

Characterizing High Rate GNSS Velocity Noise for Synthesizing a GNSS Strong Motion Learning Catalog

Timothy Dittmann *, Jade Morton ¹, Brendan Crowell ³, Diego Melgar ⁴, Jensen DeGrande ³, David Mencin ²

¹Ann and H. J. Smead Aerospace Engineering Sciences Department, University of Colorado Boulder, Boulder, CO, USA, ²EarthScope Consortium, USA, ³Department of Earth and Space Sciences, University of Washington, Seattle, Washington, USA, ⁴Department of Earth Sciences, University of Oregon, Eugene, U.S.A

Author contributions: *Conceptualization:* All. *Methodology:* DiegoM, TD. *Software:* TD, BC, DiegoM. *Validation:* TD, JD. *Formal Analysis:* TD, JM. *Writing - original draft:* TD, JM. *Writing - Review & Editing:* All authors. *Funding acquisition:* JM, DavidM.

Abstract Data-driven approaches to identify geophysical signals have proven beneficial in high dimensional environments where model-driven methods fall short. GNSS offers a source of unsaturated ground motion observations that are the data currency of ground motion forecasting and rapid seismic hazard assessment and alerting. However, these GNSS-sourced signals are superposed onto hardware-, location- and time-dependent noise signatures influenced by the Earth's atmosphere, low-cost or spaceborne oscillators, and complex radio frequency environments. Eschewing heuristic or physics based models for a data-driven approach in this context is a step forward in autonomous signal discrimination. However, the performance of a data-driven approach depends upon substantial representative samples with accurate classifications, and more complex algorithm architectures for deeper scientific insights compound this need. The existing catalogs of high-rate ($\geq 1\text{Hz}$) GNSS ground motions are relatively limited. In this work, we model and evaluate the probabilistic noise of GNSS velocity measurements over a hemispheric network. We generate stochastic noise time series to augment transferred low-noise strong motion signals from within 70 kilometers of strong events ($\geq M_W 5.0$) from an existing inertial catalog. We leverage known signal and noise information to assess feature extraction strategies and quantify augmentation benefits. We find a classifier model trained on this expanded pseudo-synthetic catalog improves generalization compared to a model trained solely on a real-GNSS velocity catalog, and offers a framework for future enhanced data driven approaches.

Non-technical summary Global Navigation Satellite System (GNSS) signals are a source of valuable earthquake ground motion data that is traditionally sourced from inertial-based instruments. Inertial-based instruments include a class of sensors that use Newton's first law to directly measure ground velocity or acceleration. Routine noise of GNSS is more complex than the inertial-based instruments, which in turn has limited the scope of adoption of GNSS in earthquake monitoring. Machine learning applied to the scientific domain has shown that it can separate signal from noise and offer deeper scientific insights, but our existing datasets are relatively limited. Implementing an effective machine learning model for any scientific objective depends on having a sufficiently large, accurately labeled dataset for training and validating the model. We present an expanded "pseudo-synthetic" catalog comprised of transferred real-world signals added to synthetic GNSS velocity noise generated from real world noise analysis. We demonstrate how training a model on our expanded synthetic dataset outperforms training on limited real data and can support more sophisticated learning objectives offering deeper understanding.

1 Introduction

Distributed observations of coseismic ground motions are the backbone of accurate ground motion models, finite fault modeling, and early warning. If available in real-time, GNSS-derived high rate time differenced carrier phase (TDCP) velocities (GRAAS and SOLOVIEV, 2004) applied to seismology (Colosimo et al., 2011) are an additional source of these intrinsic measurements (Parameswaran et al., 2023) that are traditionally sourced from dedicated inertial sensor networks. If available in near-real time or post processing, GNSS

velocities can contribute to catalogs of ground motion measurements used for empirical regional and local ground motion models (Crowell et al., 2023). GNSS spatially complements or substitutes existing inertial ground motion observations (Crowell, 2021), especially valuable in sparse networks (Grapenthin et al., 2017). Furthermore, GNSS expands the dynamic range of inertial measurements, and contributes to magnitude estimation (Murray et al., 2023) when inertial sensors saturate (Melgar et al., 2013) during the largest, most destructive events.

However, ambient GNSS velocity noise remains well above the noise floor of inertial sensors, largely due to

Production Editor:
Gareth Funning
Handling Editor:
Mathilde Radiguet
Copy & Layout Editor:
Hannah F. Mark

Received:
May 9, 2023
Accepted:
September 11, 2023
Published:
October 5, 2023

*Corresponding author: stdi2687@colorado.edu

sources of uncertainty related to ranging of space-based weak radio frequency signals. Analysis of high rate positioning noise (Genrich and Bock, 2006), carrier phase noise (Wang et al., 2021), and TDCP velocities (Shu et al., 2018; Crowell et al., 2023) has shed valuable insight into the factors that influence the ambient noise floor of these GNSS velocities. To date, the GNSS velocity noise frequency spectrum has not been evaluated across sufficiently large temporal and spatial scales to statistically report on the ambient noise across a network. Ambient noise characterization methods developed in the seismic community offer a statistical approach to represent ambient noise frequency content for sensor network monitoring and calibration. The probabilistic spectrum of GNSS velocity noise illuminates the limit of seismic signal detection in GNSS.

Improved classification of seismic signals within GNSS noise will expand the range in which GNSS contributes material ground motion observations with minimal false alerting for denser in situ observations and early warning integrity. Methods for addressing this signal to noise (SNR) challenge exist: variations on a short term average over long term average (STA/LTA) detection adopted from inertial seismic sensors resolve static offsets (e.g. Allen and Ziv, 2011; Colombelli et al., 2013) but filter valuable dynamics encoded in the waveforms; threshold based detection methods (e.g. Crowell et al., 2009; Hodgkinson et al., 2020; Dittmann et al., 2022a) capture dynamics but struggle to balance sensitivity with false alerting, and must mitigate false alerts with external dependencies such as spatially correlating or temporally windowing from seismic triggers. Machine learning (ML) models combine a range of feature inputs to improve the decision confidence in separating seismic signal from noise (e.g. Meier et al., 2019; Dittmann et al., 2022b) in stand-alone mode. However, the generalization performance of any such classifier or deeper ML model will ultimately be limited by the model selection and optimization, the extent of the labeled catalog for training, and the quality of the labels.

Previous GNSS seismic catalogs illustrate how limited the observed long-tail, larger magnitude GNSS seismic events datasets are (Ruhl et al., 2018). For example, the EarthScope/UNAVCO continuous geodetic archive began archiving lower sampling rate GNSS observations in 1993 and 5Hz high rate data retrieval in 2006. Decreased hardware costs coupled with commercial and scientific demand only relatively recently allowed for global high-rate network expansion. Additional geodetic networks (e.g. INGV: Italy, GEONET: NZ) complement EarthScope's high rate catalog worthy of inclusion on the order of doubling, not the order(s) of magnitude needed for deeper learning to answer more sophisticated questions. One solution to this small data challenge is synthesizing waveforms using kinematic finite fault ruptures and Green's functions ("FakeQuakes," Melgar et al., 2016; Williamson et al., 2020). This model-driven approach is invaluable for the largest, most destructive events, where a data-driven strategy for these infrequent events is inherently insufficient. However, this method is not yet practical for generalizing across global rupture scenarios and great

care must be taken to not bias results with *unknown unknowns* of fault models and ground motion propagation of future events. This is an area of active research.

An intermediate real-world-data driven alternative is to transfer samples from a separate source of our signals of interest (Hoffmann et al., 2019). Inertial sensors have existed at more locations for far longer than the first positioning satellite was launched. Event catalogs of zero-baseline inertial measurements offer low-noise ground motion velocities to be transferred as our truth waveforms of accurately labeled samples. The GNSS noise probabilistic power spectral density (PPSD) characterization offers the necessary information to superpose stochastic noise for training over a range of noise conditions. The final component to improved generalization are the learning training decisions, including model selection and feature engineering. With appropriately applied domain knowledge to increasingly larger data volumes, the revolution of transferable classification and regression model algorithm development is readily adaptable to earth science questions (Bergen et al., 2019; Kong et al., 2018).

To improve our understanding of GNSS velocity sensitivity relative to ambient noise, expand the quantity of available labeled training data, and improve detection classification performance in a highly variable noise environment, we characterized the GNSS velocity noise frequency spectrum from which we augmented transferred inertial velocity waveforms observed over 80 years with synthetic GNSS velocity time series. This manuscript presents a framework for expanding the available, accurately labeled GNSS velocity waveforms and evaluates the improved signal detection gained from learning on such a catalog. Finally, we present the expanded catalog to support evolving, deeper learned models to train on.

2 Materials and Methods

2.1 Lightweight GNSS Velocity Processing

A GNSS receiver generates precise relative phase estimations by tracking the signal carrier wave using a phase lock loop. To achieve absolute positioning using carrier phase measurements, a suite of measurement error source models must be estimated to account for thermal noise, satellite and receiver oscillators, multipath reflections, atmospheric and ionospheric effects from a 20,000 kilometer signal propagation path, and unknown carrier cycle integer offsets (Teunissen, 2020). These correction models incur costs, computationally, potentially monetarily, and in performance for resolving carrier phase ambiguities to estimate absolute position. In past and current implementations of using geodetic measurements for capturing earthquakes, absolute positions are differenced from an a priori position to extract relative topocentric motion, the signal of interest. TDCP or variometric processing (GRAAS and SOLOVIEV, 2004) differences these precise carrier phase measurements in consecutive epochs to remove temporally correlated error sources and consistent integer ambiguities. TDCP uses the precision of these mea-

surements to its advantage, by foregoing absolute positioning in exchange for precise relative velocity measurements while still benefiting from multi-signal observability across a visible satellite constellation. In this context, TDCP advantageously does not require ambiguity resolution convergence, lacks complex error models which in turn minimizes measurement noise, and reduces computational requirements. These factors combined with the simplicity of the algorithmic inputs makes it ideal for seismic ground deformation applications (Colosimo et al., 2011; Benedetti et al., 2014; Hohensinn and Geiger, 2018; Grapenthin et al., 2018; Parameswaran et al., 2023) at higher rates and potentially on the network edge.

We use the SNIVEL processing method (Crowell, 2021) for estimating 5Hz GPS TDCP. This method uses the narrow lane GPS-only L1/L2 phase combination, the Klobuchar ionospheric correction, the Niell tropospheric correction, and broadcast satellite ephemeris. Observations are weighted as a function of satellite elevation angle with a seven degree elevation mask. While development accommodating precise orbits (Shu et al., 2020), multi-GNSS, cycle slip detection/mitigation (Fratarcangeli et al., 2018), and higher order noise source mitigation is ongoing and warranted, the current method is capable of capturing ground motions of nearfield M4.9 and larger sources at teleseismic distances (Crowell, 2021; Dittmann et al., 2022b).

