

The temporal limits of predicting fault failure

Kun Wang¹, Christopher Johnson², Kane Bennett¹, and Paul Johnson¹

¹Los Alamos National Laboratory

²Los Alamos National Laboratory, Los Alamos National Laboratory

November 26, 2022

Abstract

Machine learning models using seismic emissions can predict instantaneous fault characteristics such as displacement in laboratory experiments and slow slip in Earth. Here, we address whether the acoustic emission (AE) from laboratory experiments contains information about near-future frictional behavior. The approach uses a convolutional encoder-decoder containing a transformer layer. We use as input progressively larger AE input time windows and progressively larger output friction time windows. The attention map from the transformer is used to interpret which regions of the AE contain hidden information corresponding to future frictional behavior. We find that very near-term predictive information is indeed contained in the AE signal, but farther into the future the predictions are progressively worse. Notably, information for predicting near future frictional failure and recovery are found to be contained in the AE signal. This first effort predicting future fault frictional behavior with machine learning will guide efforts for applications in Earth.

The temporal limits of predicting fault failure

Kun Wang^{1,2*}, Christopher W. Johnson^{1*}, Kane C. Bennett¹, Paul A. Johnson¹

¹Los Alamos National Laboratory, Geophysics Group, Los Alamos, N.M.

²Center for Nonlinear Studies, Los Alamos National Laboratory, Los Alamos, NM, 87545, USA

Key Points:

- The limits of machine learning models are explored for predicting the future frictional behavior of laboratory earthquake slip
- A deep learning model is applied to acoustic emission data broadcast from the fault to predict future friction
- Acoustic emissions broadcast from the fault are found to carry information used for predicting future fault friction

Statement:

This manuscript is a non-peer reviewed preprint submitted to ESSOAr.

*NOTE: KW and CWJ contributed equally to this paper

Corresponding author: Paul A. Johnson, paj@lanl.gov

Abstract

Machine learning models using seismic emissions can predict instantaneous fault characteristics such as displacement in laboratory experiments and slow slip in Earth. Here, we address whether the acoustic emission (AE) from laboratory experiments contains information about near-future frictional behavior. The approach uses a convolutional encoder-decoder containing a transformer layer. We use as input progressively larger AE input time windows and progressively larger output friction time windows. The attention map from the transformer is used to interpret which regions of the AE contain hidden information corresponding to future frictional behavior. We find that very near-term predictive information is indeed contained in the AE signal, but farther into the future the predictions are progressively worse. Notably, information for predicting near future frictional failure and recovery are found to be contained in the AE signal. This first effort predicting future fault frictional behavior with machine learning will guide efforts for applications in Earth.

Plain Language Summary

Over the last 5 years it has been shown that machine learning is a powerful tool to learn about the inside of faults from the acoustic emissions that the fault broadcasts. For instance, the emissions can inform us of the instantaneous fault displacement and friction in addition to timing of an upcoming earthquake in laboratory experiments as well as for the phenomenon of slow slip in Earth. Here we show that by applying a deep learning approach that has shown striking results in natural language processing and computer vision, the seismic emissions also contain foreshadowing information about the immediate future of fault friction. Remarkably, the emissions can tell us if the laboratory fault is about to fail, and how it will begin to recover.

Introduction

Applying machine learning techniques to fault slip from laboratory shear experiments has demonstrated that continuous acoustic/seismic signals emanating from an active fault contain rich information regarding the instantaneous and future characteristics of that system. The fault acoustic emissions are imprinted with information regarding its current state and where it is in the slip cycle. Indeed, the statistical features of the continuous seismic signal emitted from the fault and identified by the machine learning model allow the prediction of instantaneous fault friction, displacement, fault gouge thickness, slip velocity and the timing of the next laboratory earthquake ('labquake'), e.g., (Rouet-Leduc et al., 2017, 2018; Lubbers et al., 2018; Hulbert et al., 2019, 2019; P. A. Johnson et al., 2021; Wang et al., 2021). The labquake magnitude can also be predicted to a lesser degree, and with considerably less precision and accuracy than quantities related to timing (Hulbert et al., 2019).

Other model inputs such as seismic wave velocity, wave amplitude (Shokouhi et al., 2021), and geodetic measurements (Corbi et al., 2019) have also been successfully applied to predict instantaneous and future fault characteristics. Since the seminal paper first describing machine learning model predictions by Rouet-LeDuc and colleagues (Rouet-Leduc et al., 2017), many works now exist describing these behaviors—too many to reference here. In addition, a Kaggle competition based on laboratory earthquake prediction was also recently held where many hundreds of competitors took part in developing machine learning models for predicting timing of the upcoming slip event (P. A. Johnson et al., 2021).

