# GWAK: Gravitational-Wave Anomalous Knowledge with Recurrent Autoencoders

**Ryan Raikman**[1,5*]**, Eric A. Moreno**[2]**, Ekaterina Govorkova**[2]**, Ethan J Marx**[1,2]**, Alec Gunny**[1,2]**, William Benoit**[3]**, Deep Chatterjee**[1,2]**, Rafia Omer**[3]**, Muhammed Saleem**[3]**, Dylan S Rankin**[4]**, Michael W Coughlin**[3]**, Philip C Harris**[2]**, Erik Katsavounidis**[1,2]

[1]MIT LIGO Laboratory, USA,
[2]Massachusetts Institute of Technology, USA,
[3]University of Minnesota, USA,
[4]University of Pennsylvania, USA,
[5]Carnegie Mellon University, USA

E-mail: *rraikman@mit.edu

July 2023

**Abstract.** Matched-filtering detection techniques for gravitational-wave (GW) signals in ground-based interferometers rely on having well-modeled templates of the GW emission. Such techniques have been traditionally used in searches for compact binary coalescences (CBCs), and have been employed in all known GW detections so far. However, interesting science cases aside from compact mergers do not yet have accurate enough modeling to make matched filtering possible, including core-collapse supernovae and sources where stochasticity may be involved. Therefore the development of techniques to identify sources of these types is of significant interest. In this paper, we present a method of anomaly detection based on deep recurrent autoencoders to enhance the search region to unmodeled transients. We use a semi-supervised strategy that we name *"Gravitational Wave Anomalous Knowledge"* (GWAK). While the semi-supervised nature of the problem comes with a cost in terms of accuracy as compared to supervised techniques, there is a qualitative advantage in generalizing experimental sensitivity beyond pre-computed signal templates. We construct a low-dimensional embedded space using the GWAK method, capturing the physical signatures of distinct signals on each axis of the space. By introducing signal priors that capture some of the salient features of GW signals, we allow for the recovery of sensitivity even when an unmodeled anomaly is encountered. We show that regions of the GWAK space can identify CBCs, detector glitches and also a variety of unmodeled astrophysical sources.

## 1. Introduction

Since the original observation of gravitational waves (GW) [1] by Advanced LIGO [2] and Advanced VIRGO [3], and with the recent introduction of KAGRA [4], more than 90 gravitational-wave events [5–7] have been catalogued to date, fundamentally transforming our way of observing the universe. While all of the detected signals thus far correspond to the coalescence of binary black hole (BBH), binary neutron star (BNS), or black hole - neutron star (BHNS) mergers [8–12], there is growing interest in the detection of unmodeled signals, which are not described by any known theoretical waveforms, computationally prohibitive to simulate, or are stochastic in nature. These signals may originate from sources such as core-collapse supernovae (CCSN) [13], as well as exotic sources such as cosmic strings [14, 15], axion stars [16], neutron star glitches, or primordial black holes [17, 18]. Detection of such sources could lead to new insights into the fundamental physics of the universe.

Typical methods for detecting GWs rely on matched filtering techniques [19], which compare the observed data to a known signal template. These methods require precise knowledge of the signal prior, such as the waveform and the parameters of the source, in order to detect the signal. While matched filtering is a well-established and powerful method for detecting gravitational waves, it has certain limitations. For example, matched filtering is sensitive only to signals that match the known templates, and may miss signals with different waveforms or parameters. This makes matched filtering unsuitable for unmodeled signal searches.

The GW community has developed several unmodeled approaches that do not rely on a specific waveform model. Some of these currently used by the International Gravitational-wave Network (IGWN) include cWB [20, 21], which searches for and reconstructs GW transient signals without relying on a specific waveform model. Another framework, oLIB [22] uses the Q transform to decompose GW strain data into several time-frequency planes of constant quality factors. The pipeline flags data segments containing excess power and searches for clusters of these segments to identify possible GW candidate events. MLy (read as "Emily") [23] is a machine-learning-based search for generic sub-second-duration transient GW signals in the 20 to 500 Hz frequency band; MLy works by utilizing convolutional neural networks (CNNs) trained to recognize signals that are simultaneous and coherent between detectors.

In this paper, we explore the use of a method introduced in Ref. [24] deployed within the High Energy Physics community for the development of an anomaly detection pipeline for data collected by GW observatories. We introduce *"Gravitational Wave Anomalous Knowledge"* (GWAK), a strategy for anomaly search that combines deep learning (DL) techniques with prior information on potential signals to improve the sensitivity of detection. The GWAK algorithm is based on the intuition that unknown transient sources should loosely resemble known signals, as well as be coherent between present GW detectors. We apply the GWAK method to GW datasets, by introducing signal priors that capture some of the salient features of GW signatures, allowing for the

recovery of sensitivity even when the observed signal mismatches the known priors.

Because GWAK does not rely on precise prior knowledge of the signal and can detect signals with unknown waveforms or parameters by matching incoming data streams with salient features (cross-correlation, oscillations, etc.) that are generic to broad types of GWs, GWAK is more robust and powerful for the detection of unknown signals. As such, it can be used as a complementary approach to methods such as matched filtering for detecting transient GWs, and it has the potential to improve the sensitivity of GW detection systems to sources of this type.

This paper is organized as follows: in Section 2, we provide a brief review of deep learning in GW detection and previous works in machine learning anomaly detection. In Section 3.2, we describe the data used for this study. In Section 3, we present the GWAK method for constructing embedded spaces for anomalous searches and the autoencoder and transformer architectures used to build these embedded spaces. In Section 4, we discuss the performance of the GWAK method on real GW data. Conclusions and next steps are provided in Section 5.

The code used to analyze data and generate results and plots can be found at ‡.

## 2. Related work

Deep learning approaches for GW detection are well explored [25–35]. However, these methods typically rely on supervised learning techniques, which provide competitive efficiency by exploiting neural network nonlinearity and the information provided by ground-truth labels. By construction, these methods rely on a realistic simulation of the signal generated by a specific kind of source, which is assumed upfront. With supervised DL approaches, there is no guarantee of generalizability to an out-of-training event, what we refer to as "anomalies."

Autoencoders are commonly implemented for a variety of GW applications. They can be used for non-linear subtraction of noise from incoming time-series of GW strain [36–38] in real-time [39]. Additionally, autoencoders are used in generative models, speeding up the computation of GW waveforms relative to very computationally expensive numerical relativity simulations [40]. Finally, due to the transient nature of single-detector artifacts, "glitches", autoencoders can be used for glitch classification in an unsupervised manner [41].

There have also been explorations into unsupervised GW detection to enhance detection capabilities beyond signal templates and simulation. Initial studies with a Convolutional Neural Network (CNN) based autoencoder [42] and Long-Short Term Memory (LSTM) based autoencoder [43] show that unsupervised detection is possible. Both studies rely on learning the typical features of the background. Once a model is trained, it is then used to evaluate the similarity between background priors and new unseen data. To do this, the algorithms use reconstruction loss, computed by comparing the original signal and the signal outputted by the model trained on the

‡ https://github.com/ML4GW/gwak

background prior, as a detection statistic. By comparing the reconstruction error of new data with a threshold corresponding to an allowed false alarm rate (FAR), data points that deviate enough from the normal pattern are identified as anomalous. This relies on the autoencoder's inability to reconstruct any potential signal that deviates from the background, or generally the prior on which the model was trained, triggering a high reconstruction loss. This approach has been shown to effectively detect out-of-training anomalies in GW datasets and has the potential to improve the performance of anomaly detection systems. However, for a specific signal, an algorithm trained with an unsupervised procedure on unlabeled data is typically less accurate than a supervised classifier trained on labeled data.

We build upon this approach by training several autoencoders; each with a different signal prior to enhance signal sensitivity. The signals that were chosen to be used in this paper are described in Sec. 3.2. Unlike an earlier study [43] which used simulated Gaussian background as a proof of concept, our method is trained with real background data. This is significantly more challenging than just simulated Gaussian noise, and therefore no direct comparison of the results presented here to the previous ones can be made.

Although the autoencoder architecture can implicitly learn detector correlation, a related method explicitly leverages the correlation across detectors [44]. The method trains detector correlation with a white-noise burst (WNB) [45, 46] signal prior, which should be harder to detect than any other type of signal given its lack of a distinctive morphology. Similarly to [44], we choose not to rely on our autoencoders to learn the correlation between the two detector sites. Instead, we directly compute and include the correlation in our final metric that is used to find signal events as another axis in the GWAK embedded space. As such, the present version of our method is applicable to the case of aligned GW detectors. However, more off-plane GW detectors are being proposed, and a future area of work is generalizing the method to an arbitrary detector network.

## 3. GWAK algorithm

GWAK (reads: *guac*) builds on the concept of semi-supervision, pulling on concepts from both supervised and unsupervised learning. The semi-supervised method is manifested by using simulated signals as approximations for anomalous-unmodeled signals in the GWAK embedded space. We use five classes of datasets that can help us to build an informative space to search for these unmodeled signals. A separate unsupervised autoencoder network is then trained on each class of samples separately, resulting in a low-dimensional "GWAK" space consisting of the coherence metrics between the autoencoder inputs and outputs for each autoencoder, which is then used to search for anomalous signals. This approach is particularly useful for detecting new physics phenomena, where the signal prior is unknown but the simulation of some potential signal pattern, like BBH and sine-Gaussians, is available. The GWAK method results

in classes of anomalous signals inhabiting different regions of the continuous GWAK space, each being reconstructed differently by the five autoencoders. Searches can then be performed on the lower-dimensional embedded space in the regions that anomalies are expected to inhabit.
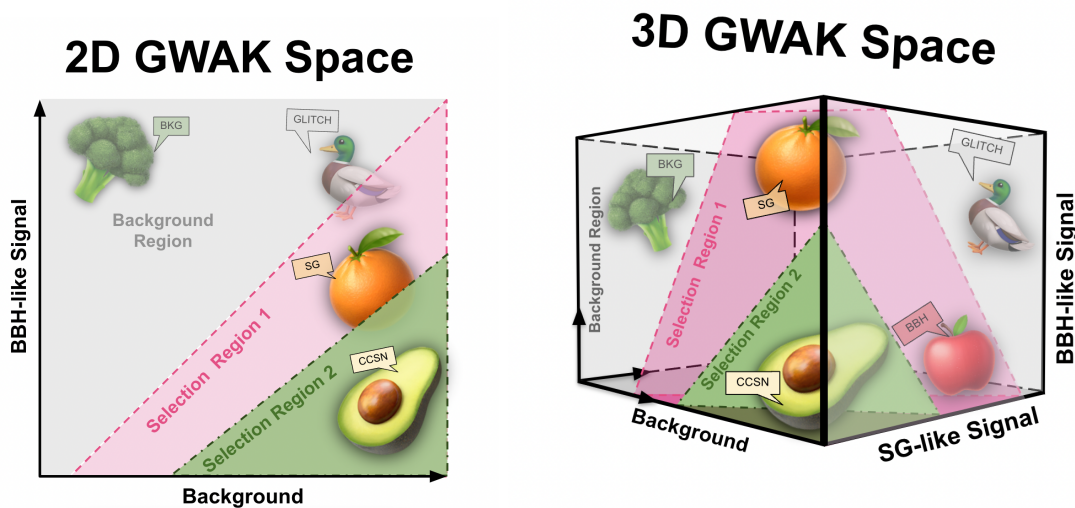


Figure 1: The 2D and 3D schematic depicting the main idea behind the GWAK algorithm. For plotting purposes we only show one background and two signal axes, while in this paper we are actually using 5 (2 for background and 3 for signals), and the main idea remains the same. All the fruits (orange, avocado and apple) are "signals" that we would like to select, broccoli represents background that we would like to suppress and the duck represents true anomalies, such as detector malfunctioning, and glitches. Shaded pink and green region depicts selection regions that would correspond to different FARs. For higher FAR (pink shading) more signal-like anomalies would be picked up (orange, avocado and apple), while for a selection that corresponds to a lower FAR (green shading), only "avocado signal" would be picked up.

### 3.1. Network architectures

One of the main advantages of LSTM autoencoders is their ability to handle sequential data with temporal dependencies. This makes them suitable for anomaly detection in time-series data, such as GWs, speech signals, and sensor data. The LSTM autoencoder consists of an encoder and a decoder, where the encoder maps the input sequence to a fixed-length vector, and the decoder maps the vector back to the original sequence. We use a similar architecture to that optimized in [43].

