

Spatio-Temporal Graph Convolutional Networks for Earthquake Source Characterization

Xitong Zhang¹, Will Reichard-Flynn¹, Miao Zhang², Matthew Hirn³, and Youzuo Lin¹

¹Los Alamos National Laboratory (DOE)

²Dalhousie University

³Michigan State University

November 21, 2022

Abstract

Accurate earthquake location and magnitude estimation play critical roles in seismology. Recent deep learning frameworks have produced encouraging results on various seismological tasks (e.g., earthquake detection, phase picking, seismic classification, and earthquake early warning). Most existing machine learning earthquake location methods utilize waveform information from a single station. However, multiple stations contain more complete information for earthquake source characterization. Inspired by recent successes in applying graph neural networks in graph-structured data, we develop a Spatio-Temporal Graph Convolutional Neural Network (STGCN) for estimating earthquake locations and magnitudes. Our graph neural network leverages geographical and waveform information from multiple stations to construct graphs automatically and dynamically by an adaptive feature integration process. Given input waveforms collected from multiple stations, the neural network constructs different graphs and fuses spatial-temporal consistency effectively from various stations based on graphs' edges. Using a recent graph neural network and a fully convolutional neural network as baselines, we apply STGCN to earthquakes cataloged by Southern California Seismic Network from 2000 to 2019 and induced earthquakes collected in Oklahoma from 2014 to 2015. STGCN yields more accurate earthquake locations than those obtained by the baseline models and performs comparably in terms of depth and magnitude prediction, though the ability to predict depth and magnitude remains weak for all tested models. Our work demonstrates the potential of using graph neural networks and multiple stations for better automatic estimation of earthquake epicenters.

Spatio-Temporal Graph Convolutional Networks for Earthquake Source Characterization

*

Xitong Zhang^{1,2,*}, Will Reichard-Flynn^{1,*}, Miao Zhang³, Matthew Hirn^{2,4,5},
and Youzuo Lin^{1,◇}

¹Geophysics Group, Earth and Environmental Sciences Division, Los Alamos National Laboratory

²Department of Computational Mathematics, Science and Engineering, Michigan State University

³Department of Earth and Environmental Sciences, Dalhousie University

⁴Department of Mathematics, Michigan State University

⁵Center for Quantum Computing, Science and Engineering, Michigan State University

Key Points:

- We design multiple station-based graph deep neural networks for earthquake source characterization.
- Both geographic distance and waveform feature similarity are considered in a graph convolutional neural network for feature combination.
- Our network is capable of automatically selecting and combining relevant seismic stations to characterize earthquake source parameters.

*: X. Zhang and W. Reichard-Flynn equally contribute to this work.

Corresponding author: ◇ Y. Lin, ylin@lanl.gov

Abstract

Accurate earthquake location and magnitude estimation play critical roles in seismology. Recent deep learning frameworks have produced encouraging results on various seismological tasks (e.g., earthquake detection, phase picking, seismic classification, and earthquake early warning). Most existing machine learning earthquake location methods utilize waveform information from a single station. However, multiple stations contain more complete information for earthquake source characterization. Inspired by recent successes in applying graph neural networks in graph-structured data, we develop a Spatio-Temporal Graph Convolutional Neural Network (STGCN) for estimating earthquake locations and magnitudes. Our graph neural network leverages geographical and waveform information from multiple stations to construct graphs automatically and dynamically by an adaptive feature integration process. Given input waveforms collected from multiple stations, the neural network constructs different graphs and fuses spatial-temporal consistency effectively from various stations based on graphs' edges. Using a recent graph neural network and a fully convolutional neural network as baselines, we apply STGCN to earthquakes cataloged by Southern California Seismic Network from 2000 to 2019 and induced earthquakes collected in Oklahoma from 2014 to 2015. STGCN yields more accurate earthquake locations than those obtained by the baseline models and performs comparably in terms of depth and magnitude prediction, though the ability to predict depth and magnitude remains weak for all tested models. Our work demonstrates the potential of using graph neural networks and multiple stations for better automatic estimation of earthquake epicenters.

Plain Language Summary

Machine learning-based approaches have recently become prevalent in seismological tasks such as earthquake source characterization, which is the interest of this paper. The location and magnitude of an earthquake can be best determined by relating the motion recorded at multiple stations in a network. Therefore, it would be beneficial to

45 combine the waveforms from multiple seismic stations for source characterization. Be-
46 cause of the irregular spatial distribution of seismic stations, graph convolutional neu-
47 ral networks (a deep learning architecture which handles graph-structured data) have
48 great potential in combining both spatial and temporal information from different seis-
49 mic stations. In this work, waveforms recorded at multiple stations are passed through
50 neural networks with connective links based on the similarity of waveform features and
51 geographic locations. The model is tested on two datasets and compared to two pub-
52 lished baselines (graph convolutional neural network and fully convolutional network).
53 Compared with the baselines, STGCN achieves improved accuracy for epicenter estima-
54 tion and comparable accuracy for depth and magnitude estimation.

55 **1 Introduction**

56 Earthquake source characterization plays a fundamental role in various seismic stud-
57 ies, including earthquake early-warning, hazard assessment, subsurface energy exploration,
58 etc. (L. Li et al., 2020). Characterization of an earthquake source can be posed as a clas-
59 sical inverse problem. Its purpose is to infer the source information (location, magnitude,
60 etc) from seismic recordings. Various approaches have been developed to characterize
61 earthquake sources, the most well-established being travelttime-based inversion (Z. Zhang
62 et al., 2017; Z. Li & van der Baan, 2016; Lin et al., 2015; H. Zhang & Thurber, 2003)
63 and waveform-based inversion (Beskardes et al., 2018; Zhebel & Eisner, 2015; Pesicek
64 et al., 2014; Gajewski et al., 2007). Travelttime-based methods implement a multi-step
65 process, in which the arrival times of P and S waves are determined through phase de-
66 tection and then associated to specific earthquakes; earthquake locations are estimated
67 as an inversion process given arrival times, station locations, and a velocity model. Mag-
68 nitudes are calculated based on waveform amplitudes. Though travelttime-based meth-
69 ods are commonly used in seismic applications, they are susceptible to noise-related er-
70 rors, particularly when estimating low-magnitude events, and fail to utilize abundant phase
71 and amplitude information in the complete waveform. In contrast, waveform-based in-

72 version integrates all phase and amplitude information recorded in seismographs, result-
73 ing in high quality source characterization, however, which is computationally expen-
74 sive. Both methods require domain expertise to properly tune parameters in the inver-
75 sion process. Deep learning for source characterization provides a data-driven alterna-
76 tive, where integrated location and magnitude predictions extract full-waveform features
77 with less computational expense than waveform inversion.

78 Advances in algorithms and computing, and the availability of large, high-quality
79 datasets have allowed machine learning techniques to attain spectacular success in seis-
80 mological applications (Kong et al., 2019; Bergen et al., 2019) including phase picking
81 (Zhu & Beroza, 2019), seismic discrimination (Z. Li et al., 2018), waveform denoising (Zhu
82 et al., 2019), phase association (Ross et al., 2019), earthquake location (Perol et al., 2018),
83 as well as magnitude estimation (Mousavi & Beroza, 2020b). Although machine learn-
84 ing has long been applied to seismic event detection (J. Wang & Teng, 1995; Tiira, 1999),
85 the first work to leverage recent advances in deep learning was developed by Perol et al.
86 (2018), where convolutional neural networks (CNN's) were trained to detect earthquakes
87 from single station recordings and predict the source locations from among six regions.
88 Though successful in establishing foundational research in machine learning for earth-
89 quake location, the CNN model is restricted to waveforms from a single seismic station
90 and can only classify earthquakes into broad geographic groups without providing spe-
91 cific location information. Since then, more advanced single-station approaches have been
92 developed to improve location accuracy. Mousavi and Beroza (2020a) build Bayesian neu-
93 ral networks to learn epicenter distance, P-wave travel time, and associated uncertainty
94 from single-station data.

95 Recently, multi-station based machine learning methods have shown promising re-
96 sults. For instance, Kriegerowski et al. (2019) develop a CNN structure that combines
97 three-component waveforms from multiple stations to predict hypocenter locations, re-
98 sulting in more accurate source parameters than single station methods. X. Zhang et al.

99 (2020) developed an end-to-end fully convolutional network (FCN) to predict the prob-
100 ability distribution of earthquake location directly from input data recorded at multi-
101 ple stations, which was extended to determine earthquake locations and magnitudes from
102 continuous waveforms for earthquake early warning (X. Zhang et al., 2021). Shen and
103 Shen (2021) also adopt a CNN framework, extracting the location, magnitude, and ori-
104 gin time from continuous waveforms collected across a seismic network.

105 Though multiple-station approaches improve upon single-station methods, the use
106 of standard convolutional layers is limited in several ways: (1) CNN’s are designed to
107 function on evenly-spaced grids (i.e. photographs) where information is exclusively shared
108 between adjacent cells, and (2) CNN’s require the input of station locations to be static
109 (i.e. recordings from station 01 must always be found at position 01 of the input file) in
110 order to learn positional mapping. These assumptions are inappropriate for seismic net-
111 works, which are not regularly-spaced and may record information related to non-adjacent
112 stations. Additionally, station outages, the addition/removal of stations to seismic net-
113 works, and the ability to select a localized array for the detection of small-magnitude events
114 makes dynamic station input highly desirable for source characterization.

115 Münchmeyer et al. (2020) developed an attention-based transformer model for earth-
116 quake early warning, which was extended to predict hypocenters and magnitudes of events
117 in Münchmeyer et al. (2021). While this model is successful in implementing a multi-
118 station approach that allows for dynamic inputs, high computational complexity restricts
119 inputs to a relatively small number of stations. Another method for implementing flex-
120 ible, multi-station input that avoids high complexity for large networks is through graph
121 convolution. This method is implemented by van den Ende and Ampuero (2020), who
122 develop a multi-station source characterization model. This model regards features as
123 nodes on an edgeless graph, implementing single-station convolution and global pooling.
124 However, global pooling may not sufficiently extract all useful information from multi-
125 ple seismic stations, as the pooling layer is ideally applied after global features are ob-

126 tained by feature fusion along the spatial dimension. Yano et al. (2021) introduce a multi-
127 station technique in which edges are manually constructed. While this technique allows
128 for more meaningful features to be constructed than in global pooling, manually-selected
129 edges require station inputs to be fixed during training and implementation, introduc-
130 ing the same limitation inherent to CNN's. Similarly, McBrearty and Beroza (2022) pro-
131 poses a GCN framework using multiple pre-defined graphs constructed on both labels
132 and station locations. The model requires the arrival time and is evaluated by a synthetic
133 dataset.