In previous works, only predictions of the instantaneous characteristics and/or timing (magnitude) of the next slip event were attempted. The question we address here is the following—is there information contained in the continuous seismic signal regarding the

near future fault behavior including and beyond the next slip event? One might expect for instance, that due to a slip event the granular gouge is reset due to the intense, system wide perturbation associated with the labquake. We have observed previously that intense, externally applied dynamic perturbations—dynamic earthquake triggering—can have persistent impact through multiple slip cycles (P. A. Johnson et al., 2008). Such perturbations are considerably stronger than the seismic emission associated with the slip event however. Therefore it is unclear whether or not fault slip physical information can be carried within the fault gouge through a slip event into the next slip cycle.

Predicting future physical behavior in the form of bulk friction from current seismic emissions is the focus of this work. Such work will ultimately guide efforts in Earth to determine if there is near-term predictive information regarding fault physical characteristics such as displacement and earthquake timing contained in the continuous seismic emissions from active faults.

In the following, we describe the laboratory experiment and the accompanying data, and the deep learning modeling applied to those data. We then describe the prediction results and close with discussion of the results and conclusions.

Methods

Deep Learning Model

The deep learning model developed here for the future activity in the laboratory fault experiments is based on the Transformer architecture (Vaswani et al., 2017) designed as an encoder-decoder, or sequence-to-sequence, model that has shown remarkable results in natural language processing, computer vision, and time series modeling applications (Dosovitskiy et al., 2020; Li et al., 2019; Zhou et al., 2021). This particular model type was selected because the self-attention mechanism greatly improves the performance on sequential data compared to earlier designs that combine an attention model with recurrent neural networks such as Long Short-Term Memory or Gated Recurrent Units (Bahdanau et al., 2014; Luong et al., 2015). The success of the future prediction task relies on the ability of the Transformer model to learn a long-range dependency. Prior applications using Transformer models for long-sequence time-series forecasting have shown good results up to about 1000 future time points. These improvements rely on more efficient self-attention mechanisms such as LogSparse (Li et al., 2019) and ProbSparse (Zhou et al., 2021) that reduce computational complexity and memory usage.

In this work, the deep learning model is designed to predict the future friction coefficient as model output using as input acoustic emission (AE) obtained from laboratory biaxial shear data. Our approach uses a simplified version of the Transformer model that takes a time series vector as input and encodes it to a high dimensional representation of the signal that is then fed to the latent space tensor. The coefficients are passed through the Transformer and output to a second latent space. A decoder branch outputs friction coefficients at the corresponding future time step.

The first step is to design the input and output branches of the model as a convolutional encoder-decoder (CED) model and independently pre-train as an autoencoder (Figure 1). The input-branch decoder is needed to obtain the embedding of the unit window of signals for a functional latent space. The output-branch encoder is needed for generating the target μ vectors for the teacher-forcing training stage of the Transformer. Note that the input decoder and the output encoder (the blocks in dashed lines in Figure 1) are not used in the Transformer model application, but only for model pre-training. Following the pre-training procedure, the parameters of the input and output CED models are fixed and the Transformer alone is trained using the future friction prediction data.

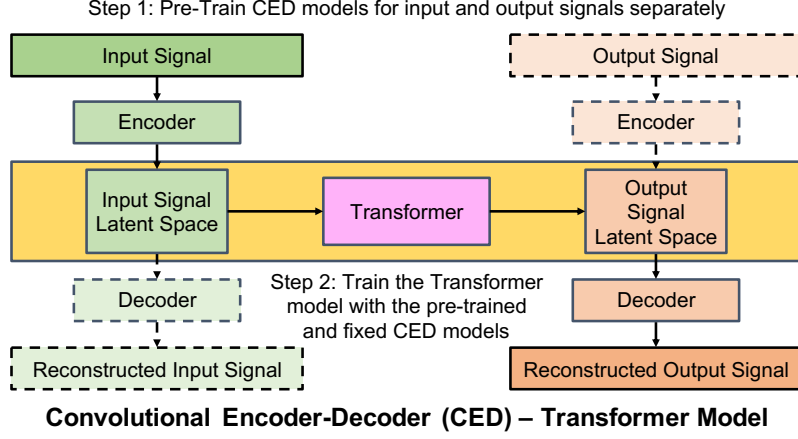


Figure 1: The deep learning model architecture used for pre-training and application to future predictions. The input signal (green column) and output signal (pink column) show the convolutional encoder-decoders that are connected by the Transformer via the latent space. The encoder consists of convolutional layers that take the input acoustic emission time series and map it into a sequence of vectors in a smaller space. The vectors are fed into a Transformer (yellow box) to predict output latent vectors and then to a decoder which transforms them into an output sequence with the target fault friction. The model components noted by the dashed lines are used in Step 1 for pre-training as an autoencoder to obtain the embedding vectors and not used in Step 2 with the Transformer for future predictions.

