# Bayesian-Deep-Learning Estimation of Earthquake Location from Single-Station Observations

S.Mostafa Mousavi[1*] & Gregory C. Beroza[1†]

[1]*Geophysics Department, Stanford University, Stanford, California, USA*

**We present a deep learning method for single-station earthquake location, which we approach as a regression problem using two separate Bayesian neural networks. We use a multi-task temporal-convolutional neural network to learn epicentral distance and P travel time from 1-minute seismograms. The network estimates epicentral distance and P travel time with absolute mean errors of 0.23 km and 0.03 s respectively, along with their epistemic and aleatory uncertainties. We design a separate multi-input network using standard convolutional layers to estimate the back-azimuth angle, and its epistemic uncertainty. This network estimates the direction from which seismic waves arrive to the station with a mean error of 1 degree. Using this information, we estimate the epicenter, origin time, and depth along with their confidence intervals. We use a global dataset of earthquake signals recorded within 1 degree (~112 km) from the event to build the model and to demonstrate its performance. Our model can predict epicenter, origin time, and depth with mean errors of 7.3 km, 0.4 second, and 6.7 km respectively, at different locations around the world. Our approach can be used for fast earthquake source characterization with a limited number of observations, and also for estimating location of earthquakes that are sparsely recorded - either**

---
*E-mail: mmousavi@stanford.edu
†E-mail: beroza@stanford.edu

**because they are small or because stations are widely separated.**

## 1   Introduction

Recent years have seen a renewed surge of interest in applying machine-learning techniques to seismic signal processing and earthquake monitoring. Promising results from multiple studies indicate neural-network-based models can outperform traditional algorithms in tasks such as: earthquake signal detection [1,2], phase picking [3–6], first-motion polarity determination [7,8], denoising [9,10], discrimination [11,12], and association [13–15].

However, earthquake location remains a challenging task. Perol et.al. [1] trained a convolutional neural network to simultaneously learn classification (event vs. noise) and to group earthquakes into 6 clusters initially defined by K-means in Oklahoma. Lomax et.al. [16] expanded this approach by classifying seismic waveforms into a larger number of classes including: event/noise (1 class), station–event distance (50 classes), station–event azimuth (36 classes, each 10 degrees), event magnitude (20 classes), and event depth (20 classes); however, their model did not generalize well and suffered from high error rates. On the other hand multi-station approaches (e.g. [17]) result in a better performance by learning the move out patterns for specific station configuration at a local region.

Neither these, nor other, neural networks applied to earthquake data quantify uncertainty in

their output. Machine learning can be thought of as inferring plausible models that explain data and can be used to make predictions about unseen data. Uncertainty plays a key role in that process of quantifying the reliability of those predictions. Data can be consistent with different models and the question of what model is appropriate based on such data is uncertain. Predictions using future data are also uncertain [18]. Typical deep learning models do not capture uncertainties in the output. In regression models, output is a single vector that regresses to the mean of the data, but in classification models, and the output probability is not equivalent to model confidence. That is, a model can be uncertain in its predictions even with a high softmax output [19]. With model confidence we can treat uncertain input and special cases properly. In the case of classification models for earthquake signal detection and arrival time measurement, for example, we might pass uncertain cases to a human analyst for expert analysis.

Neural networks with model uncertainty are known as Bayesian neural networks. They offer a probabilistic interpretation of models usually by placing prior probability distributions over the network weights.

In this paper, we approach single-station earthquake location as a regression problem using two separate Bayesian neural networks. For learning epicentral distance and P travel time we designed a multi-task temporal convolutional network. The network consists of causal dilated convolutions and residual connections that estimate epicentral distance and P travel time simultaneously along with their epistemic and aleatory uncertainties. We use a separate multi-input network with standard convolutional layers to estimate the back-azimuth and its epistemic uncertainty. Using

this information, we estimate the epicenter, origin time, and depth along with their confidence intervals. We use a global data set for building the model and for demonstrating its performance. The proposed approach can be used for rapid earthquake source characterization using a limited number of observations. This can have many different applications including in earthquake early warning systems [20] and in locating earthquakes that are sparsely recorded.

## 2   Results

### Dataset

We use the STanford EArthquake Dataset (STEAD) [21] for the training and testing of the models. STEAD is a global dataset of labeled 3-component seismic waveforms (earthquake and non-earthquake). Here, we only use earthquake waveforms recorded at epicentral distances of less than 110 km with signal-to-noise ratio of 25 decibels and higher. We only use stations for which north-south and east-west components are properly aligned to their correct geographic orientations. Based on these criteria, we select ~ 150,000 waveforms to be used for the training (% 80) and testing (% 20) of the networks. We show the geographical distribution of the events associated with these waveforms and their characteristics Figure 1 and Figure 2. Waveforms are 1 minute in duration with a sampling rate of 100 HZ and were band-passed filtered from 1-45 HZ.

### Network Architecture

We designed two separate networks, one for predicting the epicentral distance and P travel

4

time (dist-PT network) and the other for back-azimuth estimation (BAZ network).

The dist-PT network is a multi-task temporal convolutional network consisting of 1D convolutional layers where convolutions are causal and dilated Figure 3. The input to the network is a $6000 \times 4$ matrix where the first three rows are 3-component waveforms (each 1 minute long with 100 samples per second) and the last row is a binary vector where values between P and S arrival times are set to 1 and the rest to zero. The last vector is similar to the labeling in output layers of detector/picker networks; however, here we use it as the input to highlight the part of the waveforms that contain the most important information for the regression tasks.

The main body of the dist-PT network consists of 11 dilational convolution layers (dilation rate doubles for each layer) each with a relu activation function and 20 kernels of size 6. At the end, network has two fully connected layers each with a linear activation function and two neurons. The network has 58,500 trainable parameters. The full description of the optimization function is given in the method section. We applied dropout to every dilated convolutional layer in the network and trained the model with a dropout rate of 0.20. The aleatory uncertainties are implicitly learned from a customized loss function during the training without a need for uncertainty labels. This loss function acts as an intelligent regression function that makes the model robust to noisy data. We sample the posterior distribution over the weights to obtain the posterior probabilities and to estimate uncertainties.

