A Little Data goes a Long Way: Automating Seismic Phase Arrival Picking at Nabro Volcano with Transfer Learning

Sacha Lapins^{1,1}, Berhe Goitom^{1,1}, J.-Michael Kendall^{2,2}, Maximilian J. Werner^{1,1}, Katharine V. Cashman^{1,1}, and James Hammond^{3,3}

¹University of Bristol ²University of Oxford ³Birkbeck, University of London

November 30, 2022

Abstract

Supervised deep learning models have become a popular choice for seismic phase arrival detection. However, they don't always perform well on out-of-distribution data and require large training sets to aid generalization and prevent overfitting. This can present issues when using these models in new monitoring settings. In this work, we develop a deep learning model for automating phase arrival detection at Nabro volcano using a limited amount of training data (2498 event waveforms recorded over 35 days) through a process known as transfer learning. We use the feature extraction layers of an existing, extensively-trained seismic phase picking model to form the base of a new all-convolutional model, which we call U-GPD. We demonstrate that transfer learning reduces overfitting and model error relative to training the same model from scratch, particularly for small training sets (e.g., 500 waveforms). The new U-GPD model achieves greater classification accuracy and smaller arrival time residuals than off-the-shelf applications of two existing, extensively-trained baseline models for a test set of 800 event and noise waveforms from Nabro volcano. When applied to 14 months of continuous Nabro data, the new U-GPD model detects 31,387 events with at least four P-wave arrivals and one S-wave arrival, which is more than the original base model (26,808 events) and our existing manual catalogue (2,926 events), with smaller location errors. The new model is also more efficient when applied as a sliding window, processing 14 months of data from 7 stations in less than 4 hours on a single GPU.

1	A Little Data goes a Long Way: Automating Seismic Phase Arrival Picking at Nabro	
2	Volcano with Transfer Learning	
3		
4	Enter authors here: Sacha Lapins ¹ , Berhe Goitom ¹ , J-Michael Kendall ² , Maximilian J.	
5	Werner ¹ , Katharine V. Cashman ¹ and James O. S. Hammond ³	
6	¹ School of Earth Sciences, University of Bristol, UK.	
7	² Department of Earth Sciences, University of Oxford, UK.	
8	³ Department of Earth and Planetary Sciences, Birkbeck, University of London, UK.	
9		
10	Corresponding author: Sacha Lapins (sacha.lapins@bristol.ac.uk)	
11		
12	Key Points:	
13	• Transfer learning using existing model trained on California earthquake data produces	
14	effective new model for monitoring at Nabro volcano	
15	• Nabro transfer learning model shows improved S-wave picking resulting in smaller	
16	location errors than even manual phase picks	
17	• Changing task from classification to segmentation results in more efficient model	
18	processing 14 months of data from 7 stations in 4 hours	
19		

20 Abstract

21

Supervised deep learning models have become a popular choice for seismic phase arrival 22 detection. However, they don't always perform well on out-of-distribution data and require large 23 training sets to aid generalization and prevent overfitting. This can present issues when using these 24 models in new monitoring settings. In this work, we develop a deep learning model for automating 25 phase arrival detection at Nabro volcano using a limited amount of training data (2498 event 26 waveforms recorded over 35 days) through a process known as transfer learning. We use the 27 feature extraction layers of an existing, extensively-trained seismic phase picking model to form 28 the base of a new all-convolutional model, which we call U-GPD. We demonstrate that transfer 29 learning reduces overfitting and model error relative to training the same model from scratch, 30 particularly for small training sets (e.g., 500 waveforms). The new U-GPD model achieves greater 31 classification accuracy and smaller arrival time residuals than off-the-shelf applications of two 32 existing, extensively-trained baseline models for a test set of 800 event and noise waveforms from 33 Nabro volcano. When applied to 14 months of continuous Nabro data, the new U-GPD model 34 detects 31,387 events with at least four P-wave arrivals and one S-wave arrival, which is more 35 than the original base model (26,808 events) and our existing manual catalogue (2,926 events), 36 with smaller location errors. The new model is also more efficient when applied as a sliding 37 window, processing 14 months of data from 7 stations in less than 4 hours on a single GPU. 38

39

40 Plain Language Summary

41

Seismic monitoring increasingly relies on automated signal processing as the rate of data 42 acquisition grows. Supervised deep learning models have proven to be effective for detecting and 43 characterizing seismic events, but training such highly parameterized models generally requires 44 large amounts of manually labelled data. Once trained, however, these models extract general 45 seismic waveform features that can be used to train new models with more limited training data. 46 In this work, we use the generalized knowledge of seismic data from a model trained on millions 47 of earthquakes in California to train a new model for detecting volcanic earthquakes at Nabro 48 volcano, Eritrea, a recently active and, prior to its 2011 eruption, poorly monitored volcano. Using 49

a small training set of waveforms, the new model more accurately detects phase arrivals and noise than off-the-shelf applications of two baseline models. The new model is efficient, processing 14 months of data in less than 4 hours. It is also effective, detecting more volcanic events and showing improved levels of S-wave arrival picking. The result is smaller event location errors than even our manual picks. This level of efficiency and consistency highlights the role that machine learning can play in volcano-seismic monitoring.

56

57 1 Introduction

58

Seismic monitoring plays a fundamental part in mitigating hazards at volcanoes. During 59 periods of unrest, thousands of earthquakes can occur each day, producing a diverse range of 60 seismic signals that reflect a multitude of interlinked volcanic processes (e.g., migrating fluids, 61 fault movement, explosions, rockfalls). These earthquakes are generally recorded by broadband 62 seismometers, which are highly sensitive to ground motion across a wide range of frequencies and 63 record signals at high sample rates (typically 100 times or more per second). This level of detail, 64 however, comes at the cost of generating vast amounts of data. Many seismic networks utilize tens 65 or even hundreds of seismometers at a given time (e.g., Hansen & Schmandt, 2015), making real-66 time manual inspection of these time series practically infeasible. Previous seismic deployments 67 have also generated extensive legacy datasets that can offer insights into historical volcanic activity 68 69 and opportunities to further our understanding of volcanic processes. The main challenge is therefore to identify and characterize volcanic earthquakes in a robust and timely manner so as to 70 71 provide vital clues regarding the state of a volcano and the likelihood or impact of an eruption or 72 hazard, as well as be able to accurately and efficiently process large existing datasets for further analysis within a reasonable timeframe. 73

74

Identifying earthquake phase arrivals, particularly the initial primary (P-) and secondary/shear (S-) wave arrivals, forms the basis of most seismic processing tasks (e.g., determining locations, magnitudes and source parameters). Manually identifying these phase arrivals yields greater accuracy and estimates of arrival time uncertainty than automated approaches but is extremely time-consuming. Alternatively, most automated approaches are orders

of magnitude quicker but typically require clear phase arrivals, existing 'templates' of previously 80 catalogued earthquakes (e.g., Gibbons & Ringdal, 2006; Lengliné et al., 2016; Shelly et al., 2007), 81 or pre-processing / feature extraction steps calibrated for a small set of earthquake characteristics 82 (e.g., trigger algorithms based on the ratio of short-term average to long-term average signal 83 amplitude, STA/LTA; Withers et al., 1998). A challenge for application to volcanology is that 84 volcanic earthquakes can exhibit widely varying time-frequency characteristics, often with low 85 amplitudes or obscured phase arrivals, and new phases of unrest can produce previously unseen 86 seismic signals that differ from existing earthquake templates. Furthermore, methods based on 87 existing seismic catalogues are unsuitable for new seismic deployments where a catalogue of 88 events has not been collected. 89

90

91 A recently successful approach for seismic phase arrival detection is the use of supervised deep learning models (e.g., Dokht et al., 2019; Mousavi et al., 2019; Ross et al., 2018b; Woollam 92 93 et al., 2019; Zhu & Beroza, 2019). These methods are based on convolutional neural networks (CNN), a variant of classical neural networks that employ convolution operations, as opposed to 94 95 matrix multiplication, in at least part of the model. These operations are employed in 'hidden' convolutional layers that allow the network to learn a large set of filters to extract useful features 96 97 from the input data and map them to a desired output (e.g., to identify phase arrivals in earthquake waveforms; Fig 1). Typically, multiple convolutional layers are applied in succession and in 98 99 combination with other operations, such as non-linear 'activation', down-sampling and normalization, to extract complex patterns from the data using a hierarchy of simpler filter kernels. 100 These extracted features can then be fed into a standard fully-connected neural network or other 101 machine learning architecture for classification, segmentation, regression, clustering or inference 102 (e.g., Mousavi et al., 2019; Ross et al., 2018b; van den Ende & Ampuero, 2020). As such, the 103 104 'convolutional' part of CNNs act as the model's feature extraction system. With each successive convolutional layer, the extracted features move from lower-level, general signal features 105 (resembling, for example, long/short period wavelets in seismological waveform models; Fig 1A 106 inset) to more task specific, high-level features (Yosinski et al., 2014). The final 'classification' 107 layers of the model map these features to the desired output and can be considered the most task 108 specific part of the model, empirically tuned to the distribution of the training data (Yosinski et 109 al., 2014). 110

Such an approach gives supervised deep learning models a strong advantage over 112 traditional algorithms that require considerable manual intervention or rely on a small set of 113 manually determined characteristics and simple threshold criteria. In general, however, these 114 models require substantial amounts of labelled data during training to generalize to out-of-sample 115 data (the amount dependent on various factors, such as network architecture, number of network 116 parameters and training hyperparameters; e.g., D'souza et al., 2020; He et al., 2019; Sun et al., 117 2017). In the case of seismological supervised models, these models can demonstrate impressive 118 levels of generalization to phase arrival detection in other geographic and tectonic settings, if 119 trained with sufficient data (e.g., Mousavi et al., 2020; Tan et al., 2021). However, as with 120 practically any deep learning model, they can also suffer significant loss in performance when 121 122 faced with data that differs in source or distribution from their training data (e.g., Barbedo, 2018; Zech et al., 2018; Fig 7). As such, the requirement for extensive training sets can place the 123 traditional paradigm of supervised learning (i.e., using a large amount of hand-labelled data to 124 train a single model for a desired domain or problem) out of reach for many real-world 125 126 applications.

