



**HAL**  
open science

# Automated Seismic Source Characterization Using Deep Graph Neural Networks

Martijn P.A. van den Ende, J.-p. Ampuero

► **To cite this version:**

Martijn P.A. van den Ende, J.-p. Ampuero. Automated Seismic Source Characterization Using Deep Graph Neural Networks. *Geophysical Research Letters*, 2020, 47 (17), pp.e2020GL088690. 10.1029/2020GL088690 . hal-02931963

**HAL Id: hal-02931963**

**<https://hal.science/hal-02931963>**

Submitted on 20 Apr 2021

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

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

# Geophysical Research Letters

## RESEARCH LETTER

10.1029/2020GL088690

### Key Points:

- We propose a deep learning approach for automated earthquake location and magnitude estimation based on graph neural network theory
- This new approach processes multistation waveforms and incorporates station locations explicitly
- Including station locations improves the accuracy of epicenter estimation compared to models that are location agnostic

### Supporting Information:

- Text S1

### Correspondence to:

M. P. A. van den Ende,  
 martijn.vandenende@geoazur.unice.fr

### Citation:

van den Ende, M. P. A., & Ampuero, J.-P. (2020). Automated seismic source characterization using deep graph neural networks. *Geophysical Research Letters*, 47, e2020GL088690. <https://doi.org/10.1029/2020GL088690>

Received 2 MAY 2020

Accepted 24 AUG 2020

Accepted article online 30 AUG 2020

## Automated Seismic Source Characterization Using Deep Graph Neural Networks

M. P. A. van den Ende<sup>1</sup>  and J.-P. Ampuero<sup>1</sup> 

<sup>1</sup>IRD, CNRS, Observatoire de la Côte d'Azur, Université Côte d'Azur, Géoazur, France

**Abstract** Most seismological analysis methods require knowledge of the geographic location of the stations comprising a seismic network. However, common machine learning tools used in seismology do not account for this spatial information, and so there is an underutilized potential for improving the performance of machine learning models. In this work, we propose a graph neural network (GNN) approach that explicitly incorporates and leverages spatial information for the task of seismic source characterization (specifically, location and magnitude estimation), based on multistation waveform recordings. Even using a modestly-sized GNN, we achieve model prediction accuracy that outperforms methods that are agnostic to station locations. Moreover, the proposed method is flexible to the number of seismic stations included in the analysis and is invariant to the order in which the stations are arranged, which opens up new applications in the automation of seismological tasks and in earthquake early warning systems.

**Plain Language Summary** To determine the location and size of earthquakes, seismologists use the geographic locations of the seismic stations that record the ground shaking in their data analysis workflow. By taking the distance between stations and the relative timing of the onset of the shaking, the origin of the seismic waves can be accurately reconstructed. In recent years, machine learning (a subfield of artificial intelligence) has shown great potential to automate seismological tasks, such as earthquake source localization. Most machine learning methods do not take into consideration the geographic locations of the seismic stations, and so the usefulness of these methods could still be improved by providing the locations at which the data were recorded. In this work, we propose a method that accounts for geographic locations of the seismic stations, and we show that this improves the machine learning predictions.

## 1. Introduction

Seismic source characterization is a primary task in earthquake seismology and involves the estimation of the epicentral location, hypocentral depth, and moment of the seismic source. Particularly for the purposes of earthquake early warning, emergency response, and timely information dissemination, an estimate of the seismic source characteristics needs to be produced rapidly, preferably without the intervention of an analyst. One computational tool that satisfies these requirements is machine learning, making it a potential candidate to address the challenge of rapid seismic source characterization.

Recently, attempts have been made to apply machine learning to seismic source characterization (Käuffel et al., 2014; Kriegerowski et al., 2019; Lomax et al., 2019; Mousavi & Beroza, 2020b, 2020a; Perol et al., 2018). In the ConvNetQuake approach of Perol et al. (2018), a convolutional neural network was adopted to distinguish between noise and earthquake waveforms and to determine the regional earthquake cluster from which each event originated. This study was later extended by Lomax et al. (2019) to global seismicity. Mousavi and Beroza (2020b) employed a combined convolutional-recurrent neural network to estimate earthquake magnitudes, without location estimates. It is noteworthy that these studies focussed on the analysis of single-station waveforms as an input for the models, which goes against the common intuition that at least three seismic stations are required to triangulate and locate a seismic source. One possible explanation for the performance of these methods is that they rely on waveform similarity (Perol et al., 2018) and differences in phase arrival times (Mousavi & Beroza, 2020b). Unfortunately, since the parametrization through high-dimensional machine learning methods does not carry a clear physical meaning, this hypothesis is not easily tested.

Alternatively, a multistation approach would take as input for each earthquake all the waveforms recorded by the seismic network. One compelling argument in favor of single-station approaches is that for each earthquake there are as many training samples as there are stations, whereas in the multistation approach there is only one training sample per earthquake (the concatenated waveforms from the whole network). Since the performance of a deep learning model tends to benefit from larger volumes of data available for training, the model predictions may not improve when combining multiple station data into a single training sample. Second, microearthquakes are usually not recorded on multiple seismic stations if the seismic network is sparse, warranting further development of single-station methods. Lastly, concatenating data from multiple stations in a meaningful way are nontrivial. If the seismic network has a Euclidean structure, that is, if it is arranged in a regular pattern like for uniformly spaced seismic arrays or fiber-optic distributed acoustic sensing, the data can be naturally arranged into, for example, a 2-D image, where the distance between each pixel is representative of the spatial sampling distance. Unfortunately, most seismic networks are not arranged in a regular structure, so that the geometry of the network needs to be learned implicitly, as was attempted by Kriegerowski et al. (2019). Even though this approach yielded acceptable hypocenter location estimates, it remains an open question whether better results could be achieved when the non-Euclidean nature of the seismic network is better accounted for. Moreover, the seismic stations comprising the network may not be continuously operational over the period of interest (due to (de)commissioning, maintenance, or temporary campaigning strategies), leading to gaps in the fixed Euclidean data structure. Rather, seismic networks are better represented by a time-varying *graph* structure.

