

manuscript submitted to Reviews of Geophysics

1 **Deep Learning for Geophysics: Current and Future Trends**

2 **Siwei Yu¹ and Jianwei Ma²**

3 ¹ Center of Geophysics, Institute of Artificial Intelligence, and School of Mathematics, Harbin
4 Institute of Technology, Harbin, China.

5 ² School of Earth and Space Sciences, Peking University, Beijing, China.

6 Corresponding author: Jianwei Ma (jwm@pku.edu.cn)

7 **Key Points:**

- 8 • The concept of deep learning and classical architectures of deep neural networks are
9 introduced.
- 10 • A review of state-of-the-art deep learning methods in geophysical applications is
11 provided.
- 12 • The future directions for developing new deep learning methods in geophysics are
13 discussed.

14 **Abstract**

15 Recently, a new data-driven technique, i.e., deep learning (DL), has attracted significantly
16 increasing attention in the geophysical community. The collision of DL and traditional methods
17 has brought opportunities as well as challenges. DL was proven to have the potential to predict
18 complex system states accurately and relieve the “curse of dimensionality” in large temporal and
19 spatial geophysical applications. We address the basic concepts, state-of-the-art literature, and
20 future trends by reviewing DL approaches in various geosciences scenarios. Exploration
21 geophysics, earthquakes, and remote sensing are the main focuses. More applications, including
22 Earth structure, water resources, atmospheric science, and space science, are also reviewed.
23 Additionally, the difficulties of applying DL in the geophysical community are stressed. The
24 trends of DL in geophysics in recent years are analyzed. Several promising directions are
25 provided for future research involving DL in geophysics, such as unsupervised learning, transfer
26 learning, multimodal DL, federated learning, uncertainty estimation, and active learning. A
27 coding tutorial and a summary of tips for rapidly exploring DL are presented for beginners and
28 interested readers of geophysics.

29 **Plain Language Summary**

30 With the rapid development of artificial intelligence (AI), students and researchers in the
31 geophysical community would like to know what AI can bring to geophysical discoveries. We
32 present a review of deep learning, a popular AI technique, for geophysical readers to understand
33 recent advances, open problems, and future trends. This review aims to pave the way for more
34 geophysical researchers, students, and teachers to understand and use deep learning techniques.

35 **1 Introduction**

36 Geophysics is a discipline that uses physical principles and methods to investigate and
37 characterize the Earth, from the Earth’s core to the Earth’s surface. Modern geophysics extends
38 to outer space, from the outer layers of the Earth’s atmosphere to other planets. The general
39 methods of geophysics consist of data observation, processing, modeling, and prediction.
40 Observation is an essential means by which humans come to understand unknown geophysical
41 phenomena. Data observation uses mainly noninvasive techniques such as seismic waves,

manuscript submitted to Reviews of Geophysics

42 gravity fields, and remote sensing. Processing the recovery of clean data from raw observations
43 includes denoising, reconstruction, etc. Modeling uses mathematical and physical knowledge to
44 characterize geophysical phenomena and laws. Predictions provide the unknown based on the
45 known data and models. Spatial predictions are used to uncover the Earth's interior, such as in
46 exploration geophysics, which images the physical properties of the subsurface. Temporal
47 predictions provide the historical or future states of the Earth, such as in weather forecasting.

48 With the development of observation equipment, the amount of observed data is
49 increasing at an impressive speed. Processing, modeling and prediction with such a large amount
50 of observed data and solving bottlenecks in geophysics are significant problems. Taking
51 modeling as an example, one of the most challenging tasks in modeling is to characterize the
52 Earth with a high resolution. However, there is an unfortunate contradiction in traditional
53 methods that prevents the simultaneous achievement of both a high resolution and a wide range
54 of data observation due to hardware limitations. Therefore, it is nearly impossible to obtain a
55 high resolution model of the Earth, either spatially or temporally, since the Earth has an
56 extremely large spatial and temporal scale. An Earth system numerical simulation facility in
57 China, called EarthLab, can at most provide a resolution of 25 km for the atmosphere and 10 km
58 for oceans based on a high-performance computation device with 15 P FLOPs (floating-point
59 operations per second). Several specific difficult tasks in geophysics are listed in Table 1.

60 To illustrate the bottlenecks in processing and prediction, we use exploration geophysics
61 as an example. Exploration geophysics aims to observe Earth's subsurface or other planets with
62 data collected at the surface, such as seismic fields and gravity fields. The main process of
63 exploration geophysics includes pre-processing and imaging, where imaging means predict the
64 subsurface structures. In the geophysical signal pre-processing stage, the simplest assumption
65 regarding the shape of underground layers is that the reflective seismic records are linear in small
66 windows ([Spitz 1991](#)). Further assumptions include that the data are sparse under certain
67 transforms ([Donoho and Johnstone 1995](#)), such as the curvelet domain ([Herrmann and](#)
[Hennenfent 2008](#)) or the time-frequency domain ([Mousavi and Langston 2016](#), [Mousavi et al.](#)
[2016](#), [Mousavi and Langston 2017](#)), and that the data are low-rank after the Hankel transform
[\(Oropeza and Sacchi 2011\)](#), among others. However, the predesigned linear assumption or sparse

manuscript submitted to Reviews of Geophysics

transform assumption is not adaptive to different types of seismic data and may lead to low denoising or interpolation quality for data with complex structures. In the geophysical imaging stage, wave equations are fundamental tools to govern the kinematics and dynamics of seismic wave propagation. Acoustic, elastic, or viscoelastic wave equations introduce an increasing number of factors into the wave equations, and the generated wave field records can precisely estimate real scenarios. However, as the wave equation becomes increasingly complex, the numerical implementation of the equation becomes nontrivial, and the computational cost increases considerably for large-scale scenarios.

Different from traditional model-driven methods, machine learning (ML) is a type of data-driven approach that trains a regression or classification model through a complex nonlinear mapping with adjustable parameters based on a training dataset. The comparison of model-driven and data-driven approaches is summarized in Figure 1. For decades, ML methods have been widely adopted in various geophysical applications, such as exploration geophysics ([Poulton 2002](#), [Lim 2005](#), [Huang et al. 2006](#), [Helmy et al. 2010](#), [Zhang et al. 2014](#), [Jia and Ma 2017](#)), earthquake localization ([Mousavi et al. 2016](#)), aftershock pattern analysis ([DeVries et al. 2018](#)), and Earth system analysis ([Reichstein et al. 2019](#)). A review article about ML in solid Earth geoscience was recently published in Science ([Bergen et al. 2019](#)). The topic includes a variety of ML techniques, from traditional methods, such as logistic regression, support vector machines, random forests and neural networks, to modern methods, such as deep neural network and deep generative models. The article stresses that ML will play a key role in accelerating the understanding of the complex, interacting and multiscale processes of Earth's behavior.

In the ML community, an artificial neural network (ANN) is one such regression or classification model that is analogous to the human brain and consists of layers of neurons. An ANN with more than one layer, i.e., a deep neural network (DNN), is the core of a recently developed ML method, named deep learning (DL) ([LeCun et al. 2015](#)). DL mainly encompasses supervised and unsupervised approaches depending on whether labels are available or not, respectively. Supervised approaches train a DNN by matching the input and labels and are usually used for classification and regression tasks. Unsupervised approaches update the parameters by building a compact internal representation and then are used for clustering or

manuscript submitted to Reviews of Geophysics

100 pattern recognition. In addition, DL also contains semi-supervised learning where partial labels
101 are available and reinforcement learning where a human-designed environment provides
102 feedback for the DNN. Figure 2 summarizes the relationship from artificial intelligence to DL
103 and the classification of DL approaches. DL has shown potential in overcoming the limitations
104 of traditional approaches in various areas. The performance of DL is even superior to the
105 performance of the human brain in specific tasks, such as image classification (5.1% versus
106 3.57% with respect to the top-5 classification errors, [He et al. 2016](#)) and the game Go.

107 The geophysical community has shown a great interest in DL in recent years. Figure 3
108 show the published papers related to artificial intelligence in two major geophysical unions, i.e.,
109 society of exploration geophysics (SEG) and American geophysical union (AGU). A clear
110 exponential growth is observed in both libraries due to the use of DL techniques. Moreover, DL
111 has also provided several astonishing results to the geophysical community. For instance, on the
112 STanford EAरthquake Dataset (STEAD), the earthquake detection accuracy is improved to 100%
113 compared to 91% accuracy of the traditional STA/LTA (short time average over long time
114 average) method ([Mousavi et al. 2019](#), [Mousavi et al. 2020](#)). DL makes characterizing the earth
115 with high resolution on a large scale possible ([Chattopadhyay et al. 2020](#), [Chen et al. 2019](#),
116 [Zhang et al. 2020](#)). DL can even be used for discovering physical concepts ([Iten et al. 2020](#)).

117 Our review introduces DL-related literature covering a variety of geophysical
118 applications, from deep to the Earth's core to distant outer space, and mainly focuses on
119 exploration geophysics, earthquake science and a geophysical data observation method for
120 remote sensing. This review intends to first provide a glance at the most recent DL research
121 related to geophysics, along with an analysis of the changes and challenges DL brings to the
122 geophysical community, and then discuss the and future trends. Figure 4 gives a glance at the
123 topics included in this review. In addition, we provide a cookbook for beginners who are
124 interested in DL, from geophysical students to researchers.

125 The review part consists of three sections. The second section contains concepts, and we
126 introduce the basic idea of DL (S2). The third section review DL applications in geophysical

127 areas (S3). A discussion of future trends directions (S4) are given as extensions of this review.
128 S5 summarizes this review. A tutorial section for beginners is given in the appendix.

129 **2 The theory of deep learning**

130 Readers who are already familiar with general theory in DL may skip to Section 3. We
131 denote scalars by italic letters, vectors by bold lowercase letters and matrices by bold uppercase
132 letters. In geophysics, a large number of regression or classification tasks can be reduced to,

$$\mathbf{y} = \mathbf{L}\mathbf{x}, \quad (1)$$

133 where \mathbf{x} stands for unknown parameters, \mathbf{y} stands for observation which we partially know, and
134 \mathbf{L} is a forward or degraded operator in geophysical data observation, such as noise contamination,
135 subsampling, or physical response. However, \mathbf{L} is usually ill-conditioned or not invertible, or
136 even not known. The inverse of \mathbf{L} is mainly approximately achieved by two routines. First, an
137 optimization objective loss function is established with an additional constraint, such as sparsity
138 constraint in dictionary learning. Second, given an extensive training set, a mapping between \mathbf{x}
139 and \mathbf{y} is established by training, as done in DL, which is especially suitable for situations where
140 \mathbf{L} is not precisely known.

141 To bring the reader into DL gradually, this paper first introduces another approach, i.e.,
142 dictionary learning ([Aharon et al. 2006](#)), since the theoretical frameworks of dictionary learning
143 and DL are similar. In dictionary learning, an adaptive dictionary is learned as a representation of
144 the target data. The key features of dictionary learning are single-level decomposition,
145 unsupervised learning, and linearity. Single-level decomposition means that one dictionary is
146 used to represent a signal. Unsupervised learning means no labels are provided during dictionary
147 learning. Besides, only the target data are used without an extensive training set. Linearity
148 implies that the data decomposition on the dictionary is linear. The above features make the
149 theory of dictionary learning simple. This review will help readers transfer existing knowledge
150 on dictionary learning to DL.

151 2.1 Dictionary learning

152 To solve Equation (1), an optimization function $E(\mathbf{x};\mathbf{y})$ with a regularization term R is
 153 constructed:

$$E(\mathbf{x};\mathbf{y}) = D(\mathbf{Lx},\mathbf{y}) + R(\mathbf{x}) \quad (2)$$

154 where D is a similarity measurement function. Typically, the L_2 -norm $\|\mathbf{Lx} - \mathbf{y}\|_2$ is used under
 155 the assumption of Gaussian distribution for the error. Tikhonov regularization ($R(\mathbf{x}) = \|\mathbf{x}\|_2^2$) and
 156 sparsity are two popular regularization terms. In sparsity regularization, $R(\mathbf{x}) = \|\mathbf{Wx}\|_1$, where \mathbf{W}
 157 is a sparse transform with several vectorized bases. \mathbf{W} is also termed as the dictionary. The goal
 158 of dictionary learning is to train an optimized sparse transform \mathbf{W} , which is used for the sparse
 159 representation of \mathbf{x} . The objective function of dictionary learning involves learning \mathbf{W} via matrix
 160 decomposition with constraints R_w and R_v on the dictionary \mathbf{W} and coefficient \mathbf{v} ,

$$E(\mathbf{W},\mathbf{v}) = D(\mathbf{W}^T \mathbf{v}, \mathbf{x}) + R_w(\mathbf{W}) + R_v(\mathbf{v}) \quad (3)$$

161 where \mathbf{W} and \mathbf{v} are optimized alternatively, i.e., dictionary updating and sparse coding. Here we
 162 introduce two dictionary learning approaches: K-SVD and data-driven tight frame (DDTF).

163 K-SVD (where SVD is singular value decomposition) ([Aharon et al. 2006](#)) regularizes
 164 the sparsity of \mathbf{v} and normalizes the energy of \mathbf{W} . K-SVD uses orthogonal matching pursuit for
 165 sparse coding and several tricks in dictionary updating. First, one component of the dictionary is
 166 updated at a given time, and the remaining terms are fixed. Second, a rank-1 approximation SVD
 167 algorithm is used to obtain the updated dictionary and coefficients simultaneously, thereby
 168 accelerating convergence and reducing computational memory. K-SVD is applied in geophysics
 169 with extensions to improve efficiency ([Nazari Siahsar et al. 2017](#)).

170 Despite the success of K-SVD in signal enhancement and compression, dictionary
 171 updating is still time-consuming regarding high-dimensional and large-scale datasets, such as 3D
 172 prestack data in seismic exploration. K-SVD includes one SVD step to update one dictionary
 173 term. Can the entire dictionary be updated by one SVD for efficient improvement? A data-driven
 174 tight frame (DDTF) ([Cai et al. 2014](#), [Liang et al. 2014](#),) was proposed by enforcing a tight frame
 175 constraint on the dictionary \mathbf{W} . The tight frame condition is a slightly weaker condition than

manuscript submitted to Reviews of Geophysics

176 orthogonality, for which the perfect reconstruction property holds. With the tight frame property,
 177 dictionary updating in DDTF is achieved with one SVD, which is hundreds of times faster than
 178 K-SVD. DDTF has been applied in high dimensional seismic data reconstruction (Yu et al. 2015,
 179 Yu et al. 2016). An example of a learned dictionary with 3D DDTF for a seismic volume is
 180 shown in Figure 5.

181 2.2 Deep learning

182 Unlike dictionary learning, DL treats geophysical problems as classification or regression
 183 problems. A DNN F is used to approximate \mathbf{x} from \mathbf{y} ,

$$\mathbf{x} = F(\mathbf{y}; \Theta) \quad (4)$$

184 where Θ is the parameter set of the DNN. In classification tasks, \mathbf{x} is a one-hot encoded vector
 185 representing the categories. Θ is obtained by building a high-dimension approximation between
 186 two sets $\mathbf{X} = \{\mathbf{x}_i, i = 1 \dots N\}$ and $\mathbf{Y} = \{\mathbf{y}_i, i = 1 \dots N\}$, i.e., the labels and inputs. The
 187 approximation is achieved by minimizing the following loss function to obtain an optimized Θ :

$$E(\Theta; \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^N \|\mathbf{x}_i - F(\mathbf{y}_i; \Theta)\|_2^2 \quad (5)$$

188 If F is differentiable, a gradient-based method can be used to optimize Θ . However, a
 189 large Jacobi matrix is involved when calculating $\nabla_\Theta E$, making it infeasible for large-scale
 190 datasets. A back-propagation method (Rumelhart et al. 1986) is proposed to compute $\nabla_\Theta E$ and
 191 avoid calculating the Jacobi matrix. In unsupervised learning, the label \mathbf{x} is not known, such that
 192 additional constraints are required, such as making \mathbf{x} identical to \mathbf{y} .

193 The relations of DL and dictionary learning are as follows: the depth of decomposition,
 194 the amount of training data, and the nonlinear operators. Dictionary learning is usually a single-
 195 level matrix decomposition problem. A double sparsity (DS) dictionary learning was proposed
 196 to explore deep decomposition (Rubinstein et al. 2010). The motivation of DS is that the learned
 197 dictionary atoms still share several underlying sparse pattern for a generic dictionary. In other
 198 words, the dictionary is represented with a sparse coefficient matrix multiplied by a fixed

manuscript submitted to Reviews of Geophysics

199 dictionary, as in discrete cosine transform. Inspired by DS dictionary learning, can we propose
200 triple, quadruple or even centuple dictionary learning? We know cascading linear operators are
201 equivalent to a single linear operator. Therefore, using more than one fixed dictionary does not
202 improve the signal representation ability compared to that ability of one fixed dictionary if no
203 additional constraints are provided. In DL, nonlinear operators are combined in such a deep
204 structure. An ANN with one hidden layer and nonlinear operators can represent any complex
205 function with a sufficient number of hidden neurons. To fit ANN with many hidden neurons, we
206 need an extensive training set, while dictionary learning involves only one target data. To
207 compare the learned features of dictionary learning in Figure 5, the hierarchical structures of
208 filters in DL are shown in Figure 6.

209 The theory of DL can be penetrated from different angles except for dictionary learning
210 (Figure 7). DL can be treated as an ultra-high dimensional nonlinear mapping from data space to
211 the feature space or the target space, where the nonlinear mapping is represented by a DNN.
212 Therefore, DL is basically a high-dimensional nonlinear optimization problem. Recurrent neural
213 networks (RNNs) are basically a solution of the ordinary differential equation with the Euler
214 method ([Chen et al. 2018](#)). A generative adversarial network (Goodfellow et al. 2014, Creswell
215 et al. 2018) (GAN) can be interpreted by the theory of optimal transportation, since the targets of
216 GAN are mainly manifold learning and probability distribution transformation, i.e.,
217 transformation between the given white noise and the data distribution ([Lei et al. 2020](#)). RNNs
218 and GANs are two specific DNNs and will be introduced in the next subsection.

219 2.3 Deep neural network architectures

220 The key components of DL are the training set, network architectures and parameter
221 optimization. The architectures of DNNs vary in different applications; here, we introduce
222 several commonly used architectures.

223 A fully connected neural network (FCNN) (Figure 8a) is an ANN composed of fully
224 connected layers where the inputs of one layer are connected to every unit in the next layer. The
225 weighted summation of the inputs passes through a nonlinear activation function f in one unit.
226 The typical f in DL are rectified linear unit (ReLU), sigmoid and tanh functions, as shown in

manuscript submitted to Reviews of Geophysics

227 Figure 9a. The number of layers in a FCNN has a significant effect on the fitting and
228 generalization abilities of the model. However, FCNNs were restricted to a few layers due to the
229 computational capacity of the available hardware, the vanishing and explosion gradient problem
230 during optimization, etc. With the development of hardware and optimization algorithms, ANNs
231 tend to become deeper. On the other hand, if a raw dataset is the input directly into the FCNN,
232 massive parameters are required since each pixel corresponding to one feature, especially for
233 high dimensional inputs. FCNN requires preselected features as inputs into the neural network
234 with full reliance on experience and ignores the structure of the input entirely. Automated feature
235 selection algorithms are proposed ([Qi et al. 2020](#)), but require high computational resources. To
236 reduce the number of parameters in an FCNN and consider local coherency in an image,
237 convolutional neural networks (CNN) (Figure 8b) were proposed to share network parameters
238 with convolutional filters.

239 CNNs have developed rapidly since 2010 for image classification and segmentation, and
240 several popular CNNs include VGGNet ([Simonyan and Zisserman 2015](#)) and AlexNet
241 ([Krizhevsky et al. 2017](#)). CNNs are also used in image denoising ([Zhang et al. 2017](#)) and super-
242 resolution tasks ([Dong et al. 2014](#)). A CNN uses original data rather than selected features as an
243 input set and use convolutional filters to restrict the inputs of a neural network to within a local
244 range. The convolutional filters are shared by different neurons in the same layer. As shown in
245 Figure 9b, one typical block in CNN consist of one convolutional layer, one nonlinear layer, one
246 batch normalization and one pooling layer. Convolutional layers and nonlinear layers provide the
247 basis components of CNN. Batch normalization layers prevent gradient explosion and make
248 stabilize the training. Pooling layers subsamples the input to extract key features. The simplest
249 CNNs are named as vanilla CNNs, which are CNNs with simple sequential structures (the same
250 for vanilla FCNN). Vanilla CNNs are reliable for most applications in geophysics, such as
251 denoising, interpolation, velocity modeling, and data interpretation, if many training samples and
252 labels are available.

253 More DL network architectures have been proposed for specific tasks based on vanilla
254 FCNNs or CNNs. A deep convolutional autoencoder (CAE, Figure 8c) is a type of CNN
255 consisting of an encoder and a decoder. The encoder uses convolutional layers and pooling

manuscript submitted to Reviews of Geophysics

256 layers to extract critical features in a latent space from the inputs, resulting in a contracting path.
257 The decoder uses deconvolutional layers and unpooling layers to decode the features into the
258 original data space, resulting in an expanding path. Here deconvolution and unpooling are
259 transpose operators corresponding to convolution and pooling. In a generalized CAE, the middle
260 of the network can also have larger dimension than the two ends. If the outputs are the same as
261 the inputs, a CAE works in an unsupervised way, and the latent features are used for other tasks,
262 such as clustering. The learned latent features can also be used for dimension reduction in large-
263 scale tasks. If labels are provided as outputs, the network architecture of CAE can also work in a
264 supervised way.

265 U-Nets ([Ronneberger et al. 2015](#)) (Figure 8d) have U-shaped structures and skip
266 connections. The skip connections bring low-level features to high levels. U-Net was first
267 proposed for image segmentation and has been applied in seismic data processing, inversion, and
268 interpretation. The U-shape structure with a contracting path and expanding path makes every
269 data point in the output contain all information from the input, such that the approach is suitable
270 for mapping data in different domains, such as inverting velocity from seismic records. The input
271 size of the test set must be the same as that in the training set for a trained U-Net. The data need
272 processed patch-wisely if the size is not identical to the requirement of U-Net.