134 In this study, to harness the full functionality of Graph Convolutional Neural Net-
135 works (GCN's) while maintaining flexibility in the location and number of seismic sta-
136 tions, we design a data-driven framework, spatio-temporal graph convolutional neural
137 network (STGCN), that creates edges automatically to combine waveform features and
138 spatial information. In order to evaluate the performance of our approach, we compare
139 STGCN to two baselines: the GCN model designed by van den Ende and Ampuero (2020)
140 and the Fully Convolutional Network (FCN) designed by X. Zhang et al. (2020). We ap-
141 ply all three models to the two datasets upon which the baselines were originally tested
142 and trained: (1) regional $2.5 < M < 6$ earthquakes recorded by 185 seismic stations
143 in Southern California from 2000 to 2019 (van den Ende & Ampuero, 2020), and (2) lo-
144 cal $0 < M < 4$ earthquakes recorded by 30 seismic stations in Oklahoma from 2014
145 to 2015 (X. Zhang et al., 2020). Next, model stability is evaluated with different hyper-
146 parameters. Finally, we examine the transferrability of STGCN to seismic networks out-
147 side of the training domain.

148 The layout of this article is as follows. In Section 2, we describe the fundamentals
149 of graph-based CNN models and STGCN. In Section 3, we introduce the field data, train-
150 ing procedures, and experimental results. In Section 4, we discuss the mechanisms which
151 enhance and inhibit the performance of STGCN in the context of previous work. Finally,
152 in Section 5, we present concluding remarks and discuss future work.

2 Methodology

In this section, we describe our framework and the major components of our STGCN. A graph is constructed by a set of nodes and edges. Our proposed framework constructs input-dependent graphs automatically, in which a node represents a seismic station and the edge connecting two nodes denotes that extracted features from these two nodes will be combined during convolution. The input to the network is collection of three-channel waveforms from each seismic station, along with the latitude and longitude of the recording stations. The output is the earthquake magnitude and location denoted by latitude, longitude and depth.

2.1 Overview



Figure 1: The overview of STGCN. There are three major components in STGCN: (1) Waveform feature extraction for obtaining time domain feature from each station independently. (2) Spatial feature fusion for time domain feature integration from different stations based on their geographic locations and extracted feature similarity. (3) Earthquake location and magnitude prediction given spatial features from the previous step.

Graph convolutional neural networks (GCN's) are designed to handle graphical data, or data that can be represented by vertices connected by edges. In GCN's, convolution and pooling operates along connecting edges. In CNN's, on the other hand, convolution and pooling operates on regions closest together on a Euclidean grid, meaning that in-

167 put order directly impacts information-sharing and featurization. This is not the case
168 for GCN's, in which edges are not restricted to Euclidean grids but may instead be con-
169 structed by any criteria. Two major advantages of GCN architectures are that they do
170 not require a fixed input order, and can handle graphs with different sets of vertices. These
171 properties of GCN's fit well in seismic data analysis with inputs from multiple stations.
172 It is common for stations in a seismic network to be added, removed, or repositioned,
173 or for the recording quality of individual stations to fluctuate over time due to opera-
174 tion and/or equipment issues. It is therefore beneficial to dynamically select relevant seis-
175 mic stations for source characterization. We therefore propose a dynamic GCN frame-
176 work as the basis for STGCN.

177 Inspired by Y. Wang et al. (2019), our graph convolutions follow the design of Edge-
178 Conv layers to automatically generate edges between nodes. Instead of manually con-
179 structing fixed edges or implementing an edgeless graph, our framework learns to com-
180 bine useful information from multiple stations implicitly during the training process. Our
181 framework consists of three major components as shown in Figure 1:

- 182 • Waveform feature extraction: We first extract time-domain features from the wave-
183 form recorded at each seismic station using a CNN-based encoder. The three-channel
184 seismic recordings are reduced to a low dimensional representation.
- 185 • Spatial feature fusion: We then represent the seismic station network as a graph,
186 in which each node (i.e. station) is connected to other nodes by automatically gen-
187 erated edges. Through iterative steps of edge generation and convolution, the per-
188 ceptive field is gradually enlarged. The model integrates and fuses features from
189 different stations to obtain a high-order view of the recorded wavefield over the
190 seismic network. The graph convolutional architecture considers both geographic
191 locations and waveform feature similarity among multiple seismic stations.

- Prediction: The last component is the prediction module. A fully-connected neural network outputs four normalized scalars corresponding to latitude, longitude, depth and magnitude based on features learned from the previous steps.

2.2 Graph Convolutional Layers

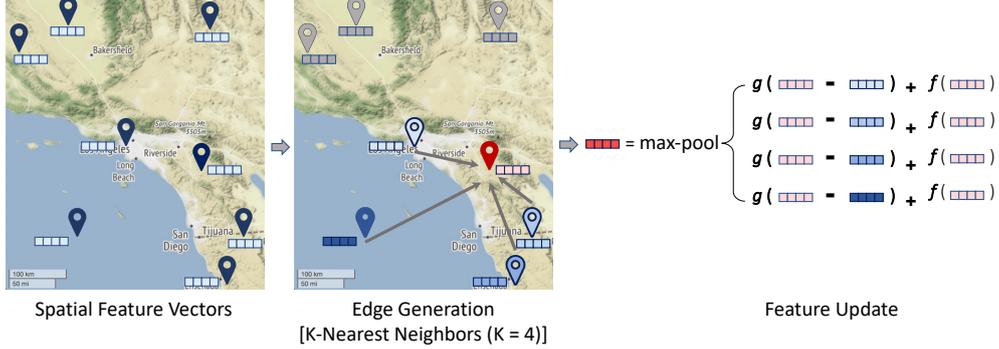


Figure 2: The overview of a graph convolutional layer. Each graph convolutional layer consists of two parts: (1) Edge generation among different stations. (2) Feature updating for each station based on the generated links. In the figure, the feature of the red station is updated based on four nearby blue stations. $g(\cdot)$ and $f(\cdot)$ represent two learnable networks.

The spatial feature fusion process is the most important component and consists of four graph convolution layers. The goal of each graph convolution layer is to enlarge the perceptive field by combining the extracted feature of each seismic stations and auto-selected neighbor stations. As shown in the Figure 2, each graph convolution layer can be broken down into two steps:

- Edge generation: Each station node is connected to several other station nodes which show maximum similarity to the node. Similarity measurements are based on two criteria:

204 1. *Geographic distance*: The geographic distance is the intuitive choice, since ad-
 205 jacent stations tend to record related signals due to similar wave paths. Addi-
 206 tionally, events are more likely to be mutually recorded by stations in close prox-
 207 imity, especially in the case of small-magnitude events.

208 2. *Feature similarity*: As the same earthquake event can be recorded by distant
 209 stations in a large area, waveform similarity provides a complimentary perspec-
 210 tive to geographic distance. We compare l_2 distance of features from station i
 211 and j directly by $\|x_i - x_j\|_2^2$, and thus we can combine two waveform features
 212 from two stations further away, where x_i and x_j are the extracted feature vec-
 213 tors.

214 In edge generation, we link every station with its K-nearest neighbors based on
 215 their similarity, where K is a tunable hyperparameter. In our framework, both ge-
 216 ographic proximity and waveform feature similarity are considered. In practice,
 217 the similarity between waveforms can also be affected by other factors, such as wave
 218 path and signal to noise ratio. By training with a large amount of samples with
 219 different sets of seismic stations with distinct spatial distributions, the network
 220 will learn to embed these implicit and complex factors to low dimensional features
 221 automatically, in order to minimize the misfit between labels and predictions.

- Feature update: Given the edges, we update the features of each stations by

$$\tilde{x}_i = \max_{j \in \mathcal{N}_{\text{distance}}(i)} g(x_i - x_j) + f(x_i) + \max_{j' \in \mathcal{N}_{\text{feature}}(i)} g(x_i - x_{j'}) + f(x_i), \quad (1)$$

222 where the max operation refers to the element-wise max-pooling. x_i , x_j and $x_{j'}$
 223 are features of station i , j and j' , respectively. j is a neighbor of i based on ge-
 224 ographic distance and j' is a neighbor of i by measuring feature similarity from
 225 the previous edge generation step. $g(\cdot)$ and $f(\cdot)$ are two trainable fully connected
 226 neural networks. \tilde{x}_i is the updated feature of station i . Max pooling is conducted
 227 along the constructed edges to combine information from the K-nearest neighbors
 228 of i . The update is asymmetric for station i and j to encourage the update pro-

229 cesses of i and j to be different, as it is possible that only one of the stations records
 230 the event.

231 **2.3 Architecture**

232 A graphical illustration of the architecture is presented in Figure 3. Time domain
 233 waveform features are extracted from each station independently using an encoder with
 234 eleven convolutional layers. The extracted features are used in spatial feature fusion, in
 235 which time domain features are concatenated to station locations before each graph con-
 236 volution. Our STGCN uses four groups of graph convolutional layers to obtain spatially
 237 hierarchical features. Two graphs are generated within each group: one in which edges
 238 are generated based on geographic distance, and one in which edges are generated based
 239 on waveform feature similarity. After convolution, the features obtained from both graphs
 240 are summed together prior to max pooling. For graphs in which geographic distance dic-
 241 tates edges, two scalars containing station coordinates are concatenated to each updated
 242 feature before each convolution. After all four groups of convolutions, the features from
 243 each group are concatenated together as a hierarchical representation for final source char-
 244 acterization regression.

After all feature outputs are concatenated, the features are individually processed
 with a final CNN layer. The output is then regressed to scalar predictions of latitude,
 longitude, depth, and magnitude using a fully-connected neural network. The objective
 function is

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N \frac{1}{4} \|y_i - \hat{y}_i\| \quad (2)$$

245 where \hat{y}_i and y_i are the prediction and ground truth values of i th sample, respectively,
 246 represented as vectors of latitude, longitude, depth and magnitude.