The deep learning tools most commonly used in seismology, convolutional neural networks (CNNs), and multilayer perceptrons (MLPs) (see also supporting information Text S1; Fukushima, 1980; LeCun et al., 2015; Rosenblatt, 1957; Rumelhart et al., 1986; Schramowski et al., 2020) are well suited to Euclidean data structures but are not optimal for graph data structures. One important characteristic of graphs is that they are not defined by the ordering or positioning of the data but only by the relations between data. As such, valid operations on a graph need to be invariant to the data order. This is not generally the case for CNNs, which exploit ordering as a proxy for spatial distance, nor for MLPs, which rely on the constant structure of the input features. Fortunately, much progress has been made in the field of *graph neural networks* (GNNs; Gori et al., 2005; Scarselli et al., 2009; Zhou et al., 2019), providing a robust framework for analyzing non-Euclidean data using existing deep learning tools.

In this contribution, we will demonstrate how GNNs can be applied to seismic source characterization using data from multiple seismic stations simultaneously. The method does not require a fixed seismic network configuration, and so the number of stations to be included in each sample is allowed to vary over time. Moreover, the stations do not need to be ordered geographically or as a function of distance from the seismic source. This makes the proposed method suitable for earthquake early warning and disaster response applications, in which the number and location of stations on which a given event is recorded is not known a priori.

## 2. Methods

### 2.1. Basic Concepts of Graph Neural Networks

Over the past few years, numerous deep learning techniques have been proposed that allow for the analysis of non-Euclidean data structures (Bronstein et al., 2017; Zhou et al., 2019), which has found applications in point cloud data (Qi et al., 2017; Wang et al., 2019), curved manifolds (Monti et al., 2017), and  $N$ -body classical mechanics (Sanchez-Gonzalez et al., 2019), among many others. As a subclass of non-Euclidean objects, graphs highlight relations between objects, typically represented as nodes connected by edges. Commonly studied examples of graph-representable objects include social networks (Hamilton et al., 2017), molecules (Duvenaud et al., 2015), and urban infrastructures (Cui et al., 2019). Owing to the lack of spatial ordering of graph structures, mathematical operations performed on graphs need to be invariant to the order in which the operations are executed. Moreover, nodes and relations between them (i.e., the edges) may not be fixed, and so the graph operations need to generalize to an arbitrary number of nodes and/or edges (and potentially the number of graphs) at any given moment. In essence, suitable graph operations are those that can be applied to the elements of a *set* of unknown cardinality. These can be simple mathematical operations such as taking the mean, maximum, or sum of the set, or they can involve more expressive aggregation (Battaglia et al., 2018) and message passing (Gilmer et al., 2017) operations.

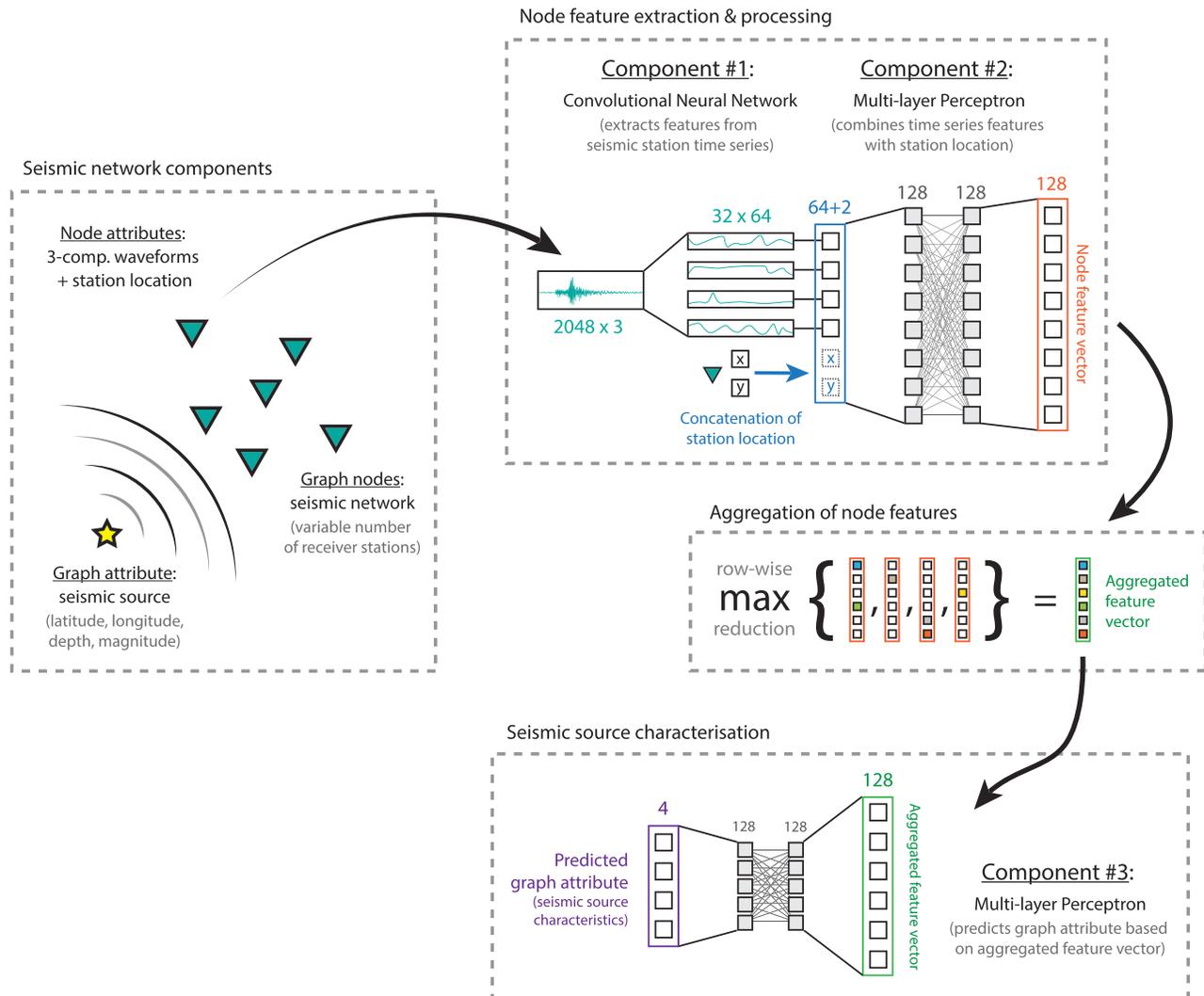
To make the above statement more concrete, we represent a seismic network by an edgeless graph in which each seismic station is a node. In the context of seismic source characterization, information travels from the seismic source to each individual receiver station independently of the relative positions between the stations. Since no information is transmitted from one station to another, it is not intuitive to include, for example, the relative distance between two stations. While local site amplifications could play an important role in the seismic source characterization process, such information should be encoded in the absolute location of each station rather than the relative location. Hence, for the task of seismic source characterization, the relations between individual stations are not physically meaningful, and so we do not include edges connecting the nodes in the analysis, reducing the graph to an unordered set. While a graph with no edges may seem ludicrous, the existence of edges is not a requirement for defining a graph, and basic architectural principles (e.g., Battaglia et al., 2018) still apply. Naturally, in cases where the relation between stations is relevant, for example, in seismic array beamforming (which relies on relative locations and arrival times), edge information should be included. Each node in our graph carries two attributes: a three-component seismic waveform time series and a geographic location. The graph itself carries four attributes: the latitude, longitude, depth, and magnitude of the seismic source. Through suitable processing and aggregation of the node attributes, the objective for the GNN is to predict the graph attributes.