For the signal classes, we used the LSTM autoencoder as described above, with the bottleneck size of 4, 8 and 8 for binary black hole (BBH), low and high frequency sine-gaussian (SG) signals, respectively. For the background classes, we used a fully connected dense model. This was done as the signal classes have temporal behavior

to exploit, whereas glitches have smaller-scale, localized features. The total number of trainable parameters for the BBH LSTM AE is 510324, for SG 64–512 Hz and SG 512–1024 Hz are 511672, for background AE 243352 and for glitches AE is 241302.
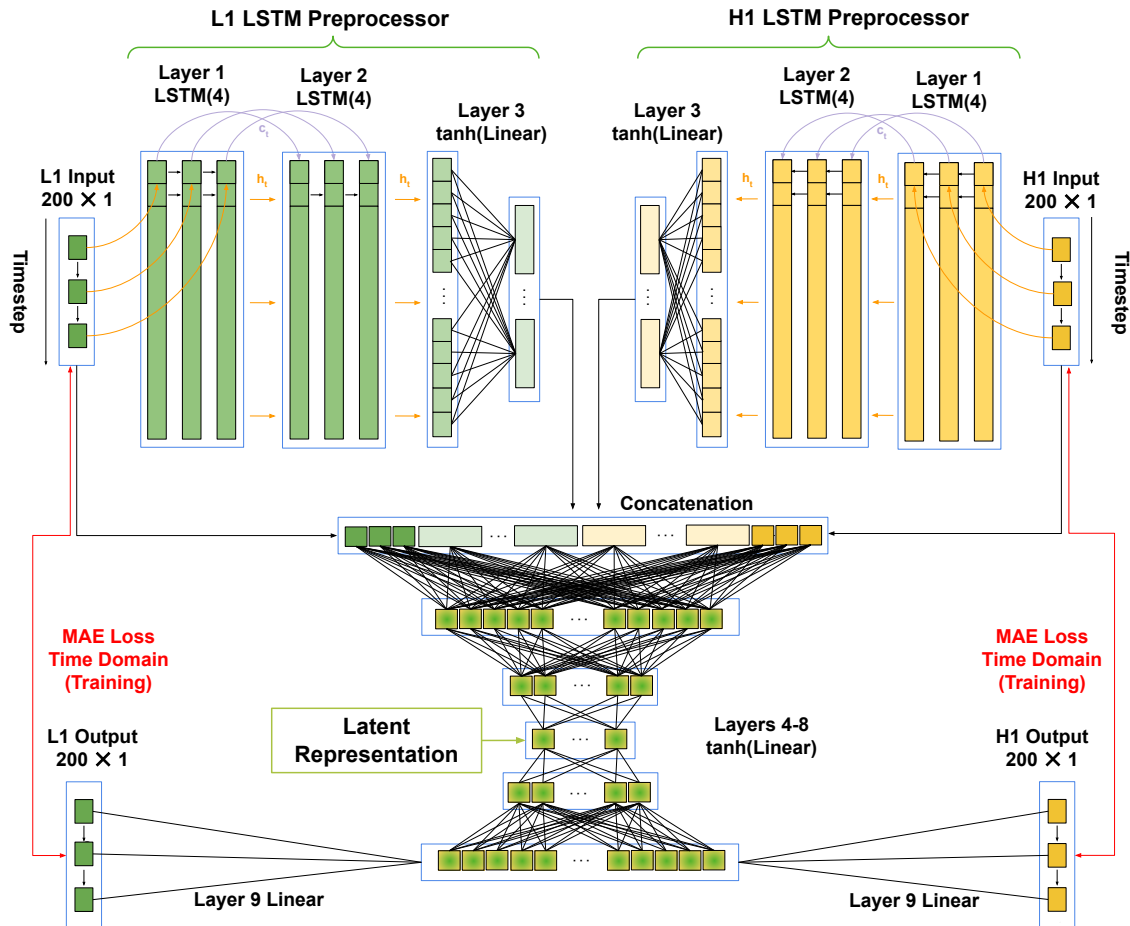


Figure 2: Graphical illustration of LSTM autoencoder, preserving the traditional encoder and decoder structure, which allows for a reconstruction loss detection statistic.

### 3.2. Data samples

The dataset used in this study was collected by the LIGO Hanford (H1) and LIGO Livingston (L1) [47] gravitational-wave detectors during the first half of the third observing run (O3a), which took place between 1st April 2019 and 1st Oct 2019. We specifically used publicly available data between GPS times of 1238166018 and 1238170289, right at the beginning of the run. Next, the time-series data were downsampled from 16384 Hz to 4096 Hz, and processed to remove and create a separate dataset of transient instrumental artifacts (glitches) using the excess power identification algorithm Omicron [48]. We used $Q_{min} = 3.3166$, $Q_{max} = 108$, and $f_{min} = 32$ for the Omicron algorithm. We then took 4 s segments of the data without noise artifacts to
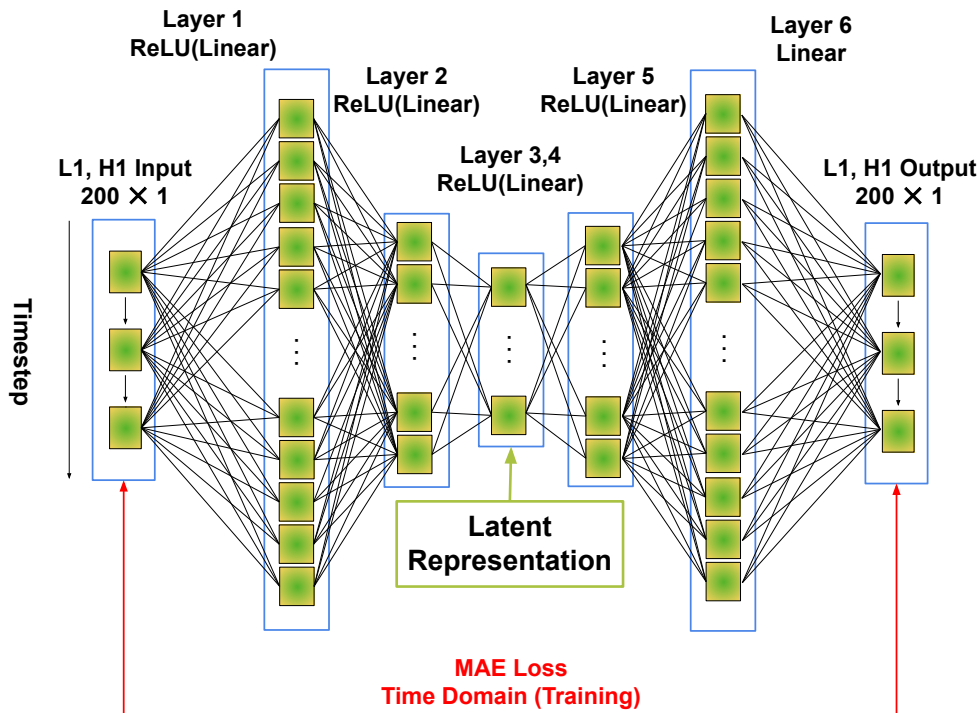
Figure 3: Graphical illustration of DNN autoencoder, preserving the traditional encoder and decoder structure, which allows for a reconstruction loss detection statistic.

serve as the baseline for the injection of signals. As such, we created five different classes of data as proxies for signals and background signatures:

- **Binaries** – Simulated BBH signals using IMRPhenomPv2 [49–51] injected into the real background noise, as shown in Fig. 5(top left). The simulation parameters are given in Table 1.

- **Background** – Background from O3a with DQsecDB § state flag `DCS-ANALYSIS_READY_C01:1` applied and excess power glitches [48] and known GW-events removed, as shown in Fig. 4(bottom).

- **Sine-Gaussian 64 − 512 Hz** – Generic low frequency signal model used to simulate generic GW sources, as shown in Fig. 5 (middle left).

- **Sine-Gaussian 512 − 1024 Hz** – Generic high frequency signal model used to simulate generic GW sources, as shown in Fig. 5 (bottom left).

- **Glitches** – Transient instrumental glitches (often of unknown origin) flagged by Omicron [48] as having excess power, as shown in Fig. 4 (top).

To create samples of BBH and SG signals, we used numerical simulations to generate

|        | Parameter | Prior | Limits | Units |
|--------|-----------|-------|--------|-------|
| BBH | $m_1$ | - | $(5, 100)$ | $M_\odot$ |
| | $m_2$ | - | $(5, 100)$ | $M_\odot$ |
| | Mass ratio $q$ | Uniform | $(0.125, 1)$ | - |
| | Chirp mass $M_c$ | Uniform | $(25, 100)$ | $M_\odot$ |
| | Tilts $\theta_{1,2}$ | Sine | $(0, \pi)$ | rad. |
| | Phase $\phi$ | Uniform | $(0, 2\pi)$ | rad. |
| | Right Ascension | Uniform | $(0, 2\pi)$ | rad. |
| | Declination $\delta$ | Cosine | $(-\pi/2, \pi/2)$ | rad. |
| sine-Gaussian | $Q$ | Uniform | $(25, 75)$ | - |
| | Frequency | Uniform | $(64, 512)$ and $(512, 1024)$ | Hz |
| | Phase $\phi$ | Uniform | $(0, 2\pi)$ | rad. |
| | Right Ascension | Uniform | $(0, 2\pi)$ | rad. |
| | Declination $\delta$ | Cosine | $(-\pi/2, \pi/2)$ | rad. |
| | Eccentricity | Uniform | $(0, 0.01)$ | - |
| | $\Psi$ | Uniform | $(0, 2\pi)$ | rad. |

Table 1: Sampling parameters and priors for BBH (top) and sine-Gaussian (bottom) injections.

$h_+$ and $h_\times$ polarization modes. We then sampled sky localizations uniformly in the sky, projected the polarization modes onto the sky location, and injected the projected modes into the two LIGO detectors.

The BBH sample is generated with the parameters and priors [52] as shown in Table 1, taken from `bilby` processing spins BBH prior, and the sine-Gaussian samples are generated with the parameters and priors as shown in Table 1.

We employed a series of digital signal processing techniques to prepare data for training autoencoders in the context of GW detection. Specifically, we first applied a whitening filter to normalize the data with respect to one hour of surrounding background data. This filter effectively suppressed frequency regions dominated by noise and reduced effects from spectral lines ‖ ¶. Moreover, we implemented a band-pass filter within the frequency range of 30–1500 Hz to further attenuate noise outside of the most sensitive frequency range of GW instruments. After applying these filters, we removed 1 s intervals from each end of the data samples to eliminate any edge effects from preprocessing. The remaining 2 s samples, each containing either an injected signal, pure background, a low/high frequency sine-Gaussian, or a glitch artifact, were used to generate training data.

To obtain a set of windows suitable for training, we extracted 200 data-points (total duration of 50 ms sampled at 4096 Hz) from each sample. Our experimentation revealed

‖ https://gwosc.org/s6speclines/
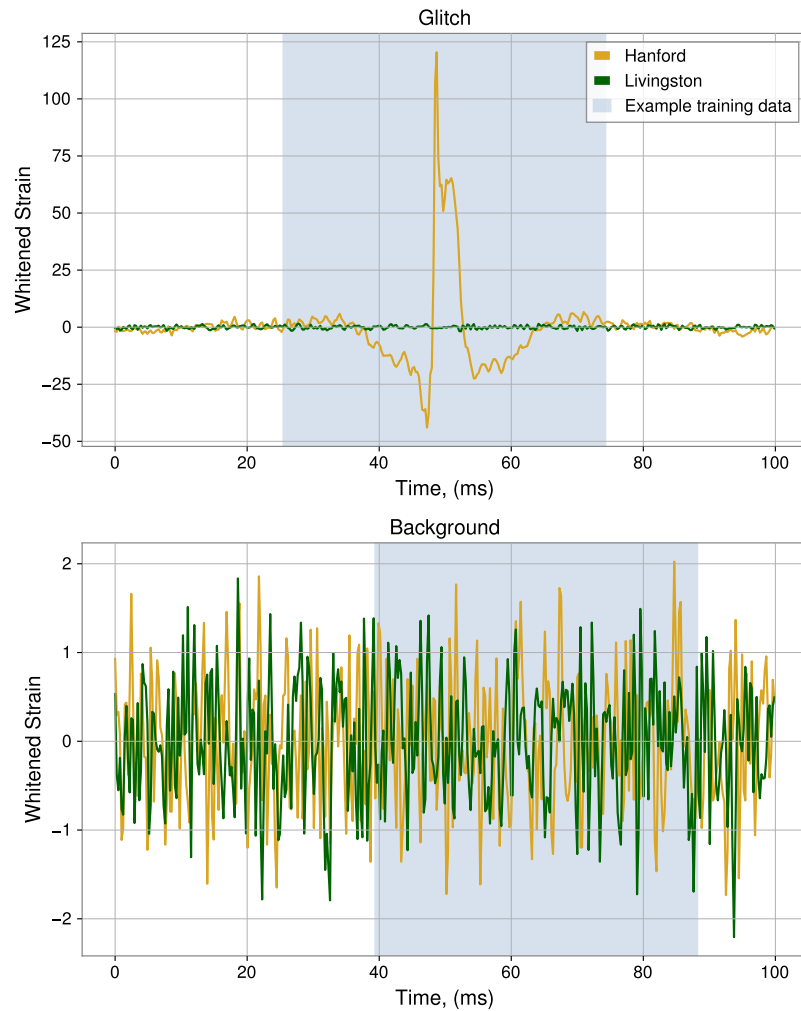¶ https://dcc.ligo.org/LIGO-T1500415/public

Figure 4: Example of GWAK classes: glitch (top) and background (bottom) strains. The light blue shading highlights an example region that is passed as input to the autoencoders for training.

that training the autoencoders with input lengths greater than 200 data-points would result in degraded performance, due to the difficulty to capture long-term behavior in an LSTM, especially on an evolving signal. While reducing the number of data-points below 200 could enhance computational efficiency, it reduces the autoencoder's ability to learn the evolution of a shorter-duration signal.

