Built on
**LDRD**
Laboratory Directed Research
and Development

**SAND20XX-XXXXR**
**LDRD PROJECT NUMBER**: 209745
**LDRD PROJECT TITLE**: Seismic Phase Identification with Speech
Recognition Algorithms
**PROJECT TEAM MEMBERS**: Timothy J. Draelos, Stephen Heck, Jennifer
Galasso, Ronald Brogan[1]

## ABSTRACT:

Seismic signals are composed of the seismic waves (phases) that reach a sensor, similar to the way speech signals are composed of phonemes that reach a listener's ear. Large/small seismic events near/far from a sensor are similar to loud/quiet speakers with high/low-pitched voices. We leverage ideas from speech recognition for the classification of seismic phases at a seismic sensor. Seismic Phase ID is challenging due to the varying paths and distances an event takes to reach a sensor, but there is consistent structure of the makeup (e.g. ordering) of the different phases arriving at the sensor.

Current Phase ID techniques do not take into account the global spectrotemporal structure of a waveform that includes a phase arrival for an event. Together with scalar value measurements of seismic signal detections, we use the seismogram and its spectrogram as inputs to a merged deep neural network with convolutional and recurrent layers to learn the frequency structure over time of different phases. Our best results come from the use of a Long Short-Term Memory network merged with horizontal slowness, amplitude, SNR of signal detections, and the time since the previous signal detection. The classification performance of First-P phases versus non-First-P (95.6% class average accuracy) and suggests a significant impact on the reduction of false and missed events in seismic signal processing pipelines.

## INTRODUCTION:

Since 2012, neural-inspired deep neural networks have revolutionized speech recognition [1], an extremely competitive multimillion dollar industry. Because of the similarities of seismic and acoustic signals and sensors, speech recognition algorithms are well-suited to seismic Phase ID, the classification of seismic phases at a seismic sensor. Seismic Phase ID is challenging due to the varying paths and distances seismic events take to reach a sensor, but there is discriminative spectrotemporal structure of the different phases arriving at the sensor and consistent ordering of the phases in time based on geophysical properties.

Seismic signals are composed of seismic waves reaching a sensor, similar to speech signals being composed of phonemes reaching a listener's ear. Moreover, large/small seismic events near/far from a sensor are similar to loud/quiet speakers with high/low-pitched voices. Neural-

---

[1] ENSCO, Inc.

Sandia National Laboratories

U.S. DEPARTMENT OF
**ENERGY**

inspired deep neural networks have revolutionized speech recognition, a competitive multimillion dollar industry. The Deep Speech 2 speech recognition algorithm [2] captures the frequency structure over time of an input acoustic signal and identifies words as sequences of phonetic labels. Figure 1 shows how the Deep Speech 2 architecture can be applied to seismic Phase ID. For both acoustic and seismic signals, spectrograms are valuable computational steps to capture spectrotemporal content. The early layers of Deep Speech 2 utilize convolutional layers to identify larger geometric structures in spectrograms. Deeper recurrent layers perform sequence learning on the spectrotemporal structure from the convolutional layers. In theory, the Connectionist Temporal Classification (CTC) loss function will be able to apply a sequence of phase labels to the input signal without segmentation of individual phases and with some gaps between phases from the same event.
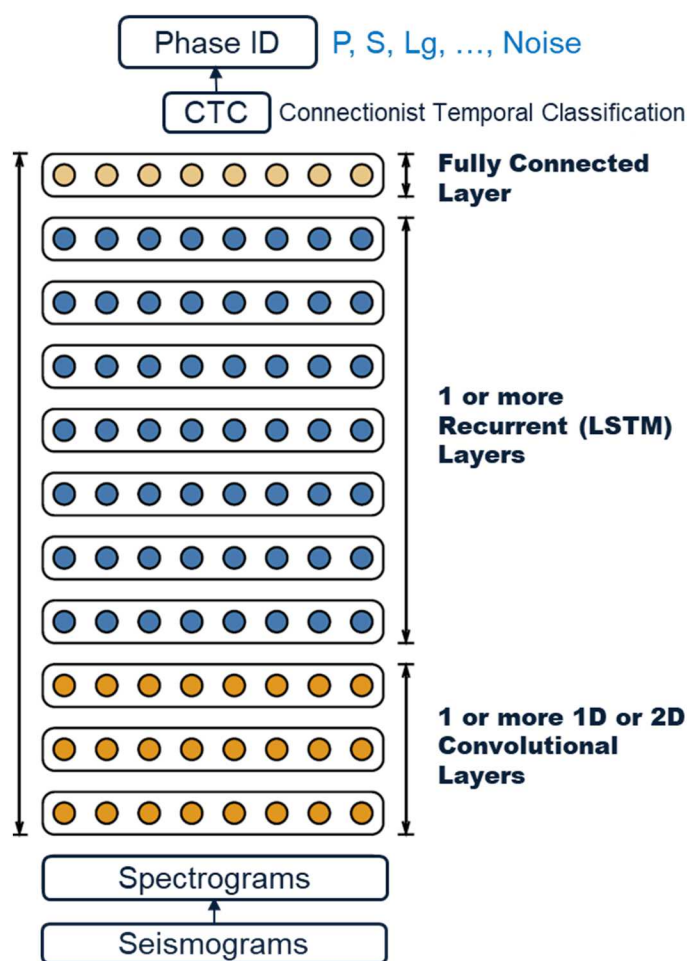