## 2.2. Model Architecture

The model architecture employed in this work consists of three components that operate sequentially—see Figure 1 and Text S2 for details (Hu et al., 2020; Saxe et al., 2014; Tompson et al., 2015). Firstly, we analyze the waveforms of a given station using a CNN. This CNN processes the three-component waveform (comprising  $N_f$  time samples) and extracts a set of  $N_f$  features. The geographic location (latitude/longitude) of the seismic station is then appended to produce a feature vector of size  $N_f + 2$ . This feature vector serves as an input for the second component: An MLP that recombines the time series features and station location into a final station-specific feature vector of size  $N_q$ . This process is repeated for all  $N_s$  stations in the network using the same CNN and MLP components (i.e., the exact same operations are applied to each station individually). The convolution operations are performed only along the time axis. The output of the CNN after concatenation with each station location is then of size  $N_s \times (N_f + 2)$ , and the output of the MLP is of size  $N_s \times N_q$ .

After processing of the node attributes (the waveforms and locations of each station), the output of the MLP is max reduced over all stations to yield a graph feature vector. Empirically we have found that a max reduce yields better results than averaging or summation. The extracted features carry no physical meaning, and the information content of the feature vectors adapts to the type of aggregation during training. Hence, the most suitable type of aggregation needs to be determined experimentally. Finally, the graph feature vector is fed into a second MLP to predict the graph attributes, being the latitude, longitude, depth, and magnitude of the seismic source. Each of these source attributes is scaled so that they fall within the continuous range of  $-1 < x < +1$ , enforced by a tanh activation function in the last layer in the network. In contrast to previous work (Lomax et al., 2019; Perol et al., 2018), no binning of the source characteristics is performed. Moreover, we do not perform event detection, as this has already been done in numerous previous studies (Dysart & Pulli, 1990; Li et al., 2018; Mousavi et al., 2019; Wu et al., 2019, and others) and is essentially a solved problem. Instead, we focus on the characterization of a given seismic event. Note that the procedure above is intrinsically invariant to the number and ordering of the seismic stations: The feature extraction and recombination with the geographic location is performed for each node individually and does not incorporate information from the other stations in the network. The aggregation and the resulting graph feature vector are also independent of the number and ordering of stations. Finally, the seismic source characteristics are predicted from this invariant graph feature vector and are hence completely independent of the network input ordering and size.

To regularize the learning process, we include dropout regularization (Srivastava et al., 2014) with a dropout rate of 15% between each layer in each model component. Since the mechanics of convolutional layers are different from “dense” layers (those defining the MLPs), we use *spatial dropout* regularization (Tompson et al., 2015) that randomly sets entire feature maps of a convolutional layer to zero (as opposed to individual elements in the feature maps). The use of dropout regularization is dually motivated: First of all it reduces overfitting on the training set, as the model cannot rely on a single layer output (which could be randomly set to zero), promoting redundancy and generalization within the model. Secondly, by randomly perturbing the



**Figure 1.** Synoptic overview of the adopted model architecture. The three-component waveforms from a receiver station are fed into a CNN, after which the extracted features are combined with the station's geographic location and further processed by an MLP. The resulting node feature vector of all the stations is aggregated, and this aggregated feature vector is passed through a second MLP that predicts the seismic source characteristics.

data flow within the neural networks, the model output becomes probabilistic. The probability distribution of the model predictions for a given event can be acquired by evaluating a given input multiple times at inference time, with the variability produced by the dropout regularization. This technique is commonly referred to as Bayesian dropout (Gal & Ghahramani, 2016), as it yields a posterior distribution and hence provides a means to estimate the epistemic uncertainty for the predictions.

### 2.3. Data Description and Training Procedure

To construct a training set, we use ObsPy (Beyreuther et al., 2010) to download the broadband station inventory and earthquake catalog of the Southern California Seismic Network (SCSN; Hutton et al., 2010) over the period 2000–2015. For both the seismic station and event locations, we limit the latitude range from  $32^\circ$  to  $36^\circ$ , and the longitude range from  $-120^\circ$  to  $-116^\circ$ . The lower earthquake magnitude limit is set to 3 with no depth cut-off. In total, 1,377 events and 187 stations are included in the data set. After downloading the three-component waveforms and removing the instrument response, we filter the waveforms to a 0.1–8 Hz bandpass and interpolate onto a common time base of  $1 \leq t \leq 101$  seconds after the event origin time, over 2,048 evenly spaced time samples ( $\approx 20$  Hz sampling frequency). For an average  $P$  wave speed of  $6 \text{ km s}^{-1}$ , this time interval allows the stations at the far ends of the domain (roughly  $440 \times 440 \text{ km}$  in size) to record the event while keeping the data volume compact. The lower limit of the frequency band is chosen below

the corner frequency of the earthquakes in this analysis ( $M_w < 6$ , with corresponding corner frequency  $f_c > 0.2$  Hz Madariaga, 1976) such that information regarding the seismic moment is retained. The upper frequency limit acknowledges the common notion that attenuation and scattering rapidly reduce the signal spectrum at higher frequencies. Although the start time of all selected waveforms is fixed relative to their event origin time, the shift-equivariance of the convolution layers ensures that the extracted features are not sensitive to their timing with respect to the origin. Subsequent aggregation over the time-axis renders the features strictly time-invariant. As a result, selecting a different start of the data time window (which is inevitable when the event origin time is unknown) does not affect the model performance. The processed waveforms are then scaled by their standard deviation and stored in a database which includes the locations of the seismic stations that have recorded the events. Note that not all stations are operational at the time of a given event and hence the number of stations with recordings of the event varies.