To optimize the data processing and facilitate learning by the network, the data are normalized to have a standard deviation of one on a sample-per-sample basis. This normalization was undertaken mainly for the reason that the neural networks struggled to learn with unnormalized samples. The strongest example of this is with the glitch dataset, where strain magnitude can reach amplitudes hundreds to thousands of times above the background.

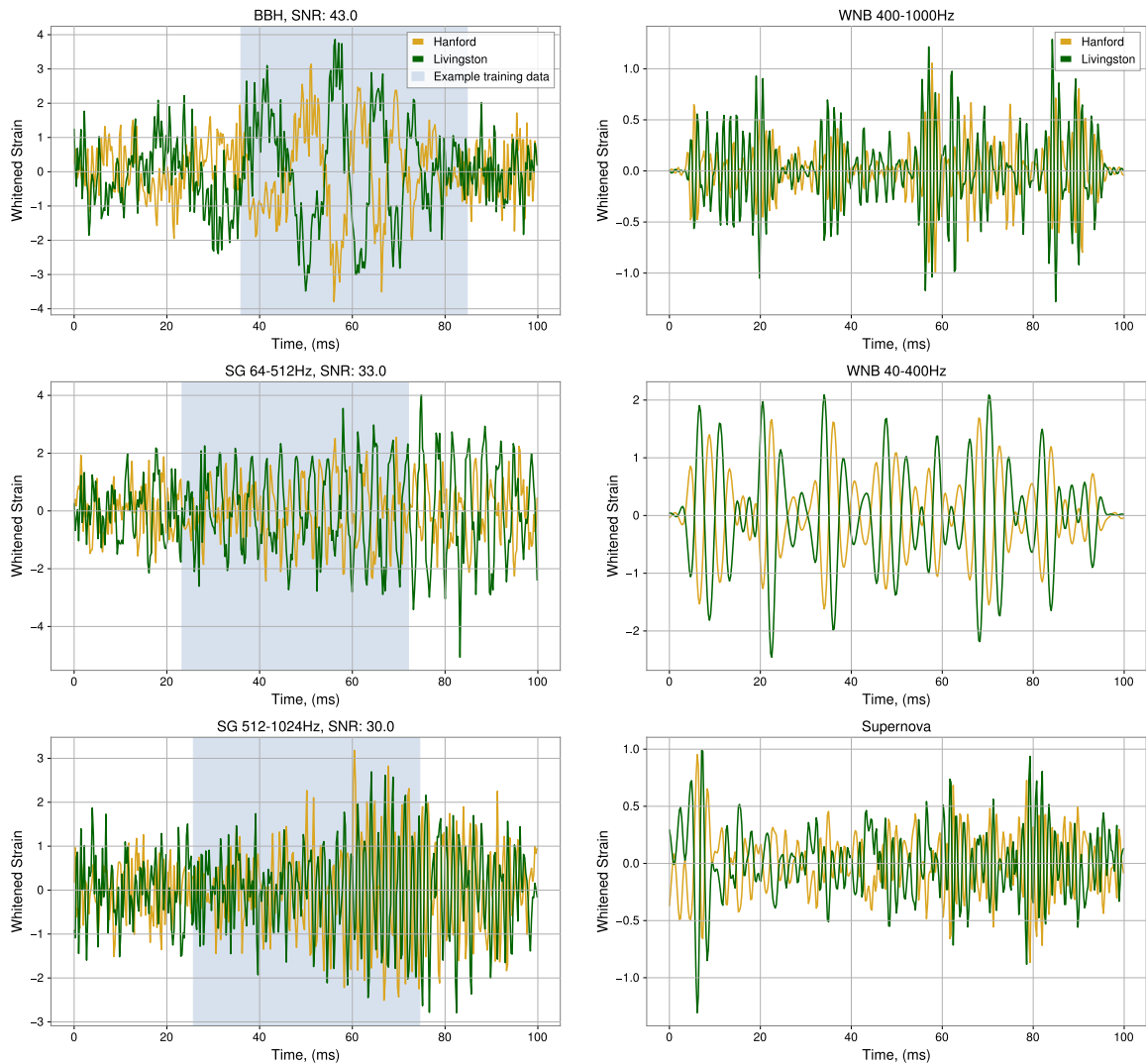For each of the non-coherent classes (background and glitch), the dataset is split

Figure 5: Example of signal-like classes: BBH (top left), WNB (top, middle right), sine-gaussian (middle, bottom left) and Supernova (bottom right) strains. The light blue shading highlights an example region that is passed as input to the autoencoders for training.

into three parts: 80,000 training samples (80%), 10,000 validation samples (10%), and 10,000 test samples (10%). For each signal class (BBH, low frequency SG and high frequency SG), we generate 5 sub-datasets, each with a specified SNR injection range, as shown in Fig. 6. Each of these sub-datasets has an identical splitting procedure: 80,000 training samples (80%), 10,000 validation samples (10%), and 10,000 test samples (10%). The training and validation datasets were used for the autoencoders to create loss values on which to update and score the networks respectively. The test samples were used for the recreation plots, as in Fig. 7, as well as to show the GWAK feature space, Fig. 10.

Signal events needed to build the GWAK space were created injecting simulated GWs on top of artifact-free detector noise. This provides an analogous situation to a real

GW, in which the strain from the incoming wave is recorded in combination with the detector noise. We do not explore the case of coincident detector aritfacts with transient astrophysical signals. The injection of signal also accounts for the difference in GW time-of-arrival at each detector owing to light travel time from the sky localization of the signal, which is significant at a sampling rate of 4096 Hz, corresponding to a maximum of 40 samples.

### 3.3. Autoencoder Training

To train an autoencoder, the input sequence is passed through the encoder and the decoder, and the model is optimized to minimize the reconstruction error between the input and the output sequence. Once the model is trained, it can be used to identify data points that deviate from the normal pattern by comparing the reconstruction error of new data with a threshold.

   To train our five autoencoders, each corresponding to one of the data classes, we used two different schemes. To train the glitch and background classes, we use the same dataset for all 200 epochs, the Adam [53] optimizer, and Mean Absolute Error (MAE) loss, computed between the autoencoder input and autoencoder reconstruction. To train the binary black hole, sine-Gaussian low frequency and high frequency classes, we used a curriculum scheme. Over 200 epochs, we used 5 different datasets, each with identical injections but different uniform SNR priors: $U[192, 384]$ (epochs 1-40), $U[96, 192]$ (epochs 41-80), $U[8, 96]$ (epochs 81-120), $U[24, 48]$ (epochs 121-160), $U[12, 24]$ (epochs 161-200), as shown in Fig. 6. Here, we used the Adam optimizer, reset after each step of the curriculum, and computed error between the autoencoder output of a noisy input (injection into the real background with specified SNR) and the "clean," noise-less template. This was intended to train the autoencoders to learn to reconstruct the signal itself, without any noise. To compute the validation loss at each epoch, we used a subset of the $U[12, 24]$ SNR dataset for each curriculum. This was intended to provide a fair computation of the validation loss across each curriculum. The MAE loss $\mathbf{L}$ between original data $\mathbf{D}$ and reconstruction $\mathbf{R}$ is given by

$$\mathbf{L} = \frac{1}{N}\frac{1}{2}\frac{1}{200}\sum_{i=1}^{N}\sum_{j=1,2}\sum_{k=1}^{200}|\mathbf{D}_i[j,k] - \mathbf{R}_i[j,k]|$$

$\mathbf{D}_i[j,k]$ represents the i-th sample of the original data, taken from the j-th detector at the k-th timestep, and likewise $\mathbf{R}_i[j,k]$ represents the i-th sample of the autoencoder output, taken from the j-th detector at the k-th timestep. The loss curves are shown in Fig. 6. The example of autoencoder reconstruction obtained with pre-trained autoencoders is shown in Fig. 7, and recreation samples of other training classes are shown at the end of the paper.
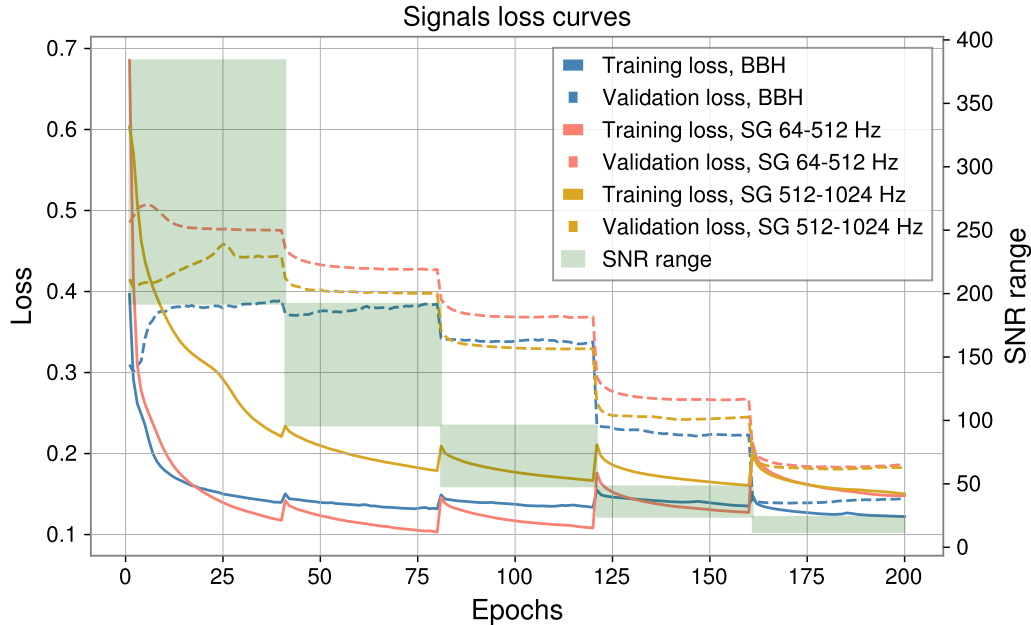
Figure 6: Autoencoder training and validation losses for signal classes, using curriculum learning to progressively reduce validation loss. The validation loss for each training SNR range is computed on the validation data from the last SNR step. In solid/dashed colours are the training/validation losses for BBH (blue), SG 64–512 Hz (salmon) and SG 512–1024 Hz (dark yellow). A light green shaded region depicts the SNR range for each step of training, spanning the range of injected SNR corresponding to each curriculum. For example, the first curriculum contains injections in the SNR range [192, 384]. At the transition between curricula, we see a small spike in the training loss, especially for the SG 512–1024 Hz model. This is due to the transition to a "more difficult" training dataset, as the lower SNR leads to a less distinguishable signal. With the transition between curricula, we see a sharp drop in validation loss over the course of a few epochs. This is due to the fact that the validation set is maintained for each autoencoder through training as belonging to the lowest - $U[12, 24]$ SNR range. With the transition to a new curriculum with a lower SNR range, the training data for that curriculum will more closely match the validation set, explaining the rapid drop.