Back-azimuth estimation is a continuous orientation prediction, which prohibits the direct use of a typical L2 loss function because the angle is in a non-Euclidean space. To handle this

5

problem, we represent back-azimuth angles, *baz*, as points on a unit circle $baz = (\cos \theta, \sin \theta)$ during the training, and convert the predicted 2D points to the back-azimuth angles during testing.

The back-azimuth network (Figure 4) primarily consists of 1D convolutional layers and has two inputs: 1) a $(150 \times 3)$ matrix (0.5 second before and 1 second after the P arrival); 2) a $(7 \times 3)$ matrix consisting of the covariance matrix, eigenvalues, and the eigenvectors derived from the 3-component waveforms for the same time window. We convolve the two input matrices with 4 and 1 convolutional layers and feed them into two fully connected layers with 100 and 2 neurons respectively to predict the coordinates of the back-azimuth angle on the unit circle. All the other layers have relu activations, except for the last fully connected layer. The kernel size used in all convolutional layers is 3 while the number of kernels varies between 20 and 64. Overall, the network is very light and only has ~46,000 trainable parameters. Using multiple inputs and point estimations prevent us from a stable estimation of aleatory uncertainty using the intelligent loss function (Equation (2)). This loss function estimates uncertainties for each output ($\cos \theta$ and $\sin \theta$) separately. Moreover, it is hard to assign the estimated uncertainties to the corresponding input. Our attempts to estimate a single aleatory uncertainty for both components did not in a stable optimization. However, we estimate the epistemic uncertainty using the Monte Carlo dropout sampling procedure described in the method section and as done for the previous network. We use dropout rates of 0.1 and 0.3 after the convolutional and fully connected layers respectively.

**Regression Results**

The regression results for the test set are presented in Figure 5. The best coefficient of

determination is obtained for P travel time estimation. The network is able to estimate P travel time with a standard deviation of 0.66 second. The mean error for epicentral distance estimates is 0.23 km with a standard deviation of 5.42 km. Compared to these results, back-azimuths estimates are more uncertain. This is mainly due to the complication of estimating orientation; however, a coefficient of determination of 0.87 for the regression results and a mean error rate of ~ 1 degree can be considered good results due to the fact that only 1.5 seconds of the waveforms are used for these estimates.

There is some positive correlation between the estimated uncertainties and prediction errors for both distance and P travel time (Figure 6), as expected. This correlation is slightly stronger for the epicentral estimates, which suggests that estimated uncertainties might be used in practical applications to identify uncertain predictions. In both cases, the aleatory uncertainties reflect the errors better than the estimated epistemic uncertainties. This indicates a lesser role for model errors in the final output.

We do not estimate the aleatory uncertainty for the back-azimuth estimates due to technical complications; however, we still observe a weak positive correlation between the estimated model uncertainty and the errors in the predictions (Figure 7).

Errors in both distance and P travel time estimates seem to increase slightly with the increase in station-event distance; however, such a correlation is not apparent for the back-azimuth estimates (Figure 8). The increase in distance has a clear effect on both estimated aleatory and epistemic uncertainties where the effect on aleatory uncertainty follows a power law, but is near-linear for

7

the epistemic uncertainty.

While the effects of signal-to-noise ratio of waveforms on regression errors are visible, the correlation with the estimated uncertainties is less clear. Events with magnitudes larger than 2.5 have higher uncertainties. This is likely due to the fact that we have less training data for larger magnitudes (Figure 8).

An interesting observation is the negative bias in distance estimation for deeper events (where the predicted distances are larger than the actual distances) (Figure 8). This explains the over estimations in the upper left side of the regression line in Figure 5-a. Both epistemic and aleatory uncertainties are higher for the deeper events, which are fewer in number in the training dataset. We see a similar (though weaker) trend of higher epistemic uncertainties for the back-azimuth estimates, but most of the errors for the back-azimuth estimates are caused by shallow events.

**Location Results**

We used the distance and back-azimuth predictions to estimate the epicenter of the associated event for each observation in the test set. We calculate the error ellipse for each epicentral location based on estimated uncertainties for distance and back-azimuth and their projections onto the reference Earth model. We use the P travel time estimates to calculate origin times and to provide a rough estimate of earthquake depth. For the depth estimation we assume that the P waves follow a straight-line path between source and station. We assumed a velocity of 5.6 km/s for the P wave and calculated the incident angle using the estimated distance the P wave has traveled together

with the estimate of epicentral distance. Estimated locations and associated error ellipses for 16 events are shown in Figures 9 and 10.

We estimate locations and errors for each observation (station) and averaged for each event based on the number of available observations. To demonstrate generalization of the model, the examples are selected from different regions in Asia, Africa, Central US, Nevada, San Juan island, Southern California, and Alaska. The uncertainties for the horizontal location, depth, and origin time for the cataloged events are presented if they have been reported by the monitoring networks. We selected examples randomly to reflect an unbiased representation of the model's performance. Figure 9 presents some examples with moderate errors. Errors in the origin time estimates are the lowest in general, which is consistent with the regression results. From the error ellipses, we can see a higher contribution of uncertainties in the back-azimuth estimations, which is again in agreement with the regression results. Location estimate errors are in a reasonable range considering the reported uncertainties for the catalog locations (e.g. Figure 9 -f, g, and h); however, we note that our location estimates are based on only single-station observations and without the use of any velocity model. Figure 9 -c, d, and e suggest that if we include the location uncertainties for individual estimates into the averaging process, the final location estimates might improve for the cases with multiple observations.

A surprising result in these examples are the relatively good estimates of earthquake depths. Depth estimation is a persistently difficult and uncertain, yet important factor in source characterization. Our depth estimate results are within an acceptable range considering that we do not

train our network directly for end-to-end learning of depth (because of the high uncertainty and inhomogeneity in reported depths by different networks around the world that make the labeling and learning process challenging). Instead, our estimates are derived by combining the predicted values for the P travel distance and epicentral distance.

