

# MLGWSC-1: The first Machine Learning Gravitational-Wave Search Mock Data Challenge

Marlin B. Schäfer<sup>1,2</sup>, Ondřej Zelenka<sup>3,4</sup>, Alexander H. Nitz<sup>1,2</sup>, He Wang<sup>5</sup>, Shichao Wu<sup>1,2</sup>, Zong-Kuan Guo<sup>5</sup>, Zhoujian Cao<sup>6</sup>, Zhixiang Ren<sup>7</sup>, Paraskevi Nousi<sup>8</sup>, Nikolaos Stergioulas<sup>9</sup>, Panagiotis Iosif<sup>10,9</sup>, Alexandra E. Koloniari<sup>9</sup>, Anastasios Tefas<sup>8</sup>, Nikolaos Passalis<sup>8</sup>, Francesco Salemi<sup>11,12</sup>, Gabriele Vedovato<sup>13</sup>, Sergey Klimenko<sup>14</sup>, Tanmaya Mishra<sup>14</sup>, Bernd Brügmann<sup>3,4</sup>, Elena Cuoco<sup>15,16,17</sup>, E. A. Huerta<sup>18,19</sup>, Chris Messenger<sup>20</sup>, Frank Ohme<sup>1,2</sup>

<sup>1</sup>*Max-Planck-Institut für Gravitationsphysik, Albert-Einstein-Institut, D-30167 Hannover, Germany*

<sup>2</sup>*Leibniz Universität Hannover, D-30167 Hannover, Germany*

<sup>3</sup>*Friedrich-Schiller-Universität Jena, D-07743 Jena, Germany*

<sup>4</sup>*Michael Stifel Center Jena, D-07743 Jena, Germany*

<sup>5</sup>*CAS Key Laboratory of Theoretical Physics, Institute of Theoretical Physics, Chinese Academy of Sciences, Beijing 100190, China*

<sup>6</sup>*Department of Astronomy, Beijing Normal University, Beijing 100875, China*

<sup>7</sup>*Peng Cheng Laboratory, Shenzhen, 518055, China*

<sup>8</sup>*Department of Informatics, Aristotle University of Thessaloniki, GR-54124 Thessaloniki, Greece*

<sup>9</sup>*Department of Physics, Aristotle University of Thessaloniki, GR-54124 Thessaloniki, Greece*

<sup>10</sup>*GSI Helmholtz Center for Heavy Ion Research, Planckstraße 1, 64291 Darmstadt, Germany*

<sup>11</sup>*Università di Trento, Dipartimento di Fisica, I-38123 Povo, Trento, Italy*

<sup>12</sup>*INFN, Trento Institute for Fundamental Physics and Applications, I-38123 Povo, Trento, Italy*

<sup>13</sup>*INFN, Sezione di Padova, I-35131 Padova, Italy*

<sup>14</sup>*Department of Physics, University of Florida, PO Box 118440, Gainesville, FL 32611-8440, USA*

<sup>15</sup>*European Gravitational Observatory (EGO), I-56021 Cascina, Pisa, Italy*

<sup>16</sup>*Scuola Normale Superiore, Piazza dei Cavalieri 7, I-56126 Pisa, Italy*

<sup>17</sup>*INFN, Sezione di Pisa, Largo Bruno Pontecorvo, 3, I-56127 Pisa, Italy*

<sup>18</sup>*Data Science and Learning Division, Argonne National Laboratory, Lemont, Illinois 60439, USA*

<sup>19</sup>*Department of Computer Science, University of Chicago, Chicago, Illinois 60637, USA and*

<sup>20</sup>*SUPA, School of Physics and Astronomy, University of Glasgow, Glasgow G12 8QQ, United Kingdom*

We present the results of the first Machine Learning Gravitational-Wave Search Mock Data Challenge (MLGWSC-1). For this challenge, participating groups had to identify gravitational-wave signals from binary black hole mergers of increasing complexity and duration embedded in progressively more realistic noise. The final of the 4 provided datasets contained real noise from the O3a observing run and signals up to a duration of 20 seconds with the inclusion of precession effects and higher order modes. We present the average sensitivity distance and runtime for the 6 entered algorithms derived from 1 month of test data unknown to the participants prior to submission. Of these, 4 are machine learning algorithms. We find that the best machine learning based algorithms are able to achieve up to 95% of the sensitive distance of matched-filtering based production analyses for simulated Gaussian noise at a false-alarm rate (FAR) of one per month. In contrast, for real noise, the leading machine learning search achieved 70%. For higher FARs the differences in sensitive distance shrink to the point where select machine learning submissions outperform traditional search algorithms at FARs  $\geq 200$  per month on some datasets. Our results show that current machine learning search algorithms may already be sensitive enough in limited parameter regions to be useful for some production settings. To improve the state-of-the-art, machine learning algorithms need to reduce the false-alarm rates at which they are capable of detecting signals and extend their validity to regions of parameter space where modeled searches are computationally expensive to run. Based on our findings we compile a list of research areas that we believe are the most important to elevate machine learning searches to an invaluable tool in gravitational-wave signal detection.

## I. INTRODUCTION

The first gravitational-wave (GW) observation on September 14, 2015 [1] achieved by the LIGO and Virgo Collaboration [2, 3] started the era of GW astronomy. During the first observing run (O1) two more GWs from coalescing binary black holes (BBHs) were detected. The second observing run (O2) saw  $\mathcal{O}(10)$  additional confident BBH detections as well as the first detection of a binary neutron star (BNS) merger [4–9]. The third observing run (O3) was split into two parts, O3a and O3b.

During O3a a further  $\mathcal{O}(40)$  BBHs as well as a second BNS merger were reported [10–12]. O3b added another  $\mathcal{O}(40)$  BBH events as well as finding the first two confident detections where the component masses are consistent with the merger of a neutron star black hole system (NSBH) [13, 14]. The fourth observing run (O4) is scheduled to begin in spring of 2023 and is expected to significantly increase the volume from which sources can be detected [15, 16].

GW signals are commonly identified in the background noise of the detectors using matched filtering [13, 17–19]. Matched filtering compares pre-computed models

of expected signals, known as templates, with the data from the detectors [20]. When a model matches the data to a pre-defined degree and data-quality requirements are met, a candidate detection is reported. Loosely modelled searches [21, 22], which look for coherent excess power in multiple detectors, are also employed by the LIGO-Virgo-KAGRA collaboration (LVK) to find potential signals.

The rate of detections has drastically increased from O1 to O3. This increase was enabled by continued detector upgrades at the two advanced LIGO observatories in Hanford and Livingston [2], as well as sensitivity improvements for the advanced Virgo detector [3]. With the entry into service of Kagra [23] a fourth observatory joined the network of ground based GW-detectors towards the end of O3. The rate of detections is expected to further increase during O4 as the sensitivity of the detectors improves and the volume from which sources can be detected grows.

With an increasing rate of detections, it is likely that systems with unexpected physical properties will be observed more frequently in the future. Optimally searching for these is a challenge for matched filtering based searches, where the computational cost scales linearly with the number of templates used. The inclusion of effects such as precession, eccentricity, or higher order modes requires millions of templates to not miss potential signals [24–26] and thus are computationally prohibitive, especially when real-time alerts should be issued. Loosely modeled searches are inherently capable of detecting arbitrary sources at a fixed computational cost but are prone to miss more signals due to their lower sensitivity in parameter regions where accurate models exist.

In recent years, machine learning has been applied in many scientific fields to enable or improve research into computationally expensive topics [27]. Some examples include the prediction of protein structure used in pharmaceutical studies [28], improvements to material composition and synthesis [29], or event reconstruction at the Large Hadron Collider [30]. There is also ongoing research into using neural networks to discover closed form expressions from raw data [31] or optimizing machine learning algorithms to take advantage of physical symmetries of the underlying problem [32–34].