### 3.4. Feature extraction

While the MAE loss was used in training, at evaluation time we opted for a frequency-domain based features, computed between the input and autoencoder output. The intuition for choosing to compute features in the frequency domain is based on the fact that signal autoencoders were trained with clean, noiseless targets. Since the reconstruction is a noise-less signal, and one 50 ms window of a signal does not exist in a wide range of frequencies, the reconstruction will be a localized peak in frequency space. On the contrary, the original signal, especially one of low SNR, will contain both the
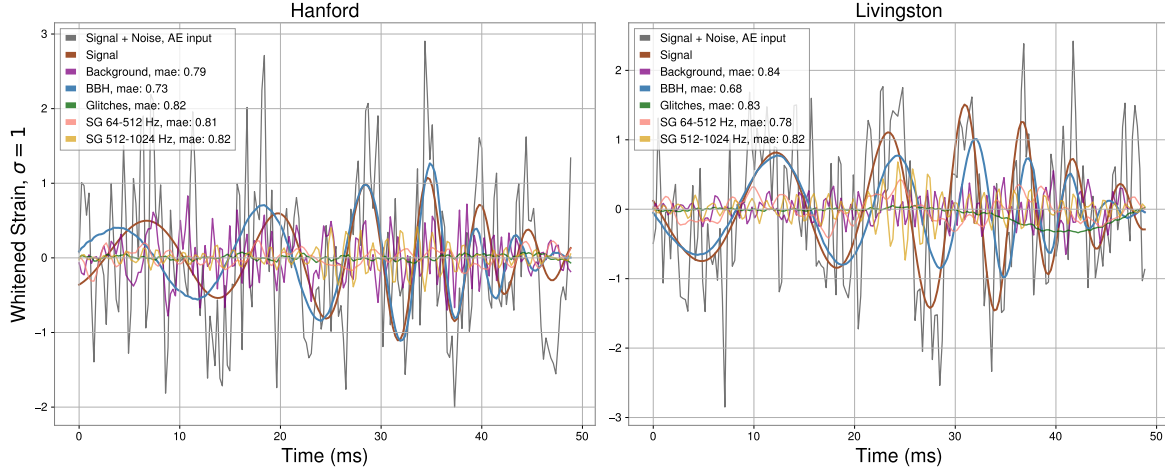
Figure 7: Example of recreation on injected BBH signal, with the noise-less template also shown. The recreation of the BBH autoencoder (in blue) follows closely the original signal injection (in brown) as expected, since it was trained to reconstruct the noiseless input. While background (in purple), glitches (in green), SG 64–512 Hz and SG 512–1024 Hz clearly fail to reconstruct the injected BBH signal, as expected.

narrow-bandwidth frequency feature along with noise distributed approximately evenly throughout the entire frequency range. When computing the MAE between the original input and reconstructed output, the presense of high frequency noise will inflate the loss, since the autoencoder, by design, does not fit noise. To bypass this, we chose to compute a dot product in frequency space. In the frequency regime where the true signal exists, both original and reconstructed signals will be similar, and as such the dot product will yield a high value. In the other frequency regimes where the signal is not present, the reconstructed signal is close to zero, and as such these noisy contributions get removed.

We choose two features per autoencoder to be the following: Let $H_O, L_O, H_R, L_R$ correspond to the original Hanford and Livingston signals and reconstructed Hanford and Livingston signals respectively, from a single autoencoder. Each are 200-datapoint segments, sampled at 4096 Hz. We then take the Fourier transform of each, yielding $\widetilde{H_O}, \widetilde{L_O}, \widetilde{H_R}, \widetilde{L_R}$. The two features, per autoencoder, are $|\widetilde{H_O} \cdot \widetilde{H_R}|$ and $|\widetilde{L_O} \cdot \widetilde{L_R}|$. We also used a general "frequency space correlation" feature to compliment the Pearson correlation 3.5, defined by $|\widetilde{H_O} \cdot \widetilde{L_O}|$. $|\widetilde{A} \cdot \widetilde{B}|$ represents the magnitude of the dot product of two complex vectors, namely the Fourier transforms of $A$ and $B$. Training an autoencoder to optimize directly on the $|\widetilde{H_O} \cdot \widetilde{H_R}|$, $|\widetilde{L_O} \cdot \widetilde{L_R}|$ features proved to be too simple of a task. Each autoencoder would consistently learn the largest feature in Fourier space, leading autoencoders generalizing too well, i.e. not being specific enough to their respective training class. By scoring the network's performance in the time domain, i.e. with MAE between the input and reconstructed output, the network must accurately recreate more details, such as temporal offsets, signal evolutions, as well as the corresponding frequency components, forcing the specificity to the class. Finally, by

training with MAE, it provides us with a nice visual picture of the autoencoder output which we can easily compare against the input, but this would not necessarily be the case of using the frequency-domain features directly for loss.

### 3.5. Pearson cross-correlation

To derive a comprehensive metric for inference, we incorporated information regarding the cross-correlation between the two detector sites, in conjunction with the GWAK information. Given that any astrophysical signal will invariably manifest in both detector sites, the correlation of the measured strains is of critical significance for the signal search procedure. Although we utilized information from both detectors during the GWAK space training phase, we opted to directly incorporate cross-correlation information in our final metric.

To accomplish this, we employed the Pearson correlation coefficient [54]. Specifically, we computed the Pearson correlation coefficient between the Hanford and inverted Livingston sites by selecting the maximum correlation coefficient from all possible time shifts for a 200 datapoint window. Since the physical separation between the detectors is about 3000 km, corresponding to a time of flight of 10 ms, we must iterate over all possible temporal shifts within 10 ms to contain the correct time delay. At a sampling rate of 4096 Hz, this corresponds to a shift of 40 datapoints in either direction. This iteration over all possible time shifts is an advantage of the Pearson correlation coefficient over the frequency-domain correlation coefficient, and as such we chose to include both in our GWAK space. The source of this time delay is due to the sky location of the gravitational wave source. The Livingston detector is inverted to account for the detectors' relative orientations [55]. The Pearson correlation coefficient is a widely used statistical measure that provides a measure of the strength of a linear relationship between two variables, in this case, the strain measurements from the two sites. The coefficient ranges from $-1$ to 1, with values close to 1 indicating a strong positive correlation, values close to $-1$ indicating a strong negative correlation, and values close to 0 indicating a lack of correlation between the two variables.

Given two detector streams, the Pearson cross-correlation at time t was computed via

$$\mathbf{P} = \max_{\Delta} \frac{\sum_{k=t-100}^{t+100}(H_k - <H>) \cdot (L_{k-\Delta} - <L>) \cdot (-1)}{\sqrt{\left(\sum_{i=t-100}^{t+100}(H_i - <H>)^2\right)\left(\sum_{j=t-100-\Delta}^{t+100-\Delta}(L_j - <L>)^2\right)}},$$

$$<H> = \frac{1}{200}\sum_{i=t-100}^{t+100} H_i, \quad <L> = \frac{1}{200}\sum_{j=t-100-\Delta}^{t+100-\Delta} L_j$$

The presence of a multiplicative factor of $-1$ serves as the inversion of the Livingston detector due to the relative orientation. The range of $\Delta$ is within the maximum time of flight in units of datapoints, so $\Delta \in [-40, 40]$. In addition to the autoencoder frequency domain features, this yields 5 (autoencoders) $\cdot$2 (features per autoencoder) $+1$ (pearson) $+1$ (frequency-domain correlation) $= 12$ overall features.

## 3.6. Artificial background with timeslides

To create background data that we used a technique called timeslides. This involved temporally shifting the data from one detector relative to another by at least 10 ms, the light travel time between detectors, guaranteeing that there is no astrophysical correlation present. As a background dataset for 3.7, we computed 8 hours worth of timeslides. To evaluate our entire algorithm and report false alarm rates for detections, we computed 1 year worth of timeslides.

## 3.7. Linear combination optimization

To construct a final metric which combines the information from the above 12 features, we opted for a linear combination of those values to produce one final metric value. Given that we are trying to generalize to unknown anomalous signals, opting for a simple linear model aims to reduce any bias towards known signal regions. To find optimal values for the parameters of the linear classifier, we used a simple Linear Support Vector Machine (SVM), which aims to optimize the binary classification of background and signal classes. The classification by a Linear SVM is simply given by $\vec{W}^T\vec{X} + b$, where $\vec{W}$ represents the learned weight vector, of the same dimensionality as the GWAK vector $\vec{X}$, which contains the 12 GWAK features. $b$ represents a bias term. For the background dataset, we used 8 hours worth of timeslides as described in 3.6. Also using this stretch of timeslides, we computed a set of normalization coefficients. For each of the 12 features, these were the mean value and standard deviation of that feature across 8 hours of analysis. These were used as to rescale each feature to mean zero and standard deviation one independently. This helped with training the linear classifier, as while the Pearson correlation coefficient is O(1), the frequency domain coefficients can be O(1000). For the signal dataset, we generated a new dataset, comprising of 6 classes - BBH, SG (64 – 512 Hz), SG ( 512 – 1024Hz), low-frequency white noise burst (40-400Hz), high frequency white noise burst (400-1000Hz), and supernova [56] (85 $M_\odot$ progenitor mass, SFHo equation of state), each with a SNR prior of $U(10, 100)$.

We trained this linear fit for 5000 epochs, with a learning rate of 0.01 and using the Adam optimizer. Fig. 8 shows the learned coefficients for each of the 12 GWAK features. This serves as a sanity check that the impact of each feature on the final metric score is as we expect. Signals are classified via a more negative value, so the features corresponding to the signal autoencoders should have more negative values, as well as the correlation features, as a higher correlation score is more representative of an astrophysical event. On the contrary, autoencoders corresponding to non-astrophysical data (background or glitch) should have more positive values, since a stronger relationship to those classes will be less indicative of the signal. These intuitions are all demonstrated via Fig. 8, so we pass this sanity check.
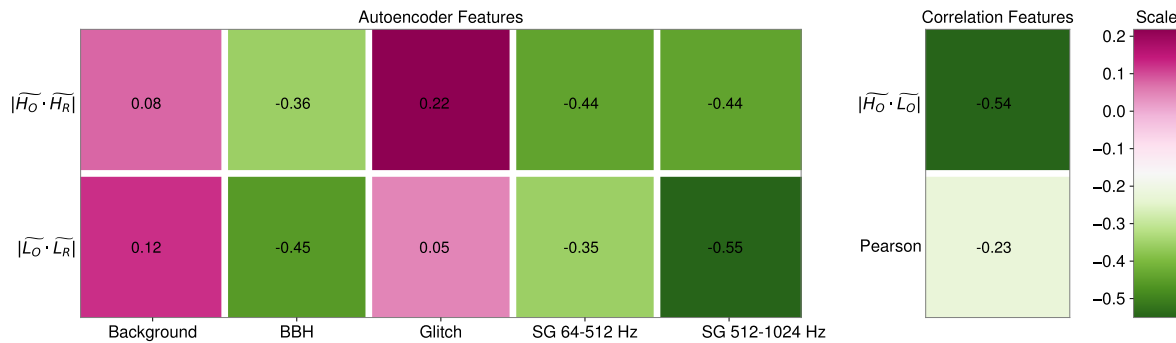
Figure 8: Learned coefficients of a linear combination of the 12 GWAK features. As expected, the features corresponding to the signal autoencoders have more negative values, as well as the correlation features, as a higher correlation score is more representative of an astrophysical event. While background and glitch features have more positive values since a stronger relationship to those classes will be less indicative of the signal.

### 3.8. Smoothing for the final metric

As the GWAK algorithm assigns a single metric value to each $50\,\mathrm{ms}$ window, it does not naively carry sensitivity to signals longer then the $50\,\mathrm{ms}$ window. A longer signal will have multiple evaluation points, but to compute FAR only the lowest metric value is taken. To increase sensitivity by specifying the GWAK algorithm to signals greater than one window in length, we convolved the timeseries of final metric evaluations with uniform kernels of varying size, with the idea being that the sensitivity to a given anomalous signal is maximized when the kernel length is of order the signal length. We present the detection efficiency using this method in Fig. 15.

## 4. Results

This section describes how our proposed methodology can be used to discover anomalous events. Additionally, we evaluate the effectiveness of our semi-supervised approach in detecting several potential sources of GWs, without using information about these signals during the training phase.

### 4.1. Background and glitch mitigation

As the goal of this method is to identify anomalous signals in the background, the mitigation of non-astrophysical data being identified as signal-like improves sensitivity to real signals by reducing the corresponding false alarm rate. In particular, detector artifacts or "glitches" can often pose a problem as they can have signal-like morphology [57,58] as well as possess high SNR in a single detector. This serves as the motivation for training a glitch autoencoder, as it should be able to recognize single detector artifacts reliably. Upon "recognition" of a glitch, via the glitch autoencoder frequency domain features, an automatic down-weight will be given to the event in question. In addition, as

glitches are uncorrelated between detectors, occurring locally, the use of both correlation features also helps to mitigate false alarms caused by glitches. As there is no correlation between observations in one detector stream and another within the maximum light travel distance during a glitch, those features will correspondingly have small values, and similarly down-weight the glitch to have a less signal-like score.