The latent vectors of the input AE signals are obtained by the pre-trained encoder model branch. For example, a total AE history of 2.56 s is encoded into a sequence of 10 input latent vectors, each corresponding to 0.256 s windows of the input signal, with an embedding dimension of 64. The self-attention is the learned score representing the time-dependent relative importance of the input signal to the model for each prediction. The latent vectors are computed using a multi-head attention layer (Vaswani et al., 2017) in the Transformer’s encoder. These self-attention processed AE vectors are the keys and values to the multi-head attention in the Transformer’s decoder and the embedding vectors of the output μ generated in the past are the queries. The Transformer decodes to the next output embedding vector to match the target μ vector.

The approach applied here uses categorical predictions that represent different embedding, or latent, vectors from the groups of 256 signal data. This is analogous to translating input texts to output token IDs for natural language processing with a Transformer model except that here the dictionary terms are replaced with groups of 256 signal output data points. The alternative approach, which was tested during model development, is to predict embedding vectors of dimension 64 for a regression analysis. Hereafter, we refer to using the model to predict latent vectors of output friction coefficients as the regression approach and predicting category IDs of the outputs as the classification approach. The advantage of the classification approach is that the Transformer only needs to predict the probability logits, with the maximum value corresponding to the highest probability category. In the regression approach, alternatively, the Transformer needs to predict all dimensions of the latent vectors accurately to successfully estimate the output values. We obtain categories of the output embedding vectors through the Vector Quantised Variational AutoEncoder (VQ-VAE) (Oord et al., 2017) that learns a codebook of discrete vectors in the latent space. All continuous vectors after the convolutional encoder are compared against the codebook and the discrete latent vector with the smallest euclidean distance

is fed to the convolutional decoder. The codebook is updated by an exponential moving average (Oord et al., 2017) during the training stage. In a classification model the target categories are the discrete latent vectors of the codebook.

The model training is performed in 2 steps. In Step 1 the input and output CED branches are trained independently using the mean squared error (MSE) loss function between the target and reconstructed signal, plus the loss associated to the commitment loss in the VQ-VAE and the L2 regularization. In Step 2 the Transformer is trained using the classification approach for imbalanced data which focuses more on the categories that are difficult for the model to predict (Lin et al., 2017), since the number of values associated with the friction drops is much smaller than when the friction increases or is stable. Additionally, Step 2 is tested by training the Transformer using the regression approach with the MSE loss function between the target latent vectors of the output friction and the predicted vectors. In both training steps a batch size of 8 is used with the Adam optimizer and a custom learning rate scheduler. The training is terminated when the reconstruction loss on the validation data does not diminish for 50 epochs and the lowest validation loss is used for the final model.

Hyperparameter selection and model design optimization

Model hyperparameters and architecture are optimized for laboratory experiment p4677 data using a Bayesian optimizer (scikit-optimize package; Head et al., 2021) to train thousands of models with slightly different designs. Optimization is performed for the two CED models for embedding the input AE and output frictional coefficient, respectively, and the Transformer model for predicting future output embeddings. For the output CED model, the number of convolutional layers n_{out}^l , filter size f_{out} , kernel size k_{out} , the number of codes in the VQ-VAE’s codebook n_{out}^{vq} , the activation function act_{out} and the L2 regularization weight $l2_{out}$ are optimized. Similarly for the input CED model, n_{in}^l , f_{in} , k_{in} , n_{in}^{vq} , act_{in} , and $l2_{in}$ are optimized. The hyperparameters for the Transformer include the number of encoder layers n_{enc}^{trsf} , the number of decoder layers n_{dec}^{trsf} , the number of attention heads n_{attn}^h , and the dropout rate r_{drop} (see Table S1). The search space is comprised of a range of values for each variable listed. The objective function for hyperparameter minimization is the MSE between the target future friction coefficients and the predicted frictions on the validation data of experiment No. p4677 using AE history of 2.56 s to predict 0.256 s into the future. Each optimization attempt is run for 500 iterations, with the first 100 iterations using random parameter selections. The procedure is distributed onto 16 GPUs with 5 processes running per GPU (80 models trained simultaneously) to rapidly test the search space. The procedure is iterated to adjust the search spaces until the final values are not at the boundary limits. The optimized hyperparameters for the final model design are provided in Supporting Information Table S1.

