

University of Massachusetts Dartmouth

Department of Mathematics

An Exploration of Statistical Methods in Determining the Gaussianity of LIGO Detector Data

A Data Science project

by

Ashish Thomas Mathew
(02101455)

Project Advisors:

Dr. Sarah Caudill

University of Massachusetts Dartmouth

Dr. Melissa Lopez

National Institute for Subatomic Physics (Nikhef)

Submitted in Partial Fulfillment of the
Requirements for the Degree of
Master of Science

May, 2025

Abstract

The Advanced Laser Interferometer Gravitational-wave Observatory (LIGO) is one of the most sophisticated instruments ever built, capable of measuring motion 10,000 times smaller than a proton. It was the first device to detect a gravitational wave (GW) event. While analyzing the GW data from LIGO, the detector noise is approximated to be stationary and Gaussian. However, due to its high sensitivity, a common issue it faces is its susceptibility to glitches: transient, non-Gaussian noise bursts that contaminate the time series data obtained from the detector. This project aims to develop a tool to check if the distribution of the amplitudes of the detector noise is Gaussian or not as a way to assess whether it is clean up to a defined threshold, hence improving the quality of the analysis performed. Taking the null hypothesis that the distribution is Gaussian, we study statistical approaches of normality testing, namely the Shapiro-Wilk test, the Kolmogorov-Smirnov test, and the Anderson-Darling test, to test this null hypothesis against the alternative of the distribution being non-Gaussian. We found that of the three tests, the Shapiro-Wilk test performed the best on the complete frequency range in detecting non-Gaussian samples with an accuracy of 73.86% and a precision of 91.67%. In comparison, the Kolmogorov-Smirnov test had an accuracy of 64.43% and a precision of 100%, while the Anderson-Darling test had an accuracy of 63.9% and a precision of 100% when detecting non-Gaussianity.

1 Introduction

Nearly a century after Einstein's prediction of the existence of gravitational waves (GWs) in 1916, the first direct gravitational wave detection was achieved by the Advanced Laser Interferometer Gravitational-Wave Observatory (aLIGO) [1] and Virgo [2] interferometers during the binary black hole merger event GW150914 on September 14, 2015. The LIGO Livingston (L1) interferometer underwent several upgrades following this run - effectively a massive redesign - enhancing its sensitivity by 15% to 25% [3]. This improvement was evident during the O2 observing run during which L1, in conjunction with its Hanford (H1) counterpart and Virgo (V1), detected eleven new gravitational wave signals [4]. Subsequently, during the O3 run, all three detectors operated at their best possible sensitivity, leading to the first single-detector GW detection, GW190425, achieved by LIGO Livingston [5]. There have been over 90 GW events recorded at a high level of confidence since the LIGO-Virgo collaboration's inception to date [6].

A consequence of these detectors' sensitivities is their susceptibility to noise. During observation runs, various noise sources, such as seismic noise, suspension thermal noise, and sensing noise, affect the data collected by the interferometers [7]. Among these, the most problematic are **glitches**: transient events caused by non-astrophysical phenomena such as anthropogenic noise, weather conditions, or instrument malfunctions [8, 9]. Glitches manifest as localized bursts of excess power in interferometer time series data and often do not have well-defined sources. They can occur at energy levels and frequencies that overlap with GW signals, thereby mimicking them and increasing the number of false positive detections. During the first half of the third observing run (O3a) H1 and L1 recorded glitch rates of 0.29 min^{-1} and 1.10 min^{-1} respectively, which rose to 0.32 min^{-1} and 1.17 min^{-1} during the second half (O3b) [6]. A notable example of glitches posing such an issue was during the detection of GW170817 [10], where instrumental noise transients during the signal's detection time complicated its identification and subsequent analysis.

This project aims to develop a tool to check if the amplitudes of detector noise in a sample of data from the detector show the presence of a glitch through tests of **Gaussianity**. Gaussianity

is the property of a sample distribution of data being normally distributed, also known as a Gaussian distribution. A normal/Gaussian distribution is a continuous probability distribution for a random variable where most of the data points are located around the mean, with lesser amounts of data further away. The Gaussian distribution is symmetric about the mean μ , which is equal to zero, with the standard deviation σ being equal to 1. This distribution is important as it can approximate many real world phenomena, such as the distribution of heights and weights of a country's population. It is also the basis for many statistical tests and methods. Any variation from this distribution is considered non-Gaussian. Figure 1 shows an example of a Gaussian distribution.

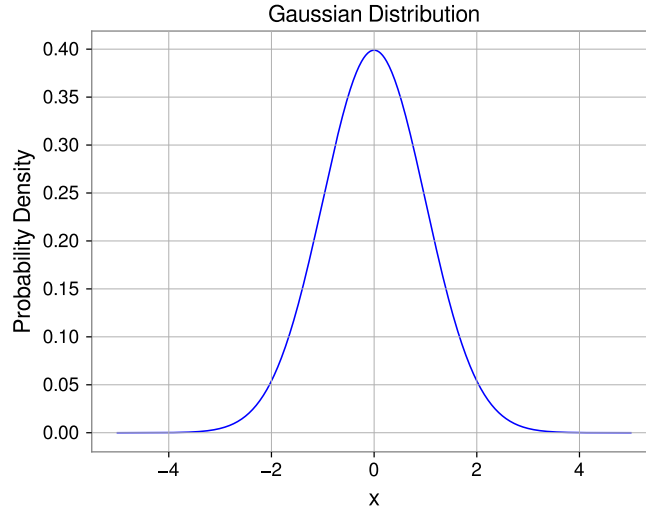


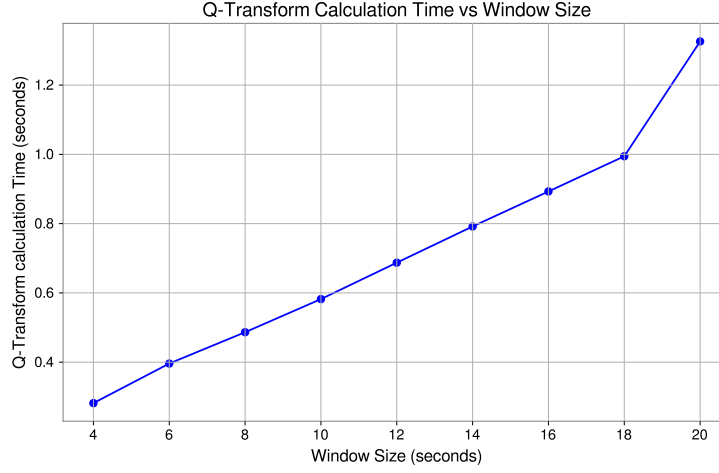
Figure 1: Example of a Gaussian distribution

All the noise sources in a detector collectively produce time series signals which can be treated as stochastic processes with their corresponding joint probability distributions and statistical properties [8]. In gravitational wave data analysis the noise floor of the detectors is approximated to be stationary and Gaussian. In the event of a glitch, the noise has some structure with a high signal-to-noise ratio (SNR), making it *non-Gaussian*.

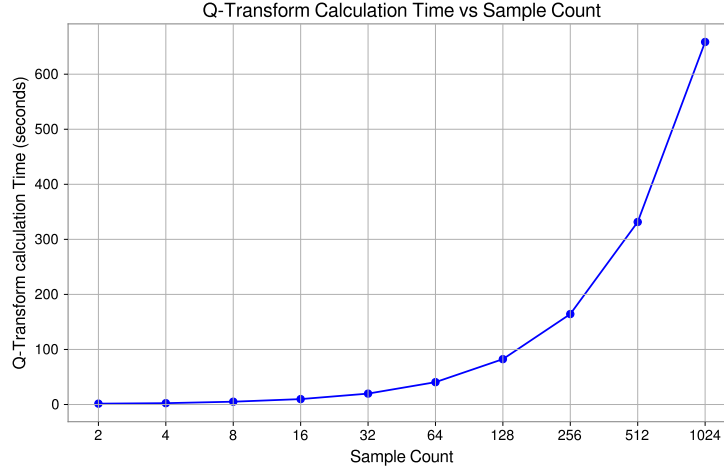
This study uses three tests of Gaussianity, namely, (1) the Shapiro-Wilk, (2) Kolmogorov-Smirnov, and (3) Anderson-Darling tests on preprocessed samples of the L1 interferometer data to determine their normality and assess how well they differentiate between Gaussian and non-Gaussian samples. This experiment also includes studying the performance of these tests at different frequency ranges, with particular emphasis on their effectiveness across various glitch types.

The main motivation for this project is to gauge the use of statistical hypothesis testing on time-amplitude domain data as a faster alternative to the current glitch detection methods working in the time-frequency domain. Detecting and mitigating the effects of glitches in interferometric strain data remains an active area of research within astrophysical data analysis, with several techniques proposed and implemented for the same [11–13]. However, many of these are computationally expensive, slowing down the process of assessing of the signals for glitch activity. A lot of work around GW and glitch analysis makes use of the *Q-transform* [14–16] and modifications of it such as the Q-gram [17]. The Q-transform is a modification of the standard Fourier transform where the analysis window scales inversely with frequency. It is used as inputs to statistical tests

to obtain a p-value depicting the statistical significance of excess power in the data. Assuming the background data to be Gaussian, the excess power that does not follow a Gaussian distribution suggests the presence of a glitch. While this method is effective for post-detection run analysis, it is much slower when scaled up to work on multiple samples. Figures 2a and 2b illustrate the time taken to compute the Q-transform while varying the sample window size and the number of samples respectively.



(a) A plot of the Q-transform calculation time against the sample window for a single sample.



(b) A plot of the Q-transform calculation time against the sample count for fixed window size.

Figure 2: Plots showing how the time taken to compute the Q-transform scales with variation of the sample window size and the number of samples.

This report is organized as follows: Section 2 discusses how interferometer data is obtained and describes the properties of the time series data used in this study. It also outlines the process of acquiring clean and glitched samples and the preprocessing applied. Section 3 describes in detail the statistical tests of Gaussianity used to assess the sample data. The experimental procedures and results are presented in section 4, followed by section 5 which discusses the implications

of our findings. Finally, Section 6 summarizes the results of our experiments, discussing the strengths and limitations of the methods used. The link to this project can be found at <https://github.com/ashishmathew0297/gwpy-experiments>

2 Data Acquisition and Conditioning

The Advanced LIGO and Virgo interferometers are large-scale, heavily modified versions of the Michelson interferometer (Figure 3a), originally invented by American physicist Albert A. Michelson in 1887 [6, 18]. LIGO operates in a vacuum, using a laser beam of light split into two orthogonal parts with the help of a half-silvered mirror mounted on horizontal seismometer suspensions. The orthogonally split beams are then sent through two arms of the interferometer, known as Fabry-Pérot cavities. Each of these arms is 4 kilometers long, consisting of mirrors at each end. The first of the mirrors is fixed in place, while the second one is movable with the help of precise micrometer drives. The laser beams, when reflected by the mirrors on either end of the arms, are combined again at the beam splitter and sent to photoelectric detectors.