127

Transfer learning is based on the idea of knowledge transfer from one task to another (Pan 128 & Yang, 2010; Zhuang et al., 2020) and can be a powerful tool when we do not have sufficient 129 130 labelled data to train a reliable model from scratch, or when existing models perform poorly. At its simplest, the first *n* convolutional layers and their weights from the feature extraction part of an 131 existing model are copied to the first *n* layers of a new model for a related or similar task, with the 132 remaining layers either re-initialized with randomized weights or replaced (e.g., Razavian et al., 133 2014; Yosinski et al., 2014). These tasks need not be near-identical or even superficially related, 134 135 as long as low-level data characteristics are shared between tasks (e.g., Efremova et al., 2019; Tran et al., 2020; Zamir et al., 2018). The intuition is that generalized knowledge of data structure and 136 properties from one model trained with abundant labelled data (or 'big data') can guide a learning 137 algorithm towards a good solution for a new task with far more limited, or even no, labelled data. 138 139

In this paper, we evaluate the utility of inductive transfer learning (i.e., when labelled data are available for both the source and target tasks) for small seismic training sets and produce a

deep learning model that accurately and robustly picks phase arrivals from a deployment at Nabro 142 volcano in Eritrea, a region with little or no prior seismic monitoring. We leverage the knowledge 143 acquired from training a model on millions of seismic waveforms recorded by the Southern 144 California Seismic Network (SCSN), hereby referred to as the GPD model (Generalized seismic 145 Phase Detection; Ross et al., 2018b), and apply it to seismograms from Nabro volcano in Eritrea, 146 for which we have limited hand-labelled data (manual phase arrival picks) from the first couple of 147 months of a 14-month seismic deployment (Goitom, 2017; Hamlyn et al., 2014). The new model 148 task differs from the original GPD model task in that it is modified from one of classification 149 (assigning a single class label *P-wave*, *S-wave* or *noise* to an entire 4-second waveform; Fig 1A) 150 to one of segmentation (assigning a class label P-wave, S-wave or noise to each datapoint within 151 that 4-second waveform; Fig 1B). We achieve this by replacing the fully-connected uppermost 152 153 layers of the original GPD model with further convolutional layers, creating an all-convolutional model commonly referred to as a U-Net (Ronneberger et al., 2015). We refer to this specific model 154 155 design as the U-GPD model, utilizing GPD model weights within a U-Net architecture. The new data from Nabro volcano also exhibit differences in instrument calibration and sample rates from 156 157 the original GPD model training data, as well as differing waveform characteristics between tectonic and volcanic event types (Lahr et al., 1994; Lapins et al., 2020; McNutt & Roman, 2015). 158 159

In the following section, we introduce transfer learning and recent applications in 160 161 seismological deep learning. In Sections 3 and 4, we present our proposed transfer learning 162 method, U-GPD model architecture and seismic data recorded at Nabro volcano. In Section 5, we present a series of model comparisons. We first use common training metrics to demonstrate that 163 transfer learning reduces overfitting and model error, particularly for very small training sets (< 164 1000 waveforms), when compared with a model reinitialized with randomized weights before 165 166 training (i.e., trained from scratch with no transfer learning). We then apply these new models to a test dataset of known P-/S-wave arrivals and sections of noise and compare performance with 167 off-the-shelf applications of the base GPD model and another extensively-trained phase-picking 168 model, PhaseNet (Zhu & Beroza, 2019). We find that the U-GPD transfer learning model yields 169 improved phase arrival identification, particularly for S-waves, and false detection rate at Nabro 170 volcano. Altering the model task from classification to segmentation also improves pick time 171 residuals over the base GPD model for these test data. Finally, we apply both our new U-GPD 172

transfer learning model and the original base GPD model to the full 14-month seismic deployment 173 174 at Nabro volcano through a sliding window approach. The new U-GPD model identifies more useable S-wave arrivals than the base GPD model, yielding smaller subsequent location errors than 175 even our manual analyst's phase arrival picks. The new model also runs an order of magnitude 176 faster, processing 14 months of data from 7 broadband seismometers in less than 4 hours on a 177 single GPU. Our findings indicate that transfer learning can be extremely useful for volcano 178 seismic monitoring, even with limited computing resources and data. We conclude this paper with 179 a discussion of our findings, methodology and practical considerations of transfer learning in 180 Section 6. All data and code used throughout this paper are made fully and publicly available (see 181 Data Availability Statement). 182



Confidential manuscript submitted to JGR: Solid Earth

Figure 1. A) Model architecture for Generalized seismic Phase Detection (GPD) CNN model 186 (Ross et al., 2018b). Model can be considered as two parts: a feature extraction system 187 (convolutional layers) and classification part (fully connected layers). GPD model outputs 3 x 188 prediction values (probability of P, S or noise) for an entire 400-sample 3-component waveform 189 (i.e., output dimensions: 1 x 3). Examples of filter kernels (dashed line inset) from lowest 190 convolutional layer that extract generalized seismic waveform features determined through model 191 training on extensive SCSN dataset. These indicate that the GPD model has learnt to extract 192 193 different features from vertical and horizontal components. B) Proposed transfer learning model architecture ("U-GPD"). GPD model feature extraction system is copied to new model and fine-194 tuned with new Nabro data and low learning rate. Low learning rate ensures that useful features 195 are not 'unlearned'. New convolutional layers replace the GPD classification layers and are trained 196 197 using new Nabro data and higher learning rate. Model outputs 3 x prediction values for each datapoint in 400-sample 3-component waveform (i.e., output dimensions: 400 x 3). 198

199

200 2 Transfer Learning

201

There are many approaches to transfer learning (see Pan & Yang, 2010; Zhuang et al., 2020) 202 203 for comprehensive surveys), including using 'off-the-shelf' feature extraction systems from existing state-of-the-art CNNs (e.g., Maqsood et al., 2019; Razavian et al., 2014), learning domain-204 205 invariant or global representations across multiple tasks (e.g., Glorot et al., 2011; Li et al., 2014; Tzeng et al., 2015; Zhuang et al., 2015), applying pre-processing steps to make input data 206 representations more similar between datasets (e.g., Daumé, 2007; Sun et al., 2016) and the use of 207 domain-adversarial models (e.g., Ganin et al., 2016). Here we employ the first of these approaches 208 for P- and S-wave arrival time picking at Nabro volcano, utilizing pre-trained filters from an 209 existing, extensively trained CNN model (the GPD model; Ross et al., 2018b) to train a new model 210 with different output dimension and task type (see Section 3.1, U-GPD Model Architecture). Most 211 seismological studies that have employed transfer learning in this way have used pre-trained filters 212 from models designed for non-seismological tasks, such as image recognition. For example, filters 213 trained to recognize photographic images or handwritten characters have been used to detect 214 earthquakes and classify volcano-seismic event types from spectrograms (Huot et al., 2018; Lara 215 et al., 2020; Titos et al., 2020) and interpret seismic facies (Dramsch & Lüthje, 2018). 216

Some studies have chosen to fine-tune entire seismic deep learning models, essentially 218 updating the models with new data (or equivalently 'pre-training' the models with larger datasets, 219 depending on perspective). El Zini et al. (2020) pre-train an autoencoder with abundant unlabeled 220 data to learn compressed data representations of 2D seismic images. These model weights then 221 serve as a starting point for a model that segments seismic images, with weights fine-tuned using 222 limited labelled training data. This approach was shown to outperform the transfer of weights from 223 image recognition models and training a model from scratch. Bueno et al. (2020) fine-tune a 224 Bayesian neural network (BNN) to improve classification of volcano-seismic event characteristics 225 between datasets and time periods. They show that this approach increases model accuracy and 226 reduces epistemic uncertainty when applied to new volcanic systems or phases of activity. With 227 228 a similar aim but different approach to the work of this paper, Chai et al. (2020) utilize pre-trained weights from another existing phase arrival detection model, PhaseNet (Zhu & Beroza, 2019), to 229 pick phase arrivals from hydraulic fracturing experiments. They use the entirety of the PhaseNet 230 model and its pre-trained weights as a starting point for training and then fine-tune all model 231 232 weights equally using just 3,500 seismograms. They present improved results over the original PhaseNet model, which was trained using 700,000 seismograms of regional Californian seismicity, 233 when applied to higher sample rate data (100 kHz) from a very different setting (i.e., hydraulic 234 fracturing). Whilst these studies show that fine-tuning entire models can be an effective strategy, 235 poor hyperparameter choices (model learning rate, number of training epochs, etc.) can 236 inadvertently retrain the model (also known as 'catastrophic forgetting'; e.g., Kirkpatrick et al., 237 2017) or lead to settling on a non-global minimum within the parameter space, reopening the 238 potential for overfitting when the number of model parameters is large and the training dataset is 239 small (El Zini et al., 2020; Yosinski et al., 2014). The work in this paper differs from that of Chai 240 241 et al. (2020) in that only the weights from the feature extraction part (i.e., the first 'half') of the GPD model are transferred to our new U-GPD model. These weights are fine-tuned using a much 242 lower learning rate (weight update step size) to retain useful learned knowledge from the original 243 model but optimize cohesion with the rest of the new model, which is redesigned to reduce the 244 total number of trainable parameters, among other optimizations (see Section 3.1, Model 245 Architecture), and initialized with randomized weights (Fig 1). 246

248 **3 Proposed Model**

249 **3.1 U-GPD Model Architecture**

250

As outlined briefly above, we utilize pre-trained parameters from the convolutional layers 251 of the GPD model as a starting point for our U-GPD transfer learning model. The original GPD 252 model was trained using 4.5 million hand-labelled seismograms (1.5 million of each class P, S and 253 noise) recorded by the Southern California Seismic Network (SCSN) between the years 2000 and 254 2017. These training data were all 400-sample (4 sec) 3-component waveforms, high-pass filtered 255 above 2 Hz and (re)sampled at 100 Hz. All events had epicentral distances less than 100 km and 256 magnitudes between -0.81 and 5.7 M (various magnitude scales). The GPD model was chosen as 257 a base for our transfer learning model as these data characteristics (magnitude range, sample rate 258 259 and event distances) are comparable to those observed and recorded by volcano observatories. Furthermore, the short input length of 4 seconds (400 samples at 100 Hz sample frequency) means 260 there is less chance of erroneously labelling or missing relatively small magnitude or overlapping 261 phase arrivals. Finally, the GPD model's 'sequential' architecture, with each layer being solely 262 connected to the layers directly before and after, also means the model is more interpretable and 263 makes it easier to isolate its feature extraction system. 264

265

During model training, we fine-tune these pre-trained parameters using a very small 266 learning rate (1 \times 10⁻⁵), rather than keep them fixed (e.g., Yosinski et al., 2014). Learning rate 267 effectively controls how much model weights can change and a small learning rate will keep 268 adjustments to the pre-trained GPD feature extraction weights small. The aim of this fine-tuning 269 step is to modify any highly specific features from the source domain (particularly in the higher-270 level feature extraction layers) and overcome optimization difficulties arising from splitting the 271 272 GPD convolutional layers from co-adapted classification layers (Yosinski et al., 2014), without unlearning the important generalized waveform features we wish to exploit. 273

274

We then replace the GPD model's fully-connected layers (i.e., the task-specific classification part of the model) with further convolutional layers and up-sampling operations, combined with ReLU activation function (Nair & Hinton, 2010) and batch normalization (Ioffe &

Szegedy, 2015), to produce a model output with the same dimensions as model input (400 samples 278 x 3 channels; Fig 1B). Each of the three output channels represents the model's prediction (or 279 'probability') of a P-wave arrival, S-wave arrival or neither (hereby referred to as noise), 280 respectively, at each datapoint in the waveform. This all-convolutional approach has been adopted 281 by other phase arrival picking models (e.g., Woollam et al., 2019; Zhu & Beroza, 2019) and has 282 several distinct advantages when applied to seismic phase arrival detection: i) it provides less 283 ambiguous labelling of phase arrivals when compared to the original GPD model's approach of 284 assigning a single class prediction (P, S or noise) to an entire 400-sample 3-channel waveform; ii) 285 convolutional layers tend to have fewer parameters than fully connected neural network layers so 286 less training data is required to avoid overfitting; iii) by producing a model with input and output 287 traces of same dimension, we require less overlap when applied as a rolling window method, 288 289 producing a model that runs orders of magnitude faster on continuous sections of data.