Figure 10 presents examples where our method performed poorly; however, even in the cases of unsuccessful location, single-station predictions are more-or-less pointing toward the source location. In most of these cases estimated location uncertainties are relatively large, which could be used to distinguish between these estimates and more reliable ones. An interesting observation is that even in these cases where location errors are large (mostly due to errors in the back-azimuth estimates), estimated origin time is still very close to the ground truth. This robust estimation of origin time in addition to a rough estimation of back-azimuth angle could be used for event association across a network.

To get a broader view of the performance of our location estimates, we plotted the predicted epicenters paired with the cataloged locations for two regions of Alaska and northern California inFigure 11. We can see that the predicted locations reveal the overall pattern of seismicity correctly, and that the outliers are sparse. For example, the linear seismicity on the northern end of the creeping section of the San Andreas fault (between 36 and 37 degrees north) and the event cluster in Geysers (38.8 degrees north) is clearly recovered in (Figure 11-b). These are averaged results for the event-based estimates without removing any high-uncertainty estimates or any weighting of the results during the averaging.

The overall performance of our locations for the entire test set can be seen from the error distributions presented in Figure 12. Our method predicts epicenter, origin time, and depth with a mean error of 7.3 km, 0.4 seconds, and 6.7 km respectively. These are in agreement with our previous observations and with regression results.

To understand these errors and their potential sources further, we show each of them plotted as a function of event magnitude, depth, predicted uncertainties, and reported uncertainty in the catalog in Figure 13. Larger events tend to have larger prediction errors, which may be attributable to the more sparse training data for larger events in the dataset. The strong performance for small events indicates the sensitivity of the method. The ability to locate smaller events is important because small events are exactly those that are most likely to be detected on fewer stations, and a single-station method might be the only way to locate them. There is a strong correlation between errors in the depth estimation and the event depth. Very shallow events and events deeper than 20 km have larger errors. Predicted uncertainties seem to correlate with the errors in the predictions. The fifth row in Figure 13 suggests some part of the mismatch between the predictions and ground truth might be due to uncertainties in the reported location, origin time, or depth. The relations between estimated uncertainties (by our model) and reported uncertainties in the earthquake catalogs are presented in the last column.

## 3  Discussion and Conclusions

We present a successful application of deep learning for earthquake location based on single-station observations. The model is trained and tested using a global data set. Our test results indicate that our neural network can directly learn a general mapping function between the raw 3-component seismograms (and known P and S arrivals) and epicentral distance, P travel time, and back azimuth without the need for a local velocity model. Distance and back-azimuth predictions can be directly expanded to regional and teleseismic distances by training with waveforms from these events; however, the current approximate depth estimation procedure will not be applicable due to a combination of more complex wave propagation and the sphericity of the Earth. A distinct advantage of our approach lies in its Bayesian framework, which provides an estimate of uncertainties in data and model and allows us to estimate confidence intervals in the final estimated location.

Our strategy in using a multi-task network for closely related tasks (i.e., distance and travel time) and separate networks for the other task (i.e., back-azimuth) improve the overall performance. The designed networks are light with a relatively low number of trainable parameters, which requires fewer data in the training set. Availability of larger training sets with high-quality labels will allow us to design more flexible/powerful networks and potentially to improve the performance by reducing epistemic uncertainties. Deeper networks with more convolutional layers have the potential, for example, to reduce the sensitivity of the model to noise.

Building large training sets for these tasks is challenging because, in addition to the quality of general labeling such as existence of earthquake signal in the window, accuracy of picks, accuracy

12

of metadata used for estimate of back-azimuth, P travel time, and distance etc, station orientations are also important. Our current method does not consider station orientations that differ from the geographical orientations, and is not applicable to the borehole stations. This prevents us from using the high signal-to-noise ratio data provided by borehole instruments.

The largest uncertainty in location results is caused by uncertainties in the back-azimuth. Errors and high uncertainties in back-azimuth estimates may be due to seismic station installation and orientation error [22,23]; however, learning the orientation angle using neural networks is technically challenging as well. Utilizing more advanced methods for continuous orientation estimation (e.g. [24]) might improve the results. Moreover, previous studies [25] showed the length of the window used for polarization estimation play an important role in precision of back-azimuth estimation. Optimizing this window length based on dominate signal frequency might be a potential solution for this [26]. We note that our method might provide a novel approach to detecting mis-oriented seismometers.

Our results indicate superior performance of the neural-network based approach compared with traditional singel-station location methods (e.g. [27–29]). The proposed method can provide a rapid estimate of earthquake location directly from single instruments. This may be useful for rapid public reporting or earthquake early warning systems [30]. Lockman and Allen [31] investigated the accuracy of event parameter determination using single station for the purpose of early warning using P-wave arrival only. They concluded that estimated hypocentral distance and backazimuth with accuracy of$\pm15$ km and $\pm20°$ respectively that was obtained by the high-quality stations in

southern California are sufficient to provide useful early warning. We showed that our approach can result in a much higher accuracy for smaller events with more complicated high-frequency waveforms. Our approach has the potential to function even using a portion of waveform (e.g. P-wave only) through augmentation. On the other hand, the provided information about P and S arrival times (as an extra vector in addition to the three-component waveforms) are mainly used to speed up the training process by directing the attention of the network to the informative part of the input data and our tests indicate the network can perform well even when picks are not available.

## 4   Methods

**Bayesian deep learning**

Bayesian deep learning lies at the intersection of Bayesian statistical theory and deep learning, and makes it possible to express different aspects of uncertainty in deep learning models probabilistically. In Bayesian-deep learning, all quantitities are represented as probability distributions rather than a point estimates. Prior probability distributions are applied over model parameters and are used to represent how they relate to the data. Uncertainties in the observed data are then inferred using probability theory. Learning from a training dataset is done by transforming the prior probability distributions (defined before training on data), into posterior distributions (determined after observing data) that captures a set of plausible model parameters given the data [18].

**Types of Uncertainties**

14

Three types of uncertainty contribute to the model predictions; aleatory uncertainty, epistemic uncertainty, and ontological uncertainty.