## 4.2. Anomaly metric

We employ the linear combination method outlined in Sec. 3.7, which includes the two frequency-domain features per autoencoder, the frequency-domain correlation, and the Pearson cross-correlation presented in Sec. 3.5. The example of a full GWAK pipeline is shown in Fig. 9.



Figure 9: Example of metric calculation on a hypothetical event. The event is reconstructed with each of the 5 pre-trained autoencoders. Pearson and frequency domain correlation are computed on the given input and then each of the values is multiplied with a corresponding coefficient which arises from the pre-trained linear metric. The sum of all the features multiplied by their coefficients is referred to as the final metric and is used to make a decision on if the event is signal-like.

The resulting coefficients of the linear classifier, physically describing a hyperplane, are then used to project points from the 12-dimensional feature space down to a one-dimensional space via the dot product, or geometrically the perpendicular distance from that point to the classifying hyperplane. This one-dimensional real number is our final metric value. Since our classifier was trained to predict 0 for signals and 1 for backgrounds, a more negative final metric value corresponds to a more "signal-like" or louder input, whereas a less negative/positive metric value indicates background or no signal of interest. Once we compute the final metric value for a hypothetical signal, we would like to know the corresponding false alarm rate. This corresponds to the frequency that a non-astrophysical input (glitch or otherwise) from the gravitational-wave detectors would lead to a trigger of the same significance as the hypothetical signal. The 12

features multiplied by the corresponding coefficient and grouped by a corresponding autoencoder are shown in Fig. 10. We show each of the five data classes in this space to highlight how strongly the signals are separated and grouped in this new, learned GWAK space.

## 4.3. Evaluation on Core-Collapse Supernova

We use existing core-collapse supernovae simulations to see how our approach extends to anomalous signals which GWAK has not seen during training. We use [56] (85 $M_\odot$
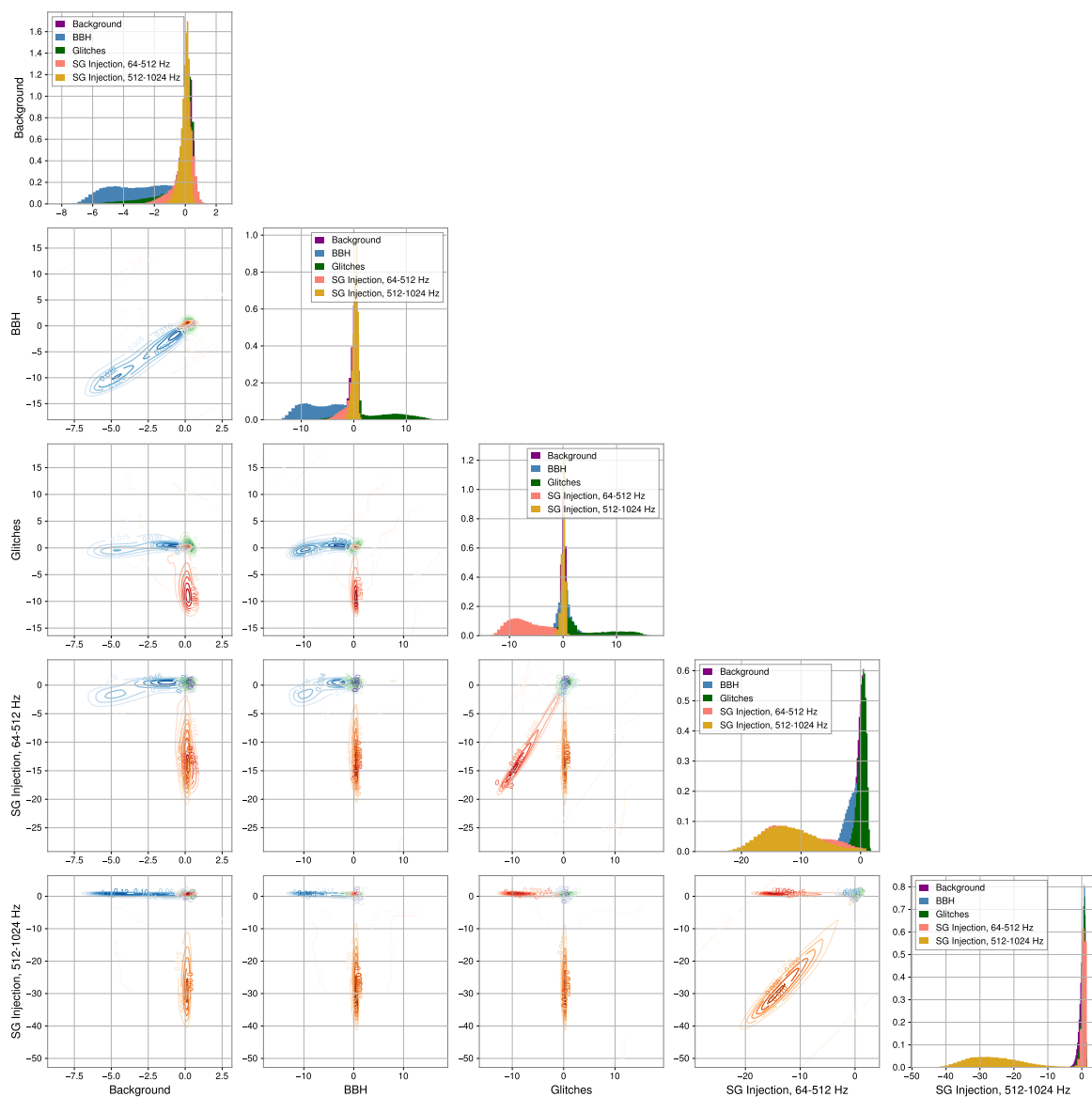


Figure 10: Trained GWAK space on 5 signal classes: BBH, background, glitch, and low and high frequency Sine-gaussian; demonstrating that regions of the GWAK space correspond to different event classes

progenitor mass, SFHo equation of state). On the top left of Fig. 11, the evaluation of GWAK axes and Pearson correlation with time and on the top right total metric value and FAR are shown as an example of the algorithm's "reaction" to unseen signals. Both BBH and sine-Gaussian losses drop at the time of the signal, which indicates that the strain at that moment is more signal-like. The Pearson correlation increases indicating a strong correlation between the two detector sides. The FAR at the moment of the event drops to a level consistent with one or fewer events per month, which means that even with strong trigger restrictions, detection of that type of event would be possible with our new proposed algorithm.

The bottom plot in Fig. 11 (left) reflects the scan of minimum metric value for different signal-to-noise ratio ranges of the injected core-collapse supernova models. As expected, with the higher signal-to-noise ratio, the total metric is lower resulting in a lower FAR. If the trigger bandwidth would allow for up to 1 false event per hour being triggered, GWAK would be able to detect core-collapse supernova events with ∼22 SNR.

### 4.4. Evaluation on White Noise Bursts

Furthermore, we assessed the performance of our method on white-noise bursts, which are signals characterized by the presence of $h^+$ and $h^\times$ polarizations that are independent time series of Gaussian noise, which is whitened over a specific frequency range and multiplied by a sigmoid envelope. The bandwidth of each injected signal is selected uniformly and randomly from a range spanning 40–400 Hz and 400–1000 Hz. The duration of the signal is chosen to be 0.1 s. Theoretically, these would be the most difficult signals to detect with our algorithm, as their lack of distinctive morphology would render the signal autoencoder features useless. However, as shown in Fig. 11, the SG autoencoders were able to generalize to the WNBs.

To evaluate the performance of our algorithm, we generated these signals with SNRs uniformly distributed between 10 and 100. The average final metric value and corresponding standard deviation for various SNR ranges are shown in Fig. 11 (right) and the lines corresponding to different false alarm rates.

The demonstration of the GWAK algorithm from strain to final metric is shown in Fig. 11. Starting with the whitened strain, we split up our data into 200 datapoint windows with a step of 5 datapoints between windows. For each window, the 10 autoencoder features are computed, along with the frequency domain correlation, and Pearson correlation, as described in Sec. 3.5. Using the learned SVM weights as shown in Fig. 8, we reduce the 12-dimensional GWAK space to 7 dimensions, and show those values at each timestep in the second panel. For the correlation features, this is simply done by multiplying the value by the corresponding weight. For the autoencoder features, the same is done, but the $|\widetilde{H_O} \cdot \widetilde{H_R}|$ and $|\widetilde{L_O} \cdot \widetilde{L_R}|$ values are combined into a single value via addition. Finally, all of those weighted features are summed yielding the final metric value, which is shown in the third panel, along with various false alarm rate thresholds.

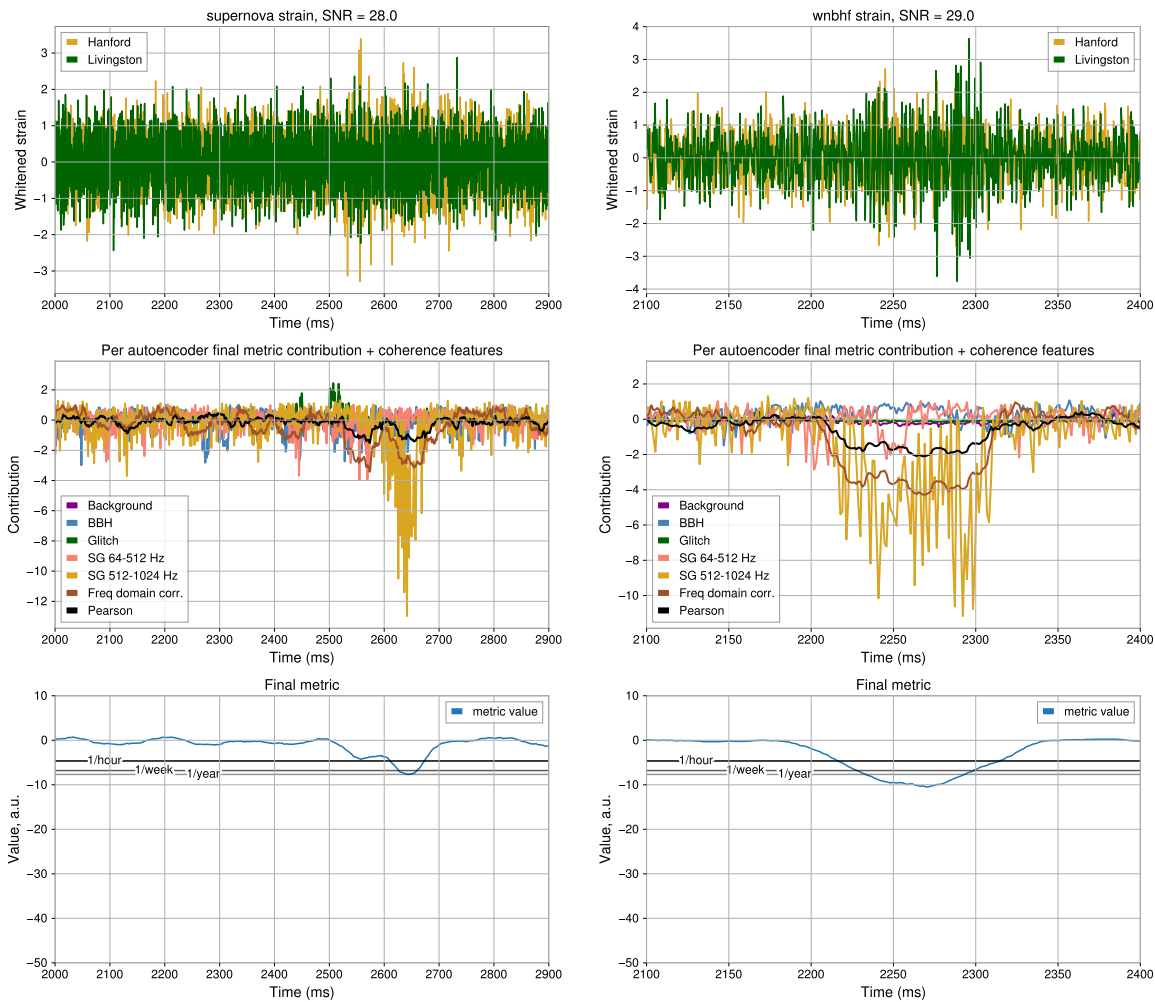The method performance on different signals at various SNR values is shown in

Figure 11: Strain (top), GWAK metric response (middle) and final metric response (bottom) for WNB (right) and Supernova (left). The evaluation of GWAK axes and Pearson correlation with time and on the top right total metric value and FAR are shown as an example of the algorithm's "reaction" to unseen signals. Both BBH and sine-Gaussian losses drop at the time of the signal, which indicates that the strain at that moment is more signal-like. The Pearson correlation increases indicating a strong correlation between the two detector sides. At the moment of the event, the FAR drops to a level consistent with or below one event per year, which means that even with strong trigger restrictions, detection of that type of event would be possible with our new proposed algorithm.