The advanced LIGO interferometer is designed such that if the beams travel equal distances, the recombination would lead to a destructive interference pattern, resulting in no light coming out of the instrument [19]. When a disturbance of either astrophysical or terrestrial origin is detected, an infinitesimally small change occurs in lengths of the detector arms, in the order of $10^{-8}m$ [20]. One of the arms is stretched, and the other arm is compressed in the perpendicular direction, altering the lengths of the reflected and combined beams, misaligning them and resulting in an interference pattern (Figure 3b). This interference pattern, coupled with results from several other sensors, provides detailed information on the *strain* of the GW detector.



Figure 3: Illustrations of a basic Michelson interferometer [21, 22] and what an interference pattern would look like.

The **GWpy** Python package is widely used in this project. It provides a suite of tools to access and condition detector strain data from the Gravitational-Wave Open Science Centre (GWOSC) database [23]. This data is in the time domain, as timestamps in the GPS time system at nanosecond precision, and records the amplitude of the calibrated noise event as a differential change in lengths of the interferometer arms. **GWpy** handles detector data using the **TimeSeries** object, which is built upon **numpy** arrays. This allows compatibility with most of the core **numpy** utilities along with custom functions for signal processing, tabular data filtering, and visualization.

We use the strain data from L1 during O3a for this project due to the high rate of occurrence of glitches. Figure 4 shows the steps taken to acquire signal data from the detector and condition it for statistical testing.



Figure 4: The Data Acquisition and Conditioning Pipeline

There are two main types of detector samples that we receive from the detectors: **Glitched detector data** and **Clean detector data** with separate modules used to obtain them.

2.1 The Glitch Segment Acquisition Module

The **Glitch Segment Acquisition Module** uses GPS times of glitch occurrences obtained from *Gravity Spy* [15] to fetch our glitch samples from the GWOSC database. GravitySpy is a large-scale citizen-science project that combines astrophysics, machine learning, and human efforts to classify glitches in GW interferometer data. The *Omicron* transient search algorithm is used by Gravity Spy to generate q-transform spectrograms and calculate the SNR of the time series samples [11]. This algorithm is crucial in identifying the most useful samples for data classification and analysis. Based on the morphological characteristics from the spectrograms, a total of 22 glitch classes were identified with an SNR above 7.5 and peak frequencies between 10 Hz and 2048 Hz.

In the O3a data for L1, Fast_Scattering and Tomte glitches are the most prevalent, while Chirp, 1080Lines and Wandering_line glitches have fewer samples, as shown in Table 1. The No_Glitch class represents glitch samples that lack significant traits or energy levels and do not fit in with the other classes morphologically. Hence, for our study, we do not consider this glitch class.

Glitch Class	Count	Glitch Class	Count
Fast_Scattering	21749	Whistle	896
Tomte	18708	Low_Frequency_Lines	788
Blip_Low_Frequency	7549	Scratchy	207
Scattered_Light	5398	Repeating_Blips	164
No_Glitch	5358	Violin_Mode	164
Extremely_Loud	4319	Paired_Doves	155
Koi_Fish	4268	Light_Modulation	72
1400Ripples	2363	Helix	21
Blip	1947	Wandering_Line	20
Power_Line	1189	1080Lines	9
Low_Frequency_Burst	1187	Chirp	6

Table 1: Glitch counts per class for LIGO Livingston (L1) during the O3a run.

The time series samples of the glitches are obtained using the `TimeSeries.fetch_open_data()` API provided by GWpy. When querying the API, a sample rate of 4096 Hz is used and a padding of 10 seconds is taken on either side of each of the glitch GPS time to obtain the beginning and end times of the window for preprocessing the glitch. This is done to ensure that the sample being used is large enough for data conditioning.

2.2 The Clean Segment Acquisition Module

The **Clean Segment Acquisition Module** is used to both, find the start and end times of the gaps of clean detector noise as well as obtain the corresponding time series data. This is done using `gwtrigfind`, a package developed to search for event triggers files from GW detectors, in conjunction with `EventTable` and `TimeSeries.fetch_open_data()`, provided by GWpy.

Taking the start and end GPS times of the O3a run, `gwtrigfind` is used to find the file path containing Omicron triggers from the *L1:GDS-CALIB-STRAIN* channel of the L1 detector. `EventTable` is then used to load all the trigger data, containing information on the start and end GPS times of the events. Taking the time frames between successive events, we obtain the start and end times of the gaps between glitches/GW triggers, which do not have a significant amount of noise activity. The sizes of the gaps are clamped between 7 and 30 seconds because too short of a gap could lead to the inclusion of glitches in the sample, as the time series data may not have enough time to stabilize after a glitch. Additionally, a short time segment would affect the PSD calculation and subsequent preprocessing. Too long of a time gap could be indicative of detector malfunction or a period of time when it is not operational.

The start and end GPS times of the clean segments are then used to fetch the time series data from the GWOSC database at a sample rate of 4096 Hz. Then their Q-grams, and subsequently, the p-values of each of them are computed [17].

The p-value calculation is done by taking the energy values from the Q-grams, dividing them into N bins and fitting a monotonic logarithmic curve to the maximum energy values of the bins ϵ_0 to obtain the parameters A and t . Then the p-value is calculated as follows

$$P(\text{Q-gram max } \epsilon = \epsilon_0) = 1 - \exp \left[-\frac{AN}{t} e^{-\epsilon_0 t} \right] \quad (1)$$

The samples with a p-value greater than 0.05 are considered to be clean segments of data.

The clean GPS times from the clean segment acquisition module are saved in a CSV file, and the queried time series data is also stored locally. This allows them to be used for procuring the time series samples without having to query the GWOSC database during multiple stages of analysis.

2.3 Data Preprocessing & Conditioning

The strain readings obtained from a GW detector are usually a combination of the GW signal and detector noise [7, 24, 25]. Given our focus on the noise, it is important to bring out its characteristics such that it can be studied. Considering the strain output to be $s(t)$ with $n(t)$ as the noise and $h(t)$ being the possible GW signal received, the detector output is given by:

$$s(t) = n(t) + h(t) \quad (2)$$

The preprocessing in this project, which was mentioned in previous section and will be discussed in sections 2.3.1 and 2.3.2, involves the use of Fourier transforms to divide the signal into its frequency components. Given the signal $x(t)$ in the time domain we get its corresponding Fourier transform in the frequency domain $\tilde{x}(f)$ as

$$\tilde{x}(f) = \mathcal{F}\{x(t)\}(f) = \int_{-\infty}^{\infty} x(t)e^{-2\pi ift} dt \quad (3)$$

and the inverse Fourier Transform given by

$$x(t) = \mathcal{F}^{-1}\{\tilde{x}(f)\}(t) = \int_{-\infty}^{\infty} \tilde{x}(f)e^{2\pi ift} df \quad (4)$$

2.3.1 Power Spectral Density and Amplitude Spectral Density

One of the methods to study the detector noise is the **Power Spectral Density (PSD)**. The PSD, provides a representation of how power is distributed across different frequency bands in a noise signal[26], allowing for the identification of dominant frequency components and their corresponding amplitudes. Assuming the noise sources to be non-stationary and taking the frequency domain representation of the noise signal $n(f)$, the ensemble average of the Fourier components is given by [7]

$$\langle \tilde{n}(f)\tilde{n}^*(f') \rangle = \frac{1}{2}\delta(f - f')S_n(f). \quad (5)$$

Where, $\langle \cdot \rangle$ represents the ensemble average, $\tilde{\cdot}$ represents the Fourier transform from 3, \cdot^* represents the complex conjugate, and $S_n(f)$ is the PSD. Aside from PSD, another way to characterize the noise of a detector is with the **Amplitude Spectral Density (ASD)** values, given by

$$ASD(f) = \sqrt{S_n(f)} \quad (6)$$

The ASD measures the amplitude of the signal at each frequency, and it is often used to compare the noise levels of different detectors or to assess a detector's sensitivity to specific frequencies. It is normally plotted on a logarithmic scale, with frequency on the x-axis and the ASD on the y-axis, allowing for easy identification of noise peaks and their corresponding frequencies.

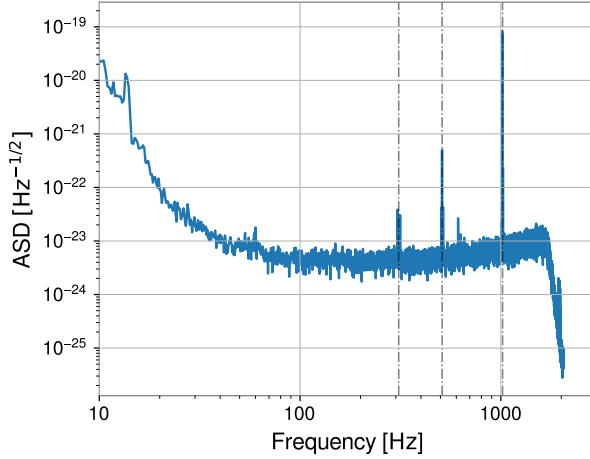


Figure 5: ASD plot for a Tomte Glitch. The gray dotted lines coincide with strong spectral lines at around the 300, 512, and 1024 Hz ranges, marked by dotted lines. There is a big chance that these are of instrumental origin.

From figure 5, which shows an example Tomte glitch ASD plot, to obtain the root-mean-square strain noise at a given frequency band, we integrate over the squares of the ASD readings over the frequency band of interest and take its square root. The problem with ASD plots, however, is that it does not visually capture the glitch signal well enough because they are relatively weak, highly transient and easily overpowered by the instrument noise. To better visualize the glitch, we would instead choose to **whiten** the time series data and experiment with it in the time-amplitude domain.

2.3.2 Whitening

Since the detector noise $n(t)$ comes from a combination of sources from both, the detector and the environment, the signals that we want to study are buried under coloured detector noise [7]. This makes the calculation of the PSD difficult. The process of **whitening** helps with this issue. Whitening makes the data Gaussian-like with a uniform variance and no correlation of the noise [7].