*Figure 1. Conceptual application of Deep Speech 2 neural network to seismic Phase ID. Input to the network is at the bottom and output from the network is at the top.*

Current Phase ID techniques do not take into account the global spectrotemporal structure of phase or event seismograms. Seismic Phase ID can follow a similar approach as Deep Speech 2 to classify the elements of an event's seismogram that a station records, but it is nonetheless

extremely challenging and exploratory. The primary challenges lie in the difference between acoustic signals and seismic signals.

- Seismic events create different types of energy that reach seismic sensor stations via differing paths through the earth. The distance between the event and the sensor affects the spectrotemporal content of an event's waveform recorded at a station.
- Depending on the size of the event and distance between station and event, seismic signals can have very low signal-to-noise ratios (SNR), unlike most speech recognition applications.
- When a seismic phase is detected at a station, identifying the phase can help determine where and when the event occurred and which phases are associated with the event.

Figure 2 shows a record section of an event detected at many stations, illustrating the propagation of seismic waves in distance versus time. The same event produces multiple phases recorded at each station. The distance from a station to the event and the type of energy received from the event determines the labels of phases received at each station. The predictable sequence of phases in time suggests that a sequence learning approach to Phase ID is worth exploring.
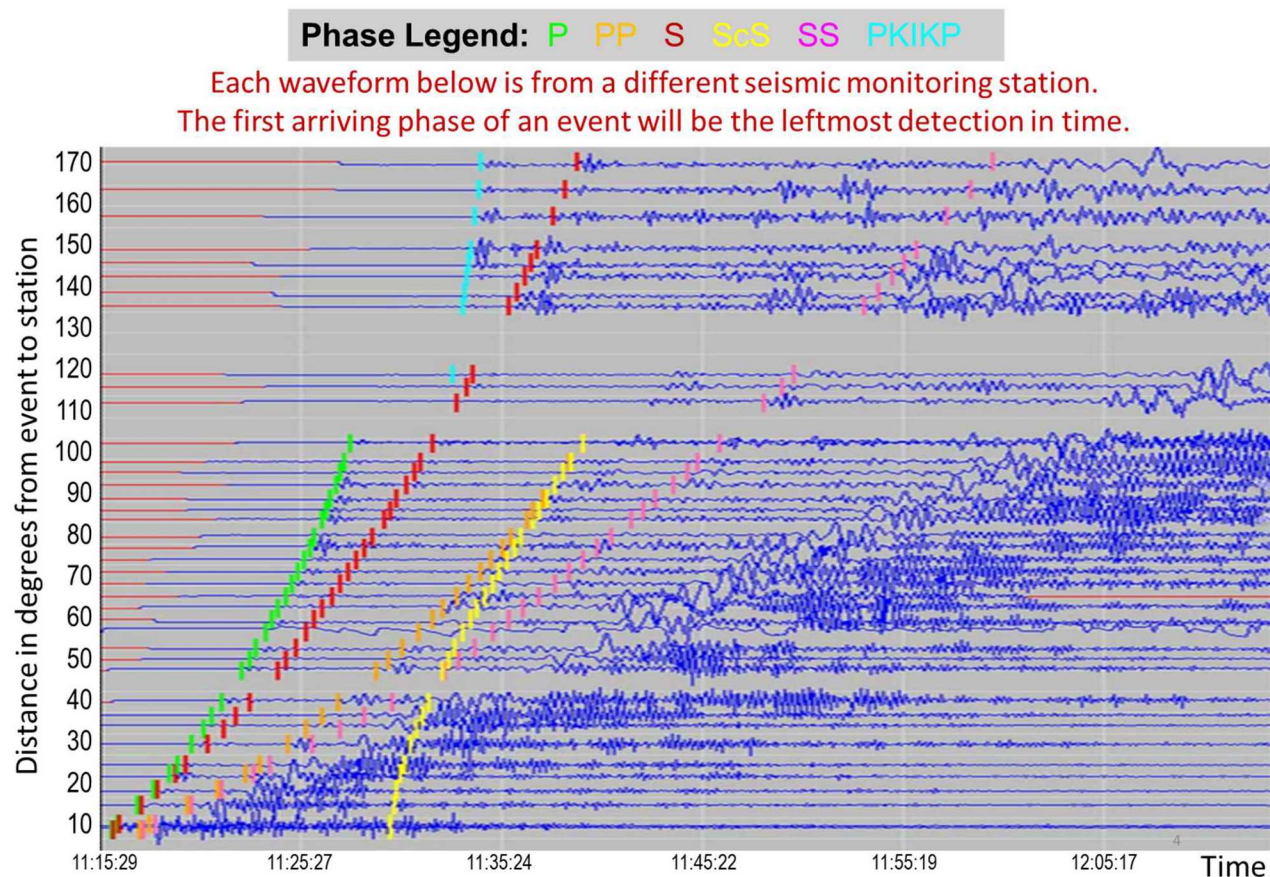