The aleatory (also known as irreducible, inherent, stochastic, or type-A) uncertainty captures uncertainty with respect to information that our data cannot explain. Aleatory uncertainty can be divided into two sub-categories of Heteroscedastic (data-dependent) and Homoscedastic (task-dependent) uncertainties.

**Uncertainty Estimation**

Let $\widehat{y}$ be the output of a neural network model with a loss function $\ell(.,.)$ (i.e. the Euclidean for regression). We denote by $\boldsymbol{W}$ the network's weight (model parameters) and by $\boldsymbol{X} = \{x_1, \ ..., \ x_N\}$ and $\boldsymbol{Y} = \{y_1, \ ..., \ y_N\}$ a set of $N$inputs and outputs respectively. The objective is to learn unknown parameters, $\boldsymbol{W}$, by minimizing the loss between predictions and actual values $(\ell(y_i, \widehat{y_i}))$. Assuming a Gaussian distribution for $\boldsymbol{Y}$the data likelihood is defined as$p(\boldsymbol{Y}|\boldsymbol{X}, \boldsymbol{W})$. Given $\boldsymbol{X}$and $\boldsymbol{Y}$, the posterior probability distribution of weight, $p(\boldsymbol{W}|\boldsymbol{X}, \boldsymbol{Y}) = p(\boldsymbol{Y}|\boldsymbol{X}, \boldsymbol{W})p(\boldsymbol{W})/p(\boldsymbol{Y}|\boldsymbol{X})$, can be inferred via Bayesian theory that provides a set of model parameters. The prior probability distribution over weight, $p(\boldsymbol{W})$, then can be used to estimate the weight uncertainty [32].

Although conceptually simple, this is difficult to perform in practice because the marginal probability, $p(\boldsymbol{Y}|\boldsymbol{X})$, can not be evaluated analytically and requires some approximations [33–35]. In these approximation techniques, a simple distribution, $q(\boldsymbol{W})$, over the network's weights is learned

by minimizing the Kullback-Leibler (KL) divergence between the approximating distribution and the full posterior, $KL(q(\boldsymbol{W})||p(\boldsymbol{W}|\boldsymbol{X},\boldsymbol{Y}))$.

Here, we use Monte Carlo dropout sampling [19,35] as a variational Bayesian approximation for this purpose. Dropout [36], is a simple regularization technique commonly used to prevent over-fitting in deep learning. Dropout randomly removes network units during the training and by doing this samples from a number of trained networks with reduced width. For test data, dropout approximates the effect of averaging the predictions of thinned networks using the weights of the un-thinned network [37]. Gal and Ghahramani [19] showed that a neural network with with dropout applied during both training and test times is mathematically equivalent to an approximation to the probabilistic deep Gaussian process [38]. This is equivalent to Monte Carlo sampling from an approximate posterior distribution over models that find an approximating distribution ($q(\boldsymbol{W})$) with minimum KL divergence to the posterior probability distribution. This technique is computationally efficient, and unlike other approximations, can be easily applied to large and complex networks.

Here, we use dropout to approximate the posterior distribution for a model, $f$, that maps an input directly to the predictive mean and variance:

$$[\widehat{y},\widehat{\sigma}^2] = f(x) \tag{1}$$

where $f$ is a Bayesian convolutional neural network that transform the input seismogram, $x$, directly to an output value, $\widehat{y} \in \mathbb{R}$, and a measure of aleatory uncertainty given by variance, $\sigma^2$. This variance (noise in the data) can be directly estimated from the data using a customized loss

function that includes both the regression and an uncertainty regularization term [39]:

$$\ell_{BNN} = \frac{1}{2N} \sum_{i=1}^{N} exp(-s_i) \; ||y_i - \widehat{y_i}||^2 + \frac{1}{2} s_i \tag{2}$$

where $s_i = \log \widehat{\sigma}_i^2$. The variance is data-dependent and is learned implicitly from the loss function during the training without the need for uncertainty labels. It represents the hetroscedastic aleatoric uncertainties and is useful in cases where input data point have different noise levels. The predictive uncertainty (a combination of epistemic and aleatory uncertainty) then can be estimated as:

$$Var(y) \approx \frac{1}{T} \sum_{t=1}^{T} \widehat{y}_t^2 - \left( \frac{1}{T} \sum_{t=1}^{T} \widehat{y}_t \right)^2 + \frac{1}{T} \sum_{t=1}^{T} \widehat{\sigma}_t^2 \tag{3}$$

where $T$ is the number of Monte Carlo dropout samples.

Simultaneous prediction of output and uncertainty by network makes the model robust to noisy data. The loss function acts as an intelligent regression function that allows the network to learn to attenuate the effect of erroneous labels by adapting the residual's weighting [39].

**Temporal Convolution Networks**

Earthquake signals contain sequential information. High-frequency compressional P-waves arrive before transverse S-waves, which in turn arrive before dispersive surface waves. These temporal dependencies among different components of an earthquake signal arise from the physics of elastic wave propagation and motivate the use of neural networks that are capable of sequence modeling. For this reason, recurrent neural nets (e.g. LSTM [40] or GRU [41]) particularly suitable for modeling earthquake signals.

It has been shown that a combination of LSTMs and CNNs can achieve good performance in learning both the local structures in a seismic signal and the temporal dependencies among these structures while reducing the computational time of sequential learning [2]; however, recent studies have also demonstrated that certain convolutional architectures can outperform recurrent architectures across a diverse range of sequential-modeling tasks and data sets, while demonstrating the same capacity with longer effective memory [42]. These convolutional architectures have the advantages of architectural simplicity, parallelism, flexible receptive field size, stable gradients, low memory requirements for training, and variable length inputs. Such temporal convolution networks, are a family of autoregressive feed-forward models with causal dilated convolutions and residual connections. Temporal convolution networks are the latest development in sequential modeling in speech recognition, time series analysis, natural language processing, and signal processing.

Our network for distance/P-travel time estimation consists of 1D convolutional layers where convolutions are causal and dilated. Causal convolution means that the information flow is only from past to future, such that to make a prediction at time point $t$, only the information up to time $t$, is used.