290

291 The new convolutional layers are initialized with completely randomized weights and trained with a higher learning rate (1×10^{-3}) than the pre-trained GPD weights. A higher 292 learning rate effectively allows the randomized weights in the new model layers to be adjusted 293 much more than the pre-trained GPD weights. The learning rates used for each part of the model 294 were determined through experimentation, insight from previous works (e.g., Ross et al., 2018a, 295 2018b), and on the basis that the learning rate for fine-tuning pre-trained weights should be orders 296 of magnitude lower than that used for tuning randomized weights (e.g., Yosinki et al., 2014). We 297 298 note that there are more formal strategies (e.g., grid/random search, Bayesian optimization, bandit strategies, gradient reversal; Bergstra & Bengio, 2012; Feurer et al., 2015; Klein et al., 2016; 299 Maclaurin et al., 2015; Snoek et al., 2015) for determining optimal model hyperparameters. Such 300 strategies, however, add significant computational cost as they generally require repeatedly 301 training models with differing hyperparameter choices, producing a much greater search space. 302 The aim in this study is not to present the absolute best possible model architecture and set of 303 304 hyperparameters specific to this deployment at Nabro (as these choices will likely be specific to application and training set size) but to illustrate how existing models can be tailored to new 305 datasets to improve performance in those settings. Furthermore, it would prove more difficult to 306 attribute any observed improvements to the use of transfer learning and U-Net architecture, as 307 308 opposed to the hyperparameter optimization strategy. We do, however, implement two further

hyperparameter choices that were found to improve performance. First, we use dilated filter 309 kernels in the new convolutional layers (e.g., van den Oord et al., 2016; Yu & Koltun, 2016) to 310 increase the size of the model's receptive field (or 'field of view') and aggregate multi-scale 311 context. Second, the new layers are subjected to spatial dropout (Tompson et al., 2015), where 312 30% of the feature maps (output of filter operations) in each convolutional layer are effectively 313 dropped (set to zero) at the start of each training epoch. This step promotes independence between 314 the features the model extracts and prevents overfitting (Tompson et al., 2015). Precise details of 315 316 U-GPD model dimensions and hyperparameters are provided in Supplementary Materials (Fig S1). 317

318

The overall network architecture outlined above is sometimes referred to as a U-Net 319 320 (Ronneberger et al., 2015). With each step through the network, the input data are progressively downsampled with an increasing number of features extracted, creating a contracting network path 321 322 that is forced to sacrifice detail and learn a more compressed, general representation of the input waveform to discriminate between classes (P, S or noise). The model then follows a symmetrically 323 324 expanding path, where the data are progressively upsampled and the number of features reduced, to regain precise temporal or spatial detail and return an output with equal dimension to the model 325 input (Ronneberger et al., 2015). Skip connections (addition operators), which act as direct, one-326 way pathways between layers in the contracting and expanding sides of the model (Fig 1B), are 327 328 used to retain precise waveform details that may be lost through this contraction/expansion process 329 and have been shown to greatly improve the likelihood of model parameters settling on the global minimum during training (Li et al., 2017). 330

331

332 **3.2 Phase Arrival Labels and Model Hyperparameters**

333

Each 3-component waveform in our training dataset has a corresponding 3-channel 'mask' that provides a ground truth label (P, S or *noise*) for each waveform datapoint. During training, the model aims to minimize the difference between its predictions and these ground truth labels. Labels are presented as binary values (0's or 1's), with P-wave arrivals indicated by a +/- 0.14 sec boxcar function, centered on the manually picked P-wave arrival time, and S-wave arrivals indicated by +/- 0.19 sec boxcar function, also centered on the manually picked S-wave arrival time. These boxcar widths provide a good balance between phase arrival detection rate and arrival
time precision and compensate for human error in the ground truth labels. Previous studies have
used Gaussian-style probability masks, with values ranging between 0 and 1, for labelling phase
arrivals (e.g., Woollam et al., 2019; Zhu & Beroza, 2019). We find that label accuracy on our test
data (e.g., Figs 5, 6 and 7) and event location error distributions from the full deployment (e.g.,
Fig 10C & D) are near-identical when using either approach but training with boxcar masks
produces a model that detects ~ 10% more events when run over continuous data.

347

As with the original GPD model, our new U-GPD model was trained using a categorical cross entropy loss function (Text S7) and the Adam optimization algorithm (Kingma & Ba, 2014). The model weights that produced lowest loss value on the validation dataset during training were selected as our final model weights. Other loss functions that address the imbalance between arrival and noise labels (as the majority of labels in any given waveform are not a phase arrival), such as a focal loss function that effectively adds weighting parameters to cross entropy loss (Lin et al., 2017), were trialed but yielded no improvement in model performance.

355

356 **4 Data**

357

Nabro volcano is one of two calderas that form the Bidu Volcanic Massif on the Eritrea-358 Ethiopia international border (Fig 2). Located in the Afar region at the northern end of the Main 359 Ethiopian Rift, it erupted unexpectedly for the first time in recorded history on 12th June, 2011, 360 disrupting continental aviation and initiating a significant humanitarian crisis (Bojanowski, 2011; 361 Donovan et al., 2018; Goitom et al., 2015). At the time, there were no seismic or other monitoring 362 networks operating in Eritrea but earthquakes were felt around the volcano several hours and days 363 prior to eruption, prompting evacuation (Goitom et al., 2015). This seismicity is the first of note 364 in global catalogues for the region (Goitom et al., 2015). Despite this fortuitous warning, at least 365 seven people were tragically killed and about 12,000 were displaced (Bojanowski, 2011; Goitom 366 et al., 2015; Hamlyn et al., 2014). The eruption is particularly notable for the vast amount of SO_2 367 emitted into the atmosphere, one of the largest eruptive SO₂ masses globally since the eruption of 368 Mount Pinatubo in 1991 (Fromm et al., 2014; Goitom et al., 2015; Theys et al., 2013), and the 369

comparative rarity of recorded historical eruptions in the region (Goitom et al., 2015; Hamlyn etal., 2014).

372

In August, 2011, approximately two months after the eruption began, eight 3-component 373 broadband seismometers (5 x Guralp CMG-6T, 3 x Guralp CMG-40T; Fig 2) were deployed 374 around the volcano to monitor ongoing activity (Hamlyn et al., 2014; Hammond et al., 2011). 375 These stations remained operational for 14 months until October, 2012. The first two months of 376 data were collected at a sample rate of 100 Hz before dataloggers were switched to a sample rate 377 of 50 Hz for the remainder of the deployment to maximize data recovery while minimizing service 378 runs. Data from the full deployment occupies 70 GB of disk space (miniSEED format). Manual 379 phase arrival picking conducted on the first four months of data (2011-08-30 to 2011-12-31; 380 Goitom, 2017; Hamlyn et al., 2014) identified a total of 2926 events, from which the first 35 days 381 of data (all 100 Hz sample rate) were quality checked and used for training and validating our 382 transfer learning model. Five subsequent days of data (2 x 100 Hz days, 3 x 50 Hz days) were 383 selected and quality checked to serve as test data. The reason to exclude 50 Hz data from model 384 385 training is to emulate data availability in the early stages of this seismic deployment and demonstrate that changes in sample rate can be overcome without compiling new training datasets 386 387 through a process known as data augmentation. The raw data for all datasets (training, validation and testing) were self-normalized, with linear trend removed, and left unfiltered. 388





Figure 2. Regional topographic map (90 m CGIAR Shuttle Radar Topography Mission and 392 GEBCO bathymetry model, grey-scale map center) and seismic deployment (30 m ALOS Digital 393 Surface Model, color map bottom right) around Nabro volcano. Red triangles (center map) indicate 394 395 Holocene volcanoes (Global Volcanism Program, 2013) with Nabro volcano highlighted in white. Inverted blue triangles (bottom right map) indicate operational broadband seismic stations 396 deployed around Nabro volcano from August 2011 to October 2012 (station NAB6, inverted 397 yellow triangle, was flooded shortly after deployment and not operational). Training and validation 398 data were taken from dark blue stations only (NAB1, NAB2, NAB3, NAB4 and NAB8). 399

A total of 2921 waveforms with labelled P- and S-wave arrivals from 978 events (2011-401 08-30 to 2011-10-03) and five stations were used as training and validation data (only five stations 402 were consistently operational during this time; dark blue stations in Fig 2 bottom right map). 403 Training and validation data were grouped and divided so that no event appeared in both datasets 404 to avoid data leakage (the model being trained on event data that also appears in validation or 405 testing). 857 events (2498 waveforms) were used for model training and 121 events (423 406 waveforms) were used for model validation, a training-validation split of approximately 85%-15%. 407 624 sections of noise (20 secs length) were manually identified across all five stations (2011-08-408 31 to 2011-09-27), with 500 sections (2500 waveforms) and 85 sections (425 waveforms) used for 409 model training and validation, respectively. Two noise waveforms were randomly dropped from 410 411 each dataset so that the training and validation noise data comprise 2498 and 423 waveforms, respectively, to match the number of event waveforms. 412

413

A separate test dataset of 400 event waveforms with labelled P- and S-wave arrivals (132 events) and 400 noise waveforms (80 sections of noise) was also produced for subsequent model testing. These data come from a different time period than those used for training and validation data, with 200 waveforms from a period where data were recorded at 100 Hz sample rate (2011-10-04 and 2011-10-05) and 200 waveforms from a period with 50 Hz sample rate (2011-10-14, 2011-10-15 and 2011-11-27) for each category. All training, validation and test data were manually identified and quality checked.

421

The success of U-Net architectures relies on an effective data augmentation strategy when 422 working with smaller datasets (Ronneberger et al., 2015). This allows the network to learn 423 424 invariance to certain changes in input signal without them needing to appear in the annotated dataset. Here we outline a data augmentation strategy that improves performance of our U-GPD 425 transfer learning model (Fig S2). First, as all stations were switched from 100 Hz sample frequency 426 to 50 Hz sample frequency part way through the seismic deployment, we randomly select subsets 427 of the training data (all originally sampled at 100 Hz) to be decimated to 50 Hz sample frequency 428 throughout training. Each training sample (i.e., each 3-component waveform) has a probability of 429 0.5 of being selected for decimation before each training epoch, with an anti-aliasing, low-pass 430

finite impulse response (FIR) filter applied and linear phase shift removed. Second, we randomly time-shift our P- and S-wave arrivals relative to the model input 'window', so that our waveforms differ slightly from epoch to epoch and the model must learn signal features that indicate arrivals rather than where they occur within the input window (i.e., arrivals don't need to occur in the center of the window for the model to detect them). With our noise data, a random 400-sample window is chosen at each training epoch from our 20-second noise sections, introducing more waveform variety between training epochs.

438

All data processing and model training/testing were performed in Python using the ObsPy
(Beyreuther et al., 2010; Krischer et al., 2015; Megies et al., 2011), TensorFlow (Abadi et al.,
2015; https://tensorflow.org) and Keras (Chollet et al., 2015; https://keras.io) libraries.

442

443 **5 Results**

444 **5.1 Training Metrics (Transfer Learning vs No Transfer Learning)**

445

To examine the impact of transfer learning and determine how much training data is 446 required to produce an effective model, we use varying sized subsets of the training data 447 throughout model training (i.e., 250, 500, 750, ..., 2000, 2250 and 2498 training samples). Figure 448 3 compares how model loss (measure of distance between model predictions and ground truth 449 labels) on training and validation data evolves throughout training between our transfer learning 450 model and the same model with completely re-initialized weights (i.e., with no transfer learning) 451 for our smallest and largest subsets of training data (250 and 2498 training samples, respectively). 452 The learning rate is set to be equal (1×10^{-3}) across the whole re-initialized model as we are no 453 longer fine-tuning existing knowledge. All other hyperparameters, including dropout rate, are kept 454 455 the same. The models trained without transfer learning (Fig 3B and D) show a much greater degree of overfitting: the model loss on the training data continues to decrease with more training while 456 the loss on validation data (data that the model does not use during training) hits an inflection point 457 and starts increasing, reflecting that the model is 'memorizing' the precise features of the training 458 459 data at the cost of generalization (Shorten & Khoshgoftaar, 2019). By contrast, the validation loss continues to decrease for the models trained with transfer learning (Fig 3A and C). Furthermore, 460

the minimum validation loss achieved by the transfer learning models for each training dataset size 461 is lower than when transfer learning is not employed (Fig 3 horizontal dashed lines). Such 462 diagnostics indicate that transfer learning is successfully preventing overfitting to the training data 463 and will likely produce a model that generalizes better to non-training data (Shorten & 464 Khoshgoftaar, 2019). The greatly improved performance on validation data using the smallest 465 subset of training data (Fig 3A and B) shows that transfer learning is particularly useful for 466 reducing overfitting and model loss when training data are very limited, but this advantage is 467 468 progressively diminished with increasing training dataset size (Figs 3 and 4).







