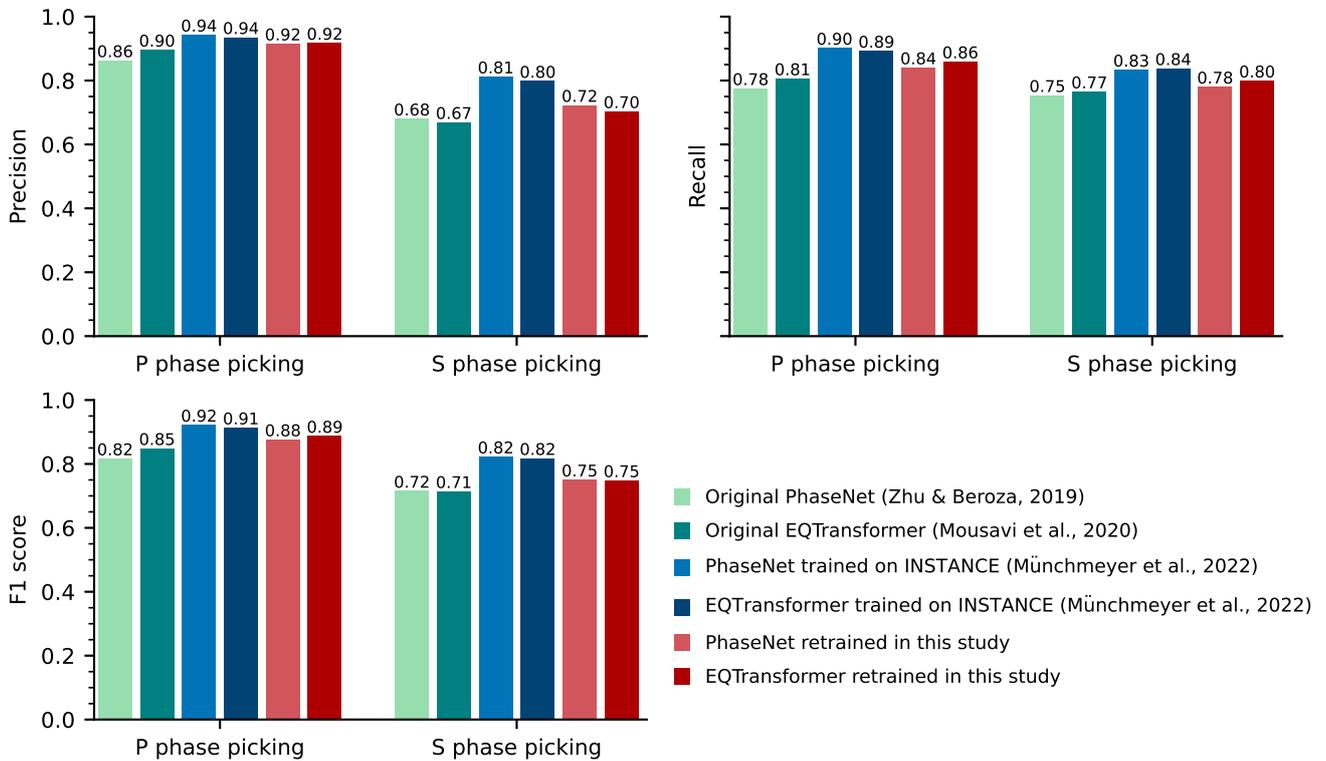


Figure S29. Evaluation of different models on the test set of the INSTANCE data set. The optimal decision thresholds are selected to maximize the F1 score on the validation set of the INSTANCE data set (Figure S28).