Dilated convolutions (also called trous, or convolution with holes) are used to learn longer history by increasing the receptive field exponentially [43,44]. They allow for having very large receptive fields and they can learn multi-scale structures without greatly increasing the number of parameters. A large receptive field with small kernel size makes it possible to combine local and

global informationFigure 14.

Residual structure [45], allows a deeper network without degradation with a larger receptive field. Each residual block consists of two layers of dilated causal convolution and rectified linear units [46]. Here we use 1D-spatial droupouts [36] for regularization and also uncertainty estimation. Since in our architecture the input and output have different widths, a $1 \times 1$ convolution is used instead of identity mapping to match the size for addition operation.

The dilation factor, $d$ is increased exponentially with the depth of the network to ensure that all of the inputs within the receptive field are convolved with some filters while allowing for a very large effective history.

1. Perol, T., Gharbi, M. & Denolle, M. Convolutional neural network for earthquake detection and location. *Science Advances* **4**, e1700578 (2018).

2. Mousavi, S. M., Zhu, W., Sheng, Y. & Beroza, G. C. CRED: A deep residual network of convolutional and recurrent units for earthquake signal detection. *Scientific reports* **9**, 10267 (2019).

3. Zhu, W. & Beroza, G. C. PhaseNet: a deep-neural-network-based seismic arrival-time picking method. *Geophysical Journal International* **216**, 261–273 (2018).

4. Pardo, E., Garfias, C. & Malpica, N. Seismic Phase Picking Using Convolutional Networks. *IEEE Transactions on Geoscience and Remote Sensing* (2019).

5. Ross, Z. E., Meier, M.-A., Hauksson, E. & Heaton, T. H. Generalized seismic phase detection with deep learning. *Bulletin of the Seismological Society of America* **108**, 2894–2901 (2018).

6. Chen, Y. Automatic microseismic event picking via unsupervised machine learning. *Geophysical Journal International* **212**, 88–102 (2017).

7. Ross, Z. E., Meier, M.-A. & Hauksson, E. P wave arrival picking and first-motion polarity determination with deep learning. *Journal of Geophysical Research: Solid Earth* **123**, 5120–5129 (2018).

8. Mousavi, S. M., Zhu, W., Ellsworth, W. & Beroza, G. Unsupervised Clustering of Seismic Signals Using Deep Convolutional Autoencoders. *IEEE Geoscience and Remote Sensing Letters* (2019).

9. Zhu, W., Mousavi, S. M. & Beroza, G. C. Seismic signal denoising and decomposition using deep neural networks. *arXiv preprint arXiv:1811.02695* (2018).

10. Zhang, C., van der Baan, M. & Chen, T. Unsupervised dictionary learning for signal-to-noise ratio enhancement of array data. *Seismological Research Letters* **90**, 573–580 (2018).

11. Mousavi, S. M., Horton, S. P., Langston, C. A. & Samei, B. Seismic features and automatic discrimination of deep and shallow induced-microearthquakes using neural network and logistic regression. *Geophysical Journal International* **207**, 29–46 (2016).

12. Nakano, M., Sugiyama, D., Hori, T., Kuwatani, T. & Tsuboi, S. Discrimination of seismic signals from earthquakes and tectonic tremor by applying a convolutional neural network to running spectral images. *Seismological Research Letters* **90**, 530–538 (2019).

13. McBrearty, I. W., Delorey, A. A. & Johnson, P. A. Pairwise association of seismic arrivals with convolutional neural networks. *Seismological Research Letters* **90**, 503–509 (2019).

14. McBrearty, I. W., Gomberg, J., Delorey, A. A. & Johnson, P. A. Earthquake Arrival Association with Backprojection and Graph TheoryEarthquake Arrival Association with Backprojection and Graph Theory. *Bulletin of the Seismological Society of America* .

15. Ross, Z. E., Yue, Y., Meier, M.-A., Hauksson, E. & Heaton, T. H. PhaseLink: A deep learning approach to seismic phase association. *Journal of Geophysical Research: Solid Earth* **124**, 856–869 (2019).

16. Lomax, A., Michelini, A. & Jozinović, D. An investigation of rapid earthquake characterization using single-station waveforms and a convolutional neural network. *Seismological Research Letters* **90**, 517–529 (2019).

17. Kriegerowski, M., Petersen, G. M., Vasyura-Bathke, H. & Ohrnberger, M. A deep convolutional neural network for localization of clustered earthquakes based on multistation full waveforms. *Seismological Research Letters* **90**, 510–516 (2018).

18. Ghahramani, Z. Probabilistic machine learning and artificial intelligence. *Nature* **521**, 452 (2015).

19. Gal, Y. & Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, 1050–1059 (2016).

20. Minson, S. E. *et al.* Crowdsourced earthquake early warning. *Science advances* **1**, e1500036 (2015).

21. Mousavi, S. M., Sheng, Y., Zhu, W. & Beroza, G. STanford EArthquake Dataset (STEAD): A Global Data Set of Seismic Signals for AI. *IEEE Access* URL 10.1109/ACCESS.2019. 2947848.

22. Ringler, A. T., Hutt, C. R., Persefield, K. & Gee, L. S. Seismic station installation orientation errors at ANSS and IRIS/USGS stations. *Seismological Research Letters* **84**, 926–931 (2013).

23. Ekström, G. & Busby, R. W. Measurements of seismometer orientation at USArray transportable array and backbone stations. *Seismological Research Letters* **79**, 554–561 (2008).

24. Hara, K., Vemulapalli, R. & Chellappa, R. Designing deep convolutional neural networks for continuous object orientation estimation. *arXiv preprint arXiv:1702.01499* (2017).

25. Roberts, R. G., Christoffersson, A. & Cassidy, F. Real-time event detection, phase identification and source location estimation using single station three-component seismic data. *Geophysical Journal International* **97**, 471–480 (1989).

26. Cichowicz, A. An automatic S-phase picker. *Bulletin of the Seismological Society of America* **83**, 180–189 (1993).

27. Magotra, N., Ahmed, N. & Chael, E. Single-station seismic event detection and location. *IEEE Transactions on Geoscience and Remote Sensing* **27**, 15–23 (1989).