273 A GAN (Figure 5e) can be applied in adversarial training with one generator to produce a
274 fake image or any other type of data and one discriminator to distinguish the produced one from
275 the real ones. When training the discriminator, the real dataset and generated dataset correspond
276 to labels one and zero, respectively. Additionally, when the generator is trained, all datasets
277 correspond to the label one. Such a game will finally allow the generative network to produce
278 fake images that the discriminative network cannot distinguish from real images. A GAN is used
279 to generate samples with similar distributions as the training set. The generated samples are used
280 for simulating realistic scenarios or expanding the training set. An extended GAN, named
281 CycleGAN, was proposed with two generators and two discriminators for signal processing ([Zhu](#)
282 [et al. 2017](#)). In CycleGAN, a two-way mapping is trained for mapping two datasets from one to
283 the other. The training set of CycleGAN is not necessarily paired as in a vanilla CNN, which
284 makes it relatively easy to construct training sets in geophysical applications.

manuscript submitted to Reviews of Geophysics

RNNs (Figure 8f) are commonly used for tasks related to sequential data, where the current state depends on the history of inputs fed into the neural network. Long short-term memory (LSTM) ([Hochreiter and Schmidhuber 1997](#)) is a widely used RNN that considers how much historical information is forgotten or remembered. LSTM can reduce the vanishing gradient problem, such that training on longer sequences is possible. Therefore, the inference accuracy of LSTM increases with the amount of historical information considered. In geophysical applications, RNNs are mainly used for predicting the next sample of a temporally or spatially sequenced dataset. RNNs are also used for seismic wavefield or earthquake signal modeling by simulating the time-dependent discrete partial differential equation.

3 DL geophysical applications

The most direct method for applying DL in geophysics is transferring geophysical tasks to computer vision tasks, such as denoising or classification. However, in certain geophysics applications, the characteristics of geophysical tasks or data are quite different from those of computer vision. For example, in geophysics, we have large-scale and high-dimensional data but fewer annotated labels. In this section, we introduce how DL approaches relieve the bottlenecks of traditional methods, what difficulties we encounter and how to solve them. The development of DL applications in exploration geophysics is first reviewed, followed by applications in earthquake science, remote sensing and other areas.

3.1 Exploration geophysics

Exploration geophysics images the Earth's subsurface by inverting collocated physical fields at the surface, among which seismic wavefields are the most commonly used. Seismic exploration uses reflective seismic waves to predict subsurface structures. The main processes of seismic exploration consist of seismic data sampling and processing (denoising, interpolation, etc.), inversion (migration, imaging, etc.), and interpretation (fault detection, facies classification, etc.). Figure 10 summarizes the procedure of exploration geophysics. Figure 11 compares traditional and DL-based methods in exploration geophysics.

311 3.1.1 Seismic data processing

312 Seismic data are contaminated by different types of noise, such as random noise from the
313 background, ground rolls that travel along the surface with high energy and mask useful signals,
314 and multiple that reflected multi-times between the interfaces. One of the long-standing
315 problems in exploration geophysics is to remove noise and improve the signal-to-noise ratio
316 (SNR) of signals. Traditional methods use handcrafted filters or regularization for denoising
317 certain kinds of noise by analyzing the corresponding features ([Herrmann and Hennenfent 2008](#)).
318 However, handcrafted filters fail when the signal and noise share a common feature space. DL
319 methods avoid feature selection when used for seismic denoising. For example, U-Net-based
320 DeepDenoiser can separate signals and noise by learning a nonlinear regression ([Zhu et al. 2019](#)).
321 Moreover, with DnCNN ([Zhang et al. 2017](#)), a CNN for denoising, the same architecture can be
322 used for three kinds of seismic noise while achieving a high SNR ([Yu et al. 2019](#)) as long as a
323 corresponding training set is constructed. However, there is still a long way to go. A DNN
324 trained on synthetic datasets does not have a good generalization ability to field data. To make
325 the network reusable, transfer learning ([Donahue et al. 2014](#)) can be used for field data denoising.
326 Sometimes the labels of clean data are difficult to obtain, and one solution is to use multiple
327 trials involving user-generated white noise to simulate real white noise ([Wu et al. 2019](#)).

328 An example of scattered ground-roll attenuation is shown in Figure 12 ([Yu et al. 2019](#)).
329 Scattered ground roll is mainly observed in the desert area, and is caused by the scattering of
330 ground roll when the near surface is laterally heterogeneous. The scattered ground roll is difficult
331 to remove because it occupies the same frequency domain as the reflected signals. DnCNN was
332 used to remove scattered ground roll successfully.

333 Due to environmental or economic limitations, seismic geophones are usually located
334 irregularly or not densely enough under the principle of Nyquist sampling. The reconstruction or
335 regularization of seismic data to a dense and regular grid is essential to improve inversion
336 resolution. In the beginning, end-to-end DNNs were proposed for the reconstruction of regularly
337 missing data ([Wang et al. 2019](#)) and randomly missing data ([Wang et al. 2020](#), [Mandelli et al.](#)
338 [2018](#)). However, the training sets are numerically synthetic, and do not generalize well to field

manuscript submitted to Reviews of Geophysics

339 data. We can borrow training data from a natural image dataset to train DnCNN and then embed
340 it in the traditional project onto a convex set (POCS, [Abma and Kabir 2006](#)) framework ([Zhang](#)
341 [et al. 2020](#)). The resulting interpolation algorithm generalized well to seismic data. Moreover, no
342 new networks were required for the interpolation of other datasets. Figure 13 gives the training
343 set and a simple interpolation result ([Zhang et al. 2020](#)).

344 First arrival picking is used to select the first jumps of useful signals and has been
345 automated but needs intense human intervention to check pickings with significant static
346 corrections, weak energy, low signal-to-noise ratios, and dramatic phase changes. DL helps
347 improve the automation and accuracy of first arrival picking on realistic seismic data. It is natural
348 to transform first arrival picking into a classification problem by setting the first arrival as ones
349 and other locations as zeros when DL is used ([Hu et al. 2019](#)). However, such a setting can cause
350 imbalanced labels. An interesting approach treats first arrival picking as an image classification
351 problem, where anything before the first arrival is set to zero, and all instances after the first
352 arrival are set to one ([Wu et al. 2019](#)). This method works well for noisy situations and field
353 datasets. After the segmentation image is obtained, a more advanced picking algorithm, such as
354 an RNN, can be applied to take advantage of the global information ([Yuan et al. 2020](#)).

355 Figure 14 shows the results of the first arrival picking based on U-Net. We used 8000
356 synthetic seismological samples. A gradient constraint was added to the loss function to enhance
357 the continuity of the selected positions. For the output, three classifications were set: zeros before
358 the first arrival, ones after the first arrival, and twos for the first arrival. The training dataset was
359 contaminated with strong noise and had missing traces. The predicted picking results were close
360 to the labels.

361 More DL-based seismic signal processing literature that does not belong to the mentioned
362 scope is summarized in this paragraph. Signal compression is essential for the storage and
363 transmission of seismic data. Traditional seismic data are stored in 32 bits per sample. With an
364 RNN to estimate the relationships among samples in a seismic trace and compress seismic data,
365 only 16 bits are needed for lossless representations, such that half storage is saved ([Payani et al.](#)
366 [2019](#)). Seismic registration aligns seismic images for tasks such as time-lapse studies. However,

manuscript submitted to Reviews of Geophysics

367 when large shifts and rapid changes exist, this task is extremely difficult. A CNN is trained with
368 two seismic images as inputs and the shift as output by learning from the concept of optical flow.
369 The method outperforms traditional methods but is dependent on the training dataset (Dhara and
370 Bagaini 2020).

371 3.1.2 Seismic data imaging

372 Seismic imaging is a challenging problem since traditional methods such as tomography
373 and full waveform imaging (FWI) suffer from several bottlenecks. 1. Imaging is time-consuming
374 due to the curse of dimensionality. 2. Imaging relies heavily on human interactions to select
375 proper velocities. 3. Nonlinear optimization needs a good initialization or low frequency
376 information, however there is a lack of low frequency energy in recorded data. DL methods help
377 relieve the bottlenecks from several angles.

378 First, end-to-end DL-based imaging methods use recorded data as inputs and velocity
379 models as outputs, which provides a totally different imaging approach. DL methods avoid the
380 mentioned bottlenecks, providing a next-generation imaging method. The first attempts at DL in
381 staking (Park and Sacchi 2019), tomography (Araya-Polo et al. 2018) and FWI (Yang and Ma
382 2019) show promising results on synthetic 2D data. One important issue is that the input is in the
383 data space and the output is in the model space, both with high dimensional parameters. U-Net is
384 used to transfer from different spaces with different dimensions, and downsampling is used to
385 reduce the parameters while training the DNN (Yang and Ma 2019). Figure 15 shows the
386 velocity inversion results from (Yang and Ma 2019).

387 However, end-to-end DL imaging also has disadvantages, such as a lack of training
388 samples and restricted input sizes due to memory limitations. An interesting work used smoothed
389 natural images as velocity models, thus producing a large number of models to construct the
390 training set (Wang and Ma 2020). Figure 16 shows how (Wang and Ma 2020) convert a three-
391 channel color image to a velocity model.

392 To make DL-based imaging applicable to large scale inputs, more works aim to
393 collaborate with traditional methods and aim to solve one of the mentioned bottlenecks, such as
394 extrapolating the frequency range of seismic data from high to low frequencies for FWI

manuscript submitted to Reviews of Geophysics

395 (Ovcharenko et al. 2019, Fang et al. 2020), and adding constraints to FWI (Zhang and Alkhalifah
396 2019). To mitigate the “curse of dimensionality” problem of global optimization in FWI, CAE is
397 used to reduce the dimension of FWI by optimizing in the latent space (Gao et al. 2019). Another
398 work aims at the high computational cost of forward modeling when the high-order finite
399 difference method is used. A GAN is used to produce a high-quality wavefield from a low-
400 quality wavefield with a lower-order finite difference in the context of surface-related multiples,
401 ghosts, and dispersion (Siahkoohi et al. 2019). U-Net can be used for velocity picking in stacking
402 (Figure 17). The inputs are seismological data, and the outputs have values of one where the
403 picks are located and values of zero elsewhere.

404 An alternative is to replace the FWI object with an RNN loss function. The structure of
405 an RNN is similar to that of finite different time evolution, and the network parameters
406 correspond to the selected velocity model. Therefore, optimizing an RNN is equivalent to
407 optimizing FWI (Sun et al. 2020). Such a strategy is extended to the simultaneous inversion of
408 velocity and density (Liu 2020). Figure 18 shows the structure of a modified RNN-based on the
409 acoustic wave equation used in (Liu 2020). The diagram represents the discretized wave equation
410 implemented in an RNN with a flow chart. The optimized method in FWI can also be learned by
411 a DNN rather than with a gradient-descent-based approach (Sun and Alkhalifah 2020). An ML-
412 descent method is proposed to consider the historical information of the gradient based on an
413 RNN rather than handcrafted features.

414 3.1.3 Seismic data interpretation and attributes analysis

415 Seismic interpretation (faults, layers, dips, etc.) or attribute analysis (impedance,
416 frequency, facies, etc.) can be used to help the extraction of subsurface geologic information and
417 locate underground sweet points. However, both tasks are time-consuming since interventions by
418 experts are required. Preliminary works show that DL has the potential to improve the efficiency
419 and accuracy in seismic interpretation or attribute analysis.

420 The localization of faults, layers, and dips in seismic interpretation is similar to object
421 detection in computer vision. Therefore, DNNs for image detection can be directly applied in
422 seismic interpretation. However, unlike the computer vision industry, it is difficult to obtain a

manuscript submitted to Reviews of Geophysics

423 public training set or to manually construct a training set for field datasets. Building realistic
424 synthetic datasets rather than handcrafted field datasets is more efficient and can produce similar
425 results. Therefore, synthetic samples are used for training. To build an approximately realistic 3D
426 training dataset, randomly choosing folding and faulting parameters in a reasonable range is
427 required ([Wu et al. 2020](#)). Then, the dataset is used to train a 3D U-Net for the seismic structural
428 interpretation of features, such as faults, layers, and dips, in field datasets. If the detected objects
429 are of a small proportion, a class-balanced binary cross-entropy loss function is used to adjust the
430 data imbalance so that the network is not trained to predict only zeros ([Wu et al. 2019](#)). An
431 alternative to a synthetic training set is a semi-automated approach that annotates the targets on a
432 coarse scale and predicts them on a fine scale ([Wu et al. 2019](#)). An example of synthetic post-
433 stack image and field data fault analysis is shown in Figure 19 ([Wu et al. 2020](#)).

434 Attribute analysis is similar to image classification, where seismic images are inputs and
435 areas with labels as different attributes are output. Therefore, DNNs for image classification can
436 be directly applied in seismic attribute analysis ([Das et al. 2019](#), [You et al. 2020](#), [Feng et al.](#)
437 [2020](#)). If the attributes cannot be directly computed from the seismic data, a DNN can work in a
438 cascaded way ([Das and Mukerji 2020](#)). If labels are not available, CAE is used for feature
439 extraction, and then a clustering method, such as K-means, is used for unsupervised clustering
440 ([Duan et al. 2019](#), [He et al. 2018](#), [Qian et al. 2018](#)). Clustering refers to grouping similar
441 attributes in an unsupervised manner. For example, we can use clustering to decide whether a
442 region contains fluvial facies or faults based on stacked sections. CAE and K-means can further
443 be optimized simultaneously for better feature extraction ([Mousavi et al. 2019](#)). To mitigate the
444 dependence of vanilla CNNs on the amount of labeled seismic data available, a 1D CycleGAN-
445 based algorithm was proposed for impedance inversion ([Wang et al. 2019](#)). The CycleGAN did
446 not require training set pairing. Only two sets with and without high fidelity are needed. To
447 consider the spatial continuity and similarity of adjacent traces, an RNN is used in facies analysis
448 ([Li et al. 2019](#)).

449 3.2 Earthquake science

450 The goal of earthquake data processing is quite different from that of exploration
451 geophysics; therefore, this section focuses on DL-based earthquake signal processing. The
452 preliminary processing of earthquake signals includes classification to distinguish real
453 earthquakes from noise and arrival picking to identify the arrival times of primary (P) and
454 secondary (S) waves. Further applications involve earthquake location and Earth tomography.
455 DL has shown promising results in these applications.

456 3.2.1 Earthquake and noise classification

457 Earthquake signal and noise classification is the most fundamental and difficult task in
458 earthquake early warning (EEW). Traditional EEW systems suffer from false and missed alerts.
459 DNN can be directly applied in signal and noise discrimination since it is a classification task.
460 With a sufficient training set, DNN can achieve up to 99.2% ([Li et al. 2018](#)) and 99.5% precision
461 ([Meier et al. 2019](#)) in different regions. To detect small and weak earthquake signals robust to
462 strong noise and non-earthquake signals, a residual network with convolutional and recurrent
463 units is developed ([Mousavi et al. 2019](#)). RNN and CNN are also used in a more challenging task
464 to distinguish between anthropogenic sources, such as mining or quarry blasts, and tectonic
465 seismicity ([Linville et al. 2019](#)). More categories of signals are required to identify in specific
466 tasks, such as in volcano seismic detection. Seismic signals can be used to detect six classes:
467 long-period events, volcanic tremors, volcano-tectonic events, explosions, hybrid events, and
468 tornados ([Malfante et al. 2018](#)). Uncertainty is also considered in volcano-seismic monitoring
469 ([Bueno et al. 2019](#)).

470 We provide an example of using the wavelet scattering transform (WST) ([Mallat 2012](#))
471 and a support vector machine for earthquake classification with a limited number of training
472 samples. The WST involves a cascade of wavelet transforms, a module operator, and an
473 averaging operator, corresponding to convolutional filters, a nonlinear operator, and a pooling
474 operator in a CNN, respectively. The critical difference between the WST and a CNN is that the
475 filters are predesigned with the wavelet transform in the WST. In our case, only 100 records
476 were used for training, and 2000 records were used for testing. We obtained a classification

manuscript submitted to Reviews of Geophysics

477 accuracy as high as 93% with the WST method. Figure 20 shows the architecture of the WST
478 algorithm.

479 3.2.2 Arrival picking

480 Arrival picking for earthquakes identifies the arrival time of P and S waves. Traditional
481 automated arrival picking algorithms, such as short-term average/long-term average method
482 (STA/LTA), are less precise than human experts and rely on thresholding setting. DL-based
483 arrival picking overcomes these shortcomings and helps illuminate the Earth structure clearly
484 ([Wang et al. 2019](#)). With a sufficiently large training set, one can achieve remarkably high
485 picking and classification accuracies higher than STA/LTA ([Zhao et al. 2019](#), [Zhou et al. 2019](#)),
486 even close to or better than human experts ([Ross et al. 2018](#), 19.4 million seismograms training
487 set). If labels are not sufficient, a GAN-based model EarthquakeGen can be used to artificially
488 expand labeled data sets ([Wang et al. 2019](#)). The detection accuracy was greatly improved by
489 performing artificial sampling for the training set. Simultaneous earthquake detection and phase
490 picking can further improve the accuracy of both tasks ([Zhou et al. 2019](#), [Mousavi et al. 2020](#)).

491 3.2.3 Earthquake location and other applications

492 Earthquake location and magnitudes estimation are important in EEW and subsurface
493 imaging. Conventional earthquake location significantly relies on a velocity model and suffer
494 from inaccurate phase picking. CNN is used for earthquake location by using received
495 waveforms at several stations as input and location map as output ([Zhang et al. 2020](#)). This
496 method worked well for earthquakes ($M_L < 3.0$) with low SNRs, for which traditional methods
497 fail. The prediction results and errors of earthquake source locations are indicated in Figure 21.
498 DL also help estimate earthquake locations and magnitudes based on signals from a single
499 station ([Mousavi and Beroza 2020](#), [Mousavi and Beroza 2020](#)). Further applications involving
500 associating seismic phases, which involves grouping the phase picks on multiple stations
501 associated with an individual event ([Ross et al. 2019](#)) and relationship analysis between a strong
502 earthquake and postseismic deformation ([Yamaga and Mitsui 2019](#)).

503 3.3 Remote sensing – a geophysical data observation means

504 Remote sensing is an important means to collect geophysical data and images by using
505 sensors in satellites or aerial crafts. Remote sensing imagery mainly includes optical images,
506 hyperspectral images, and synthetic aperture radar (SAR) images. Large-scale and high-
507 resolution satellite optical color imagery can be used for precision agriculture and urban
508 planning. To address the issue of objection rotation variations, a rotation-invariant CNN for
509 object detection in very high-resolution optical remote sensing images was proposed, where a
510 rotation-invariant layer was introduced by enforcing the training samples before and after
511 rotation to share the same features ([Cheng et al. 2016](#)). If the labels are not accurate, a two-step
512 training approach was used where first the CNN was initialized by numerous inaccurate
513 reference data and then refined on a small amount of correctly labeled data ([Maggiori et al.](#)
514 [2017](#)). To further improve the image resolution, the image contours were extracted with an edge-
515 enhancement GAN to remove the artifacts and noise in super resolution ([Jiang et al. 2019](#)).

516 Images obtained by hyperspectral sensors have rich spectral information, such that
517 different land cover categories can potentially be precisely differentiated. In recent years,
518 numerous works have explored DL methods for hyperspectral image classification ([Li et al.](#)
519 [2019](#)). To consider the spectral-spatial structure simultaneously, a 3D CNN rather than a 2D one
520 should be used to extract the effective features of hyperspectral imagery ([Chen et al. 2016](#)). The
521 extracted features are useful for image classification and target detection and open a new window
522 for future research. An alternative means to explore the relationships among different spectrum
523 channels is to use RNN, which regards hyperspectral pixels as sequential data input ([Mou et al.](#)
524 [2017](#)).

525 SAR systems artificially enlarge the aperture of radar to produce high-resolution images.
526 SAR can operate in all-weather and day-and-night conditions. CNN is used for target
527 classification in SAR images, which avoided handcrafted features and provided higher accuracy
528 ([Chen et al. 2016](#)). To consider both the amplitude and phase information of complex SAR
529 imagery, a complex-valued CNN for SAR image classification was proposed to process
530 complex-valued inputs ([Zhang et al. 2017](#)).

manuscript submitted to Reviews of Geophysics

531 3.4 Other AI geophysical applications

532 We investigate more AI geophysical applications in this section. The topics are roughly
533 arranged by the order from the Earth to outer space.

534 3.4.1 The Earth's structure

535 Understanding the structure of the Earth is a challenging task since observations are
536 mainly limited on the earth's surface. The earth is roughly divided into the surface, crustal layers,
537 mantle and core and from the surface to inside; however, the detailed structures and properties of
538 the earth are not clear. An important soil attribute, moisture, is predicted historically with high
539 fidelity from two recent years of satellite data, showing LSTM's potential for hindcasting, data
540 assimilation, and weather forecasting ([Fang et al. 2017](#), [Fang et al. 2020](#)). The high-resolution
541 3D CT data of rocks is required to determine the rock's property but results in a small field of
542 view. A CycleGAN was proposed to obtain super resolution images from low resolution one by
543 training on an unpaired dataset ([Niu et al. 2020](#)). Volcanic deformation was detected by using a
544 CNN to classify interferometric fringes in wrapped interferograms ([Anantrasirichai et al. 2018](#)).
545 The crustal thickness in eastern Tibet and the western Yangtze craton are estimated by Rayleigh
546 surface wave velocities based on DNN ([Cheng et al. 2019](#)). The mantle thermal state of
547 simplified model planets was predicted based on DL with an accuracy of 99% for both the mean
548 mantle temperature and the mean surface heat flux compared to the calculated values ([Shahnas
549 and Pysklywec 2020](#)).