2.2 Observed High Rate GNSS Velocity Noise Model and Synthetic Noise

Understanding GNSS noise is imperative to applying GNSS observations to answer complex geophysical questions. Such investigations range from low frequency estimation of secular plate velocities (Williams et al., 2004) to higher frequency (>1Hz) signals, including structural monitoring (SHEN et al., 2019; Hohensinn et al., 2020), space weather (Yang et al., 2017), and deformation monitoring (Geng et al., 2018; Avallone et al., 2011). Previous studies show that GNSS position noise is a combination of white and colored or power-law noise (Langbein and Bock, 2004). Starting from lowest frequencies, the “dam profile” of exponentially decaying noise with increased frequency is inferred to be a result of correlated signal path and processing contributions including multipath, ephemerides, clocks, and atmospheric effects. GNSS highest frequency position noise is attributed to receiver thermal noise and often presented as a white spectrum (Genrich and Bock, 2006). Receiver thermal noise is parameterized as a function of incoming signal strength and carrier phase tracking filter design, including filter bandwidth and sample integration time. These baseband signal tracking loop design choices balance dynamic stress response with thermal noise mitigation (Yang et al., 2017), and are reflected in this highest frequency noise profile (Moschas and Stiros, 2013; Häberling et al., 2015). As an aside, for these reasons a calibrated high frequency instrument response, similar to what has become the defacto standard in digital inertial instruments, has been proposed (Ebinuma and Kato, 2012). We note this as worthy of fur-

ther investigation for future efforts integrating TDCP velocity noise into monitoring but have not yet observed an instrument bias with respect to capturing seismic strong motion signals in 5Hz velocities.

The EarthScope geodetic archive captures 5Hz data of stations recording concurrent with larger magnitude earthquakes. This includes at least 1 hour of “ambient” 5Hz data antecedent to the hour in which the event takes place. We process with SNIVEL all available 5Hz pre-event hour long windows for our ambient GNSS velocity dataset. This dataset consists of 1507 hours from 904 stations since 2007 distributed from the Caribbean to Alaska. We use this sample space to be representative of GNSS velocity distributions both spatially and temporally.

We evaluated the spectrum of GNSS TDCP noise over this sample set by adopting a seismic ambient noise characterization method of McNamara and Buland (2004) modified for GNSS displacements by Melgar et al. (2020). In this approach, further modified for 5Hz GNSS velocities, we calculated the power spectral density of 10 minute 5Hz single component velocity windows. We evaluated power spectral densities (PSD) at periods from 205s down to 0.4s in 512 bins. PSDs were smoothed in octave intervals and then stacked across 73 aligned frequency bins over all available PSD segments. The result is a probabilistic power spectral density (PPSD), or distribution of power spectral densities over the samples included. These PPSDs have been adopted for seismic network monitoring (Casey et al., 2018) and offer valuable insight for anticipated signal sensitivity. We combined horizontal topocentric components into a single PPSD and then estimated an independent vertical PPSD, given GNSS vertical noise is approximately 3-5 times larger.

We stored 19 distribution slices (every 5th percentile from 5% to 95%) of the real-world noise quantiles from which to generate synthetic stochastic noise time series (See the pre-event time window of Figure 2). We adopted the approach of Melgar et al. (2020) for GNSS position displacements, first proposed by Boore (1983) and further developed by Graves and Pitarka (2010). In this approach, we were able to maintain the frequency content of the noise at respective reference levels while randomizing the phase for generating unique time series. We accommodated amplitude loss in the domain transformations with linear scaling. For additional context of this strategy, Lin et al. (2021) demonstrated an ML application leveraging the Melgar et al. (2020) approach for generating displacement noise time series superposed on synthesized FakeQuake displacements to train a deep learning model estimating Chilean subduction zone moment magnitudes.

2.3 Strong Motion Observations and Augmentation

Our signals of interest are velocity waveforms from medium to larger earthquakes (>M5.0) which GNSS velocities are sensitive to (Dittmann et al., 2022a). The Next Generation Attenuation for Western United States 2.0 (NGAW2) project (Ancheta et al., 2014) is a database

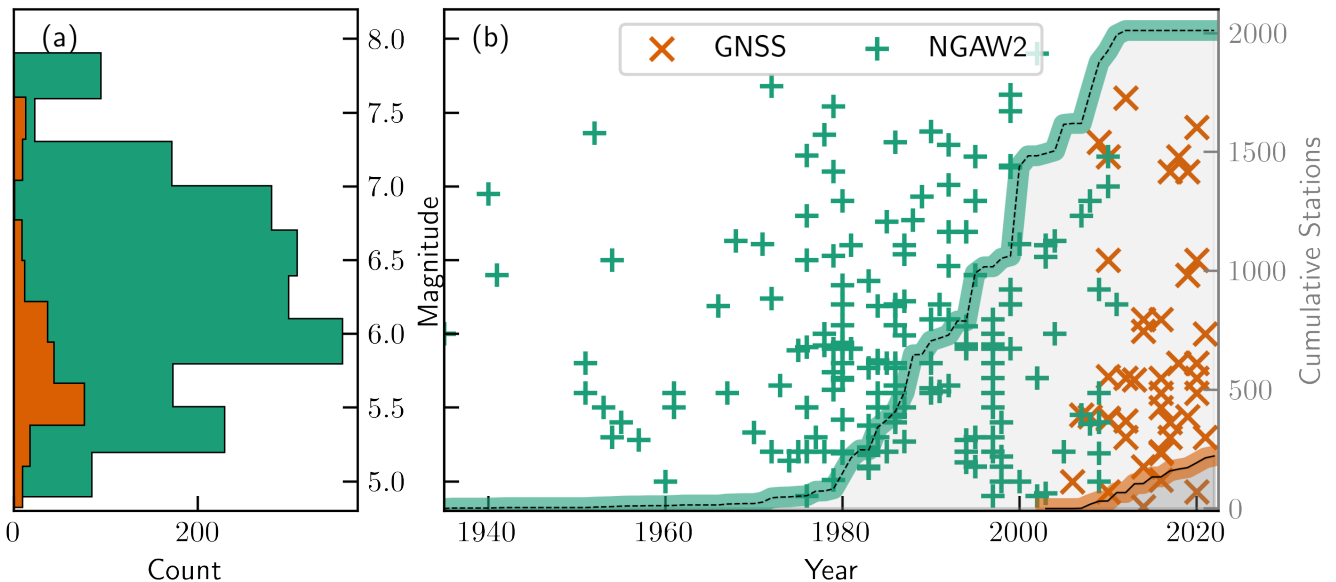


Figure 1 (a) A histogram comparing the EarthScope 5Hz GNSS catalog (“GNSS”) with the NGA West-2 database (“NGAW2”) for events observed by stations within 70 kilometers and sensitivity radii. The scatter plot in (b) shows the individual event magnitudes as a function of time, and the secondary axis line plot is the cumulative station count over time observing the events. In the cumulative line plot, the dashed line is the “NGAW2” and the solid line is “GNSS”.

of global strong motion measurements and response spectral ordinates from “shallow crustal earthquakes in active tectonic regimes” spanning over 75 years including 21,339 three component records from 599 events ranging M3.0 to M7.9. Global seismic networks contribute strong motion accelerograms or broadband velocity measurements that are processed by the NGAW2 project into acceleration, velocity, and displacement waveforms. The processing consists of an acausal Butterworth filter to reduce high- and low-frequency noise and an instrument response correction; further information regarding processing is given by [Ancheta et al. \(2014\)](#). The records were visually inspected for corner frequency determination, quality, and completeness, making the catalog an ideal source of low-noise larger ground motion measurements. A primary application of such a catalog is for ground motion prediction research to inform earthquake engineering. We use the processed velocity waveforms as our noise-free signal. It is worth noting that the seismic community has coalesced around several extensive labeled datasets to benchmark and facilitate rapid growth of deep learning models for a variety of applications ([Mousavi and Beroza, 2022](#)). We considered the several existing curated seismic ML catalogs ([Woollam et al., 2022](#)), but found these predominantly emphasized weaker signals. This is logical given the signal-to-noise challenges from inertial measurements looking to ML for use in seismology, but provides insufficient amplitudes for detection in synthesized GNSS strong motion observations.

We focus our effort on the portion of the database containing nearfield (≤ 70 km radius) observations of M5.0 to M7.9 within expected sensitivity radius of 1cm/s peak ground velocity given the scaling laws of [Fang et al. \(2020\)](#) for rapid hazard applications. Future work is extensible to the limits of detection above the noise

floor (>1000 km). We collected 2007 waveforms from 217 events (Figure 1). The processed velocity time series are offered at either 100 or 200 Hz sampling rate. We low pass filtered these waveforms with a filter corner frequency of 2.5Hz and then downsampled to 5Hz. We adopted a recursive short-term average over long-term average (STA/LTA) detection algorithm to label ground motion on each individual component. We found this is a sufficient automatic detector given its performance ([Withers et al., 1998](#)) in these relatively strong signals and factoring in the subsequent noise injected into our system. We used a 5 second short-term window and 10 second long-term window with a detection threshold ratio of 1.5. This metric was chosen through trial and error for its sensitivity for our larger strong motion signals of interest ([Trnkoczy, 2012](#)).