Fig. 12. For each signal type, we generate 10,000 waveforms using the SNR prior $U[5, 50]$. We then run the full GWAK algorithm and obtain the minimum metric value achieved in each injection, as shown in Fig. 11. We then group the injections into SNR bins of width 5, and show the average metric value as well as the $1\sigma$ range for each bin, via the solid line and filled region respectively. The final metric values corresponding to certain false alarm thresholds are also shown. We first see that the three training classes
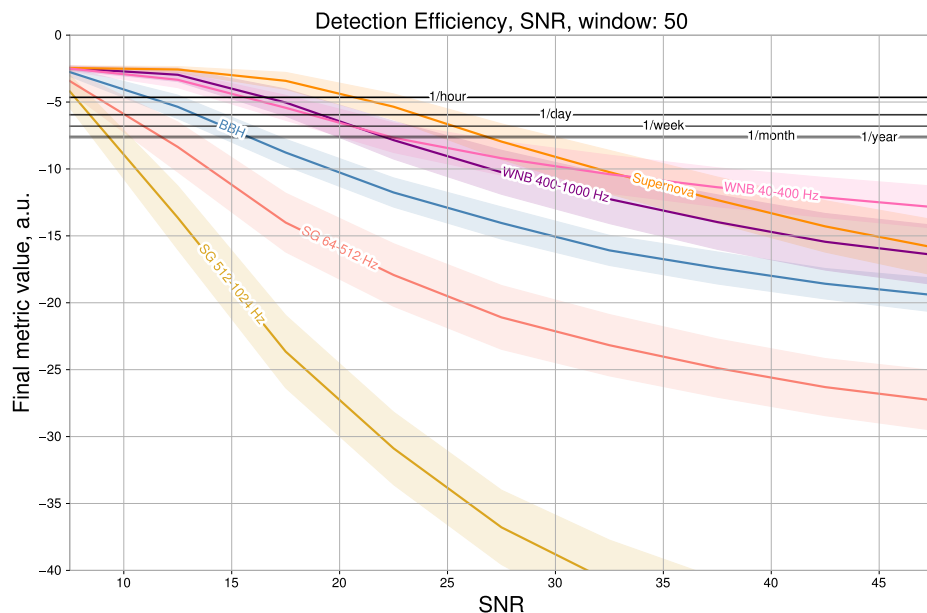
Figure 12: The final metric as a function of SNR for GWAK axes training signals, BBH (blue), SG 64–512 Hz (yellow), SG 512–1024 Hz (salmon) and for potential anomalies, WNB 40–400 Hz (pink), WNB 400–1000 Hz (purple), and Supernova (orange). The black lines of varied width correspond to different FARs, from the FAR of 1 per hour to 1 per year. For each of the lines, the events that are below that line would be detected. As expected, the signals used to train the GWAK axes are on average detected better, eg more events with a smaller SNR are detected for a given FAR threshold.

- BBH, SG 64-512 Hz, SG 512-1024 Hz are most efficiently detected, which is expected as they have specific autoencoder models. Moving on to the anomalous signals - WNB 40-400 Hz, WNB 400-1000 Hz, and Supernova, we see that they are not detected as readily as the training signals, but are still able to achieve satisfactory false alarm rates around 1-2/month around 20-25 SNR.

We show the detection efficiency using a receiver operating characteristic (ROC) curve in Fig. 13 in order to compare it with other ML techniques. Here, we pick a fixed false alarm rate threshold of 1/year and compute what fraction of injections, for each signal, at each SNR, have detection statistics below the threshold. Similar to Fig. 12, we see that the signals corresponding to training classes can be more readily detected at lower SNRs, and the anomalous signals require slightly higher SNRs to be detected. This graph can be compared to the one presented by MLy [23], while for the BBH and SGs we achieve similar performance, the efficiencies for WNBs and Supernovae are lower in our case. This is expected since the MLy algorithm was trained in a supervised manner, using WNBs during the training, while we only used those signals for finding the linear coefficients of the final metric but not as the GWAK axes.
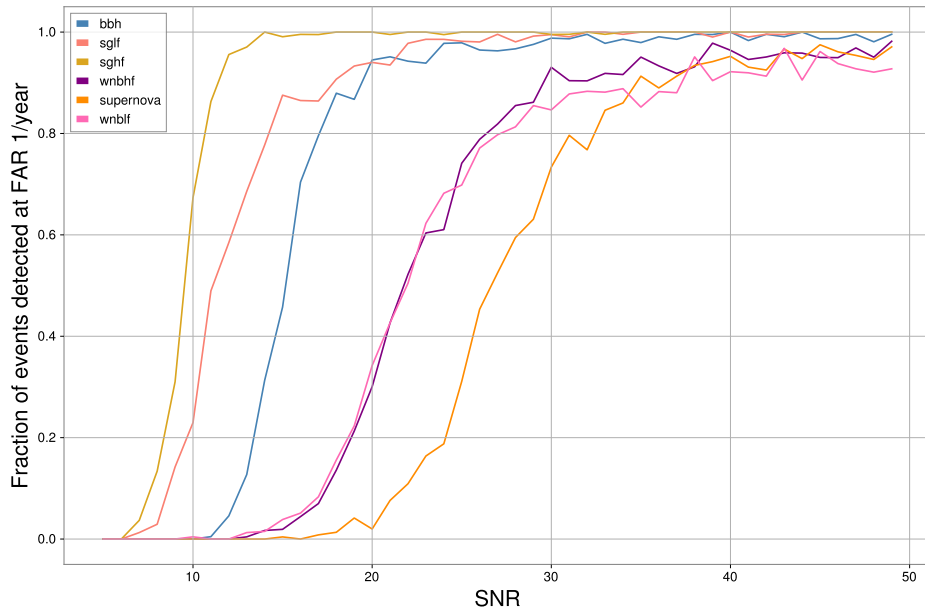
Figure 13: Detection efficiency as a function of SNR for GWAK axes training signals, and for potential anomalies. We see that the signals corresponding to training classes can be more readily detected at lower SNRs, and the anomalous signals require slightly higher SNRs to be detected.

## 4.5. Comparison to a supervised search

To quantify the loss of efficiency of the unsupervised GWAK method in comparison to a supervised search, we perform the following study. Firstly, we use pre-trained GWAK axes for the BBH search by using a small, dense network instead of a simple linear combination for the final metric. We train this network in a supervised manner, using BBH as the signal class and timeseries as the background class. In this way, we can quantify how much signal efficiency is lost when using a general linear combination instead of overfitting on a specific signal. The results are shown in Fig. 14. We observe that BBH detection efficiency surpasses that achieved with the linear final metric. However, as anticipated, the detection efficiency for all other signals significantly decreases. We omitted the use of a smoothing window, as it was determined to be the most efficient for BBH-supervised search. Thus, we must compare it to the BBH ROC without the application of a smoothing window, as depicted in the Fig. 15.

We demonstrate that, in general, enhancing the detection efficiency for a specific signal through supervised training is feasible. However, this improvement often incurs a noticeable decline in performance for other signals. Given our objective to develop an algorithm capable of detecting unknown signals, we lack the necessary information to train it in a supervised manner. Nevertheless, in follow-up works, we may consider exploring the adoption of a more sophisticated final metric function, though we must exercise caution to prevent overfitting to the signals employed in optimizing this metric.
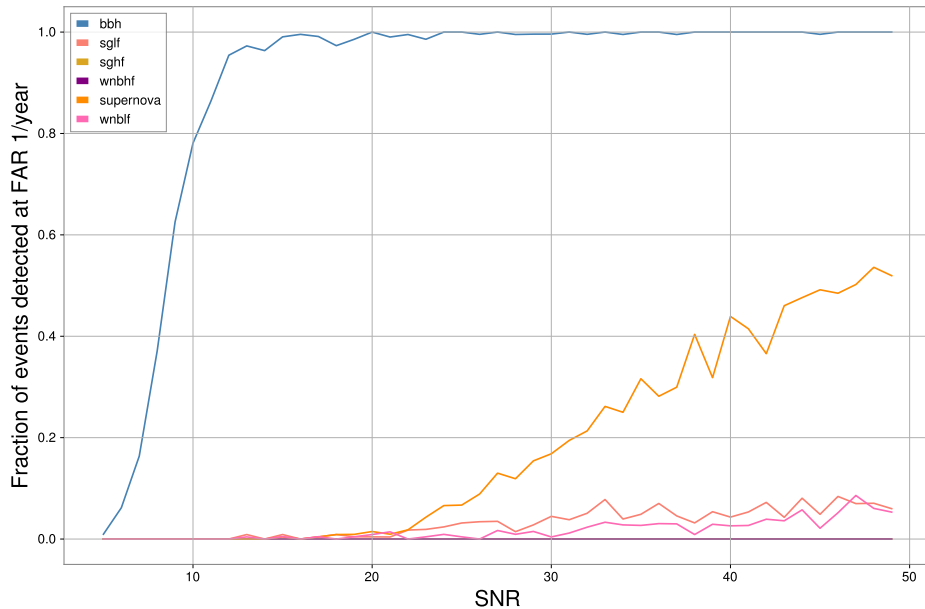
Figure 14: Detection efficiency for BBH and other signals and anomalies obtained using the final metric trained in a supervised manner on BBH signals pre-processed by the GWAK algorithm. This demonstrates that the GWAK method can approach a supervised search when given this specific task.

### 4.6. Comparison of Different Smoothing Windows

In Fig. 15, we illustrate that the optimization of detection efficiency, both for known and anomalous signals, occurs when employing smoothing kernels with sizes approximately equal to the signal length. These findings affirm that utilizing smoothing within the final metric space is an effective strategy to adapt the GWAK algorithm to diverse signal lengths