More relevant to this work, machine learning algorithms have also started to be explored as alternative algorithms for many GW data-analysis tasks. These include detector glitch classification [35–37], parameter estimation [38–42], continuous GW detection [43–49], enhancements for existing pipelines [50–57], surrogate waveform models [58–60], as well as various signal detection algorithms [61–80]. For a summary of many methods we refer the reader to [81, 82]. In this work we focus solely on detection algorithms for BBH GW signals, which have been the most commonly observed type of sources to date [10–12]. These signals are the easiest to detect for machine learning algorithms due to their short duration.

Many of the works considering the usage of machine

learning for GW signal detection are difficult to cross-compare. Most algorithms target different datasets and derived metrics are often motivated more by machine learning practices than by state-of-the-art GW searches. It is, therefore, hard to pinpoint exactly how capable machine learning search algorithms currently are and where the main difficulties arise. To achieve the goal of an objective characterization of machine learning GW search capabilities, a common ground for comparison is required.

Here we present the results of the first Machine Learning Gravitational-Wave Search Mock Data Challenge (MLGWSC-1). In an attempt to provide a common ground of comparison for different algorithms and in preparation of O4, we have calculated sensitive distances from 6 different submissions calculated on datasets of one month duration to collect and compare a suite of searches. We want to motivate the utilization of machine learning based searches in a production setting by providing a definitive resource to allow for easy comparison between different algorithms, be it machine learning based, matched filtering based, or completely unmodeled. This challenge is the first of its kind<sup>1</sup> and hopefully more will be held in the future, expanding to more difficult scenarios.

The mock data used in this challenge consists of 4 datasets containing noise of increasing realism and signals with increasing complexity for the two detectors LIGO Hanford and LIGO Livingston [2]. The final dataset challenges participants to identify GWs from spinning BBHs with a duration of up to 20s added to real detector noise from O3a. The signals also take precession effects and higher order modes into account.

Submissions are evaluated on mock data of one month duration for each of the four datasets. We calculate sensitive distances for each algorithm and estimate the computational efficiency based on the runtime. The final dataset should provide an accurate picture of the possible real-world performance these algorithms can achieve. However, we note that direct comparison of the runtime performance of the different algorithms is complicated by differing hardware usage and optimization.

We find that machine learning algorithms are already competitive with state-of-the-art searches on simulated data containing injections drawn from the limited parameter space covered by this challenge. The most sensitive machine learning algorithm manages to retain  $\geq 93\%$  of the sensitive distance measured for the PyCBC pipeline [14] on Gaussian background data down to a false-alarm rate (FAR) of 1 per month. For higher FARs the separation between the approaches generally shrinks.

---

<sup>1</sup> There has previously been a public Kaggle challenge [83]. First in the sense of this paper refers to our setup of providing continuous data.

Most machine learning searches, as tested here, are less sensitive on real noise than on simulated data. The traditional algorithms handle this transition better. As a consequence, the most sensitive machine learning algorithm retains  $\geq 70\%$  of the sensitive distance of the PyCBC search down to a FAR of 1 per month. However, the sensitivity achieved of machine learning algorithms on real data is still substantial and shows that they are capable of rejecting non-Gaussian noise artifacts without any hand-tuned glitch classification.

From the evaluation of the different datasets we conclude that the main difficulties for current machine learning algorithms are the ability to analyze the consistency of detected signals between detectors and the maximum duration of signals that can be detected. Solving these issues would allow for better performance at FARs  $< 1$  per month and enable a fast detection of potentially electromagnetic bright sources such as BNS or NSBH mergers.

All code used in this challenge is open source and available at [84]. Therein we also collect the individual submissions by groups that have given their consent, provide the analysis results, and make available all plots used in this paper for all submissions.

This paper is structured as follows. In [section II](#) we provide the details on the challenge, the datasets, as well as the evaluation process. All submissions are briefly introduced in [section III](#). The results of the challenge and a brief discussion can be found in [section V](#). We conclude and give an outlook into possible future work in [section VI](#).

## II. METHODS

All submissions described in [section III](#) are evaluated on the same datasets, and all machine learning submissions are evaluated under the same conditions. Below we describe the provided material from the challenge, the requirements for the submitted algorithms, as well as the evaluation process.

### A. Challenge Resources

In this challenge participants are asked to identify GW signals submerged in detector noise. To provide grounds of comparison, all submissions are evaluated on the same datasets. To allow for optimization of the submitted algorithms for the task at hand, participants had access to code that allowed them to generate arbitrary amounts of data equivalent to that used during the final evaluation of this challenge. All code used for data generation and algorithm evaluation is open source and can be found at [84].

In particular, participants had access to the code that was used to generate the final challenge sets, but not the specific seed that was used. The specifics of the datasets are described in [section IIB](#). They were also provided

with the code that was used to generate the metrics we provide in this paper. Details on the metrics can be found in [section IIC](#).

### B. Test Data

The challenge provides a script to generate semi-continuous raw test data for any of the four datasets described below. It allows the user to choose a specific seed and a total duration of the output data. The code subsequently generates up to three files; the first containing pure noise, the second containing the same noise with injected GW signals, and the third containing the parameters of the injected signals.

The files containing the pure noise and the noise with additive signals are of the same structure. They are HDF5 [85] files with two groups named “H1” and “L1” containing data from the two detectors LIGO Hanford and LIGO Livingston, respectively. Each group consists of  $N$  HDF5-datasets, each holding the detector data of a single segment, as well as information on the GPS starting time of the segment, and its sampling rate. Each segment has a minimum duration of 2 h, is sampled at 2048 Hz, and contains continuous data. The files also contain information on the meta-data used to create the file. This meta-data is removed in the final challenge sets.

We chose to split data into smaller segments of uncorrelated noise for two reasons. First, real detectors are not equally sensitive for months at a time and data quality differs to an extent where certain data cannot be used for analyses. As such, any algorithm should be able to handle gaps in the data. Second, the noise characteristic varies over time. Segmenting simulated data allows us to easily incorporate different models for the power spectrum over the duration of the data. Subsequently, the noise model can be increased in complexity for the four datasets.

Minimal pre-processing is done on the data that is handed to the submitted algorithms. We only apply a low-frequency cutoff of 15 Hz which is used to enable a reduction in file size for real-detector data that has to be downloaded. The low-frequency cutoff reduces the dynamic range of the data, which allows us to scale the data and cast it to lower numerical precision. Any other pre-processing is left to the algorithms and is factored into the performance evaluation. The scaling is inverted during data loading.

A larger index of the dataset signifies a greater complexity and realism of the dataset. Participants may choose to optimize for any of the 4 datasets but are only allowed to submit a single algorithm, which is subsequently tested with all 4 datasets. We do this to test the ability of the search to generalize to slightly varying conditions.

Many parameters of the injected signals are drawn from the same distributions irrespective of the dataset. A summary of these distributions can be found in [Ta-](#)

Parameter	Uniform distribution
Coalescence phase	$\Phi_0 \in (0, 2\pi)$
Polarization	$\Psi \in (0, 2\pi)$
Inclination	$\cos \iota \in (-1, 1)$
Declination	$\sin \theta \in (-1, 1)$
Right ascension	$\varphi \in (-\pi, \pi)$
Chirp-Distance	$d_c^2 \in (130^2, 350^2) \text{ Mpc}^2$

TABLE I. A summary of the distributions shared between all datasets from which parameters are drawn.

ble I. All signals are generated using the waveform model `IMRPhenomXPHM` [86] with a lower frequency cutoff of 20 Hz. The waveform model was chosen for its ability to simulate both precession and higher-order modes. The merger times of two subsequent signals are separated by a random time between 24 s to 30 s to avoid any overlap. We apply a taper to the start of each waveform.

In Figure 1 we show an overview of the intrinsic parameters used in this challenge and compare it to the parameter space searched by state of the art searches [13, 14].