We exploited our “noise-free” signal waveforms and realistic stochastic noise generation by adopting data augmentation of transferred signals. Data augmentation is a form of regularization in which the size of a data catalog is artificially increased by creating augmented copies of our original waveforms ([Zhu et al., 2020](#)). Augmentation not only expands extents of a data catalog, valuable in relatively limited event datasets such as ours, but also improves generalization ([Bishop, 1995](#)). Successful augmentation trains a classifier to learn features or patterns in the presence of a larger range of authentic noise factors ([Iwana and Uchida, 2021](#)). In our application, we injected a synthetic noise time series derived from a single reference level of noise spectrum with unique phase values (Figure 2). We did this at seven noise reference levels on equivalent intervals from the 5th to 95th percentile to augment each strong motion waveform, a form of magnitude augmentation or jitter. We also buffered each augmented waveform with a random number of samples to misalign the samples in time

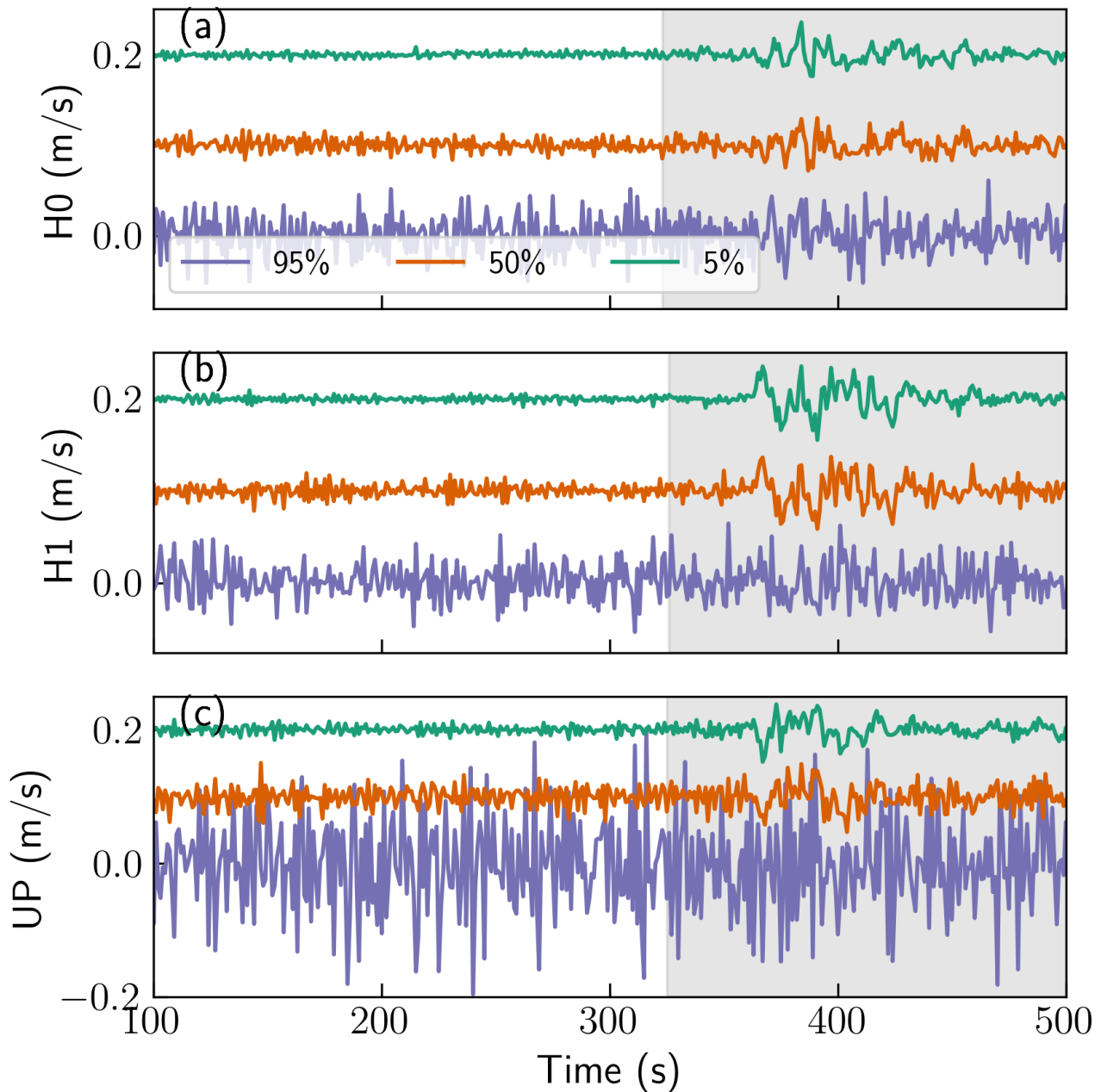


Figure 2 Example of three component waveforms from a single event NGAW2 waveform from Chi-Chi, Taiwan (2003, M6.2 50km radius) with three levels of synthetic noise added (5%, 50%, and 95% quantiles). (a) and (b) are the horizontal components, H0 and H1 respectively. (c) is the vertical waveform and noise component.

relative to each waveform replica. This resulted in seven different pseudo-synthetic observational waveforms for each station-event pair. This approach minimized overfitting in our models by training on a range of noise for a given signal at different offsets in each feature window, and expanded our catalog seven-fold from 2,007 strong motion waveforms to 14,049 pseudo synthetic GNSS velocity waveforms (Figure 1).

Additionally, we included the ambient catalog used in creation of the PPSDs to ensure the classifier is both trained on and tested against real-world GNSS velocity noise. This strategy was particularly important for potential disturbances not captured by the ambient synthetic noise generation process, such as the most infre-

quent events that might get statistically removed from the stochastic power spectrum but could result in detrimental false alerts if their signature is not learned. For example, the lowest frequency offsets from processing artifacts are infrequent enough to barely impact the probabilistic spectrum, but if not these are not included in training they could present as a synchronized event. We validated the performance of training a classifier on this synthetic catalog against the previously labeled EarthScope 5Hz GNSS velocities (see *Data and code availability*). For description of this dataset, please refer to [Dittmann et al. \(2022b\)](#). This curated catalog of GNSS velocity waveforms was processed identically as the noise catalog of this work; but one fundamental difference

is it is labeled through visual inspection instead of a known “truth” of our lowest noise inertial waveforms.

2.4 Model Selection, Feature Engineering and Training

First we validated the performance of a classifier trained on our strong motion waveforms relative to our previous GNSS velocity catalog approach. We used a random forest classifier (Breiman, 2001) for our detection model. Random forest is an ensemble method of decision trees. A decision tree is an algorithm that splits inputs along features to classify samples. A single decision tree can be biased by the initial features selected to seed the splitting; random forest overcomes this potential bias by running an ensemble of decision trees and having each cast a vote, where the majority eventually rules. We set up a binary classification that is demonstrated to have high accuracy and balance of sensitivity and false alerting in GNSS velocities. By keeping our model consistent with our previous work, we validated the newly formed catalog.

For validation comparison, we preserved our strategy from Dittmann et al. (2022b) of 30s overlapping windows. Future work will further optimize this sampling strategy with respect to sensitivity and real-time performance. From each window sample, we extracted a series of features to test their performance for our signal detection classification. In the time domain, we extract metrics akin to the traditional thresholding methods, including the four largest amplitudes, the median, and the median absolute deviation. In the frequency domain we included the entire PSD range over the 5Hz sampling of 30s windows, which includes periods from 1 second to 30 seconds. Variations on both of these time and frequency metrics were evaluated in our previous work, with the lower frequency (3s-15s period) horizontal PSD the most influential for the classifier model. However, while the overall performance over the entire catalog was a marked improvement from the current, variability in the false positive rate of the ambient dataset combined with missed detections of nearfield smaller magnitude events warrants further investigation.

Each sample consisted of one or a combination of these features for 30 second windows for all three components (Figure 3). STA/LTA labels were reduced to a single positive or negative outcome from 450 samples (150 samples per window x 3 components). Given our knowledge of signal relative to noise in this synthetic dataset, we also assigned a SNR metric for each sample, which was the peak single difference between signal power and noise power across all frequency bins. We employed a similar nested cross validation approach to our previous work for comparison and validation. Because the number of discrete events is still relatively small, we wished to minimize the potential bias from random validation and testing set selection.

In nested cross validation (Bishop and Nasrabadi, 2007), we ran 10 different testing scenarios, where each scenario keeps aside a different subset of one tenth of the events. Within each fold, we also ran an inner loop

of 5 fold cross validation across a grid search of hyperparameters. This technique further minimized overfitting hyperparameters by cross validating across a range of sample subsets. Our hyperparameters included the depth of nodes, or the number of decision splits, the number of estimators or decision trees, class weighting, a strategy that can assist with imbalanced datasets such as ours, and finally a SNR training threshold. This last hyperparameter was uniquely available to this pseudo-synthetic dataset; we generated the noise added to the signal, and so with this information we can accurately quantify the relative detectability. Using this as a hyperparameter allowed us to optimize training sets to include the largest extent of low signal-to-noise samples that benefit the model, while avoiding degrading model performance with undetectable low SNR.

In cross validation, we optimized the model on F1 scores, a balance of precision and recall. F1 is the harmonic mean of precision and recall. Precision is equal to the number of true positives (TP) over the sum of TP and false positives (FP), and recall is the number of TP over the sum of TP and false negatives (FN). Dittmann et al. (see 2022b).

3 Results and Discussion

3.1 Noise Characteristics

In TDCP velocity noise, we observe a V-shaped noise spectral profile in the PPSD (Figure 4). Periods longer than 6s follow a power law profile, likely reflecting correlated errors such as multipath and atmospheric effects not completely removed in the time differencing. This result is aligned with Melgar et al. (2020), which identified 1Hz PPP displacement noise as a red noise with a dam profile down to their Nyquist frequency. They infer that multipath and troposphere are the primary sources of the PPP “random walk” correlated noise signature (5s-200s period), and anticipate a spectral flattening to white noise around their maximum resolvable frequency (0.5 Hz) (Melgar et al., 2020). 1Hz PPP PPSD had a corner around 3 seconds, while in TDCP the lower frequency power law corner is at 6 seconds period. Another notable difference with TDCP processing reflected in this profile is the absence of absolute atmospheric models. In TDCP, the single slant path phase differences with first order corrections remove all but higher order gradients. Unfiltered time-differenced velocities will not accumulate error from potentially biased corrections models, a challenge of PPP. Shu et al. (2020) noted that inclusion of precise satellite clocks and orbits can significantly reduce longer period drifts existing in displacements derived from GNSS variometric velocities that otherwise must be detrended.

At approximately 4-6s period the noise spectrum inflects and begins increasing at a mirrored power law exponential to the lower frequencies. In TDCP at higher rates (>1Hz), Crowell et al. (2023) observes in multiple sample rates from a single receiver that TDCP velocities have increased noise in the time domain, roughly a factor of 7 of standard deviation from 1 Hz to 10 Hz veloci-

	Number of Station-Event Waveforms	Number of Samples	Labeling Strategy
GNSS Event Catalog (<70km) (Dittmann et al., 2022b)	247	5,187	visual inspection
Ambient Noise Training	1,507	88,893	assumed event-free
Ambient Noise Testing	1,507	85,806	assumed event-free
NGAW2	2,007	60,330	zero-noise truth labels
NGAW2 with Augmentation	14,049	422,309	zero-noise truth labels

Table 1 Extent and strategy of catalogs used in this research of noise and M5+ events within expected detectability and 70km radius.

ties. In the frequency domain these velocities present as a reverse power law of increasing noise as frequency increases, flattening at a corner around 0.2s period (5Hz). We observe a similar spectral shape in our PPSDs. Furthermore, Shu et al. (2018) processed up to 50Hz and identified a spectral “knee” around 3.5Hz; the highest frequencies observed in our study terminated at this “knee”. We infer this highest frequency (>1Hz) correlated noise to be predominantly influenced by receiver thermal noise, and likely receiver baseband design dependent (Moschas and Stiros, 2013). Crowell et al. (2023) also finds that the lowest noise power in the frequency domain exists in the 1-10s periods of the highest sample rate observations (20Hz in their study), notable given this intersects the spectral region of the seismic ground motion waveforms of interest. Given the spectrum at higher sampling rates, there is likely potential for improved screening of TDCP velocities for our signals of interest to reduce temporal aliasing (Hohensinn et al., 2020; Crowell et al., 2023).