*Figure 2. Example record section displaying seismograms recorded for an event at multiple seismic monitoring stations at varying distances from the event.*
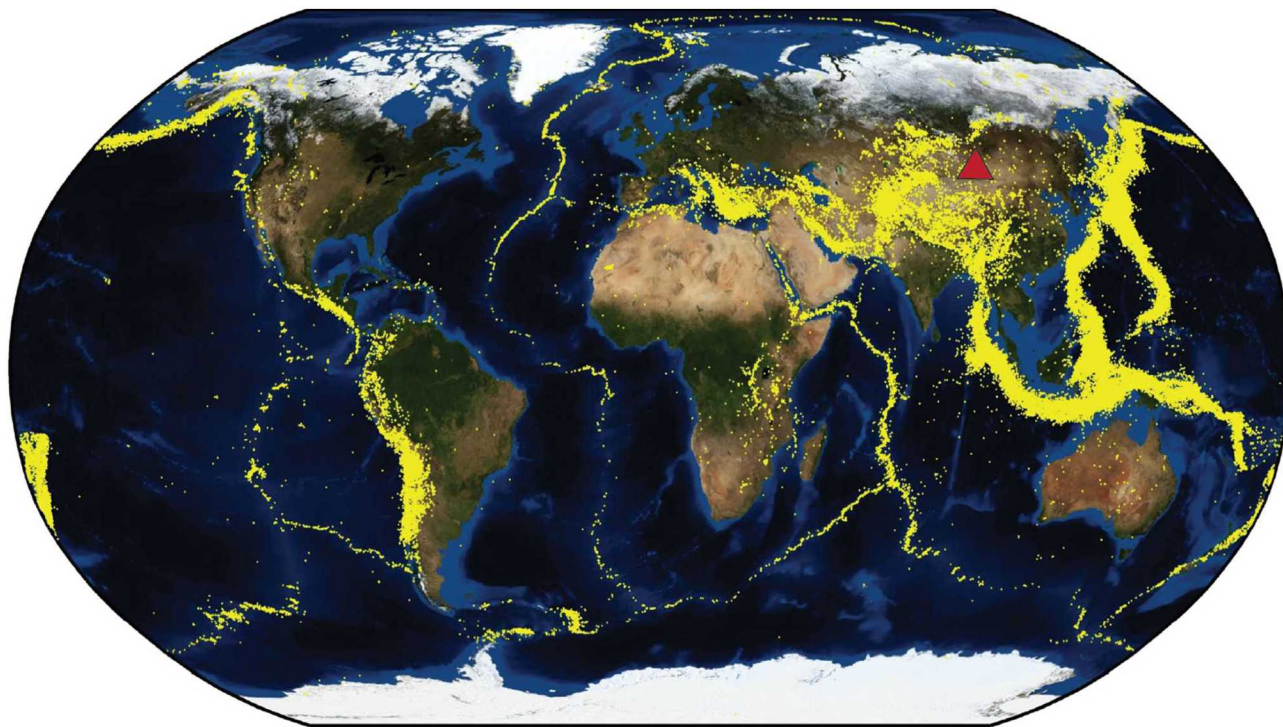
Sequence learning, utilizing a geophysics model of phase arrivals in time, is the ultimate goal of our work on Phase ID, but current results apply to the classification of segmented phases based on their arrival times.

## DETAILED DESCRIPTION OF EXPERIMENT/METHOD:

Phases detected at the Makanchi seismic array (MKAR) in Khazakstan were used in this research. MKAR has been a primary station in the international monitoring system (IMS) since January, 2002, and is sensitive to detect phases from seismic events around the world. Figure 3Figure 3 shows a map with the location of MKAR and detected events. The type of phase detected at MAKR is a function of the type of energy from the event, the distance from the event, the depth of the event, and the magnitude of the event.



*Figure 3. Events detected by MKAR (red triangle). Yellow dots indicate event locations.*

Table 1 lists 30 different seismic phases detected at MKAR since 2002 that were used in our experiments.

*Table 1. Descriptions of seismic phases used in Phase ID experiments. Phases in Bold are considered First-P. The number of each phases listed in the IMS LEB Association table between 2002 through 2017 is in the Count column.*