550 3.4.2 Water resources

551 Water on Earth has a great impact on ecosystems and natural disasters. DL can help
552 address several major challenges in water sciences ([Shen 2018](#)). DL can predict the loop current
553 in the ocean by learning the pattern in sea surface height (SSH). An LSTM was proposed to
554 predict SSH and current loop in the Gulf of Mexico within 40 kilometers nine weeks in advance
555 ([Wang et al. 2019](#)). Due to the limit of memory, the region of interest is split into different sub-
556 regions. Further works directly reconstruct SSH on a large and spatial and temporal space based
557 on sparsely sampled data with CNN ([Manucharyan et al. 2021](#)). By using observation from

manuscript submitted to Reviews of Geophysics

558 satellite and coastal stations simultaneously, GAN can be used to reconstruct the SSH of the
559 whole North-Sea ([Zhang et al. 2020](#)). DL also help estimate the iceberg in the pan-Antarctic
560 near-coastal zone that covers the whole Antarctic continent for monitoring ice melt and sea level
561 increasing ([Barbat et al. 2019](#)), and coastal inundation for a better understanding of the
562 geospatial and temporal characteristics of coastal flooding ([Liu et al. 2019](#)).

563 In addition to oceans, water is stored in different forms, such as rivers, lakes, rain, and
564 snow. DL has found its roles in estimating groundwater storage ([Sun et al. 2019](#)), global water
565 storage in the US ([Sun et al. 2020](#)), measuring accurate river widths by super resolution ([Ling et](#)
566 [al. 2019](#)), predicting the temperature of lake water ([Read et al. 2019](#)), predicting rainfall and
567 runoff ([Akbari Asanjan et al. 2018](#)), and prediction water vapor retrieval from remote sensing
568 data ([Acito et al. 2020](#)).

569 3.4.3 Atmospheric science

570 Atmospheric science observes and predicts climate, weather and atmospheric
571 phenomena. Global observation of global atmospheric parameters is difficult since the earth is
572 extremely large and sensor locations are limited. Researchers chose a CNN-based inpainting
573 algorithm to reconstruct missing values in global climate datasets such as HadCRUT4 ([Kadow et](#)
574 [al. 2020](#), Figure 22). Air pollution is damaging both the earth's environment and human health.
575 Researchers used DL to estimate ground-level PM2.5 or PM10 levels by using satellite
576 observations and station measurements ([Li et al. 2017](#), [Shen et al. 2018](#), [Tang et al. 2018](#)). DL
577 also helps improve the accuracy of weather forecasting, which is a long-standing challenge in
578 atmospheric science ([Scher and Messori , Bonavita and Laloyaux 2020](#)). The tracks of typhoons
579 were predicted with a GAN based on satellite images ([Rüttgers et al. 2019](#)). A six-hour-advance
580 track with an average error of 95.6 km was produced. Flow-dependent typhoon-induced sea
581 surface temperature cooling was estimated by a DNN and used for improving typhoon
582 predictions ([Jiang et al. 2018](#)).

manuscript submitted to Reviews of Geophysics

583 3.4.4 Space science

584 Global space parameter estimation and prediction are long-standing tasks in space
585 science. Researchers used a DNN to predict short-term and long-term 3D dynamic electron
586 densities in the inner magnetosphere ([Chu et al. 2017](#)). This network can obtain the
587 magnetospheric plasma density at any time and for any location. A regularized GAN is used to
588 reconstruct dynamic total electron content (TEC) maps ([Chen et al. 2019](#)). Several existing maps
589 were used as references to interpolate missing values in some regions, such as the oceans. The
590 TEC maps can also be predicted two hours in advance with an LSTM ([Liu et al. 2020](#)) or one
591 day in advance with a GAN ([Lee et al. 2021](#)). Further, a DNN is used to estimate the relationship
592 between electron temperature and electron density in small regions ([Hu et al. 2020](#)). Therefore,
593 the global electron density is easily measured and used to predict the global electron temperature.
594 The geomagnetic storm can be predicted with LSTM with uncertainty estimation ([Tasistro - Hart
595 et al. 2020](#)), providing confidence in the output.

596 An aurora is an astronomical phenomenon commonly observed in polar areas. Auroras
597 are caused by disturbances in the magnetosphere caused by the solar wind. Auroral classification
598 is important for polar and solar wind research. Researchers used DNN to classify auroral images
599 ([Clausen and Nickisch 2018](#), Figure 23). The classification results can further be used to produce
600 an auroral occurrence distribution ([Zhong et al. 2020](#)). To handle the situation where limited
601 images were annotated, a CycleGAN model was used to extract key local structures from all-sky
602 auroral images ([Yang et al. 2019](#)).

603 **4 Future trends directions for deep learning in geophysics**

604 4.1 The development trends of DL in geophysics

605 The landmark achievements of DL appeared after 2015, such as VGGNet ([Simonyan and
606 Zisserman 2015](#)), ResNet ([He et al. 2016](#)), AlexNet ([Krizhevsky et al. 2017](#)) and AlphaGo in
607 2016. The first attempts to apply DL in subjects related to geophysics focused on remote sensing
608 in 2016 and 2017 ([Chen et al. 2016](#), [Chen et al. 2016](#), [Maggiori et al. 2017](#), [Li et al. 2017](#)), since
609 remote sensing is a common technique widely used in many areas. In 2018 and 2019, more

manuscript submitted to Reviews of Geophysics

610 geophysical areas, such as exploration geophysics ([Araya-Polo et al. 2018](#)) and earthquake
611 studies ([Mousavi, Zhu et al. 2019](#)), started to employ DL.

612 The first attempts started with simple FCNN methods followed by complex networks,
613 such as CNN, RNN, and GAN models. With respect to the training set, early works used end-to-
614 end training borrowed from the computer vision area, which requires a large number of
615 annotated labels, while recent works have started to consider unsupervised learning ([He et al.](#)
616 [2018](#)) and the combination of DL with a physical model ([Wu and McMechan 2019](#),
617 [Chattopadhyay et al. 2020](#)). In 2020, more works are focused on the uncertainty of DL methods
618 ([Grana et al. 2020](#), [Cao et al. 2020](#), [Mousavi and Beroza 2020](#)). More examples are listed in
619 Table 2. From these trends, we can conclude that an increasing number of researchers are trying
620 to develop DL methods that are specifically designed for geophysical tasks to make DL methods
621 more practical. In the next subsection, we introduce these future trends in detail.

622 4.2 Future directions for deep learning in geophysics

623 DL, as an efficient artificial intelligence technique, is expected to discover geophysical
624 concepts and inherit expert knowledge through machine-assisted mathematical algorithms.
625 Despite the success of DL in some geophysical applications such as earthquake detectors or
626 pickers, their use as a tool for most practical geophysics is still in its infancy. The main problems
627 include a shortage of training samples, low signal-to-noise ratios, and strong nonlinearity.
628 Among these issues, the critical challenge is the lack of training samples in geophysical
629 applications compared to those in other industries. Several advanced DL methods have been
630 proposed related to this challenge, such as semi-supervised and unsupervised learning, transfer
631 learning, multimodal DL, federated learning, and active learning. We suggest that a focused be
632 placed on the subjects below for future research in the coming decade.

633 4.2.1 Semi-supervised and unsupervised learning

634 In practical geophysical applications, obtaining labels for a large dataset is time-
635 consuming and can even be infeasible. Therefore, semi-supervised or unsupervised learning is
636 required to relieve the dependence on labels. [Dunham et al. 2019](#) focused on the application of

manuscript submitted to Reviews of Geophysics

637 semi-supervised learning in a situation in which the available labels were scarce. A self-training-
638 based label propagation method was proposed, and it outperformed supervised learning methods
639 in which unlabeled samples were neglected. Semi-supervised learning takes advantage of both
640 labeled and unlabeled datasets. The combination of AE and K-means is an efficient unsupervised
641 learning method ([He et al. 2018](#), [Qian et al. 2018](#)). An autoencoder is used to learn low-
642 dimensional latent features in an unsupervised way, and then K-means is used to cluster the
643 latent features.

644 4.2.2 Transfer learning

645 Usually, we must train one DNN for a specific dataset and a specific task. For example, a
646 DNN may effectively process land data but not marine data, or a DNN may be effective in fault
647 detection but not in facies classification. Transfer learning ([Donahue et al. 2014](#)) is suggested to
648 increase the reusability of a trained network for different datasets or different tasks.

649 In transfer learning with different datasets, the optimized parameters for one dataset can
650 be used as initialization values for learning a new network with another dataset; this process is
651 called fine-tuning. Fine-tuning is typically much faster and easier than training a network with
652 randomly initialized weights from scratch. In transfer learning involving different tasks, we
653 assume that the extracted features should be the same in different tasks. Therefore, the first
654 layers in a model trained for one task are copied to the new model for another task to reduce the
655 training time. Another benefit of transfer learning is that with a small number of training samples,
656 we can promptly transfer the learned features to a new task or a new dataset. Diagrams of these
657 two transfer learning methods are shown in Figure 24. Further topics in transfer learning include
658 the relationship between the transferability of features ([Yosinski et al. 2014](#)) and the distance
659 between different tasks and different data sets ([Oquab et al. 2014](#)).

660 4.2.3 Combination of DL and traditional methods

661 Can we combine traditional and DL approaches to combine geophysical mechanics and
662 DL? Intuitively, such a combination can produce a more precise result than traditional methods
663 and a more reliable result than DL methods.

manuscript submitted to Reviews of Geophysics

How can DL be incorporated into traditional methods? In a traditional iteration optimization algorithm, the thresholding-based denoiser can be replaced by a DL denoiser (Zhang et al. 2017) such that the reconstructed results are improved. On the other hand, different tasks use the same denoiser without training a new denoiser. Another technique, DIP, uses a DNN architecture as a constraint on the data and ensembles traditional physical models for different tasks (Lempitsky et al. 2018). Similar to the idea of DIP, Wu and McMechan 2019 showed that a DNN generator can be added to an FWI framework. First, a U-Net-based generator $F(\mathbf{v}; \Theta)$ with random input \mathbf{v} was used to approximate a velocity model \mathbf{m} with high accuracy. Then, $\mathbf{m} = F(\mathbf{v}; \Theta)$ was inserted into the FWI objective function,

$$E_{\text{FWI}}(\Theta) = \frac{1}{2} \|P(F(\mathbf{v}; \Theta)) - \mathbf{d}_r\|_2^2 \quad (6)$$

where \mathbf{d}_r is the seismic record and P is the forward wavefield propagator. The gradient of E_{FWI} with respect to network parameters Θ is calculated with the chain rule. U-Net is only used for regularizing the velocity model. After training, one forward propagation of the network will produce a regularized result.

Traditional optimization methods also benefit from the autodifference mechanism in DL, which makes optimization more efficient by replacing conjugate gradient descent or LBGFS with DL optimization methods, such as SGD and Adam (Sun et al. 2020, Wang et al. 2020). DL also inspired new directions in the study of traditional nonlinear optimization algorithms, such as ML-descent (Sun and Alkhalifah 2020) and DL-based adjoint state methods (Xiao et al.).

How can traditional methods be incorporated into DL? With an additional physical constraint on DL methods, fewer training samples are required to obtain a more generalized inference than those of traditional methods. Raissi et al. 2019 proposed a physically informed neural network (PINN) that combines training data and physical equation constraints for training. Taking wave modeling as an example, the wavefield was represented with a DNN, $u(x, t) = F(x, t; \Theta)$, such that the acoustic wave equation was:

$$u_{tt} = c^2 \Delta u \xrightarrow{u(x, t) = F(x, t; \Theta)} F_{tt}(x, t; \Theta) = c^2 \Delta F(x, t; \Theta) \quad (7)$$

manuscript submitted to Reviews of Geophysics

688 How can DL and traditional methods cooperate? Another benefit of combining data-
689 driven and model-driven approaches is that we can obtain high-resolution solutions on a large
690 scale. The process on a large scale was numerically solved with a low-resolution grid based on
691 physical equations. On a small scale, the process was solved by data-driven DL methods
692 ([Chattopadhyay et al. 2020](#)). Therefore, the high computational demand on a fine scale is
693 avoided. DL can also be used for discovering physical concepts ([Iten et al. 2020](#)).

694 It is more common to hear someone ask, “Does machine learning have a real role in
695 hydrological modeling?” rather than, “What role will hydrological science play in the age of
696 machine learning?” ([Nearing et al. 2020](#)). As the authors claim, DL has uncovered the principles
697 in large-scale rainfall-runoff simulations, which cannot be explained by physical models. DL has
698 a great impact on traditional methods, causing a collision between new and old ideas. We believe
699 that DL and physical-based methods will be used together to move science forward for a long
700 time.

701 4.2.4 Multimodal deep learning

702 To improve the resolution of inversion, the joint inversion of data from different sources
703 has been a popular topic in recent years ([Garofalo et al. 2015](#)). One of the advantages of DNNs is
704 that they can fuse information from multiple inputs. In multimodal DL ([Ngiam et al. 2011](#),
705 [Ramachandram and Taylor 2017](#)), inputs are from different sources, such as seismic data and
706 gravity data. Collecting data from different sources can help relieve the bottleneck of a limited
707 number of training samples. Besides, using multimodal datasets can increase the quality and
708 reliability of DL methods ([Zhang et al. 2020](#)). [Feng et al. 2020](#) used data integration to forecast
709 streamflow where 23 variables were used integrated, such as precipitation, solar radiation, and
710 temperature. Figure 25 shows an illustration of multimodal DL.

711 4.2.5 Federated learning

712 To provide a practical training set in DL for geophysical applications, collecting available
713 datasets from different institutes or corporations might be a possible solution. However, data
714 transfer via the internet is time-consuming and expensive for large-scale geophysical datasets.

manuscript submitted to Reviews of Geophysics

715 Besides, most datasets are protected and cannot be shared. Federated learning was first proposed
716 by Google ([Mcmahan et al. 2017](#), [Li et al. 2020](#)) to train a DNN with user data from millions of
717 cellphones without privacy or security issues. The encrypted gradients from different clients are
718 assembled in a central server, thus avoiding data transfer. The server updates the model and
719 distributes information to all clients (Figure 26). In a simple federated learning setting, the clients
720 and the server share the same network architecture. We give a possible example of federated
721 learning in geophysics based on the concept that some corporations do not share the annotations
722 of first arrivals; however, they can benefit from federated learning by training a DNN together
723 for first arrival picking.

724 4.2.6 Uncertainty estimation

725 One of the remaining questions associated with applying DL in geophysics is related to
726 whether the results of DL-based model-driven methods with a solid theoretical foundation can be
727 trusted. DL-based uncertainty analysis methods include Monte Carlo dropout ([Gal and](#)
728 [Ghahramani 2016](#)), Markov chain Monte Carlo (MCMC) ([de Figueiredo et al. 2019](#)), variational
729 inference ([Subedar et al. 2019](#)), etc. For example, in Monte Carlo dropout, dropout layers are
730 added to each original layer to simulate a Bernoulli distribution. With multiple realizations of
731 dropout, the results are collected, and the variance is computed as the uncertainty. DL with
732 uncertainty estimation in inference is reported in areas such as volcano-seismic monitoring
733 ([Bueno et al. 2019](#)), geomagnetic storm forecasting ([Tasistro-Hart et al. 2020](#)), weather
734 forecasting ([Scher and Messori , Bonavita and Laloyaux 2020](#)), soil moisture predictions ([Fang](#)
735 [et al. 2020](#)) and earthquake locations estimation ([Mousavi and Beroza 2020](#)).

736 4.2.7 Active learning

737 To train a high-precision model using a small amount of labeled data, active learning is
738 proposed to imitate the self-learning ability of human beings ([Yoo and Kweon 2019](#)). An active
739 learning model selects the most useful data based on a sampling strategy for manual annotation
740 and adds this data to the training set; then, the updated dataset is used for the next round of
741 training (Figure 27). One of the sampling strategies is based on the uncertainty principle, i.e., the
742 samples with high uncertainty are selected. Taking fault detection as an example, if a trained

manuscript submitted to Reviews of Geophysics

743 network is not sure whether a fault exists at a given location, we can annotate the fault manually
744 and add the sample to the training set.

745 **5 Summary**

746 DL methods have created both opportunities and challenges in geophysical fields.
747 Pioneering researchers have provided a basis for DL in geophysics with promising results; more
748 advanced DL technologies and more practical problems must now be explored. To close this
749 paper, we summarize a roadmap for applying DL in different geophysical tasks based on a three-
750 level approach.

- 751 • Traditional methods are time-consuming and require intensive human labor and
752 expert knowledge, such as in first-arrival selection and velocity selection in
753 exploration geophysics.
- 754 • Traditional methods have difficulties and bottlenecks. For example, geophysical
755 inversion requires good initial values and high accuracy modeling and suffers from
756 local minimization.
- 757 • Traditional methods cannot handle some cases, such as multimodal data fusion and
758 inversion.

759 With the development of new artificial intelligence models beyond DL and advances in
760 research into the infinite possibilities of applying DL in geophysics, we can expect intelligent
761 and automatic discoveries of unknown geophysical principles soon.

762 **6 Appendix: a deep learning tutorial for beginners**

763 **6.1 A coding example of a DnCNN**

764 The implementation of DL algorithms in geophysical data processing is quite simple
765 based on existing frameworks, such as Caffe, Pytorch, Keras, and TensorFlow. Here, we provide
766 an example of how to use Python and Keras to construct a DnCNN for seismic denoising. The
767 code requires 12 lines for dataset loading, model construction, training, and testing. The dataset

manuscript submitted to Reviews of Geophysics

768 is preconstructed and includes a clean subset and a noisy subset; the overall dataset includes
 769 12800 samples with a size of 64×64 (available at <https://bit.ly/33SyXPO>).

```

770 1. import h5py
771 2. from tensorflow.keras.layers import Input,Conv2D,BatchNormalization,ReLU,Subtract
772 3. from tensorflow.keras.models import Model
773 4. ftrain = h5py.File('noise_dataset.h5','r')
774 5. X, Y = ftrain['/X'][()], ftrain['/Y'][()]
775 6. input = Input(shape=(None,None,1))
776 7. x = Conv2D(64, 3, padding='same',activation='relu')(input)
777 8. for i in range(15):
778 9.     x = Conv2D(64, 3, padding='same',use_bias = False)(x)
779 10.    x = ReLU()(BatchNormalization(axis=3, momentum=0.0,epsilon=0.0001)(x))
780 11. x = Conv2D(1, 3, padding='same',use_bias = False)(x)
781 12. model = Model(inputs=input, outputs=Subtract()([input, x]))
782 13. model.compile(optimizer="rmsprop", loss="mean_squared_error")
783 14. model.fit(X[:-1000], Y[:-1000], batch_size=32, epochs=50, shuffle=True)
784 15. Y_ = model.predict(X[-1000:])

```

785 Any appropriate plotting tool can be used for data visualization. The training takes less
 786 than one hour on an NVidia 2080Ti graphics processing unit. The readers can try this code in
 787 their own areas as long as a training set is compatibly constructed.

788 6.2 Tips for beginners

789 We introduce several practical tips for beginners who want to explore DL in geophysics
 790 from the perspective of the three most critical steps in DL: data generation, network construction
 791 and training. Though exploration geophysics is used as example, the tips for data generation and
 792 network training are generally applicable to most areas. Network construction generally depends
 793 on the task.

794 6.2.1 Data generation

795 As noted by [Poulton 2002](#), “training a feed-forward neural network is approximately 10%
 796 of the effort involved in an application; deciding on the input and output data coding and creating
 797 good training and testing sets is 90% of the work”. In DL, we advise that the percentages of the

manuscript submitted to Reviews of Geophysics

798 effort for network construction and dataset preparation should be approximately 40% and 60%.
799 First, most DL approaches use an original data set as the input, thus reducing coding decision
800 efforts. Second, a wider variety of network architectures and parameters can be used in DL
801 compared to those in traditional neural networks. Overall, constructing a proper training set plays
802 a more prominent role in DL.

803 Synthetic datasets can be used effectively in DL, which is advantageous since labeled real
804 datasets are sometimes difficult to obtain. First, to assess the applicability of DL in a specific
805 geophysical application, using synthetic datasets is the most convenient method. Second, if a
806 satisfactory result is obtained with synthetic datasets, a few annotated real datasets can be used
807 for transfer learning via parameter tuning. Third, if the synthetic datasets are sufficiently
808 complicated, i.e., if the most important factors are considered when generating the datasets, the
809 trained network may be able to process realistic datasets directly ([Wu et al. 2020](#) and [Wu et al.](#)
810 [2019](#)).

811 A synthetic training set should be diverse. First, we suggest using an existing synthetic
812 dataset with an open license, instead of generating a dataset. For specific tasks, such as FWI, a
813 dataset may need to be generated based on a wave equation. Second, data augmentation methods,
814 such as rotation, reflection, scaling, translation, and adding noise, missing traces, or faults to
815 clean datasets, can be used to expand the training set. The goal is to generate extremely large
816 synthetic datasets that are as close to realistic datasets as possible.

817 To generate realistic datasets, we suggest using existing methods to generate labels that
818 should then be checked by a human. For example, in first-arrival picking, an automatic picking
819 algorithm is used to preprocess the datasets, and the results are then provided to an expert who
820 identifies the outliers. We also suggest using active learning ([Yoo and Kweon 2019](#)) to provide a
821 semiautomated labeling procedure. First, all datasets with machine annotation are used to train a
822 DNN, and the samples with high predicted uncertainty are required to be manually annotated.

823 6.2.2 Network construction for different tasks

824 Beginners are suggested to use a DnCNN or U-Net for testing. DnCNNs are available for
825 most tasks in which the input and output share the same domain, such as denoising, interpolation,

manuscript submitted to Reviews of Geophysics

826 and attribute analysis. The input size of a DnCNN can vary since there are no pooling layers
827 involved. However, each output data point is determined by a local field from the input rather
828 than from the entire input set. Additionally, U-Net contains pooling layers, and all input points
829 are used to determine an output point. U-Nets are available for tasks even when the inputs and
830 outputs are in different domains, such as in FWI. However, the input size of U-Net is fixed once
831 trained and the data need processed patch-wisely.

832 Combining a CAE and K-means is suggested for unsupervised clustering tasks, such as
833 attribute classification. We do not suggest CycleGAN for geophysical tasks since the training
834 process is extremely time-consuming and the results are not stable. An RNN provides a high-
835 performance framework for time-dependent tasks, such as forward wave modeling and FWI.
836 RNNs are also used for regression and classification tasks involving temporal or spatial
837 sequential datasets, such as in the denoising of a single trace.