### 1. Dataset 1

The noise from the first dataset is purely Gaussian and simulated from the PSD model `aLIGOZeroDetHighPower` [87] for both detectors. This means that the PSD used to generate the data contains no sharp peaks originating from factors such as the power grid, is the same for all segments, and is known to the participants.

Injected signals are non-spinning and no higher-order modes are simulated. The component masses are uniformly drawn from  $10 M_\odot$  to  $50 M_\odot$ . We enforce the condition that the primary mass has to be equal or larger than the secondary mass. With this mass range, at a lower frequency cutoff of 20 Hz, and for non-spinning systems the signal duration is on the order of 1 s.

The first dataset represents a solved problem, as it has already been excessively studied in the past [61, 63, 74]. It is meant as a starting point where people new to the field can refer to existing literature to get off the ground initially. We expected many of the algorithms to perform equally well on this set.

The final challenge set for dataset 1 was generated with the seed 1 068 209 514 and a start time of 0.

### 2. Dataset 2

The noise for the second dataset is also purely Gaussian and simulated. However, in contrast to the first dataset the PSDs were derived from real data from O3a and as such contain power peaks at certain frequencies and are noisy. We generated a total of 20 PSDs for each detector. The PSDs used to generate the noise are randomly chosen from these lists and as such are unknown

to the search algorithm. The lists themselves are known to the participants. The PSDs in both detectors are independent of each other but do not change over time.

Signals are now allowed to have a spin aligned with the orbital angular momentum with a magnitude between  $-0.99$  and  $0.99$ . Additionally, the mass range is adjusted to draw component masses from the range  $7 M_\odot$  to  $50 M_\odot$ . This change increases the maximum duration of the signals at a lower frequency cutoff of 20 Hz to  $\approx 20$  s. No higher-order modes are simulated for this dataset and due to the aligned spin requirement no precession effects are present in the waveform.

The second dataset was intended to pose a considerable increase in difficulty to the first dataset. Using an unknown PSD which was derived from real data requires participants to estimate it during the analysis, if the algorithm requires it. However, we expected that increasing the signal duration to up to 20 s would be the more prominent reason for an increase in difficulty as many previous machine learning algorithms have had trouble when dealing with large inputs [39, 46, 68, 88]. Finally, we did not expect a large increase in the difficulty of the dataset due to the inclusion of aligned spins.

The final challenge set for dataset 2 was generated with the seed 2 743 406 703 and a start time of 0.

### 3. Dataset 3

The noise for the third dataset is also simulated and purely Gaussian. The increase in difficulty of the noise comes from varying the PSDs over time. Instead of choosing a single random PSD from the list of 20 PSDs per detector described in section II B 2 and generating all noise with that one PSD, the PSD for dataset 3 is randomly chosen for each segment.

The mass range from  $7 M_\odot$  to  $50 M_\odot$  and subsequently the maximum signal duration of 20 s is unchanged compared to section II B 2. However, instead of requiring the spins to be aligned with the orbital angular momentum, their orientation is isotropically distributed with a magnitude between 0 and 0.99. As a consequence, precession effects are now present in the waveforms. Additionally, we also model all higher-order  $(l, m)$ -modes available in `IMRPhenomXPHM`, which are:  $(2, 2)$ ,  $(2, -2)$ ,  $(2, 1)$ ,  $(2, -1)$ ,  $(3, 3)$ ,  $(3, -3)$ ,  $(3, 2)$ ,  $(3, -2)$ ,  $(4, 4)$ ,  $(4, -4)$  [86].

The main challenge of this dataset was intended to be the inclusion of precession effects. While these are not as impactful for short duration, high mass systems, they can substantially alter the signal morphology for lower mass systems. Adding higher-order modes can also substantially increase signal complexity. Both of these effects are currently not modeled in any production search relying on accurate signal models, as their inclusion requires an increase in size of the filter bank to include millions of templates [24, 25]. As such, we expected many if not all of the submitted algorithms to struggle with this dataset. On the other hand, any machine learning based

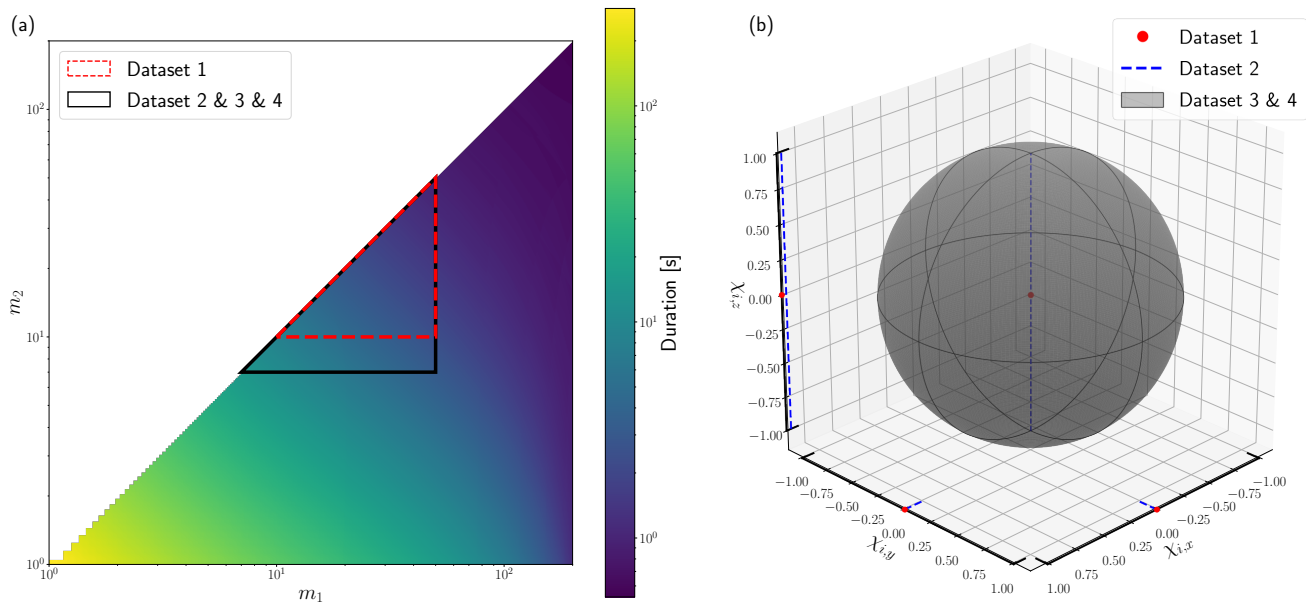


FIG. 1. An illustration of the range for the intrinsic parameters covered by this challenge. The left panel (a) shows a typical range for the component masses used by state of the art searches [14]. The color indicates the duration of the waveform from 20 Hz. The triangles show the parameter regions covered by this challenge. The right panel (b) shows the component-spin  $\chi_i$  distribution of the different datasets in this challenge.

algorithm that operates successfully on this dataset may motivate the utilization of machine learning in production searches in the future by extending the searchable parameter space.

The final challenge set for dataset 3 was generated with the seed 470 182 217 and a start time of 0.

#### 4. Dataset 4

Dataset 4 is the only dataset that contains real detector noise obtained from the Gravitational Wave Open Science Center (GWOSC) [89]. All noise was sampled from parts of O3a that had the “data” quality flag and none of the flags “CBC\_CAT1”, “CBC\_CAT2”, “CBC\_HW\_INJ”, or “BURST\_HW\_INJ” were active. We consider only segments where the data from both LIGO Hanford and LIGO Livingston clear the above conditions and excluded 10 s around any detection listed in GWTC-2 [10]. Afterwards we discarded any segments shorter in duration than 2 h. To allow for different noise realizations, we shift the data from LIGO Livingston by a random time from 0 s to 240 s while keeping the data from LIGO Hanford fixed. The time shifts are independent for each segment and to avoid any possible overlap between neighbouring segments, we consider each segment on its own.