| Phase | Description | Count |
|---|---|---|
| **P** | A longitudinal wave, bottoming below the uppermost mantle; also an upgoing longitudinal wave from a source below the uppermost mantle. | 271220 |
| **Pn** | Any P wave bottoming in the uppermost mantle or an upgoing P wave from a source in the uppermost mantle. | 34579 |
| **Pg** | At short distances, either an upgoing P wave from a source in the upper crust or a P wave bottoming in the upper crust. At larger distances also arrivals caused by multiple P-wave reverberations inside the whole crust with a group velocity around 5.8 km/s. | 2979 |
| **PKP** | Unspecified P wave bottoming in the core. | 17430 |
| **PKPab** | P wave bottoming in the upper outer core; ab denotes the retrograde branch of the PKP caustic. | 5836 |
| **PKPbc** | P wave bottoming in the lower outer core; bc denotes the prograde branch of the PKP caustic. | 11795 |
| Lg | A wave group observed at larger regional distances and caused by superposition of multiple S-wave reverberations and SV to P and/or P to SV conversions inside the whole crust. The maximum energy travels with a group velocity around 3.5 km/s. | 20426 |
| S | A shear wave, bottoming below the uppermost mantle; also an upgoing shear wave from a source below the uppermost mantle. | 2747 |
| Sn | Any S wave bottoming in the uppermost mantle or an upgoing S wave from a source in the uppermost mantle. | 13554 |
| PP | Free surface reflection of P wave leaving a source downwards. | 2594 |
| pP | P resulting from reflection of upgoing P at the free surface. | 8101 |
| sP | P resulting from converted reflection of upgoing S at the free surface. | 614 |
| PcP | P reflection from the core-mantle boundary. | 17326 |
| ScP | S to P converted reflection from the core-mantle boundary. | 5557 |
| PKiKP | P wave reflected from the inner core boundary. | 7068 |
| PKKPbc | PKKP bottoming in the lower outer core. | 3329 |
| pPKPbc | PKPbc resulting from reflection of upgoing P at the free surface. | 1698 |
| Pdiff | P diffracted along the core-mantle boundary in the mantle. | 1509 |
| PKP2 | Free surface reflection of PKP. | 1387 |
| PKhKP | a precursor to PKPdf due to scattering near or at the core-mantle boundary. | 1367 |
| pPKP | PKP resulting from reflection of upgoing P at the free surface. | 1061 |
| SKPbc | SKP bottoming in the lower outer core. | 834 |
| PKP2bc | Free surface reflection of PKP; bc denotes the prograde branch of the PKP caustic. | 652 |
| PKKPab | PKKP bottoming in the upper outer core. | 585 |
| PKKP | Unspecified P wave reflected once from the inner side of the core-mantle boundary. | 488 |
| SKP | Unspecified S wave traversing the core and then the mantle as P. | 404 |
| P4KPbc | P wave reflected 3 times from inner side of the CMB; bc denotes the prograde branch of the PKP caustic. | 356 |
| pPKPab | PKPab resulting from reflection of upgoing P at the free surface. | 255 |
| SKKPbc | SKKP bottoming in the lower outer core. | 217 |
| P3KPbc | P wave reflected 2 times from inner side of the CMB; bc denotes the prograde branch of the PKP caustic. | 200 |

## Preliminary Phase ID

The simplest phase classification experiment we conducted was between two classes of similar phases: Pg, Pn, and P vs. S, Sn, and Lg. This level of Phase ID is potentially the simplest test of

our classification approach. The most impactful Phase ID experiment conducted is classification of First-P phases vs. non-First-P phases, given all the different phases in Table 1.

## Impactful Phase ID

The Probabilistic Event Detection, Association, and Location (PEDAL) algorithm uses seismic signal detections to create probabilistically realistic hypothetical events on a 3-D grid covering the earth, including depth. It determines event location and origin time by assuming all signal detections are first-arriving compressional waves (P phases) at each station and finds the grid point where the pairwise combinations of detection observations (arrival time, azimuth, and horizontal slowness) at different stations compared with predictions is highest (assuming an event originating at the grid point). First-P phases are defined from 0 to 180 degrees and include Pg at local distances, Pn at regional distances, and P, then PKP, then, PKPbc and PKPab at teleseismic distances. In practice, many automated detections are not First-P phases and can lead to the false events by the PEDAL algorithm. If a phase classifier can filter all but First-P phases for PEDAL during event detection and location, the number of false and missed events could be significantly reduced and the quality of valid events improved. Therefore, the primary Phase ID experiment conducted in this work is that of First-P vs. not-First-P classification.

## Classification Features

The data used to discriminate between classes of phases are described in Table 2 and include beamed waveform data extracted 5 seconds before the arrival time of the signal detection through 10 seconds after the arrival time, spectrograms computed from the waveforms, and scalar measurement values from the signal detection. Spectrograms were computed on 601 samples of the seismograms using 256-sample Fast-Fourier Transforms (FFTs) with 248-sample overlap.