838 To adjust the hyperparameters of a DNN and optimization algorithms, we suggest using
839 an autoML toolbox, such as Autokeras, instead of manually adjusting the values. The basic
840 objective is to search for the best parameter combination within a given sampling range. Such a
841 search is exceptionally time-consuming, and a random search strategy may accelerate the tuning
842 process. Moreover, for most applications, the default architecture gives reasonable results.

843 6.2.3 Training, validation, and testing

844 The available dataset should be split into three subsets: one training set, one validation set,
845 and one test set to optimize the network parameters. The proportions of the subsets depend on
846 the overall size of a dataset. For datasets with 10K-50K samples, the proportions are suggested to
847 be 60%, 20%, and 20%, respectively. For larger datasets (for instance, those larger than 1M),
848 much smaller portions are often used for validation and test (approximately 1% to 5%) since the
849 alternative can result in using unnecessarily large test/validation sets and wasting the data that
850 can be used for training and building a better model. In a classification task, we suggest using
851 one-hot coding in training. The validation set is used to test the network during training. Then,
852 the model with the best validation accuracy is selected rather than the final trained model. If the
853 validation accuracy does not improve or decrease after some saturation during training, an early

manuscript submitted to Reviews of Geophysics

854 stopping strategy is suggested to avoid overfitting. Network hyperparameters should be tuned
855 according to the validation accuracy. The validation set is used to guide training, and the test set
856 is used to test the model based on unseen datasets; however, this set should not be used for
857 hyperparameter tuning.

858 Two commonly seen issues during training are as follows: the validation loss is less than
859 the training loss, and the loss is not a number. Intuitively, the training loss should be less than the
860 validation loss since the model is trained with a training dataset. Several potential reasons for this
861 issue are as follows: 1. regularization occurs during training but is ignored during validation,
862 such as in the dropout layer; 2. the training loss is obtained by averaging the loss of each batch
863 during an iteration, and the validation loss is obtained based on the loss after one iteration; and 3.
864 the validation set may be less complicated than the training set, especially when only the training
865 set has been augmented. The potential reasons for NaN loss are as follows: 1. the learning rate is
866 too high; 2. in an RNN, one should clip the gradient to avoid gradient explosion and 3. zero is
867 used as a divisor, negative values are used in logarithm, or an exponent is assigned too large of a
868 value.

869 **Glossary**

870 AE: Autoencoder; an ANN with the same inputs and outputs.

871 AI: Artificial Intelligence; Machines are taught to think like humans.

872 ANN: Artificial neural network; a computing system inspired by biological neural networks
873 that constitute animal brains.

874 Aurora: A natural light display in the earth's sky; disturbances in the magnetosphere caused
875 by the solar wind.

876 BNN: Bayesian neural network; the network parameters are random variables instead of
877 regular variables.

878 CAE: Convolutional autoencoder; an AE with shared weights.

879 CNN: Convolutional neural network; a DNN with shared weights.

880 DDTF: Data-driven tight frame; A dictionary learning method using a tight frame constraint
881 for the dictionary.

manuscript submitted to Reviews of Geophysics

- 882 Deblending: In seismic exploration, several explosion sources are shot very close in time to
883 improve efficiency. Then, the seismic waves from different sources are blended. The
884 recorded dataset first needs to be deblended before further processing.
- 885 Dictionary: A set of vectors used to represent signals as a linear combination.
- 886 DIP: Deep image prior; the architecture of a DNN is used as a prior constraint for an image.
- 887 DL: Deep learning; a machine learning technology based on a deep neural network.
- 888 DnCNN: Denoised convolutional neural network.
- 889 DNN: Deep neural network; an ANN with many layers between the input and output layers.
- 890 DS: Double sparsity; the data are represented with a sparse coefficient matrix multiplied by
891 an adaptive dictionary. The adaptive dictionary is represented by a sparse coefficient matrix
892 multiplied by a fixed dictionary.
- 893 Event: In exploration geophysics, a seismic event means reflected waves with the same
894 phase. In seismology, an event means a happened earthquake.
- 895 Facies: A seismic facies unit is a mapped, three-dimensional seismic unit composed of
896 groups of reflections whose parameters differ from adjacent facies units.
- 897 Fault: a discontinuity in a volume of rock across which there has been significant
898 displacement as a result of rock-mass movement.
- 899 FCN: Fully convolutional network; an FCN is a network that contains no fully connected
900 layers. Fully connected layers do not share weights.
- 901 FCNN: Fully connected neural network; an FCNN is a network composed of fully connected
902 layers.
- 903 FWI: Full waveform inversion; full waveform information is used to obtain subsurface
904 parameters. FWI is achieved based on the wave equation and inversion theory.
- 905 GAN: Generative adversarial network; GANs are used to generate fake images. A GAN
906 contains a generative network and a discriminative network. The generative network tries to
907 produce a nearly real image. The discriminative network tries to distinguish whether the
908 input image is real or generated. Therefore, such a game will eventually allow the generative
909 network to produce fake images that the discriminative network cannot distinguish from real
910 images.
- 911 Graphics processing unit (GPU): A parallel computing device. GPUs are widely used for
912 training neural works in deep learning.
- 913 HadCRUT4: Temperature records from Hadley Centre (sea surface temperature) and the
914 Climatic Research Unit (land surface air temperature).

manuscript submitted to Reviews of Geophysics

- 915 K-means: A classical clustering algorithm, where K is the number of clusters.
- 916 K-SVD: A dictionary learning method using SVD for dictionary updating.
- 917 LSTM: long short-term memory; LSTM considers how much historical information is
918 forgotten or remembered with adaptive switches.
- 919 Magnetosphere: Range of the magnetic field surrounding an astronomical object where
920 charged particles are affected.
- 921 ML: Earthquake local magnitude; a method for measuring earthquake scale.
- 922 Patch: In dictionary learning, an image is divided into many patches (blocks) that are the
923 same size as the atoms in a dictionary.
- 924 PINN: Physical informed neural network; A physical equation is used to constrain the neural
925 network.
- 926 PM: Particulate matter. PM10 are coarse particles with a diameter of 10 micrometers or less;
927 PM2.5 are fine particles with a diameter of 2.5 micrometers or less.
- 928 ResNet: Residual neural network; ResNets contain skip connections to jump over several
929 layers. The output of a residual block is the residual between the input and the direct output.
- 930 RNN: Recurrent neural network; in time-sequenced data processing applications, RNNs use
931 the output of a network as the input of the subsequent process to consider the historical
932 context.
- 933 SAR: Synthetic aperture radar; the motion of a radar antenna over a target is treated as an
934 antenna with a large aperture. The larger the aperture is, the higher the image resolution will
935 be.
- 936 Solar wind: A stream of charged particles released from the upper atmosphere of the Sun.
- 937 Sparse coding: Input data are represented in the form of a linear combination of a dictionary
938 where the coefficients are sparse.
- 939 Sparsity: The number of nonzero values in a vector.
- 940 SVD: Singular value decomposition; a matrix factorization method. $\mathbf{A}=\mathbf{U}\mathbf{S}\mathbf{V}$, where \mathbf{U} and \mathbf{V}
941 are two orthogonal matrices, \mathbf{S} is a diagonal matrix whose elements are the singular values of
942 \mathbf{A} . SVD is used for dimension reduction by removing the smaller singular values. SVD is
943 also used for recommendation systems and natural language processing.
- 944 Tight frame: A frame provides a redundant, stable way of representing a signal, similar to
945 dictionary. A tight frame is a frame with the perfect reconstruction property; i.e., $\mathbf{W}^T\mathbf{W}=\mathbf{I}$.
- 946 Tomography: Inversion of the subsurface velocity based on travel time information.

manuscript submitted to Reviews of Geophysics

947 U-Net: U-shaped network; U-Nets have U-shaped structures and skip connections. The skip
948 connections bring low-level features to high levels.

949 Wave equation: A partial differential equation that controls wave propagation.

950 WST: Wavelet scattering transform; a transform involves a cascade of wavelet transforms, a
951 module operator, and an averaging operator.

952 **Acknowledgments**

953 The work was supported in part by the National Key Research and Development Program
954 of China under grant nos. 2017YFB0202902 and 2018YFC1503705 and NSFC under grant nos.
955 41625017 and 41804102. We thank Society of Exploration Geophysicists, Nature Research, and
956 American Association for the Advancement of Science for allowing us to reuse the original
957 figures from their journals.

958 **Data Availability Statement**

959 Data were not used, nor created for this research.

960

961 **References**

- 962 Abma, R. and N. Kabir (2006). 3D interpolation of irregular data with a POCS algorithm. *Geophysics*. **71**(6): 91-97.
- 963 Acito, N., M. Diani and G. Corsini (2020). Cwv-net: A deep neural network for atmospheric column water vapor
964 retrieval from hyperspectral vnir data. *IEEE Transactions on Geoscience and Remote Sensing*. **58**(11): 8163-8175.
- 965 Aharon, M., M. Elad and A. Bruckstein (2006). K-SVD: An algorithm for designing overcomplete dictionaries for
966 sparse representation. *IEEE Transactions on Signal Processing*. **54**(11): 4311-4322.
- 967 Akbari Asanjan, A., T. Yang, K. Hsu, S. Sorooshian, J. Lin and Q. Peng (2018). Short - term precipitation forecast
968 based on the persiann system and lstm recurrent neural networks. *Journal of Geophysical Research: Atmospheres*.
969 **123**(22): 12-563.
- 970 Anantrasirichai, N., J. Biggs, F. Albino, P. Hill and D. Bull (2018). Application of machine learning to classification
971 of volcanic deformation in routinely-generated InSAR data. *Journal of Geophysical Research*. **123**(8): 6592-6606.
- 972 Araya-Polo, M., J. Jennings, A. Adler and T. Dahlke (2018). Deep-learning tomography. *The Leading Edge*. **37**(1):
973 58-66.
- 974 Barbat, M. M., T. Rackow, H. H. Hellmer, C. Wesche and M. M. Mata (2019). Three years of near-coastal antarctic
975 iceberg distribution from a machine learning approach applied to SAR imagery. *Journal of Geophysical Research-Oceans*. **124**(9): 6658-6672.
- 976 Bergen, K. J., P. A. Johnson, M. V. de Hoop and G. C. Beroza (2019). Machine learning for data-driven discovery
977 in solid earth geoscience. *Science*. **363**(6433): 1-10.
- 978 Bonavita, M. and P. Laloyaux (2020). Machine learning for model error inference and correction. *Journal of
979 Advances in Modeling Earth Systems*. **12**(12): e2020MS002232.
- 980 Bueno, A., C. Benitez, S. De Angelis, A. D. Moreno and J. M. Ibanez (2019). Volcano-seismic transfer learning and
981 uncertainty quantification with bayesian neural networks. *IEEE Transactions on Geoscience and Remote Sensing*.
982 **58**(2): 892-902.
- 983 Cai, J., H. Ji, Z. Shen and G. Ye (2014). Data-driven tight frame construction and image denoising. *Applied and
984 Computational Harmonic Analysis*. **37**(1): 89-105.
- 985 Cao, R., S. Earp, S. A. L. de Ridder, A. Curtis and E. Galetti (2020). Near-real-time near-surface 3D seismic
986 velocity and uncertainty models by wavefield gradiometry and neural network inversion of ambient seismic noise.
987 *Geophysics*. **85**(1): KS13-KS27.
- 988 Chattopadhyay, A., A. Subel and P. Hassanzadeh (2020). Data - driven super - parameterization using deep
989 learning: Experimentation with multiscale lorenz 96 systems and transfer learning. *Journal of Advances in Modeling
990 Earth Systems*. **12**(11): e2020MS002084.
- 991 Chen, R. T., Y. Rubanova, J. Bettencourt and D. Duvenaud (2018). Neural ordinary differential equations. arXiv
992 preprint arXiv:1806.07366.
- 993 Chen, S., H. Wang, F. Xu and Y. Jin (2016). Target classification using the deep convolutional networks for SAR
994 images. *IEEE Transactions on Geoscience and Remote Sensing*. **54**(8): 4806-4817.
- 995 Chen, Y., H. Jiang, C. Li, X. Jia and P. Ghamisi (2016). Deep feature extraction and classification of hyperspectral
996 images based on convolutional neural networks. *IEEE Transactions on Geoscience and Remote Sensing*. **54**(10):
997 6232-6251.
- 998 Chen, Z., M. Jin, Y. Deng, J.-S. Wang, H. Huang, X. Deng and C.-M. Huang (2019). Improvement of a deep
999 learning algorithm for total electron content maps: Image completion. *Journal of Geophysical Research*. **124**(1):
1000 790-800.
- 1001 Cheng, G., P. Zhou and J. Han (2016). Learning rotation-invariant convolutional neural networks for object
1002 detection in VHR optical remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*. **54**(12):
1003 7405-7415.
- 1004 Cheng, X., Q. Liu, P. Li and Y. Liu (2019). Inverting rayleigh surface wave velocities for crustal thickness in eastern
1005 Tibet and the western Yangtze craton based on deep learning neural networks. *Nonlinear Processes in Geophysics*.
1006 **26**(2): 61-71.
- 1007 Chu, X., J. Bortnik, W. Li, Q. Ma, R. Denton, C. Yue, V. Angelopoulos, R. M. Thorne, F. Darrouzet, P. Ozhogin, C.
1008 A. Kletzing, Y. Wang and J. Menietti (2017). A neural network model of three-dimensional dynamic electron
1009 density in the inner magnetosphere. *Journal of Geophysical Research*. **122**(9): 9183-9197.
- 1010 Clausen, L. B. N. and H. Nickisch (2018). Automatic classification of auroral images from the oslo auroral themis
1011 (OATH) data set using machine learning. *Journal of Geophysical Research-Space Physics*. **123**(7): 5640-5647.
- 1012 Das, V. and T. Mukerji (2020). Petrophysical properties prediction from prestack seismic data using convolutional
1013 neural networks. *Geophysics*. **85**(5): N41-N55.

manuscript submitted to Reviews of Geophysics

- 1015 Das, V., A. Pollack, U. Wollner and T. Mukerji (2019). Convolutional neural network for seismic impedance
1016 inversion. *Geophysics*. **84**(6): R869-R880.
- 1017 de Figueiredo, L. P., D. Grana, M. Roisenberg and B. B. Rodrigues (2019). Gaussian mixture markov chain Monte
1018 Carlo method for linear seismic inversion. *Geophysics*. **84**(3): R463-R476.
- 1019 DeVries, P. M. R., F. Viegas, M. Wattenberg and B. J. Meade (2018). Deep learning of aftershock patterns
1020 following large earthquakes. *Nature*. **560**(7720): 632-634.
- 1021 Dhara, A. and C. Bagaini (2020). Seismic image registration using multiscale convolutional neural networks.
1022 *Geophysics*. **85**(6): V425-V441.
- 1023 Donahue, J., Y. Jia, O. Vinyals, J. Hoffman and T. Darrell (2014). DeCAF: A deep convolutional activation feature
1024 for generic visual recognition. International Conference on Machine Learning: 647-655.
- 1025 Dong, C., C. C. Loy, K. He and X. Tang (2014). Learning a deep convolutional network for image super-resolution.
1026 European Conference on Computer Vision: 184-199.
- 1027 Donoho, D. L. and I. M. Johnstone (1995). Adapting to unknown smoothness via wavelet shrinkage. *Journal of the
1028 american statistical association*. **90**(432): 1200-1224.
- 1029 Duan, Y., X. Zheng, L. Hu and L. Sun (2019). Seismic facies analysis based on deep convolutional embedded
1030 clustering. *Geophysics*. **84**(6): IM87-IM97.
- 1031 Dunham, M. W., A. Malcolm and J. Kim Welford (2019). Improved well-log classification using semisupervised
1032 label propagation and self-training, with comparisons to popular supervised algorithms. *Geophysics*. **85**(1): O1-O15.
- 1033 Fang, J., H. Zhou, Y. Elita Li, Q. Zhang, L. Wang, P. Sun and J. Zhang (2020). Data-driven low-frequency signal
1034 recovery using deep-learning predictions in full-waveform inversion. *Geophysics*. **85**(6): A37-A43.
- 1035 Fang, K., D. Kifer, K. Lawson and C. Shen (2020). Evaluating the potential and challenges of an uncertainty
1036 quantification method for long short - term memory models for soil moisture predictions. *Water Resources
1037 Research*. **56**(12): e2020WR028095.
- 1038 Fang, K., C. Shen, D. Kifer and X. Yang (2017). Prolongation of SMAP to spatiotemporally seamless coverage of
1039 continental U.S. Using a deep learning neural network. *Geophysical Research Letters*. **44**(21).
- 1040 Feng, D. P., K. Fang and C. P. Shen (2020). Enhancing streamflow forecast and extracting insights using long-short
1041 term memory networks with data integration at continental scales. *Water Resources Research*. **56**(9):
1042 e2019WR026793.
- 1043 Feng, R., T. Mejer Hansen, D. Grana and N. Balling (2020). An unsupervised deep-learning method for porosity
1044 estimation based on poststack seismic data. *Geophysics*. **85**(6): M97-M105.
- 1045 Gal, Y. and Z. Ghahramani (2016). Dropout as a Bayesian approximation: Representing model uncertainty in deep
1046 learning. International Conference on Machine Learning.
- 1047 Gao, Z., Z. Pan, J. Gao and Z. Xu (2019). Building long-wavelength velocity for salt structure using stochastic full
1048 waveform inversion with deep autoencoder based model reduction. *SEG Technical Program Expanded Abstracts:*
1049 1680-1684.
- 1050 Garofalo, F., G. Sauvin, L. V. Socco and I. Lecomte (2015). Joint inversion of seismic and electric data applied to
1051 2D media. *Geophysics*. **80**(4): EN93-EN104.
- 1052 Grana, D., L. Azevedo and M. Liu (2020). A comparison of deep machine learning and Monte Carlo methods for
1053 facies classification from seismic data. *Geophysics*. **85**(4): WA41-WA52.
- 1054 He, K., X. Zhang, S. Ren and J. Sun (2016). Deep residual learning for image recognition. *IEEE Conference on
1055 Computer Vision and Pattern Recognition*: 770-778.
- 1056 He, Y., J. Cao, Y. Lu, Y. Gan and S. Lv (2018). Shale seismic facies recognition technology based on sparse
1057 autoencoder. *International Geophysical Conference*.
- 1058 Helmy, T., A. Fatai and K. Faisal (2010). Hybrid computational models for the characterization of oil and gas
1059 reservoirs. *Expert Systems with Applications*. **37**(7): 5353-5363.
- 1060 Herrmann, F. J. and G. Hennenfent (2008). Non-parametric seismic data recovery with curvelet frames. *Geophysical
1061 Journal International*. **173**(1): 233-248.
- 1062 Hochreiter, S. and J. Schmidhuber (1997). Long short-term memory. *Neural Computation*. **9**(8): 1735-1780.
- 1063 Hu, A., B. Carter, J. Currie, R. Norman, S. Wu and K. Zhang (2020). A deep neural network model of global topside
1064 electron temperature using incoherent scatter radars and its application to gnss radio occultation. *Journal of
1065 Geophysical Research*. **125**(2): 1-17.
- 1066 Hu, L., X. Zheng, Y. Duan, X. Yan, Y. Hu and X. Zhang (2019). First-arrival picking with a U-net convolutional
1067 network. *Geophysics*. **84**(6): U45-U57.
- 1068 Huang, K., J. You, K. Chen, H. Lai and A. Don (2006). Neural network for parameters determination and seismic
1069 pattern detection. *SEG Technical Program Expanded Abstracts*: 2285-2289.