Laboratory experiment

We use laboratory data obtained from a bi-axial shear device. The device is a double-direct shear apparatus comprised of three blocks with two fault gouge layers, e.g., (Marone, 1998; P. A. Johnson et al., 2013a; Rouet-Leduc et al., 2017; P. A. Johnson et al., 2021). The three block system is held in place by a fixed horizontal load, while the center block is driven by a piston oriented perpendicular to the horizontal load. The shear stress and friction on the gouge layers, from the center drive block and measured with a load cell, is an important experimental output. The simulated fault gouge comprises class IV spheres (dimension from 105–149 μm). The initial layer thickness is 2×4 mm (two layers), and the roughened interfaces with the drive block have dimensions 10 cm \times 10 cm. The drive block vertical displacement rate is 5 $\mu\text{m/s}$, corresponding to a strain rate of approximately 1.2×10^{-3} /s. The apparatus is servo-controlled so that constant normal stress and displacement rate of the drive block are maintained at ± 0.1 kN and ± 0.1 $\mu\text{m/s}$, respectively. The apparatus is monitored via computer to record load on the drive block and drive block displacement

at 1 kHz. Once the system has been sheared to the point where it is in steady state conditions, the laboratory fault gouge layers fail in quasi-periodic cycles of stick and slip.

Laboratory Data Analysis

For the present analysis, experiment No. p4677 is used to train and validate the model, in which a fixed normal stress on the blocks of 2.5 MPa is maintained. Continuous AE broadcast from the fault zone were recorded via a Verasonics multichannel recording system and used as the input signal to the deep learning model, as described above. The coefficient of friction (μ) was used as the target output signal. The entire experiment p4677 signals are split into 60% training and 40% validation data. These segments of signals are further split into sliding windows with the window size for inputs l_{in} , the window size for outputs l_{out} , the step size for the input windows l_{step} , and the lag between the start of the input and output windows l_{lag} . For the first step of training the CED models with instantaneous prediction data, $l_{in} = l_{out} = l_{step}$, $l_{lag} = 0s$, and l_{in} is equal to the unit window length (0.256 s in the Results) of signals embedded into one latent vector. For the second step of training the Transformer in the latent space with future prediction data while keeping the pre-trained CED fixed, l_{in} and l_{out} are no longer required to be the same, $l_{lag} = l_{in}$ and $l_{step} = l_{out}$ to ensure zero overlapping between input and output windows. The values of 0.256 s, 0.512 s, 1.024 s, 1.536 s, 2.048 s, 2.56 s are tested for l_{in} and l_{out} in order to study the window size effect on the prediction accuracy of the future fault slips. The size of the datasets depends on the above length hyperparameters. For example, for future prediction data with $l_{in} = 2.56s$, $l_{out} = 0.256s$, $l_{lag} = 2.56s$ and $l_{step} = 0.256s$, there are 694 pairs of input and output signal windows in the training dataset and 459 pairs in the validation dataset. Experiment No. p4581 is used as the testing dataset, in which progressively larger normal loads were applied and the steady states are at 3, 4, 5, 6, 7, 8 MPa normal stresses. In the following we focus on the 3MPa load analysis.

All input and output signals are normalized by subtracting the mean and dividing by the standard deviation using the statistics extracted from the training signal data. For the p4677 data, the statistics from the training signals are 8.878 ± 21.018 for the input AE signals and 0.652 ± 0.0413 for the output μ . When making predictions using experiment p4581 data with increasing normal loads, the statistics are extracted from the first 40% of the signals at each normal stresses for AE and 0.435 ± 0.024 for μ .

Results and Discussion

We first compared the future prediction accuracy between the classification approach and the regression approach, as well as between using the full waveform and AE magnitude signals as model inputs. We found classification outperforms regression and therefore will focus on these results. The optimal model we found was using the classification approach with the AE magnitudes as inputs after hyperparameter optimization (Table S1).

The future prediction accuracy using progressively increasing input and output window sizes on the validation and testing signals are tabulated in the matrix shown in Figure 2. A larger matrix showing more window sizes is provided in Figure S1 of Supporting Information. The Mean Absolute Percent Error (MAPE) values are noted in the matrix figures for validation and testing, respectively. We note that average MAPE/MSE scores are not a perfect indicator of model performance; however, they do provide an indication of prediction accuracy (and are commonly applied metrics). The validation data set uses a portion of experiment p4677 data the model had not seen during training. The trained model is then applied to the p4581 testing data set for an experiment that were conducted at a slightly different load level (3 MPa vs. 2.5 MPa). The results show for the validation set that near-term predictions are about the same with the exception of the shortest window length. In contrast, the testing data shows degrading prediction as the input window length

increased beyond 2.048 s. Similarly, predictions further into the future degrade with window length. The best results for the training data are in the range of 1.024 - 2.048 seconds input and 0.256 s prediction window length. Predictions degrade with longer window lengths into the future. In summary, the model predictions are best for short future time sequences, and predictions beyond approximately 2.5 s are less good.