*Table 2. Data elements derived from a signal detection used for Phase ID.*

| Description | Notation | Data Type | Normalization |
|---|---|---|---|
| Horizontal Slowness | Sh | Scalar | If Sh < 0, exclude phase<br>If Sh > 35, set Sh = 35<br>Sh = Sh / 35 |
| Detection Amplitude | Amp | Scalar | If Amp < 0, exclude phase<br>If Amp > 100, set Sh = 100<br>Amp = (9 * Amp / 100) + 1<br>Amp = log10(Amp) |
| Detection Signal to Noise Ratio | SNR | Scalar | If SNR > 100, set SNR = 100<br>SNR = (9 * SNR / 100) + 1<br>SNR = log10(SNR) |
| Time since prior detection | T | Scalar | If first phase, set T = 0<br>If T > 1600, set T = 1600<br>T = log10(T) / log10(1601) |
| Seismogram | Se | 601-element vector | Se = 2 * (Se-min(Se) / (max(Se)-min(Se)) - 1 |
| Spectrogram | Sp | Array of 36 time steps x 76 frequency steps | Sp = log10(9*Sp + 1)<br>Sp = Sp / max(Sp) |

The classifier used for seismic Phase ID can utilize all the features described in Table 2 simultaneously in a merged deep neural network (DNN) shown in Figure 4. The signal detection measurements for each phase detected by MKAR (based on arrival ID in the IMS LEB association table) are taken from the arrival table and waveform samples from 5 seconds before the phase arrival time through 10 seconds after are extracted from which spectrograms are computed. All or subsets of this information can be fed into the DNN simultaneously together with the phase's ground truth label during training. During testing the same input can be input and the output of the DNN can be compared against ground truth.



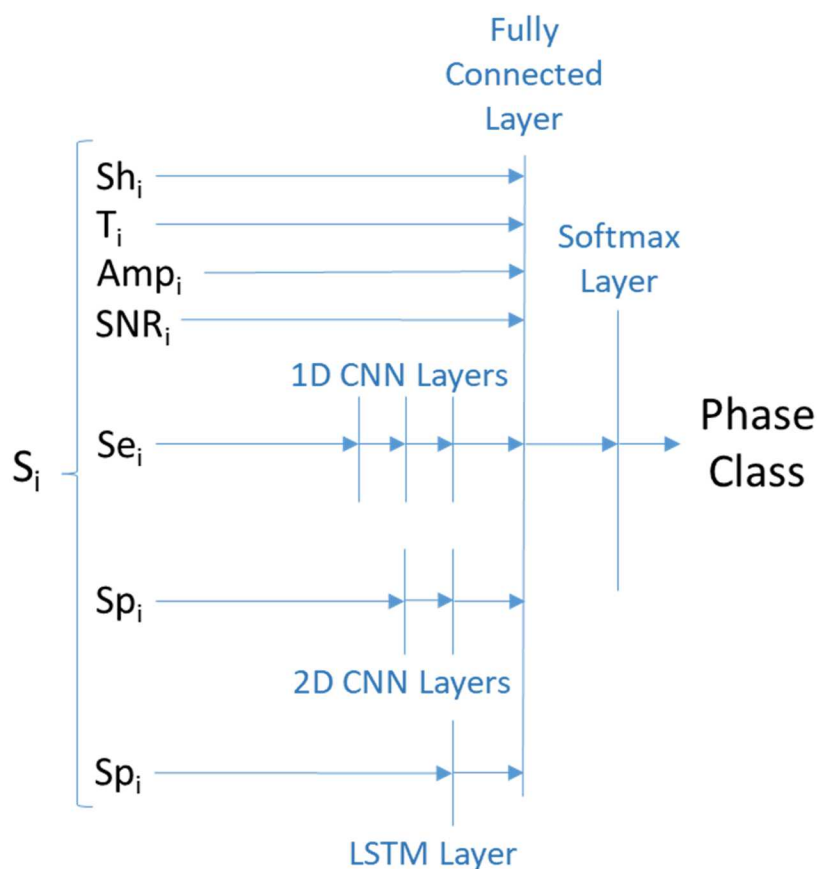*Figure 4. Merged deep neural network for Phase ID. Input to the network is on the left and output from the network is on the right.*

## RESULTS:

Results of our merged DNN solution to Phase ID were compared against 2 other approaches
1. The horizontal slowness of a detection.
2. The Phase ID field (iphase) in the arrival table for each detection in the IMS that comes from an automated Phase ID system.

It took us much of the project to arrive at a system that produced acceptable results. We don't show quantitative results except for our final system, but below is a list of challenges and break-throughs of the R&D process, all of which involves data issues.

1. The first attempt at training a phase classifier used from just the single 3-component element of the Makanchi array. The hope was that using all three components would help distinguish between phases. This may be the case, but we used just one filter band for each component (0.5 Hz – 6 Hz) and the SNR for much of the data was too low to see any identifiable signal in waveforms and spectrograms. At this point in the project, we also had a bug in our code of extracting waveforms 5 seconds before the arrival time through 10 seconds after. Therefore, our negative result is not conclusive.

2. To increase the SNR of waveforms and subsequent spectrograms, we performed beam-forming of all waveforms leveraging the location of events relative to MKAR. Although in real-life, the location of an event is not known when a detection arrives at a station, it is common practice to beam-form in multiple directions from the station, choosing the waveform with the highest SNR for further processing, in this case Phase ID. This process was a significant step in providing data with discernible structure in waveforms and spectrograms.