Whitening is done by dividing the Fourier coefficients of the original time series data by the estimated ASD [8].

$$\tilde{n}_{\text{whitened}}(f) = \frac{\tilde{n}(f)}{ASD} = \frac{\tilde{n}(f)}{\sqrt{S_n(f)}} \quad (7)$$

$$n_{\text{whitened}}(t) = \mathcal{F}^{-1} \{ \tilde{n}_{\text{whitened}}(f) \} \quad (8)$$

This data from this process allows for a clearer view of the glitch signal, and can be used for further analysis, such as statistical testing or machine learning classification.

In figure 6 we see a Tomte glitch signal before and after whitening. In the first plot we do not clearly see the effects of the glitch on the time series data. However, after whitening the data,

we can clearly see the glitch signal in the time-amplitude plot. The Q-transform allows a view of the non-stationary portion of the glitch signal, which shows a clear peak in the normalized energy coinciding with the glitch event. For a clean signal we would not see any significant peaks in the whitened time series or non-stationary parts in the Q-transform.

In this project, the data from the clean and glitch segment acquisition modules are whitened using the method. After whitening, the glitch samples are cropped to a 0.5-second time window around the GPS time of the glitch event, effectively giving us a 1-second glitch sample to work with. The clean data, in turn, is divided into 1-second segments, each being treated as a separate clean sample for statistical testing. The 1-second time frame is selected, as it provides us 4096 data points to work with, ensuring that the tests give us the most reliable results possible.

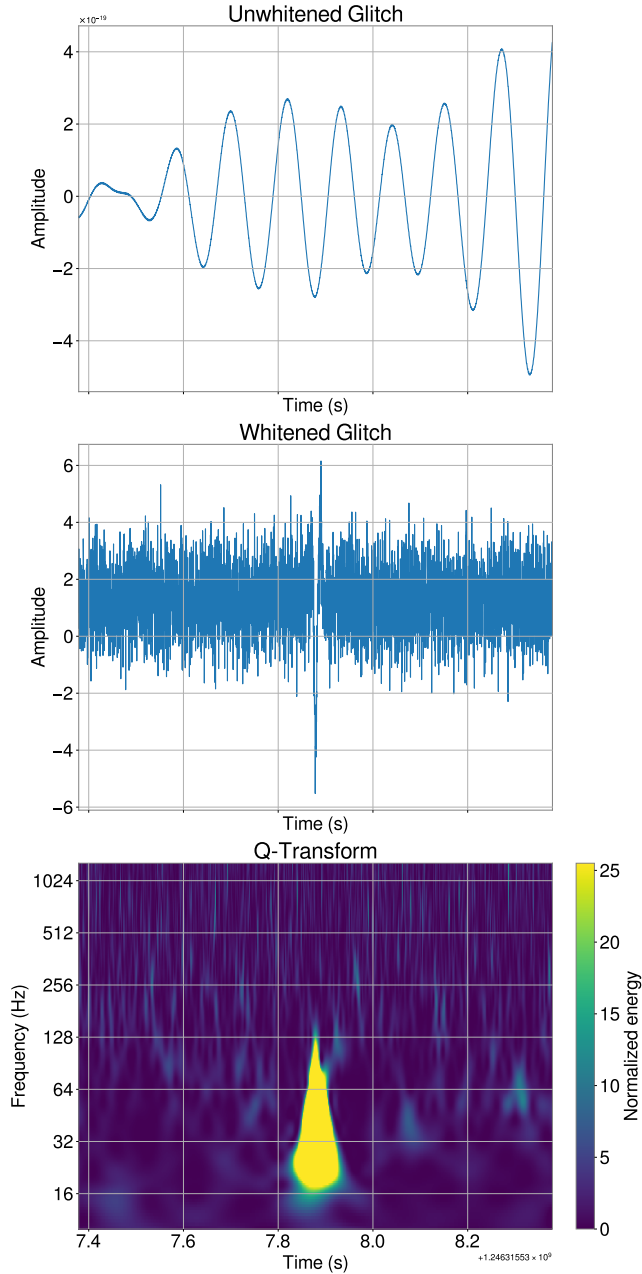


Figure 6: Example of a Tomte glitch signal before and after whitening.

3 Methods

Our main objective is to determine the Gaussianity of time series data obtained from L1 and use this as a criterion to determine the presence of a glitch in it.

If a sample is found to be non-Gaussian, this indicates an excess of power at certain frequencies due to which there is a high chance of a glitch being present. In the previous section, we have seen how whitening the time series data brings out the glitch characteristics. Taking an example of a random Blip glitch, if we were to treat its whitened amplitude values as a population of points and plot the distribution, we notice that it contains heavier tails and a higher peak, hence making it non-Gaussian [7](#). Similarly, when observing a sample of clean data, the signal is closer to that of a Gaussian distribution. Hence, we propose a pair of hypotheses to determine the presence of a glitch in a sample of data from the detector.

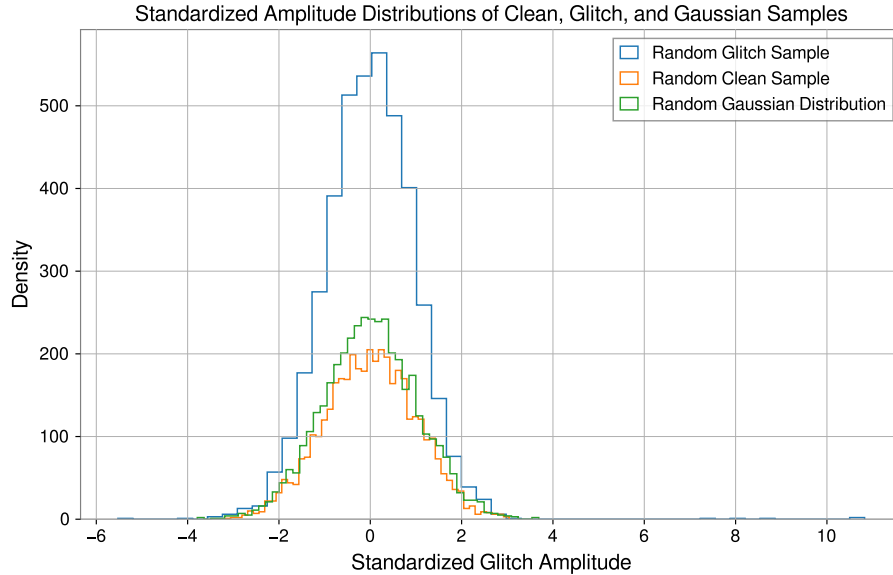


Figure 7: The distributions of standardized amplitudes for a glitch and a clean signal compared to Gaussian distributions.

Null Hypothesis (H_0):

Considering a sample distribution D of whitened time series detector data, we assume the distribution of the amplitude values to be Gaussian

$$D = \mathcal{N}(\mu, \sigma^2)$$

Hence, the sample does not contain a glitch, i.e. it is considered to be clean.

Alternative Hypothesis (H_1):

Considering a sample distribution D of whitened time series detector data, we find that the distribution of the amplitudes is non-Gaussian,

$$D \neq \mathcal{N}(\mu, \sigma^2)$$

Hence, the sample contains a glitch.

We employ three statistical tests of normality to test these hypotheses: the Shapiro-Wilk test, the Kolmogorov-Smirnov test, and the Anderson-Darling test. Each of these have their own strengths and limitations, and are used to assess the normality of the data in different ways. These tests give us a statistic along with a p-value, the latter of which will be used to determine the significance of the result. The p-value tells us how strongly the statistic rejects the null hypothesis based on a significance level α . For our tests, we have the significance level at 0.05. If the p-value computed by our tests is less than 0.05, we reject the null hypothesis and conclude that the sample data is non-Gaussian, hence containing a glitch. Any samples with a p-value greater than or equal to 0.05 are considered to be clean.

3.1 The Shapiro-Wilk Test

The Shapiro-Wilk test [27] is a parametric test, meaning that it assumes the data follows a specific distribution. This parametric test follows the idea of calculating a test statistic W based on the ratio of the best estimate of the sample data's variance to its actual variance. This was the first test that could detect departures from normality due to either the skewness or the kurtosis of the data.

Skewness is a measure of asymmetry of a distribution. It gives us information how much the data deviates either to the left or right from a symmetric distribution. A skewness of 0 shows a perfectly symmetric distribution, while positive and negative values show a right and left skew respectively. **Kurtosis**, on the other hand, measures the tails of the distribution, indicating how much of the data is contained in the tails compared to a normal distribution. A normal distribution has a kurtosis of 3, while values greater than 3 and less than 3 indicate heavier and lighter tails respectively.

If we were to consider x_1, x_2, \dots, x_n , to be a collection of ordered sample points from a population, ranked smallest to largest, with the sample mean \bar{x} and a_1, a_2, \dots, a_n being constants computed from the expected values of the order statistics of a normal distribution, then the test statistic W is given by:

$$W = \frac{\left(\sum_{i=1}^n a_i x_{(i)}\right)^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (9)$$

Taking $m = (m_1, m_2, \dots, m_n)$ to be the expected values of the order statistics of a normal distribution such that $m_i = \mathbb{E}[z_i]$ where z_i is the i -th order statistic of the normal distribution, and V the covariance matrix of the ordered statistics, the constants a_i , which make up the **optimal weight vector** $a \in \mathbb{R}^n$, can be computed as

$$a_i = \frac{m^T V^{-1}}{\|V^{-1}m\|} \quad (10)$$

$$= \frac{m^T V^{-1}}{(m^T V^{-1} V^{-1} m)^{1/2}} \quad (11)$$

The test statistic W lies in the range $[0, 1]$, with values closer to 1 indicating that the sample data is more likely to be Gaussian. The null hypothesis H_0 is rejected if W is significantly less than 1, indicating that the sample data is non-Gaussian.

The p-value of the statistic is usually used to determine the significance of the result. If the p-value is less than the significance level (α) defined, the null hypothesis is rejected, indicating that the sample data is non-Gaussian and likely contains a glitch.

The Shapiro-Wilk test, though powerful, has some limitations. It is sensitive to the sample size and may not perform well with small or large samples. It works ideally for sample sizes $n \leq 2000$. The `shapiro()` function from `scipy.stats` works well for $n \leq 5000$ after which the p-value calculations' accuracy drops. Additionally, for this test to work the data must be continuous and univariate.