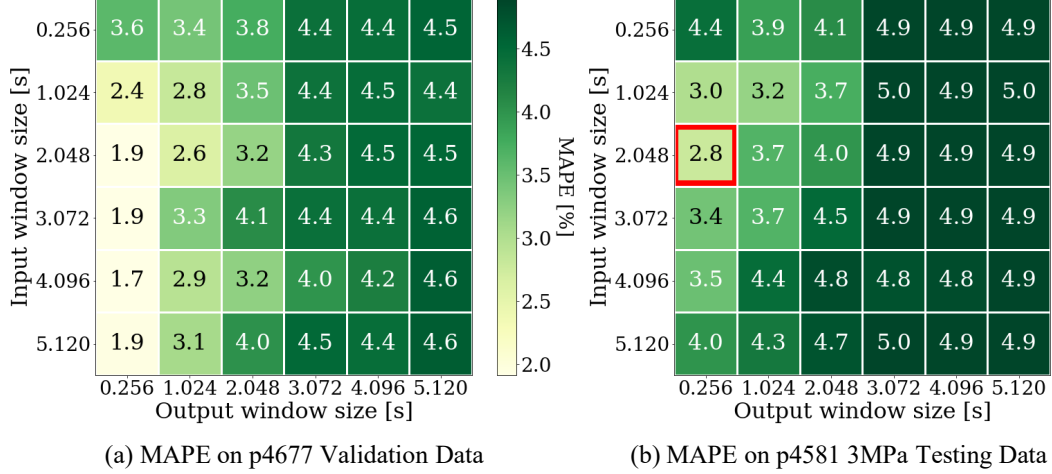


Figure 2: Grid search for 36 different input and prediction window lengths, showing MAPE of future predictions using different input/output window sizes. The validation data (left) is from experiment p4677, 2.5 MPa applied load data with 60%/40% train/validation split, and the testing data (right) are from the experiment p4581 dataset conducted at 3 MPa. The window highlighted with the red rectangle has the smallest testing data error and these model predictions are plotted in Figure 3. Supplementary matrices of MAPE results with finer window size (0.256 s step size) until 2.56 s is provided in Figure S1.

Figure 3 shows the predictions for the time window exhibiting the best MAPE—the matrix element denoted with the red box in Figure 2. Figure 4 shows predictions for the diagonal elements shown in the right panel in Figure 2. In Figure 3 one can see that the predictions are reasonably good overall. Many frictional failures are predicted, as are the succeeding recoveries. For some slip events, the event is not well predicted and the onset of recovery is therefore poorly predicted, e.g., events at approximately 44 and 54 sec. Large precursors may be predicted (e.g., approximately 120 sec) or missed entirely (e.g., 15 sec). While a slip event may be predicted, the full friction failure magnitude is not as was observed in previous works applying many different models, e.g., (P. A. Johnson et al., 2021). In general the predictions improve through the slip cycle. The fact that the model is able to predict the failure event just before it takes place is encouraging. The onset of frictional recovery post-failure is also predicted in many cases suggesting there is knowledge contained in the AE regarding the onset of the future labquake cycle. The intense acoustic emission associated with the failure event apparently does not entirely erase fault system memory.

From the self-attention calculated in the Transformer latent-space, the relative importance of different parts of the AE input for future predictions can be analyzed by plotting the attention scores to visualize their distribution with respect to the AE history. The top panel of Figure 3 shows the mean attention scores and their standard deviation, computed from successive sliding input windows. The attention scores represent the average of the attention heads for a given input window. An expanded view of eight stress cycles is shown

in the lower panel to better visualize the variance in the signal and the accuracy of the friction predictions. Where the attention scores are largest indicates time intervals when the AE signal is most important for the future friction predictions. In general, the self attention scores indicate the model predictions are better in time intervals where impulsive precursors occur immediately preceding failure but also where the scores fluctuate the most.