3. Finally, given that the type of phase is unknown for a detection nor even the distance from the station to the event, the use of multiple filter bands was used as follows.
   - 0.5 Hz – 1.5 Hz
   - 1.0 Hz – 2.0 Hz
   - 1.5 Hz – 3.0 Hz
   - 2.0 Hz – 4.0 Hz
   - 4.0 Hz – 8.0 Hz

   The combination of beamed waveforms and multiple filter bands provided spectrograms with discriminative structure for our final Phase ID results.

Results of two phase classification experiments are given below, one a simple demonstrative experiment and the other an experiment suggestive of high impact on improved event detection and location with PEDAL. All experiments used a classifier with the following elements.

- A spectrogram fed into a single bidirectional LSTM layer with 40 nodes and Scaled Exponential Linear Unit (SELU) activation functions. The spectrogram data input to the classifier is from beam-formed waveforms with 5 filter channels.
- The output of the LSTM network merged with four detection measurements (slowness, amplitude, SNR, and time since previous phase) in a fully connected layer of 20 nodes with SELU activation functions.
- The output of the FC layer fed into another FC layer with 40 nodes with SELU activation functions.
- A softmax layer with 2 nodes as the final output.

Figure 5 shows plots of seismograms for three different phases next to their corresponding spectrograms. Note that the arrival time of a phase should be 1/3 from the left of each plot, but

some phases have more gradual onsets or lower SNR, making it harder to discern or the arrival time may always be accurate. Another issue in need of further exploration is the optimal amount of time to use before and after the arrival time of each phase.

# Seismograms          Spectrograms



*Figure 5. Seismograms and spectrograms for P, Lg, and Sn phases.*

Figure 6 shows tables of binary classification results on the preliminary Phase ID experiment. Using the horizontal slowness value for each detection does very well in classifying between a class of Pn, Pg, and P phases and a class of Lg, Sn, and S phases. The merged DNN, however, does a better job by reducing misclassifications of Lg,Sn, and S.

Sandia National Laboratories                    U.S. DEPARTMENT OF **ENERGY**

## Merged Network

- **Accuracy on held-out Test Set**
  - 99.7% Overall
  - 99.3% Class Average

|  |  | Network Predictions | | |
|---|---|---|---|---|
|  |  | **Pn,Pg,P** | **Lg,Sn,S** | **Accuracy** |
| **Ground Truth** | Pn,Pg,P | 1121 | 2 | 99.8% |
|  | Lg,Sn,S | 2 | 172 | 98.9% |

## Simple Slowness (sh) Test

Pn,Pg,P: sh ≤ 17.3

Lg,Sn,S: sh > 17.3

- **Accuracy on held-out Test Set**
  - 98.1% Overall
  - 93.8% Class Average

|  |  | Network Predictions | | |
|---|---|---|---|---|
|  |  | **Pn,Pg,P** | **Lg,Sn,S** | **Accuracy** |
| **Ground Truth** | Pn,Pg,P | 1120 | 3 | 99.7% |
|  | Lg,Sn,S | 21 | 153 | 87.9% |

*Figure 6. Simple binary phase classification (Pn, Pg, and P vs. Lg, Sn, and S) results. The top table shows results from a merged DNN. The bottom table shows results using a horizontal slowness threshold.*

Figure 7 shows tables of binary classification results on the First-P Phase ID experiment. The merged DNN offers superior performance than the other two approaches. The accuracy numbers suggest that it could have a significant impact in reducing missed and false event detection in a seismic signal processing pipeline.

**Merged Network**

- Accuracy on held-out Test Set
  - 97.0% Overall
  - 95.6% Class Average

| Ground Truth | | Network Predictions | | |
|---|---|---|---|---|
| | | Not First-P | First-P | Accuracy |
| | Not First-P | 440 | 27 | 94.2% |
| | First-P | 29 | 1393 | 98.0% |

**Simple Slowness (sh) Test**

First-P: sh ≤ 15

- Accuracy on held-out Test Set
  - 83.7% Overall
  - 58.1% Class Average

| Ground Truth | | Network Predictions | | |
|---|---|---|---|---|
| | | Not First-P | First-P | Accuracy |
| | Not First-P | 189 | 278 | 40.5% |
| | First-P | 30 | 1392 | 97.9% |

**Automated Phase ID (iphase)**

- Accuracy on held-out Test Set
  - 65.7% Overall
  - 51.9% Class Average

| Ground Truth | | Network Predictions | | |
|---|---|---|---|---|
| | | Not First-P | First-P | Accuracy |
| | Not First-P | 114 | 353 | 24.4% |
| | First-P | 294 | 1128 | 79.3% |

*Figure 7. First-P vs. Not First-P phase classification results. The top table shows results from a merged DNN, the middle table shows results using a horizontal slowness threshold, and the bottom table shows results from the automated IMS phase labeler.*

## DISCUSSION:

Seismic Phase ID is a difficult problem because of the differing paths of identical phases, low SNR of some signal detections, and the additional information used by human analysts beyond individual waveform inspection and detection measurements. Much of the R&D was spent wrangling with data in such a way that DNNs could learn discriminative features for use in Phase ID. Whereas machine learning projects historically involved significant investment in feature engineering, much practical application of deep learning, particularly in uncharted waters, is largely about data engineering. For the foreseeable future and likely forever, machine learning will require investment and guidance from experts in the data related to the application.