3.2 The Kolmogorov-Smirnov Test

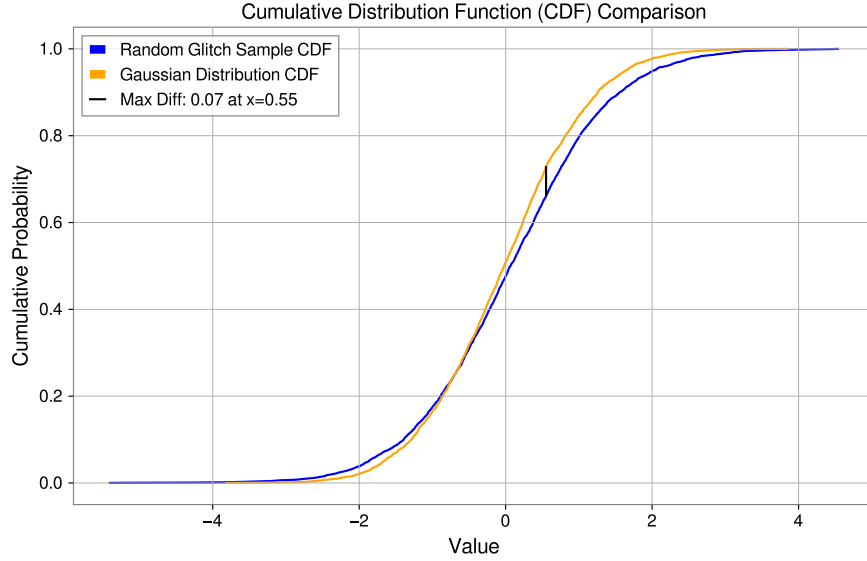


Figure 8: The KS statistic calculation for a sample glitch against a Gaussian distribution. The black line is the maximum absolute difference between the sample glitch CDF and the Gaussian CDF, which is the KS statistic D_n .

The Kolmogorov-Smirnov (KS) test is a non-parametric test, which means that it does not have any assumption that our distribution follows a specific probability distribution. Originally proposed in the 1930s, this test is used to decide whether a sample comes from a particular type of distribution [28, 29]. The idea of this test is to check the difference in shape between the distributions being studied, which is done by measuring the maximum vertical distance between the **empirical cumulative distribution function (ECDF)** of the sample and the **cumulative distribution function (CDF)** of the reference distribution, which in this case is a random Gaussian distribution 8. This test has two versions: the **one-sample KS test** which is used to compare an empirical sample distribution with a reference theoretical distribution, and the **two-sample KS test** which compares the empirical distributions of two samples. We use the one-sample KS test in this project.

Taking a sample of ordered points x_1, x_2, \dots, x_N as N ordered points from a population, the Empirical CDF of the sample, is given by

$$F_N = \frac{n(i)}{N} \quad (12)$$

Where $n(i)$ is the number of points less than or equal to x_i in the sample. This is in the form of a step function increasing at steps of $1/N$ at each ordered point [30]. Taking the theoretical CDF of the distribution being tested as $F(x)$, a Gaussian distribution in our case, the KS test statistic D_n is defined as

$$D_n = \sup_{1 \leq i \leq N} \left(F(x_i) - \frac{i-1}{N}, \frac{i}{N} - F(x_i) \right) \quad (13)$$

Here, sup refers to the supremum or maximum of the differences between the theoretical CDF and the ECDF. If D_n is significantly larger than the defined significance level, it indicates that the sample data significantly differs from a Gaussian distribution, hence rejecting the null Hypothesis. In our implementation, we use the p-value of the statistic to determine the significance of the result as it is directly proportional to D_n .

3.2.1 Note: Scaling the data

The Kolmogorov-Smirnov test is highly affected by the values of the data. We need to ensure that the comparisons being done in the test are not affected by the absolute values of the data points. Hence, this requires for the data to be scaled down to the same range of a standard normal distribution. This is achieved by scaling the data to have a mean of 0 and a standard deviation of 1. This is also known as **standardization** or **z-score normalization**.

To achieve this the `StandardScaler` class from the `sklearn.preprocessing` module is used. This module works by transposing the data to have a zero mean, calculates the standard deviation, and scales the data points down such that the standard deviation is 1. The **standard score** or **z-score** of a sample point, which tells us how many standard deviations it is from the mean, is given by

$$z = \frac{x - \mu}{\sigma}. \quad (14)$$

Where x is the original data point, μ is the mean, and σ is the standard deviation of the data. The z value obtained can be treated as the new data point after scaling and is particularly useful in this case as it does not change the position of the data point in its calculation.

3.3 The Anderson Darling test

The Anderson-Darling (AD) test is a non-parametric goodness-of-fit test used to determine whether a sample of data is drawn from a specific distribution, most commonly a normal distribution. It is a modification of the Kolmogorov-Smirnov and Cramer-von Mises tests with more sensitivity to differences in the tails of the distribution [30, 31].

Similar to the Kolmogorov-Smirnov test, this test compares the ECDF of the sample data with the CDF of a theoretical reference distribution by taking a distance function between them. The main difference here is that the full range of the data is used to calculate the distance function rather than just the largest distance value [31]. Here, taking n to be the number of elements in the sample and $w(x)$ to be the weighting function, the Anderson-Darling statistic distance A^2 , given by the squared area between the sample ECDF $F_n(x)$ and the theoretical reference CDF $F(x)$, is calculated as

$$A^2 = n \int_{-\infty}^{\infty} (F_n(x) - F(x))^2 w(x) dF(x) \quad (15)$$

Taking the weighting function $w(x)$ is defined as

$$w(x) = \frac{1}{F(x)(1 - F(x))} \quad (16)$$

We get A^2 [32] as

$$A^2 = n \int_{-\infty}^{\infty} \frac{(F_n(x) - F(x))^2}{F(x)(1 - F(x))} dF(x) \quad (17)$$

To obtain random samples from a distribution function, the CDF is to be calculated and random samples x are to be taken from a uniform distribution $x \in [0, 1]$. Hence, the Anderson-Darling test statistic can be rewritten as

$$A^2 = -N - S \quad (18)$$

where S is given by

$$S = \frac{1}{N} \sum_{i=1}^N (2i - 1) (\ln F(x_i) + \ln(1 - F(x_{N-i+1}))) \quad (19)$$

hence, giving us the final form of the Anderson-Darling statistic as

$$A^2 = -N - \sum_{i=1}^N \left(\frac{2i - 1}{N} \right) (\ln F(x_i) + \ln(1 - F(x_{N-i+1}))) \quad (20)$$

This is effectively a weighted cross-product of the samples, allowing for a comparison of the sample data with the reference distribution.

In the case of the mean and standard deviation of the sample distribution being known with the distribution size, N being greater than 5, the α levels are 1%, 2.5%, 5%, 10% and 15%, with the critical values for each significant level being pre-computed through Monte-Carlo simulations. The null hypothesis is rejected if the calculated statistic is greater than the critical value at a selected α , indicating that the sample data is not normally distributed.

3.4 Using the tests for glitch detection

In our experiments, we treat the statistical tests as classifiers, where a positive class indicates the presence of a glitch in the sample, irrespective of the type, and a negative class represents a clean signal. This is an approach usually taken with machine learning models, which predict an outcome by learning patterns from the data. In our case, the statistical tests are based on fixed sets of criteria, hence forgoing the learning aspect and directly providing a result based on the input data. This provides a straightforward way of evaluating the performance of each test through the evaluation metrics of the results.

The evaluation metrics used in this project are based on the following definitions:

- **True Positives (TP)**: The number of samples correctly identified as glitches by the test.
- **False Negatives (FN)**: The number of samples incorrectly identified to be clean/Gaussian when they contained glitches.
- **False Positives (FP)**: The number of samples incorrectly identified as glitches when they were actually clean.
- **True Negatives (TN)**: The number of samples correctly identified as clean by the test.

The following evaluation metrics are what we consider for our tests:

- **Accuracy**: The number of correct predictions out of the total number of samples. This is the most common metric used to determine the overall performance of the model or test being studied. However, it can be misleading when taken at face value as it does not account for dataset imbalances or model/test bias.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (21)$$

- **Recall or True Positive Rate (TPR)**: The number of actual positive samples correctly identified as positive over the total number of actual positive samples gives us the *Recall*. This tells us how trustworthy the positive predictions made by the test are.

$$\text{Recall} = \frac{TP}{TP + FN} \quad (22)$$

- **False Negative Rate (FNR)**: This is the number of actual positive samples incorrectly identified as negative over the total number of actual positive samples. We would ideally want this value to be low to ensure not missing out on identifying glitches in our noise samples.

$$\text{FNR} = \frac{FN}{TP + FN} \quad (23)$$

- **Specificity or True Negative Rate (TNR)**: This is the number of actual negative samples correctly identified as negative over the total number of actual negative samples. This tells us how trustworthy the negative predictions made by the test are.

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (24)$$

- **False Positive Rate (FPR):** This is the number of actual negative samples incorrectly identified as positive over the total number of actual negative samples. It is ideal to keep this value low to ensure that our tests do not incorrectly label clean signals as glitches.

$$\text{FPR} = \frac{FP}{FP + TN} \quad (25)$$

- **Precision:** The number of predicted positive samples that were actually positive over the total number of predicted positive samples gives us the *Precision*. The precision tells us how trustworthy the positive predictions made by the test are.

$$\text{Precision} = \frac{TP}{TP + FP} \quad (26)$$

- **F1 Score:** This is another measure of accuracy which makes use of the Precision and Recall. As opposed to Accuracy, which computes how many times the model makes a correct prediction, it sees how well the model performs on each class. Learning models play a balancing act between the Precision and Recall, as the precision places more emphasis on getting the positive predictions right, while recall focuses on getting as many positive predictions as possible. Getting both these values high would be an indicator that the model or test is performing well. The F1 score combines precision and recall into a single metric by taking their harmonic mean to ensure that both are maximized.

$$\text{F1 Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (27)$$

Along with studying how effective these tests are at determining the presence of a glitch in detector data, we will also be validating the claim in [33], which states that the Shapiro-Wilk test is the most powerful, followed by the Anderson darling test and the Kolmogorov-Smirnov test.

4 Experimentation

The implementation of our statistical tests on the whitened time series data is fairly straightforward, with the `scipy.stats` package providing dedicated functions for each. Once the time series data for the clean and glitch samples are obtained and conditioned, their amplitude values are directly used as inputs to the following functions to obtain the relevant statistics and p-values.