manuscript submitted to Reviews of Geophysics

- 1070 Iten, R., T. Metger, H. Wilming, L. Del Rio and R. Renner (2020). Discovering physical concepts with neural
1071 networks. *Phys Rev Lett.* **124**(1): 010508.
- 1072 Jia, Y. and J. Ma (2017). What can machine learning do for seismic data processing? An interpolation application.
1073 *Geophysics.* **82**(3): V163-V177.
- 1074 Jiang, G.-q., J. Xu and J. Wei (2018). A deep learning algorithm of neural network for the parameterization of
1075 typhoon-ocean feedback in typhoon forecast models. *Geophysical Research Letters.* **45**(8): 3706-3716.
- 1076 Jiang, K., Z. Wang, P. Yi, G. Wang, T. Lu and J. Jiang (2019). Edge-enhanced GAN for remote sensing image
1077 superresolution. *IEEE Transactions on Geoscience and Remote Sensing.* **57**(8): 5799-5812.
- 1078 Kadow, C., D. M. Hall and U. Ulbrich (2020). Artificial intelligence reconstructs missing climate information.
1079 *Nature Geoscience.* **13**(6): 408-413.
- 1080 Krizhevsky, A., I. Sutskever and G. E. Hinton (2017). Imagenet classification with deep convolutional neural
1081 networks. *Communications of the ACM.* **60**(6): 84-90.
- 1082 LeCun, Y., Y. Bengio and G. Hinton (2015). Deep learning. *Nature.* **521**(7553): 436-444.
- 1083 Lee, S., E. Y. Ji, Y. J. Moon and E. Park (2021). One - day forecasting of global tec using a novel deep learning
1084 model. *Space Weather.* **19**(1): 2020SW002600.
- 1085 Lei, N., D. An, Y. Guo, K. Su, S. Liu, Z. Luo, S.-T. Yau and X. Gu (2020). A geometric understanding of deep
1086 learning. *Engineering.* **6**(3): 361-374.
- 1087 Lempitsky, V., A. Vedaldi and D. Ulyanov (2018). Deep image prior. *IEEE Conference on Computer Vision and
1088 Pattern Recognition:* 9446-9454.
- 1089 Li, L., Y. Lin, X. Zhang, H. Liang, W. Xiong and S. Zhan (2019). Convolutional recurrent neural networks based
1090 waveform classification in seismic facies analysis. *SEG Technical Program Expanded Abstracts:* 2599-2603.
- 1091 Li, S., W. Song, L. Fang, Y. Chen, P. Ghamisi and J. A. Benediktsson (2019). Deep learning for hyperspectral
1092 image classification: An overview. *IEEE Transactions on Geoscience and Remote Sensing.* **57**(9): 6690-6709.
- 1093 Li, T., A. K. Sahu, A. Talwalkar and V. Smith (2020). Federated learning: Challenges, methods, and future
1094 directions. *IEEE Signal Processing Magazine.* **37**(3): 50-60.
- 1095 Li, T., H. Shen, Q. Yuan, X. Zhang and L. Zhang (2017). Estimating ground-level PM2.5 by fusing satellite and
1096 station observations: A geo-intelligent deep learning approach. *Geophysical Research Letters.* **44**(23): 11,985-
1097 911,993.
- 1098 Li, Z., M. A. Meier, E. Hauksson, Z. Zhan and J. Andrews (2018). Machine learning seismic wave discrimination:
1099 Application to earthquake early warning. *Geophysical Research Letters.* **45**(10): 4773-4779.
- 1100 Liang, J., J. Ma and X. Zhang (2014). Seismic data restoration via data-driven tight frame. *Geophysics.* **79**(3): V65-
1101 V74.
- 1102 Lim, J. S. (2005). Reservoir properties determination using fuzzy logic and neural networks from well data in
1103 offshore korea. *Journal of Petroleum Science and Engineering.* **49**(3-4): 182-192.
- 1104 Ling, F., D. Boyd, Y. Ge, G. M. Foody, X. Li, L. Wang, Y. Zhang, L. Shi, C. Shang, X. Li and Y. Du (2019).
1105 Measuring river wetted width from remotely sensed imagery at the subpixel scale with a deep convolutional neural
1106 network. *Water Resources Research.* **55**(7): 5631-5649.
- 1107 Linville, L., K. Pankow and T. Draelos (2019). Deep learning models augment analyst decisions for event
1108 discrimination. *Geophysical Research Letters.* **46**(7): 3643-3651.
- 1109 Liu, B., X. Li and G. Zheng (2019). Coastal inundation mapping from bitemporal and dual-polarization SAR
1110 imagery based on deep convolutional neural networks. *Journal of Geophysical Research-Oceans.* **124**(12): 9101-
1111 9113.
- 1112 Liu, L., S. Zou, Y. Yao and Z. Wang (2020). Forecasting global ionospheric tec using deep learning approach. *Space
1113 Weather.* **18**(11): e2020SW002501.
- 1114 Liu, S. (2020). Multi-parameter full waveform inversions based on recurrent neural networks. Dissertation for the
1115 Master Degree in Science, Harbin Institute of Technology.
- 1116 Maggiori, E., Y. Tarabalka, G. Charpiat and P. Alliez (2017). Convolutional neural networks for large-scale remote-
1117 sensing image classification. *IEEE Transactions on Geoscience and Remote Sensing.* **55**(2): 645-657.
- 1118 Malfante, M., M. Dalla Mura, J. I. Mars, J. P. Metaxian, O. Macedo and A. Inza (2018). Automatic classification of
1119 volcano seismic signatures. *Journal of Geophysical Research-Solid Earth.* **123**(12): 10645-10658.
- 1120 Mallat, S. (2012). Group invariant scattering. *Communications on Pure and Applied Mathematics.* **65**(10): 1331-
1121 1398.
- 1122 Mandelli, S., F. Borra, V. Lipari, P. Bestagini, A. Sarti and S. Tubaro (2018). Seismic data interpolation through
1123 convolutional autoencoder. *SEG Technical Program Expanded Abstracts:* 4101-4105.

manuscript submitted to Reviews of Geophysics

- 1124 Manucharyan, G. E., L. Siegelman and P. Klein (2021). A deep learning approach to spatiotemporal sea surface
1125 height interpolation and estimation of deep currents in geostrophic ocean turbulence. *Journal of Advances in*
1126 *Modeling Earth Systems*. **13**(1).
- 1127 McMahan, H. B., E. Moore, D. Ramage, S. Hampson and B. A. Y. Arcas (2017). Communication-efficient learning
1128 of deep networks from decentralized data. *International Conference on Artificial Intelligence and Statistics*.
- 1129 Meier, M. A., Z. E. Ross, A. Ramachandran, A. Balakrishna, S. Nair, P. Kundzicz, Z. F. Li, J. Andrews, E.
1130 Hauksson and Y. S. Yue (2019). Reliable real-time seismic signal/noise discrimination with machine learning.
1131 *Journal of Geophysical Research-Solid Earth*. **124**(1): 788-800.
- 1132 Mou, L., P. Ghamisi and X. X. Zhu (2017). Deep recurrent neural networks for hyperspectral image classification.
1133 *IEEE Transactions on Geoscience and Remote Sensing*. **55**(7): 3639-3655.
- 1134 Mousavi, S. M. and G. C. Beroza (2020). Bayesian-deep-learning estimation of earthquake location from single-
1135 station observations. *IEEE Transactions on Geoscience and Remote Sensing*. **58**(11): 8211-8224.
- 1136 Mousavi, S. M. and G. C. Beroza (2020). A machine-learning approach for earthquake magnitude estimation.
1137 *Geophysical Research Letters*. **47**(1).
- 1138 Mousavi, S. M., W. L. Ellsworth, W. Zhu, L. Y. Chuang and G. C. Beroza (2020). Earthquake transformer—an
1139 attentive deep-learning model for simultaneous earthquake detection and phase picking. *Nature communications*.
1140 **11**(1): 1-12.
- 1141 Mousavi, S. M., S. P. Horton, C. A. Langston and B. Samei (2016). Seismic features and automatic discrimination
1142 of deep and shallow induced-microearthquakes using neural network and logistic regression. *Geophysical Journal*
1143 *International*. **207**(1): 29-46.
- 1144 Mousavi, S. M. and C. A. Langston (2016). Hybrid seismic denoising using higher - order statistics and improved
1145 wavelet block thresholding. *Bulletin of the Seismological Society of America*. **106**(4): 1380-1393.
- 1146 Mousavi, S. M. and C. A. Langston (2017). Automatic noise-removal/signal-removal based on general cross-
1147 validation thresholding in synchrosqueezed domain and its application on earthquake data. *Geophysics*. **82**(4):
1148 V211-V227.
- 1149 Mousavi, S. M., C. A. Langston and S. P. Horton (2016). Automatic microseismic denoising and onset detection
1150 using the synchrosqueezed continuous wavelet transform. *Geophysics*. **81**(4): V341-V355.
- 1151 Mousavi, S. M., W. Zhu, W. Ellsworth and G. Beroza (2019). Unsupervised clustering of seismic signals using deep
1152 convolutional autoencoders. *IEEE Geoscience and Remote Sensing Letters*. **16**(11): 1693-1697.
- 1153 Mousavi, S. M., W. Zhu, Y. Sheng and G. C. Beroza (2019). CRED: A deep residual network of convolutional and
1154 recurrent units for earthquake signal detection. *Scientific Reports*. **9**(1): 1-14.
- 1155 Nazari Siahzar, M. A., S. Gholtashi, A. R. Kahoo, W. Chen and Y. Chen (2017). Data-driven multitask sparse
1156 dictionary learning for noise attenuation of 3D seismic data. *Geophysics*. **82**(6): V385-V396.
- 1157 Nearing, G. S., F. Kratzert, A. K. Sampson, C. S. Pelissier, D. Klotz, J. M. Frame, C. Prieto and H. V. Gupta (2020).
1158 What role does hydrological science play in the age of machine learning? *Water Resources Research*:
1159 e2020WR028091.
- 1160 Ngiam, J., A. Khosla, M. Kim, J. Nam, H. Lee and A. Y. Ng (2011). Multimodal deep learning. *International*
1161 *Conference on Machine Learning*.
- 1162 Niu, Y., Y. D. Wang, P. Mostaghimi, P. Swietojanski and R. T. Armstrong (2020). An innovative application of
1163 generative adversarial networks for physically accurate rock images with an unprecedented field of view.
1164 *Geophysical Research Letters*. **47**(23): e2020GL089029.
- 1165 Oquab, M., L. Bottou, I. Laptev and J. Sivic (2014). Learning and transferring mid-level image representations using
1166 convolutional neural networks. *IEEE Conference on Computer Vision and Pattern Recognition*.
- 1167 Oropeza, V. and M. Sacchi (2011). Simultaneous seismic data denoising and reconstruction via multichannel
1168 singular spectrum analysis. *Geophysics*. **76**(3): V25-V32.
- 1169 Ovcharenko, O., V. Kazei, M. Kalita, D. Peter and T. Alkhalifah (2019). Deep learning for low-frequency
1170 extrapolation from multioffset seismic data. *Geophysics*. **84**(6): R989-R1001.
- 1171 Park, M. J. and M. D. Sacchi (2019). Automatic velocity analysis using convolutional neural network and transfer
1172 learning. *Geophysics*. **85**(1): V33-V43.
- 1173 Payani, A., F. Fekri, G. Alregib, M. Mohandes and M. Deriche (2019). Compression of seismic signals via recurrent
1174 neural networks: Lossy and lossless algorithms. *SEG Technical Program Expanded Abstracts 2019*: 4082-4086.
- 1175 Poulton, M. M. (2002). Neural networks as an intelligence amplification tool: A review of applications. *Geophysics*.
1176 **67**(3): 979-993.
- 1177 Qi, J., B. Zhang, B. Lyu and K. Marfurt (2020). Seismic attribute selection for machine-learning-based facies
1178 analysis. *Geophysics*. **85**(2): O17-O35.

manuscript submitted to Reviews of Geophysics

- 1179 Qian, F., M. Yin, X. Liu, Y. Wang, C. Lu and G. Hu (2018). Unsupervised seismic facies analysis via deep
1180 convolutional autoencoders. *Geophysics*. **83**(3): A39-A43.
- 1181 Raissi, M., P. Perdikaris and G. E. Karniadakis (2019). Physics-informed neural networks: A deep learning
1182 framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of
1183 Computational Physics*. **378**: 686-707.
- 1184 Ramachandram, D. and G. W. Taylor (2017). Deep multimodal learning: A survey on recent advances and trends.
1185 *IEEE Signal Processing Magazine*. **34**(6): 96-108.
- 1186 Read, J. S., X. Jia, J. Willard, A. P. Appling, J. A. Zwart, S. K. Oliver, A. Karpatne, G. J. A. Hansen, P. C. Hanson,
1187 W. Watkins, M. Steinbach and V. Kumar (2019). Process-guided deep learning predictions of lake water
1188 temperature. *Water Resources Research*. **55**(11): 9173-9190.
- 1189 Reichstein, M., G. Camps-Valls, B. Stevens, M. Jung, J. Denzler, N. Carvalhais and Prabhat (2019). Deep learning
1190 and process understanding for data-driven earth system science. *Nature*. **566**(7743): 195-204.
- 1191 Ronneberger, O., P. Fischer and T. Brox (2015). U-net: Convolutional networks for biomedical image segmentation.
1192 *Medical Image Computing and Computer Assisted Intervention*: 234-241.
- 1193 Ross, Z. E., M.-A. Meier and E. Hauksson (2018). Pwave arrival picking and first-motion polarity determination
1194 with deep learning. *Journal of Geophysical Research: Solid Earth*. **123**(6): 5120-5129.
- 1195 Ross, Z. E., Y. S. Yue, M. A. Meier, E. Hauksson and T. H. Heaton (2019). Phaselink: A deep learning approach to
1196 seismic phase association. *Journal of Geophysical Research-Solid Earth*. **124**(1): 856-869.
- 1197 Rubinstein, R., M. Zibulevsky and M. Elad (2010). Double sparsity: Learning sparse dictionaries for sparse signal
1198 approximation. *IEEE Transactions on Signal Processing*. **58**(3): 1553-1564.
- 1199 Rumelhart, D. E., G. E. Hinton and R. J. Williams (1986). Learning representations by back-propagating errors.
1200 *Nature*. **323**(6088): 533-536.
- 1201 Rütgers, M., S. Lee, S. Jeon and D. You (2019). Prediction of a typhoon track using a generative adversarial
1202 network and satellite images. *Scientific reports*. **9**(1): 1-15.
- 1203 Scher, S. and G. Messori Ensemble methods for neural network - based weather forecasts. *Journal of Advances in
1204 Modeling Earth Systems*: e2020MS002331.
- 1205 Shahnas, M. H. and R. N. Pysklywec (2020). Toward a unified model for the thermal state of the planetary mantle:
1206 Estimations from mean field deep learning. *Earth and Space Science*. **7**(7).
- 1207 Shen, C. (2018). A transdisciplinary review of deep learning research and its relevance for water resources scientists.
1208 *Water Resources Research*. **54**(11): 8558-8593.
- 1209 Shen, H., T. Li, Q. Yuan and L. Zhang (2018). Estimating regional ground-level PM2.5 directly from satellite top-
1210 of-atmosphere reflectance using deep belief networks. *Journal of Geophysical Research*. **123**(24).
- 1211 Siahkoohi, A., M. Louboutin and F. J. Herrmann (2019). The importance of transfer learning in seismic modeling
1212 and imaging. *Geophysics*. **84**(6): A47-A52.
- 1213 Simonyan, K. and A. Zisserman (2015). Very deep convolutional networks for large-scale image recognition.
1214 *International Conference on Learning Representations*.
- 1215 Spitz, S. (1991). Seismic trace interpolation in the F-X domain. *Geophysics*. **56**(6): 785-794.
- 1216 Subedar, M., R. Krishnan, P. L. Meyer, O. Tickoo and J. Huang (2019). Uncertainty-aware audiovisual activity
1217 recognition using deep Bayesian variational inference. *International Conference on Computer Vision*: 6300-6309.
- 1218 Sun, A. Y., B. R. Scanlon, H. Save and A. Rateb (2020). Reconstruction of grace total water storage through
1219 automated machine learning. *Water Resources Research*: e2020WR028666.
- 1220 Sun, A. Y., B. R. Scanlon, Z. Zhang, D. Walling, S. N. Bhanja, A. Mukherjee and Z. Zhong (2019). Combining
1221 physically based modeling and deep learning for fusing grace satellite data: Can we learn from mismatch? *Water
1222 Resources Research*. **55**(2): 1179-1195.
- 1223 Sun, B. and T. Alkhailah (2020). MI-descent: An optimization algorithm for full-waveform inversion using
1224 machine learning. *Geophysics*. **85**(6): R477-R492.
- 1225 Sun, J., Z. Niu, K. A. Innanen, J. Li and D. O. Trad (2020). A theory-guided deep-learning formulation and
1226 optimization of seismic waveform inversion. *Geophysics*. **85**(2): R87-R99.
- 1227 Tang, G., D. Long, A. Behrangi, C. Wang and Y. Hong (2018). Exploring deep neural networks to retrieve rain and
1228 snow in high latitudes using multisensor and reanalysis data. *Water Resources Research*. **54**(10): 8253-8278.
- 1229 Tasistro - Hart, A., A. Grayver and A. Kuvshinov (2020). Probabilistic geomagnetic storm forecasting via deep
1230 learning. *Journal of Geophysical Research: Space Physics*: e2020JA028228.
- 1231 Wang, B., N. Zhang, W. Lu and J. Wang (2019). Deep-learning-based seismic data interpolation: A preliminary
1232 result. *Geophysics*. **84**(1): V11-V20.
- 1233 Wang, J., Z. Xiao, C. Liu, D. Zhao and Z. Yao (2019). Deep learning for picking seismic arrival times. *Journal of
1234 Geophysical Research: Solid Earth*. **124**(7): 6612-6624.

manuscript submitted to Reviews of Geophysics

- 1235 Wang, J. L., H. Zhuang, L. M. Chérubin, A. K. Ibrahim and A. M. Ali (2019). Medium-term forecasting of loop
 1236 current eddy cameron and eddy darwin formation in the gulf of mexico with a divide-and-conquer machine learning
 1237 approach. *Journal of Geophysical Research*. **124**(8): 5586-5606.
- 1238 Wang, N., H. Chang and D. Zhang (2020). Deep - learning - based inverse modeling approaches: A subsurface flow
 1239 example. *Journal of Geophysical Research: Solid Earth*: e2020JB020549.
- 1240 Wang, T., Z. Zhang and Y. Li (2019). Earthquakegen: Earthquake generator using generative adversarial networks.
 1241 SEG Technical Program Expanded Abstracts: 2674-2678.
- 1242 Wang, W. and J. Ma (2020). Velocity model building in a crosswell acquisition geometry with image-trained
 1243 artificial neural network. *geophysics*. **85**(2): U31-U46.
- 1244 Wang, W., G. A. McMechan and J. Ma (2020). Elastic full-waveform inversion with recurrent neural networks. Seg
 1245 technical program expanded abstracts 2020, Society of Exploration Geophysicists: 860-864.
- 1246 Wang, Y., Q. Ge, W. Lu and X. Yan (2019). Seismic impedance inversion based on cycle-consistent generative
 1247 adversarial network. SEG Technical Program Expanded Abstracts: 2498-2502.
- 1248 Wang, Y., B. Wang, N. Tu and J. Geng (2020). Seismic trace interpolation for irregularly spatial sampled data using
 1249 convolutional autoencoder. *Geophysics*. **85**(2): V119-V130.
- 1250 Wu, H., B. Zhang, F. Li and N. Liu (2019). Semiautomatic first-arrival picking of microseismic events by using the
 1251 pixel-wise convolutional image segmentation method. *Geophysics*. **84**(3): V143-V155.
- 1252 Wu, H., B. Zhang, T. Lin, D. Cao and Y. Lou (2019). Semiautomated seismic horizon interpretation using the
 1253 encoder-decoder convolutional neural network. *Geophysics*. **84**(6): B403-B417.
- 1254 Wu, H., B. Zhang, T. Lin, F. Li and N. Liu (2019). White noise attenuation of seismic trace by integrating
 1255 variational mode decomposition with convolutional neural network. *Geophysics*. **84**(5): V307-V317.
- 1256 Wu, X., Z. Geng, Y. Shi, N. Pham, S. Fomel and G. Caumon (2020). Building realistic structure models to train
 1257 convolutional neural networks for seismic structural interpretation. *Geophysics*. **85**(4): WA27-WA39.
- 1258 Wu, X., L. Liang, Y. Shi and S. Fomel (2019). FaultSeg3D: Using synthetic data sets to train an end-to-end
 1259 convolutional neural network for 3D seismic fault segmentation. *Geophysics*. **84**(3): IM35-IM45.
- 1260 Wu, X., Y. Shi, S. Fomel, L. Liang, Q. Zhang and A. Z. Yusifov (2019). Faultnet3d: Predicting fault probabilities,
 1261 strikes, and dips with a single convolutional neural network. *IEEE Transactions on Geoscience and Remote Sensing*.
57(11): 9138-9155.
- 1263 Wu, Y. and G. A. McMechan (2019). Parametric convolutional neural network-domain full-waveform inversion.
 1264 *Geophysics*. **84**(6): R881-R896.
- 1265 Xiao, C., Y. Deng and G. Wang Deep - learning - based adjoint state method: Methodology and preliminary
 1266 application to inverse modelling. *Water Resources Research*: e2020WR027400.
- 1267 Yamaga, N. and Y. Mitsui (2019). Machine learning approach to characterize the postseismic deformation of the
 1268 2011 Tohoku-Oki earthquake based on recurrent neural network. *Geophysical Research Letters*. **46**(21): 11886-
 1269 11892.
- 1270 Yang, F. and J. Ma (2019). Deep-learning inversion: A next-generation seismic velocity model building method.
 1271 *Geophysics*. **84**(4): R585-R584.
- 1272 Yang, F. and J. Ma (2019). Deep-learning inversion: A next-generation seismic velocity model building method.
 1273 *Geophysics*. **84**(4): R583-R599.
- 1274 Yang, Q., D. Tao, D. Han and J. Liang (2019). Extracting auroral key local structures from all-sky auroral image by
 1275 artificial intelligence technique. *Journal of Geophysical Research-Space Physics*. **124**(5): 3512-3521.
- 1276 Yoo, D. and I. S. Kweon (2019). Learning loss for active learning. *IEEE Conference on Computer Vision and*
 1277 *Pattern Recognition*.
- 1278 Yosinski, J., J. Clune, Y. Bengio and H. Lipson (2014). How transferable are features in deep neural networks.
 1279 *Neural Information Processing Systems*.
- 1280 You, N., Y. E. Li and A. Cheng (2020). Shale anisotropy model building based on deep neural networks. *Journal of*
 1281 *Geophysical Research: Solid Earth*. **125**(2): e2019JB019042.
- 1282 Yu, S., J. Ma and S. Osher (2016). Monte Carlo data-driven tight frame for seismic data recovery. *Geophysics*. **81**(4):
 1283 V327-V340.
- 1284 Yu, S., J. Ma and W. Wang (2019). Deep learning for denoising. *Geophysics*. **84**(6): V333-V350.
- 1285 Yu, S., J. Ma, X. Zhang and M. Sacchi (2015). Interpolation and denoising of high-dimensional seismic data by
 1286 learning a tight frame. *Geophysics*. **80**(5): V119-V132.
- 1287 Yuan, P., S. Wang, W. Hu, X. Wu, J. Chen and H. Van Nguyen (2020). A robust first-arrival picking workflow
 1288 using convolutional and recurrent neural networks. *Geophysics*. **85**(5): U109-U119.
- 1289 Zhang, C., C. Frogner, M. Araya-Polo and D. Hohl (2014). Machine-learning based automated fault detection in
 1290 seismic traces. *76th EAGE Conference and Exhibition 2014*. **2014**(1): 1-5.