To reduce the amount of data that has to be downloaded by participants we pre-selected the suitable parts of the O3a data. We then applied a low frequency cutoff of 15 Hz and scaled the data by a factor of  $\approx 2^{69}$ . Finally, the data was converted to single precision and stored in

a compressed format. This allowed us to provide a download link to a single file of 94 GB size containing enough data to generate up to 7 024 699  $s \approx 81$  d of coincident real noise for both detectors. The data was scaled by the constant factor to avoid the loss of dynamic range due to the conversion from double precision to single precision. When generating test data, the data is converted back to double precision and the scaling is inverted. The code used to downsample the data is also open source and available at [84].

The signals are generated equivalently to the signals in dataset 3, i.e. masses are uniformly drawn from  $7 M_{\odot}$  to  $50 M_{\odot}$ , spins are isotropically distributed with a magnitude from 0 to 0.99, and all higher-order modes available in IMRPhenomXPHM are generated. Consequently, precession effects are simulated.

This dataset is intended to be indicative of a real-world application of the search in parameter regions which are currently sparsely searched. Given that many machine learning searches have proven to generalize well from Gaussian noise to real detector noise at higher FARs in the past [62, 64, 66, 67] we expected that machine learning algorithms that do well on dataset 3 will also be competitive for dataset 4. However, it was expected that handling short glitches may prove difficult for certain searches, especially those focusing most on the merger and ringdown.

The final challenge set for dataset 4 was generated with the seed 2 514 409 456 and a start time of 0.

### C. Evaluation

All submissions are evaluated on the challenge sets, which are generated with a seed unknown to the participants at the time of submission. The evaluation is run on the Atlas computing cluster at the Albert-Einstein-Institut (AEI), Hannover. Groups that submitted an algorithm had no direct access to the evaluation stage<sup>2</sup> and final results presented in this work were only communicated back to the groups after the submission deadline had passed.

We compute two metrics for every submission and dataset. These are the wall-clock time required by the algorithm at hand to analyze one month of data as well as the sensitive distance of the search as a function of the false-alarm rate. In essence, the sensitivity as a function of the false-alarm rate is a receiver operating characteristic (ROC) curve that factors in the varying signal strengths of the injected GWs. It is a common measure of search sensitivity for production GW-searches [90] and thus allows for easy comparisons. We do not compute the ROC curve directly, for two reasons. First, it requires the number of a negative samples in the data. Since our data is continuous and the evaluation is left to the groups, defining a negative sample is not possible. Second, the ROC curve can be changed by choosing a different signal population. For instance, the ROC curve can be driven to zero by choosing a population of signals that are excessively far from the detectors. The sensitive distance normalizes the data by the injected population.

For the calculation of the sensitive distances we use two challenge sets for each of the 4 datasets. The first contains pure noise and we will call it the background set from here on out. The second contains the same noise as the background set but adds GW signals into it. This second set will be called the foreground set from here on out. As described in section II D any search algorithm is expected to process these files and return lists of events, where an event is a combination of a GPS time, a ranking statistic-like quantity, and a value for the timing accuracy. We will call these events background or foreground events when they have been derived from the background or foreground set, respectively. For the remainder of this section we will refer to the ranking statistic-like quantity simply as ranking statistic, to simplify our statements.

To calculate the sensitivity as a function of the false-alarm rate, we need to determine the false-alarm rate as a function of the ranking statistic. Next we can also determine the sensitivity as a function of the ranking statistic. Finally, we can combine the two, by evaluating both at the same values of the ranking statistic.

We use the ranking statistic of all background events as points where both the FAR as well as the sensitivity is evaluated. Each of these is certain to be a false positive and thus ensures that the FAR is unique at each threshold, as long as the search does not return identical ranking statistics for multiple background events.

To calculate the FAR at a given ranking statistic we count the number of background events with a ranking statistic greater than this threshold. We, subsequently, turn that into a rate by dividing the number of false-positives by the duration of the background data, i.e. 2592000 s. With  $N_{\text{FP},\mathcal{R}}$  the number of false-positives at a given ranking statistic  $\mathcal{R}$  and  $T$  the time spanned by the background set, the FAR  $\mathcal{F}$  can be calculated by

$$\mathcal{F} = \frac{N_{\text{FP},\mathcal{R}}}{T}. \quad (1)$$

The sensitive volume of a search at FAR  $\mathcal{F}$  can be calculated by [90]

$$V(\mathcal{F}) = \int d\mathbf{x}d\mathbf{\Lambda} \epsilon(\mathcal{F}; \mathbf{x}, \mathbf{\Lambda}) \phi(\mathbf{x}, \mathbf{\Lambda}), \quad (2)$$

where  $\mathbf{x}$  are the spatial coordinates of the injection,  $\mathbf{\Lambda}$  are the injection parameters,  $\epsilon(\mathcal{F}; \mathbf{x}, \mathbf{\Lambda})$  is the efficiency of the search at FAR  $\mathcal{F}$ , and  $\phi(\mathbf{x}, \mathbf{\Lambda})$  is the distribution of the injection parameters  $\mathbf{x}$  and  $\mathbf{\Lambda}$ .

When injections are performed uniformly in volume up to a maximum distance  $d_{\text{max}}$ , Equation 2 can be approximated by [90]

$$V(\mathcal{F}) \approx V(d_{\text{max}}) \frac{N_{I,\mathcal{F}}}{N_I}, \quad (3)$$

where  $V(d_{\text{max}})$  is the volume of a sphere with radius  $d_{\text{max}}$ ,  $N_{I,\mathcal{F}}$  is the number of found injections at a FAR of  $\mathcal{F}$ , and  $N_I$  is the total number of injections performed. An injection is found if there is at least one foreground event that is within  $\pm\Delta t$  of the injection, where  $\Delta t$  is the time variance assigned to the event by the search algorithm. The number of found injections at a given FAR considers only those foreground events where the ranking statistic assigned to the specified event is greater than the ranking statistic corresponding to the FAR. In machine learning terms Equation 3 is the recall at a given threshold on the network output multiplied by the volume of a sphere with radius  $d_{\text{max}}$ , assuming that each injection corresponds to exactly one true positive.

However, the injections in the datasets are not performed uniformly in volume, as we sample over the chirp-distance instead of the luminosity distance. The chirp-distance is given by [91]

$$d_c = d \left( \frac{\mathcal{M}_{c,0}}{\mathcal{M}_c} \right)^{5/6}, \quad (4)$$

where  $d$  is the luminosity distance,  $\mathcal{M}_c = (m_1 m_2)^{3/5} / (m_1 + m_2)^{1/5}$  is the chirp-mass, and  $\mathcal{M}_{c,0} = 1.4/2^{1/5} M_\odot$  is a fiducial chirp-mass used as a

<sup>2</sup> This excludes submissions by the organization group. However, no member of the organization group accessed the challenge-data before the submission deadline or altered their algorithm after the submission deadline.

basis for calculation. Note that in contrast to [91] we use the luminosity distance instead of the effective distance as our basis.