28. Böse, M. *et al.* A probabilistic framework for single-station location of seismicity on Earth and Mars. *Physics of the Earth and Planetary Interiors* **262**, 48–65 (2017).

29. Abercrombie, R. E. Earthquake locations using single-station deep borehole recordings: Implications for microseismicity on the San Andreas fault in southern California. *Journal of Geophysical Research: Solid Earth* **100**, 24003–24014 (1995).

30. Kanamori, H. Quantification of earthquakes. *Nature* **271**, 411 (1978).

31. Lockman, A. B. & Allen, R. M. Single-station earthquake characterization for early warning. *Bulletin of the Seismological Society of America* **95**, 2029–2039 (2005).

32. MacKay, D. J. A practical Bayesian framework for backpropagation networks. *Neural computation* **4**, 448–472 (1992).

33. Graves, A. Practical variational inference for neural networks. In *Advances in neural information processing systems*, 2348–2356 (2011).

34. Blundell, C., Cornebise, J., Kavukcuoglu, K. & Wierstra, D. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424* (2015).

35. Gal, Y. & Ghahramani, Z. Bayesian convolutional neural networks with Bernoulli approximate variational inference. *arXiv preprint arXiv:1506.02158* (2015).

36. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research* **15**, 1929–1958 (2014).

37. Kendall, A., Badrinarayanan, V. & Cipolla, R. Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding. *arXiv preprint arXiv:1511.02680* (2015).

38. Damianou, A. & Lawrence, N. Deep gaussian processes. In *Artificial Intelligence and Statistics*, 207–215 (2013).

39. Kendall, A. & Gal, Y. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, 5574–5584 (2017).

40. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural computation* **9**, 1735–1780 (1997).

41. Chung, J., Gulcehre, C., Cho, K. & Bengio, Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555* (2014).

42. Bai, S., Kolter, J. Z. & Koltun, V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271* (2018).

43. van den Oord, A. *et al.* Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499* (2016).

44. Yu, F. & Koltun, V. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122* (2015).

45. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 770–778 (2016).

46. Nair, V. & Hinton, G. E. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, 807–814 (2010).

## Acknowledgments

## Authors' contributions

S.M.M. designed the project, networks, implemented the software, performed the training and tests, and wrote the manuscript. G.C.B. lead the project and helped with the manuscript. All authors contributed ideas to the project.

## Competing interests

The authors declare any competing interests.

Figure 1: Geograpical distribution of events used in this study.

Figure 2: Characteristics of the dataset used for training and testing.

Figure 3: The dist-PT network for estimating the epicentral distance, P travel time, and their aleatory uncertainties. A detailed description of residual dilational units are presented in the method section and figure Figure 14.

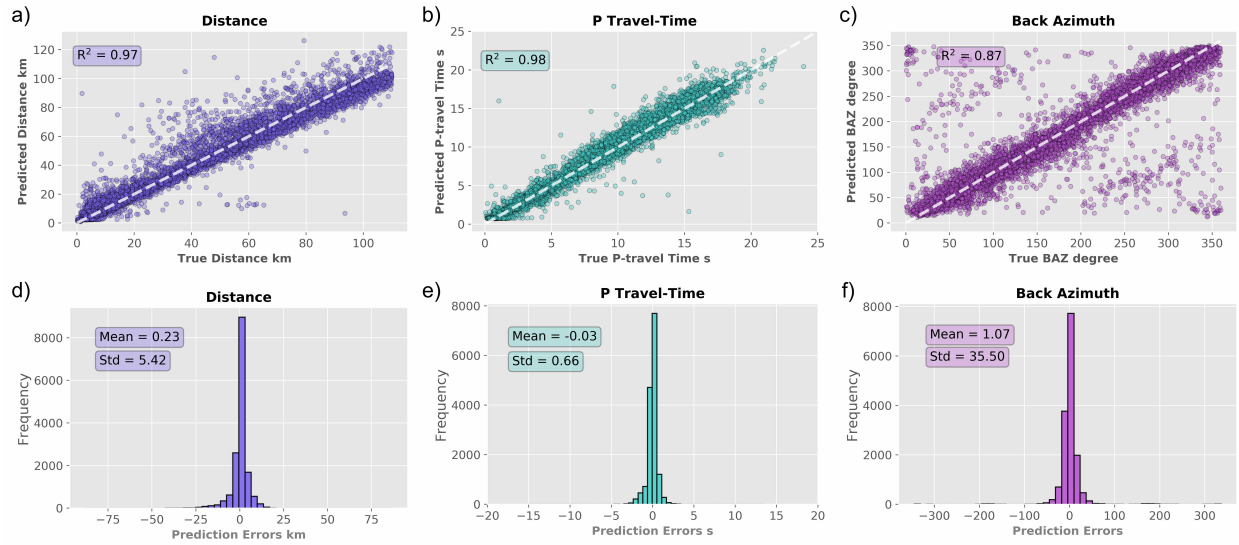Figure 4: Architecture of the BAZ network (for estimation of the back azimuth angle (BAZ)).

Figure 5: Test results for distance, P travel time, and back azimuth.

Figure 6: Prediction errors as a function of estimated uncertainties for epicentral distance and P travel time.

Figure 7: Estimated model uncertainties for test set and its relation with the prediction errors of back azimuth.