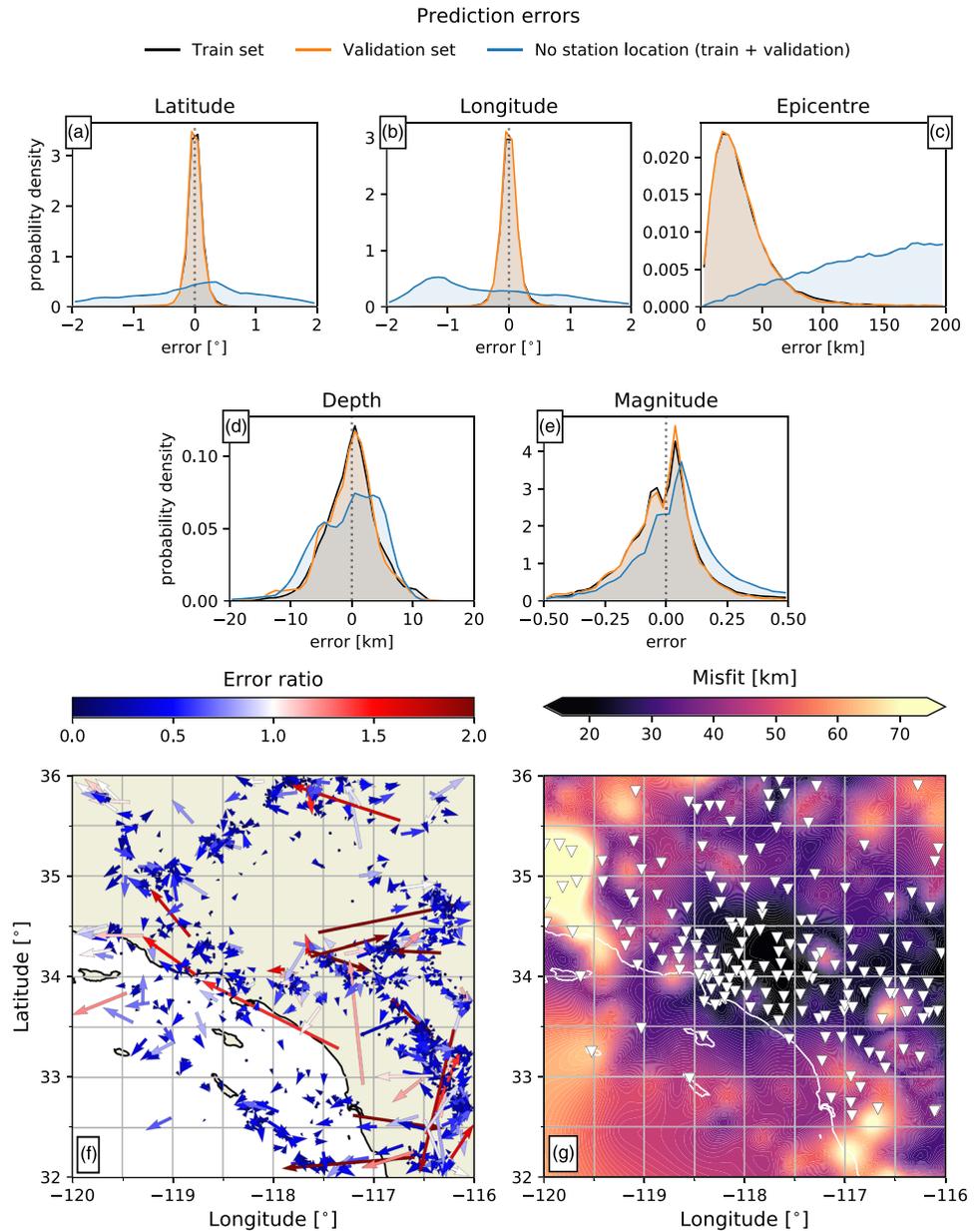
After processing the waveforms, the locations of the stations and seismic source are scaled by the minimum and maximum latitude/longitude, so that the rescaled locations fall in the range of  $\pm 1$ . Such normalization is generally considered good practice in deep learning. Similarly, the source depth is scaled to fall in the same range by taking a minimum and maximum source depth of 0 and 30 km, respectively. The earthquake magnitude is scaled taking a minimum and maximum of 3 and 6. The full data set is then randomly split 80–20 into a training set and a validation set, respectively. A batch of training samples is generated on the fly between training epochs by randomly selecting 16 training events and 50 randomly selected stations associated with each event, which we consider to strike a good balance between data volume and memory consumption. When a given event was recorded by fewer than 50 stations, the absent recordings are replaced by zeros (which do not contribute to the model performance). The model performance is evaluated through a mean absolute error loss between the predicted and target seismic source characteristics (scaled between  $\pm 1$ ), and training is performed by minimization of the loss using the ADAM algorithm (Kingma & Ba, 2017). Training is continued for 500 epochs, at which point the model performance has saturated. On a single nVidia Tesla K80, the training phase took about 1 hr in total. Once trained, evaluation of 1,377 events with up to 50 stations each takes less than 5 s of computation time (including data transfer overhead) or 3.5 ms per event.

### 3. Results and Discussion

#### 3.1. Reference Model Performance

We evaluate the performance of the trained model on both the training and validation data sets separately (Figures 2a–2e and Figure S2). The model posterior is estimated by maintaining dropout regularization at inference time (as discussed in the previous section) and performing the inference 100 times on each event in the training and validation catalogs and calculating the corresponding mean and standard deviation. Overall, the performance is similar for either data set, which indicates that overfitting on the training set is minimal. The mean absolute difference between the catalog values and the model predictions is less than  $0.11^\circ$  ( $\approx 13$  km in distance) for the latitude and longitude (which amounts to a mean epicentral location error of 18 km), 3.3 km for the depth, and 0.13 for the event magnitude. While these predictions are not as precise as typical nonrelocated estimates for Southern California (Powers & Jordan, 2010), they are obtained without phase picking or waveform amplitude modeling nor is a crustal velocity models explicitly provided (though it is implicitly encoded in the catalog hypocenter locations). Hence, the method provides a reasonable first-order estimate of location and magnitude that can serve as a starting point for subsequent refinement based on traditional seismological tools.