A future PPSD product from continuous single station measurements would enable quantitative comparisons of the ambient noise levels from one station to another for monitoring and performance analysis. These noise levels, presented in a domain familiar scheme, are a meaningful proxy for the relative sensitivity to observe ground motions. Routine outliers can be observed and correlated to disturbances or events, a potentially valuable tool for network monitoring. In this study, without continuous 5Hz observations, it is not possible to assess time or spatially related variability outside the semi-arbitrary windows currently available.

3.2 Pseudo Synthetic Model Performance

We evaluated three different feature selection strategies by deploying three independent scenarios of random forest hyperparameter tuning and model fitting on identical training and testing splits. An advantage of our pseudo-synthetic approach is our knowledge at the individual waveform level about discrete true signals relative to artificial noise across our synthetic catalog. Our feature sets were time, frequency, and a combined time and frequency set “psd-t”. Overall, we found the highest

performance from the largest feature vector of all available features (Figure 5). We found the PSD-only performance similar to the “psd-t” combined feature vectors, which aligns with our feature importances from Dittmann et al. (2022b).

The overall F1 scores of Figure 5(a) indicate the optimal classifier will include both sets of information, but the PSD-only and time-only F1 scores suggest that the frequency domain information is most valuable for its stand alone performance relative to time only features. A benefit of our random forest model is readily extracted feature importance information (Figure 6). When our random forest model was presented with the time and frequency information, the trained model distributed feature importances across spatial components and features. The horizontal components (East/North) contributed more than the vertical, consistent with previous findings aligned with increased vertical GNSS noise relative to signals (Figure 4, Genrich and Bock, 2006). Contrary to the stand alone performance of Figure 5(a), discrete time domain features have considerably more importance than the frequency domain features. However, the sum of all frequency features in Figure 6(b) is greater than the cumulative time domain features for each respective component.

Within the frequency domain, the most valuable features are in the 2-5s period range. This shape is distinct from our previous classifier Dittmann et al. (2022b), where the most valuable features were the lowest frequency power spectra (6-30s).

From a model explainability perspective, we interpret that this importance distribution reflects the strength of the ensemble decision tree algorithm to distribute its decisions across all features with encoded information to optimize performance. An equivalent algorithm would be difficult to implement and generalize using traditional thresholds or filtering of this combined information. From a domain interpretability perspective, the relative value of signal amplitudes and signal frequency content is comparable after factoring in the distribution of frequency importances across significantly more bins. A model trained on these combined features gets the best-of-all-worlds benefits that traditional approaches (e.g. STA/LTA, threshold) lack. Additionally,

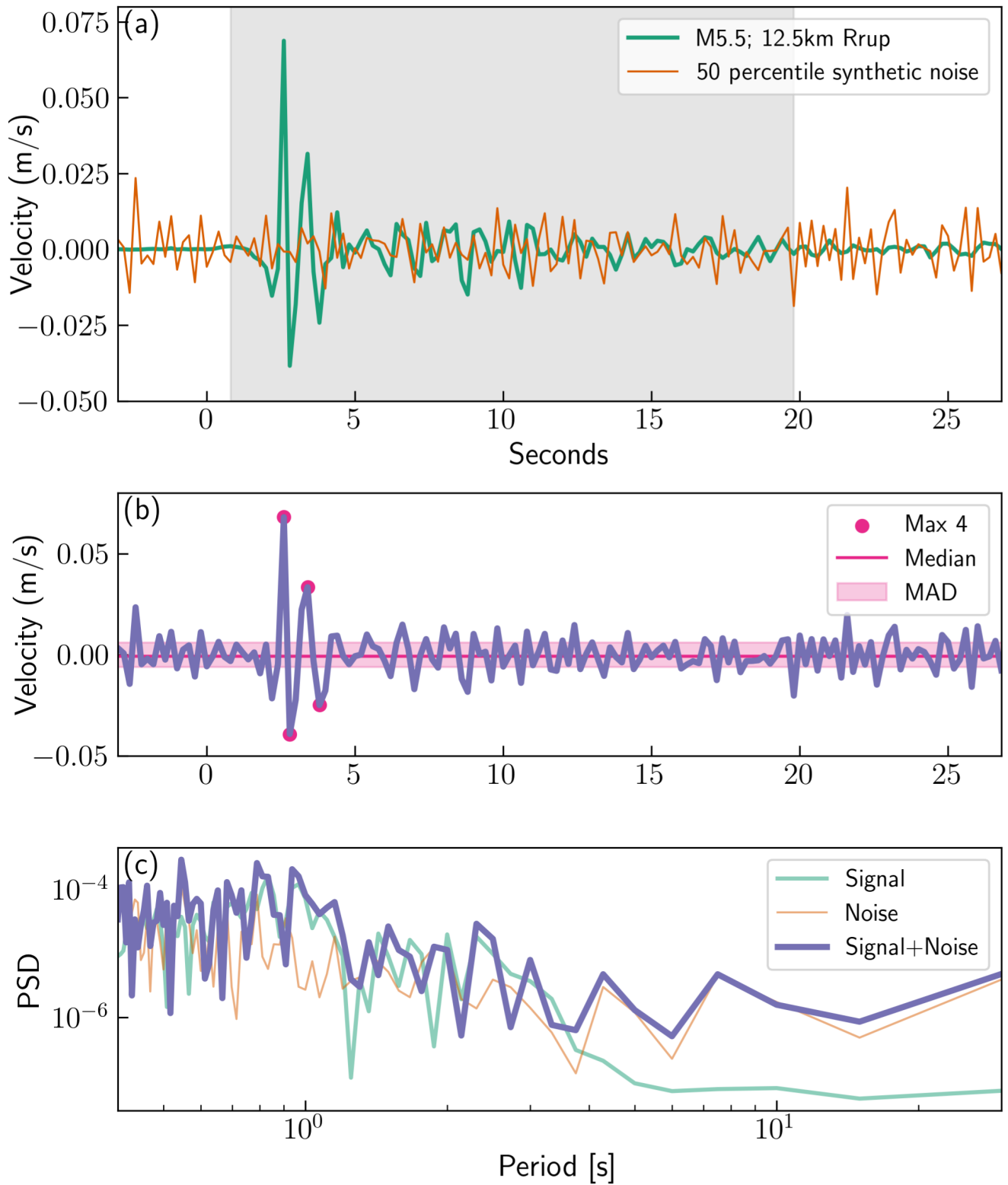


Figure 3 Demonstration of waveforms, noise and feature selection. The green timeseries in (a) is a downsampled NGAW2 waveform of a relatively weak signal for our application (a M5.5 at 12.5 km). The orange is a randomly generated noise time series using the 50th percentile noise spectrum. The gray shading is the region of detection triggered by the recursive STA/LTA. The sum of these time series (b) is then used as our observation. In the time domain (b) the features selected include the 4 largest amplitudes (solid magenta circles), the median, and the median absolute deviations, all indicated for this waveform in magenta. Finally, we also compute the power spectral density using a periodogram (in purple) and extract the power at each frequency bin. The original signal and noise periodograms are shown as well, for reference, though they are not included in the feature extraction.

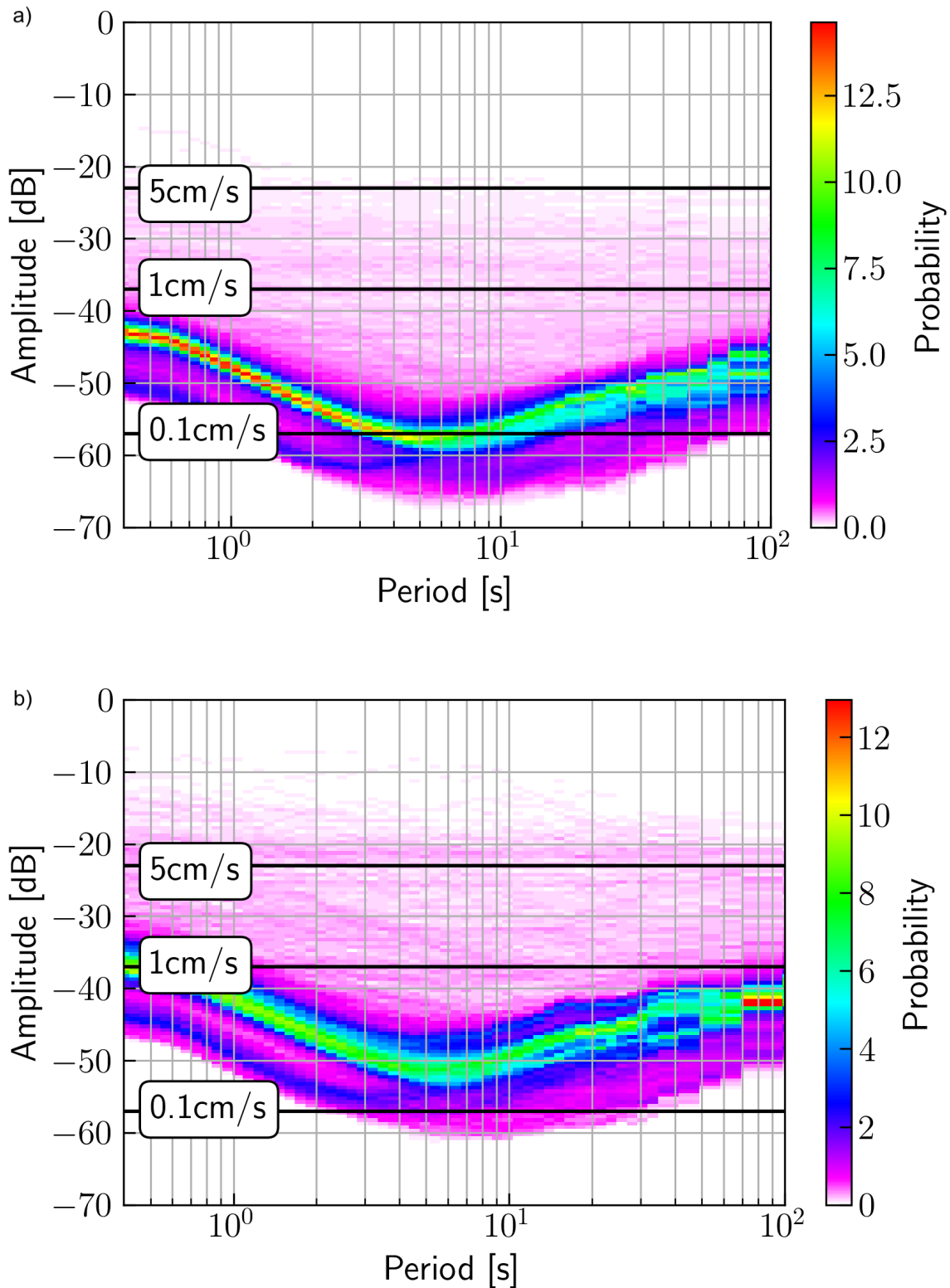


Figure 4 GNSS velocity PPSDs. Panel (a) is the combined horizontal components, panel (b) is the vertical component. Horizontal black lines are references for white noise timeseries of 3 respective standard deviations (5 cm/s, 1cm/s, 0.1cm/s).