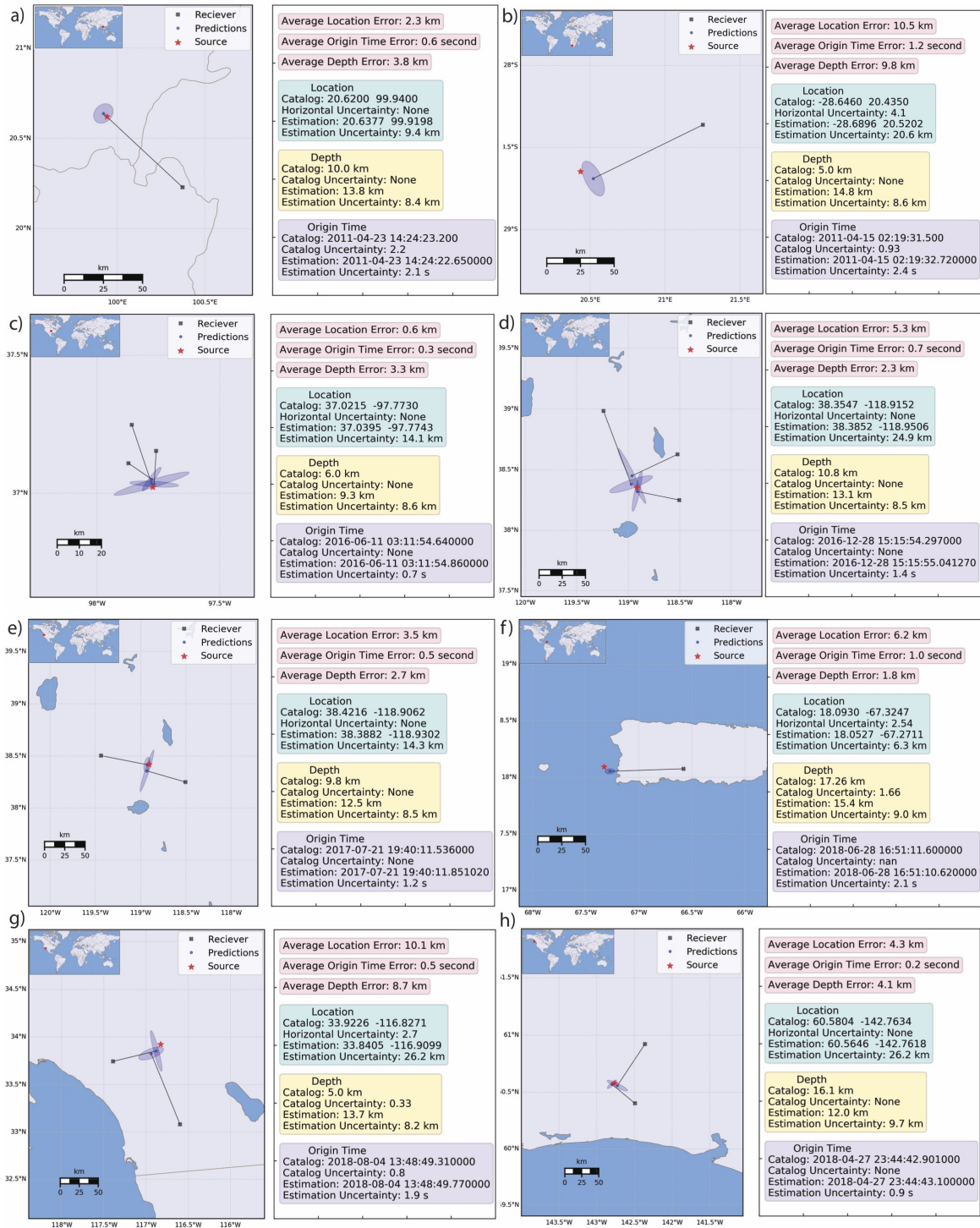Figure 8: Relations between errors and uncertainties with different characteristics of events and waveform.

Figure 9: Single-station location estimates and associated error based on predicted back-azimuth and epicentral distance. For each event, we estimate errors based on averaged results (without weighting) for multiple observations. a) is an $M_L$ 1.7 earthquake in Myanmar. b) is an $M_L$ 1.6 in south Africa. c) an $M_L$ 2.1 in Southern Kansas. d and e) are respectively an $M_L$ 2.3 and $M_L$ 1.3 south of Reno, Nevada. f) is an $m_d$ 2.19 in San Juan island. g) is an $M_L$ 3.1 northwest of Palm Springs in south California. h) is an $M_L$ 1.3 in Alaska.
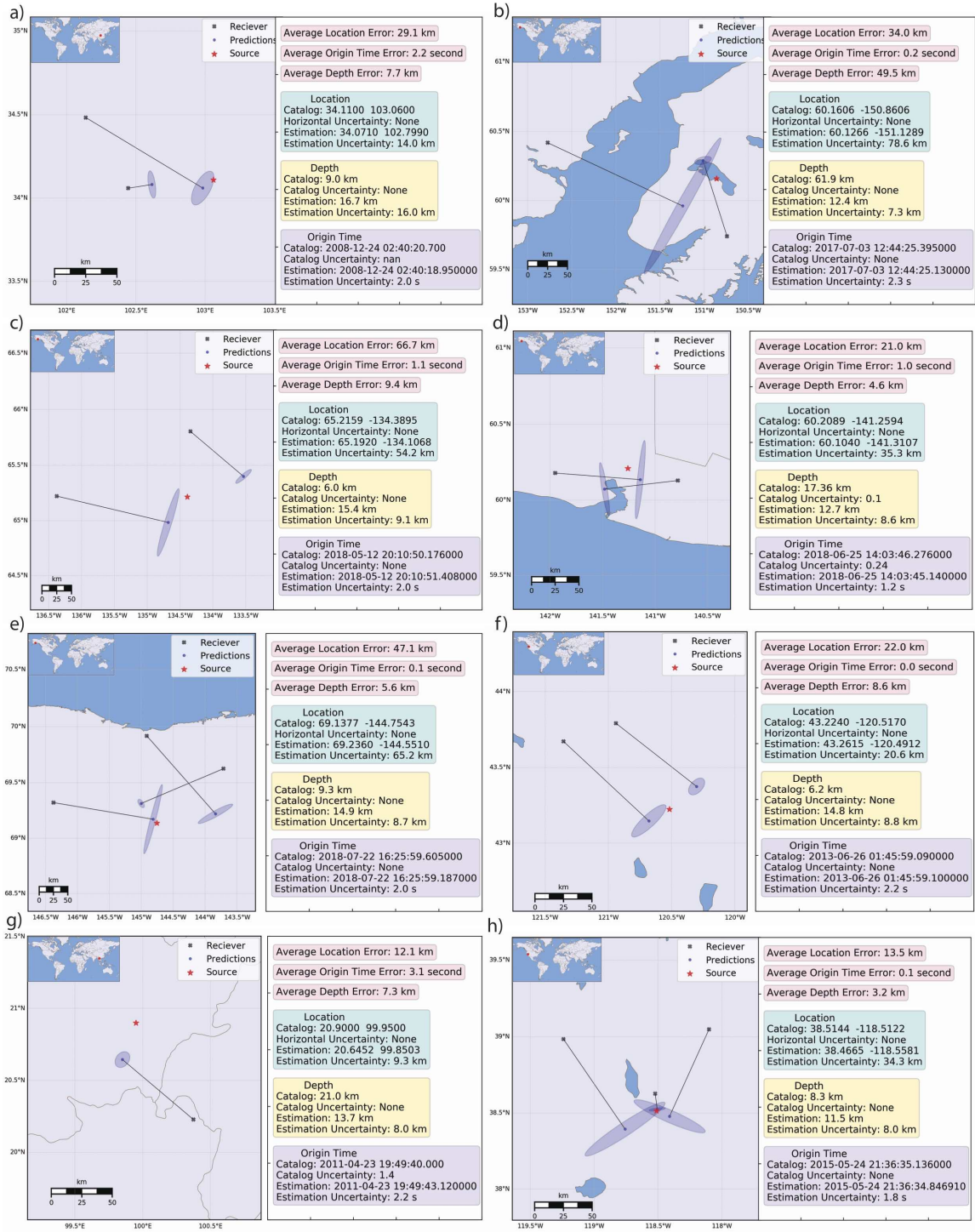
34