## 5. Conclusions

In this study, we utilized the Gravitational Wave Anomalous Knowledge (GWAK) method to identify anomalies in datasets acquired by ground-based GW observatories. The GWAK method relies on the notion of introducing alternative signal priors that capture some of the salient features of new physics signatures, enabling the restoration of sensitivity even when the alternative signal is incorrect. We separately trained five unsupervised autoencoders on a dataset consisting of normal background noise, glitches, and a collection of simulated signals that incorporate the physical characteristics of a potential new physics signature. We then established a 12-dimensional GWAK space, comprising two reconstruction losses for each of the detector sides for each of the five autoencoders and two features representing the correlation between the detector sides and searched regions of the GWAK space for anomalous signals. We then combined all
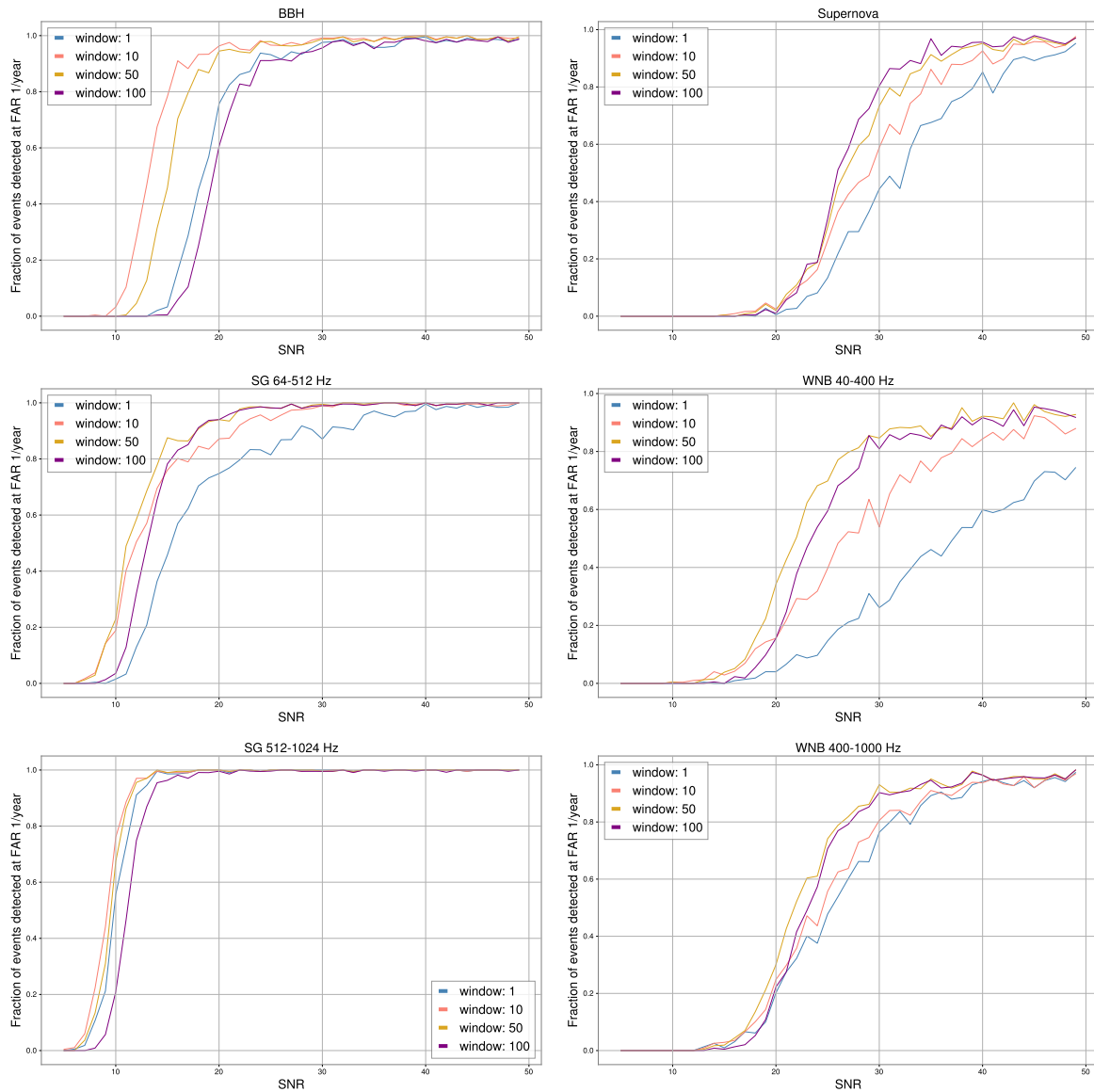
Figure 15: Different smoothing windows for training signals (left column) and anomalous signals (right column).

the 12 features into one final metric by multiplying with a corresponding coefficients from Fig. 8 and summing the result.

Our findings indicate that the GWAK method efficiently detected anomalies in the GW datasets. In particular, unmodelled sources like core-collapse supernovae and white noise bursts. Additionally, the GWAK method could differentiate signal-like anomalies from anomalous events, such as those resulting from detector glitches.

Our proposed method demonstrates promising results, detecting these sources with high accuracy and without prior knowledge of their characteristics. These findings underscore the potential value of our approach in detecting new and unexpected sources of GW signals, while simultaneously reducing the dependence on labeled training data.

In addition to serving as an unsupervised search, the machine-learning based approach allows for the implementation of the GWAK algorithm as a low-latency search tool. By detecting anomalies rapidly, alerts can be sent to electromagnetic telescopes for follow-up.

In future work, there are many potential avenues to explore. One is using the recreation as a denoising tool instead of just a detection statistic, allowing for rapid parameter estimation. This is especially important with electromagnetic follow-up, as the telescopes need information on the source location for detailed observation. On the detection side of things, an idea is to use normalizing flows to learn the high dimensional manifolds on which signals lie and use the probability value as a score alternate to the less intuitive engineered autoencoder frequency domain features. There is also potential to improve the search by changing the way that the 12 GWAK features are reduced to a single value. While the linear fit served as a simple and explainable method and did not really possess the potential to overfit the known signals and therefore decrease sensitivity to anomalies, it suffers the drawback that potentially useful patterns between features are not exploited. A simple example is with the features $|\widetilde{H_O} \cdot \widetilde{H_R}|$ and $|\widetilde{L_O} \cdot \widetilde{L_R}|$ for a single autoencoder. For a glitch event, you would expect only one of these values to increase, corresponding to the detector in which the glitch occurred, so you would expect an asymmetry between these features for a non-astrophysical event. The opposite is true for a BBH event, for example, as the fact that it is present in both detector channels means that there should be symmetry in the BBH autoencoder features. While this is just one intuitive example, more complex relations could certainly exist. Yet another area to explore is modifying the network architecture to allow for a longer signal length to generalize the algorithm to various signal durations.

To conclude, the GWAK method displays potential as a powerful tool for detecting anomalies in GW datasets and has the potential to enhance the performance of GW anomaly detection systems.

**Acknowledgments**

## 6. Appendix

### *6.1. Backgrounds training curve*

For completeness, in Fig. 16 we show the training and validation losses for the background and glitch autoencoders.
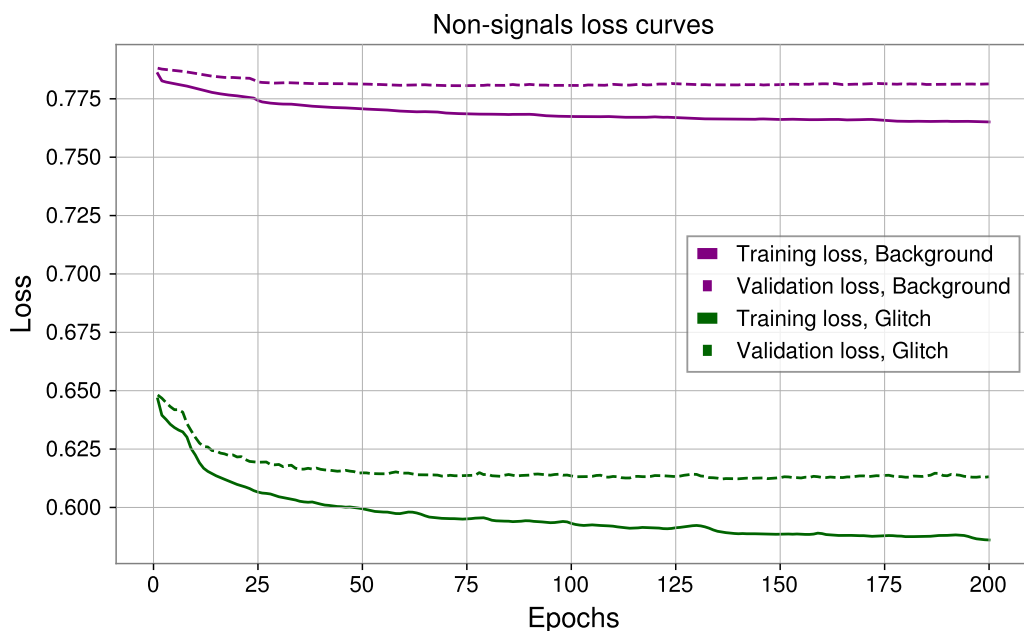


Figure 16: Autoencoder training and validation losses for background classes. The training/validation losses are in bold/dashed for background (in purple) and for glitches (in green).

## 6.2. Recreations

For completeness, in Fig. 17, 18, 19 and 20 we show recreation plots for the low frequency SG, high frequency SG, glitches and background correspondingly. For each of the dataset types we can see that the corresponding autoencoder is the best in reconstructing the input, as expected.
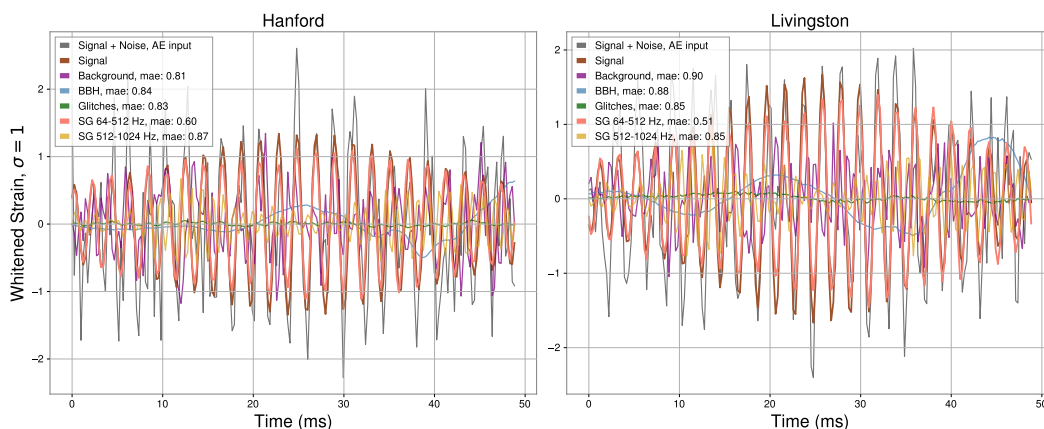


Figure 17: Example of recreation plots obtained with all five pre-trained autoencoders on low frequency sine-gaussian dataset.
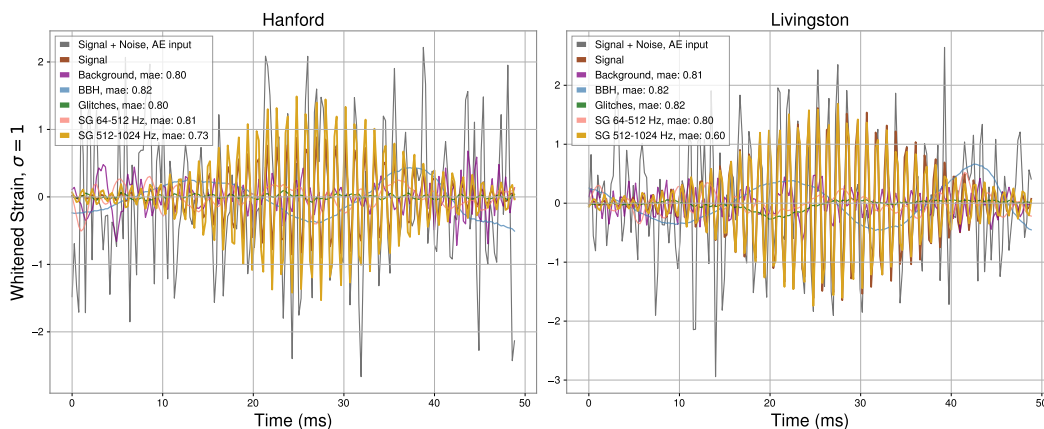


Figure 18: Example of recreation plots obtained with all five pre-trained autoencoders on high frequency sine gaussian dataset.
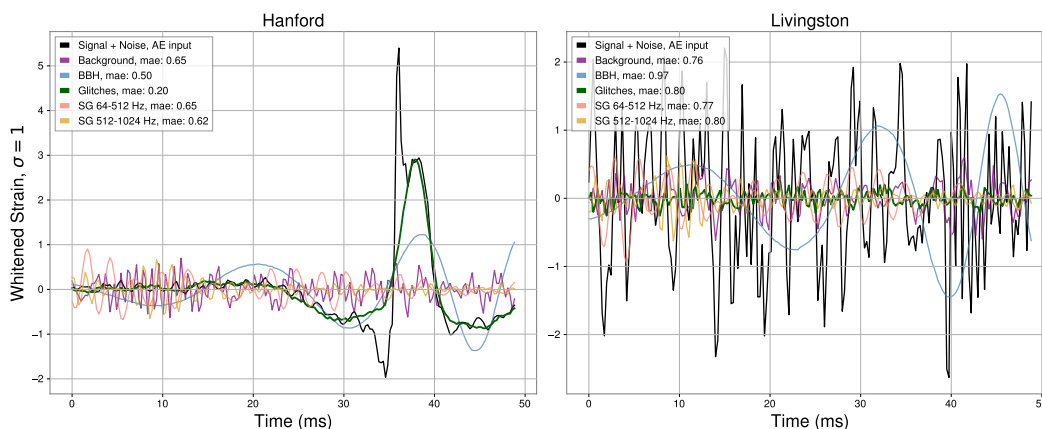
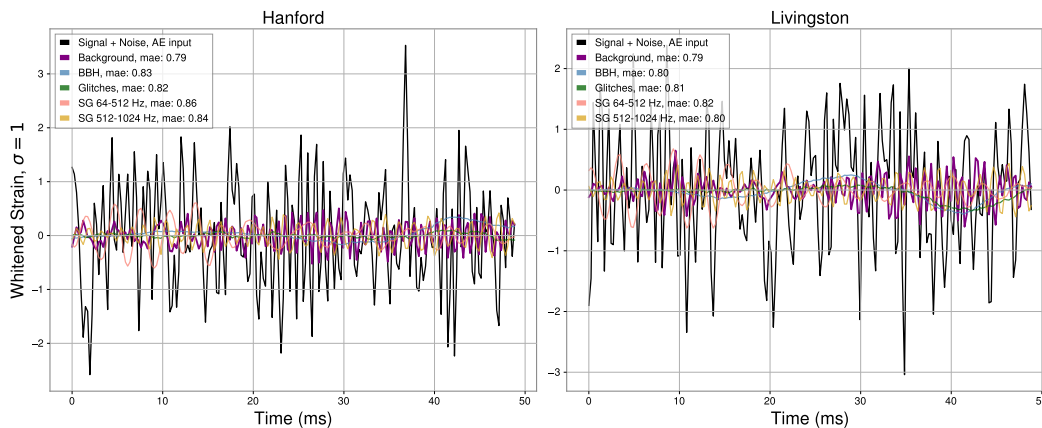Figure 19: Example of recreation plots obtained with all five pre-trained autoencoders on glitch dataset.