the difference between the importances of this classifier and the previous classifier we infer is due to the nature of the labeling; these pseudo synthetic waveforms are labeled with low-noise “truth” models, so higher frequency, including more pulse-like signals, are more readily labeled. This is in contrast to the visual inspec-

tion, in which the human eye is inherently drawn to and presumably biased by longer period coherent signals. We will further evaluate in the validation section that training on augmented pseudo synthetic waveforms outperforms human-level classification performance.

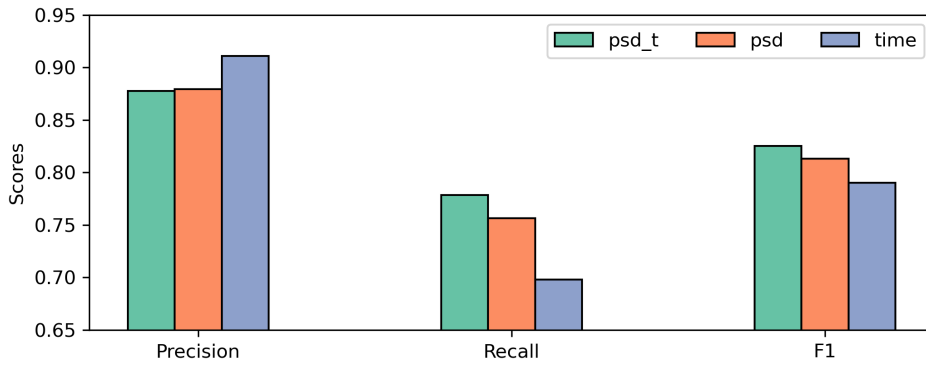


Figure 5 Testing feature extraction strategies across the NGAW2 synthetic dataset. Precision, recall, and F1 scores are presented as a function of feature extraction strategies across the entire catalog in 10 fold nested cross validation. “PSD” are the frequency domain features, “time” are the time domain features, and “psd-t” are the same time and frequency features concatenated into a single feature vector.

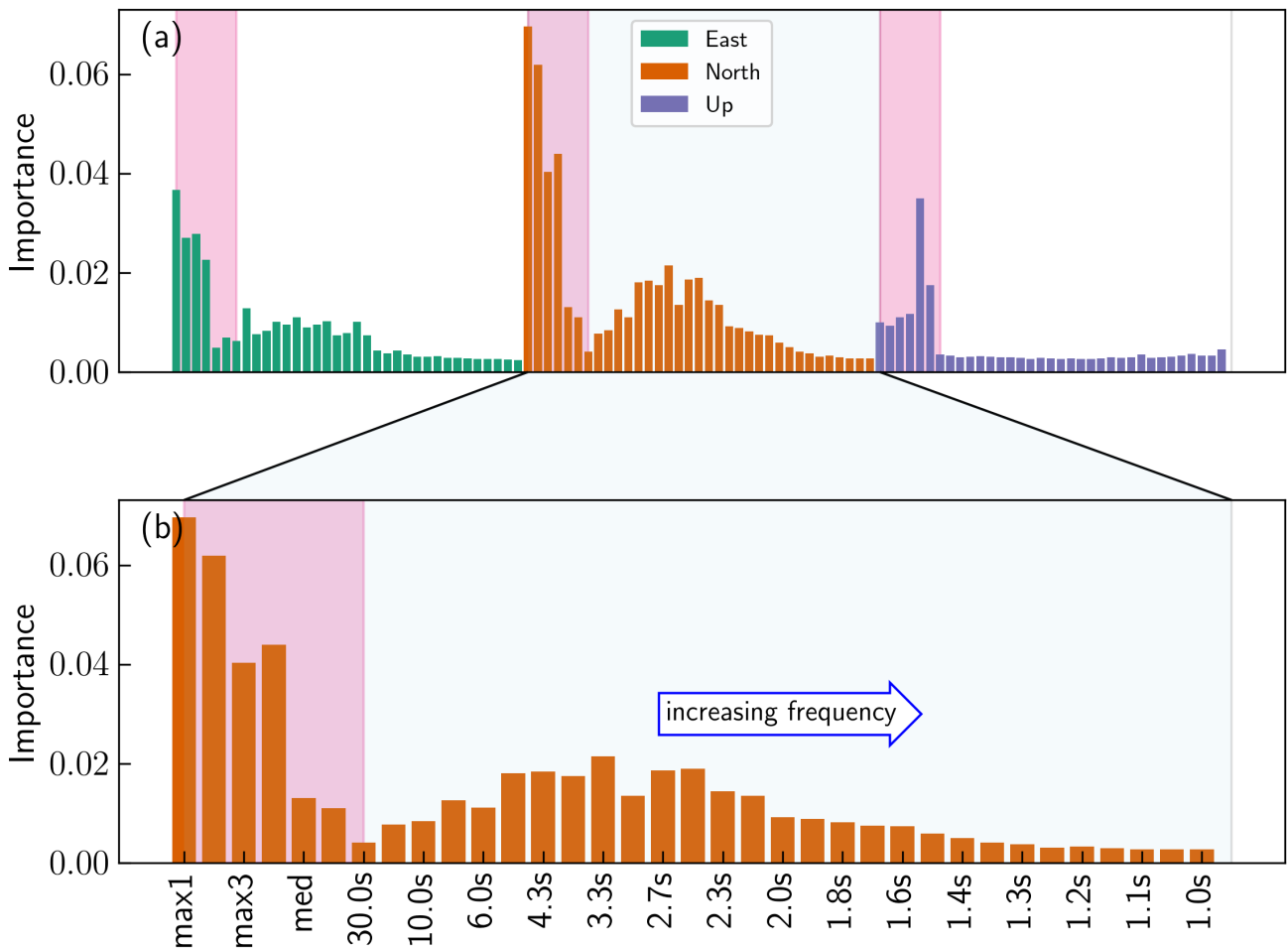


Figure 6 Feature importances for our random forest classification model cross validated and trained on the entire NGA West 2 synthetic GNSS dataset and the ambient noise dataset. Panel (a) shows the concatenated importances for all features across all components when a model is trained on all the features at once; the pink shading represents time domain features, the unshaded section are the frequency domain features. The second panel (b) is a close up of the North component features, with the same background shading schema. Every other feature is labeled for reference, max1 is largest amplitude, max2 is second largest, ..., 1.0s is the power in the 1s period bin of the PSD; for a single window example, see 3.

3.3 Quantifying Augmentation

Figure 1 and Table 1 make evident that transferred signals with data augmentation significantly expanded the

GNSS velocity catalog with respect to the number of unique waveforms. Additionally, data augmentation is an opportunity to expand sample feature space by lever-

aging our knowledge of the signals relative to the noise to train high quality labels with elevated noise environments (Zhu et al., 2020). We quantified the performance impact of augmentation by comparing models trained with and without augmentation. We ran identical complete nested cross validation testing scenarios using two different training tactics. In the first, we allowed the model to train on all 7 replicas of each waveform. In the second, we only provided the lowest noise waveform in training. Panels (a-b) of Figure 7 are from the first training scenario with augmentation. We tested on all replicas of the testing set waveforms, but for visualization purposes the left panel (a) is the performance of the 20th percentile median noise waveforms, and the right panel (b) is the performance of the 80th percentile high noise waveforms. The 20th or 80th percentiles are chosen to represent the “high” and “low” noise levels. SNR metrics were derived from the known noise time series and known signal periodograms. With data augmentation, we observed decreasing SNR for the same catalog while testing against increased noise levels (from panel a to b or c to d), with an overall true positive rate from 90% to 84%. When we compared the 20% noise levels with and without data augmentation (panels a, c), we notice a similar drop in performance without augmentation. Finally, when we looked at the highest noise samples without augmentation, we see a dramatic decrease in performance despite testing on the identical waveforms with the same SNR, from 90% to 75%.

3.4 Validation with Observed High Rate GNSS Velocity Event Waveforms

Finally, to validate our synthesis of GNSS velocity waveforms against real-world GNSS velocities, we reran a nested cross validation experiment with the entire real-world GNSS velocity catalog of Dittmann et al. (2022b) as a reference to compare the synthetically generated model. Similar in testing design to the previous comparison of data augmentation, we evaluated the performance of two classification models against the same semi-random testing subsets in the nested cross validation loops and reported on the mean performance. In this testing split scenario, one model was fit on the remaining ‘real’ data using hyperparameters extracted from k-fold cross validation for each training set, while the other model was fit on the entire synthetic GNSS velocities catalog. All other feature engineering strategies were held consistent and both models were evaluated against the same ‘real’ testing sets. The synthetic GNSS trained model yielded better performance metrics, including increased precision, recall, and F1 (Figure 8). This performance can best be explained by the extent of training sets: the synthetic model was trained on 14,049 waveforms, where the “real” model was trained on ~200, depending on the nested cross validation run testing slice. The added extent and density of information in the transferred and augmented training data improved model generalization for unseen events.

Additionally, we ran an ambient test where we take the best fit model from each dataset and applied it to a yet unseen ambient noise dataset (for dataset descrip-

tion, see Table 1). We found the GNSS velocity trained model had a nearly identical false positive rate, where false positive rate is one minus the true negative rate (Figure 8). This further validates that our noise training and augmentation strategy was effective in improving performance in difficult noise conditions, as our performance improvement in the event catalog did not come at the expense of ambient performance.