Figure 10: Single-station location estimates and associated error based on predicted back-azimuth and epicentral distance. For each event, errors are calculated based on averaged results (without weighting) for multiple observations. a) is an $M_L$ 2.5 event in central China, b) is an $M_L$ 1.6 in southern Alaska, c) is an $M_L$ 2.5 in northwest Canada. d) is an $M_L$ 2.1 in southeast Alaska. e) is an $M_L$ 1.8 in northern Alaska. f) is an $M_L$ 2.3 in Oregon, g) is an $M_L$ 2.2 in Myanmar, h) is an $M_L$ 1.3 in Nevada.

Figure 11: Single-station location predictions connected to their associated ground truth (locations in the catalogs) for Alaska (a) and Northern California (b).
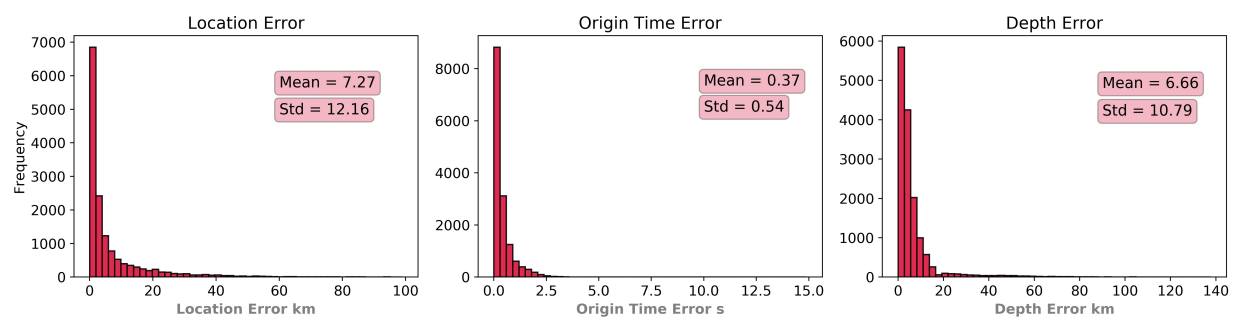
Figure 12: Statistics of location, origin time, and depth estimations for all the data (globally distributed) in the test set. Results are based on single-station estimates averaged for each event.

Figure 13: The relationships between prediction errors and event magnitude, depth, predicted uncertainties, and reported uncertainty in the catalog.
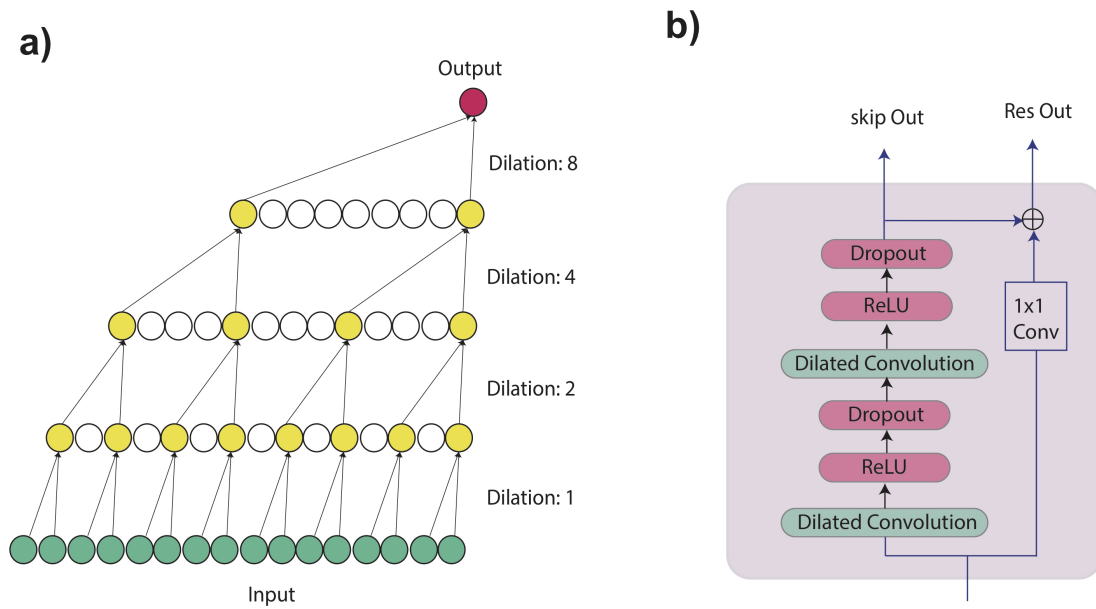
Figure 14: a) demonstration of dilational convolution operation for dilations 1, 2, 4, and 8. b) residual structure used over the dilational network.