The exploration of deep neural networks to capture spectrotemporal structure in phase waveforms in combination with a few signal detection measurements in a merged classifier led to a successful automated solution to Phase ID in a controlled, yet realistic and meaningful setting. The potential impact of higher quality event detections (fewer missed and false events) on operational seismic signal processing pipelines is significant. Since only one seismic monitoring station in a global network was used in this R&D, additional investigation is warranted into different kinds of monitoring stations used in different networks for different applications than global nuclear explosion monitoring. Human analysts often inspect waveforms

Sandia National Laboratories

U.S. DEPARTMENT OF ENERGY

from multiple stations to determine event detections and phase labels. We did not explore that here, but as long as there are examples with some ground truth of the steps that analysts take, merged DNNs may likely offer some assistance.

Typical machine learning projects are dominated with data issues and this project did not disappoint. In addition to the importance of improving SNR via beam-forming and the use of multiple filter bands, seismic Phase ID, particularly with the IMS data used, has severe class imbalances. Table 1 shows that P phases are about an order of magnitude more common than the next most common phase, Pn, and three orders of magnitude more common than some infrequent phases. This issue must be addressed while training a classifier so that it doesn't learn to predict everything as a P phase and get high accuracies numbers. We developed a custom minibatch sampler to present to the classifier during each iteration of training. The minibatch sampler is capable of collecting specific numbers of training sample by phase and/or by class. In this way, we get good class average accuracy results.

I presented initial results of the Phase ID work at the 2018 Machine Learning and Deep Learning Workshop and presented more details of final results at the University of New Mexico Meeting on Machine Learning Applications to Seismology. It was suggested by an attendee to use particle motion detection (e.g., the rectilinear detection measurement) for 3-component (non-array) stations. Exploration of other detection measurements is warranted in the future since our focus was on evaluating the value of the spectrotemporal structure in waveforms for Phase ID.

Although the First-P phase classifier performed admirably, discriminating between individual phases (see Table 1) will likely require contextual information, such as treating a collection of phases as a sequence and leveraging a geophysics model of the ordering of phases in time. Table 3 lists 22 events and the associated list of phases (sometimes just one) for each event detected by MKAR. Each sequence is like a spoken word in speech recognition and utilizing a sequence learning approach to individual Phase ID will likely bear fruit.

*Table 3. Phase sequences for multiple events.*

| Origin ID | Phase List |
|-----------|------------|
| 15148126 | P |
| 15148127 | P |
| 15148128 | P |
| 15148166 | PKP |
| 15148216 | P,pP,ScP |
| 15148226 | P |
| 15148268 | Pn |
| 15148273 | P,S,PKKPbc |
| 15148278 | PKP,SKPbc |
| 15148281 | Pn,Sn,Lg |

| | |
|---|---|
| 15148323 | P,PcP |
| 15148325 | P |
| 15148366 | P |
| 15149933 | Pn |
| 15149965 | P |
| 15149971 | Pn,Sn,Lg |
| 15150034 | P |
| 15150038 | P |
| 15150039 | P,PcP |
| 15150043 | P |
| 15150070 | P,PcP,ScP |
| 15150076 | PKP |

Figure 8 illustrates how the CTC loss function can perform supervised learning on sequence data without requiring alignmfent between input data and labels or segmentation of signal elements.



Figure 8. The use of Connectionist Temporal Classification for speech recognition. On the bottom is a sequence of spectrograms in time. On the top are multiple examples of the same words spoken in different time durations and with different delays.

Figure 9 shows how CTC might be used for Phase ID on a sequence of detections or a streaming waveform data.
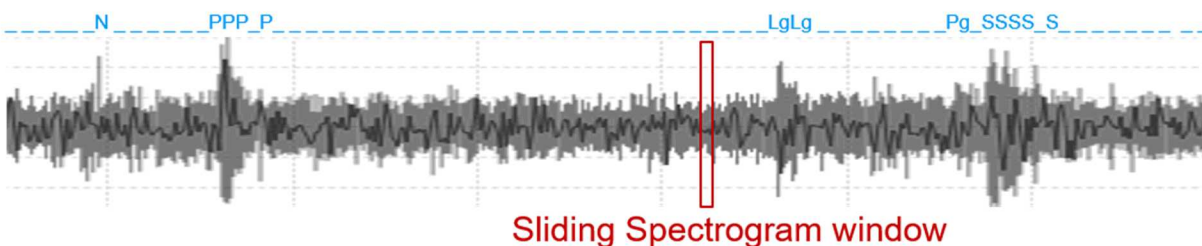
*Figure 9. Streaming Phase ID.*

## ANTICIPATED OUTCOMES AND IMPACTS:

The results from the exploratory R&D clearly demonstrate superior performance over standard automated methods of Phase ID on controlled experiments. This suggests positive impacts on the following applications.