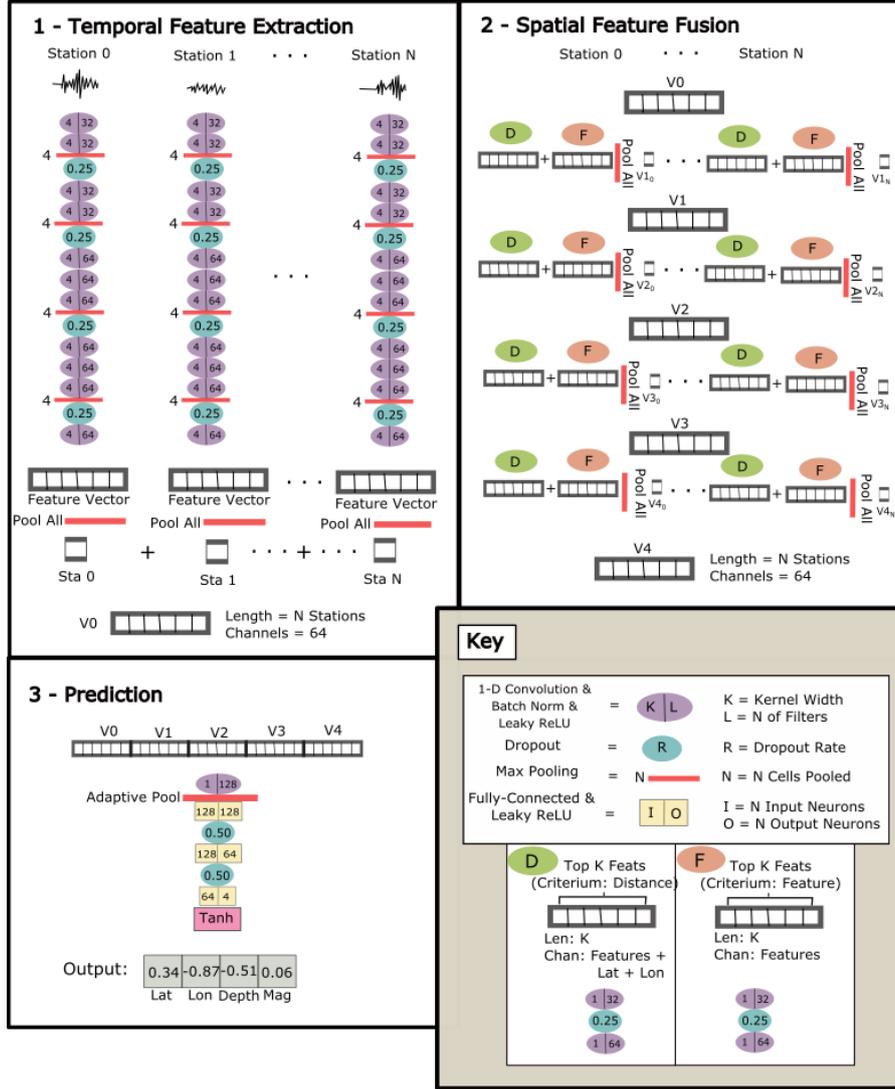


Figure 3: Overview of STGCN. STGCN includes three components, following the framework outlined in Figure 1.

247 3 Experiments and Results

248 In this section, the data, experiment settings, and results are discussed. We eval-
249 uate STGCN with three major experiments: (1) performance on two datasets compared
250 to GCN and CNN baselines, (2) stability analysis of STGCN with various settings, and
251 (3) the transferrability of STGCN to regions outside of the training domain.

3.1 Data Description

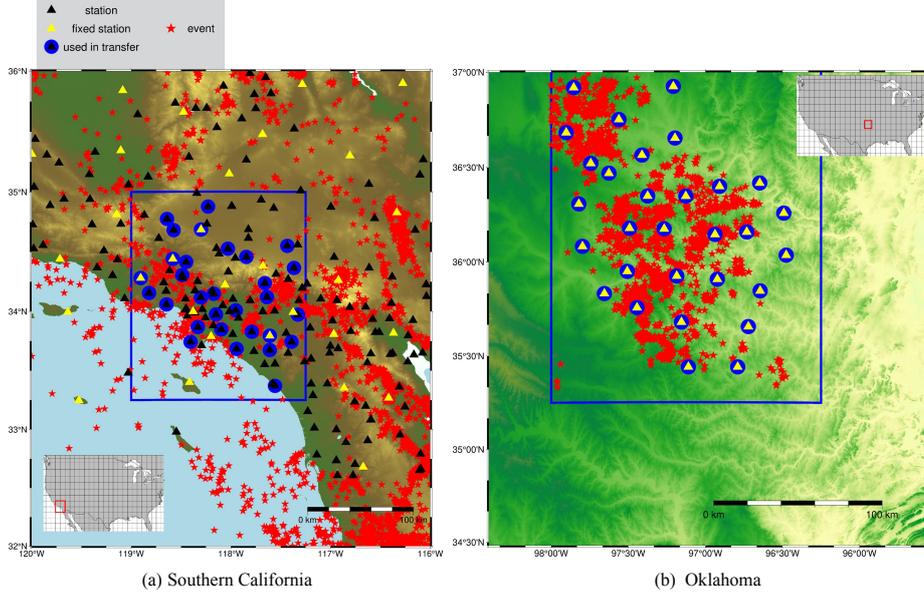


Figure 4: Maps of the two target regions used in this study: (a) Southern California and (b) Oklahoma. The distribution of all seismic stations (black triangles) and earthquakes (red stars) are shown. The areas used selected the transferability study are contained within the blue squares. In the map of Southern California, the 30 stations selected for fixed input testing are yellow triangles, and the 30 stations selected for the transferability study are surrounded by a blue circle.

Consistent with target regions in the GCN baseline (van den Ende & Ampuero, 2020) and the FCN baseline (X. Zhang et al., 2020), we correspondingly collected earthquake datasets from Southern California and Oklahoma. The former data including station inventory, earthquake catalogue and waveforms are downloaded from the Southern California Seismic Network (SCSN) and Southern California Earthquake Data Center (SCEDC) (Hutton et al., 2010) from January 2000 to June 2019 and accessed using ObsPy (Beyreuther et al., 2010). STGCN makes predictions by outputting values between -1 and 1. Thus we constrain our the labels of events to fit within a normalized range. We limit stations and

261 events to a geographic subset from 32° to 36° latitude, and from -120° to -116° lon-
 262 gitude (van den Ende & Ampuero, 2020). We select events from a depth range of 0-30
 263 km and a magnitude range of $2.5 < M < 6$. The final dataset contains 2,209 events
 264 recorded by 185 broadband seismic stations. On average, 48 seismic stations are func-
 265 tional for all events. The maximum number of functional seismic stations we can down-
 266 load raw waveforms from is 142. The spatial distribution of events and stations is illus-
 267 trated in Figure 4. After removing the instrument response, the signals are bandpass fil-
 268 tered from 1–8 Hz. In the second target region, we collect induced earthquake dataset
 269 in Oklahoma from March 2014 to July 2015 (Nanometrics Seismological Instruments,
 270 2013). We limit the dataset to events between 34.482° to 37° latitude, and from -98.405°
 271 to -95.527° longitude with depths from 0-12 km (X. Zhang et al., 2020). Magnitudes
 272 range from $1.5 < M < 4$. The instrument response is removed, and waveforms are band-
 273 pass filtered from 1 – 8 Hz. The final dataset contains 3,456 events recorded from 30
 274 stations.

275 An arbitrary scaling factor of $1e7$ is multiplied across both datasets to raise the ex-
 276 tremely small amplitudes to an acceptable range without eliminating magnitude infor-
 277 mation. Each recording contains 200 sec of seismic displacement collected by three or-
 278 thogonal channels, which is interpolated into 4,096 evenly spaced samples, resulting in
 279 a sampling rate of approximately 20 Hz. We use a sliding window to handle the uncer-
 280 tainty of the arrival time that would occur in practical use by cropping shorter time seg-
 281 ments from longer raw waveforms at different positions in time. Thus, the actual arrival
 282 signal can locate at different time steps and the model will learn to extract proper rep-
 283 resentation from raw seismic waveforms that have different arrival times during train-
 284 ing. In the end, we use a sliding window with a length of 100 sec and a stride of 5 sec
 285 to create ten 100 sec samples from each 200 sec recording. Each sample is associated with
 286 a label containing latitude, longitude, depth and magnitude values normalized from -1
 287 to 1.

288 One advantage of our GCN over CNN’s or GCN’s with fixed edges is its ability to
 289 make predictions using dynamic inputs (i.e., the selected stations and their order in the
 290 input file are not necessarily the same for each sample). To demonstrate this ability, we
 291 perform tests with STGCN and the GCN baseline using Southern California data with
 292 dynamic inputs, in which functioning stations are randomly selected for each event. How-
 293 ever, to make a fair comparison between STGCN and the FCN baseline, the same sta-
 294 tions must occupy the same position in each input. Using the Southern California and
 295 Oklahoma datasets, we train STGCN as well as both baselines on thirty fixed stations
 296 to compare the performance of all methods. The GCN models can be adaptively trained
 297 to make predictions given any number of input stations. If the number of functioning
 298 stations is less than the target number of stations for any given event, the input is padded
 299 with zeroed channels and the coordinates of the missing stations are set to $(-1, -1)$. For
 300 the two datasets, events are omitted where < 25 stations are functioning. In the South-
 301 ern California dataset where phase reports from SCEDC are available, only events with
 302 > 5 stations recording available P and/or S phases are kept, considering the sparse cov-
 303 erage of the stations in a large region. Overall, each event in the Southern California dataset
 304 is detected by an average of 31 stations.

3.2 Training Procedure

306 In the experiments, we use AdamW as the optimizer with a learning rate of $3e-4$.
 307 The l_2 regularization term λ is $1e-4$. Models are trained for 400 epochs with early stop-
 308 ping after 50 epochs without validation error improvement, from which we select the model
 309 with the best validation performance. We use a 20-80 split to divide each dataset into
 310 testing and training data, and reserve 20% of the training data for validation. The datasets
 311 are not randomly shuffled, but rather separated by time in which training data precedes
 312 testing data. This approach avoids potential information leakage (Kaufman et al., 2012)
 313 which might occur from spatially and temporally localized swarms. This method of split-
 314 ting data also better simulates a real-use case, in which historic earthquakes would be

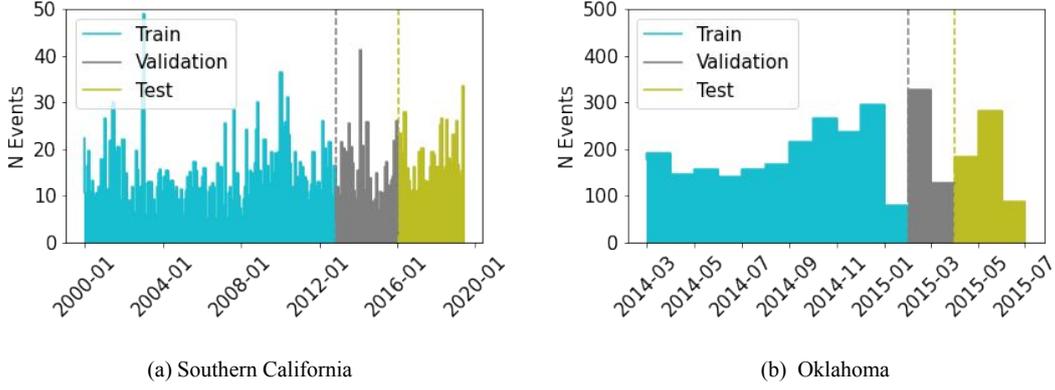


Figure 5: The monthly earthquake frequency distribution for (a) Southern California and (b) Oklahoma. The temporal boundaries between the training, validation, and testing data are indicated by color.