Figure 3. Model loss vs. training epoch number. A) Transfer learning model and 250 training 472 samples of each class (P, S or neither). B) Model trained without transfer learning (i.e., initially 473

Confidential manuscript submitted to JGR: Solid Earth

randomized weights) and 250 training samples of each class. C) Transfer learning model and full 474 training dataset (2498 training samples of each class). D) Model trained without transfer learning 475 (i.e., initially randomized weights) and full training dataset. Blue curve shows model loss for 476 training data, red curve shows model loss for validation data (not seen during training). A lower 477 model loss on training data (blue) than validation data (red) means the model shows signs of 478 overfitting. The degree of overfitting (gap between blue and red curves) is much greater for the 479 models without transfer learning (**B** and **D**) with validation loss hitting an inflection point then 480 increasing whilst training loss continues to decrease. The transfer learning models also achieve a 481 smaller minimum validation loss (horizontal dashed line) for each training set size. 482

483

Figure 4 shows the highest model accuracy (the proportion of labels the model classifies 484 485 correctly) and lowest model loss achieved by our transfer learning and re-initialized models on validation data when trained using each subset size of training data. The transfer learning model 486 487 achieves lower model loss regardless of training dataset size (Fig 4B). As training dataset size increases, the difference between the lowest loss achieved by the two models (gap between red 488 489 circles and red triangles, Fig 4B) decreases and the advantages of transfer learning diminish. Generally, loss is considered a more robust metric than accuracy for model performance on future 490 491 data as it measures the distance between model predictions and ground truth labels, whereas accuracy simply measures a binary true/false score. However, accuracy still provides useful 492 493 information regarding model performance. In particular, the transfer learning model shows a stable relationship between maximizing model accuracy and minimizing model loss (gap between black 494 and red circles is very small for all training subset sizes), where the training strategy of minimizing 495 model loss appears to achieve the same goal as maximizing model accuracy, again a sign of 496 reduced overfitting. The re-initialized model (black and red triangles), on the other hand, shows a 497 498 much less stable relationship in this regard, with diverging training scores (Fig 4) indicating that high model accuracy comes at the cost of higher model loss and low model loss comes at the cost 499 of lower model accuracy for these small training set sizes when transfer learning is not employed. 500 The increased model loss for model weights with highest model accuracy (black triangles) also 501 502 suggests that the model has become overconfident in its predictions (it has large errors on the small proportion of labels it gets wrong) and is therefore likely to perform worse on out-of-distribution 503

data, with more false or missed phase arrival detections (e.g., a phase arrival being labelled as 504 noise with high model confidence, or vice versa). 505

506

Model performance between the two approaches (transfer learning vs re-initialization) 507 converges as training set size increases, indicating that the need for transfer learning decreases 508 with increased training set size, as expected. In fact, model performance with transfer learning 509 appears to plateau, or possibly even degrade, at training subset sizes of more than 1500 samples. 510 511 This suggests that, with enough training data, transfer learning could potentially inhibit the model's ability to learn useful features in the new data that are absent in the original GPD training data. 512 This apparent variance in performance may also simply be a result of the stochasticity arising from 513 training using randomized weights in the new part of our transfer learning model. 514





517

Figure 4. Model accuracy (A) and loss (B) for various subsets of training data. Open red circles 518 are transfer learning model weights from epoch that achieves lowest validation loss (e.g., dashed 519 horizontal lines in Fig 3), open black circles are transfer learning model weights from epoch that 520 achieves highest validation accuracy, solid red triangles are re-initialized model (no transfer 521 learning) weights from epoch that achieves lowest validation loss, and solid black triangles are re-522 initialized model weights from epoch that achieves highest validation accuracy. 523

525 5.2 Test Dataset (Known Arrival Times)

526

Following model training, we test the above models (i.e., new model with and without 527 transfer learning) and two baseline models (GPD and PhaseNet) using the test dataset outlined in 528 Section 4. We examine the proportion of correct class predictions (Fig 5) and the residuals between 529 model and manually determined phase arrival pick times (Fig 6). Due to differences in model task 530 types (classification vs segmentation), we apply all models as sliding windows over 1000-sample 531 waveforms (note that the PhaseNet model takes a 3000-sample waveform as input so we examine 532 only the middle 1000 samples for this model). To account for human picking error in collating our 533 test set, we define a true positive for each phase arrival type (P or S) as the model prediction 534 exceeding a given threshold value for that arrival type within 0.5 secs of the manually determined 535 arrival, such that predicted arrival times very close to the manually determined arrival time are 536 considered accurate. A true positive for sections of noise is defined as no phase arrival prediction 537 exceeding a given threshold value at any point within that section of data. The test data are pre-538 processed as per the training data for each model (i.e., GPD model tested on 2 Hz high-pass filtered 539 data and all other models, including PhaseNet, tested on raw data; all detrended and self-540 541 normalized).

542

The GPD model is tested using four different threshold values (Fig 5A - D) as this value 543 strongly controls the number of false or missed phase arrival detections generated by this model. 544 545 When the threshold is set to be whichever class label (P, S or N) has the highest predicted value 546 for a given waveform, nearly all P- and S-wave arrivals are detected by the GPD model (99.75 % and 95 % detection rate, respectively; Fig 5A). However, this threshold criterion makes the GPD 547 model extremely prone to false phase arrival detections in sections of noise, with 44 % of 1000-548 sample noise waveforms in our test dataset containing at least one false phase arrival detection 549 (Fig 5A, bottom right square) and many of our 1000-sample event waveforms containing multiple 550 phase arrival triggers (e.g., Fig 7B & E). When this threshold criterion is applied to continuous 551 sections of data from Nabro, the number of false phase arrival detections overwhelmingly 552 outweighs the number of true phase arrival detections and becomes unmanageable in terms of 553

correctly associating phases, identifying true events and processing the data within computational
 memory constraints.

556

One way to lower the number of false phase arrival detections is to use a higher threshold 557 value for P- and S-wave predictions. Figure 5B shows the GPD model's performance on our test 558 data using a 0.9 threshold value (i.e., a P or S prediction 'probability' must exceed 0.9 to be 559 included). The number of false detections in sections of noise is greatly reduced (down from 44 % 560 of waveforms to 10 % of waveforms) but at the cost of reduced true phase arrival detections (~ 561 95% and ~82% of P- and S-wave arrivals, respectively). Part of this performance dip is 562 undoubtedly due to the difference in sample rates between one half of the test data (50 Hz) and the 563 GPD model's training data (all 100 Hz). When the threshold value is increased further (i.e., P or S 564 prediction must exceed 0.95 or 0.99; Fig 5C and D), the GPD model yields even fewer false phase 565 arrival detections in noise sections but at the cost of fewer P- and S-wave arrivals. 566

567

Figure 5E shows the performance of the PhaseNet model on our test dataset. This model is 568 included as it adopts the same U-Net segmentation approach as our new model and is trained on 569 data from a variety of instrument types, although the training data is still exclusively from 570 California. The PhaseNet model is much less prone to false phase arrival detections than the GPD 571 model (Fig 5E, bottom right square); as such, a much lower threshold value (0.4) can be used to 572 maximize the number of true phase arrival detections. This model accurately identifies ~ 89% and 573 ~ 83 % of P- and S-wave arrivals in our test dataset, which is better than the GPD model with a 574 threshold value that achieves a similar false detection rate (e.g., Fig 5D), but detects fewer phase 575 arrivals than our transfer learning and reinitialized models trained with Nabro data (Fig 5F – I). 576







Figure 5. Confusion matrices for base GPD model (A - D), PhaseNet model (E), U-GPD transfer learning model (**F**, 500 training samples, and **G**, 2498 training samples) and re-initialized model (**H**, 500 training samples, and **I**, 2498 training samples). Values in matrices are proportion of ground truth phase arrivals (test set) assigned by each model to a given class (values of 1 along diagonal from top left to bottom right means all phase arrivals and sections of noise correctly identified).

586

When trained using a subset of just 500 training samples for each class (P/S/N) and 587 evaluated using a prediction threshold value of 0.4, the transfer learning approach correctly detects 588 \sim 93% and \sim 94% of P- and S-wave arrivals with very few false phase arrival detections in sections 589 of noise (~1 %; Fig 5F), a clear improvement over our model trained with re-initialized weights 590 and the same training subset (Fig 5H). When our full training dataset is used (2498 samples for 591 each class), model performance converges between transfer learning (Fig 5G) and re-initialization 592 (Fig 5I), with a similar number of correctly identified phase arrivals and false detections in noise, 593 although the transfer learning model still performs marginally better, particularly on sections of 594 noise. In essence, the transfer learning model strikes a better balance between high phase arrival 595 detection rate (~97 - 98% for each phase arrival type; Fig 5G, top left and center squares) and low 596

false detection rates in sections of noise (~ 1%; Fig 5G, bottom right square) on our test data from Nabro volcano than any of the existing baseline models (Fig 5A – E) or training a model from scratch (Fig 5I).

600

Figure 6 shows the residuals for each model between their predicted phase arrival times 601 and the original manual pick times for these test waveforms. Predicted phase arrival times were 602 determined using a simple trigger algorithm (e.g., Withers et al., 1998) on each model's probability 603 time series with the time series index that yields maximum predicted value chosen as the pick time 604 for a given phase arrival type (Fig 7). The models that employ semantic segmentation (i.e., 605 PhaseNet, our U-GPD transfer learning model and our re-initialized model; Fig 6B - F) show 606 comparable pick time precision (root mean square deviation [RMSD] of 0.036, 0.038 and 0.044 607 seconds, respectively, for each model's P-wave predictions and RMSD of 0.053, 0.053 and 0.065 608 seconds, respectively, for each model's S-wave predictions). The GPD model (Fig 6A), by 609 comparison, has a more diffuse range of phase arrival pick times (RMSD of 0.217 seconds for P-610 waves and 0.188 seconds for S-waves), with some model picks made more than 1 second before 611 612 or after the manually determined arrival time. This is almost certainly a result of its more ambiguous class labelling (Fig 1) and the broad phase arrival probability peaks it generates (Fig 613 614 7).



617

Figure 6. Model phase pick residuals vs. manual phase picks for base GPD model (A), PhaseNet 618 model (B), U-GPD transfer learning model (C, 500 training samples, and D, 2498 training 619

samples), and reinitialized model (\mathbf{E} , 500 training samples, and \mathbf{F} , 2498 training samples). The models based on semantic segmentation ($\mathbf{B} - \mathbf{F}$) yield smaller phase pick residuals.

622

Figure 7 shows three example waveforms from the test set with corresponding model 623 predictions for the U-GPD transfer learning, GPD and PhaseNet models. These waveforms were 624 chosen as they have low SNR phase arrivals. Prediction labels for the U-GPD model resemble the 625 boxcar labels of the training set (Fig 7A, D & G), whereas prediction labels produced by the 626 PhaseNet model resemble the model's truncated Gaussian-style training labels (Fig 7C, F & I). 627 Despite these boxcar shapes, the U-GPD model's maximum predicted value for each phase arrival 628 consistently and accurately picks both P- and S-wave arrivals (Fig 7A, D & G). On the other hand, 629 the base GPD model's prediction labels are considerably broader and noisier (Fig 7B, E & H). The 630 U-GPD model appears to have benefitted from retraining using Nabro-specific data, as it performs 631 much better than the existing models on these challenging waveforms. 632



Figure 7. Three example waveforms from our test set. Phase arrival prediction trigger thresholds 635 (horizontal dashed lines) are 0.4, 0.9 and 0.4 for U-GPD (left), GPD (center) and PhaseNet (right), 636 respectively. (A - C), Test waveform with substantial high frequency background noise. All 637 models accurately detect P- and S-wave arrivals but GPD model makes multiple phase detections. 638 $(\mathbf{D} - \mathbf{F})$, Test waveform with low amplitude P-wave arrival. Existing GPD and PhaseNet models 639 incorrectly label S-wave arrival as close combination of P-wave arrival and S-wave arrival. U-640 GPD transfer learning model correctly detects both P- and S-wave arrivals. (G - I), Test waveform 641 with substantial low frequency background noise. P-wave arrival prediction is below trigger 642

643 thresholds for both GPD and PhaseNet models, although GPD model accurately detects S-wave

arrival. U-GPD transfer learning model correctly identifies both P- and S-wave arrivals.