- `scipy.stats.shapiro()` for the Shapiro-Wilk test
- `scipy.stats.ks_1samp()` for the one-sample Kolmogorov-Smirnov test
- `scipy.stats.anderson()` for the Anderson-Darling test

For our experiment we take up to 101 GPS times for each glitch type and 700 GPS time pairs of clean signals, which after querying and preprocessing gives us the following sample counts.

Glitch Class	Count	Glitch Class	Count
Clean_Signal	1824	Repeating_Blips	101
1400Ripples	101	Scattered_Light	101
Air_Compressor	101	Power_Line	100
Blip_Low_Frequency	101	Koi_Fish	99
Blip	101	Low_Frequency_Burst	98
Fast_Scattering	101	Light_Modulation	72
Extremely_Loud	101	Scratchy	45
Whistle	101	Helix	21
Violin_Mode	101	Wandering_Line	9
Low_Frequency_Lines	101	Chirp	6
Paired_Doves	101	1080Lines	6
Tomte	101		

Table 2: Count of the clean and glitch samples taken for the experiment.

It is evident that due to the variation in the number of glitches, the data we are dealing with is skewed. Due to this we will need to study how each of our tests perform when detecting each type of glitch. Additionally, some glitches, such as 1400 Ripples, 1080 Lines, and Scattered Light, occur at specific frequency ranges. Due to this, we also experiment with using a band pass filter on our data between 10 Hz to 512 Hz and 512 Hz to 1024 Hz, and applying our statistical tests to this filtered data.

4.1 Experimenting on the complete frequency range

Taking the whole frequency range into account, we get the following results for the tests (Table 3).

Test	TP	FN	FP	TN	Accuracy	TPR	TNR	FPR	FNR	Precision	F1 Score
Shapiro	913	856	83	1741	0.74	0.52	0.95	0.05	0.48	0.92	0.66
KS	491	1278	0	1824	0.64	0.28	1.00	0.00	0.72	1.00	0.43
Anderson	472	1297	0	1824	0.64	0.27	1.00	0.00	0.73	1.00	0.42

Table 3: Glitch detection results for the full frequency range from 10 Hz to 1024Hz. All these tests were performed on the amplitude values of the whitened time series samples with an α of 0.05.

At a first glance, the Shapiro-Wilk test has the best accuracy of the three at 0.74, with the accuracies of the Kolmogorov-Smirnov and Anderson Darling tests being approximately 10% lower than it. This is also seen in the F1 score for the Shapiro-Wilk test being higher than the other two tests. The false positives being low for the Shapiro-Wilk test and zero for the other two tests tells us that all of them excel at determining whether a sample is clean. The recall of 0.52 tells us that the Shapiro-Wilk test is able to capture slightly more than half of the cases where a glitch is present whereas the Kolmogorov-Smirnov and Anderson-Darling tests struggle with the same. Looking at the number of false negatives and the false negative rates gives us more insight into their performance.

The Shapiro-Wilk test has a false negative rate and true positive rate that are almost equal, indicating that though it is able to effectively identify glitches, it also has a significant number of false negatives which would give reason to question its reliability. The Kolmogorov-Smirnov and Anderson-Darling tests in comparison have an even higher number of False Negatives, hence making them much less effective for glitch detection.

To further understand the effectiveness of these tests, we can analyze the results on a per-glitch level.

Label	TP	FN	TPR	FNR
1080Lines	4	2	0.67	0.33
1400Ripples	72	29	0.71	0.29
Air_Compressor	7	94	0.07	0.93
Blip	97	4	0.96	0.04
Blip_Low_Frequency	15	86	0.15	0.85
Chirp	2	4	0.33	0.67
Extremely_Loud	101	0	1.00	0.00
Fast_Scattering	3	98	0.03	0.97
Helix	19	2	0.90	0.10
Koi_Fish	99	0	1.00	0.00
Light_Modulation	67	5	0.93	0.07
Low_Frequency_Burst	8	90	0.08	0.92
Low_Frequency_Lines	2	99	0.02	0.98
Paired_Doves	43	58	0.43	0.57
Power_Line	7	93	0.07	0.93
Repeating_Blips	98	3	0.97	0.03
Scattered_Light	10	91	0.10	0.90
Scratchy	8	37	0.18	0.82
Tomte	55	46	0.54	0.46
Violin_Mode	86	15	0.85	0.15
Wandering_Line	9	0	1.00	0.00
Whistle	101	0	1.00	0.00

Table 4: Shapiro-Wilk test results for each glitch class at the full frequency range from 10 Hz to 1024Hz.

Looking at the results on a per-glitch level in tables 4, 5, and 6, we see that the tests’ performances vary significantly across each of the glitch classes. All the three tests achieve perfect accuracy and F1 Scores for Extremely_Loud glitches, which is to be expected due to these glitches being highly distinctive over a large frequency range with a high SNR. In the case of the Blip, Repeating_Blips, whistle and Violin_Mode glitches, the Shapiro-Wilk test far surpasses the other two tests, with a fixed pattern of the Shapiro-Wilk being the best, followed by the Kolmogorov-Smirnov test and then the Anderson-Darling test. In a few unique cases, such as 1400Ripples, Air_Compressor and Fast_Scattering glitches, the Shapiro-Wilk test is able to detect a few of the glitch samples, while the other two tests completely fail at detecting these glitches.

When looking at the results for clean signals the Shapiro-Wilk test has a slightly higher number of false positives than the other two tests, which is to be expected due to its sensitivity to the distribution of the data. The Kolmogorov-Smirnov and Anderson-Darling tests have no false positives, indicating that they are able to effectively identify clean signals without misclassifying them as glitches.

Label	TP	FN	TPR	FNR
1080Lines	0	6	0.00	1.00
1400Ripples	0	101	0.00	1.00
Air_Compressor	0	101	0.00	1.00
Blip	37	64	0.37	0.63
Blip_Low_Frequency	2	99	0.02	0.98
Chirp	0	6	0.00	1.00
Extremely_Loud	101	0	1.00	0.00
Fast_Scattering	0	101	0.00	1.00
Helix	12	9	0.57	0.43
Koi_Fish	98	1	0.99	0.01
Light_Modulation	60	12	0.83	0.17
Low_Frequency_Burst	0	98	0.00	1.00
Low_Frequency_Lines	0	101	0.00	1.00
Paired_Doves	20	81	0.20	0.80
Power_Line	1	99	0.01	0.99
Repeating_Blips	62	39	0.61	0.39
Scattered_Light	3	98	0.03	0.97
Scratchy	3	42	0.07	0.93
Tomte	14	87	0.14	0.86
Violin_Mode	18	83	0.18	0.82
Wandering_Line	5	4	0.56	0.44
Whistle	55	46	0.54	0.46

Table 5: Kolmogorov-Smirnov test results for each glitch class at the full frequency range from 10 Hz to 1024Hz.

Label	TP	FN	TPR	FNR
1080Lines	0	6	0.00	1.00
1400Ripples	0	101	0.00	1.00
Air_Compressor	0	101	0.00	1.00
Blip	35	66	0.35	0.65
Blip_Low_Frequency	2	99	0.02	0.98
Chirp	0	6	0.00	1.00
Extremely_Loud	101	0	1.00	0.00
Fast_Scattering	0	101	0.00	1.00
Helix	8	13	0.38	0.62
Koi_Fish	98	1	0.99	0.01
Light_Modulation	58	14	0.81	0.19
Low_Frequency_Burst	0	98	0.00	1.00
Low_Frequency_Lines	0	101	0.00	1.00
Paired_Doves	17	84	0.17	0.83
Power_Line	1	99	0.01	0.99
Repeating_Blips	60	41	0.59	0.41
Scattered_Light	2	99	0.02	0.98
Scratchy	3	42	0.07	0.93
Tomte	12	89	0.12	0.88
Violin_Mode	17	84	0.17	0.83
Wandering_Line	5	4	0.56	0.44
Whistle	53	48	0.52	0.48

Table 6: Anderson-Darling test results for each glitch class at the full frequency range from 10 Hz to 1024Hz.

From our observations, we can conclude that the Shapiro-Wilk detects a wider variety of glitches than the other two tests. The Kolmogorov-Smirnov and Anderson-Darling tests have a lower overall detection rate, particularly for glitches with lower SNR such as Scattered_Light.

Overall, the Shapiro-Wilk test is the best of the three tests for detecting glitches on the full frequency range, followed by the Kolmogorov-Smirnov and Anderson-Darling tests respectively, with both of them being significantly less effective. This is in line with the claim made in [33].

4.2 Experimenting with a band pass filter between 10 Hz to 512 Hz

We now perform the same analysis as before on the time series data with a band pass filter applied between 10 Hz and 512 Hz before the preprocessing and conditioning steps. The results of the tests on this filtered data are shown in Table 7.

Test	TP	FN	FP	TN	Accuracy	TPR	TNR	FPR	FNR	Precision	F1 Score
Shapiro	1252	517	735	1089	0.65	0.71	0.60	0.40	0.29	0.63	0.67
KS	630	1139	24	1800	0.68	0.36	0.99	0.01	0.64	0.96	0.52
Anderson	589	1180	0	1824	0.67	0.33	1.00	0.00	0.67	1.00	0.50

Table 7: Glitch detection results for a frequency range of 10 Hz to 512 Hz. All these tests were performed on the amplitude values of the whitened time series samples with an α of 0.05.

Compared to the results of the previous analysis in Table 3, we see an overall increase in true positives for all three tests with a trade-off being the increase in the number of false positives. The Shapiro-Wilk test shows the highest increase in the number of false positives followed by the Kolmogorov-Smirnov test and Anderson-Darling tests respectively. There is an overall decrease in the accuracies for all three tests, with the Shapiro-Wilk test having the lowest accuracy at 0.65.

The increase in true positives at lower frequencies could point towards the possibility that the tests are more sensitive to glitches in this frequency range. However, given that the data has been filtered to only include the low frequency content of the data, it could also mean that removing the higher frequencies has affected the sample distributions, making them imbalanced. This could be a reason for the increase in false positives in the tests.

In this case, the F1 score is a more robust to use than accuracy as it gives us an idea of how well the tests balance the Precision and Recall. In terms of F1 score, we notice that the tests in this case perform better than on the complete frequency range. The Shapiro-Wilk test still performs the best in this scenario, followed by the Kolmogorov-Smirnov and Anderson-Darling tests respectively.