Since we can compute the posterior distribution for each event, we can compare the confidence intervals given by the posterior with the true epicenter location error. In Figure 2f, we plot the residual vectors between the predicted epicenter locations and those in the catalog. To visualize the model uncertainty, we compute an error ratio metric as the distance between the predicted and cataloged epicenters, normalized by the 95% confidence interval obtained from the model posterior. Hence, values less than 1 indicate that the true epicenter location falls within the 95% confidence interval, while values greater than 1 indicate the converse. Most of the predictions have an error ratio  $< 1$ . This assessment of the uncertainty in the predictions only addresses epistemic uncertainties but does not immediately address aleatoric uncertainties (errors or bias on the SCSN catalog). The epicentral errors reported for the SCSN catalog are approximately 2 km, even though an in-depth analysis of these errors suggests that this error assessment is somewhat overestimated



**Figure 2.** (a)–(e) Prediction error distributions for the trained model, for (a) latitude, (b) longitude, (c) epicentral distance, (d) depth, and (e) magnitude of each event. The model performance when including the station geographic locations is evaluated separately for the train and validation data sets, showing minimal overfitting. When the station locations are omitted, the performance is evaluated on the combined data set. (f) Residuals of the epicentral locations. Each arrow represents one cataloged event, starting at the predicted epicenter and pointing toward the catalog epicenter. The colors indicate the ratio of the misfit over the 95% confidence interval of the model posterior. Hence, blue colors indicate that the catalog epicenter falls within the 95% confidence interval and red colors that the epicenter falls outside of it. (g) Overlay of the locations of seismic stations on the interpolated prediction error (in km).

(Powers & Jordan, 2010). The expected aleatoric uncertainties are therefore much smaller than the epistemic uncertainties given by the model posterior distribution.

The spatially interpolated prediction error seems partly correlated with the local density of seismic stations (Figure 2g), as regions with the highest station density also exhibit a low prediction error. The largest systematic errors are found in the northwest and southeast corners of the selected domain, where the station density is low and where the model seems unable to achieve the bounding values of latitude and longitude. This observation can be explained by the behavior of the tanh activation function, which asymptotically

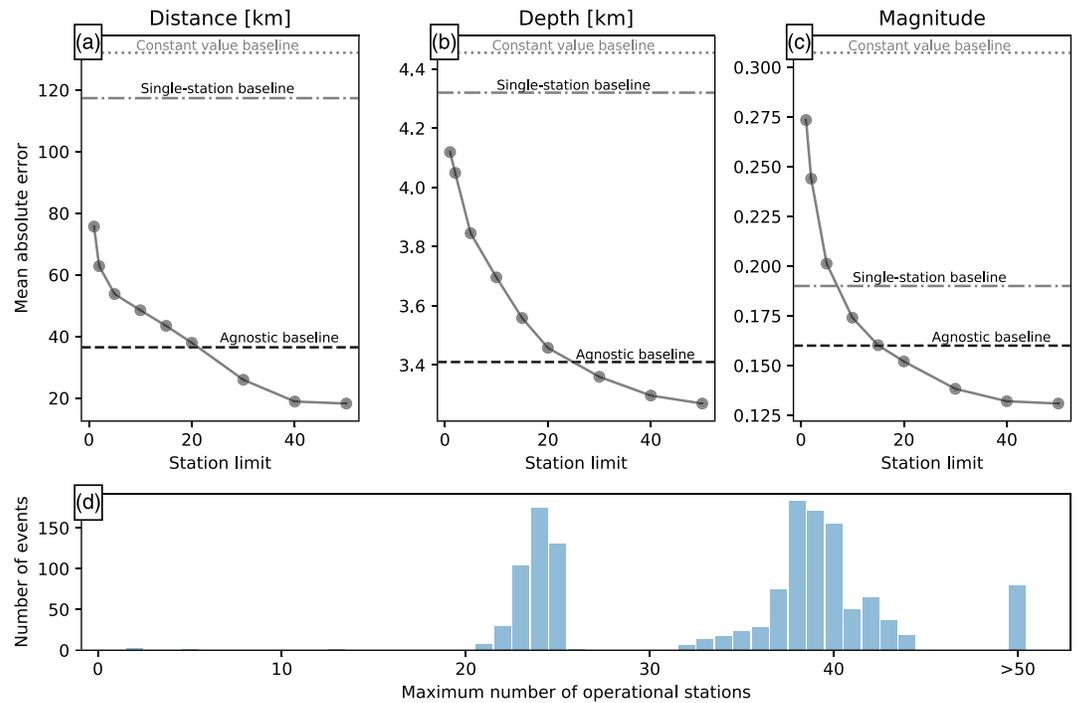
approaches its range of  $\pm 1$ , corresponding with the range of latitudes and longitudes of the training samples. Hence, increasingly larger activations are required to push the final location predictions toward the boundaries of the domain, biasing the results toward the interior. This highlights a fundamental trade-off between resolution (prediction accuracy) in the interior of the data domain, and the maximum amplitude of the predictions (which also applies to linear activation functions).

Lastly, we perform additional analyses of the sensitivity of the predictions to the signal-to-noise ratio, waveform preprocessing, and epicenter location (Figures S4–S6). These analyses show that the predictions are rather robust to the event magnitude (as a first-order proxy for signal-to-noise ratio) and insensitive to instrument corrections. Moreover, preliminary tests, in which we adopted a filter passband of 0.5–5 Hz, indicated that the choice for the prefiltering frequency band had little influence on the model performance. When the model is provided with waveforms belonging to an event with an epicenter outside of the selected training domain, the model predictions for the epicenter location collapse to an average value around the center of the domain (Figure S6). Fortunately, the uncertainty of the predictions (inferred from the posterior distribution of each event) is also much larger than for events that are located within the domain. Thus, exterior events can be distinguished from interior events through the inferred precision.

### 3.2. Influence of Geographic Information on Location Accuracy

A direct test to assess whether the station geographic location information is actually used in making the predictions (and therefore holds predictive value), we perform inference on the full data set but set the station coordinates to a fixed mean value of ( $34^\circ$ ,  $-118^\circ$ )—see Figures 2a–2e and S3. While the predictions for the event magnitude remain mostly unchanged, the estimation of the epicenter location deteriorates and becomes broadly distributed (typical for random predictions). This clearly indicates that the station location information plays an important role in estimating the epicenter locations. Thus, the adopted GNN approach, in which station location information is provided explicitly, holds an advantage over station-location agnostic methods. Interestingly, the event magnitude is almost as well resolved as when the station coordinates are included, which suggests that the model relies on the waveform data but not on station locations to estimate the magnitude. This was also observed by Mousavi and Beroza (2020b), who proposed that the relative timing of the *P* and *S* wave arrivals may encode epicentral distance information. Combined with the amplitude of the waveforms, this may implicitly encode magnitude information. Additional analysis of the contributions of (parts of) the waveforms and station locations to each of the prediction components (e.g., using Grad-CAM; Selvaraju et al., 2019) will shed more light on this.