From these improved classification results we infer that transferred, augmented “synthetic” waveforms are not only a valid substitute for high-rate GNSS measurements to partially overcome modern, smaller GNSS seismic datasets, but may outperform human-level classification performance. A future deployed classifier will be trained on the combination of data catalogs to achieve the best generalization performance for yet-to-occur events. This real-world versus pseudo-synthetic comparison and validation result also suggests that evolved transfer learning across measurement domains, including exploration of fine-tuning of more mature seismic deep ML models with GNSS velocities, could further advance GNSS seismology challenges.

4 Conclusions

We find the ambient GNSS velocity noise distribution’s shape to be consistent with previous high-rate GNSS positioning noise analysis and spectral amplitudes, and find the noise distribution to be useful for signal sensitivity, synthetic noise generation, and future network monitoring. We find that frequency, time, and combined feature extraction strategies vary slightly under different SNR regions and that data augmentation boosts overall performance by training a model in higher noise settings.

Finally, we find that a model trained on these pseudo-synthetic waveforms, with the full suite of augmentation, outperforms the model trained on strictly GNSS velocity waveforms over the magnitudes (M_W 5.0-8.0) and hypocentral distances (≤ 70 km) tested in this analysis. Augmentation improves detection around the noise-signal boundaries. The immediate benefit is an improved classification model from an expanded catalog that can be retrained on the combined pseudo-synthetic and real catalog for unseen events. Such a classifier will be embedded in enhanced network operations and hazard monitoring for automated, stand-alone event detection. The subsequent benefit is an expanded training catalog (Dittmann et al., 2023) and framework that supports deeper learning models that are “data hungry” (Mousavi and Beroza, 2022). This includes expanding functional learning outputs, such as denoising, regression for magnitude inversion, and forecasting. With respect to future training of the largest events using this catalog, we identify possible limitations of this approach for specific experimental hypotheses due to the potential for introducing magnitude saturation of inertial instruments into our model training, a phenomena we are explicitly avoiding by using GNSS as a source. Similarly, more sophisticated source-dependent learning (e.g. forecasting) will need to consider the distribu-

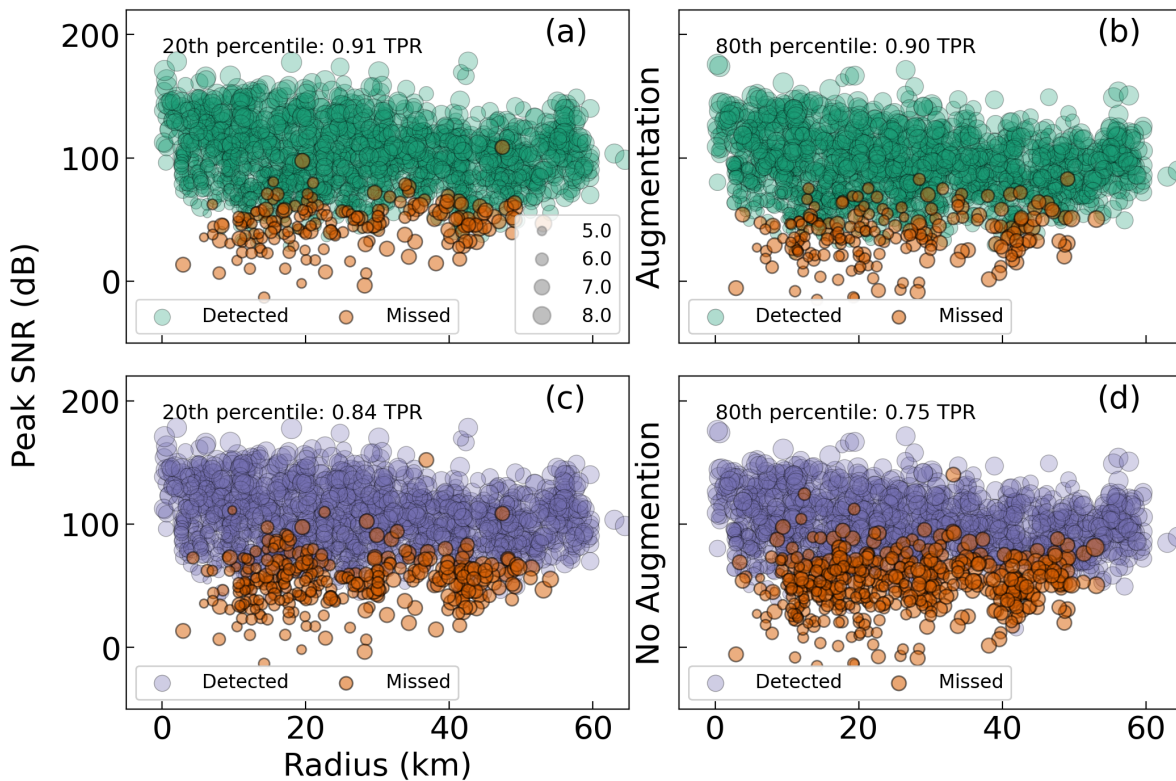


Figure 7 Comparing event detection with training on augmented noise samples across noise levels. Each panel includes the peak SNR of the waveform for each event as a function of radius from the event. This SNR metric is the peak of signal power to noise power for any frequency bin of periodograms calculated for all samples for all components for a given station-event waveform. The plot marker radius is determined by the event magnitude. The top panels (a-b) are testing the 20th noise model and 80th noise model of each station-event waveform using a classifier trained on all augmented samples. 20th and 80th are chosen to represent “low” and “high” noise. The markers are colored by a binary detected/not detected. The bottom panels (c-d) are testing the 20th/80th noise model waveforms with no data augmentation. This illustrates the value of augmentation for detection in noise, in addition to the approximate threshold of detection given our knowledge of signal and noise in this pseudo synthetic dataset. “TPR” - True Positive Rate.

tion of the NGAW2 source catalog used, specifically accounting for subduction events. Further investigations using this framework, perhaps paired with fully synthetic methods, is warranted. A loose ML integration of stand-alone inertial waveforms and this expanded GNSS-sourced waveforms enables fine-tuning (Yosinski et al., 2014) or transfer of existing inertial-based seismic detection ML models, such as Mousavi et al. (2020); Seydoux et al. (2020). Tighter amalgamation of stand-alone sensor sources benefiting from improved classification could include GNSS-sourced velocity waveforms directly in ground motion catalogs (Crowell et al., 2023) and operational monitoring systems. Such approaches would further blur distinctions between inertial and GNSS seismic signal sources, shifting from representations of different fields of earth sciences towards independent observational inputs with complementary dynamic ranges and respective noise models.

Data and code availability

The inertial seismic records are available from the Pacific Earthquake Engineering Research Center (PEER) Next Generation Attenuation for Western United States 2.0 (<https://ngawest2.berkeley.edu/>, Ancheta et al., 2014).

The 5Hz GNSS data used for TDCP processing in the study are available from the Geodetic Facility for the Advancement of Geoscience (GAGE) Global Navigation Satellite Systems (GNSS) archives as maintained by EarthScope Inc. (previously UNAVCO, Inc). The data are available in RINEX (v.2.11) format at <https://data.unavco.org/archive/gnss/highrate/5-Hz/rinex/>. SNIVEL code used for TDCP velocity processing is developed openly at <https://github.com/crowellbw/SNIVEL> (Accessed December 2021) (Crowell, 2021). SNIVEL 5Hz velocity timeseries used in this study are preserved at (Dittmann, 2022). Labeled 5Hz GNSS velocity samples and pseudo synthetic samples are preserved at (Dittmann et al., 2023).

Version 1.0.1 of the scikit-learn software used for ran-

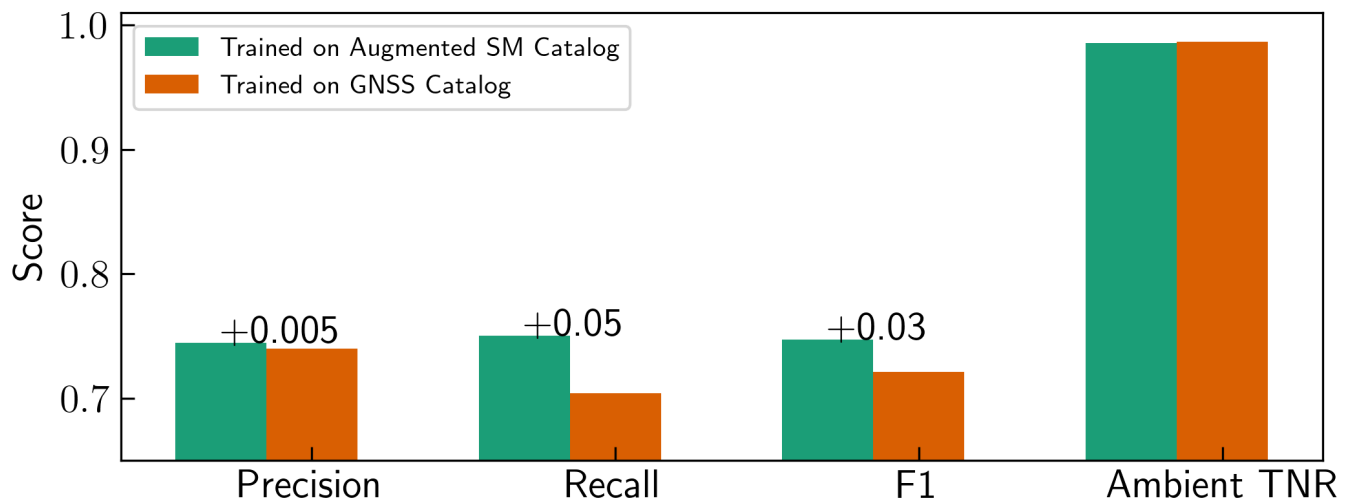


Figure 8 Testing performance of real GNSS velocity events as a function of training catalog used. These are the mean scores across the 10 testing folds of the nested cross validation. The purple results are from a model generated using cross validation of the remaining real gns dataset; the green results are from the bulk model fit to the entire NGAW2 synthetic dataset. Each uses the “psd-t” feature extraction method (combined time and frequency features). The ambient true negative rate (TNR) is estimated using a separate dataset of unseen ambient data. TNR is (true negatives) / (true negatives + false positives), or equivalent to one minus the false positive rate. The annotated text is the difference between the two approaches.

dom forest classification is preserved at (Grisel et al., 2021) and developed openly at <https://github.com/scikit-learn/scikit-learn> (Pedregosa et al., 2011). Figures were made with Matplotlib version 3.5.1 (Caswell et al., 2021), available under the Matplotlib license at <https://matplotlib.org/>.