645

646 **5.3 Full 14-Month Deployment (Unknown Arrival Times)**

647

Whilst evaluating model performance on individual, manually scrutinized waveforms is useful for benchmarking and yielding estimates of model efficacy, the model's performance in a 'real-world' setting is ultimately of most importance to seismic analysts. Evaluating such performance is inherently more challenging, however, as the number of events in long sections of monitoring data and their respective phase arrival times are unknown, and other considerations, such as computational time and resources (e.g., memory requirements and availability of optimized hardware), affect model feasibility as a monitoring tool.

655

In this section, we present results of our best performing model in the prior section (U-656 657 GPD transfer learning model trained with full training dataset of 2498 samples of each class) and the original base GPD model when run over the full 14-month Nabro seismic deployment (Fig 8). 658 As with the test dataset in Section 5.2, phase arrivals are detected at individual stations through a 659 simple trigger algorithm, where an arrival is detected if the probability assigned to that class label 660 661 (P or S) exceeds a given threshold (e.g., 0.4 for our U-GPD transfer learning model). The phase arrival time is determined as the waveform sample with the highest probability for that phase (Fig 662 663 7).

664

665 The U-GPD transfer learning model was applied to the data as a sliding window with 50 % overlap (i.e., applied at 'time shifts' of 200 samples) over 24-hour sections of data from each 666 individual station. The model takes 5 seconds to process 24 hours of 3-component data at 100 Hz 667 sample rate (or 3 seconds per day at 50 Hz sample rate) on a single graphics processing unit (GPU; 668 669 NVIDIA GeForce RTX 2080 Ti), a rate many orders of magnitude faster than 'real-time' even when run on hundreds of stations. To avoid poor predictions due to window edge effects, only the 670 middle 200 sample predictions out of 400 from each window are used to predict phase arrivals and 671 are concatenated to produce one long continuous prediction trace without overlap or gaps and with 672 the same sample rate as that of the input signal (i.e., 100 or 50 Hz). With all other processing steps 673

674 (e.g., software initialization, data read/write, signal windowing, running trigger algorithm, etc.),

the U-GPD transfer learning model picks phase arrivals at all 7 available stations from the full 14-

- 676 month deployment in less than 4 hours using a single GPU (greatly reduced when parallelized over
- 677 multiple GPUs), indicating that it could easily be used within real-time monitoring constraints.
- 678

Conversely, as the GPD model produces only one class prediction per window (Fig 1A), 679 we apply this model with much greater overlap (97.5 %; every 10 samples of data) and with 680 varying threshold values (0.9, 0.95 and 0.99) for phase arrival detection triggering. This generates 681 a prediction trace with a much coarser sample rate than the original input signal (i.e., from 100 or 682 50 Hz to 10 or 5 Hz, respectively) and takes 26 seconds per 24 hours' 3-component data at 100 Hz 683 sample rate (or 15 seconds per day at 50 Hz sample rate) on the same NVIDIA GPU, approximately 684 685 a five-fold increase in computational time with a tenth of the temporal detail. With all other processing steps, the GPD model took almost 50 hours to run over the full 14-month deployment 686 687 using a single GPU, more than a ten-fold increase in computational time over the transfer learning model, due to more (pre-)processing required (e.g., more signal windows generated and 688 689 subsequent processing). Assuming a linear increase in computational time, running the model as a sliding window over every sample of data would take ~ 260 seconds per 24 hours' 3-component 690 691 data at 100 Hz sample rate and ~ 500 hours (nearly 3 weeks) for the full 14-month deployment and 7 stations. While this is still faster than real-time, these timescales for a single or limited number 692 693 of station(s) could become limiting when applied at hundreds of stations, particularly without high 694 performance computing resources.

695

696 **5.3.1 Phase Association Method**

697

Both models detect P- and S-wave phase arrivals but do not associate them to the same event. To assess the number of locatable events detected, we group P-wave phase arrival triggers into 4-second bins and keep only bins with arrivals detected at four or more stations. This bin size was chosen to encompass the maximum plausible travel time between any two stations. If multiple arrivals were detected at the same station within a 4-second bin, the detection threshold was increased for all arrivals in that particular bin to retain only the highest probability phase picks. If

any of these bins now had arrivals at less than four stations, as a result of removing lower 704 probability phase picks, they were discarded as there would be too few stations to constrain event 705 location. If there were still multiple arrivals present at any given station, only the arrivals with 706 highest probability for each station were kept. Finally, if phase arrival bins intersected (a subset of 707 one bin was contained in another), the bin with highest mean probability was kept. This association 708 method is clearly quite crude, and only works for small, very local arrays, but allows a broad 709 evaluation of model performance at detecting phase arrivals. Use of a more rigorous phase 710 association method (e.g., Ross et al., 2019; Yeck et al., 2019) would obviously be better at 711 eliminating false arrival picks or identifying multiple events within a 4 second window, which is 712 a common feature of seismicity during volcanic unrest. However, this will mask underlying model 713 performance; e.g., the inclusion of false arrival picks is likely to generate greater estimated location 714 errors (Fig 10). 715

716

We associate S-wave arrivals to their corresponding P-wave arrivals by first locating events 717 using NonLinLoc (e.g., Lomax et al., 2000), a widely used software package for probabilistic 718 719 earthquake location, using the P-wave arrival bins outlined above (Fig 8A) and a simple 1D linear gradient velocity model from previous seismic studies at Nabro (Table S8; Goitom et al., 2015; 720 721 Hamlyn et al., 2014). The difference between P-wave arrival and event origin times were used to predict which S-wave arrival detections should be associated with each P-wave arrival using a 722 723 Vp/Vs ratio of 1.76 (Goitom et al., 2015) and S-wave travel time error of 0.25 (25%). S-wave arrival triggers that lay within this error bound for each detected P-wave arrival were associated to 724 that event. S-wave arrivals at stations without a detected P-wave arrival were not included. All 725 events were then located again in NonLinLoc using all included phase arrivals (Fig 8B). 726





Figure 8. U-GPD transfer learning model event locations (total no. of events = 33,950) using
automated phase association strategy. A) P-wave phase arrival triggers are grouped into 4 second
bins and these groupings are used to obtain initial event hypocenters and origin times. B) S-wave
phase arrival triggers are associated to P-waves in (A) using initial origin times, a Vp/Vs ratio of
1.76 and a travel-time error of 25 %. Events are then located again using all included P-wave and
S-wave arrivals.

737 5.3.2 Detected Events and Location Errors

738

Figure 9 shows the cumulative number of events detected by the U-GPD transfer learning 739 740 model (threshold value of 0.4; black solid line) and the original GPD model (threshold values of 0.9, 0.95 and 0.99; grey lines). The cumulative number of events from an existing manual 741 catalogue for this deployment (Goitom, 2017; Hamlyn et al., 2014), some of which provided the 742 transfer learning model training data, is also given for reference. Event locations for each model 743 744 and the manual catalogue are provided in Supplementary Materials (Figs S3 - S6). When only Pwave arrivals are used (Fig 9A), the GPD model with detection threshold of 0.9 appears to detect 745 the most events (total no. of events detected by GPD model = 41,007; total no. of events detected 746 by transfer learning model = 33.950). A threshold of 0.95 also detects more events than the transfer 747 learning model until shortly after the switch in instrument sample rates from 100 Hz to 50 Hz. 748

However, when we consider events with at least one associated S-wave arrival, the transfer 749 learning model detects more events overall (Fig 9B; no. of events detected by transfer learning 750 model = 31,387; no. of events detected by GPD model with 0.9 threshold = 26,808). This is 751 consistent with the results from our test dataset in Section 5.2, with the proportion of S-wave 752 arrivals accurately detected by the GPD model at these threshold values much lower than the 753 proportion of P-wave arrivals detected (Fig 5B – D). Furthermore, 6 % of noise waveforms and 754 16% of S-wave arrivals from our test data were mislabeled by the GPD model (0.9 threshold value) 755 as P-wave arrivals (Fig 5B), a higher rate of false detections or labels than the transfer learning 756 model (1 % of noise sections and 0.5% of S-waves, respectively; Fig 5G). This means that a higher 757 proportion of the P-wave groupings detected by this model with 0.9 threshold value are likely to 758 include mislabeled S-waves or false arrivals, which is reflected in subsequent event location errors 759 760 (Fig 10C – D).



Figure 9. Cumulative number of events detected by GPD model (various thresholds, grey lines)
and transfer learning model trained on full Nabro dataset (2498 samples of each class, 0.4
threshold, black line). Blue dashed line is existing manual catalogue (Goitom, 2017). All training
/ validation waveforms are from dates before switch in sample frequency (vertical dashed line).

Confidential manuscript submitted to JGR: Solid Earth

A) Cumulative number of events detected using P-wave arrivals only (see main text for event
binning procedure). B) Cumulative number of events with at least one associated S-wave arrival.

- To scrutinize these results further, we examine the number of stations with P- and S-wave 771 arrival detections per event (Fig 10A – B). In general, the events detected and picked by the U-772 GPD transfer learning model include more stations and considerably more S-wave arrivals than 773 those picked by the GPD model, although the number detected by the GPD model may have been 774 reduced by using a coarser prediction trace (every 10 samples, a requirement to reduce model run 775 time to a reasonable timeframe). This increase in the number of stations and S-wave arrivals per 776 event will constrain event locations, as seen in the location errors derived from the models' phase 777 arrival picks (Fig 10C - D). 778
- 779

Location errors are estimated by NonLinLoc using multi-dimensional Gaussian estimators 780 and subsequent confidence intervals (e.g., Lomax et al., 2000). The horizontal errors (Fig 10C) for 781 the locations produced using the transfer learning model pick times are comparable to the existing 782 783 manually picked events. Furthermore, vertical (depth) errors are much improved over the manual catalogue (Fig 10D), likely reflecting more consistency in S-wave arrival picking than that of a 784 785 manual analyst. The GPD model, by comparison, produces a more diffuse range of horizontal and vertical errors, which is likely to be a combination of coarser prediction trace, poorer pick precision 786 787 (Fig 6A), lack of S-wave arrivals (Fig 10B) and false/mislabeled P-wave arrival detections (Fig 5B). This interpretation is further supported when we look at the number of event locations lying 788 within the array (i.e., event locations lying within the convex hull of station coordinates) for each 789 model: NonLinLoc locates more events within the array using the transfer learning picks (n = 790 791 23,859) than using the GPD model with 0.9 threshold value (n = 22,826). While we expect many events to occur outside of the array (e.g., at neighboring faults or volcanic centres), this metric 792 shows that a much larger proportion of event locations detected by the GPD model lie away from 793 the volcanic edifice, which may reflect poorer pick precision, false/mislabeled arrivals or coarser 794 prediction trace, but may also reflect the event types (i.e., regional tectonic) that the original model 795 796 was trained on.