315 used to train a model to detect more recent events on a network where station config-
 316 uration and seismic characteristics may evolve over time. Figure 5 shows the monthly
 317 event frequency distribution in the training and testing dataset.

318 When testing transferability, models are tuned using a learning rate of $3e-5$ for
 319 2,000 epochs with an early stopping cutoff of 100 epochs without validation improvement.
 320 All weights in the model were permitted to retrain.

321 3.3 Performance Comparison

To evaluate our developed framework, we compare the testing mean absolute error (MAE) of our proposed model (referenced as STGCN) with the baseline model by van den Ende and Ampuero (2020) (referenced as GCN) when applied to 100 randomly-selected stations from the Southern California dataset. MAE is calculated as:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|, \quad (3)$$

322 where \hat{y}_i is the model’s prediction, y_i is the true value, and n is the total number of pre-
 323 dictions. In graph convolution, seven edges ($K=7$) were generated between the each sta-

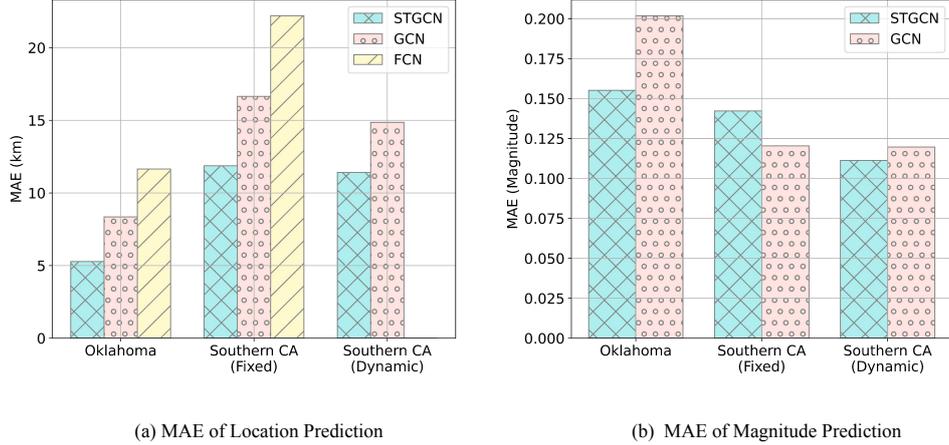


Figure 6: (a) MAE of each tested model where the location error is measured in km. Location error refers to the euclidean distance between the predicted location and the true event location. (b) MAE of the magnitude predictions from the graph convolutional neural networks when applied to the Oklahoma dataset with 30 fixed stations, the Southern California dataset with 30 fixed stations, and the Southern California dataset with 100 dynamically selected stations.

324 tion and the most similar nodes. The number of edges K is a tunable parameter, the im-
 325 pact of which we evaluate in the following section. Both models make predictions nor-
 326 malized between -1 and 1. The values are first reverted from the normalized scalars to
 327 degrees of latitude and longitude, kilometers of depth, and magnitude values. For dis-
 328 tance error calculations, degrees of latitude and longitude are converted to kilometers
 329 using conversions of 110 km/degree and 92 km/degree, respectively. The previous anal-
 330 ysis examines the performance of STGCN when applied to dynamically selected stations
 331 from a large network. To further demonstrate STGCN’s capabilities, we extend our tests
 332 to two different datasets (Oklahoma and Southern California), tested in comparison to
 333 two baselines (a GCN baseline (van den Ende & Ampuero, 2020) and a FCN baseline
 334 (X. Zhang et al., 2020)). As the FCN baseline requires a fixed input consisting of 30 sta-

Latitude	MAE (km)	MSE (10^2 km)	R^2
STGCN	7.788 ± 10.849	1.783 ± 12.697	0.977
GCN	10.201 ± 11.791	2.431 ± 12.438	0.969

Longitude	MAE (km)	MSE (10^2 km)	R^2
STGCN	7.563 ± 9.408	1.457 ± 8.209	0.982
GCN	10.095 ± 12.086	2.480 ± 11.865	0.970

Depth	MAE (km)	MSE (10^2 km)	R^2
STGCN	3.486 ± 2.958	0.209 ± 0.377	0.256
GCN	3.837 ± 3.166	0.247 ± 0.399	0.120

Magnitude	MAE	MSE	R^2
STGCN	0.111 ± 0.115	0.0257 ± 0.0824	0.837
GCN	0.120 ± 0.126	0.0302 ± 0.105	0.807

Table 1: Performance of the STGCN model proposed in this paper and the GCN baseline when applied to the Southern California dataset with dynamic inputs. MAE refers to the mean absolute error (Equation 3) and MSE refers to the mean squared error (Equation 4), where a lower value indicates less error. The R^2 value (Equation 5) is a measure of how strongly variation in the predicted values are related to variation in the ground truth value, where a value close to 1 is indicative of high accuracy.

335 tions, the 30 stations active for the greatest number of events in the Southern Califor-
336 nia dataset were used as the inputs for all samples. The selected stations are highlighted

337 in Figure 4. As the Oklahoma network consists of only 30 stations, all 30 stations were
 338 used. The performance overview is shown in Figure 6, which clearly shows that our pro-
 339 posed model achieve higher localization accuracy than baselines for all datasets. The FCN
 340 baseline doesn't support magnitude prediction, and two GCN-based models achieve com-
 341 parable performance.

342 STGCN makes predictions with an average of 8.3 km less location error, a 49% im-
 343 provement across all tested datasets when compared to the FCN baseline, and has the
 344 ability to predict magnitude as well as location. Across all datasets, STGCN makes pre-
 345 dictions with an average of 3.8 km less location error than the GCN baseline, a 28% im-
 346 provement. While magnitude does not improve for every individual dataset, STGCN shows
 347 an overall improvement in magnitude when all tested datasets are considered.

The detailed evaluation results of Southern California dataset with dynamic seis-
 mic stations are shown in Table 1. In terms of MAE, our GCN model outperforms the
 GCN baseline for all predictions (latitude, longitude, depth, magnitude), with most im-
 provement achieved in latitude and longitude prediction. In addition to MAE, the mean
 squared error and R^2 values are displayed. Mean squared error is calculated as:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (4)$$

where \hat{y}_i is the model's prediction, y_i is the true value, and n is the total number of pre-
 dictions. R^2 is calculated as:

$$R^2 = 1 - \frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}, \quad (5)$$

348 where \hat{y}_i is the model's prediction, y_i is the true value, \bar{y} is the average true value, and
 349 n is the total number of predictions. STGCN demonstrates better performance with both
 350 measures of accuracy and is more consistent (smaller standard deviations in prediction
 351 accuracy). However, both STGCN and the GCN baseline demonstrate exceedingly low
 352 R^2 values for depth prediction. In terms of magnitude, STGCN and GCN perform com-
 353 parably when all measures of accuracy are considered.

Latitude	MAE (km)	MSE (10^2 km)	R^2
STGCN	4.487 ± 9.264	1.060 ± 9.484	0.947
GCN	7.166 ± 12.414	2.055 ± 14.820	0.897
FCN	9.219 ± 16.418	3.545 ± 23.070	0.822

Longitude	MAE (km)	MSE (10^2 km)	R^2
STGCN	4.151 ± 7.035	0.667 ± 5.502	0.937
GCN	5.934 ± 8.144	1.015 ± 5.547	0.904
FCN	9.308 ± 11.883	2.279 ± 8.244	0.785

Depth	MAE (km)	MSE (10^2 km)	R^2
STGCN	1.760 ± 1.473	0.053 ± 0.083	0.026
GCN	1.701 ± 1.423	0.049 ± 0.078	0.090
FCN	1.865 ± 1.546	0.059 ± 0.084	-0.086

Magnitude	MAE	MSE	R^2
STGCN	0.154 ± 0.123	0.0388 ± 0.0668	0.787
GCN	0.195 ± 0.142	0.0582 ± 0.0831	0.681

Table 2: Performance of STGCN, GCN and FCN baselines when applied to the Oklahoma dataset with fixed inputs. MAE refers to the mean absolute error (Equation 3) and MSE refers to the mean squared error (Equation 4), where a lower value indicates less error. The R^2 value (Equation 5) is a measure of how strongly variation in the predicted values are related to variation in the ground truth value, where a value close to 1 is indicative of high accuracy.

Latitude	MAE (km)	MSE (10^2 km)	R^2
STGCN	8.022 ± 9.664	1.577 ± 9.297	0.970
GCN	11.263 ± 11.696	2.637 ± 8.010	0.949
FCN	14.415 ± 21.827	6.842 ± 34.697	0.869

Longitude	MAE (km)	MSE (10^2 km)	R^2
STGCN	7.840 ± 11.645	1.971 ± 19.095	0.972
GCN	11.485 ± 12.199	2.807 ± 10.252	0.960
FCN	16.369 ± 24.872	8.865 ± 47.323	0.874

Depth	MAE (km)	MSE (10^2 km)	R^2
STGCN	3.869 ± 3.380	0.264 ± 0.411	-0.016
GCN	4.264 ± 3.384	0.296 ± 0.403	-0.141
FCN	4.105 ± 3.324	0.279 ± 0.431	-0.074

Magnitude	MAE	MSE	R^2
STGCN	0.142 ± 0.117	0.0340 ± 0.0624	0.796
GCN	0.120 ± 0.118	0.0283 ± 0.0880	0.830

Table 3: Performance of STGCN, GCN and FCN baselines when applied to the Southern California dataset with fixed inputs. MAE refers to the mean absolute error (Equation 3) and MSE refers to the mean squared error (Equation 4), where a lower value indicates less error. The R^2 value (Equation 5) is a measure of how strongly variation in the predicted values are related to variation in the ground truth value, where a value close to 1 is indicative of high accuracy.