When sampling the injections from the distributions defined in Table I using the chirp-distance, effectively the maximum luminosity distance  $d$  is selected based on the chirp-mass; the smaller the chirp mass, the smaller the maximum luminosity distance at which injections are placed. This allows us to increase the number of detectable low mass systems and, subsequently, make statistically meaningful statements about the sensitivity for these systems without requiring a large increase in the amount of data that needs to be analyzed. However, when considering a fixed chirp mass, injections are still placed uniformly within that sphere of the adjusted maximum luminosity distance. In Equation 3 we assumed that each injection was placed uniformly within the volume spanned by the sphere with volume  $V(d_{\max})$ . To adjust it for sampling over luminosity distance we have to factor in that the probed distance depends on the selected chirp mass. We, therefore, find

$$V(\mathcal{F}) \approx \frac{V(d_{\max})}{N_I} \sum_{i=1}^{N_{I,\mathcal{F}}} \frac{V\left(d_{c,\max}\left(\frac{\mathcal{M}_{c,i}}{\mathcal{M}_{c,0}}\right)^{5/6}\right)}{V\left(d_{c,\max}\left(\frac{\mathcal{M}_{c,\max}}{\mathcal{M}_{c,0}}\right)^{5/6}\right)}, \quad (5)$$

where  $\mathcal{M}_{c,i}$  is the chirp mass of the  $i$ -th found injection,  $d_{c,\max}$  is the upper limit on the injected chirp distances, and  $\mathcal{M}_{c,\max}$  is the upper limit on the injected chirp masses. This expression can be simplified to yield

$$V(\mathcal{F}) \approx \frac{V(d_{\max})}{N_I} \sum_{i=1}^{N_{I,\mathcal{F}}} \left(\frac{\mathcal{M}_{c,i}}{\mathcal{M}_{c,\max}}\right)^{5/2}, \quad (6)$$

which is the formula we use to estimate the sensitive volume of a search algorithm. Instead of quoting the volume directly we convert it to the radius of a sphere with the corresponding volume and quote that instead.

We also measure the time the algorithm requires to evaluate an entire month of test data. Since all machine learning search algorithms are running on the same hardware these values can be used to compare the speed of the different analyses on the given hardware. For a summary of the available hardware resources please refer to Table II. However, we expect the computational time to be dominated by pre-processing steps, which can in theory be heavily optimized. For this challenge, though, we did not expect many submissions to invest resources into optimizing their pre-processing and thus advise the reader to not overemphasize the provided numbers.

All runtimes are measured twice; once for the foreground set and once for the background set. In both cases the wall-time that has passed between calling the executable and it returning is measured.

Hardware type	Specification
CPU	2× Intel Xeon Silver 4215, 8(16) cores(threads) at 2.5 GHz
GPU	8× NVIDIA RTX 2070 Super (8 GB VRAM)
RAM	192 GB

TABLE II. Main hardware specifications available to each search algorithm during final testing.

#### D. Submission Requirements

All submissions are provided with the path to a single file containing the input data they have to process. In particular they have to be able to read HDF5 files, the structure of which is detailed in section II B. Importantly, no pre-processing other than the introduction of a low frequency cutoff of 15 Hz has been applied to the data. All other pre-processing has to be performed by the algorithms themselves. In addition to the path to the input data, each algorithm is provided with a second path at which it is expected to store a single HDF5 file. This file has to contain three one-dimensional datasets of equal size named “time”, “stat”, and “var”.

The “time” dataset is expected to contain the GPS times at which the algorithm predicts a GW signal to be present. These are compared to the injection times to determine which injections were found, which were missed, and how many false positives the analysis produced.

The “stat” dataset is expected to contain a ranking-statistic like quantity for every GPS time in the “time” dataset. Here, ranking-statistic like quantity means a value where larger numbers indicate a higher degree of believe for the search to have found a GW signal. Having a ranking-statistic like quantity associated to all candidate detections enables us to assign a statistical significance to any event.

The “var” dataset is expected to contain the estimated timing accuracy of the search algorithm for all GPS times in the “time” dataset. This value determines the window around the GPS time returned by the search within which an injection has had to be made in order to consider the detection a true positive and the injection to be found. This value may be constant for all times at which the search expects to have seen a signal. We allowed searches to specify this value themselves, as we felt it to be unsuitable for a signal detection challenge to require a fixed timing accuracy. In principle, this freedom can be abused by choosing an accessively high value of  $\Delta t$  and claiming all events as true positives. However, all groups have chosen values on similar scales and more importantly far shorter than the average separation of two injections.

Throughout the paper, we will refer to events returned by the search. By that we mean a single tuple  $(t, \mathcal{R}, \Delta t)$  contained in the “time”, “stat”, and “var” datasets, re-

Python package	Version
bilby	1.1.3
pycbc	efeae6
tensorflow-gpu	2.6.0
tensorflow-probability	0.14.0
torch	1.9.1+cu11

TABLE III. A selective list of the core Python packages available to algorithms during evaluation. A complete list is given at [84].

spectively.

To be able to execute all algorithms without major problems, we ask participants to either provide a single executable that can be run on the Linux command-line utilizing only the provided software stack or to provide a singularity image that we can execute. In both cases the algorithms have to accept two positional command line arguments; the path to the input data file and the path at which the output file should be stored. The main Python packages available to submitted executables are listed in Table III, for a full list refer to [84].

Each algorithm is executed by hand and closely monitored by the organization team of the challenge. Participants are not allowed to directly tune or influence the final evaluation.

To ensure that participants have submitted the correct version of their algorithm and to make sure that their algorithm behaves as expected on the evaluation hardware and software, all algorithms are first evaluated on a validation set which is generated equivalently to the final test set. The results on this validation set are then communicated back to the submitting group. Once the group has approved that their algorithm performs within the expected margin of error, the algorithm is applied to the real challenge sets. These challenge sets are the same for all participants and were kept secret until the deadline for final submissions had passed.

Since multiple members of the organization team have submitted algorithms to this challenge, the challenge datasets were only generated after the submission deadline had passed. The script to generate test data provides an option to use a random seed. This option was used to generate the final challenge datasets and ensures that no submission had knowledge of the challenge set prior to the submission deadline.

We allowed all participants to retract their submissions at any point prior to the final publication of our results. This means that participants were allowed to retract their submissions even after they were informed about the performance of their algorithm on the final challenge sets and after they have seen the performance of other entries. No group made use of this freedom and retracted their submission after results were internally published.

### III. SUBMISSIONS

In this section we briefly introduce the different algorithms. For more details on the individual submissions we refer the reader to the original works cited within each subsection. The subsections are titled by the group name and are given in order of registration to the challenge.

All algorithm preparation was performed by the individual groups using their own available hardware resources. This crucially includes training of machine learning algorithms, for which no resources were provided by the organizers of this challenge. There were no strict requirements to submit algorithms that are based on machine learning techniques. We even encouraged the submission of a few traditional algorithms to quote a point of reference. However, the available resources detailed in section II C for evaluation of the test sets are tailored to suit the needs of machine learning algorithms.

#### A. MFCNN

<sup>3</sup> The submission of the MFCNN group is based on the works from He et al. [92]. The authors of [92] refer to the model as matched-filtering convolutional neural network (MFCNN). MFCNN is a semi-coherent search model. The basic idea of the model is to use waveform templates as learnable weights in neural network layers. Analogously to the standard coincident matched-filtering searches the output of each matched-filtering layer is maximized and normalized in the unit of matched-filtering SNRs for each GW detector. However, triggers are not generated on a single detector. The remaining part of the neural network is a usual convolutional neural network that is employed afterwards to jointly analyze the output from all detectors. Finally, a SoftMax function is applied to evaluate the confidence score of a GW signal being present in the GW detector network. The architecture was designed to take the advantages of both matched-filtering and convolutional neural networks and combine them to search for real GW events in GWTC-1 [5]. To adapt to this challenge, the source code [93] of the submission was translated from the MXNet framework [94] used in the original work to a PyTorch [95] implementation.

The training data for the model is generated by the code that generates dataset 4. The training data are input into the model directly with none of the usual pre-processing such as band-pass or whitening, which is consistent with the original work [92]. In fact, the model is equipped with a whitening layer to estimate the power spectrum for each input data. The main modification

<sup>3</sup> The corresponding authors for the MFCNN submission are He Wang, Shichao Wu, Zong-Kuan Guo, Zhoujian Cao, and Zhixiang Ren



used in this challenge is to randomly sample 25 templates in the first matched-filtering layer from the same parameter space used in dataset 4 of this challenge. It performs significantly better than the original gridded and fixed template configuration. The subsequent convolution network of the model is constructed using the current excellent lightweight models MobileNetV3 [96] which give state-of-the-art results in major computer vision problems. The submission uses curriculum learning, during which the model is trained with decreasing multiples of signal amplitude. The multiplicative factor is lowered from 50 to 1 until convergence. Multiple models were randomly initialized and trained on a NVIDIA Tesla V100 GPU, from which the best was chosen for this submission.