Related to this, we investigate the effect of the (maximum) number of stations included at inference time by selecting, for each event, the stations recording the waveforms with the *M* highest standard deviations. All other waveforms are set to zero and therefore do not contribute to the predictions. If a given event was recorded by fewer than *M* stations, only the maximum number of operational stations was used with no augmentation. We perform the inference for  $M = \{1, 2, 5, 10, 15, 20, 30, 40, 50\}$  stations and compute the mean absolute error of the predictions for the epicenter location (expressed as a distance in km; Figure 3a), hypocentral depth (Figure 3b), and event magnitude (Figure 3c). For all the predicted quantities, we observe that the misfit with the catalog values rapidly decreases with the maximum number of stations included in the analysis, until the performance saturates at around  $M \geq 40$ . The reason for this saturation may lie in the distribution of the number of operational stations per event (Figure 3d). Since the majority of cataloged events is recorded by fewer than 40 stations, increasing *M* beyond 40 is only potentially beneficial only for a small number of events. For reference, we compute two performance baselines: Firstly, we take the mean value of each quantity (latitude, longitude, depth, and magnitude) over the catalog and calculate the mean absolute error relative to these. This baseline represents the performance of a “biased coin flip” (i.e., random guessing). Secondly, we train our model specifically using only a single station per training sample, through which the method specializes to single-waveform analysis (c.f. Lomax et al., 2019; Mousavi & Beroza, 2020b; Perol et al., 2018). These baselines are included in Figure 3 as horizontal dotted and dashed-dotted lines for the mean absolute error relative to the (constant value) mean and for the single-station model, respectively. Strikingly, the model that was trained on the single-station waveforms achieves worse performance in terms of the predicted hypocenter locations than the model trained on 50 stations but using only a single station at inference time. A possible explanation for this is that the single-station model may have gotten attracted to a poor local minimum in the loss landscape, after which the model started overfitting, whereas the 50-station model was able to generalize better and descended into a better local minimum.



**Figure 3.** Effect of the number of available stations on the mean absolute error of the model predictions for (a) epicentral location, (b) hypocentral depth, and (c) event magnitude. When the number of stations included at inference time is increased, the misfit between the model predictions and the catalog values decreases. The horizontal dashed and/or dotted lines in the top panels represents the baselines discussed in the text. Panel (d) displays the frequency distribution of the number of stations recording a given event.

Lastly, we compare our model performance with a model that treats the seismic network as an Euclidean object and hence has no explicit knowledge of the geographic locations of the seismic stations (“station-location agnostic”). This station-location agnostic model only features components #1 and #3 (see Figure 1 and Text S3 for details) and does not incorporate the station locations among the data features. Instead, the stations appear in a fixed order in a matrix of size  $N_s \times N_t \times 3$ , where  $N_s = 256$  denotes the total number of stations in the network (187) plus zero padding to make  $N_s$  an integer power of two. Potentially, the station-location agnostic model is able to “learn” the configuration of the seismic network and implicitly utilize station locations in predicting the seismic source characteristics. As in most traditional CNN approaches, we use a 2-D kernel of size  $k_s \times k_t$  with  $k_s = 3$  so that information from “neighboring” stations (i.e., sequentially appearing in the grid, which does not imply geographic proximity) is combined into the next layer of the model. Downsampling of the data is performed along both the temporal and station axes. Even though the number of free parameters of the station-location agnostic model is almost twice that of the graph-based model (owing to the larger convolutional kernels), and even though the model has access to all the stations simultaneously, the prediction error of the seismic source parameters is significantly larger (dashed line in Figure 3). Moreover, the station-location agnostic model required 5 times more computation time per training epoch. Hence, the GNN approach proposed here offers substantial benefits in terms of predictive power and ease of training.

### 3.3. Potential Applications

The method proposed in this study does not require the intervention of an analyst to prepare or verify the model input data (e.g., picking  $P$  and  $S$  wave first arrivals), and so it can operate autonomously. This, combined with the rapid inference time of  $\approx 3.5$  ms for 50 stations, opens up applications in automated source characterization that require a rapid response, such as earthquake early warning (EEW; Allen & Melgar, 2019), emergency response, and timely public dissemination. The aim of this study is to demonstrate the potential of incorporating seismic station locations (and possibly other node or edge attributes in a graph structure). Therefore, the model architecture was not optimized with the purpose of EEW in mind.

Nonetheless, its modular nature allows for modifications required to accommodate the real-time demands of EEW.

The first out of three components of this model consists of a CNN that analyses the waveforms of each seismic station and yields a set of station-specific features. The advantage of using a CNN is that it has immediate access to all the available information to produce a set of features optimal for the subsequent MLP components. Alternatively, a different class of deep neural networks suitable for time series analysis, the recurrent neural networks (RNN; Hochreiter & Schmidhuber, 1997; Sherstinsky, 2020), allows for online (real-time) processing of time series. Within the generalized framework of GNNs (Battaglia et al., 2018), replacing the first CNN component with an RNN produces an equally valid model architecture, although training and interpreting RNNs comes with new challenges compared with conventional CNNs. At inference time, for each new data entry the output of the first GNN component is updated for each station individually, which could be aggregated periodically to update the final model predictions, taking into account previously seen data (the “memory” of the RNN). A robust prediction will be one for which the output of the model converges to a stable estimate of hypocenter location and magnitude. Since we here employed a CNN rather than an RNN, we do not know how much time since the first ground motions is required to converge to a stable prediction, and we anticipate that this convergence depends on the quality and consistency of the data. Moreover, different components of the prediction may converge at different rates: While the hypocenter estimate may be governed by the (first) arrival of seismic energy at the various stations in the region (and therefore on the station density), the magnitude estimate is potentially controlled by the duration of the moment-rate function (Meier et al., 2017). Owing to the opacity of our deep learning method, we cannot directly assess which part of the input governs which part in the output, and so this will need to be assessed empirically.

As mentioned in section 2.2, we focused our efforts on seismic source characterization and not event detection. For any EEW task, earthquake detection is a crucial first step, which fortunately has been demonstrated to be a task suitable for machine learning methods (e.g., Dysart & Pulli, 1990; Li et al., 2018; Mousavi et al., 2019; Wu et al., 2019). In the methods proposed in the present study, earthquake detection could be performed by adding an additional graph attribute (alongside latitude, longitude, depth, and magnitude) indicating whether or not an event has been detected (similar to Lomax et al., 2019; Perol et al., 2018). Alternatively, a dedicated detection algorithm (based on machine learning or otherwise) could run in parallel and trigger the source characterization algorithm once an event has been detected. This second approach significantly reduces computational overhead. Flexibility in the number of stations included in the model input facilitates processing of an expanding data set as more seismic stations experience ground shaking after the first detection.