Label	TP	FN	TPR	FNR
1080Lines	3	3	0.50	0.50
1400Ripples	44	57	0.44	0.56
Air_Compressor	54	47	0.53	0.47
Blip	101	0	1.00	0.00
Blip_Low_Frequency	72	29	0.71	0.29
Chirp	5	1	0.83	0.17
Extremely_Loud	101	0	1.00	0.00
Fast_Scattering	41	60	0.41	0.59
Helix	21	0	1.00	0.00
Koi_Fish	99	0	1.00	0.00
Light_Modulation	71	1	0.99	0.01
Low_Frequency_Burst	43	55	0.44	0.56
Low_Frequency_Lines	51	50	0.50	0.50
Paired_Doves	86	15	0.85	0.15
Power_Line	68	32	0.68	0.32
Repeating_Blips	98	3	0.97	0.03
Scattered_Light	65	36	0.64	0.36
Scratchy	34	11	0.76	0.24
Tomte	96	5	0.95	0.05
Violin_Mode	50	51	0.50	0.50
Wandering_Line	3	6	0.33	0.67
Whistle	46	55	0.46	0.54

Table 8: Shapiro-Wilk test results for each glitch class at a frequency range of 10 Hz to 512 Hz.

The results of the tests on a per-glitch level are given in tables 8, 9, and 10. We can see that there is no change in the effectiveness of all the three tests in detecting Extremely_Loud glitches. In the case of the Koi_Fish and Repeating_Blips glitches the performance of the Shapiro-Wilk test remains the same as before while true positive rates of the other two tests increase to 1.00.

In some glitches such as Wandering_Line glitch, the Kolmogorov-Smirnov and Anderson-Darling tests completely miss the detection of the glitch at lower frequencies, while the Shapiro-Wilk test is able to detect a few of the samples, though much lesser than when using the full frequency range. For particularly low SNR glitches such as Scattered_Light, the Shapiro-Wilk test is able to detect significantly more of the samples than before, while the other tests only show a slight improvement.

The performance of the all the three tests have decreased in the case of Whistle, Violin_Mode and Wandering_Line glitches. This suggests a possibility that these glitches have a significant amount of high frequency content which is filtered out in this analysis, hence making them harder to detect. On the other hand, for glitches such as Blip, Blip_Low_Frequency, Helix, and Repeating_Blips, the tests have a significant amount of low frequency content, hence making them easier to detect.

Label	TP	FN	TPR	FNR
1080Lines	0	6	0.00	1.00
1400Ripples	0	101	0.00	1.00
Air_Compressor	1	100	0.01	0.99
Blip	74	27	0.73	0.27
Blip_Low_Frequency	11	90	0.11	0.89
Chirp	2	4	0.33	0.67
Extremely_Loud	101	0	1.00	0.00
Fast_Scattering	2	99	0.02	0.98
Helix	18	3	0.86	0.14
Koi_Fish	99	0	1.00	0.00
Light_Modulation	66	6	0.92	0.08
Low_Frequency_Burst	3	95	0.03	0.97
Low_Frequency_Lines	6	95	0.06	0.94
Paired_Doves	48	53	0.48	0.52
Power_Line	2	98	0.02	0.98
Repeating_Blips	95	6	0.94	0.06
Scattered_Light	15	86	0.15	0.85
Scratchy	14	31	0.31	0.69
Tomte	52	49	0.51	0.49
Violin_Mode	12	89	0.12	0.88
Wandering_Line	0	9	0.00	1.00
Whistle	9	92	0.09	0.91

Table 9: Kolmogorov-Smirnov test results for each glitch class at a frequency range of 10 Hz to 512 Hz.

When looking at the clean signal samples, we see that the Shapiro-Wilk test has the highest number of false positives, which may indicate that it is more sensitive to certain types of noise or artifacts in the data. The Kolmogorov-Smirnov test has a lower number of false positives, and the Anderson-Darling test shows no false positives to be present, indicating that it is the most effective at identifying clean signals at lower frequencies compared to the other tests.

Label	TP	FN	TPR	FNR
1080Lines	0	6	0.00	1.00
1400Ripples	0	101	0.00	1.00
Air_Compressor	0	101	0.00	1.00
Blip	73	28	0.72	0.28
Blip_Low_Frequency	8	93	0.08	0.92
Chirp	1	5	0.17	0.83
Extremely_Loud	101	0	1.00	0.00
Fast_Scattering	1	100	0.01	0.99
Helix	18	3	0.86	0.14
Koi_Fish	99	0	1.00	0.00
Light_Modulation	66	6	0.92	0.08
Low_Frequency_Burst	0	98	0.00	1.00
Low_Frequency_Lines	0	101	0.00	1.00
Paired_Doves	47	54	0.47	0.53
Power_Line	1	99	0.01	0.99
Repeating_Blips	91	10	0.90	0.10
Scattered_Light	10	91	0.10	0.90
Scratchy	8	37	0.18	0.82
Tomte	48	53	0.48	0.52
Violin_Mode	10	91	0.10	0.90
Wandering_Line	0	9	0.00	1.00
Whistle	7	94	0.07	0.93

Table 10: Anderson-Darling test results for each glitch class at a frequency range of 10 Hz to 512 Hz.

4.3 Experimenting with a band pass filter between 512 Hz to 1024 Hz

Similar to the previous experiments we now perform the same tests on the time series data with a band pass filter between 512 Hz and 1024 Hz applied to it prior to the preprocessing and conditioning steps. The results of the tests on this filtered data are shown in Table 11.

Test	TP	FN	FP	TN	Accuracy	TPR	TNR	FPR	FNR	Precision	F1 Score
Shapiro	555	1214	67	1757	0.64	0.31	0.96	0.04	0.69	0.89	0.46
KS	279	1490	0	1824	0.59	0.16	1.00	0.00	0.84	1.00	0.27
Anderson	268	1501	0	1824	0.58	0.15	1.00	0.00	0.85	1.00	0.26

Table 11: Glitch detection results for a frequency range of 512 Hz to 1024 Hz.

Here, we see that the accuracies and F1 scores of all the three tests are at their lowest in contrast with the previous two experiments. This suggests that the tests perform the worst due to there being lesser high frequency content in the data. We, again, notice a small spike in the false positive rate for the Shapiro-Wilk test, comparable to the full frequency range, while the Kolmogorov-Smirnov and Anderson-Darling tests have no false positives. The Shapiro-Wilk test has the highest

true positive rate, followed by the Kolmogorov-Smirnov and Anderson-Darling tests respectively, which is consistent with the results from the previous experiments.

Label	TP	FN	TPR	FNR
1080Lines	3	3	0.50	0.50
1400Ripples	6	95	0.06	0.94
Air_Compressor	2	99	0.02	0.98
Blip	47	54	0.47	0.53
Blip_Low_Frequency	7	94	0.07	0.93
Chirp	0	6	0.00	1.00
Extremely_Loud	100	1	0.99	0.01
Fast_Scattering	5	96	0.05	0.95
Helix	8	13	0.38	0.62
Koi_Fish	82	17	0.83	0.17
Light_Modulation	55	17	0.76	0.24
Low_Frequency_Burst	8	90	0.08	0.92
Low_Frequency_Lines	5	96	0.05	0.95
Paired_Doves	7	94	0.07	0.93
Power_Line	7	93	0.07	0.93
Repeating_Blips	58	43	0.57	0.43
Scattered_Light	2	99	0.02	0.98
Scratchy	4	41	0.09	0.91
Tomte	3	98	0.03	0.97
Violin_Mode	55	46	0.54	0.46
Wandering_Line	9	0	1.00	0.00
Whistle	82	19	0.81	0.19

Table 12: Shapiro-Wilk test results for each glitch class at a frequency range of 512 Hz to 1024 Hz.

Looking at the per-glitch results in tables 12, 13, and 14, we finally see a deviation in the results of the tests for Extremely Loud glitches compared to the previous two experiments. We observe a slight decrease in the true positives for the Kolmogorov-Smirnov and Anderson-Darling tests, with the Shapiro-Wilk test still being able to detect most of the samples. This points to the fact that the Extremely Loud glitch has a significant amount of power in both the lower and higher frequency ranges, hence making it the easiest to detect over a wide range of frequencies. On the other hand, in some the cases such as those of, Scattered Light and Tomte glitches, the detection is at a lower rate by the Shapiro-Wilk test, while almost not being detected at all by the other two tests at higher frequencies. The Kolmogorov-Smirnov and Anderson-Darling tests have similar true positive rates to one another at high frequencies in the case of Koi Fish, Light Modulation, and Repeating Blips glitches. However, their true positive rates are lower than the Shapiro-Wilk test.

The Shapiro-Wilk test is able to detect all the samples of the Wandering_Line glitch, similar to the experiment with the full frequency range, while the other two tests detect a slightly larger number of samples than at the lower frequency range. The Whistle and Violin_Mode glitches show a decent number of detections in the case of the Shapiro-Wilk test, but a lower number for the other two tests at both higher and lower frequency ranges. This suggests that all of these glitches have a significant amount content in both the high and low frequency bands, with the Shapiro-Wilk test being the most effective of the three at detecting them.

Label	TP	FN	TPR	FNR
1080Lines	0	6	0.00	1.00
1400Ripples	0	101	0.00	1.00
Air_Compressor	0	101	0.00	1.00
Blip	21	80	0.21	0.79
Blip_Low_Frequency	0	101	0.00	1.00
Chirp	0	6	0.00	1.00
Extremely_Loud	97	4	0.96	0.04
Fast_Scattering	0	101	0.00	1.00
Helix	0	21	0.00	1.00
Koi_Fish	66	33	0.67	0.33
Light_Modulation	33	39	0.46	0.54
Low_Frequency_Burst	0	98	0.00	1.00
Low_Frequency_Lines	0	101	0.00	1.00
Paired_Doves	0	101	0.00	1.00
Power_Line	0	100	0.00	1.00
Repeating_Blips	31	70	0.31	0.69
Scattered_Light	0	101	0.00	1.00
Scratchy	2	43	0.04	0.96
Tomte	0	101	0.00	1.00
Violin_Mode	7	94	0.07	0.93
Wandering_Line	7	2	0.78	0.22
Whistle	15	86	0.15	0.85

Table 13: Kolmogorov-Smirnov test results for each glitch class at a frequency range of 512 Hz to 1024 Hz.

Overall, the Shapiro-Wilk test continues to show the best performance in terms of true positives and F1 score followed by the Kolmogorov-Smirnov and Anderson-Darling tests. The effectiveness of these tests are depend on the specific characteristics of each glitch class, including their frequency content and the amount of noise present in the signal.