Version 1.2.1 of the obspy software used for seismic data handling and PPSD generation is preserved at (Team, 2020) and developed openly at <https://github.com/obspy/obspy> (Krischer et al., 2015). Software used to generate psuedo synthetic waveforms and train, test and validate models is available at https://github.com/timdittmann/psuedosynth_gnss_velocities

Acknowledgements

We would like to express gratitude for the reviews by editor Mathilde Radiguet and two anonymous reviewers that improved this manuscript.

This material is based on services provided by the GAGE Facility, operated by EarthScope Consortium, with support from the National Science Foundation, the National Aeronautics and Space Administration, and the U.S. Geological Survey under NSF Cooperative Agreement EAR-1724794.

High-rate processing and ML for geoscience and hazards research is supported by NSF OAC-1835791.

Competing interests

The authors have no competing interests with this manuscript.

References

- Allen, R. M. and Ziv, A. Application of real-time GPS to earthquake early warning. *Geophysical Research Letters*, 38(16), Aug. 2011. doi: 10.1029/2011gl047947.
- Ancheta, T. D., Darragh, R. B., Stewart, J. P., Seyhan, E., Silva, W. J., Chiou, B. S.-J., Wooddell, K. E., Graves, R. W., Kottke, A. R., Boore, D. M., Kishida, T., and Donahue, J. L. NGA-West2 Database. *Earthquake Spectra*, 30(3):989–1005, Aug. 2014. doi: 10.1193/070913eqs197m.
- Avallone, A., Marzario, M., Cirella, A., Piatanesi, A., Rovelli, A., Alessandro, C. D., D’Anastasio, E., D’Agostino, N., Giuliani, R., and Mattone, M. Very high rate (10 Hz) GPS seismology for moderate-magnitude earthquakes: The case of the Mw 6.3 L’Aquila (central Italy) event. *Journal of Geophysical Research*, 116(B2), Feb. 2011. doi: 10.1029/2010jb007834.
- Benedetti, E., Branzanti, M., Biagi, L., Colosimo, G., Mazzoni, A., and Crespi, M. Global Navigation Satellite Systems Seismology for the 2012 Mw 6.1 Emilia Earthquake: Exploiting the VADASE Algorithm. *Seismological Research Letters*, 85(3):649–656, May 2014. doi: 10.1785/0220130094.
- Bergen, K. J., Johnson, P. A., de Hoop, M. V., and Beroza, G. C. Machine learning for data-driven discovery in solid Earth geoscience. *Science*, 363(6433), Mar. 2019. doi: 10.1126/science.aau0323.
- Bishop, C. M. Training with Noise is Equivalent to Tikhonov Regularization. *Neural Computation*, 7(1):108–116, Jan. 1995. doi: 10.1162/neco.1995.7.1.108.
- Bishop, C. M. and Nasrabadi, N. M. Pattern Recognition and Machine Learning. *J. Electronic Imaging*, 16:049901, 2007.
- Boore, D. M. Stochastic simulation of high-frequency ground motions based on seismological models of the radiated spectra. *Bulletin of the Seismological Society of America*, 1983. doi: <https://doi.org/10.1785/BSSA07306A1865>.
- Breiman, L. Random Forests. *Machine Learning*, 45(1):5–32, 2001. doi: 10.1023/a:1010933404324.
- Casey, R., Templeton, M. E., Sharer, G., Keyson, L., Weertman, B. R., and Ahern, T. Assuring the Quality of IRIS Data with MUSTANG.

- Seismological Research Letters*, 89(2A):630–639, Feb. 2018. doi: 10.1785/0220170191.
- Caswell, T. A., Droettboom, M., Lee, A., de Andrade, E. S., Hoffmann, T., Hunter, J., Klymak, J., Firing, E., Stansby, D., Varoquaux, N., Nielsen, J. H., Root, B., May, R., Elson, P., Seppänen, J. K., Dale, D., Lee, J.-J., McDougall, D., Straw, A., Hobson, P., Hannah, Gohlke, C., Vincent, A. F., Yu, T. S., Ma, E., Silvester, S., Moad, C., Kniazev, N., Ernest, E., and Ivanov, P. matplotlib/matplotlib: v3.5.1, 2021. <https://zenodo.org/record/5773480>. doi: 10.5281/ZENODO.5773480.
- Colombelli, S., Allen, R. M., and Zollo, A. Application of real-time GPS to earthquake early warning in subduction and strike-slip environments. *Journal of Geophysical Research: Solid Earth*, 118(7):3448–3461, July 2013. doi: 10.1002/jgrb.50242.
- Colosimo, G., Crespi, M., and Mazzoni, A. Real-time GPS seismology with a stand-alone receiver: A preliminary feasibility demonstration. *Journal of Geophysical Research: Solid Earth*, 116(B11), Nov. 2011. doi: 10.1029/2010jb007941.
- Crowell, B., DeGrande, J., Dittmann, T., and Ghent, J. Validation of Peak Ground Velocities Recorded on Very-high rate GNSS Against NGA-West2 Ground Motion Models. *Seismica*, 2(1), Mar. 2023. doi: 10.26443/seismica.v2i1.239.
- Crowell, B. W. Near-Field Strong Ground Motions from GPS-Derived Velocities for 2020 Intermountain Western United States Earthquakes. *Seismological Research Letters*, 92(2A): 840–848, Jan. 2021. doi: 10.1785/0220200325.
- Crowell, B. W., Bock, Y., and Squibb, M. B. Demonstration of Earthquake Early Warning Using Total Displacement Waveforms from Real-time GPS Networks. *Seismological Research Letters*, 80(5): 772–782, Sept. 2009. doi: 10.1785/gssrl.80.5.772.
- Dittmann, T. High Rate GNSS Velocities for Earthquake Strong Motion Signals, 2022. doi: 10.5281/ZENODO.6588601.
- Dittmann, T., Hodgkinson, K., Morton, J., Mencin, D., and Mattioli, G. S. Comparing Sensitivities of Geodetic Processing Methods for Rapid Earthquake Magnitude Estimation. *Seismological Research Letters*, 93(3):1497–1509, Jan. 2022a. doi: 10.1785/0220210265.
- Dittmann, T., Liu, Y., Morton, Y., and Mencin, D. Supervised Machine Learning of High Rate GNSS Velocities for Earthquake Strong Motion Signals. *Journal of Geophysical Research: Solid Earth*, 127(11), Oct. 2022b. doi: 10.1029/2022jb024854.
- Dittmann, T., Morton, J., Crowell, B., Melgar, D., DeGrande, J., and Mencin, D. Real and Pseudosynthetic timeseries used in “Characterizing High Rate GNSS Velocity Noise for Synthesizing a GNSS Strong Motion Learning Catalog”, 2023. <https://zenodo.org/record/7909327>. doi: 10.5281/ZENODO.7909327.
- Ebinuma, T. and Kato, T. Dynamic characteristics of very-high-rate GPS observations for seismology. *Earth, Planets and Space*, 64(5):369–377, May 2012. doi: 10.5047/eps.2011.11.005.
- Fang, R., Zheng, J., Geng, J., Shu, Y., Shi, C., and Liu, J. Earthquake Magnitude Scaling Using Peak Ground Velocity Derived from High-Rate GNSS Observations. *Seismological Research Letters*, 92(1):227–237, Oct. 2020. doi: 10.1785/0220190347.
- Fratarcangeli, F., Ravanelli, M., Mazzoni, A., Colosimo, G., Benedetti, E., Branzanti, M., Savastano, G., Verkhoglyadova, O., Komjathy, A., and Crespi, M. The variometric approach to real-time high-frequency geodesy. *Rendiconti Lincei. Scienze Fisiche e Naturali*, 29(S1):95–108, May 2018. doi: 10.1007/s12210-018-0708-5.
- Geng, J., Pan, Y., Li, X., Guo, J., Liu, J., Chen, X., and Zhang, Y. Noise Characteristics of High-Rate Multi-GNSS for Subdaily Crustal Deformation Monitoring. *Journal of Geophysical Research: Solid Earth*, 123(2):1987–2002, Feb. 2018. doi: 10.1002/2018jb015527.
- Genrich, J. F. and Bock, Y. Instantaneous geodetic positioning with 10–50 Hz GPS measurements: Noise characteristics and implications for monitoring networks. *Journal of Geophysical Research: Solid Earth*, 111(B3), Mar. 2006. doi: 10.1029/2005jb003617.
- GRAAS, F. V. and SOLOVIEV, A. Precise Velocity Estimation Using a Stand-Alone GPS Receiver. *Navigation*, 51(4):283–292, Dec. 2004. doi: 10.1002/j.2161-4296.2004.tb00359.x.
- Grapenthin, R., West, M., and Freymueller, J. The Utility of GNSS for Earthquake Early Warning in Regions with Sparse Seismic Networks. *Bulletin of the Seismological Society of America*, June 2017. doi: 10.1785/0120160317.
- Grapenthin, R., West, M., Tape, C., Gardine, M., and Freymueller, J. Single-Frequency Instantaneous GNSS Velocities Resolve Dynamic Ground Motion of the 2016 Mw 7.1 Iniskin, Alaska, Earthquake. *Seismological Research Letters*, 89(3):1040–1048, Apr. 2018. doi: 10.1785/0220170235.
- Graves, R. W. and Pitarka, A. Broadband Ground-Motion Simulation Using a Hybrid Approach. *Bulletin of the Seismological Society of America*, 100(5A):2095–2123, Sept. 2010. doi: 10.1785/0120100057.
- Grisel, O., Mueller, A., Lars, Gramfort, A., Louppe, G., Prettenhofer, P., Blondel, M., Niculae, V., Nothman, J., Joly, A., Fan, T. J., Vanderplas, J., manoj kumar, Lemaitre, G., Qin, H., Hug, N., Estève, L., Varoquaux, N., Layton, R., Metzen, J. H., Jalali, A., (Venkat) Raghav, R., Schönberger, J., du Boisberranger, J., Yurchak, R., Li, W., la Tour, T. D., Woolam, C., Eren, K., and Eustache. scikit-learn/scikit-learn: scikit-learn 1.0.1, 2021. doi: 10.5281/ZENODO.5596244.
- Hodgkinson, K. M., Mencin, D. J., Feaux, K., Sievers, C., and Mattioli, G. S. Evaluation of Earthquake Magnitude Estimation and Event Detection Thresholds for Real-Time GNSS Networks: Examples from Recent Events Captured by the Network of the Americas. *Seismological Research Letters*, 91(3):1628–1645, Mar. 2020. doi: 10.1785/0220190269.
- Hoffmann, J., Bar-Sinai, Y., Lee, L. M., Andrejevic, J., Mishra, S., Rubinstein, S. M., and Rycroft, C. H. Machine learning in a data-limited regime: Augmenting experiments with synthetic data uncovers order in crumpled sheets. *Science Advances*, 5(4), Apr. 2019. doi: 10.1126/sciadv.aau6792.
- Hohensinn, R. and Geiger, A. Stand-Alone GNSS Sensors as Velocity Seismometers: Real-Time Monitoring and Earthquake Detection. *Sensors*, 18(11):3712, Oct. 2018. doi: 10.3390/s18113712.
- Hohensinn, R., Häberling, S., and Geiger, A. Dynamic displacements from high-rate GNSS: Error modeling and vibration detection. *Measurement*, 157:107655, June 2020. doi: 10.1016/j.measurement.2020.107655.
- Häberling, S., Rothacher, M., Zhang, Y., Clinton, J. F., and Geiger, A. Assessment of high-rate GPS using a single-axis shake table. *Journal of Geodesy*, 89(7):697–709, Apr. 2015. doi: 10.1007/s00190-015-0808-2.
- Iwana, B. K. and Uchida, S. An empirical survey of data augmentation for time series classification with neural networks. *PLOS ONE*, 16(7):e0254841, July 2021. doi: 10.1371/journal.pone.0254841.
- Kong, Q., Trugman, D. T., Ross, Z. E., Bianco, M. J., Meade, B. J., and Gerstoft, P. Machine Learning in Seismology: Turning Data into Insights. *Seismological Research Letters*, 90(1):3–14, Nov. 2018. doi: 10.1785/0220180259.
- Krischer, L., Megies, T., Barsch, R., Beyreuther, M., Lecocq, T., Caudron, C., and Wassermann, J. ObsPy: a bridge for seismology into the scientific Python ecosystem. *Computational Science & Discovery*, 8(1):014003, May 2015. doi: 10.1088/1749-4699/8/1/014003.
- Langbein, J. and Bock, Y. High-rate real-time GPS network at