Figure 10. A) Number of P-wave arrival picks per event for transfer learning model (grey) and 800 base GPD model (gold). B) Number of S-wave arrival picks per event. C) Histogram of Gaussian 801 horizontal location errors (1 standard deviation) for events picked by transfer learning model (grey) 802 and base GPD model (gold), and those in the existing manual catalogue (blue). D) Histogram of 803 Gaussian vertical (depth) location errors (1 standard deviation). 804

805

6 Discussion 806

807

Transfer learning using existing seismological deep learning models can be a highly 808 effective strategy to automate phase arrival picking in settings with little or no prior monitoring. 809 We demonstrate that, with a limited number of hand-labelled waveforms (on the order of hundreds 810 to low thousands) and a few minutes of training time, one can produce a consistent and effective 811

deep learning model for phase arrival detection that requires no other manual intervention or tuning
and can process years of data in a matter of hours.

814

For small training datasets, the use of pre-existing, generalized CNN filters greatly reduces 815 model overfitting (i.e., model parameters 'memorizing' the training data) when compared with 816 training a model from scratch (Fig 3) and yields a more stable relationship between maximizing 817 model accuracy and minimizing model error (Fig 4). Furthermore, when combined with a good 818 data augmentation strategy, transfer learning can also address the issue of processing data when 819 instrument sample rates differ from those used to train existing models. When applied to data from 820 Nabro volcano, augmenting our training set with decimated waveforms greatly improves model 821 performance on lower sample rate data (Fig S2). As such, hand-labelled training data from the first 822 823 35 days of the deployment (all 100 Hz sample rate) were sufficient to detect phase arrivals throughout the duration of the deployment, even after instrument sample rates were switched to 824 50 Hz (Fig 9). Without this data augmentation step, model performance on lower sample rate data 825 declines dramatically (Fig S2). This shows that where sample rates are altered or new instruments 826 827 added during a seismic deployment, data augmentation can overcome the cost of collecting further hand-labelled data and allow models to be adapted cheaply and quickly throughout the 828 829 deployment.

830

831 The introduction of new, task-specific data and the change in model task from one of classification to one of segmentation also improves our U-GPD model pick time precision (Fig 6), 832 the number of stations per detected event (Fig 10A), the number of S-wave arrivals detected (Figs 833 5 and 10B) and computational efficiency over the original base GPD model, as well as potentially 834 reducing the number of false/mislabeled P-wave detections (Fig 5) and increasing the number of 835 836 identified events that relate directly to volcanic activity (evidenced by the increased number of events located within the array). Without manual intervention or sophisticated phase association, 837 phase arrival picks from the U-GPD transfer learning model produce locations with smaller depth 838 errors than the base GPD model and even manually determined phase arrival times (Fig 10D). This 839 is likely a result of more consistent picking and labelling, particularly for S-wave arrivals, which 840 is difficult even for manual analysts to perform consistently, and suggests that very few of the 841 events detected are false. 842

Given the greatly improved computational time over the base GPD model, the small 844 number of training events required and the use of a high-level, user-focused programming library 845 (Keras), this approach is well within the reach of volcano observatories and research groups. 846 Previous studies that analyze the pre-, syn- and post-eruptive periods at Nabro volcano have relied 847 on manually-produced seismic catalogues comprising hundreds of events (e.g., Goitom et al., 848 2015; Hamlyn et al., 2014; the latter locating 658 events over 38 days, a rate of < 18 events per 849 day). Our U-GPD transfer learning model yields a seismic catalogue that is order of magnitudes 850 larger (33,950 events over 396 days, a rate of > 85 events per day; Figs 8 and 9), with smaller 851 location errors (Fig 10), in a matter of hours. Furthermore, as the model processes 1D waveform 852 data, as opposed to 2D spectrogram images in some other existing models (e.g., Dokht et al., 2019; 853 854 Lara et al., 2020; Titos et al., 2020), it runs quickly on high resolution data without using a GPU optimized for deep learning frameworks (32 secs per 24 hours of 100 Hz data on an Intel Core i7 855 856 desktop CPU) and so could easily be deployed for real-time monitoring with limited computing resources or at much larger arrays. The methods and computational times in this paper have relied 857 858 on standard, generic libraries (ObsPy, TensorFlow and Keras); the use of more optimized, compiled code or higher-performance / lower-level languages (e.g., Julia and C) could greatly 859 improve computational times further. 860

861

Beyond phase arrival picking, the generalized waveform features extracted by existing, 862 extensively trained models, such as the GPD model (Fig 1A), could serve as a useful feature 863 extraction system for models designed for other waveform processing tasks. For example, 864 information regarding frequency content and orientation of seismic energy extracted by the GPD 865 model (Fig 1A inset) could reasonably provide useful features for a new model designed to 866 automatically classify volcano seismic event types (e.g., Bueno et al., 2020; Hibert et al., 2017; 867 Lara et al., 2020), particularly when available annotated datasets are small or unbalanced. 868 However, with larger datasets, there is the potential for transfer learning to inhibit learning of new, 869 useful features, particularly if the source and target tasks or data distributions differ considerably. 870

871

The number of seismological studies to date that employ transfer learning is relatively low (e.g., Bueno et al., 2020; Chai et al., 2020; El Zini et al., 2020; Huot et al., 2018; Titos et al., 2020). This is undoubtedly, in part, due to the lack of extensively trained, well-documented, publicly available seismological models. However, the number is likely to grow as more extensive datasets and models are developed and released into the public domain. We credit the availability of the GPD model in the public domain and use of a popular, user-focused machine learning framework (Keras) as the foundation of the work presented in this paper. Such availability facilitates adaptation and experimentation; development of other publicly available models and extensive datasets would aid progress in the field of seismological machine learning.

881

Whilst the application of transfer learning can overcome the perception that deep learning 882 models require a 'large upfront cost' in terms of data and computational resources, the 883 development and benchmarking of large-scale, extensive models and datasets are still imperative 884 885 to push the field of seismological machine learning forwards and extend applications to all aspects of seismic processing and inference. However, it is hoped that applications such as the one 886 887 presented in this paper will motivate the initial investment in the development of such models so that the cost of producing effective task-specific models (e.g., through transfer learning) is 888 889 progressively reduced.

890

891

892 Acknowledgments

893

The seismic data were collected with funding from the Natural Environment Research Council 894 (NERC) project NE/J012297/1 ("Mechanisms and implications of the 2011 eruption of Nabro 895 896 volcano, Eritrea"). The UK seismic instruments and data management facilities were provided under loan number 976 by SEIS-UK at the University of Leicester. The facilities of SEIS-UK are 897 supported by NERC under Agreement R8/H10/64. Author SL was supported by a GW4+ Doctoral 898 Training Partnership studentship from the Natural Environment Research Council (NERC) 899 900 [NE/L002434/1]. Author BG was funded by the Engineering and Physical Sciences Research Council (EPSRC) and the School of Earth Sciences at the University of Bristol. Author MJW was 901 902 funded by UKRI GCRF EP/P028233/1 ("PREPARE") and NERC NE/R017956/1 ("EQUIPT4RISK"). Author JMK was funded by NERC grant NE/R018006/1. Author KVC 903

was supported by the AXA Research Fund. We gratefully acknowledge support from the sponsors 904 of the Bristol University Microseismicity ProjectS (BUMPS) and the NERC Centre for the 905 Observation and Modelling of Earthquakes, volcanoes and Tectonics (COMET). We also 906 gratefully acknowledge the cooperation we received from the Eritrea Institute of Technology, 907 Eritrean government, Southern and Northern Red Sea Administrations, local sub-zones and village 908 administrations. We thank the Department of Mines, Ministry of Energy and Mines for their 909 continued support throughout the project. Special thanks go to Zerai Berhe, Mebrahtu Fisseha, 910 Michael Eyob, Ahmed Mohammed, Kibrom Nerayo, Asresehey Ogbatsien, Andemichael 911 Solomon and Isaac Tuum. We thank Alem Kibreab and Prof. Ghebrebrhan Ogubazghi for their 912 vital help in facilitating the fieldwork. 913

- 914
- 915

916 Data Availability Statement

917

918 All seismic Nabro Urgency Array (Hammond 2011; data from the et al., https://doi.org/10.7914/SN/4H 2011) publicly available through IRIS 919 are Data Services (http://service.iris.edu/fdsnws/dataselect/1/). IRIS Data Services are funded through the Seismological 920 Facilities for the Advancement of Geoscience (SAGE) Award of the National Science Foundation under 921 922 Cooperative Support Agreement EAR-1851048. See Hammond et al. (2011) for further details on 923 waveform data access and availability. Model training, validation and test sets / metadata are archived and available through Zenodo (Lapins et al., 2021; https://doi.org/10.5281/zenodo.4498549). Full code to 924 reproduce our U-GPD transfer learning model, perform model training, run the U-GPD model over 925 continuous sections of data and use model picks to locate events in NonLinLoc (Lomax et al., 2000) are 926 927 available at https://github.com/sachalapins/U-GPD, with the release (v1.0.0) associated with this paper also 928 archived and available through Zenodo (Lapins, 2021; https://doi.org/10.5281/zenodo.4558121).

929

930

931 References

932

Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., ... Zheng, X. (2015). TensorFlow: Large-Scale
Machine Learning on Heterogeneous Distributed Systems (software available from

935 https://www.tensorflow.org).

- 936 Barbedo, J. G. A. (2018). Impact of dataset size and variety on the effectiveness of deep learning and transfer
- 937 learning for plant disease classification. Computers and Electronics in Agriculture, 153, 46-53.
- 938 https://doi.org/10.1016/j.compag.2018.08.013
- 939 Bergstra, J, Bengio, Y. (2012). Random Search for Hyper-Parameter Optimization. Journal of Machine Learning 940 Research, 13(10), 281-305.
- 941 Beyreuther, M., Barsch, R., Krischer, L., Megies, T., Behr, Y., & Wassermann, J. (2010). ObsPy: A python toolbox 942 for seismology. Seismological Research Letters, 81(3), 530–533. https://doi.org/10.1785/gssrl.81.3.530
- 943 Bojanowski, A. (2011). Volcano mix-up. Nature Geoscience, 4(8), 495. https://doi.org/10.1038/ngeo1222
- Bueno, A., Benitez, C., De Angelis, S., Diaz Moreno, A., & Ibanez, J. M. (2020). Volcano-Seismic Transfer 944 945 Learning and Uncertainty Quantification with Bayesian Neural Networks. IEEE Transactions on Geoscience 946 and Remote Sensing, 58(2), 892–902. https://doi.org/10.1109/TGRS.2019.2941494
- 947 Chai, C., Maceira, M., Santos-Villalobos, H. J., Venkatakrishnan, S. V., Schoenball, M., Zhu, W., ... Thurber, C. 948 (2020). Using a Deep Neural Network and Transfer Learning to Bridge Scales for Seismic Phase Picking. 949
- Geophysical Research Letters, 47(16). https://doi.org/10.1029/2020GL088651
- 950 Chollet, F., & and others. (2015). Keras (software available from https://keras.io).
- 951 D'souza, R. N., Huang, P.-Y., & Yeh, F.-C. (2020). Structural Analysis and Optimization of Convolutional Neural 952 Networks with a Small Sample Size. Scientific Reports, 10(834). https://doi.org/10.1038/s41598-020-57866-2
- 953 Daumé, H. (2007). Frustratingly easy domain adaptation. In 45th Annual Meeting of the Association for 954 Computational Linguistics (ACL) (pp. 256-263). https://arxiv.org/abs/0907.1815
- 955 Dokht, R. M. H., Kao, H., Visser, R., & Smith, B. (2019). Seismic Event and Phase Detection Using Time-
- 956 Frequency Representation and Convolutional Neural Networks. Seismological Research Letters, 90(2A), 481-957 490. https://doi.org/10.1785/0220180308
- 958 Donovan, A., Blundy, J., Oppenheimer, C., & Buisman, I. (2018). The 2011 eruption of Nabro volcano, Eritrea:
- 959 perspectives on magmatic processes from melt inclusions. Contributions to Mineralogy and Petrology, 173(1), 960 1-23. https://doi.org/10.1007/s00410-017-1425-2
- 961 Dramsch, J. S., & Lüthje, M. (2018). Deep-learning seismic facies on state-of-the-art CNN architectures. In SEG 962 Technical Program Expanded Abstracts (pp. 2036–2040). Anaheim, CA, USA.
- 963 Efremova, Di. B., Sankupellay, M., & Konovalov, D. A. (2019). Data-Efficient Classification of Birdcall through 964 Convolutional Neural Networks Transfer Learning. In 2019 Digital Image Computing: Techniques and
- 965 Applications (DICTA) (pp. 1–8). Perth, Australia. https://doi.org/10.1109/DICTA47822.2019.8946016
- El Zini, J., Rizk, Y., & Awad, M. (2020). A Deep Transfer Learning Framework for Seismic Data Analysis: A Case 966 967 Study on Bright Spot Detection. IEEE Transactions on Geoscience and Remote Sensing, 58(5), 3202–3212. https://doi.org/10.1109/TGRS.2019.2950888 968
- Feurer, M., Klein, A., Eggensperger, K., Springenberg, J., Blum, M., & Hutter, F. (2015). Efficient and Robust 969 970 Automated Machine Learning. In Advances in Neural Information Processing Systems 28 (NIPS 2015), 2962-971 2970.
- 972 Fromm, M., Kablick III, G., Nedoluha, G., Carboni, E., Grainger, R., Campbell, J., & Lewis, J. (2014). Correcting