manuscript submitted to Reviews of Geophysics

- 1291 Zhang, H., X. Yang and J. Ma (2020). Can learning from natural image denoising be used for seismic data
1292 interpolation? *Geophysics*. **85**(4): WA115-WA136.
- 1293 Zhang, K., W. Zuo, Y. Chen, D. Meng and L. Zhang (2017). Beyond a Gaussian denoiser: Residual learning of deep
1294 CNN for image denoising. *IEEE Transactions on Image Processing*. **26**(7): 3142-3155.
- 1295 Zhang, K., W. Zuo, S. Gu and L. Zhang (2017). Learning deep CNN denoiser prior for image restoration. *IEEE*
1296 Conference on Computer Vision and Pattern Recognition: 2808-2817.
- 1297 Zhang, X., J. Zhang, C. Yuan, S. Liu, Z. Chen and W. Li (2020). Locating induced earthquakes with a network of
1298 seismic stations in oklahoma via a deep learning method. *Scientific Reports*. **10**(1): 1941.
- 1299 Zhang, Z. and T. Alkhalifah (2019). Regularized elastic full-waveform inversion using deep learning. *Geophysics*.
1300 **84**(5): R741-R751.
- 1301 Zhang, Z., E. V. Stanev and S. Grayek (2020). Reconstruction of the basin - wide sea - level variability in the north
1302 sea using coastal data and generative adversarial networks. *Journal of Geophysical Research: Oceans*. **125**(12):
1303 e2020JC016402.
- 1304 Zhang, Z., H. Wang, F. Xu and Y. Jin (2017). Complex-valued convolutional neural network and its application in
1305 polarimetric SAR image classification. *IEEE Transactions on Geoscience and Remote Sensing*. **55**(12): 7177-7188.
- 1306 Zhao, M., S. Chen, L. Fang and A. Y. David (2019). Earthquake phase arrival auto-picking based on U-shaped
1307 convolutional neural network. *Chinese Journal of Geophysics*. **62**(8): 3034-3042.
- 1308 Zhong, Y., R. Ye, T. Liu, Z. Hu and L. Zhang (2020). Automatic aurora image classification framework based on
1309 deep learning for occurrence distribution analysis: A case study of all-sky image datasets from the yellow river
1310 station. *Journal of Geophysical Research*.
- 1311 Zhou, Y., H. Yue, Q. Kong and S. Zhou (2019). Hybrid event detection and phase-picking algorithm using
1312 convolutional and recurrent neural networks. *Seismological Research Letters*. **90**(3): 1079-1087.
- 1313 Zhu, J., T. Park, P. Isola and A. A. Efros (2017). Unpaired image-to-image translation using cycle-consistent
1314 adversarial networks. *International Conference on Computer Vision*: 2242-2251.
- 1315 Zhu, W., S. M. Mousavi and G. C. Beroza (2019). Seismic signal denoising and decomposition using deep neural
1316 networks. *IEEE Transactions on Geoscience and Remote Sensing*. **57**(11): 9476-9488.
- 1317

manuscript submitted to Reviews of Geophysics

Tables

Table 1 Examples of data-driven tasks in Geophysics

Examples of data-driven Tasks in Geophysics	
Modeling	Modeling the Earth with high spatial and temporal resolution
Spatial prediction	<p>Reconstruction</p> <ul style="list-style-type: none"> Global climate information based on limited measurements All-sky information from limited astronomy observation stations Both high resolution and large scale measurement in remote sensing <p>Inversion</p> <ul style="list-style-type: none"> High resolution subsurface structure using active seismic sources in exploration geophysics The Earth's structure based on passive earthquake measurements
Temporal prediction	<p>Forward predictiton</p> <ul style="list-style-type: none"> Rain fall nowcasting Typhoon track prediction Other natural disasters prediction in small time window <p>Backward prediction</p> <ul style="list-style-type: none"> The evolution of the Earth and the Universe in very large time window The drift of the continental
Detection	<p>Earthquake detection</p> <ul style="list-style-type: none"> Microearthquake detection Earthquake early warning <p>Pond coverage on Arctic sea ice, Coastal inundation mapping</p>
Classification	<p>Large spatial scale remote sensing imagery classification, Optical, Hyper-spectrum, SAR,</p> <p>Auroal classification</p>

manuscript submitted to Reviews of Geophysics

1319 Table 2 Examples of literature that use different network architectures for tasks beyond end-to-end training. Here
 1320 optimization oriented means using DNNs to optimize the traditional model-driven objective functions.

	CNN	CAE	U-Net	GAN	RNN
Supervised (End-to-end)	<u>Yu et al. 2019</u> <u>Dhara and Bagaini 2020</u>	<u>Wang et al. 2020</u>	<u>Yang and Ma 2019</u> <u>Wu et al. 2019</u>	<u>Siahkoohi et al. 2019</u>	<u>Yuan et al. 2020</u> <u>Linville et al. 2019</u>
Semi/ unsupervised		<u>Mousavi et al. 2019</u> <u>Duan et al. 2019</u>		<u>Niu et al. 2020</u>	
Optimization Oriented	<u>Xiao et al.</u>	<u>Sun and Alkhalifah 2020</u>			<u>Sun et al. 2020</u> <u>Wang et al. 2020</u>
Physical constraint	<u>Zhang et al. 2020</u>		<u>Wu and McMechan 2019</u>		
Uncertainty estimation	<u>Mousavi and Beroza 2020</u>				<u>Tasistro - Hart et al. 2020</u> <u>Grana et al. 2020</u>

1321

Figures

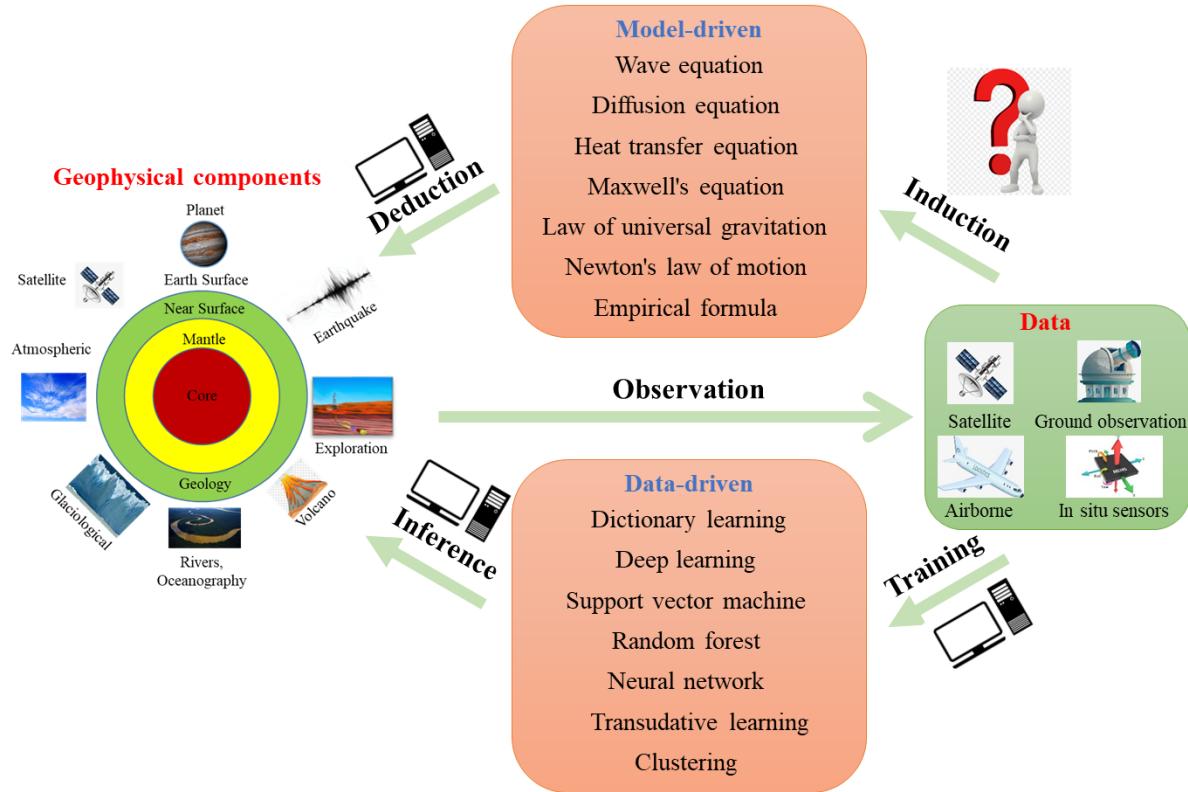


Figure 1 An illustration of model-driven and data-driven methods. On the left are the research topics in geophysics ranging from the Earth's core to the outer space. One the right is the observation means used at present. In the middle are examples of model-driven and data-driven methods. In model-driven methods, the principles of geophysical phenomena are induced from a large amount of observed data based on physical causality, then the models are used to deduct the geophysical phenomena in the future or in the past. In data-driven methods, the computer first inducts a regression or classification model without considering physical causality. Then, this model will perform tasks such as classification on incoming datasets.

1322
1323

manuscript submitted to Reviews of Geophysics

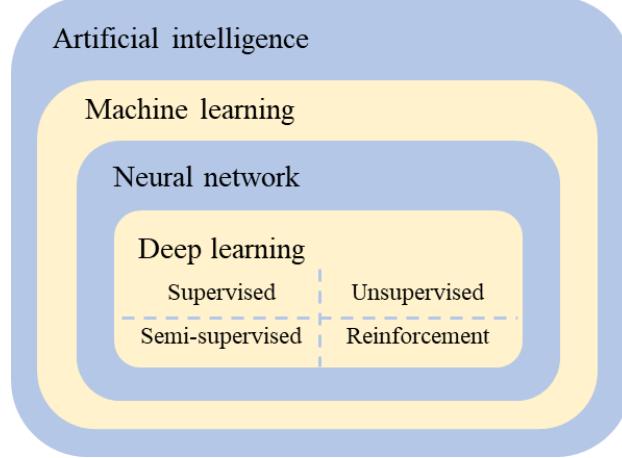


Figure 2 The containment relationship among artificial intelligence, machine learning, neural network and deep learning, and the classification of deep learning approaches.

1324

1325

manuscript submitted to Reviews of Geophysics

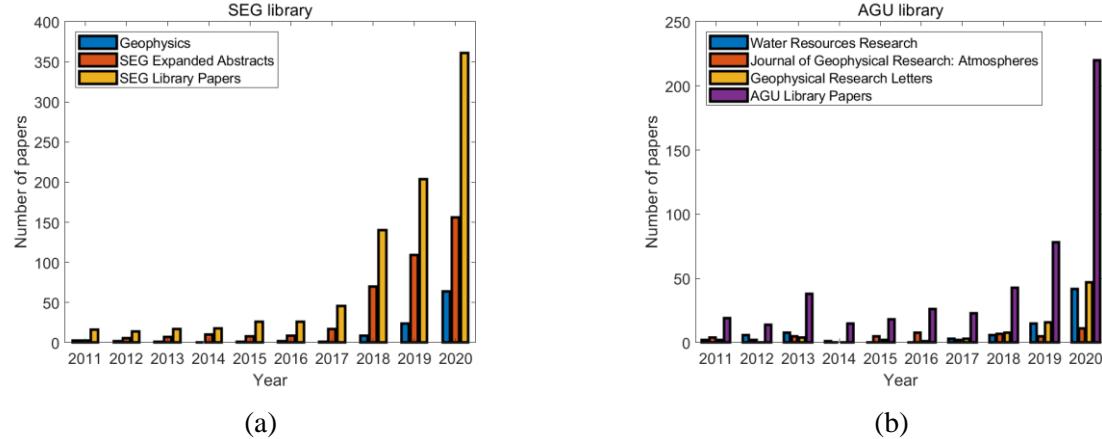


Figure 3 (a) and (b) are statics of AI-related papers in SEG Library and AGU Library. In (a), Geophysics means the flagship journal of SEG. SEG Expanded Abstracts means the Expanded Abstracts from SEG annual meeting. SEG Library papers mean the papers founded in the SEG digital library. In (b), the first three captions in the legend are the names of top journals in AGU. The fourth caption in the legend represents the papers founded in the AGU digital library.

1326

1327

manuscript submitted to Reviews of Geophysics

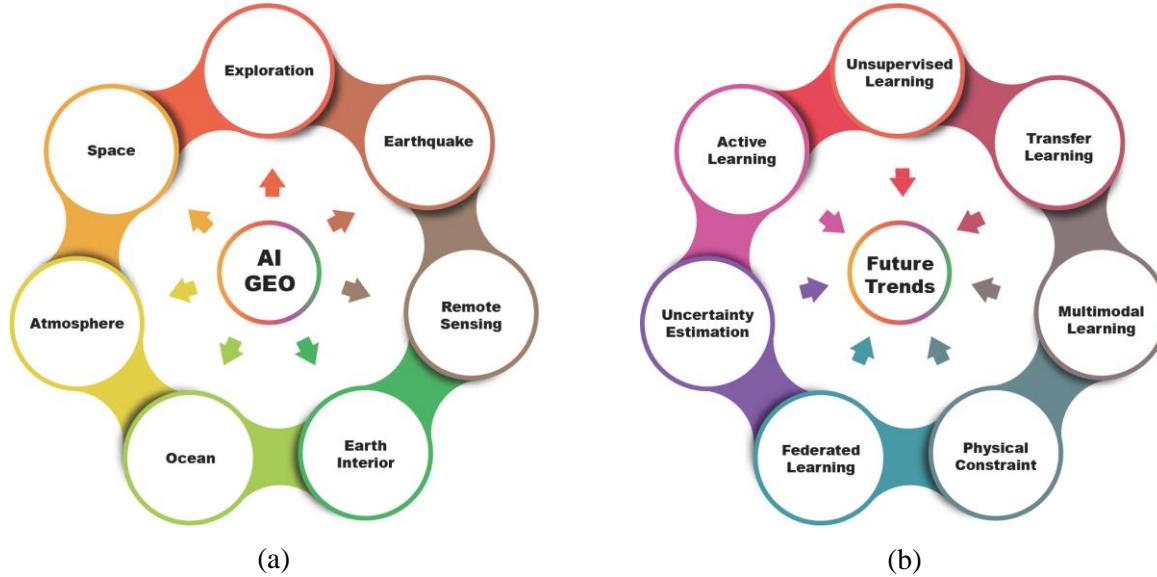


Figure 4 The topics included in this review. (a) DL-based geophysical applications. (b) The future trends of applying DL in geophysics.

1328

1329

1330

manuscript submitted to Reviews of Geophysics

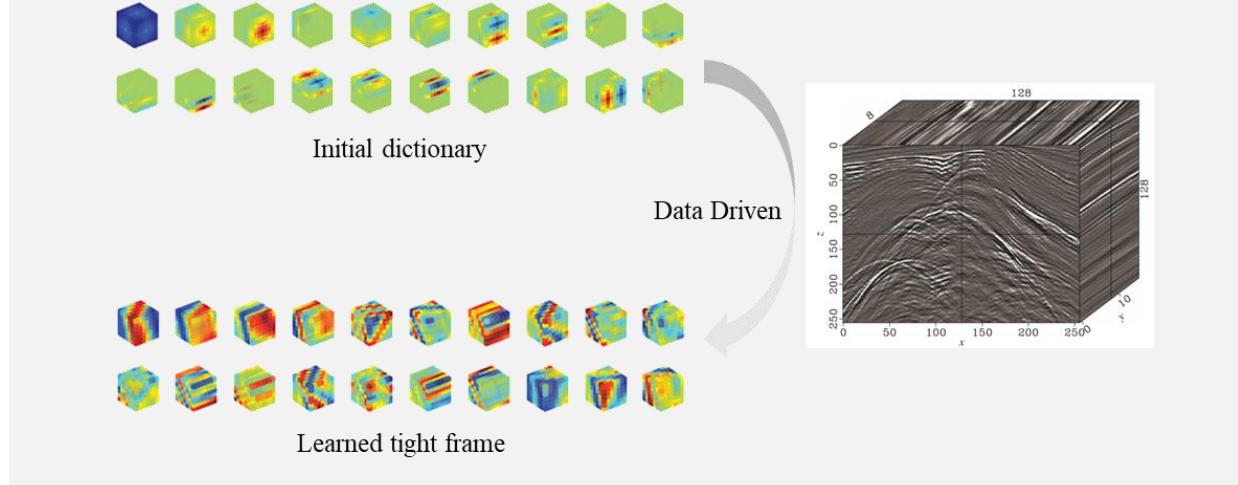


Figure 5. An illustration of dictionary learning: data-driven tight frame. The dictionary is initialized with a spline framelet. After training based on a post-stack seismic dataset, the trained dictionary exhibits apparent structures.

manuscript submitted to Reviews of Geophysics

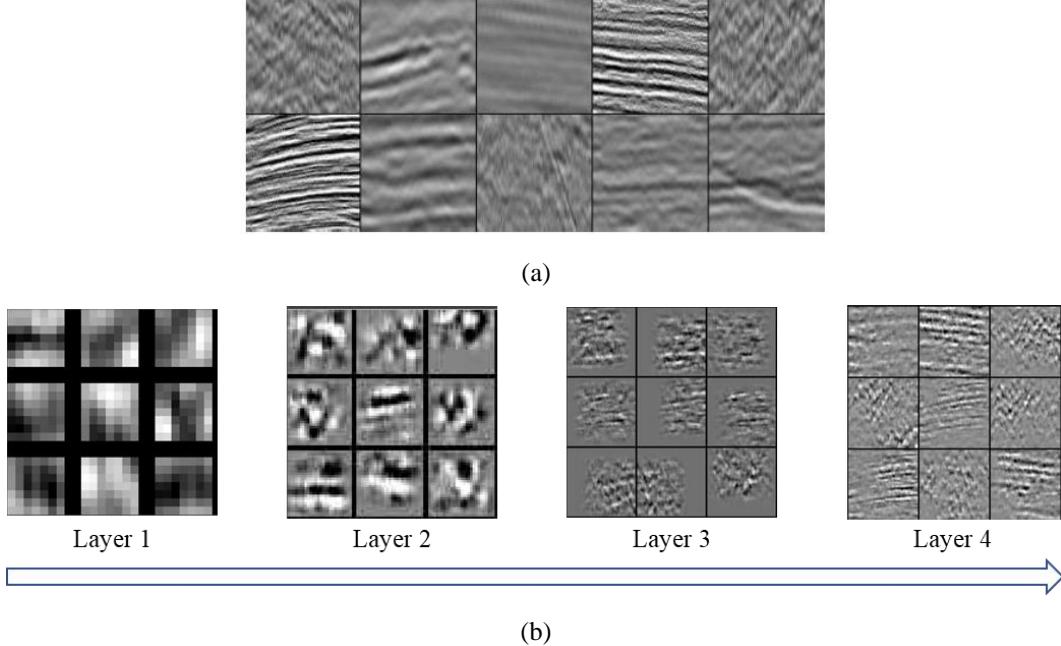


Figure 6. The learned features in deep learning. (a) Training samples. (b) In each layer, nine of the learned filters are shown. A great number of hierarchical structures are observed in different layers. Layer 1 exhibits edge structures, layer 2 shows small structures of seismic events, and layer 3 shows small portions of seismic sections. The filters in layer 2 and 3 are blank near edges, which may be caused by the boundary effect of the convolutional filter. Layer 4 gives larger seismic portions, which are approximations to the training data. The filters in layer 4 look more similar to each other than training datasets because DNN tries to learn the similar and hierarchical patterns which compose the data.

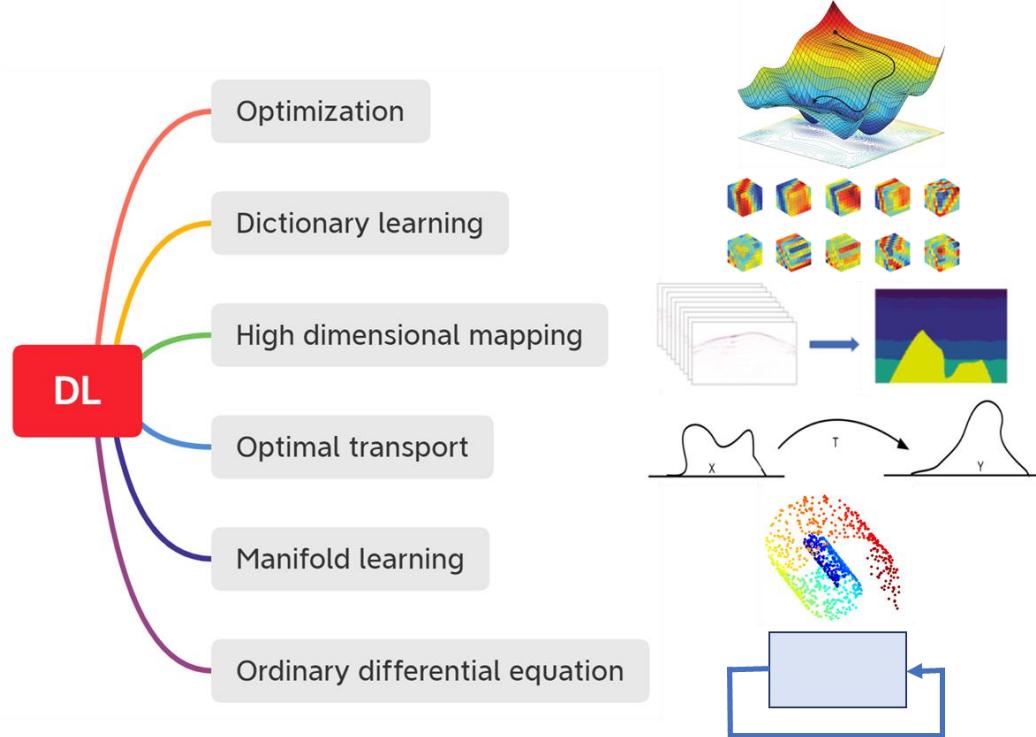


Figure 7. Understanding DL from different perspectives. Optimization: DL is basically a nonlinear optimization problem which solves for the optimized parameters to minimize the loss function of the outputs and labels. Dictionary learning: The filter training in DL is similar to that in dictionary learning. High dimensional mapping: DNN in DL is basically a high-dimensional mapping from the input to the labels. Optimal transport: a generative adversarial network can be interpreted by the theory of optimal transportation, which involves transformation between the given white noise and the data distribution. Manifold learning: The representation of training samples in the latent space of a DNN is similar to that learning a low dimensional manifold which contains all the data samples. Ordinary differential equation: a recurrent neural networks is basically a solution of an ordinary differential equation with the Euler method.