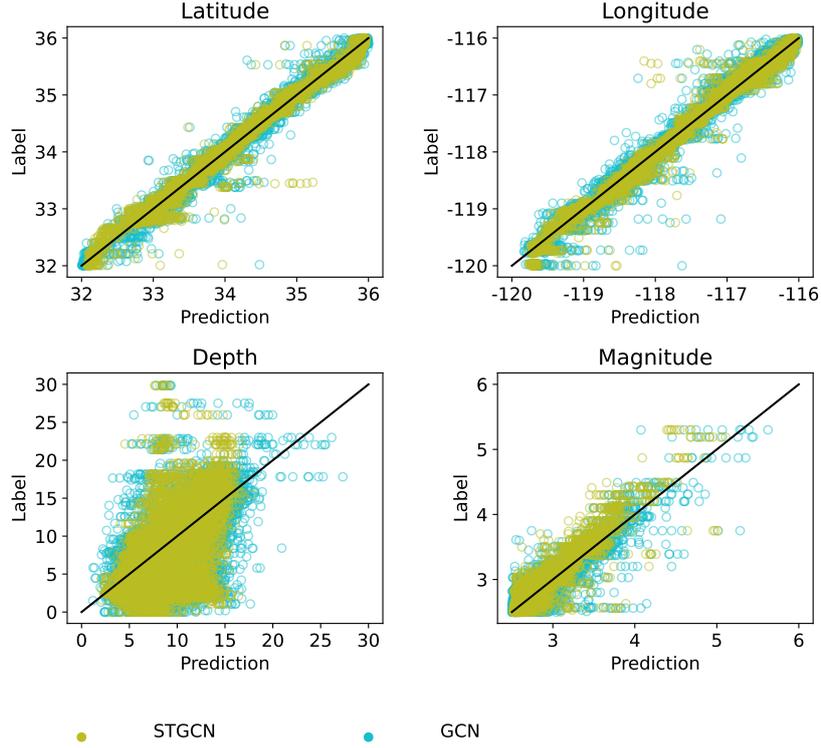


Figure 7: Testing comparison on 100 dynamically selected stations from the Southern California dataset. “STGCN” and “GCN” denote the performance of our framework and the published baseline, respectively. In the scatter plot, each point represents an event, and a position on the diagonal line corresponds to perfect agreement between the predicted value (x-axis) and the true value (y-axis). Latitude and longitude values are displayed in degrees and depth values are displayed in kilometers

354 In both datasets with fixed seismic stations, the proposed model shows significant
 355 improvement over both baselines in terms of location error, with most improvement aris-
 356 ing from latitude and longitude predictions. This improvement is supported by several
 357 performance metrics (Table 2 and Table 3). Overall location error is 5.28 km for the Ok-
 358 lahoma dataset and 11.87 km for the Southern California dataset. The higher loss for
 359 the Southern California dataset may be attributable to the larger size of the region. As
 360 locations in both the smaller and larger regions are normalized to values between -1 and

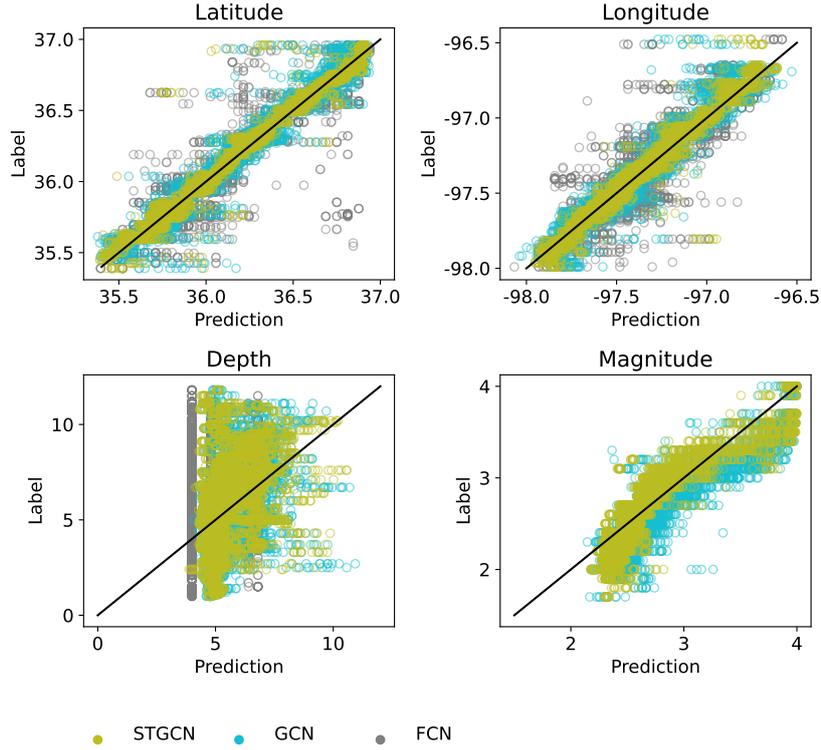


Figure 8: Testing comparison on 30 fixed stations from the Oklahoma dataset.

“STGCN”, “GCN”, and “FCN” denote the performance of our framework, the published GCN baseline, and the published FCN baseline, respectively. In the scatter plot, each point represents an event, and a position on the diagonal line corresponds to perfect agreement between the predicted value (x-axis) and the true value (y-axis). Latitude and longitude values are displayed in degrees and depth values are displayed in kilometers. Magnitude is omitted for the FCN, as this model makes only location predictions

361 1, errors in the initial prediction will result in larger errors when converted to kilome-
 362 ters in larger regions. In addition, larger regions may include a greater range of struc-
 363 tural complexity that may be more challenging for the model to learn.

364 Figure 7, 8 and 9 plot all predictions to give a richer understanding of model ca-
 365 pacity beyond individual quality metrics. Observation of individual predictions makes
 366 it clear that while both models succeed in learning a meaningful mapping to latitude and

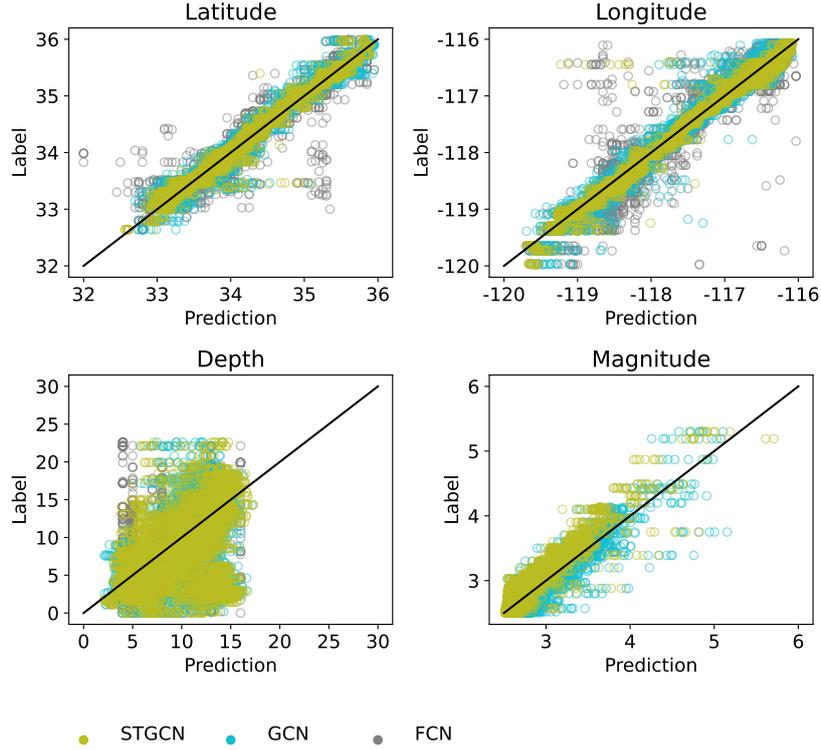


Figure 9: Testing comparison on 30 fixed stations from the Southern California dataset. “SPCGN”, “GCN”, and “FCN” denote the performance of our framework, the published GCN baseline, and the published FCN baseline, respectively. In the scatter plot, each point represents an event, and a position on the diagonal line corresponds to perfect agreement between the predicted value (x-axis) and the true value (y-axis). Latitude and longitude values are displayed in degrees and depth values are displayed in kilometers. Magnitude is omitted for the FCN, as this model makes only location predictions.

367 longitude predictions, depth predictions are highly scattered and are little better than
 368 predictions of the mean.

369 While our proposed model predicts magnitude with less error than the GCN base-
 370 line on the Oklahoma dataset, the model has greater magnitude errors when applied to
 371 the Southern California dataset (Table 1). All models perform extremely poorly when
 372 predicting depth. Therefore STGCN does not improve depth or magnitude prediction,

373 where it remains comparable to the baseline models. However, STGCN substantially im-
 374 proves latitude and longitude predictions, resulting in higher quality location estimations
 375 .