For the applications of emergency response and information dissemination, the real-time requirements are less stringent, so that some response time may be sacrificed in favor of prediction accuracy, maintaining the CNN component #1. Our method can be readily applied to automated earthquake catalog generation in regions where large volumes of raw data exist but which have not been fully processed. This typically arises in aftershock campaigns with stations that were not telemetered, for instance Ocean Bottom Seismometers. Given the relatively small size of the GNN employed here, retraining a pretrained model on data from a different region is relatively inexpensive. Out of the 110,836 trainable parameters, less than half (42,244) reside in the second and third components of the network. The first CNN component is completely agnostic to any spatial or regional information, as it only extracts features from time series of individual stations. Hence, if the waveforms in the target region are similar to those in the initial training region, the first component requires no retraining. This leaves only the smaller second and third MLP components to be retrained and adapted to the characteristics of the target region. As such, fewer training seismic events than employed for the initial training will be required for fine-tuning of the model. It is crucial to realize here that the second and third components potentially encode the crustal velocity structure and local site amplifications and are therefore specific to the domain that was selected during training (Southern California). Direct application of the trained model to other regions without retraining is unwarranted. The scaling of the retrained model performance with the number of stations will need to be assessed empirically, as it may be sensitive to station redundancy, and spatial coverage and density.

Aside from automatically providing an earthquake catalog, the estimates of the seismic source locations can offer a suitable starting point for additional seismological analyses. With the retrained model, the predicted hypocenter locations yield approximate phase arrival times at the various stations in the seismic network, which serve as a basis to set the windows for cross-correlation time-delay estimation and subsequent double-difference relocation. Grid-search based inversion efforts could be directed to a region around the predicted hypocenter location, rather than expanding the search of candidate source locations to a much larger (regional) domain. Even though the model predictions for the epicentral locations are larger than what conventional seismological techniques can achieve, there is merit in deep learning based automated source characterization to expedite current seismological workflows.

Lastly, we point out that the GNN approach is rather general and that it may be adopted in other applications such as seismic event detection or classification that benefit from geographic or relational information of the seismic network. Aside from predicting “global” graph attributes, like was done in this study, GNNs can also be employed to predict node or edge attributes. Examples of such attributes include site amplification factors and event detections for the nodes (seismic stations) and phase associations for the edges. Since many geophysical data are inherently non-Euclidean, graph-based approaches offer a natural choice for the analysis of these data and permit creative solutions to present-day challenges.

#### 4. Conclusions

In this study we propose a method to incorporate the geometry of a seismic network into deep learning architectures using a graph neural network (GNN) approach, applied to the task of seismic source characterization (earthquake location and magnitude estimation). By incorporating the geographic location of stations into the learning and prediction process, we find that the deep learning model achieves superior performance in predicting the seismic source characteristics (epicentral latitude/longitude, hypocentral depth, and event magnitude) compared to a model that is agnostic to the layout of the seismic network. In this way, multistation waveforms can be incorporated while preserving flexibility to the number of available seismic stations, and invariance to the ordering of the station recordings. The GNN-based approach warrants the exploration of new avenues in earthquake early warning and rapid earthquake information dissemination, as well as in automated earthquake catalog generation or other seismological tasks.

#### Data Availability Statement

Python codes and the pretrained model are available online (<https://doi.org/10.6084/m9.figshare.12231077>).

#### Acknowledgments

We thank the Associate Editor and two anonymous reviewers for their thoughtful comments on the manuscript. MvdE is supported by French government through the UCA<sup>JEDI</sup> Investments in the Future project managed by the National Research Agency (ANR) with the reference number ANR-15-IDEX-01. The authors acknowledge computational resources provided by the ANR JCJC E-POST project (ANR-14-CE03-0002-01JCJC E-POST).

#### References

- Allen, R. M., & Melgar, D. (2019). Earthquake early warning: Advances, scientific challenges, and societal needs. *Annual Review of Earth and Planetary Sciences*, 47(1), 361–388. <https://doi.org/10.1146/annurev-earth-053018-060457>
- Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., & Pascanu, R. (2018). Relational inductive biases, deep learning, and graph networks. arXiv:1806.01261 [cs, stat].
- Beyreuther, M., Barsch, R., Krischer, L., Megies, T., Behr, Y., & Wassermann, J. (2010). ObsPy: A Python toolbox for seismology. *Seismological Research Letters*, 81(3), 530–533. <https://doi.org/10.1785/gssrl.81.3.530>
- Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., & Vandergheynst, P. (2017). Geometric deep learning: Going beyond Euclidean data. *IEEE Signal Processing Magazine*, 34(4), 18–42. <https://doi.org/10.1109/MSP.2017.2693418>
- Cui, Z., Henrickson, K., Ke, R., & Wang, Y. (2019). Traffic graph convolutional recurrent neural network: A deep learning framework for network-scale traffic learning and forecasting. In *IEEE Transactions on Intelligent Transportation Systems* (pp. 1–2). IEEE. <https://doi.org/10.1109/TITS.2019.2950416>
- Duvenaud, D. K., Maclaurin, D., Iparraguirre, J., Bombarell, R., Hirzel, T., Aspuru-Guzik, A., & Adams, R. P. (2015). Convolutional networks on graphs for learning molecular fingerprints. In C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 28, pp. 2224–2232). Red Hook, NY: Curran Associates, Inc.
- Dysart, P. S., & Pulli, J. J. (1990). Regional seismic event classification at the NORESS array: Seismological measurements and the use of trained neural networks. *Bulletin of the Seismological Society of America*, 80(6B), 1910–1933.
- Fukushima, K. (1980). Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biological Cybernetics*, 36(4), 193–202. <https://doi.org/10.1007/BF00344251>
- Gal, Y., & Ghahramani, Z. (2016). Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, 48, 1050–1059.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., & Dahl, G. E. (2017). Neural message passing for quantum chemistry. *Proceedings of the 34th International Conference on Machine Learning*, 70, 1263–1272.
- Gori, M., Monfardini, G., & Scarselli, F. (2005). A new model for learning in graph domains, *Proceedings 2005 IEEE International Joint Conference on Neural Networks* (Vol. 2, pp. 729–734). Montreal, Quebec, CA: IEEE. <https://doi.org/10.1109/IJCNN.2005.1555942>
- Hamilton, W. L., Ying, R., & Leskovec, J. (2017). Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (pp. 1025–1035). Long Beach, California, USA: Curran Associates Inc.

- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Hu, W., Xiao, L., & Pennington, J. (2020). Provable benefit of orthogonal initialization in optimizing deep linear networks. arXiv:2001.05992 [cs, math, stat].
- Hutton, K., Woessner, J., & Hauksson, E. (2010). Earthquake monitoring in southern California for seventy-seven years (1932–2008). *Bulletin of the Seismological Society of America*, 100(2), 423–446. <https://doi.org/10.1785/0120090130>
- Käufel, P., Valentine, A. P., O’Toole, T. B., & Trampert, J. (2014). A framework for fast probabilistic centroid-moment-tensor determination—Inversion of regional static displacement measurements. *Geophysical Journal International*, 196(3), 1676–1693. <https://doi.org/10.1093/gji/ggt473>
- Kingma, D. P., & Ba, J. (2017). Adam: A method for stochastic optimization. arXiv:1412.6980 [cs].
- Kriegerowski, M., Petersen, G. M., Vasyura-Bathke, H., & Ohrnberger, M. (2019). A deep convolutional neural network for localization of clustered earthquakes based on multistation full waveforms. *Seismological Research Letters*, 90(2A), 510–516. <https://doi.org/10.1785/0220180320>
- LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436–444. <https://doi.org/10.1038/nature14539>
- Li, Z., Meier, M. A., Hauksson, E., Zhan, Z., & Andrews, J. (2018). Machine learning seismic wave discrimination: Application to earthquake early warning. *Geophysical Research Letters*, 45, 4773–4779. <https://doi.org/10.1029/2018GL077870>
- Lomax, A., Michelini, A., & Jozinović, D. (2019). An investigation of rapid earthquake characterization using single station waveforms and a convolutional neural network. *Seismological Research Letters*, 90(2A), 517–529. <https://doi.org/10.1785/0220180311>
- Madariaga, R. (1976). Dynamics of an expanding circular fault. *Bulletin of the Seismological Society of America*, 66, 639–666.
- Meier, M. A., Ampuero, J. P., & Heaton, T. H. (2017). The hidden simplicity of subduction megathrust earthquakes. *Science*, 357(6357), 1277–1281. <https://doi.org/10.1126/science.aan5643>
- Monti, F., Boscaini, D., Masci, J., Rodolà, E., Svoboda, J., & Bronstein, M. M. (2017). Geometric deep learning on graphs and manifolds using mixture model cnns. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 5425–5434). <https://doi.org/10.1109/CVPR.2017.576>
- Mousavi, S. M., & Beroza, G. C. (2020a). Bayesian-deep-learning estimation of earthquake location from single-station observations. In *IEEE Transactions on Geoscience and Remote Sensing* (pp. 1–14). <https://doi.org/10.1109/TGRS.2020.2988770>
- Mousavi, S. M., & Beroza, G. C. (2020b). A machine-learning approach for earthquake magnitude estimation. *Geophysical Research Letters*, 47, e2019GL085976. <https://doi.org/10.1029/2019GL085976>
- Mousavi, S. M., Zhu, W., Sheng, Y., & Beroza, G. C. (2019). CRED: A deep residual network of convolutional and recurrent units for earthquake signal detection. *Scientific Reports*, 9(1), 1–14. <https://doi.org/10.1038/s41598-019-45748-1>
- Perol, T., Gharbi, M., & Denolle, M. (2018). Convolutional neural network for earthquake detection and location. *Science Advances*, 4(2), e1700578. <https://doi.org/10.1126/sciadv.1700578>
- Powers, P. M., & Jordan, T. H. (2010). Distribution of seismicity across strike-slip faults in California. *Journal of Geophysical Research*, 115, B05305. <https://doi.org/10.1029/2008JB006234>
- Qi, C. R., Yi, L., Su, H., & Guibas, L. J. (2017). Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (pp. 5105–5114). Long Beach, California, USA: Curran Associates Inc.
- Rosenblatt, F. (1957). The perceptron: A perceiving and recognizing automaton (*Tech. Rep. No. 85-60-1*). Buffalo, New York: Cornell Aeronautical Laboratory.
- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, 323(6088), 533–536. <https://doi.org/10.1038/323533a0>
- Sanchez-Gonzalez, A., Bapst, V., Cranmer, K., & Battaglia, P. (2019). Hamiltonian graph networks with ode integrators. arXiv:1909.12790 [physics].
- Saxe, A., McClelland, J. L., & Ganguli, S. (2014). Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *International Conference on Learning Representations*.
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., & Monfardini, G. (2009). The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1), 61–80. <https://doi.org/10.1109/TNN.2008.2005605>
- Schramowski, P., Stammer, W., Teso, S., Brugger, A., Luigs, H. G., Mahlein, A. K., & Kersting, K. (2020). Right for the wrong scientific reasons: Revising deep networks by interacting with their explanations. arXiv:2001.05371 [cs, stat].
- Selvaraju, R. R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., & Batra, D. (2019). Grad-cam: Visual explanations from deep networks via gradient-based localization. *International Journal of Computer Vision*, 128, 336–359. <https://doi.org/10.1007/s11263-019-01228-7>
- Sherstinsky, A. (2020). Fundamentals of recurrent neural network (RNN) and long short-term memory (LSTM) network. *Physica D: Nonlinear Phenomena*, 404, 132306. <https://doi.org/10.1016/j.physd.2019.132306>
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56), 1929–1958.
- Tompson, J., Goroshin, R., Jain, A., LeCun, Y., & Bregler, C. (2015). Efficient object localization using convolutional networks. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 648–656). <https://doi.org/10.1109/CVPR.2015.7298664>
- Wang, Y., Sun, Y., Liu, Z., Sarma, S. E., Bronstein, M. M., & Solomon, J. M. (2019). Dynamic graph cnn for learning on point clouds. Association for Computing Machinery.
- Wu, Y., Lin, Y., Zhou, Z., Bolton, D. C., Liu, J., & Johnson, P. (2019). DeepDetect: A cascaded region-based densely connected network for seismic event detection. *IEEE Transactions on Geoscience and Remote Sensing*, 57(1), 62–75. <https://doi.org/10.1109/TGRS.2018.2852302>
- Zhou, J., Cui, G., Zhang, Z., Yang, C., Liu, Z., Wang, L., & Sun, M. (2019). Graph neural networks: A review of methods and applications. arXiv:1812.08434 [cs, stat].