To search for triggers and evaluate the performance of the model, a sliding window approach is implemented. The evaluation data is divided into overlapping segments corresponding to the input size of the model. Subsequently, all segments are passed through the model resulting in a sequence of predictions and a table of SNR peaks from the 25 sorted matched filters. The step size is 1 second and a threshold of 0.5 is set on the network output as in [92]. The “time”-, “var”- and “stat”- dataset of the output file described in section IID are derived from the table of SNR peaks associated with directly filtering the templates with the data. The GPS time and time variance of each trigger are designated as the median value and the interquartile range of SNR peaks from the nearby segments, respectively. We count the coincident SNR peaks between two detectors to quantify the ranking-statistic. Other experiments are still in progress and are supposed to be published alongside further details in a standalone paper.

The final version of the algorithm submitted by the MFCNN group was provided after the submission deadline had past. A vital flaw in their original contribution was discovered and was allowed to be fixed.

## B. PyCBC

<sup>4</sup> The *PyCBC* submission is based on a standard configuration of the PyCBC-based archival search for compact-binary mergers [14]. The search infrastructure was used, in addition to cWB, for the first detection of gravitational waves, GW159014 [1], in production analyses by multiple groups to produce gravitational-wave catalogs [13, 14] and targeted analyses [97]. A similar low-latency PyCBC-Live analysis is also based around the same toolkit [18, 98]. The analysis uses matched filtering to identify candidate observations in combination with a bank of predetermined waveform templates

that correspond to the expected gravitational-wave signals [20]. Matched filtering is known to be the optimal linear filter for stationary, Gaussian noise. To account for the potential non-Gaussian noise transients [99–101], each candidate and the surrounding data are checked for consistency with the expected signal [102, 103]. In addition, the properties of candidates, such as their time of arrival, amplitude, and phases in each detector are checked for consistency with an astrophysical population [104].

The empirically measured noise distribution and the consistency with the expected gravitational-wave signal are combined to calculate a ranking statistic for each potential candidate [104, 105]; this ranking statistic is used as the “stat” value of dataset output, along with its associate trigger time in “time”. The “var” dataset is set to a constant of 0.25 s. Two template banks are used for the submitted results. For dataset 1, a template bank of non-spinning waveform templates, using the IMRPhenomD [106] model, is created using stochastic placement. Datasets 2, 3, and 4 were evaluated with a common template bank that includes templates that account for spin which is aligned with the orbital angular momentum. Furthermore, only the dominant mode of the gravitational-wave signal was used and effects such as precession were not accounted for. In both cases, the mass boundaries of the template bank conform to the challenge set parameters.

The final version of the algorithm submitted by the PyCBC group was provided after the submission deadline had past. A vital flaw in their original contribution was discovered and was allowed to be fixed. Furthermore, the PyCBC submission strictly speaking uses a different algorithm for dataset 1 than for all other datasets, as the template banks are not the same. The change in template banks was accepted, as this work does not focus on a runtime analysis.

## C. CNN-Coinc

<sup>5</sup> This submission is based on the works from Gabbard et al. [63] and Schäfer et al. [74]. It utilizes the network architecture presented in [63] with a prepended batch-normalization layer [107]. As such the network processes 8192 input samples, which corresponds to 4 s at a sampling rate of 2 kHz. The network is trained only once and applied to the data from both detectors individually. Afterwards the outputs are correlated to find coincident events as detailed in [74]. The source code for training the network and applying it to test data of the format used in this challenge is open source and can be found at [108]. The algorithm was designed to enable an easy and efficient estimation of the search background by applying time shifts between the individual detectors data.

---

<sup>4</sup> The corresponding author for the PyCBC submission is Alexander H. Nitz.

---

<sup>5</sup> The corresponding author for the CNN-Coinc submission is Martin B. Schäfer.

While this feature cannot be utilized in this challenge, the original paper [74] highlights the advantages of this approach.

The network is trained on parts of the real O3a noise from the Hanford detector as provided in this challenge. Signals are generated using the waveform approximant IMRPhenomXPHM [86] from the same parameter distribution used in datasets 3 and 4 in this challenge. Merger times of the signals are varied between 2.9 s to 3.1 s from the start of the input window of the network. The signals are pre-whitened by one of the provided Hanford PSDs used in datasets 2 and 3. Noise samples are non-overlapping parts taken from the real noise data provided by this challenge, where each segment is whitened by an estimate of the PSD on that segment. The network was trained for 100 epochs using the loss and optimizer settings provided in [74] on a single NVIDIA RTX 2070. The epoch with the greatest binary accuracy on a single training run was chosen for this challenge.

During evaluation the network is applied to the challenge-data using a sliding window approach. Each data segment is whitened by an estimate of the PSD of that segment obtained by Welch’s method [20, 109]. All data is whitened before the network is applied for computational efficiency. Subsequently, the network is applied to the data via a sliding window with a step size of 204 samples  $\approx 0.1$  s. Afterwards a threshold of 3.4 is applied on the unbounded Softmax replacement (USR) output, which was introduced in [73]. Coincident events are calculated using the same procedure and parameters as outlined in [74]. The “time”- and “stat”-dataset of the output file described in section IID list the coincident event times and ranking statistic values, respectively. The time variance of the “var”-dataset is set to a constant value of 0.3 s.

#### D. TPI FSU Jena

<sup>6</sup> This submission closely followed the method of [73], which is itself based on [63], with several modifications to adapt to the specifics of the challenge. The core of the algorithm is a convolutional neural network that accepts a  $2 \times 2048$  input tensor corresponding to 1 second of data from 2 detectors sampled at 2048 Hz. Its architecture is derived from that of [73] and deviates from the original network by a larger size of the individual layers and a doubled number of convolutional layers. These modifications are the result of a hyperparameter variation experiment which found these settings to be optimal. A standalone publication on this submission giving further details on the methodology is in preparation. The final layer of the network is a Softmax layer over two inputs

which is used for training and removed using the USR [73] during evaluation.

The network is trained on a dataset constructed by whitening a randomly chosen part of the real noise file and slicing it to produce 1-second noise samples and injecting whitened IMRPhenomXPHM-generated BBH waveforms into half the noise samples at SNRs uniformly drawn between 7 and 20. The waveform parameters are drawn from the same distributions as are used in dataset 4 of this challenge. The training dataset consists of  $10^6$  samples and the validation set of  $2 \cdot 10^5$  samples.

During evaluation, each segment in the input file is whitened separately using the estimated PSD and sliced into 1-second segments at 0.1-second spacing. These are fed to the network with the USR applied. First-level triggers are selected by applying a threshold of -8, which are then clustered into events. For each event, the “time” and “stat” in the output file are the values of the highest ranking statistic first-level trigger of each cluster, and “var” is set to 0.2 seconds. The algorithm is implemented using the PyTorch framework [95] and spawns child processes to whiten individual segments. The network evaluation is performed by the parent process.

#### E. Virgo-AUTH

<sup>7</sup> This submission is based on a simple per-dataset binary classification scheme. Interestingly, it was found that training a model only on dataset 2 or only on dataset 4 can yield impressive results on the other datasets as well. Specifically, training samples from dataset 2 can generalize well to dataset 3 and 1 and not so well on dataset 4, whereas training samples from dataset 4 can generalize well on datasets 1, 2 and 3. Thus, training samples were only generated from dataset 4. An adaptive normalization mechanism [110] was used instead of batch normalization as the first layer, to handle non-stationary timeseries. For the neural network architecture a deep, ResNet-like model [111] with a depth of 54 layers was used.