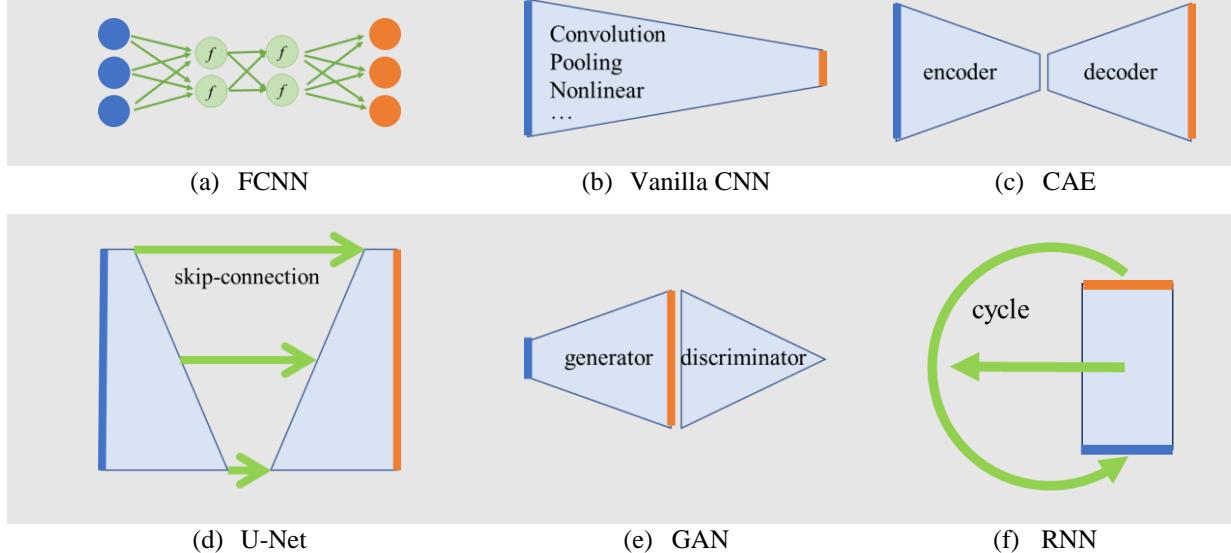


Figure 8. Sketches of DNNs. The blue lines indicate inputs, and the orange lines indicate outputs. The length of the blue and orange lines represents the data dimension. The green lines indicate intermedia connections. (a) In FCNN, the inputs of one layer are connected to every unit in the next layer. f stands for a nonlinear activation function. In (b)-(f), we omit the details of the layers and maintain the shape of each network architecture. (b) Vanilla CNN is cascaded by convolutional layers, pooling layers, nonlinear layer, and etc. In CNN, the outputs of the convolutional layers are either the same or smaller than the input depending on the strides used for convolution. Pooling layers will reduce the size of the extracted features. In regression or classification tasks, the output usually has the same dimension or a smaller dimension than the input (where (b) shows the latter situation). The difference between regression and classification is that the outputs are continuous variables in regression tasks and discrete variables representing categories in classification tasks. The dimension of the latent feature space in the CAE may be either larger or smaller than that of the data space, where (c) shows the latter. (d) Skip connections in U-Net are used to bring the low-level features to a high level. (e) In a GAN, low-dimensional random vectors are used to generate a sample from the generator, and then the sample is classified as true or false by the discriminator. (f) In an RNN, the output or hidden state of the network is used as input in a cycle.

1332
1333

manuscript submitted to Reviews of Geophysics

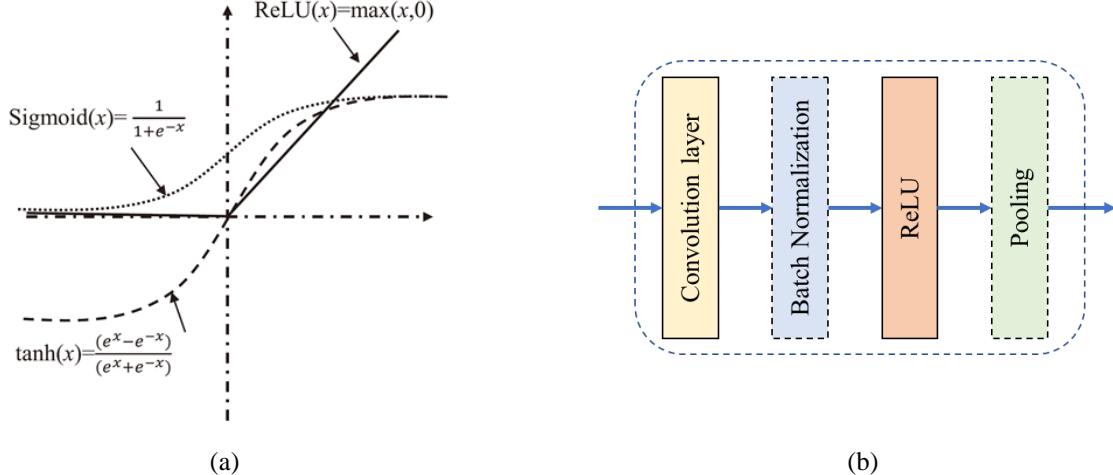


Figure 9. Details in DNN architectures. (a) Activation functions in the nonlinear layer. ReLU is commonly used since its gradient is easily computed and can avoid gradient vanishing. (b) A typical block in CNN. The convolutional layer and ReLU layer (nonlinear layer) are the basic components of one CNN block. The batch normalization layer can avoid gradient explosion. The pooling layer can extract features by subsampling the input.

1334

1335

1336

manuscript submitted to Reviews of Geophysics

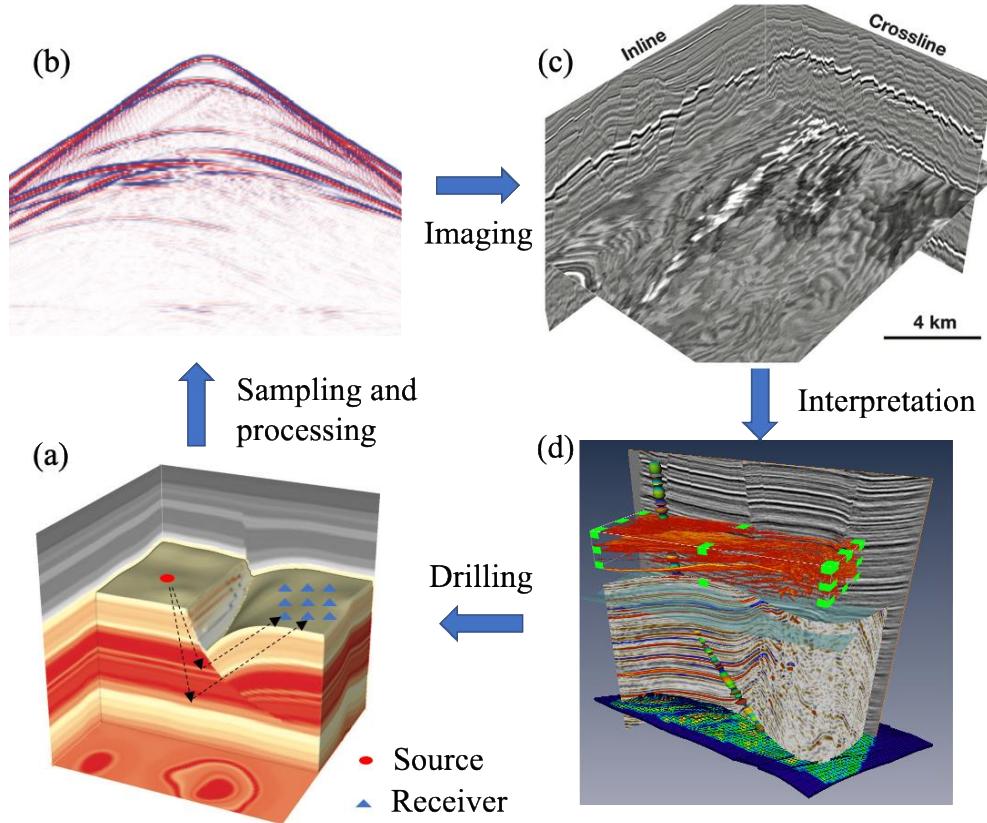
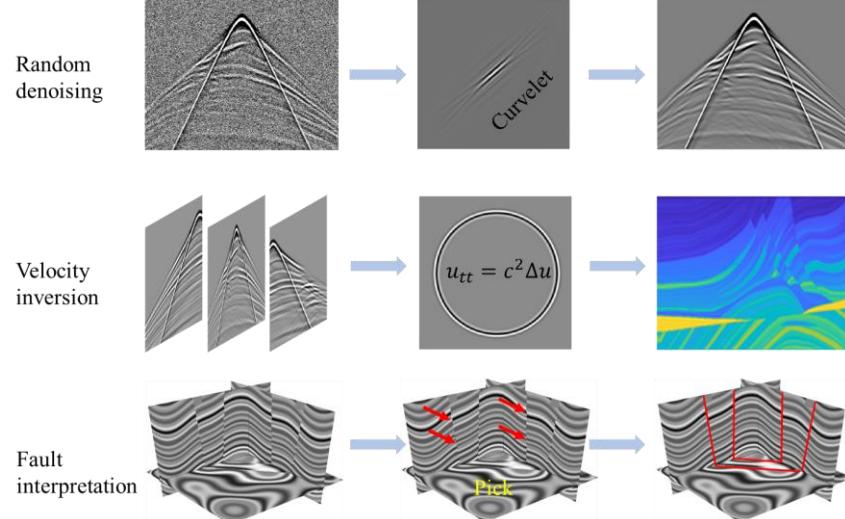
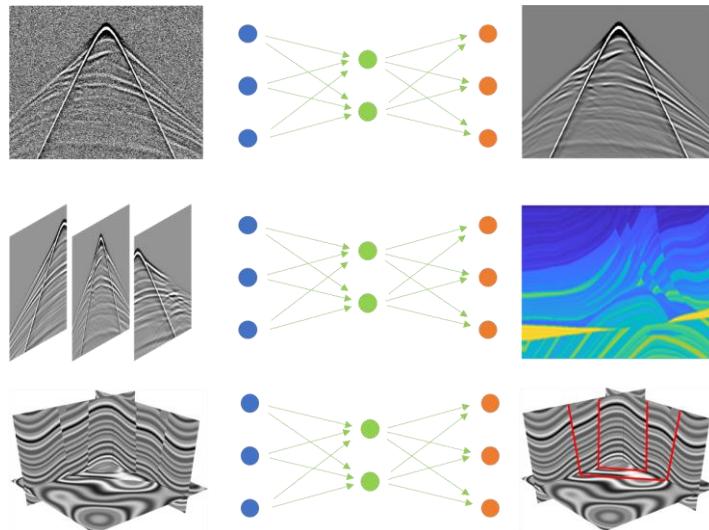


Figure 10. The procedure of exploration geophysics. (a) The subsurface structures. The seismic wave is excited at sources (red point) and propagates downward to the reflector and then propagates upwards until recorded by the receivers (blue points). (b) The seismic records are after processing. (c) The seismic imaging result, where the lines stand for the reflectors. (d) Underground properties are interpreted to determine where the reservoir locates.

manuscript submitted to Reviews of Geophysics



(a) Traditional exploration geophysics methods



(b) DL-based exploration geophysics methods

Figure 11. Comparison of traditional and DL-based methods in exploration geophysics. (a) In random denoising tasks, the curvelet denoising method (Herrmann and Hennenfent 2008) assumes that the signal is sparse under curvelet transform, and a matching method is used for denoising. In velocity inversion tasks, full-waveform inversion based on the wave equation is used for forward and adjoint modeling in the optimization algorithm. In fault interpretation tasks, faults are picked by interpreters. (b) The mentioned tasks are treated as regression problems that are optimized with neural networks. Different tasks may require different neural network architectures.

manuscript submitted to Reviews of Geophysics

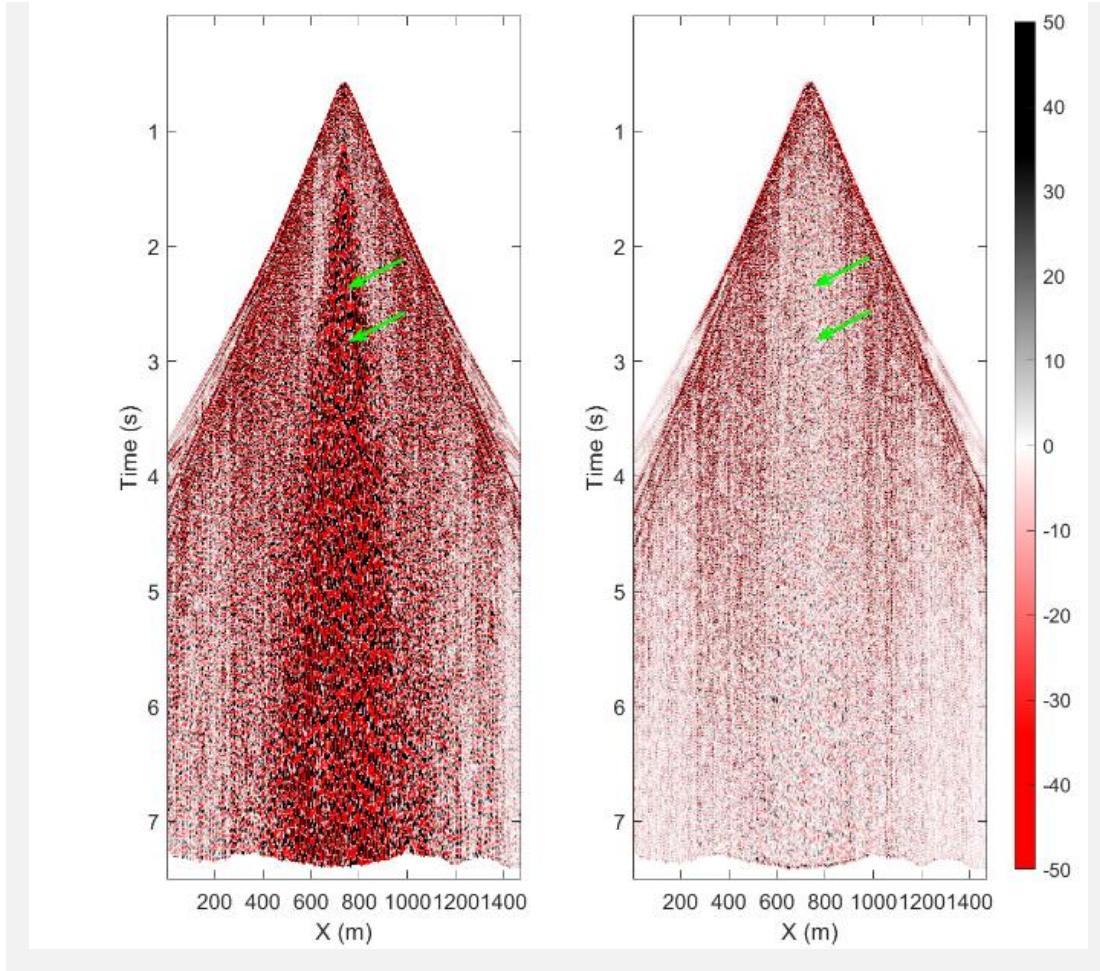


Figure 12. Deep learning for scattered ground-roll attenuation. On the left is the original noisy dataset. On the right is the denoised dataset. The scattered ground roll marked by the green arrows is removed.

manuscript submitted to Reviews of Geophysics

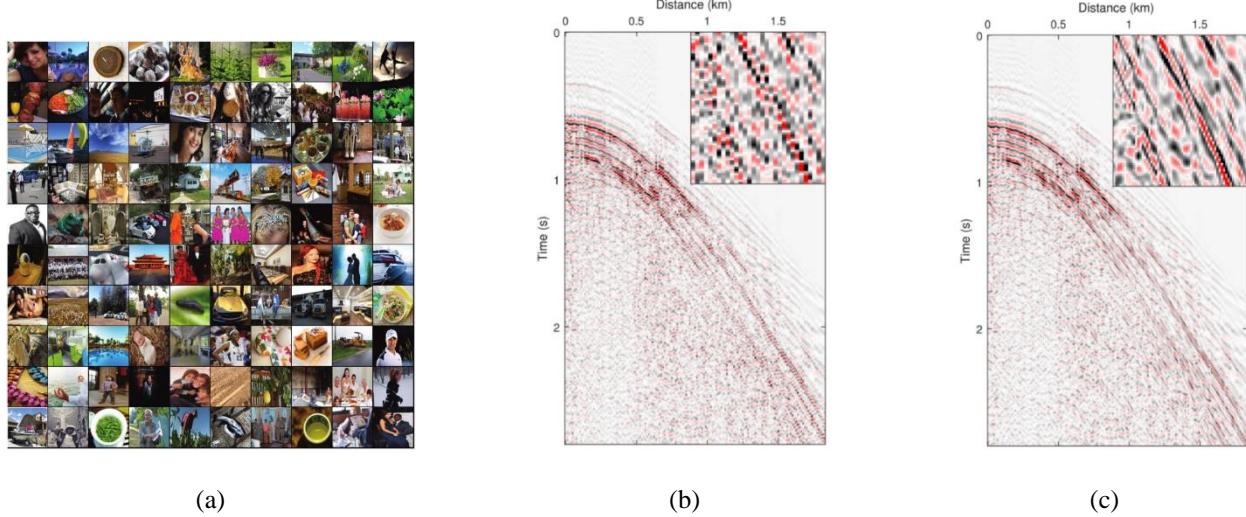


Figure 13. The training set and seismic interpolation result (Zhang et al. 2020). (a) A subset of the natural image dataset. The natural image dataset was used to train a network for seismic data interpolation. (b) An under-sampled seismic record. (c) The interpolated record corresponding to (b). The regions 1.6-1.88 s and 1.0-1.375 km are enlarged at the top-right corner.

manuscript submitted to Reviews of Geophysics

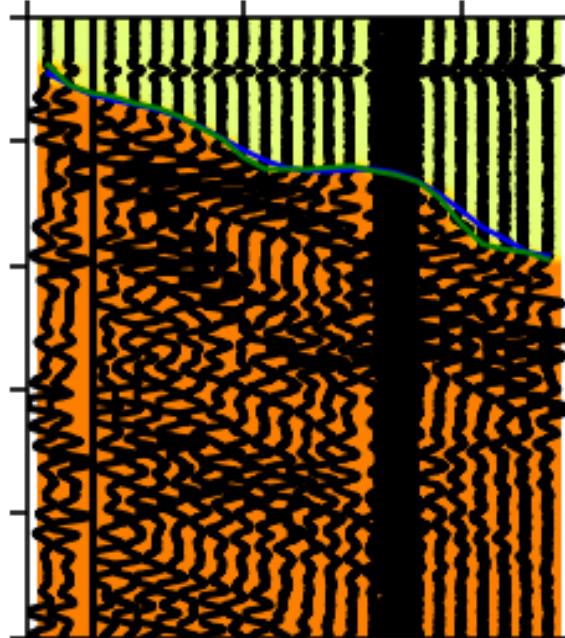


Figure 14. Phase picking based on U-Net. The inputs are seismological data. The outputs are zeros above the first arrival in the green area, ones below the first arrival in the yellow area, and twos for the first arrival on the blue line. The green line indicates the predicted first arrival. This experiment was performed based on the modified code from https://github.com/DaloroAT/first_break_picking.

manuscript submitted to Reviews of Geophysics

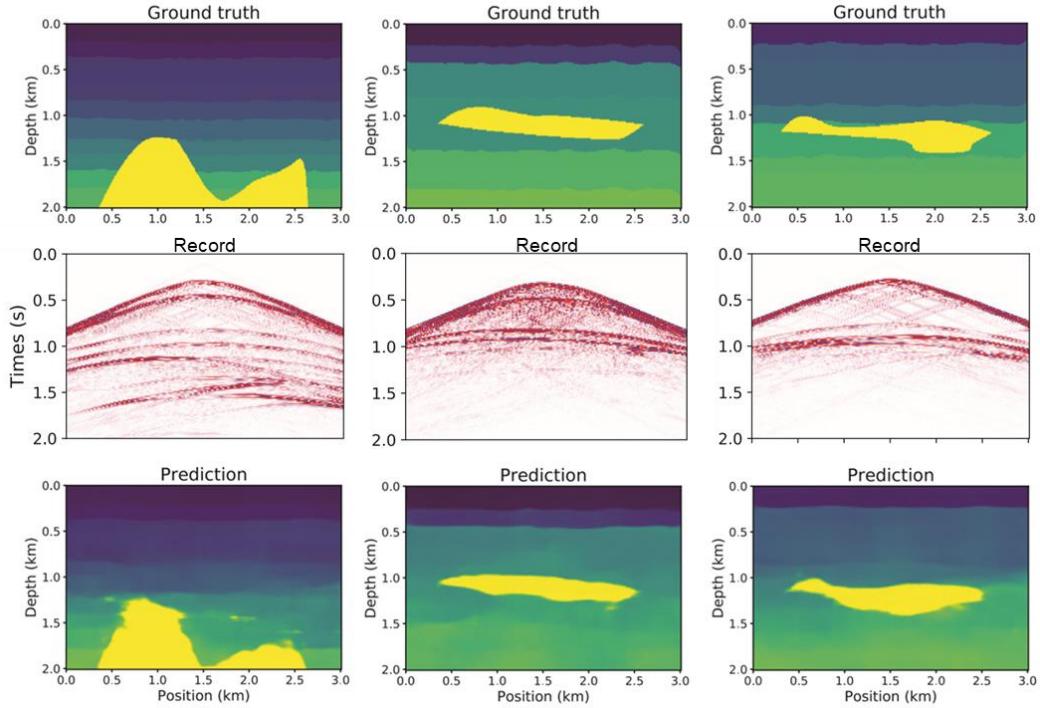


Figure 15. Predicting the velocity model with U-Net from raw seismological data ([Yang and Ma 2019](#)). The columns indicate different velocity models. From top to bottom are the ground truth velocity models, generated seismic records from one shot, and the predicted velocity models.