Figure 20: Example of recreation plots obtained with all five pre-trained autoencoders on background dataset.

## References

[1] B. P. Abbott et al. Observation of gravitational waves from a binary black hole merger. *Physical Review Letters*, 116(6):061102, 2016.

[2] The LIGO Scientific Collaboration and J Aasi et al. Advanced ligo. *Classical and Quantum Gravity*, 32(7):074001, mar 2015.

[3] F Acernese et al. Advanced virgo: a second-generation interferometric gravitational wave detector. *Classical and Quantum Gravity*, 32(2):024001, dec 2014.

[4] T. Akutsu et al. Overview of KAGRA: Detector design and construction history. *PTEP*, 2021(5):05A101, 2021.

[5] B. P. Abbott, R. Abbott, and T. D. et al. Abbott. Gwtc-1: A gravitational-wave transient catalog of compact binary mergers observed by ligo and virgo during the first and second observing runs. *Phys. Rev. X*, 9:031040, Sep 2019.

[6] R. Abbott et al. GWTC-2: Compact Binary Coalescences Observed by LIGO and Virgo During the First Half of the Third Observing Run. *Phys. Rev. X*, 11:021053, 2021.

[7] R. Abbott et al. GWTC-3: Compact Binary Coalescences Observed by LIGO and Virgo During the Second Part of the Third Observing Run. 11 2021.

[8] B. P. Abbott et al. Gw151226: Observation of gravitational waves from a 22-solar-mass binary black hole coalescence. *Physical Review Letters*, 116(24):241103, 2016.

[9] B. P. Abbott et al. Gw170104: Observation of a 50-solar-mass binary black hole coalescence at redshift 0.2. *Physical Review Letters*, 118(22):221101, 2017.

[10] B. P. Abbott et al. Gw170814: A three-detector observation of gravitational waves from a binary black hole coalescence. *Physical Review Letters*, 119(14):141101, 2017.

[11] B. P. Abbott et al. Gw170817: Observation of gravitational waves from a binary neutron star inspiral. *Physical Review Letters*, 119(16):161101, 2017.

[12] B. P. Abbott et al. Gwtc-1: A gravitational-wave transient catalog of compact binary mergers observed by ligo and virgo during the first and second observing runs. *Physical Review X*, 9(3):031040, 2019.

[13] Ernazar Abdikamalov, Giulia Pagliaroli, and David Radice. Gravitational waves from core-collapse supernovae. In *Handbook of Gravitational Wave Astronomy*, pages 1–37. Springer Singapore, 2021.

[14] M B Hindmarsh and T W B Kibble. Cosmic strings. *Reports on Progress in Physics*, 58(5):477–562, may 1995.

[15] Thibault Damour and Alexander Vilenkin. Gravitational wave bursts from cosmic strings. *Physical Review Letters*, 85(18):3761–3764, oct 2000.

[16] Eric Braaten and Hong Zhang. Axion stars, 2018.

[17] G. Franciolini. Primordial black holes: from theory to gravitational wave observations, 2021.

[18] Yu N Eroshenko. Gravitational waves from primordial black holes collisions in binary systems. *Journal of Physics: Conference Series*, 1051(1):012010, jul 2018.

[19] B. Allen et al. FINDCHIRP: An Algorithm for detection of gravitational waves from inspiraling compact binaries. *Phys. Rev. D*, 85:122006, 2012.

[20] S. Klimenko et al. Method for detection and reconstruction of gravitational wave transients with networks of advanced detectors. *Phys. Rev. D*, 93:042004, Feb 2016.

[21] S. Klimenko, S. Mohanty, M. Rakhmanov, and G. Mitselmakher. Constraint likelihood analysis for a network of gravitational wave detectors. *Phys. Rev. D*, 72:122002, Dec 2005.

[22] Ryan Lynch, Salvatore Vitale, Reed Essick, Erik Katsavounidis, and Florent Robinet. Information-theoretic approach to the gravitational-wave burst detection problem. *Phys. Rev. D*, 95:104046, May 2017.

[23] Vasileios Skliris, Michael R. K. Norman, and Patrick J. Sutton. Real-time detection of unmodelled gravitational-wave transients using convolutional neural networks, 2022.

[24] Sang Eon Park, Dylan Rankin, Silviu-Marian Udrescu, Mikaeel Yunus, and Philip Harris. Quasi anomalous knowledge: searching for new physics with embedded knowledge. *Journal of High Energy Physics*, 2021(6), jun 2021.

[25] Paul T. Baker, Sarah Caudill, Kari A. Hodge, Dipongkar Talukder, Collin Capano, and Neil J. Cornish. Multivariate Classification with Random Forests for Gravitational Wave Searches of Black Hole Binary Coalescence. *Phys. Rev. D*, 91(6):062004, 2015.

[26] D. George and E.A. Huerta. Deep Neural Networks to Enable Real-time Multimessenger Astrophysics. *Phys. Rev. D*, 97(4):044039, 2018.

[27] Shasvath J. Kapadia, Thomas Dent, and Tito Dal Canton. Classifier for gravitational-wave inspiral signals in nonideal single-detector data. *Phys. Rev. D*, 96(10):104015, 2017.

[28] D. George and E.A. Huerta. Deep Learning for Real-time Gravitational Wave Detection and Parameter Estimation: Results with Advanced LIGO Data. *Phys. Lett. B*, 778:64–70, 2018.

[29] Hunter Gabbard et al. Matching matched filtering with deep networks for gravitational-wave astronomy. *Phys. Rev. Lett.*, 120(14):141103, 2018.

[30] Andrew L. Miller et al. How effective is machine learning to detect long transient gravitational waves from neutron stars in a real search? *Phys. Rev. D*, 100(6):062005, 2019.

[31] Shreejit Jadhav, Nikhil Mukund, Bhooshan Gadre, Sanjit Mitra, and Sheelu Abraham. Improving significance of binary black hole mergers in Advanced LIGO data using deep learning:

Confirmation of GW151216. *Phys. Rev. D*, 104(6):064051, 2021.

[32] E. A. Huerta et al. Accelerated, scalable and reproducible AI-driven gravitational wave detection. *Nature Astron.*, 5(10):1062–1068, 2021.

[33] Letian Jiang and Yuan Luo. Convolutional transformer for fast and accurate gravitational wave detection. In *2022 26th International Conference on Pattern Recognition (ICPR)*, pages 46–53, 2022.

[34] Chayan Chatterjee, Linqing Wen, Foivos Diakogiannis, and Kevin Vinsen. Extraction of binary black hole gravitational wave signals from detector data using deep learning. *Physical Review D*, 104(6), sep 2021.

[35] Damon Beveridge, Linqing Wen, and Andreas Wicenec. Detection of binary black hole mergers from the signal-to-noise ratio time series using deep learning, 2023.

[36] Rich Ormiston, Tri Nguyen, Michael Coughlin, Rana X. Adhikari, and Erik Katsavounidis. Noise Reduction in Gravitational-wave Data via Deep Learning. *Phys. Rev. Res.*, 2(3):033066, 2020.

[37] P. Bacon, A. Trovato, and M. Bejger. Denoising gravitational-wave signals from binary black holes with dilated convolutional autoencoder. 5 2022.

[38] Muhammed Saleem et al. Demonstration of Machine Learning-assisted real-time noise regression in gravitational wave detectors. 6 2023.

[39] Alec Gunny, Dylan Rankin, Jeffrey Krupa, Muhammed Saleem, Tri Nguyen, Michael Coughlin, Philip Harris, Erik Katsavounidis, Steven Timm, and Burt Holzman. Hardware-accelerated Inference for Real-Time Gravitational-Wave Astronomy. *Nature Astron.*, 6(5):529–536, 2022.

[40] Chung-Hao Liao and Feng-Li Lin. Deep generative models of gravitational waveforms via conditional autoencoder. *Phys. Rev. D*, 103:124051, Jun 2021.

[41] Sivaramakrishnan Sankarapandian and Brian Kulis. $\beta$-Annealed Variational Autoencoder for glitches. 7 2021.

[42] F. Morawski et al. Anomaly detection in gravitational waves data using convolutional autoencoders. *Machine Learning: Science and Technology*, 2021.

[43] Eric Moreno. Source-Agnostic Gravitational-Wave Detection with Recurrent Autoencoders: BBH Dataset, December 2021.

[44] Vasileios Skliris, Michael R. K. Norman, and Patrick J. Sutton. Real-Time Detection of Unmodelled Gravitational-Wave Transients Using Convolutional Neural Networks. 9 2020.

[45] Patrick J Sutton, Gareth Jones, Shourov Chatterji, Peter Kalmus, Isabel Leonor, Stephen Poprocki, Jameson Rollins, Antony Searle, Leo Stein, Massimo Tinto, and Michal Was. X-pipeline: an analysis package for autonomous gravitational-wave burst searches. *New Journal of Physics*, 12(5):053034, may 2010.

[46] B. P. Abbott, R. Abbott, and Others. All-sky search for short gravitational-wave bursts in the first advanced ligo run. *Phys. Rev. D*, 95:042003, Feb 2017.

[47] J. Aasi et al. Advanced LIGO. *Class. Quant. Grav.*, 32:074001, 2015.

[48] Florent Robinet, Nicolas Arnaud, Nicolas Leroy, Andrew Lundgren, Duncan Macleod, and Jessica McIver. Omicron: a tool to characterize transient noise in gravitational-wave detectors. *SoftwareX*, 12:100620, 2020.

[49] L. Santamaría, F. Ohme, P. Ajith, B. Brügmann, N. Dorband, M. Hannam, S. Husa, P. Mösta, D. Pollney, C. Reisswig, E. L. Robinson, J. Seiler, and B. Krishnan. Matching post-newtonian and numerical relativity waveforms: Systematic errors and a new phenomenological model for nonprecessing black hole binaries. *Phys. Rev. D*, 82:064016, Sep 2010.

[50] Sascha Husa, Sebastian Khan, Mark Hannam, Michael Pürrer, Frank Ohme, Xisco Jiménez Forteza, and Alejandro Bohé. Frequency-domain gravitational waves from nonprecessing black-hole binaries. i. new numerical waveforms and anatomy of the signal. *Phys. Rev. D*, 93:044006, Feb 2016.

[51] Sebastian Khan, Sascha Husa, Mark Hannam, Frank Ohme, Michael Pürrer, Xisco Jiménez Forteza, and Alejandro Bohé. Frequency-domain gravitational waves from nonprecessing black-hole binaries. ii. a phenomenological model for the advanced detector era. *Phys. Rev. D*, 93:044007,

Feb 2016.

[52] A. Nitz et al. gwastro/pycbc: Pycbc release v1.16.9, August 2020.

[53] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017.

[54] Pearson Karl. (58240–242), 1895.

[55] B. P. Abbott et al. Observation of gravitational waves from a binary black hole merger. *Physical Review Letters*, 116(6), feb 2016.

[56] Jade Powell, Bernhard Müller, and Alexander Heger. The final core collapse of pulsational pair instability supernovae. *Monthly Notices of the Royal Astronomical Society*, 503(2):2108–2122, mar 2021.

[57] B. P. Abbott et al. Characterization of transient noise in Advanced LIGO relevant to gravitational wave signal GW150914. *Class. Quant. Grav.*, 33(13):134001, 2016.

[58] Miriam Cabero et al. Blip glitches in Advanced LIGO data. *Class. Quant. Grav.*, 36(15):15, 2019.