- PEDAL – In addition to reducing the number of false and missed seismic events by filtering non-First-P phases from the event detection and location step of PEDAL, our automated Phase ID approach will likely help in the signal association step of PEDAL as well. Since this work used just one IMS station, an open question is how best to train an entire network of monitoring stations. Is a custom Phase ID classifier necessary for each station or class of stations (e.g., all array stations) or can a single classifier can be trained for an entire class of stations?
- Dynamic Networks – The ability to distinguish between noise and non-noise detections would have a significant impact on event detection with high-density local and regional seismic monitoring networks. We did not conduct and experiments to demonstrate this capability, but our results are highly suggestive of success.
- Hazard Prediction – Identifying the phases from large earthquakes as soon as possible is an important part of an early warning of potential hazards.

Although the above list relates to applications of our work, much research remains to develop fully an automated Phase ID capability.

**Data**
- Our Phase ID approach is data driven, where the more examples that are used for training, the better the system is expected to generalize. Because of the exploratory nature of this project, we used just 7 months of data out of 16 years of data from one IMS array station. Training with all 16 years of data is expected to improve results dramatically.
- In addition to the phase labels that exist in association tables, there is an opportunity to search for additional phase labels for events using "Probabilistic Labelling" since many phases are not added to analyst-reviewed bulletin if the phase was not of primary importance. Given a known event (E) and its magnitude ($m_b$), distance from a station ($\Delta$), and depth (d), if the probability of detecting a particular phase (ph) at station S, $P(ph, S \mid m_b, \Delta, d)$, is greater than a pre-specified threshold and that phase does not exist in the

Sandia National Laboratories

U.S. DEPARTMENT OF ENERGY

association table, then a highly-sensitive signal detector can be used to detect that phase in the waveform near its predicted arrival time. If the SNR of the detection is greater than a pre-specified threshold, then an additional label for that phase has been found. $P(ph, S \mid m_b, \Delta, d)$ values exist for all stations in the IMS.

- In addition to utilizing labeled data for training a phase classifier, semi-supervised learning algorithms exist that leverage unlabeled data to improve classification performance beyond that with labeled data alone. Unlabeled data is particularly plentiful in a streaming, continuous processing environment.

### Features

- The Phase ID approach demonstrated here utilized waveform data plus a few detection measurements. Many more measurements exist for each detected signal. The use of particle motion features (e.g., rectilinear motion) may offer value, particularly for 3-component stations and would be easy to add to a merged DNN.

### Sequence Learning

- Instead of classifying segmented waveforms, the ultimate Phase ID classifier treats phases in sequence with the use of a geophysics model of how different phases are expected to be ordered in time at a seismic monitoring station, similar to the way speech recognition treats phonemes in sequence with the use of a language model of how phonemes are ordered in time in words at a listener. Sequence learning algorithms are more difficult to train and we were unable to pursue this avenue in this exploratory work. Two approaches can be investigated.
  1. Use a continuous stream of waveform as input, identifying phases as they arrive. This would match the standard speech recognition approach.
  2. Use a sequence of detections as input, more easily leveraging the current work for each element of the sequence.

### Algorithms

- Advances in deep learning occur frequently. In this work, we chose well-established algorithms, but newer algorithmic approaches may offer superior performance.

### Cost-Sensitive Classification

- Cost-sensitive training can minimize the misclassifications that a user cares about most. For example, in a First-P Classifier for PEDAL, it is important to not miss First-P detections, which may result in a missed event detection, even if the system wrongly classifies more non-First-P detections.

## CONCLUSION:

Seismic Phase ID is known as an important element of event detection, location, and origin time determination as well as early warning of hazardous earthquakes. It is also known to be a difficult task to automate and involves significant human resources. Seismic applications of machine learning, particularly deep learning, abound since 1) much data (waveforms and automated measurement of signals) has been collected and for which some ground truth labeling exists and 2) although noisy at times, there is structure in the data corresponding to conclusions

important to the seismology and geoscience communities. Specifically, there are wonderful opportunities for innovate, modern solutions to automated Phase ID.

In this exploratory R&D, it was demonstrated that deep neural networks perform well on a classification problem impactful to the seismology community. It uses multiple types of data (spectrograms from beamed waveforms and detections measurements) from the IMS station MKAR to identify whether a detected signal is a First-P phase or not. Filtering all detections that are not First-P phases can lead to fewer missed and false events in a seismic signal processing pipeline. Identifying individual phases instead of important classes of phases will likely require context of each phase in time utilizing a geophysics model or how phases propagate. Sequence learning neural network methods used for speech recognition offer a potential solution. Addressing the needs of the seismological community with modern deep learning techniques on available partially-labeled seismic data will continue to be a rich domain for impactful R&D.

## BIBLIOGRAPHY:

1. Amodei, D., et al (2015). "Deep Speech 2: End-to-End Speech Recognition in English and Mandarin," CoRR abs/1512.02595. http://arxiv.org/abs/1512.02595.
2. Beaufays, F. (2015). "The Neural Networks behind Google Voice transcription," http://googleresearch.blogspot.com/2015/08/the-neural-networks-behind-google-voice.html