manuscript submitted to Reviews of Geophysics

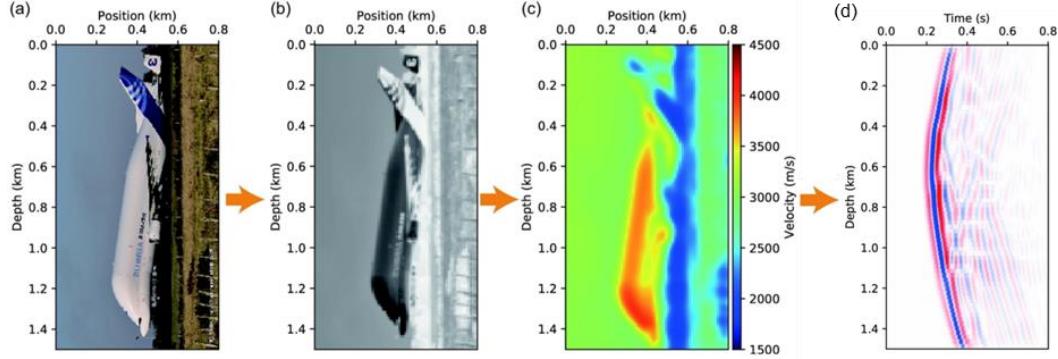


Figure 16. Converting a three-channel color image into a velocity model ([Wang and Ma 2020](#)). (a)-(c) are original color image, grayscale image, and corresponding velocity model. (d) is the seismic record generated from a cross-well geometry on (c).

manuscript submitted to Reviews of Geophysics

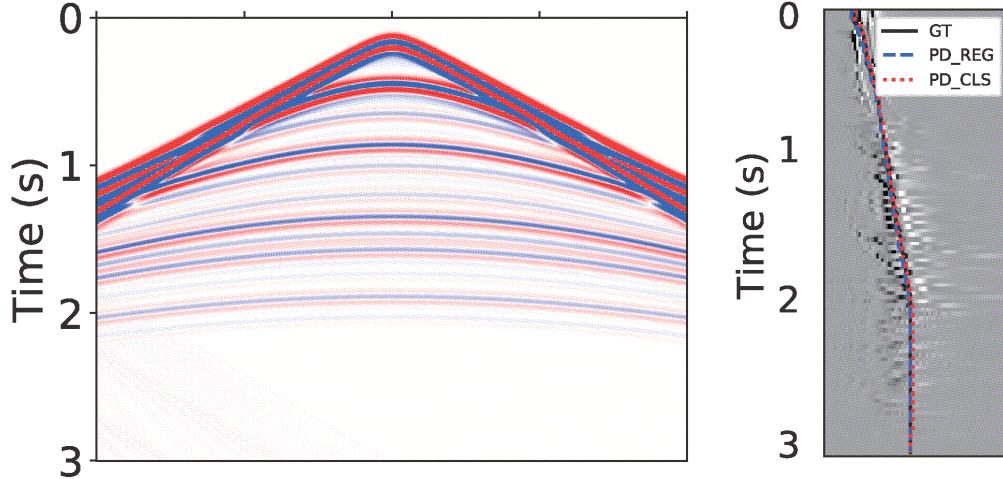


Figure 17. Velocity picking based on U-Net. The inputs are seismological data on the left. The outputs are the picking positions on the right. GT means ground truth. PD_REG and PD_CLS represent the velocity predictions of the regression network and classification network, respectively.

manuscript submitted to Reviews of Geophysics

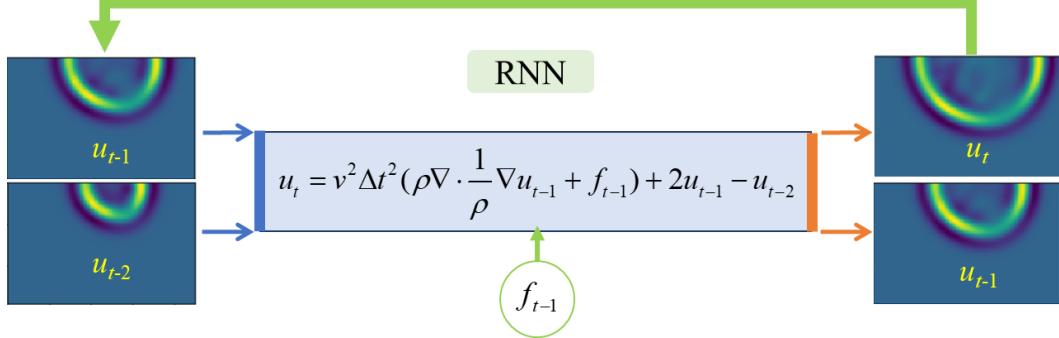


Figure 18. Modified RNN based on the acoustic wave equation for wave modeling ([Liu 2020](#)). The diagram represents the discretized wave equation implemented in an RNN. The auto-differential mechanics of a DNN help to efficiently optimize the velocity and density.

manuscript submitted to Reviews of Geophysics

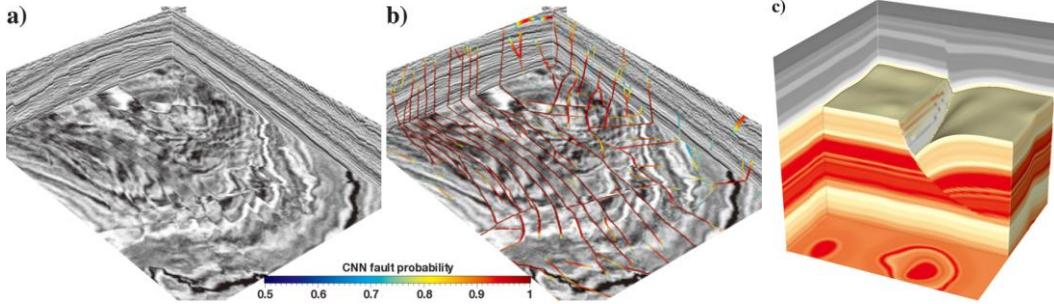


Figure 19. (a) A post-stack dataset. (b) Fault prediction result of (a). (c) A synthetic dataset ([Wu et al. 2020](#)).

manuscript submitted to Reviews of Geophysics

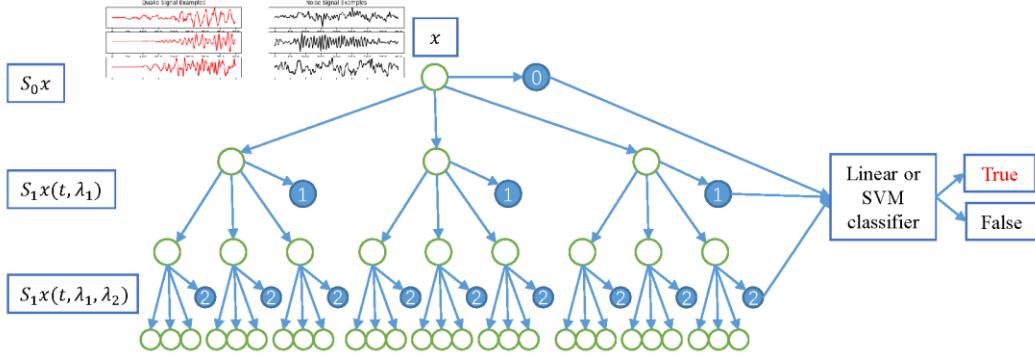


Figure 20. (a) The architecture of WST. Unlike in a CNN, the outputs of WST are combined with the outputs of each layer. Then, the outputs of WST serve as features for a classifier.

manuscript submitted to Reviews of Geophysics

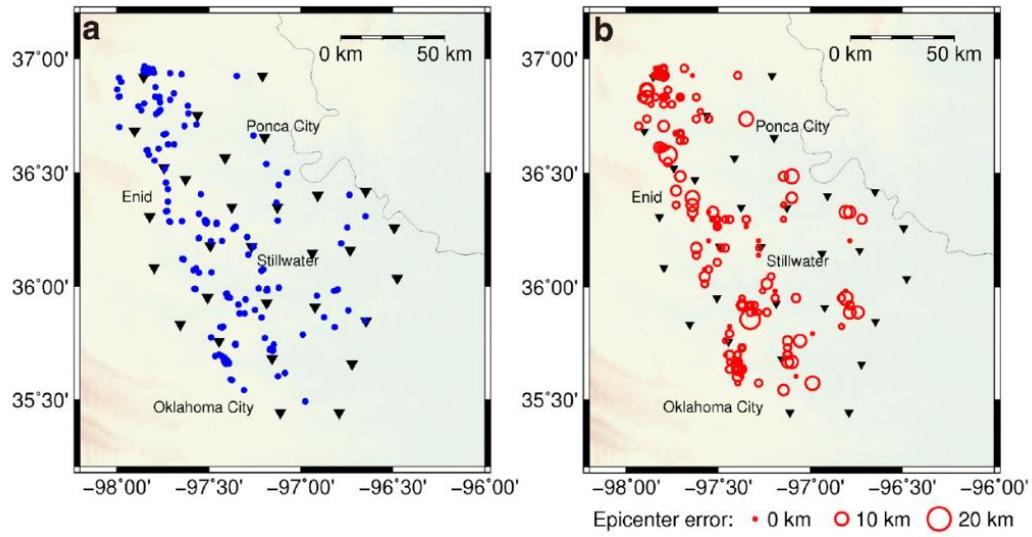


Figure 21. Locating earthquake sources with deep learning. The black triangles are stations. Left: the blue dots are the actual locations. Right: the red circles are the predicted locations. The radius of a circle represents the predicted epicenter error ([Zhang et al. 2020](#)).

manuscript submitted to Reviews of Geophysics

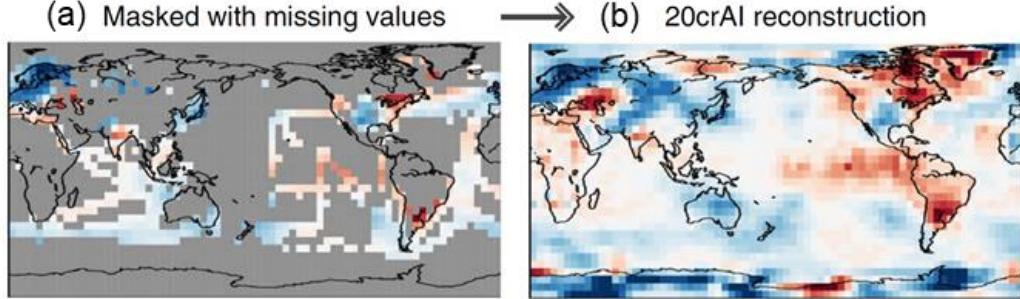


Figure 22 AI models reconstruct temperature anomalies with many missing values (Kadow et al. 2020).

manuscript submitted to Reviews of Geophysics

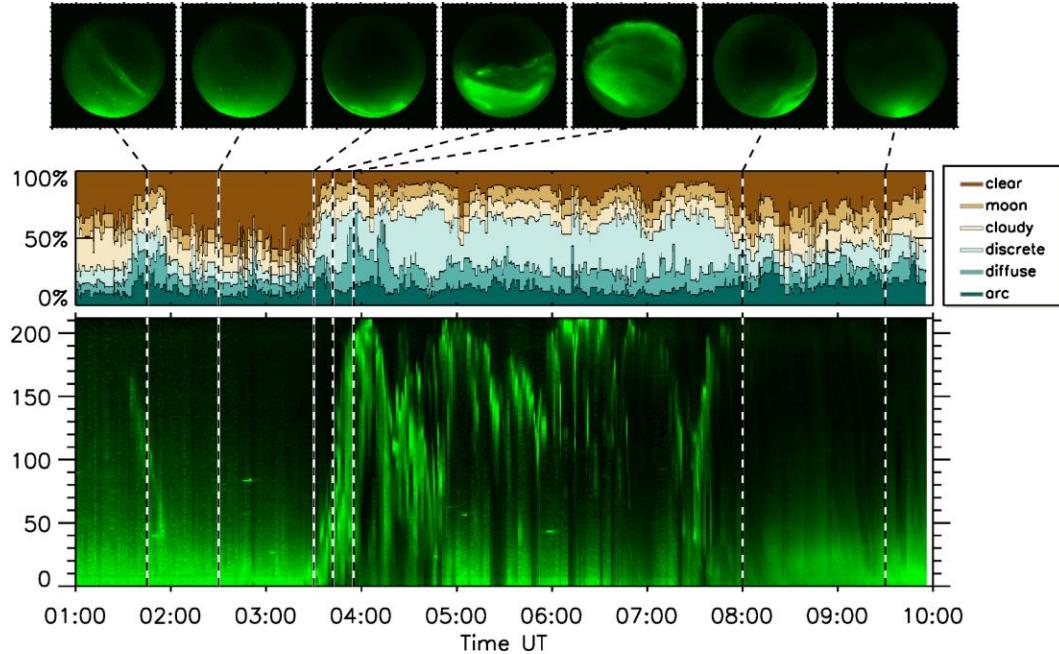


Figure 23 The bottom panel shows a keogram from auroral data collected on 21 January 2006 at Rankin Inlet. The keogram consists of a single column from the auroral images at different times. The middle panel shows the probabilities for the six categories as predicted by the ridge classifier trained with the entire training dataset. At the top are auroral images at different times. ([Clausen and Nickisch 2018](#))

manuscript submitted to Reviews of Geophysics

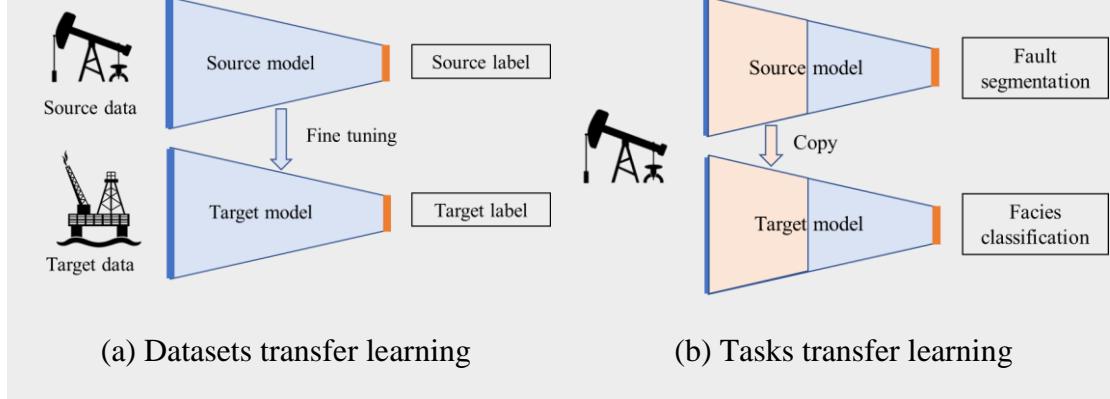


Figure 24. Diagrams of transfer learning. (a) Transfer learning between different datasets. The parameters of one trained model can be moved to another model as initialization conditions. (b) Transfer learning between different tasks. The first layers of one trained model can be copied to another model.

manuscript submitted to Reviews of Geophysics

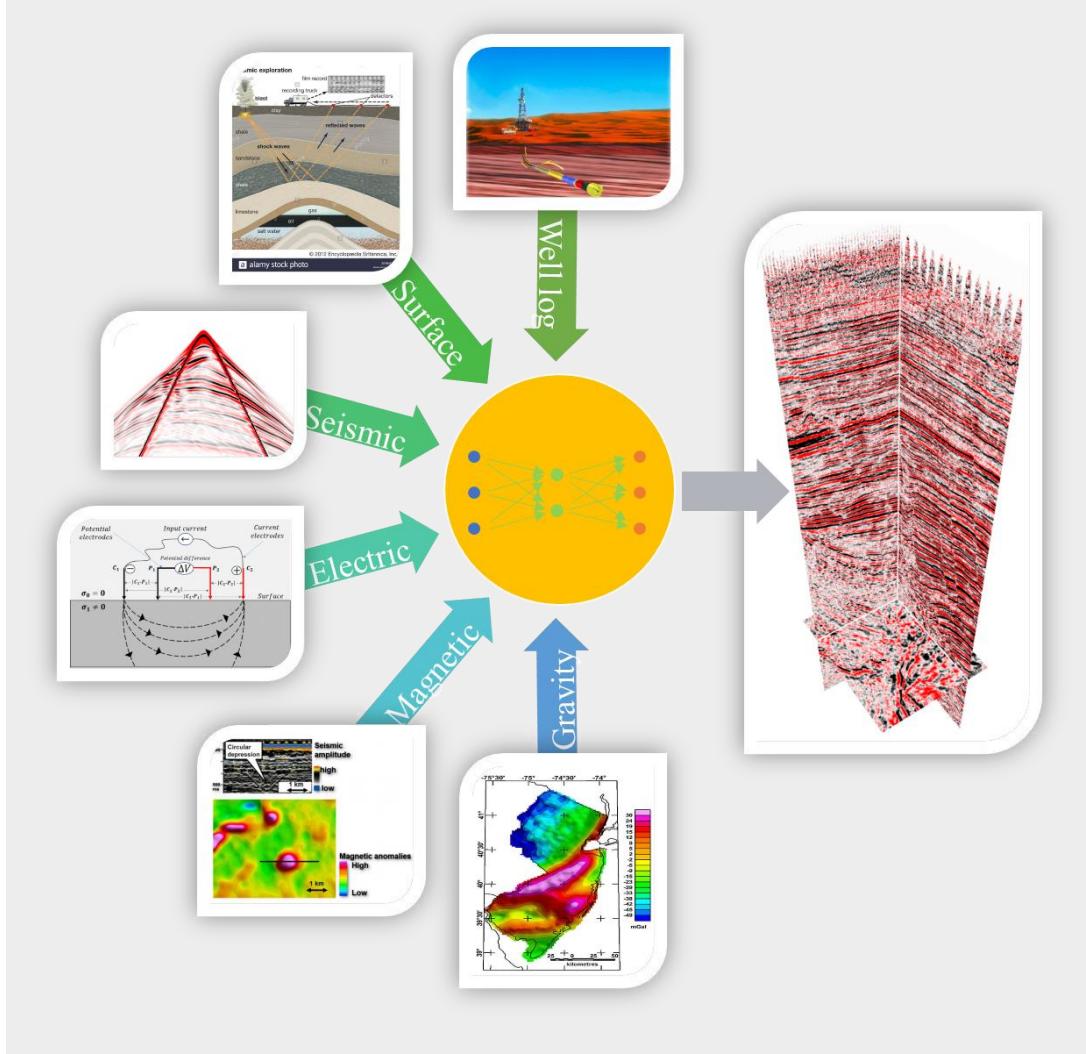


Figure 25. An illustration of multimodal deep learning

manuscript submitted to Reviews of Geophysics

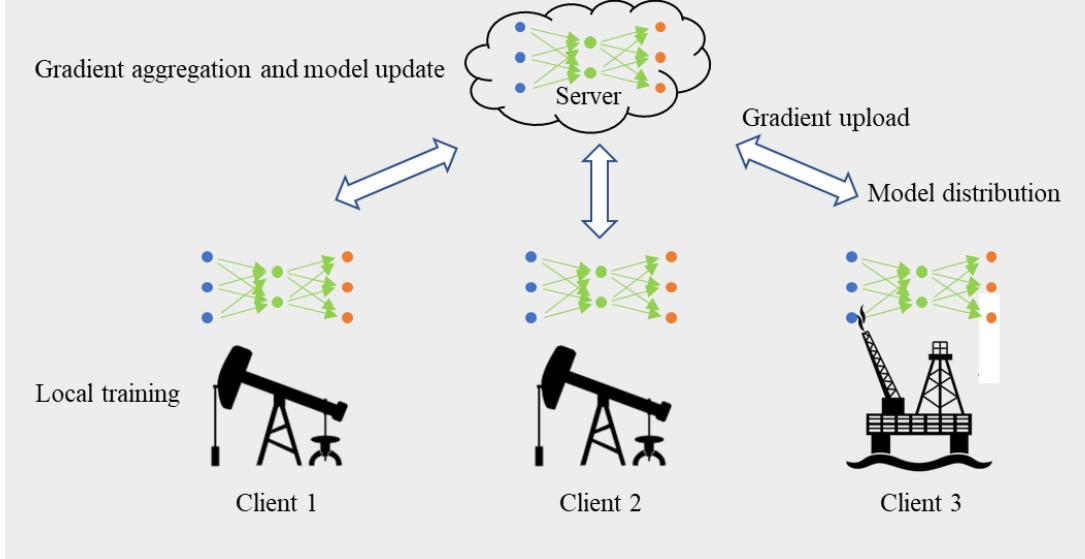


Figure 26. Federated learning. The clients train the DNN with local datasets and uploads the model gradient to the server. The server aggregates the gradients and updates the global model. Then, the updated model is distributed to all the local clients. Many rounds of training are performed until the model meets a certain accuracy requirement.

manuscript submitted to Reviews of Geophysics

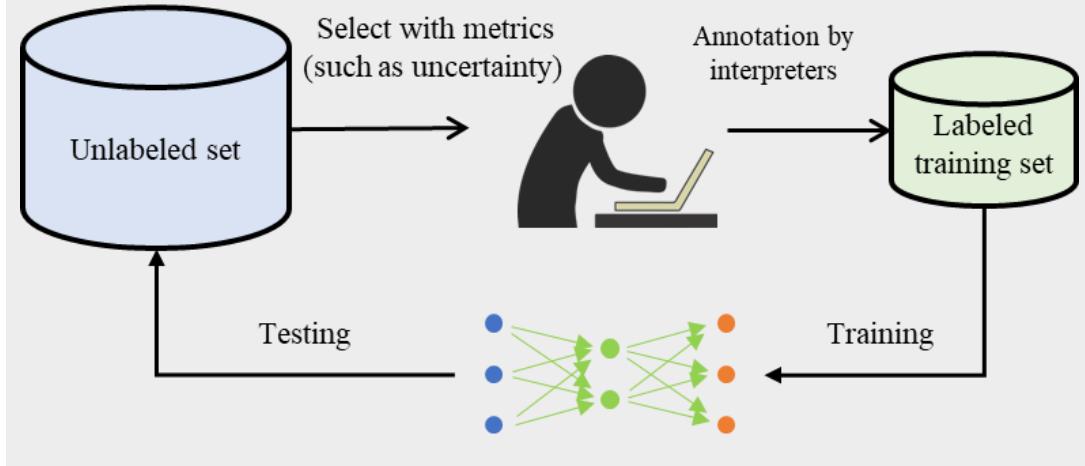


Figure 27. An illustration of active learning. We choose samples with high uncertainty and manually annotate them to serve as training samples.

1343
1344
1345