Label	TP	FN	TPR	FNR
1080Lines	0	6	0.00	1.00
1400Ripples	0	101	0.00	1.00
Air_Compressor	0	101	0.00	1.00
Blip	20	81	0.20	0.80
Blip_Low_Frequency	0	101	0.00	1.00
Chirp	0	6	0.00	1.00
Extremely_Loud	96	5	0.95	0.05
Fast_Scattering	0	101	0.00	1.00
Helix	0	21	0.00	1.00
Koi_Fish	65	34	0.66	0.34
Light_Modulation	33	39	0.46	0.54
Low_Frequency_Burst	0	98	0.00	1.00
Low_Frequency_Lines	0	101	0.00	1.00
Paired_Doves	0	101	0.00	1.00
Power_Line	0	100	0.00	1.00
Repeating_Blips	31	70	0.31	0.69
Scattered_Light	0	101	0.00	1.00
Scratchy	2	43	0.04	0.96
Tomte	0	101	0.00	1.00
Violin_Mode	6	95	0.06	0.94
Wandering_Line	6	3	0.67	0.33
Whistle	9	92	0.09	0.91

Table 14: Anderson-Darling test results for each glitch class at a frequency range of 512 Hz to 1024 Hz.

5 Discussion

To reiterate upon the criteria for our tests, our null hypothesis is that the sample of whitened detector data being studied does not contain a glitch, which we determine by the distribution of its amplitude values being Gaussian. We consider this to be a negative label in our experiment. The alternate hypothesis for our tests is that the sample of whitened data being studied contains a glitch, determined by the distribution of its amplitude values being non-Gaussian. We take this as a positive label for our experiment. The null hypothesis is rejected if the tests performed on the sample return a p-value less than 0.05, indicating that the sample is non-Gaussian.

The experiments conducted in this project suggest that the Shapiro-Wilk test is the most effective at detecting glitches across the full frequency range of 10 Hz to 1024 Hz. At the lower and higher frequencies the performance of all the three tests vary depending on the glitch classes and their characteristics. The Shapiro-Wilk test consistently shows the highest true positive rates and F1 scores of the three, indicating a fairly high level of robustness in identifying glitches. The Kolmogorov-Smirnov and Anderson-Darling tests perform similarly to one another, with the Kolmogorov-Smirnov test outperforming the Anderson-Darling test, but they are generally less effective than the Shapiro-Wilk test over all the frequency bands. This confirms the findings of the study in [33] as mentioned in section 3.4.

The Shapiro-Wilk and Kolmogorov-Smirnov tests show a significant number of false positives, particularly at lower frequencies, while the Anderson-Darling test fails to capture a substantial number of glitches. This is a cause for concern in terms of reliability of their sole use in practical applications.

Our results also highlight the fact that these tests vary in performance for each glitch type depending on their frequency content. For example, glitches such as `Extremely Loud` and `Koi Fish` are well detected in all the experiments due to their high amplitudes across the frequency spectrum, while others like `Scattered Light` and `Tomte` show varying detection rates depending on the frequency range. Some glitches, such as `Violin Mode`, `Wandering Line` and `Whistle`, are detected more effectively by the Shapiro-Wilk test compared to the other two tests, while others like `Repeating Blips` show a more balanced performance across all three tests.

Frequency Range	Test	TP	FN	FP	TN	Acc.	TPR	TNR	FPR	FNR	Prec.	F1 Score
Full Range (10 Hz - 1024 Hz)	SW	913	856	83	1741	0.74	0.52	0.95	0.05	0.48	0.92	0.66
	KS	491	1278	0	1824	0.64	0.28	1.00	0.00	0.72	1.00	0.43
	AD	472	1297	0	1824	0.64	0.27	1.00	0.00	0.73	1.00	0.42
Low Freq. (10 - 512 Hz)	SW	1252	517	735	1089	0.65	0.71	0.60	0.40	0.29	0.63	0.67
	KS	630	1139	24	1800	0.68	0.36	0.99	0.01	0.64	0.96	0.52
	AD	589	1180	0	1824	0.67	0.33	1.00	0.00	0.67	1.00	0.50
High Freq. (512 - 1024 Hz)	SW	555	1214	67	1757	0.64	0.31	0.96	0.04	0.69	0.89	0.46
	KS	279	1490	0	1824	0.59	0.16	1.00	0.00	0.84	1.00	0.27
	AD	268	1501	0	1824	0.58	0.15	1.00	0.00	0.85	1.00	0.26

Table 15: A run down of the glitch detection results across different frequency ranges for the three statistical tests being studied.

Referring to table 15, though the Shapiro-Wilk test is the best of the three at detecting glitches across all the frequency ranges in terms of accuracy and F1 Score, the accuracy being in the 64% to 74% range and F1 scores being in the 0.46 to 0.66 range alongside a significant number of false positives, especially at lower frequencies, shows us that there is still room for improvement. However, due to this being a statistical test, there is no idea of "training" unlike in the case of machine learning models. All the tests experimented with in this project are inherently limited by the nature of the data and the assumptions made about the distribution of the samples.

Though the results of these tests are not sufficient on their own, they can be used as a part of a larger system to improve the detection accuracy. A possible use for these test results would be to combine them into a feature vector to use as inputs for machine learning or deep learning models. These tests capture different aspects of the data distribution, hence making them highly viable for use as training features.

In addition, using just these three tests by themselves would not be sufficient as a feature set for training a model. These results can be used as a part of a larger feature set that includes other statistical measures, such as the variance, skewness, and kurtosis, as well as information specific to the time series data such as the SNR. Additionally, more metrics could be incorporated into the proposed feature set such as the Jarque-Bera test [34], which determines the Gaussianity of a

sample based on its Skewness and Kurtosis, and the Wasserstein distance or Wasserstein metric [35], another distance metric that can be used to compare probability distributions. This would provide a more comprehensive representation of the data, allowing us to detect the presence of glitches in our data more effectively.

6 Conclusion

In this project, we proposed a framework to detect glitches in gravitational wave detector data using statistical tests of Gaussianity on the time-amplitude domain. For our analyses, we utilized data taken from the L1 detector during the O3a run [5]. We studied the effectiveness of three statistical tests of Gaussianity, namely the Shapiro-Wilk [27], Kolmogorov-Smirnov [28], and Anderson-Darling [32] at a c , in detecting glitches in the time-amplitude domain, as opposed to the frequency-amplitude domain, which is a common approach taken in GW data analysis and glitch studies [14, 17]. We developed a preprocessing pipeline to condition the data for the tests, which included obtaining time series data for the clean and glitched data from the GWOSC database, preprocessing and conditioning the data, and applying band pass filters prior to our experiments. We then applied the three statistical tests on the preprocessed data and evaluated their performance in terms of accuracy, true positive rate, false positive rate, and F1 score.

The results of our experiments has shown that the Shapiro-Wilk test is the most effective at detecting glitches across the full frequency range of 10 Hz to 1024 Hz with an accuracy of 73.86% (rounded to 0.74 or 74%) and an F1 score of 0.66, followed by the Kolmogorov-Smirnov test at an accuracy of 64.43% (rounded to 0.64 or 64%) and F1 score of 0.43, and finally, the Anderson-Darling test at an accuracy of 63.9% (rounded to 0.64 or 64%) and F1 score of 0.42.

At lower (10 Hz to 512 Hz) and higher (512 Hz to 1024 Hz) frequency bands, the effectiveness of these tests vary, with some glitch types being detected better at higher frequencies, and others better at lower frequencies. This suggests that the success of these tests is highly dependent on the amount of power each glitch type contains at specific frequency bands. In all these cases, the tests consistently showed the same order of effectiveness as on the full frequency range, with the Shapiro-Wilk test having the highest number of glitch detections and a proportionally high number of false positives, followed by the Kolmogorov-Smirnov and Anderson-Darling tests. This hence confirms the findings of the study in [33].

The main concern with these tests is that they are limited in their capabilities for glitch detection by themselves, however they show promise as a part of a larger system, which would be a possible direction for future work.

7 Acknowledgements

This research project was conducted with the support of Dr. Sarah Caudill from the University of Massachusetts Dartmouth, and Dr. Melissa Lopez from the National Institute for Subatomic Physics (Nikhef). I am deeply grateful for their constant guidance and feedback, which was invaluable in helping me understand the nuances of not just this project but also the broader field of work I will be involved in as an aspiring PhD student. I also extend my gratitude to the faculty and staff at the University of Massachusetts Dartmouth for providing me with the resources and support to conduct this research.

This work extensively utilized computing resources from the LIGO Data Grid (LDG) and tools from the Gravitational Wave Open Science Center (<https://www.gw-openscience.org>), both being services provided by the LIGO Scientific Collaboration (LSC) and the Virgo Collaboration. I would like to thank the LIGO and Virgo Collaborations for making their data publicly available, which made this project possible.

A special thank you to my family for their constant encouragement and belief in me throughout this journey. I am the person I am today because of them, and I am forever grateful for it. Finally, I appreciate the support of my friends and peers who have been here for me through the ups and downs not only as a student, but also as someone trying to find their place in the world far away from home.

References

- [1] B. P. Abbott, R. Abbott, et al. “Observation of Gravitational Waves from a Binary Black Hole Merger”. In: *Physical Review Letters* 116.6 (Feb. 2016). ISSN: 1079-7114. DOI: [10.1103/PhysRevLett.116.061102](https://doi.org/10.1103/PhysRevLett.116.061102). URL: <http://dx.doi.org/10.1103/PhysRevLett.116.061102>.
- [2] Timothée Accadia, Fausto Acernese, et al. “Virgo: A laser interferometer to detect gravitational waves”. In: *Journal of Instrumentation* 7 (Mar. 2012), P03012. DOI: [10.1088/1748-0221/7/03/P03012](https://doi.org/10.1088/1748-0221/7/03/P03012).
- [3] Andrew Grant. “Advanced LIGO ramps up, with slight improvements”. en. In: *Physics Today* 2016.11 (Nov. 2016), p. 12128. ISSN: 19450699. DOI: [10.1063/pt.5.9074](https://doi.org/10.1063/pt.5.9074). URL: <https://pubs.aip.org/physicstoday/online/12128>.
- [4] The LIGO Scientific Collaboration, the Virgo Collaboration, 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”. In: *Physical Review X* 9.3 (Sept. 2019). arXiv:1811.12907 [astro-ph], p. 031040. ISSN: 2160-3308. DOI: [10.1103/PhysRevX.9.031040](https://doi.org/10.1103/PhysRevX.9.031040). URL: <http://arxiv.org/abs/1811.12907>.
- [5] R. Abbott et al. *GWTC-2: Compact Binary Coalescences Observed by LIGO and Virgo During the First Half of the Third Observing Run*. arXiv:2010.14527 [gr-qc]. Mar. 2021. DOI: [10.1103/PhysRevX.11.021053](https://doi.org/10.1103/PhysRevX.11.021053). URL: <http://arxiv.org/abs/2010.14527>.
- [6] The LIGO Scientific Collaboration, the Virgo Collaboration, et al. *GWTC-3: Compact Binary Coalescences Observed by LIGO and Virgo During the Second Part of the Third Observing Run*. arXiv:2111.03606 [gr-qc]. Oct. 2023. DOI: [10.1103/PhysRevX.13.041039](https://doi.org/10.1103/PhysRevX.13.041039). URL: <http://arxiv.org/abs/2111.03606>.
- [7] M. Lopez. “Exploring the Frontier of Transient Gravitational Wave Detection - Unleashing the Power of Machine Learning”. PhD thesis. Utrecht University, 2025.
- [8] The LIGO Scientific Collaboration, the Virgo Collaboration, et al. “A guide to LIGO-Virgo detector noise and extraction of transient gravitational-wave signals”. In: *Classical and Quantum Gravity* 37.5 (Mar. 2020). arXiv:1908.11170 [gr-qc], p. 055002. ISSN: 0264-9381, 1361-6382. DOI: [10.1088/1361-6382/ab685e](https://doi.org/10.1088/1361-6382/ab685e). URL: <http://arxiv.org/abs/1908.11170>.