376 **3.4 Stability Analysis**

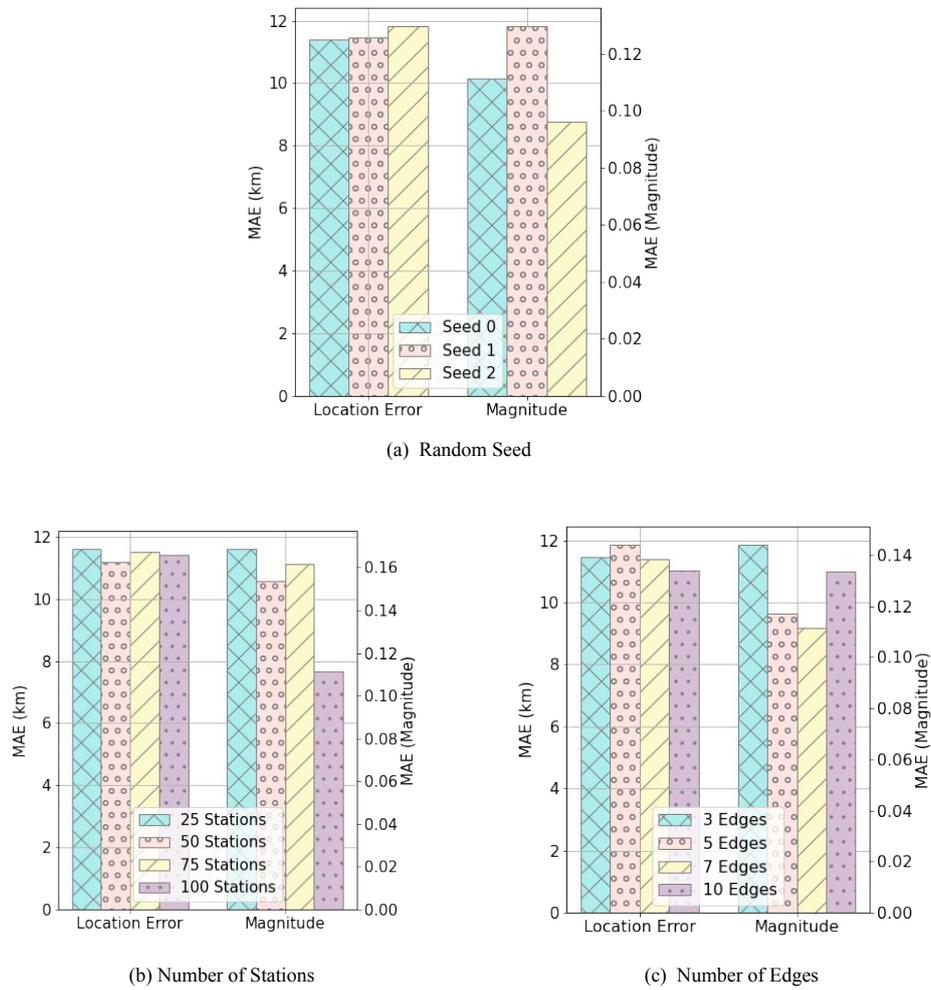


Figure 10: Stability analysis permuting (a) the random seed used to select stations for the model input, (b) the number of stations used for prediction, and (c) the number of edges used to connect nodes during graph convolution.

377 There are three critical hyper-parameters in STGCN: the number of neighbors con-
378 sidered for edge generation, the total amount of observed stations, and the random se-
379 lection of seismic stations when creating datasets. We use the Southern California Dasaset
380 to vary these hyperparameters in order to assess the stability of STGCN. The results of
381 the paramater permutation are shown in Figure 10.

382 For each prediction, a random subset of functional stations were selected. We per-
383 mutate the random seed during the selection of 100 stations, making predictions using 7
384 edges. We find that the random subsets return similar results for all predictions except
385 for magnitude, which shows a higher degree of variation. With the exception of magni-
386 tude, prediction accuracy remains similar when 25, 50, 75, or 100 stations are used. Mag-
387 nitude prediction improves substantially when 100 stations are selected. A similar pat-
388 tern is observed in the edge stability, where the number of generated edges has the great-
389 est influence on magnitude performance. Overall, the model appears to be generally sta-
390 ble, with magnitude demonstrating the greatest sensitivity to hyperparameter tuning.

391 **3.5 Transferability**

392 In many real use cases, a studied network may have a small or nonexistent cata-
393 logue of events with which to train a predictive model. It is therefore useful to test the
394 effectiveness of a pretrained model when applied to events in an unseen region. Figure 11
395 shows the performance of a model trained on the Southern California dataset and tested
396 on the Oklahoma dataset and vice versa when tuned with samples ranging from 0-250.
397 Regardless of the number of training samples, the validation and testing data remained
398 the same for each training and testing. The performance of a tuned model is compared
399 to the performance of a randomly-initialized model trained with the same number of sam-
400 ples to examine whether or not pretraining is beneficial.

401 Two equal-area regions were selected from the Oklahoma and Southern California
402 datasets. From the Southern California dataset, 30 fixed stations were selected which

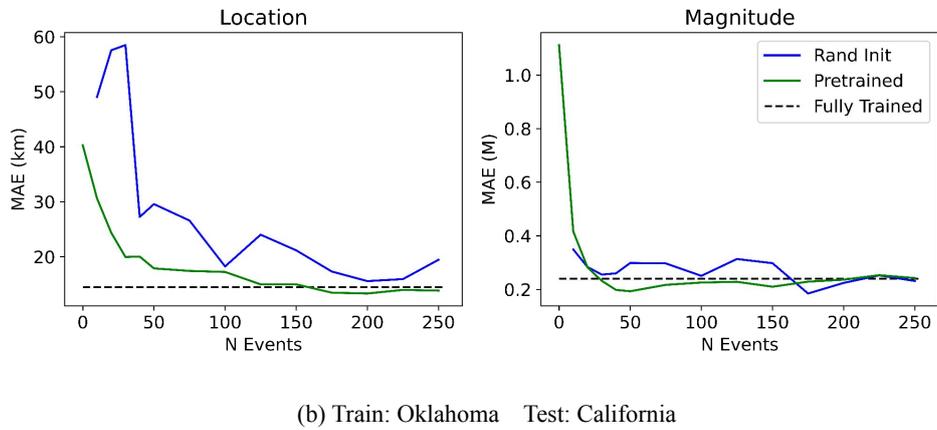
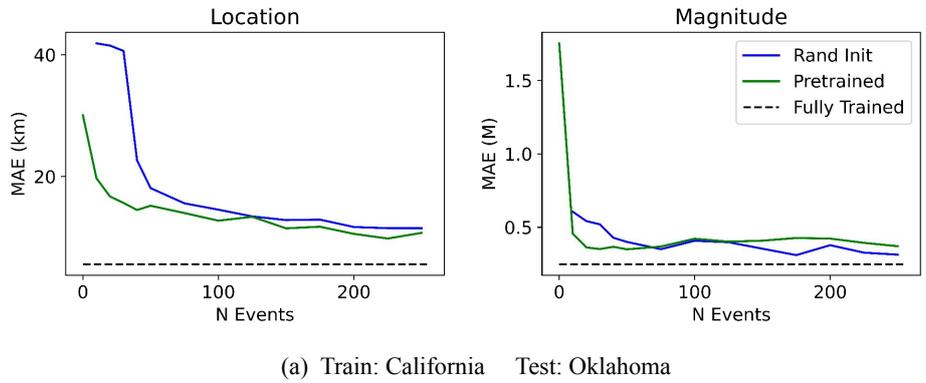


Figure 11: Transferability of (a) model trained on Southern California data and tested on Oklahoma data and (b) model trained on Oklahoma data and tested on Southern California data. The plots show the prediction error of the pretrained model (green) and randomly initialized model (blue) when a range of 0 (no retraining) to 250 events are used for training. The panels to the left show the euclidean location errors between the predicted and true hypocenter measured in km, and the panels to the right show the magnitude errors. The dashed line corresponds to the performance when randomly initialized weights are trained with all available training data from the region.

403 most closely resemble the distribution of the Oklahoma dataset with respect to minimum,
404 maximum, and mean distance between stations. The Oklahoma dataset consists of a much
405 larger training dataset comprising 2,025 events while the Southern California dataset
406 contains only 254 events. Overall, when pretrained models are applied to a new region
407 with no tuning, the models perform poorly. However, the pretrained models nonethe-
408 less predict location with greater accuracy than the models trained from random weights.
409 The benefits of transfer learning are most marked for very small datasets - after approx-
410 imately 100 events are used for training, using pretrained models has less of an advan-
411 tage over randomly initialized weights.

412 While the ranges of area and depth are equal between the two datasets, the mag-
413 nitudes of the Southern California dataset are normalized from 2.5-6 while the magni-
414 tudes of the Oklahoma dataset are normalized between 1.5-4. The tuned models were
415 able to adapt to the change in normalization given only ten events.

416 **4 Discussion**

417 Our GCN has several advantages over the FCN baseline model. One of the primary
418 advantages is the ability to make predictions on a dynamic set of inputs, allowing the
419 model to adapt to station outages, network alterations, and station subsetting. As STGCN
420 featurizes individual stations rather than an ordered network image, the model can be
421 easily trained to predict using any number of stations without architectural alteration.

422 The FCN baseline uses an image-to-image strategy, outputting a probability vol-
423 ume in which the highest values correspond to the event location. This has the advan-
424 tage of predicting a probability amplitude, which X. Zhang et al. (2020) demonstrate as
425 a useful measure of prediction uncertainty, especially in cases where earthquakes occur
426 outside the bounds of the modeled region. However, the volumetric output comes at the
427 cost of resolution limitation due to discretization. The gridded, three-dimensional out-
428 put also requires a high degree of model complexity. The FCN baseline consequently com-

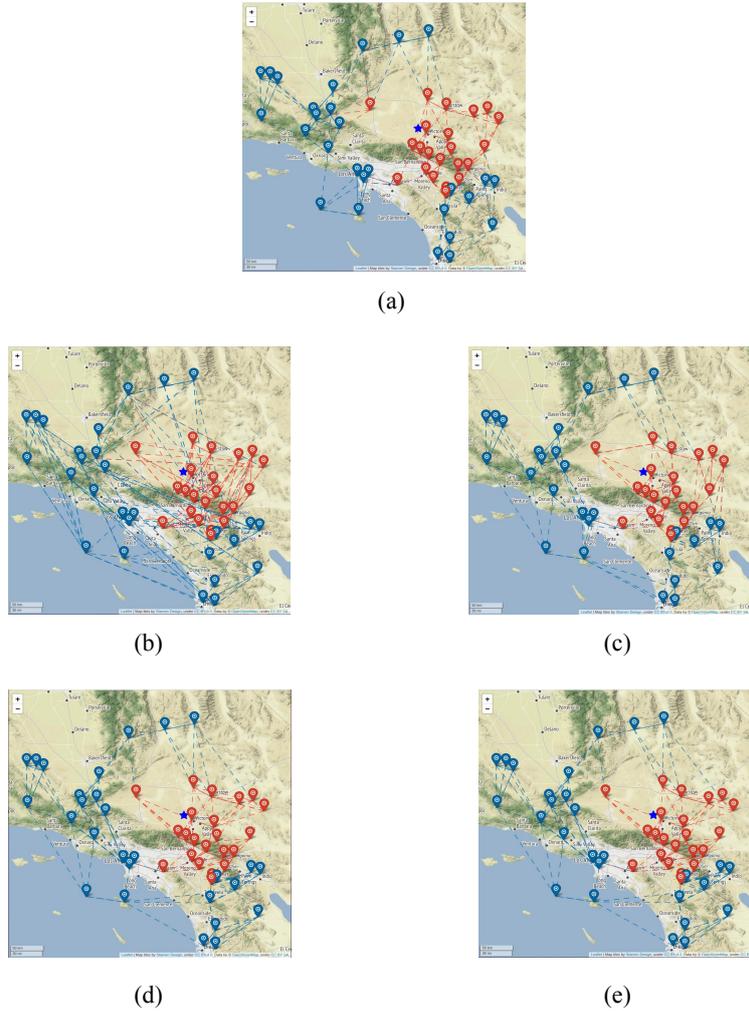


Figure 12: Graphs constructed by different layers of the graph neural network, (a) graph convolution layer based on locations of seismic stations (b) 1st, (c) 2nd, (d) 3rd and (e) 4th graph convolution layer based on the similarity of extracted features of seismic stations. Stations that detected an event in the catalogue are denoted by red symbols, while stations that did not record the event are shown in blue. Red and blue edges are generated for updating features of red and blue stations, respectively. Star represents the event location. The information from stations with the event signal are clustered in deeper layers.