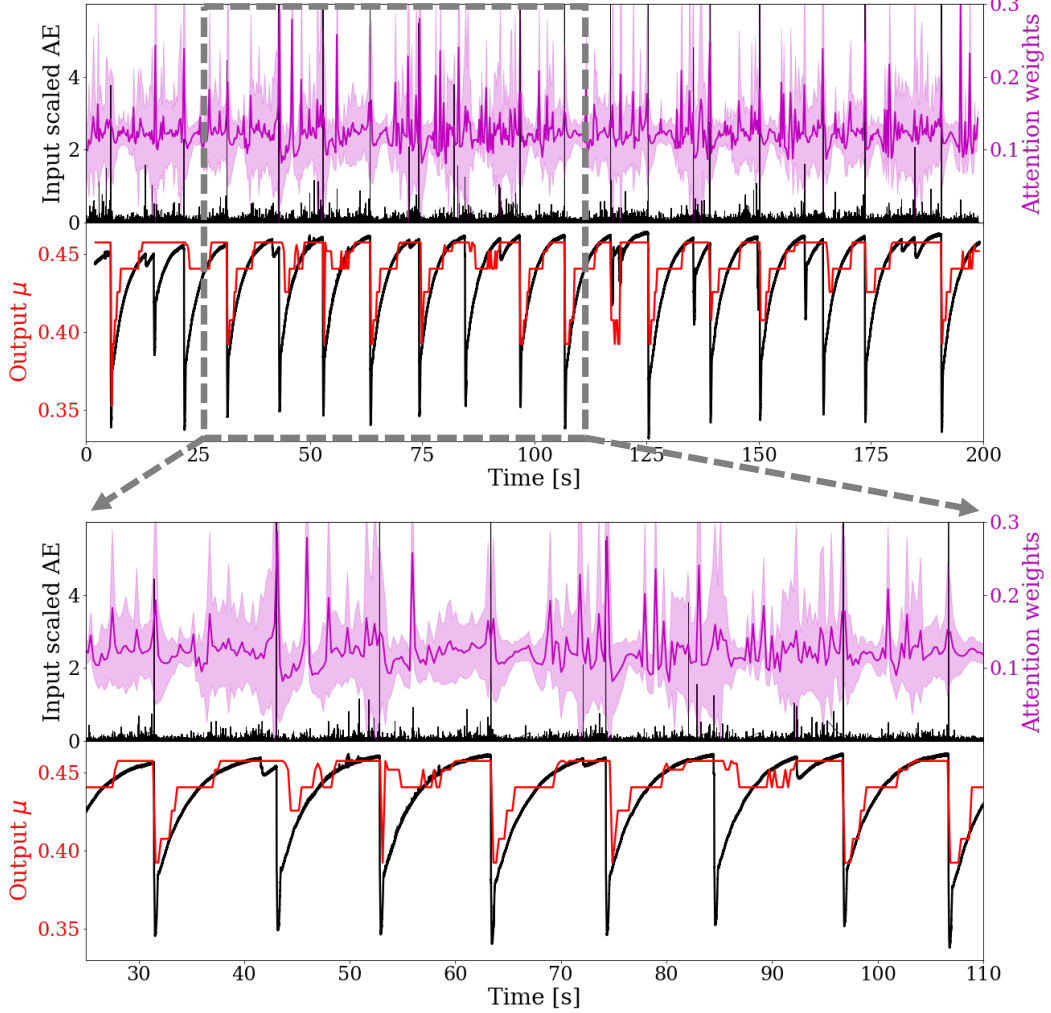


Figure 3: Attention scores shown with the AE and predicted friction values for the model shown in the red rectangle in Figure 2. In each figure the upper panel shows the mean attention score ± 1 standard deviation when predicting segments of 0.256 s window using 2.048 s history. The lower panel is the corresponding friction prediction. The bottom figure shows an expanded view from 25s to 110s to better highlight the time series details. In the 8 significant fault failures during this period, the timing of 5 failures are predicted at the correct time step, 2 are predicted after the slip event, and 1 is missed by the model.

The future predictions using the p4581 testing data for successively larger input and output windows as shown along the diagonal in the right panel in Figure 2 are presented in Figure 4. The time series data illustrate what the MAPE values listed in Figure 2 indicate. The 1.024 s input and output underscore the best prediction capability, and the prediction degrades progressively as the prediction time window size increases. Looking at the 0.256s