- [9] The LIGO Scientific Collaboration and the Virgo Collaboration. “Characterization of transient noise in Advanced LIGO relevant to gravitational wave signal GW150914”. In: *Classical and Quantum Gravity* 33.13 (July 2016). arXiv:1602.03844 [gr-qc], p. 134001. ISSN: 0264-9381, 1361-6382. DOI: [10.1088/0264-9381/33/13/134001](https://doi.org/10.1088/0264-9381/33/13/134001). URL: <http://arxiv.org/abs/1602.03844>.
- [10] The LIGO Scientific Collaboration and The Virgo Collaboration. “GW170817: Observation of Gravitational Waves from a Binary Neutron Star Inspiral”. In: *Physical Review Letters* 119.16 (Oct. 2017). arXiv:1710.05832 [gr-qc], p. 161101. ISSN: 0031-9007, 1079-7114. DOI: [10.1103/PhysRevLett.119.161101](https://doi.org/10.1103/PhysRevLett.119.161101). URL: <http://arxiv.org/abs/1710.05832>.
- [11] Florent Robinet et al. “Omicron: A tool to characterize transient noise in gravitational-wave detectors”. In: *SoftwareX* 12 (2020), p. 100620. ISSN: 2352-7110. DOI: <https://doi.org/10.1016/j.softx.2020.100620>. URL: <https://www.sciencedirect.com/science/article/pii/S2352711020303332>.
- [12] Duncan M. Macleod et al. “GWpy: A Python package for gravitational-wave astrophysics”. In: *SoftwareX* 13 (2021), p. 100657. ISSN: 2352-7110. DOI: <https://doi.org/10.1016/j.softx.2021.100657>. URL: <https://www.sciencedirect.com/science/article/pii/S2352711021000029>.
- [13] D. Davis et al. “Subtracting glitches from gravitational-wave detector data during the third observing run”. In: *Classical and Quantum Gravity* 39.24 (Dec. 2022). arXiv:2207.03429 [astro-ph], p. 245013. ISSN: 0264-9381, 1361-6382. DOI: [10.1088/1361-6382/aca238](https://doi.org/10.1088/1361-6382/aca238). URL: <http://arxiv.org/abs/2207.03429>.
- [14] S. Chatterji et al. “Multiresolution techniques for the detection of gravitational-wave bursts”. In: *Classical and Quantum Gravity* 21.20 (Oct. 2004). arXiv:gr-qc/0412119, S1809–S1818. ISSN: 0264-9381, 1361-6382. DOI: [10.1088/0264-9381/21/20/024](https://doi.org/10.1088/0264-9381/21/20/024). URL: <http://arxiv.org/abs/gr-qc/0412119>.
- [15] M Zevin et al. “Gravity Spy: integrating advanced LIGO detector characterization, machine learning, and citizen science”. In: *Classical and Quantum Gravity* 34.6 (Feb. 2017), p. 064003. ISSN: 1361-6382. DOI: [10.1088/1361-6382/aa5cea](https://doi.org/10.1088/1361-6382/aa5cea). URL: <http://dx.doi.org/10.1088/1361-6382/aa5cea>.
- [16] Siddharth Soni et al. *QoQ: a Q-transform based test for Gravitational Wave transient events*. 2023. arXiv: [2305.08257](https://arxiv.org/abs/2305.08257) [gr-qc]. URL: <https://arxiv.org/abs/2305.08257>.
- [17] Leah Vazsonyi and Derek Davis. “Identifying glitches near gravitational-wave signals from compact binary coalescences using the Q-transform”. In: *Classical and Quantum Gravity* 40.3 (Feb. 2023). arXiv:2208.12338 [astro-ph], p. 035008. ISSN: 0264-9381, 1361-6382. DOI: [10.1088/1361-6382/acafd2](https://doi.org/10.1088/1361-6382/acafd2). URL: <http://arxiv.org/abs/2208.12338>.
- [18] R. Weiss et al. *Quarterly Progress Report No. 105*. Apr. 1972. URL: <https://dspace.mit.edu/handle/1721.1/56271>.
- [19] The LIGO Scientific Collaboration. *What is an Interferometer?* URL: <https://www.ligo.caltech.edu/page/what-is-interferometer>.
- [20] Sudarshan Ghonge et al. *Assessing and Mitigating the Impact of Glitches on Gravitational-Wave Parameter Estimation: a Model Agnostic Approach*. arXiv:2311.09159 [gr-qc]. Oct. 2024. DOI: [10.48550/arXiv.2311.09159](https://doi.org/10.48550/arXiv.2311.09159). URL: <http://arxiv.org/abs/2311.09159>.

- [21] Stannered. *Simple Michelson interferometer diagram*. 2007. URL: <https://commons.wikimedia.org/wiki/File:Interferometer.svg>.
- [22] WiredSense. *The Michelson Interferometer - A Laser Lab Alignment Guide*. WiredSense Tutorials. URL: <https://www.wiredsense.com/tutorial/the-michelson-interferometer-a-laser-lab-alignment-guide>.
- [23] D. M. Macleod et al. “GWpy: A Python package for gravitational-wave astrophysics”. In: *SoftwareX* 13 (2021), p. 100657. ISSN: 2352-7110. DOI: [10.1016/j.softx.2021.100657](https://doi.org/10.1016/j.softx.2021.100657). URL: <https://www.sciencedirect.com/science/article/pii/S2352711021000029>.
- [24] Curt Cutler and Eanna Flanagan. “Gravitational Waves from Mergin Compact Binaries: How Accurately Can One Extract the Binary’s Parameters from the Inspiral Waveform?”. In: *Physical Review D* 49.6 (Mar. 1994). arXiv:gr-qc/9402014, pp. 2658–2697. ISSN: 0556-2821. DOI: [10.1103/PhysRevD.49.2658](https://doi.org/10.1103/PhysRevD.49.2658). URL: <http://arxiv.org/abs/gr-qc/9402014>.
- [25] Christopher J. Moore, Robert H. Cole, and Christopher P. L. Berry. “Gravitational-wave sensitivity curves”. In: *Classical and Quantum Gravity* 32.1 (Jan. 2015). arXiv:1408.0740 [gr-qc], p. 015014. ISSN: 0264-9381, 1361-6382. DOI: [10.1088/0264-9381/32/1/015014](https://doi.org/10.1088/0264-9381/32/1/015014). URL: <http://arxiv.org/abs/1408.0740>.
- [26] Alan V. Oppenheim and Ronald W. Schaffer. *Discrete-Time Signal Processing*. 3rd. USA: Prentice Hall Press, 2009. ISBN: 0131988425.
- [27] S. S. Shapiro and M. B. Wilk. “An Analysis of Variance Test for Normality (Complete Samples)”. In: *Biometrika* 52.3/4 (1965), pp. 591–611. ISSN: 00063444, 14643510. URL: <http://www.jstor.org/stable/2333709>.
- [28] Frank J. Massey. “The Kolmogorov-Smirnov Test for Goodness of Fit”. In: *Journal of the American Statistical Association* 46.253 (1951), pp. 68–78. ISSN: 01621459, 1537274X. URL: <http://www.jstor.org/stable/2280095>.
- [29] I. M. Chakravarti, R. G. Laha, and J. Roy. *Handbook of Methods of Applied Statistics, Volume I*. John Wiley and Sons, 1967, pp. 392–394.
- [30] William F. Guthrie. *NIST/SEMATECH e-Handbook of Statistical Methods (NIST Handbook 151)*. en. 2020. DOI: [10.18434/M32189](https://doi.org/10.18434/M32189). URL: <https://www.itl.nist.gov/div898/handbook/>.
- [31] Michael J De Smith. “StatsRef.com — The sourcebook for statistics”. In: (2025). URL: <https://www.statsref.com/>.
- [32] T. W. Anderson and D. A. Darling. “A Test of Goodness of Fit”. In: *Journal of the American Statistical Association* 49.268 (1954), pp. 765–769. DOI: [10.2307/2281537](https://doi.org/10.2307/2281537).
- [33] Nornadiah Razali and Yap Bee Wah. “Power comparisons of Shapiro-Wilk, Kolmogorov-Smirnov, Lilliefors and Anderson-Darling tests”. In: *Journal of Statistical Modeling and Analytics* 2.1 (2011), pp. 21–33.
- [34] Carlos M. Jarque and Anil K. Bera. “A Test for Normality of Observations and Regression Residuals”. In: *International Statistical Review / Revue Internationale de Statistique* 55.2 (1987), pp. 163–172. ISSN: 03067734, 17515823. URL: <http://www.jstor.org/stable/1403192>.

- [35] Victor M. Panaretos and Yoav Zemel. “Statistical Aspects of Wasserstein Distances”. In: *Annual Review of Statistics and Its Application* 6.1 (Mar. 2019), pp. 405–431. ISSN: 2326-831X. DOI: [10.1146/annurev-statistics-030718-104938](https://doi.org/10.1146/annurev-statistics-030718-104938). URL: <http://dx.doi.org/10.1146/annurev-statistics-030718-104938>.