429 prises approximately 27 million parameters while our GCN with scalar predictions com-
430 prises fewer than 0.24 million parameters.

431 The baseline GCN (van den Ende & Ampuero, 2020) implements edgeless graph
432 convolution (i.e. station-by-station convolutions with global pooling) while GCN model
433 developed in this paper implements convolution and pooling over dynamically-generated
434 edges. Figure 12 gives insight into the edge generation process. For clear visualization,
435 we select a case with 50 seismic stations with $K = 5$. In the edges generated by wave-
436 form similarity, stations that have recorded an event are generally connected to other
437 recording stations, forming different clusters than the edges generated by geographic prox-
438 imity. This indicates that the model is able to successfully extract waveform informa-
439 tion and associate stations in order to characterize an event. Moreover, the generated
440 graphs from the 3rd and 4th graph convolution layer based on the extracted feature sim-
441 ilarity converge to the same structure, indicating that the number of graph convolutional
442 layers is large enough to connect informative seismic stations together. If we only con-
443 sider the geographic proximity, one seismic station recording the earthquake will con-
444 nect to seismic stations without signal records only. It denotes that the feature similar-
445 ity is a proper complement of geographic proximity during aggregating features from dif-
446 ferent seismic stations.

447 After training in one region, STGCN does not transfer well to other regions with-
448 out retraining. This indicates that the models are encoding site-specific information such
449 as velocity structure or types of seismicity (i.e. anthropogenically induced earthquakes
450 in the Oklahoma dataset) as well as different magnitude range which affect predictions
451 in a different region. Performance improves significantly when a small amount of train-
452 ing data is used to tune the model. Using transfer learning to adapt a model from one
453 region to another is more effective than training a randomized model when a limited dataset
454 is available. However, best results are achieved when a model is trained for the region
455 of implementation using a catalogue of several hundred events.

456 While STGCN makes improvements in functionality and location error with respect
457 to the baseline models, the proposed framework faces challenges. Substantial improve-
458 ments have been made in the prediction of latitude and longitude, and an overall improve-
459 ment in magnitude is observed. However, magnitude does not improve in every dataset,
460 and depth predictions are highly inaccurate for all models. Accurate depth estimation
461 also poses a challenge for classical inversion methods (Zonno & Kind, 1984; Billings et
462 al., 1994; M. Zhang et al., 2014). As the machine learning models tested in this work are
463 trained in a purely supervised manner, the learned predictions are fundamentally lim-
464 ited by the accuracy of the training data. Errors in training data are likely to be a lead-
465 ing driver in model error in earthquake characterization, as systematically demonstrated
466 by X. Zhang et al. (2020) by observing the effects of induced label noise on models trained
467 with synthetic data.

468 We perform a similar test, training our model using synthetic data generated with
469 Pyrocko (Developers, n.d.). For each sample, receivers were placed randomly along a flat
470 surface, and a double-couple source with a random strike, dip, rake, magnitude, and lo-
471 cation was seeded. Both stations and events were placed with uniform probability in a
472 4° latitude by 4° longitude area (between 7° and 11° in the simulated volume). For events,
473 depth was constrained from $0.7 - 10$ km, strike from $0 - 180^\circ$, dip from $0 - 90^\circ$, and
474 rake from $0 - 360^\circ$ with a magnitude range of $2.5 < M < 6$. Using a precalculated
475 Green’s Function (https://greens-mill.pyrocko.org/iceland_reg.v2-453e36), wave
476 propagation was simulated through a 1-D velocity structure and recorded by the stations.
477 As the simulated waveforms have a sampling frequency of 2 Hz, the samples were dec-
478 imated to 20.24 Hz to be compatible with our model. We layered random noise over the
479 synthetic signals to make prediction more challenging. Non-detecting stations which record
480 only random noise without earthquake signal are also included in the input files. For smaller
481 events ($2.5 < M < 4$), $0 - 23\%$ of receivers were non-detecting, and for larger events
482 ($4 < M < 6$), $0 - 13\%$ of receivers were non-detecting. A total of 30 receivers were
483 included in each sample.

484 As demonstrated by Figure 13, when label error is eliminated, depth predictions
485 dramatically improve. This indicates that the inability to correctly predict depth is a
486 reflection of data quality rather than shortcomings within the model design. Note that
487 the synthetic experiment was designed for method validation and may not be applicable to
488 our field data due to different aspects (e.g., waveform frequency, velocity structure, etc).
489 Future improvement in depth prediction must therefore be solved by accounting for in-
490 correct depth labels. One solution may be to train using higher-quality datasets in which
491 meticulous relocation has been implemented. However, reliance on large quantities of re-
492 located sample events significantly restricts the areas in which supervised models can op-
493 erate. Another solution may be to avoid purely supervised methods, implementing so-
494 lutions which combine physics-based constraints with data-driven learning to overcome
495 inaccuracy in depth labels.

496 Another limitation that STGCN shares with the baselines is the ability to make
497 predictions only within a certain range of area, depth, and magnitude, which is also the
498 limitation of all machine-learning-based frameworks. The model outputs normalized val-
499 ues between -1 and 1 which correspond to a range selected at the beginning of training.
500 The spatial restrictions are similar to the bounds set in inversion-based methods and are
501 arguably less limiting, as the predictions made by our model are continuous and there-
502 fore not bound by grid-spacing. However, STGCN is more limited than non-machine learn-
503 ing methods with regard to magnitude prediction. Magnitudes falling above or below
504 the training range cannot be predicted by STGCN or the deep learning baselines. The
505 limited range of predictions adversely impacts the usefulness of the deep learning meth-
506 ods for applications such as Earthquake Early Warning, where magnitude saturation must
507 be avoided. The limitations posed by fixed prediction ranges are made less severe by STGCN's
508 ability to be tuned to new ranges with small amounts of training data. However, the fixed
509 prediction ranges nonetheless represent a weakness in our framework.

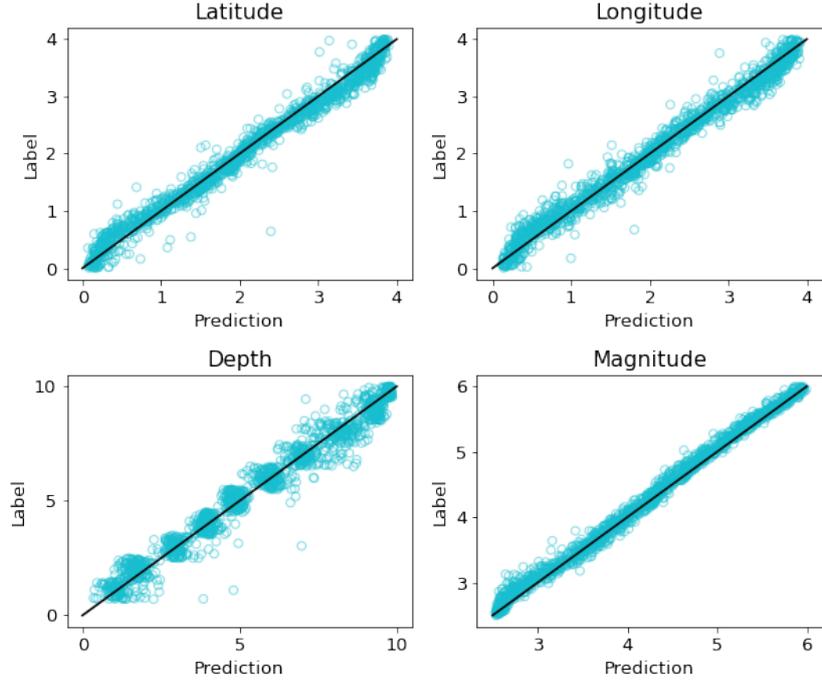


Figure 13: Testing performance of STGCN on synthetic data from 30 randomly-placed stations. In the scatter plot, each point represents an event, and a position on the diagonal line corresponds to perfect agreement between the predicted value (x-axis) and the true value (y-axis). Latitude and longitude values are displayed in degrees and depth values are displayed in kilometers.

510 **5 Conclusions and Future Work**

511 In this work, we design a graph convolutional neural network for earthquake source
 512 characterization based on waveform records from multiple stations. With experiments
 513 performed in two seismic environments, we demonstrate that STGCN outperforms both
 514 the FCN and GCN baselines, yields stable results using a range of hyperparameters, and
 515 can be applied to new datasets after retraining with a small number of events. One of
 516 the major advantages of our framework compared with other deep learning source char-
 517 acterization networks is that STGCN does not require static input or a manually gen-

518 erated graph structure. Instead, all feature generation and fusion processes are learned
519 automatically from the data to synthesize waveform features and spatial data.

520 Future improvements to our work include enhancing model capacity to predict depth,
521 a problem which synthetic tests reveal to be primarily caused by label error. This may
522 be overcome with higher-quality training data, or through methods such as physics-informed
523 machine learning. Our work thus far has focused on developing architecture to charac-
524 terize an earthquake given a discrete time series known to contain an event. Further adap-
525 tation of the core model is required to effectively process continuous waveforms in which
526 an event may not be present, or in which multiple events are contained within one win-
527 dow. An additional feature to incorporate is uncertainty quantification. Given the rel-
528 atively high degree of error in all methods for earthquake location, uncertainty is a stan-
529 dard feature in comprehensive catalogues. Uncertainty can be incorporated internally
530 (i.e. to aid in station selection) and also applied to the final predictions to identify poorly-
531 constrained events. Another interesting application is to transform the learning process
532 in an online learning manner in which a model might adaptively retrain as more recent
533 earthquakes are included in the catalogue.