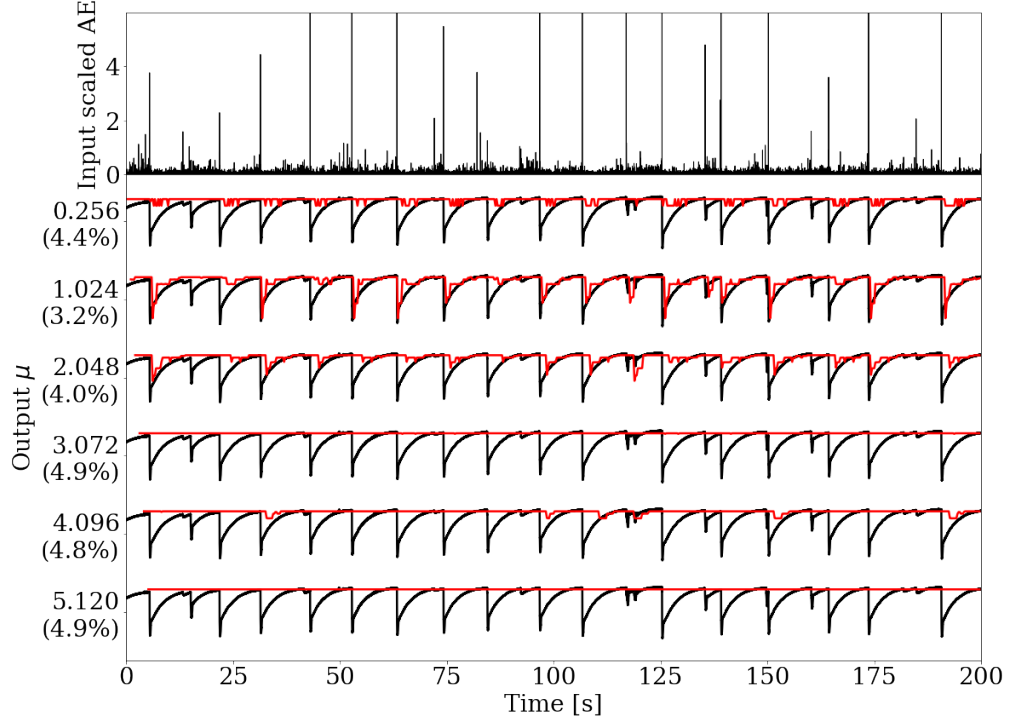


Figure 4: Test data (p4581) and the model predictions of the friction using equal length input and output window sizes (the diagonal of Figure 2b). The friction prediction is shown in red and the experiment measurements in black. The y-axis labels shows the input and output window lengths 0.256s, 1.024s, 2.048s, 3.072s, 4.096s, 5.12s and the MAPE error in parentheses. More prediction examples are provided in Supplementary materials Figure S2 and S3.

results the predictions are quite poor, possibly due to such a small input window size. For the time windows longer than 3 seconds there appears to be insufficient information contained in the AE signal to predict anything, but this also may reflect the model generalization. We note that poorer accuracy for longer window lengths could arise from the limitations of the Transformer model to capture long-term dependency in the latent space (embedding) of signals, although proven to be very successful in applications for natural language processing.

Previous work shows good instantaneous prediction capability contained in AE continuous waveforms. Signal energy was identified as key to instantaneous and time-to-failure prediction using both decision tree and deep learning models in the laboratory experiment data (e.g., Rouet-Leduc et al., 2017; P. A. Johnson et al., 2021) and for slow slip in Earth (e.g., Hulbert et al., 2019; C. W. Johnson & Johnson, 2021). The analysis described here shows there is also information contained in the continuous signal regarding the immediate future behavior. The model attention scores focus on the region preceding frictional failure where precursors occur. This is logical near failure, based on our knowledge of classical seismic precursors in the laboratory (e.g., P. A. Johnson et al., 2013b) and those observed often, but not always, in Earth (Scholz, 2002; Bouchon et al., 2013). The results also indicate that within the inter-event portion of the earthquake cycle, near term friction prediction is also possible. Perhaps most interesting is that the seismic signal contains information predictive of the fault friction immediately beyond the failure event when one might expect the system to be unstable. The attention scores do not shed light as to why this is the case, but it is truly intriguing and future endeavors might provide clues as to what features of the signal contain this information. Applying machine learning approaches such as that described in this work, to continuous seismic signals as well as other geophysical data, will continue advancing our understanding of fault activity while advancing earthquake hazards assessment.

Acknowledgements

KW and PAJ acknowledge support by the U.S. Department of Energy, Office of Science, Office of Basic Energy Sciences, Chemical Sciences, Geosciences, and Biosciences Division under grant 89233218CNA000001. KW also acknowledges support by the Center for Nonlinear Studies (CNLS) at Los Alamos National Laboratory. CWJ and KCB acknowledge Institutional Support (Laboratory Directed Research and Development) at Los Alamos National Laboratory. The authors declare no competing interests. We thank Chris Marone for the laboratory data.