- Parkfield: Utility for detecting fault slip and seismic displacements. *Geophysical Research Letters*, 31(15), 2004. doi: 10.1029/2003gl019408.
- Lin, J.-T., Melgar, D., Thomas, A. M., and Searcy, J. Early Warning for Great Earthquakes From Characterization of Crustal Deformation Patterns With Deep Learning. *Journal of Geophysical Research: Solid Earth*, 126(10), Oct. 2021. doi: 10.1029/2021jb022703.
- McNamara, D. E. and Buland, R. P. Ambient Noise Levels in the Continental United States. *Bulletin of the Seismological Society of America*, 94(4):1517–1527, Aug. 2004. doi: 10.1785/012003001.
- Meier, M.-A., Ross, Z. E., Ramachandran, A., Balakrishna, A., Nair, S., Kundzicz, P., Li, Z., Andrews, J., Hauksson, E., and Yue, Y. Reliable Real-Time Seismic Signal/Noise Discrimination With Machine Learning. *Journal of Geophysical Research: Solid Earth*, 124(1):788–800, Jan. 2019. doi: 10.1029/2018jb016661.
- Melgar, D., Bock, Y., Sanchez, D., and Crowell, B. W. On robust and reliable automated baseline corrections for strong motion seismology. *Journal of Geophysical Research: Solid Earth*, 118(3): 1177–1187, Mar. 2013. doi: 10.1002/jgrb.50135.
- Melgar, D., LeVeque, R. J., Dreger, D. S., and Allen, R. M. Kinematic rupture scenarios and synthetic displacement data: An example application to the Cascadia subduction zone. *Journal of Geophysical Research: Solid Earth*, 121(9):6658–6674, Sept. 2016. doi: 10.1002/2016jb013314.
- Melgar, D., Crowell, B. W., Melbourne, T. I., Szeliga, W., Santillan, M., and Scrivner, C. Noise Characteristics of Operational Real-Time High-Rate GNSS Positions in a Large Aperture Network. *Journal of Geophysical Research: Solid Earth*, 125(7), June 2020. doi: 10.1029/2019jb019197.
- Moschas, F. and Stiros, S. Noise characteristics of high-frequency, short-duration GPS records from analysis of identical, collocated instruments. *Measurement*, 46(4):1488–1506, May 2013. doi: 10.1016/j.measurement.2012.12.015.
- Mousavi, S. M. and Beroza, G. C. Deep-learning seismology. *Science*, 377(6607), Aug. 2022. doi: 10.1126/science.abm4470.
- Mousavi, S. M., Ellsworth, W. L., Zhu, W., Chuang, L. Y., and Beroza, G. C. Earthquake transformer—an attentive deep-learning model for simultaneous earthquake detection and phase picking. *Nature Communications*, 11(1), Aug. 2020. doi: 10.1038/s41467-020-17591-w.
- Murray, J. R., Crowell, B. W., Murray, M. H., Ulberg, C. W., McGuire, J. J., Aranha, M. A., and Hagerty, M. T. Incorporation of Real-Time Earthquake Magnitudes Estimated via Peak Ground Displacement Scaling in the ShakeAlert Earthquake Early Warning System. *Bulletin of the Seismological Society of America*, 113(3): 1286–1310, Feb. 2023. doi: 10.1785/0120220181.
- Parameswaran, R. M., Grapenthin, R., West, M. E., and Fozkos, A. Interchangeable Use of GNSS and Seismic Data for Rapid Earthquake Characterization: 2021 Chignik, Alaska, Earthquake. *Seismological Research Letters*, Mar. 2023. doi: 10.1785/0220220357.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-Learn: Machine Learning in Python. *J. Mach. Learn. Res.*, 12(null):2825–2830, 2011.
- Ruhl, C. J., Melgar, D., Geng, J., Goldberg, D. E., Crowell, B. W., Allen, R. M., Bock, Y., Barrientos, S., Riquelme, S., Baez, J. C., Cabral-Cano, E., Pérez-Campos, X., Hill, E. M., Protti, M., Ganas, A., Ruiz, M., Mothes, P., Jarrín, P., Nocquet, J., Avouac, J., and D’Anastasio, E. A Global Database of Strong-Motion Displacement GNSS Recordings and an Example Application to PGD Scaling. *Seismological Research Letters*, 90(1):271–279, Oct. 2018. doi: 10.1785/0220180177.
- Seydoux, L., Balestriero, R., Poli, P., Hoop, M. d., Campillo, M., and Baraniuk, R. Clustering earthquake signals and background noises in continuous seismic data with unsupervised deep learning. *Nature Communications*, 11(1), Aug. 2020. doi: 10.1038/s41467-020-17841-x.
- SHEN, N., CHEN, L., LIU, J., WANG, L., TAO, T., WU, D., and CHEN, R. A Review of Global Navigation Satellite System (GNSS)-based Dynamic Monitoring Technologies for Structural Health Monitoring. *Remote Sensing*, 11(9):1001, Apr. 2019. doi: 10.3390/rs11091001.
- Shu, Y., Fang, R., Li, M., Shi, C., Li, M., and Liu, J. Very high-rate GPS for measuring dynamic seismic displacements without aliasing: performance evaluation of the variometric approach. *GPS Solutions*, 22(4), Sept. 2018. doi: 10.1007/s10291-018-0785-z.
- Shu, Y., Fang, R., Liu, Y., Ding, D., Qiao, L., Li, G., and Liu, J. Precise coseismic displacements from the GPS variometric approach using different precise products: Application to the 2008 MW 7.9 Wenchuan earthquake. *Advances in Space Research*, 65(10): 2360–2371, May 2020. doi: 10.1016/j.asr.2020.02.013.
- Team, T. O. D. ObsPy 1.2.1, 2020. doi: 10.5281/ZENODO.3706479.
- Teunissen, P. J. GNSS Precise Point Positioning, Dec. 2020. doi: 10.1002/9781119458449.ch20.
- Trnkoczy, A. Topic Understanding and parameter setting of STA/LTA trigger algorithm 1 Introduction. *New Manual of Seismological Observatory Practice 2 (NMSOP-2)*, 2012. doi: https://doi.org/0.2312/GFZ.NMSOP-2_IS_8.1.
- Wang, Y., Breitsch, B., and Morton, Y. T. J. A State-Based Method to Simultaneously Reduce Cycle Slips and Noise in Coherent GNSS-R Phase Measurements From Open-Loop Tracking. *IEEE Transactions on Geoscience and Remote Sensing*, 59(10): 8873–8884, Oct. 2021. doi: 10.1109/tgrs.2020.3036031.
- Williams, S. D. P., Bock, Y., Fang, P., Jamason, P., Nikolaidis, R. M., Prawirodirdjo, L., Miller, M., and Johnson, D. J. Error analysis of continuous GPS position time series. *Journal of Geophysical Research*, 109(B3), 2004. doi: 10.1029/2003jb002741.
- Williamson, A. L., Melgar, D., Crowell, B. W., Arcas, D., Melbourne, T. I., Wei, Y., and Kwong, K. Toward Near-Field Tsunami Forecasting Along the Cascadia Subduction Zone Using Rapid GNSS Source Models. *Journal of Geophysical Research: Solid Earth*, 125(8), Aug. 2020. doi: 10.1029/2020jb019636.
- Withers, M., Aster, R., Young, C., Beiriger, J., Harris, M., Moore, S., and Trujillo, J. 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, Feb. 1998. doi: 10.1785/bssa0880010095.
- Woollam, J., Münchmeyer, J., Tilmann, F., Rietbrock, A., Lange, D., Bornstein, T., Diehl, T., Giunchi, C., Haslinger, F., Jozinović, D., Michelini, A., Saul, J., and Soto, H. SeisBench—A Toolbox for Machine Learning in Seismology. *Seismological Research Letters*, 93(3):1695–1709, Mar. 2022. doi: 10.1785/0220210324.
- Yang, R., Morton, Y., Ling, K.-V., and Poh, E.-K. Generalized GNSS Signal Carrier Tracking—Part II: Optimization and Implementation. *IEEE Transactions on Aerospace and Electronic Systems*, 53(4):1798–1811, Aug. 2017. doi: 10.1109/taes.2017.2674198.
- Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. How transferable are features in deep neural networks? 2014. doi: 10.48550/ARXIV.1411.1792.
- Zhu, W., Mousavi, S. M., and Beroza, G. C. Seismic signal augmentation to improve generalization of deep neural networks. In *Machine Learning in Geosciences*, pages 151–177. Elsevier, 2020. doi: 10.1016/bs.agph.2020.07.003.

The article *Characterizing High Rate GNSS Velocity Noise for Synthesizing a GNSS Strong Motion Learning Catalog* © 2023 by Timothy Dittmann is licensed under CC BY 4.0.