- 973 the record of volcanic stratospheric aerosol impact: Nabro and Sarychev Peak. Journal of Geophysical 974 Research: Atmospheres, 119, 10343-10364. https://doi.org/10.1002/2014JD021507 975 Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., ... Lempitsky, V. (2016). Domain-976 Adversarial Training of Neural Networks. Journal of Machine Learning Research, 17(59), 1-35. 977 https://arxiv.org/abs/1505.07818 978 Gibbons, S. J., & Ringdal, F. (2006). The detection of low magnitude seismic events using array-based waveform 979 correlation. Geophysical Journal International, 165(1), 149-166. https://doi.org/10.1111/j.1365-980 246X.2006.02865.x 981 Global Volcanism Program. (2013). Volcanoes of the World, v.4.9.3 (01 Feb 2021). Venzke, E (ed.). Smithsonian 982 Institution. Downloaded 10 Feb 2021. https://doi.org/10.5479/si.GVP.VOTW4-2013 983 Glorot, X., Bordes, A., & Bengio, Y. (2011). Domain adaptation for large-scale sentiment classification: A deep 984 learning approach. In Proceedings of the 28th International Conference on Machine Learning, ICML 2011 985 (pp. 513–520). Bellevue, WA, USA. 986 Goitom, B. (2017). The Nabro volcano, tectontic framework and seismic hazard assessment of Eritrea [doctoral 987 thesis]. University of Bristol. 988 Goitom, B., Oppenheimer, C., Hammond, J. O. S., Grandin, R., Barnie, T., Donovan, A., ... Berhe, S. (2015). First 989 recorded eruption of Nabro volcano, Eritrea, 2011. Bulletin of Volcanology, 77(85). 990 https://doi.org/10.1007/s00445-015-0966-3 991 Hamlyn, J. E., Keir, D., Wright, T. J., Neuberg, J. W., Goitom, B., Hammond, J. O. S., ... Grandin, R. (2014). 992 Seismicity and subsidence following the 2011 Nabro eruption, Eritrea: Insights into the plumbing system of an 993 off-rift volcano. Journal of Geophysical Research: Solid Earth, 119, 8267-8282. 994 https://doi.org/10.1002/2014JB011395 995 Hammond, J., Goitom, B., Kendall, J. M., Ogubazghi, G. (2011). Nabro Urgency Array [Data set]. International 996 Federation of Digital Seismograph Networks. https://doi.org/10.7914/SN/4H 2011 997 Hansen, S. M., & Schmandt, B. (2015). Automated detection and location of microseismicity at Mount St. Helens 998 with a large-N geophone array. Geophysical Research Letters, 42(18), 7390–7397. 999 https://doi.org/10.1002/2015GL064848 1000 He, T., Zhang, Z., Zhang, H., Zhang, Z., Xie, J., & Li, M. (2019). Bag of Tricks for Image Classification with 1001 Convolutional Neural Networks. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition 1002 (CVPR) (pp. 558–567). Long Beach, CA, USA. https://doi.org/10.1109/CVPR.2019.00065 1003 Hibert, C., Provost, F., Malet, J. P., Maggi, A., Stumpf, A., & Ferrazzini, V. (2017). Automatic identification of 1004 rockfalls and volcano-tectonic earthquakes at the Piton de la Fournaise volcano using a Random Forest 1005 algorithm. Journal of Volcanology and Geothermal Research, 340, 130–142. 1006 https://doi.org/10.1016/j.jvolgeores.2017.04.015 1007 Huot, F., Biondi, B., & Beroza, G. C. (2018). Jump-starting neural network training for seismic problems. In SEG 1008 Technical Program Expanded Abstracts (pp. 2191–2195).
- 1009 Ioffe, S., & Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal

- 1010 covariate shift. In *32nd International Conference on Machine Learning (ICML)* (pp. 448–456). Lille, France.
 1011 https://arxiv.org/abs/1502.03167
- Kingma, D. P., & Ba, J. (2014). Adam: A Method for Stochastic Optimization. In *3rd International Conference on Learning Representations (ICLR)* (pp. 1–15). http://arxiv.org/abs/1412.6980
- 1014 Kirkpatrick, J., Pascanu, R., Rabinowitz, N., Veness, J., Desjardins, G., Rusu, A. A., ... Hadsell, R. (2017).
- 1015 Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences*,
- 1016 *114*(13), 3521–3526. https://doi.org/10.1073/pnas.1611835114
- Klein, A., Falkner, S., Bartels, S., Henning, P., & Hutter, F. (2017). Fast Bayesian Optimization of Machine
 Learning Hyperparameters on Large Datasets. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, PMLR 54*, 528-536. https://arxiv.org/abs/1605.07079
- Krischer, L., Megies, T., Barsch, R., Beyreuther, M., Lecocq, T., Caudron, C., & Wassermann, J. (2015). ObsPy: A
 bridge for seismology into the scientific Python ecosystem. *Computational Science and Discovery*, *8*, 014003.
 https://doi.org/10.1088/1749-4699/8/1/014003
- Lahr, J. C., Chouet, B. A., Stephens, C. D., Power, J. A., & Page, R. A. (1994). Earthquake classification, location,
 and error analysis in a volcanic environment: implications for the magmatic system of the 1989–1990
- 1025 eruptions at redoubt volcano, Alaska. *Journal of Volcanology and Geothermal Research*, 62(1–4), 137–151.
 1026 https://doi.org/10.1016/0377-0273(94)90031-0
- Lapins, S. (2021). Python notebooks to accompany paper 'A Little Data goes a Long Way: Automating Seismic
 Phase Arrival Picking at Nabro Volcano with Transfer Learning' (Version v1.0.0) [Archived GitHub
 repository]. Zenodo. https://doi.org/10.5281/zenodo.4558121
- Lapins, S., Goitom, B., Kendall, J.-M., Werner, M. J., Cashman, K. V. & Hammond, J. O. S. (2021). Training,
 Validation and Test Sets for paper 'A Little Data goes a Long Way: Automating Seismic Phase Arrival
- 1032Picking at Nabro Volcano with Transfer Learning' [Dataset]. Zenodo. https://doi.org/10.5281/zenodo.4498549
- Lapins, S., Roman, D. C., Rougier, J., De Angelis, S., Cashman, K. V, & Kendall, J.-M. (2020). An examination of
 the continuous wavelet transform for volcano-seismic spectral analysis. *Journal of Volcanology and Geothermal Research*, 389, 106728. https://doi.org/10.1016/j.jvolgeores.2019.106728
- Lara, F., Lara-Cueva, R., Larco, J. C., Carrera, E. V, & León, R. (2020). A deep learning approach for automatic
 recognition of seismo-volcanic events at the Cotopaxi volcano. *Journal of Volcanology and Geothermal Research*, (xxxx), 107142. https://doi.org/10.1016/j.jvolgeores.2020.107142
- Lengliné, O., Duputel, Z., & Ferrazzini, V. (2016). Uncovering the hidden signature of a magmatic recharge at Piton
 de la Fournaise volcano using small earthquakes. *Geophysical Research Letters*, 43(9), 4255–4262.
 https://doi.org/10.1002/2016GL068383
- Li, H., Xu, Z., Taylor, G., Studer, C., & Goldstein, T. (2017). Visualizing the Loss Landscape of Neural Nets.
 Advances in Neural Information Processing Systems (NIPS), 6389–6399. https://arxiv.org/abs/1712.09913
- 1044 Li, W., Duan, L., Xu, D., & Tsang, I. W. (2014). Learning with augmented features for supervised and semi-
- 1045 supervised heterogeneous domain adaptation. *IEEE Transactions on Pattern Analysis and Machine*
- 1046 Intelligence, 36(6), 1134–1148. https://doi.org/10.1109/TPAMI.2013.167

- Lin, T.-Y., Goyal, P., Girshick, R., He, K., & Dollar, P. (2017). Focal Loss for Dense Object Detection. In 2017
 IEEE International Conference on Computer Vision (ICCV) (pp. 2999–3007). Venice, Italy: IEEE.
 https://doi.org/10.1109/ICCV.2017.324
- Lomax, A., Virieux, J., Volant, P., & Berge, C. (2000). Probabilistic earthquake location in 3D and layered models:
 Introduction of a Metropolis-Gibbs method and comparison with linear locations. In C. H. Thurber & N.
- 1052 Rabinowitz (Eds.), *Advances in Seismic Event Location* (pp. 101–134). Asterdam: Kluwer.
- Maclaurin, D., Duvenaud, D., & Adams, R. (2015). Gradient-based Hyperparameter Optimization through
 Reversible Learning. *Proceedings of the 32nd International Conference on Machine Learning, PLMR, 37,* 2113-2122. https://arxiv.org/abs/1502.03492
- Maqsood, M., Nazir, F., Khan, U., Aadil, F., Jamal, H., Mehmood, I., & Song, O. (2019). Transfer Learning
 Assisted Classification and Detection of Alzheimer's Disease Stages Using 3D MRI Scans. *Sensors*, 19(11),
 2645. https://doi.org/10.3390/s19112645
- McNutt, S. R., & Roman, D. C. (2015). Volcanic Seismicity. In *The Encyclopedia of Volcanoes* (2nd Edition, pp. 1060 1011–1034). Elsevier. https://doi.org/10.1016/B978-0-12-385938-9.00059-6
- Megies, T., Beyreuther, M., Barsch, R., Krischer, L., & Wassermann, J. (2011). ObsPy what can it do for data
 centers and observatories? *Annals of Geophysics*, 54(1), 47–58. https://doi.org/10.4401/ag-4838
- Mousavi, S. M., Ellsworth, W. L., Zhu, W., Chuang, L. Y. & Beroza, G. C. (2020). Earthquake transformer an
 attentive deep-learning model for simultaneous earthquake detection and phase picking. *Nature Communications*, *11*, 3952. https://doi.org/10.1038/s41467-020-17591-w
- Mousavi, S. M., Zhu, W., Sheng, Y., & Beroza, G. C. (2019). CRED: A Deep Residual Network of Convolutional
 and Recurrent Units for Earthquake Signal Detection. *Scientific Reports*, *9*, 10267.
 https://doi.org/10.1038/s41598-019-45748-1
- Nair, V., & Hinton, G. E. (2010). Rectified Linear Units Improve Restricted Boltzmann Machines. In *Proceedings* of the 27th International Conference on Machine Learning (pp. 807–814). Haifa, Israel.
- 1071 https://doi.org/10.5555/3104322.3104425
- Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10), 1345–1359. https://doi.org/10.1109/TKDE.2009.191
- 1074 Razavian, A. S., Azizpour, H., Sullivan, J., & Carlsson, S. (2014). CNN features off-the-shelf: An astounding
 1075 baseline for recognition. In 2014 IEEE Computer Society Conference on Computer Vision and Pattern
- 1076 *Recognition Workshops* (pp. 512–519). Columbus, OH, USA: IEEE.
- 1077 https://doi.org/10.1109/CVPRW.2014.131
- 1078 Ronneberger, O., Fischer, P., & Brox, T. (2015). U-Net: Convolutional Networks for Biomedical Image
- Segmentation. In N. Navab, J. Hornegger, W. M. Wells, & A. F. Frangi (Eds.), *Medical Image Computing and Computer-Assisted Intervention (MICCAI) 2015, Part III, Lecture Notes in Computer Science* (Vol. 9351, pp.
- 1081 234–241). Springer International Publishing. https://doi.org/10.1007/978-3-319-24574-4 28
- 1082 Ross, Z. E., Meier, M., & Hauksson, E (2018a). P Wave Arrival Picking and First-Motion Polarity Determination
- 1083 With Deep Learning. Journal of Geophysical Research: Solid Earth, 123(6), 5120-5129.