References

- Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Bouchon, M., Durand, V., Marsan, D., Karabulut, H., & Schmittbuhl, J. (2013). The long precursory phase of most large interplate earthquakes. *Nature Geoscience*, 6. doi: 10.1038/ngeo1770
- Corbi, F., Sandri, L., Bedford, J., Funicello, F., Brizzi, S., Rosenau, M., & Lallemand, S. (2019). Machine learning can predict the timing and size of analog earthquakes. *Geophysical Research Letters*, 46(3), 1303-1311. doi: 10.1029/2018GL081251
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ... others (2020). An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*.
- Head, T., Kumar, M., Nahrstaedt, H., Louppe, G., & Shcherbatyi, I. (2021, October). *scikit-optimize/scikit-optimize*. Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.5565057> doi: 10.5281/zenodo.5565057
- Hulbert, C., Rouet-Leduc, B., Johnson, P. A., Ren, C. X., Rivière, J., Bolton,

- D. C., & Marone, C. (2019). Similarity of fast and slow earthquakes illuminated by machine learning. *Nature Geoscience*, 12(1), 69–74. doi: 10.1038/s41561-018-0272-8
- Johnson, C. W., & Johnson, P. A. (2021). Learning the low frequency earthquake daily intensity on the central san andreas fault. *Geophysical Research Letters*, 48(15), e2021GL092951. doi: 10.1029/2021GL092951
- Johnson, P. A., Ferdowsi, B., Kaproth, B. M., Scuderi, M., Griffa, M., Carmeliet, J., ... Marone, C. (2013a). Acoustic emission and microslip precursors to stick-slip failure in sheared granular material. *Geophysical Research Letters*, 40(21), 5627–5631.
- Johnson, P. A., Ferdowsi, B., Kaproth, B. M., Scuderi, M., Griffa, M., Carmeliet, J., ... Marone, C. (2013b). Acoustic emission and microslip precursors to stick-slip failure in sheared granular material. *Geophysical Research Letters*, 40(21), 5627–5631. doi: 10.1002/2013GL057848
- Johnson, P. A., Rouet-Leduc, B., Pyrak-Nolte, L. J., Beroza, G. C., Marone, C. J., Hulbert, C., ... Reade, W. (2021). Laboratory earthquake forecasting: A machine learning competition. *Proceedings of the National Academy of Sciences*, 118(5). doi: 10.1073/pnas.2011362118
- Johnson, P. A., Savage, H. M., Knuth, M. W., Gombert, J., & Marone, C. (2008). Effects of acoustic waves on stick-slip in granular media and implications for earthquakes. *Nature*, 451, 57–60.
- Li, S., Jin, X., Xuan, Y., Zhou, X., Chen, W., Wang, Y.-X., & Yan, X. (2019). Enhancing the locality and breaking the memory bottleneck of transformer on time series forecasting. *Advances in Neural Information Processing Systems*, 32, 5243–5253.
- Lin, T.-Y., Goyal, P., Girshick, R., He, K., & Dollár, P. (2017). Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision* (pp. 2980–2988).
- Lubbers, N., Bolton, D., Mohd-Yusof, J., Marone, C., Barros, K., & Johnson, P. (2018). Earthquake catalog-based machine learning identification of laboratory fault states and the effects of magnitude of completeness. *Geophysical Research Letters*, 45. doi: 10.1029/2018GL079712
- Luong, M.-T., Pham, H., & Manning, C. D. (2015). Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- Marone, C. (1998). Laboratory-derived friction laws and their application to seismic faulting. *Annual Review of Earth and Planetary Sciences*, 26(1), 643–696. doi: 10.1146/annurev.earth.26.1.643
- Oord, A. v. d., Vinyals, O., & Kavukcuoglu, K. (2017). Neural discrete representation learning. *arXiv preprint arXiv:1711.00937*.
- Rouet-Leduc, B., Hulbert, C., Bolton, D. C., Ren, C. X., Riviere, J., Marone, C., ... Johnson, P. A. (2018). Estimating fault friction from seismic signals in the laboratory. *Geophysical Research Letters*, 45(3), 1321–1329. doi: 10.1002/2017GL076708
- Rouet-Leduc, B., Hulbert, C., Lubbers, N., Barros, K., Humphreys, C. J., & Johnson, P. A. (2017). Machine learning predicts laboratory earthquakes. *Geophysical Research Letters*, 44(18), 9276–9282. doi: 10.1002/2017GL074677
- Scholz, C. (2002). *The mechanics of earthquake faulting*. doi: 10.1017/CBO9780511818516
- Shokouhi, P., Girkar, V., RiviĀšre, J., Shreedharan, S., Marone, C., Giles, C. L., & Kifer, D. (2021). Deep learning can predict laboratory quakes from active source seismic data. *Geophysical Research Letters*, 48(12), e2021GL093187. doi: 10.1029/2021GL093187
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). Attention is all you need. In *Advances in neural*

- information processing systems* (pp. 5998–6008).
- Wang, K., Johnson, C. W., Bennett, K. C., & Johnson, P. A. (2021, 12). Predicting fault slip via transfer learning. *Nature Communications*, 12. doi: 10.1038/s41467-021-27553-5
- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., & Zhang, W. (2021). Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of aaai*.