534 **6 Open Research**

535 Waveform data used in this study were downloaded from the Incorporated Research
536 Institutions for Seismology (<http://ds.iris.edu/ds/nodes/dmc>) and the Southern Cal-
537 ifornia Earthquake Data Center (<https://scedc.caltech.edu/data/waveform.html>).
538 The maps in our paper were made using Generic Mapping Tools (Wessel et al., 2013)
539 and Python. The forward modelling was performed with Pyrocko (Developers, n.d.) us-
540 ing the *iceland_{reg}v2* Green’s function available at [https://greens-mill.pyrocko.org/](https://greens-mill.pyrocko.org/iceland_reg_v2-453e36)
541 *iceland_reg_v2-453e36*.

542 **Acknowledgments**

543 This work was supported by the Center for Space and Earth Science at Los Alamos
 544 National Laboratory (LANL) and by the Laboratory Directed Research and Develop-
 545 ment program of LANL under project numbers of 20210542MFR. MZ was supported by
 546 the Natural Sciences and Engineering Research Council of Canada Discovery Grant (RGPIN-
 547 2019-04297).

548 **References**

- 549 Bergen, K. J., Johnson, P. A., Maarten, V., & Beroza, G. C. (2019). Machine learn-
 550 ing for data-driven discovery in solid earth geoscience. *Science*, *363*(6433).
- 551 Beskardes, G. D., Hole, J. A., Wang, K., Michaelides, M., & Wu, Q. (2018). A com-
 552 parison of earthquake back-projection imaging methods for dense local arrays.
 553 *Geophysical Journal International*, *212*(3), 1986-2002.
- 554 Beyreuther, M., Barsch, R., Krischer, L., Megies, T., Behr, Y., & Wassermann, J.
 555 (2010). Obspy: A python toolbox for seismology. *Seismological Research*
 556 *Letters*, *81*(3), 530-533.
- 557 Billings, S., Kennett, B., & Sambridge, M. (1994). Hypocentre location: genetic al-
 558 gorithms incorporating problem-specific information. *Geophysical Journal In-*
 559 *ternational*, *118*(3), 693-706.
- 560 Developers, T. P. (n.d.). *Pyrocko: A versatile seismology toolkit for Python*. Re-
 561 trieved 2018-02-23, from <http://pyrocko.org> doi: 10.5880/GFZ.2.1.2017
 562 .001
- 563 Gajewski, D., Anikiev, D., Kashtan, B., & Tessmer, E. (2007). Localization of seis-
 564 mic events by diffraction stacking. In *Seg technical program expanded abstracts*
 565 *2007* (p. 1287-1291).
- 566 Hutton, K., Woessner, J., & Hauksson, E. (2010). Earthquake monitoring in south-
 567 ern california for seventy-seven years (1932-2008). *Bulletin of the Seismological*

- 568 *Society of America*, 100(2), 423–446.
- 569 Kaufman, S., Rosset, S., Perlich, C., & Stitelman, O. (2012). Leakage in data min-
570 ing: Formulation, detection, and avoidance. *ACM Transactions on Knowledge*
571 *Discovery from Data (TKDD)*, 6(4), 1–21.
- 572 Kong, Q., Trugman, D. T., Ross, Z. E., Bianco, M. J., Meade, B. J., & Gerstoft, P.
573 (2019). Machine learning in seismology: Turning data into insights. *Seismolog-*
574 *ical Research Letters*, 90(1), 3–14.
- 575 Kriegerowski, M., Petersen, G. M., Vasyura-Bathke, H., & Ohrnberger, M. (2019).
576 A deep convolutional neural network for localization of clustered earthquakes
577 based on multistation full waveforms. *Seismological Research Letters*, 90, 510 –
578 516.
- 579 Li, L., Tan, J., Schwarz, B., Stanek, F., Poiata, N., Shi, P., . . . Gajewski, D. (2020).
580 Recent advances and challenges of waveform-based seismic location methods at
581 multiple scales. *Reviews of Geophysics*, e2019RG000667.
- 582 Li, Z., Meier, M.-A., Hauksson, E., Zhan, Z., & Andrews, J. (2018). Machine learn-
583 ing seismic wave discrimination: Application to earthquake early warning.
584 *Geophysical Research Letters*, 45(10), 4773–4779.
- 585 Li, Z., & van der Baan, M. (2016). Microseismic event localization by acoustic time
586 reversal extrapolation. *Geophysics*, 81(3), KS123-KS134.
- 587 Lin, Y., Syracuse, E. M., Maceira, M., Zhang, H., & Larmat, C. (2015). Double-
588 difference travelttime tomography with edge-preserving regularization and a
589 priori interfaces. *Geophysical Journal International*, 201(2), 574-594.
- 590 McBrearty, I. W., & Beroza, G. C. (2022). Earthquake location and magnitude esti-
591 mation with graph neural networks. *arXiv preprint arXiv:2203.05144*.
- 592 Mousavi, S. M., & Beroza, G. C. (2020a). Bayesian-deep-learning estimation of
593 earthquake location from single-station observations. *IEEE Transactions on*
594 *Geoscience and Remote Sensing*, 1 – 14.
- 595 Mousavi, S. M., & Beroza, G. C. (2020b). A machine-learning approach for

- 596 earthquake magnitude estimation. *Geophysical Research Letters*, 47(1),
 597 e2019GL085976.
- 598 Münchmeyer, J., Bindi, D., Leser, U., & Tilmann, F. (2020, Dec). The transformer
 599 earthquake alerting model: a new versatile approach to earthquake early warn-
 600 ing. *Geophysical Journal International*, 225(1), 646–656. Retrieved from
 601 <http://dx.doi.org/10.1093/gji/ggaa609> doi: 10.1093/gji/ggaa609
- 602 Münchmeyer, J., Bindi, D., Leser, U., & Tilmann, F. (2021, Apr). Earth-
 603 quake magnitude and location estimation from real time seismic waveforms
 604 with a transformer network. *Geophysical Journal International*, 226(2),
 605 1086–1104. Retrieved from <http://dx.doi.org/10.1093/gji/ggab139> doi:
 606 10.1093/gji/ggab139
- 607 Nanometrics Seismological Instruments. (2013). *Nanometrics research network*.
 608 International Federation of Digital Seismograph Networks. Retrieved from
 609 <https://www.fdsn.org/networks/detail/NX/> doi: 10.7914/SN/NX
- 610 Perol, T., Gharbi, M., & Denolle, M. (2018). Convolutional neural network for
 611 earthquake detection and location. *Science Advances*, 4, e1700578.
- 612 Pesicek, J. D., Child, D., Artman, B., & Cieslik, K. (2014). Picking versus stacking
 613 in a modern microearthquake location: Comparison of results from a surface
 614 passive seismic monitoring array in Oklahoma. *Geophysics*, 79(6), KS61-
 615 KS68.
- 616 Ross, Z. E., Yue, Y., Meier, M.-A., Hauksson, E., & Heaton, T. H. (2019).
 617 Phaselink: A deep learning approach to seismic phase association. *Journal*
 618 *of Geophysical Research: Solid Earth*, 124(1), 856–869.
- 619 Shen, H., & Shen, Y. (2021). Array-based convolutional neural networks for au-
 620 tomatic detection and 4d localization of earthquakes in hawai ‘i. *Seismological*
 621 *Society of America*, 92(5), 2961–2971.
- 622 Tiira, T. (1999). Detecting teleseismic events using artificial neural networks. *Com-*
 623 *put. Geosci.*, 25, 929 – 938.

- 624 van den Ende, M. P., & Ampuero, J.-P. (2020). Automated seismic source char-
 625 acterisation using deep graph neural networks. *Geophysical Research Letters*,
 626 e2020GL088690.
- 627 Wang, J., & Teng, T. (1995). Artificial neural network-based seismic detector. *Bull.*
 628 *Seismol. Soc. Am.*, *85*, 308 – 319.
- 629 Wang, Y., Sun, Y., Liu, Z., Sarma, S. E., Bronstein, M. M., & Solomon, J. M.
 630 (2019). Dynamic graph cnn for learning on point clouds. *Acm Transactions On*
 631 *Graphics (tog)*, *38*(5), 1–12.
- 632 Wessel, P., Smith, W., Scharroo, R., Luis, J., & Wobbe, F. (2013). Generic mapping
 633 tools: Improved version released. *EOS*, *94*, 409-410 – 41.
- 634 Yano, K., Shiina, T., Kurata, S., Kato, A., Komaki, F., Sakai, S., & Hirata, N.
 635 (2021). Graph-partitioning based convolutional neural network for earthquake
 636 detection using a seismic array. *Journal of Geophysical Research: Solid Earth*,
 637 *126*(5), e2020JB020269.
- 638 Zhang, H., & Thurber, C. H. (2003). Double-difference tomography: The method
 639 and its application to the Hayward Fault, California. *Bulletin of the Seismolog-*
 640 *ical Society of America*, *93*(5), 1875-1889.
- 641 Zhang, M., Tian, D., & Wen, L. (2014). A new method for earthquake depth deter-
 642 mination: stacking multiple-station autocorrelograms. *Geophysical Journal In-*
 643 *ternational*, *197*(2), 1107-1116.
- 644 Zhang, X., Zhang, J., Yuan, C., Liu, S., Chen, Z., & Li, W. (2020). Locating in-
 645 duced earthquakes with a network of seismic station in Oklahoma via a deep
 646 learning method. *Scientific Report*, *10*.
- 647 Zhang, X., Zhang, M., & Tian, X. (2021). Real-time earthquake early warning
 648 with deep learning: Application to the 2016 m 6.0 central apennines, italy
 649 earthquake. *Geophysical Research Letters*, *48*(5), 2020GL089394.
- 650 Zhang, Z., Rector, J. W., & Nava, M. J. (2017). Simultaneous inversion of mul-
 651 tiple microseismic data for event locations and velocity model with bayesian

- 652 inference. *Geophysics*, 82(3), KS27-KS39.
- 653 Zhebel, O., & Eisner, L. (2015). Simultaneous microseismic event localization and
654 source mechanism determination. *Geophysics*, 80(1), KS1-KS9.
- 655 Zhu, W., & Beroza, G. C. (2019). Phasenet: a deep-neural-network-based seismic
656 arrival-time picking method. *Geophysical Journal International*, 216(1), 261–
657 273.
- 658 Zhu, W., Mousavi, S. M., & Beroza, G. C. (2019). Seismic signal denoising and
659 decomposition using deep neural networks. *IEEE Transactions on Geoscience
660 and Remote Sensing*, 57(11), 9476–9488.
- 661 Zonno, G., & Kind, R. (1984). Depth determination of north italian earthquakes
662 using grafenberg data. *Bulletin of the Seismological Society of America*, 74(5),
663 1645–1659.