- 1084 https://doi.org/10.1029/2017JB015251
- Ross, Z. E., Meier, M., Hauksson, E., & Heaton, T. H. (2018b). Generalized Seismic Phase Detection with Deep
 Learning. *Bulletin of the Seismological Society of America*, 108(5A), 2894–2901.

1087 https://doi.org/10.1785/0120180080

Ross, Z. E., Yue, Y., Meier, M., Hauksson, E., & Heaton, T. H. (2019). PhaseLink: A Deep Learning Approach to
 Seismic Phase Association. *Journal of Geophysical Research: Solid Earth*, *124*(1), 856-869.

1090 https://doi.org/10.1029/2018JB016674

- Shelly, D. R., Beroza, G. C., & Ide, S. (2007). Non-volcanic tremor and low-frequency earthquake swarms. *Nature*,
 446(7133), 305–307. https://doi.org/10.1038/nature05666
- Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on Image Data Augmentation for Deep Learning. *Journal of Big Data*, 6(1). https://doi.org/10.1186/s40537-019-0197-0
- 1095 Snoek, J., Rippel, O., Swersky, K., Kiros, R., Satish, N., Sundaram, N., Patwary, M. A., Prabhat, & Adams, R. P.

(2015). Scalable Bayesian Optimization Using Deep Neural Networks. *Proceedings of the 32nd International Conference on Machine Learning, PMLR, 37*, 2171-2180. https://arxiv.org/abs/1502.05700

- Sun, B., Feng, J., & Saenko, K. (2016). Return of frustratingly easy domain adaptation. In *30th AAAI Conference on Artificial Intelligence (AAAI)* (pp. 2058–2065). https://arxiv.org/abs/1511.05547
- Sun, C., Shrivastava, A., Singh, S., & Gupta, A. (2017). Revisiting Unreasonable Effectiveness of Data in Deep
 Learning Era. In *2017 IEEE International Conference on Computer Vision (ICCV)* (pp. 843–852). Venice,
 2017. https://doi.org/10.1109/ICCV.2017.97
- Tan, Y. J., Waldhauser, F., Ellsworth, W. L., Zhang, M., Zhu, W., Michele, M., Chiaraluce, L., Beroza, G. C. &
 Segou, M. (2021). Machine-Learning-Based High-Resolution Earthquake Catalog Reveals How Complex
- Fault Structures Were Activated during the 2016 2017 Central Italy Sequence. *The Seismic Record*, 1(1).
 https://doi.org/10.1785/0320210001
- Theys, N., Campion, R., Clarisse, L., Brenot, H., van Gent, J., Dils, B., ... Ferrucci, F. (2013). Volcanic SO2 fluxes
 derived from satellite data: a survey using OMI, GOME-2, IASI and MODIS. *Atmospheric Chemistry and Physics*, 13, 5945–5968. https://doi.org/10.5194/acp-13-5945-2013
- Titos, M., Bueno, A., García, L., Benítez, C., & Segura, J. C. (2020). Classification of Isolated Volcano-Seismic
 Events Based on Inductive Transfer Learning. *IEEE Geoscience and Remote Sensing Letters*, *17*(5), 869–873.
 https://doi.org/10.1109/LGRS.2019.2931063
- Tompson, J., Goroshin, R., Jain, A., LeCun, Y., & Bregler, C. (2015). Efficient Object Localization Using
 Convolutional Networks. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 648–656). Boston, MA, USA: IEEE. https://doi.org/10.1109/CVPR.2015.7298664
- Tran, K. T., Griffin, L. D., Chetty, K., & Vishwakarma, S. (2020). Transfer learning from audio deep learning
 models for micro-Doppler activity recognition. In *2020 IEEE International Radar Conference, (RADAR)* (pp.
 584–589). Washington, DC, USA. https://doi.org/10.1109/RADAR42522.2020.9114643
- 1119 Tzeng, E., Hoffman, J., Darrell, T., & Saenko, K. (2015). Simultaneous Deep Transfer Across Domains and Tasks.
- 1120 In 2015 IEEE International Conference on Computer Vision (ICCV) (pp. 4068–4076). Santiago, Chile: IEEE.

- 1121 https://doi.org/10.1109/ICCV.2015.463
- van den Ende, M. P. A., & Ampuero, J. P. (2020). Automated Seismic Source Characterization Using Deep Graph
 Neural Networks. *Geophysical Research Letters*, 47(17), 1–11. https://doi.org/10.1029/2020GL088690
- van den Oord, A., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., ... Kavukcuoglu, K. (2016).
 WaveNet: A Generative Model for Raw Audio. https://arxiv.org/abs/1609.03499
- 1126 Withers, M., Aster, R., Young, C., Beiriger, J., Harris, M., Moore, S., & Trujillo, J. (1998). A comparison of select
- trigger algorithms for automated global seismic phase and event detection. *Bulletin of the Seismological Society of America*, 88(1), 95–106.
- Woollam, J., Rietbrock, A., Bueno, A., & De Angelis, S. (2019). Convolutional Neural Network for Seismic Phase
 Classification, Performance Demonstration over a Local Seismic Network. *Seismological Research Letters*, 1–
 12. https://doi.org/10.1785/0220180312
- Yeck, W. L., Patton, J. M., Johnson C. E., Kragness, D., Benz, H. M., Earle, P. S., Guy, M. R., & Ambruz, N. B.
 (2019). GLASS3: A Standalone Multiscale Seismic Detection Associator. *Bulletin of the Seismological*

1134 Society of America, 109(4), 1469-1478. https://doi.org/10.1785/0120180308

- Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in deep neural networks?
 Advances in Neural Information Processing Systems, 27. https://arxiv.org/abs/1411.1792
- Yu, F., & Koltun, V. (2016). Multi-scale context aggregation by dilated convolutions. In *4th International Conference on Learning Representations (ICLR)*. https://arxiv.org/abs/1511.07122
- Zamir, A. R., Sax, A., Shen, W., Guibas, L., Malik, J., & Savarese, S. (2018). Taskonomy: Disentangling Task
 Transfer Learning. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (pp. 3712–

1141 3722). Salt Lake City, UT, USA: IEEE. https://doi.org/10.1109/CVPR.2018.00391

- Zech, J. R., Badgeley, M. A., Liu, M., Costa, A. B., Titano, J. J., & Oermann, E. K. (2018). Variable generalization
 performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study.
- 1144 PLOS Medicine, 15(11), 1–17. https://doi.org/10.1371/journal.pmed.1002683
- Zhu, W., & Beroza, G. C. (2019). PhaseNet: a deep-neural-network-based seismic arrival-time picking method. *Geophysical Journal International*, *216*(1), 261–273. https://doi.org/10.1093/gji/ggy423
- 1147 Zhuang, F., Cheng, X., Luo, P., Pan, S. J., & He, Q. (2015). Supervised representation learning: Transfer learning
- with deep autoencoders. In *IJCAI International Joint Conference on Artificial Intelligence* (pp. 4119–4125).
 Buenos Aires, Argentina.
- 1150 Zhuang, F., Qi, Z., Duan, K., Xi, D., Zhu, Y., Zhu, H., ... He, Q. (2020). A Comprehensive Survey on Transfer
- 1151 Learning. Proceedings of the IEEE, 1–34. https://doi.org/10.1109/JPROC.2020.3004555



Journal of Geophysical Research: Solid Earth

Supporting Information for

A Little Data goes a Long Way: Automating Seismic Phase Arrival Picking at Nabro Volcano with Transfer Learning

Sacha Lapins¹, Berhe Goitom¹, J-Michael Kendall², Maximilian J. Werner¹, Katharine V. Cashman¹ and James O. S. Hammond³

¹School of Earth Sciences, University of Bristol, UK; ²Department of Earth Sciences, University of Oxford, UK; ³Department of Earth and Planetary Sciences, Birkbeck, University of London, UK.

Corresponding author: Sacha Lapins (sacha.lapins@bristol.ac.uk)

Contents of this file

Figure S1 Figure S2 Figure S3 Figure S4 Figure S5 Figure S6 Text S7 Table S8

Introduction

In this supporting information document, we provide full details of our U-GPD model architecture (Fig S1), cross-entropy loss function used for model training (Text S7), velocity model used to associate detected phase arrivals and locate events (Table S8), and comparisons of U-GPD model performance on test data using varying levels of data augmentation (Fig S2). We also provide plotted locations using the U-GPD model, base GPD model (with two threshold values) and existing manual catalogue (Figs S3-S6).



Figure S1. U-GPD model architecture.

No Data Augmentation

Decimation Only







Figure S2. Results of U-GPD model, trained with different levels of data augmentation, applied to our test dataset.



Figure S3. Event locations using P- and S-wave phase arrivals detected by U-GPD model with 0.4 threshold (no. of events = 33,950).



Figure S4. Event locations using P- and S-wave phase arrivals detected by original GPD model with 0.9 threshold (no. of events = 41,007).



Figure S5. Event locations using P- and S-wave phase arrivals detected by original GPD model with 0.99 threshold (no. of events = 13,319).



Figure S6. Event locations using P- and S-wave phase arrivals from original manual catalogue (no. of events = 2,984).

Text S7.

Supervised deep learning models are optimized through minimizing a loss function that measures the distance between the model's prediction and ground truth labels. To train our U-GPD transfer learning model, we use a cross entropy loss function, defined as:

$$L(p,q) = -\sum_{i=1}^{3}\sum_{x}p_i(x)\log q_i(x)$$

where *i* is class/channel number (P, S or N), *x* are the datapoints in our waveform, $p_i(x)$ is the ground truth probability distribution for *x* and $q_i(x)$ is the model's predicted distribution for *x*. The value of this loss function decreases as the model predictions converge towards the ground truth labels.

Table S8. P-wave 1D velocity model used to associate model phase arrival picks and locate events (linear gradient between layers, assumed Vp/Vs ratio of 1.76 for S-wave velocity model).

Depth for top of layer (km b.s.l.)	Vp (km / s)
-1.5	4.10
3.0	6.10
8.0	6.80
25.0	7.40