One week of training data per dataset was generated and the generated injection parameters were used to construct all corresponding waveforms. This amounted to about 600k background segments of duration 1.25 s with a stride of 2 s between, i.e. the next sample starts 0.75 s after the end of the previous one, and about 580k waveforms, of which 300k were used for the injections. For validation, one day of data was used, resulting in about 86k noise segments and 3.2k waveforms. The noise segments and waveform segments are combined online during training, in a static manner, both for the training and for the validation sets. The input samples are whitened

<sup>6</sup> The corresponding authors for the TPI FSU Jena submission are Ondřej Zelenka, Bernd Brügmann, and Frank Ohme.

<sup>7</sup> The corresponding authors for the Virgo-AUTH submission are Paraskevi Nousi, Nikolaos Stergioulas, Panagiotis Iosif, Alexandra E. Koloniari, Anastasios Tefas, and Nikolaos Passalis.

before feeding them to the classifier. The PSD is computed online per batch of 4.25 s with a stride of 3.1 s, and each 1.25 s segment inside this duration is whitened with the same PSD. To increase speed, the Welch method for computing the PSD was implemented in PyTorch [95] and whitening is implemented as the first layer of the final detection module. Notably, this approach of computing the PSD for every 4.25 s and whitening each 1.25 s segment in a sliding window manner was found to be faster than using a precomputed PSD for every 1.25 s (about 40% faster for one month of data). After whitening, the first and last 0.125 s (0.25 s total) are removed from each sample.

The best results were obtained with a ResNet-52 type network. A Deep Adaptive Input Normalization (DAIN) layer [110] was used as the first layer after whitening, to handle distribution shifts that may be present. The final output is binary, i.e., noise plus waveforms or noise only, and the objective function used was a regularized binary cross entropy. The “var” parameter is set to 0.3 s, as the network predictions are high even when the time of coalescence is slightly outside the preset range. The “stat” parameter is set to the network confidence, i.e., a value in the  $[0, 1]$  interval corresponding to the probability that a waveform is present. Finally, 0.125 s are added to the expected time of coalescence to account for the time lost in the whitening process.

A standalone publication on the methods used in this submission is in preparation.

## F. cWB

<sup>8</sup> Coherent WaveBurst (cWB) is a waveform model-agnostic search pipeline for GW signals based on the constrained likelihood method [112–114]. The cWB pipeline has been used for the analysis of scientific data collected by the LIGO-Virgo detectors, targeting detection of signals from generic GW sources, including the compact binary mergers [13].

The cWB algorithm identifies the excess-power events in the time-frequency domain representation of strain data from multiple detectors [22, 115]. For each event, the cWB pipeline reconstructs the GW waveforms and estimates summary statistics which describe generic properties of the events like the coherence across the detector network, signal strength, and the time-frequency structure.

Recently, a boosted decision-tree algorithm, eXtreme-Gradient Boost (XGBoost) [116], was adopted and implemented within the cWB framework to automate the signal-noise classification of the cWB events [54]. Two types of input data are used for the supervised training:

signal events (from simulations) and noise events (from background estimations). For each of those, a subset of cWB summary statistics is fed to XGBoost as input features to train a signal-noise model. As in [54], the detection statistic for the machine learning-enhanced cWB algorithm is defined by:

$$\eta_r = \eta_0 \cdot W_{\text{XGB}}, \quad (7)$$

where,  $\eta_0$  is cWBs ranking statistic, and  $W_{\text{XGB}}$  is the penalty factor calculated by XGBoost ranging between 0 (noise) and 1 (signal).

This methodology has been recently used in the full reanalysis of publicly available strain data from Advanced LIGO’s Hanford and Livingston third observational run [56]: the machine learning-enhanced cWB outperforms the standard human-tuned signal-noise classification used for detection of the compact binary coalescences in the O3 run.

For this study, we chose to use machine learning-enhanced cWB; however, cWB typically rejects weak candidate triggers (i.e., with  $\text{FAR} \gg 1$  per year) at early production stages. Moreover, the whole workflow is optimized for a trigger production which saturates at  $\text{FAR} \approx 30$  to 50 per month. Therefore, we modified cWB to increase the event production rate by almost 2 orders of magnitude: the result is a cWB with sub-threshold capabilities, able to speed up computation and reduce memory allocations.

While trying to provide the most “generic” result for this study, it was decided to re-use the XGBoost model which was developed for [56]: it should be noted that the model was trained on noise and signal events sets that differ substantially from those adopted for the data sets prepared for MLGWSC-1. The noise backgrounds for dataset 3 and dataset 4 appear to be significantly quieter than O3. Also, the signals were drawn from a spin-aligned stellar-mass BBHs population model with different component mass ranges [117] and with SEOBNRv4 waveforms [118]. The above-mentioned detection statistic,  $\eta_r$ , is used as the “stat” value of dataset output, along with its associated trigger peak-time in “time”. The “var” dataset is set to a constant of 0.25 s.

The results from the cWB group were provided after the submission deadline had passed. The group assured that no tuning to the challenge set was performed.

## IV. DATA RELEASE

We provide all source code as well as the evaluation results for all submissions at [84]. The repository contains all code accessible to the participants of the challenge, which most importantly includes a script to generate data and one to produce the sensitivity statistics we provide in section V. The repository also contains code for basic visualization as part of the “contributions” folder. Adaptation of these scripts were used to create the graphics in

<sup>8</sup> The corresponding authors for the cWB submission are Francesco Salemi, Gabriele Vedovato, Sergey Klimenko, Tanmaya Mishra.

this paper. The challenge used the code of release 1.3 of the repository.

Alongside the code provided by the challenge organizers we publish the source code that was used to run the contributions for the groups PyCBC, CNN-Coinc, TPI FSU Jena, and Virgo-AUTH in the “submissions” folder of [84]. The submission code for the MFCNN group can be found at [93].

All analysis output files for all submissions created by our analysis are also publicly available and are stored in the “results” folder in [84]. For each group we make available the raw output on the foreground and the background for all 4 datasets. Additionally, all timing information is available. The exception is the cWB group, for which only results on datasets 3 and 4 are available.

The repository [84] also contains plots used in this paper for all groups, including versions we have not shown here. They can be found in the “plots” folder.

## V. RESULTS AND DISCUSSION

In this section we provide the results of our evaluation process described in section II for all 6 submissions. We calculate and discuss sensitive distances, found-missed plots, and runtimes to provide a quantitative comparison between the different submissions. We specifically focus on the difference between machine learning and traditional algorithms and reason where the core differences in performance arise.

The four datasets we use in this study were chosen to answer different questions and serve different purposes. Dataset 1 was meant as an entry point to the challenge that represents a largely solved case [61, 63, 74]. We expected most submissions to perform very similarly on this dataset. The second dataset was intended as the first major step in difficulty. We expected its main challenge to be the longer duration of the injected signals, as many machine learning algorithms target shorter durations and struggle with large analysis segments [68, 88]. Dataset 3 includes precession and higher order mode effects in the injected signals that traditional, modeled searches are not optimized for<sup>9</sup> [24–26]. We wanted to test if machine learning algorithms could get closer in performance, or even outperform, the traditional searches in these regions. The intention of dataset 4 was to provide a challenge that is representative of a realistic search on real detector data and a limited parameter space. The data contains non-Gaussian noise artifacts, that can mimic GW signals [120–123], which are strongly suppressed by sophisticated algorithms in traditional searches [17, 90, 122]. Most machine learning algorithms that target real noise do not make use of such

noise-mitigation strategies and instead rely solely on the ability of the machine learning algorithm to identify noise artifacts. This approach was reported to be effective for higher FARs in the past [62, 64, 66, 67] and we were, therefore, expecting relatively minor difference between dataset 3 and dataset 4. Furthermore, most traditional algorithms use matched filtering, which is only proven to be optimal for signal recovery when the noise is stationary and Gaussian. Since neither of the two assumptions are true for real detector data, we were also interested to test if machine learning algorithms can perform better than these searches by learning a better noise representation.

### A. Sensitivities

In this subsection we discuss the sensitive distances of the different submissions, which are a measure for how many sources can be detected at any given level of certainty, i.e. at a particular FAR. They are the core metric to determine the quality of any search. We focus on the low FAR region and truncate the plot at a FAR of  $10^3$  per month. We chose this cutoff for two reasons. First, to function as a standalone search, algorithms may only report events with low FARs. State of the art pipelines send out alerts only when the FAR is smaller than  $\mathcal{O}(1)$  per month [98]. Second, for high FARs a non-negligible number of detections originate from false associations. This means that a large number of triggers that originate from random noise coincidences are close enough to an injection to be counted as true positives.

Since all machine learning submissions chose to optimize for dataset 4, results on all prior sets also test the capability of generalizing to different signal (sub)populations. Dataset 3 is a special case, as it uses the same distribution to draw the parameters of the injected signals as dataset 4. It, therefore, differs only in the noise contents and is a good test of the performance difference of different algorithms between simulated and real noise.

The results of this challenge are summarized in Figure 2 and Table IV. The four individual panels of Figure 2 show the sensitive distances as a function of the FAR for all submissions. The panels contain the results for dataset 1 to 4 from left to right and top to bottom. In Table IV we give the numeric values for the sensitive distances at three selected FAR values of 1, 10, and 100 per month for all submissions and datasets. We also provide information on the wall-clock time used to evaluate the different sets. Due to time constraints, we only show sensitivity curves for dataset 3 and 4 for the submission from the cWB group. We also note that PyCBC used a different template bank to analyze dataset 1 than for the remaining three datasets.

We find that the machine learning algorithms from the TPI FSU Jena group presented in section IIID and the Virgo-AUTH group presented in section IIIE are very

<sup>9</sup> A full search of the entire O3 data that includes higher order modes has been performed in [119].

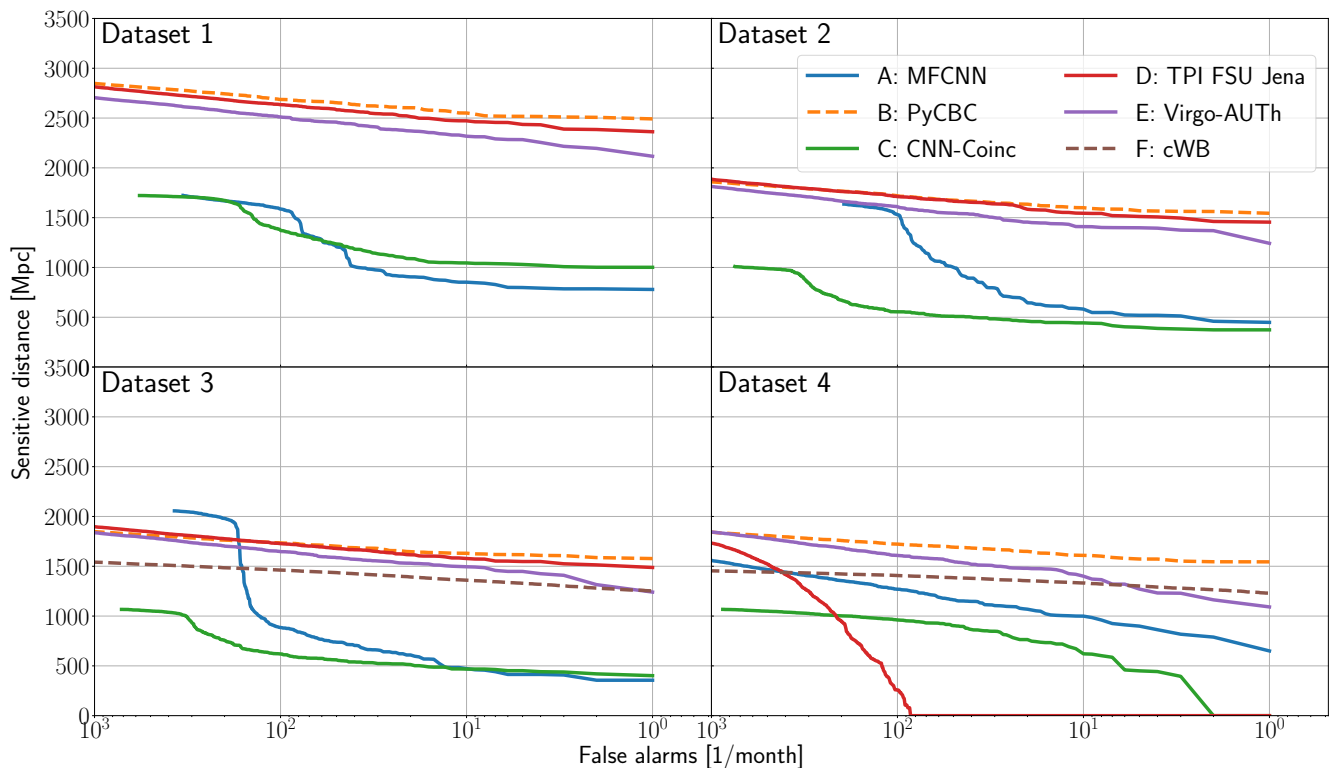


FIG. 2. The sensitive distances of all submissions and all four datasets as functions of the FAR. Submissions that made use of a machine learning algorithm at their core are shown with solid lines, others with dashed lines. The FAR was calculated on a background set that does not contain any injections.

close in sensitivity for datasets 1, 2, and 3. The submission from the TPI FSU Jena group reaches a slightly higher sensitive distance at all FARs for all of these three datasets. However, the Virgo-AUTH submission retains  $\geq 90\%$  of the sensitive distance achieved by the TPI FSU Jena submissions for FARs  $\geq 2$  per month. At lower FARs the gap widens but the individual sensitivities carry large uncertainties due to low number statistics. For higher FARs this gap narrows to a separation of roughly 4% at a FAR of 1000 per month. We suspect that the difference between the two approaches is on the order that could be explained by different initializations of the training procedure.

On dataset 4 the submission from the Virgo-AUTH group manages to maintain a stable sensitivity for the full range of tested FARs. The submission from TPI FSU Jena, on the other hand, is dominated by background triggers and seemingly struggles to adjust to the non-Gaussian noise characteristics. For high FARs the sensitivity is on a similar scale as the submission from the Virgo-AUTH group and as was observed on previous datasets, backing up the hypothesis that rejecting background triggers is the main problem. This is surprising, as both algorithms were optimized on dataset 4 but performed similarly only on datasets 1 to 3. One reason for this result may be the neural network architectures used by the different groups. The Virgo-AUTH group uses a

very deep ResNet that may be better suited to represent non-Gaussian noise artifacts. The architecture from the TPI FSU Jena group is a more straightforward convolutional architecture that may be limited in its ability to learn appropriate parameters.

The algorithms from the MFCNN group presented in [section III A](#) and the CNN-Coinc group presented in [section III C](#) also show similarities in sensitivity. Both are significantly less sensitive than the leading machine learning submission on all datasets. For datasets 1, 2, and 3, the MFCNN contribution achieves 32.5%, 30.8%, and 23.5% of the sensitive distances of the leading machine learning contribution, respectively. The CNN-Coinc submission reaches 42%, 25.5%, and 27% of the sensitivity of the leading machine learning contribution at the point of farthest separation. For dataset 4 the submission from the MFCNN and CNN-Coinc groups do comparatively better. They retain  $\geq 68\%$  and  $\geq 50\%$  of the sensitive distance of the leading machine learning submission down to a FAR of 10 per month, respectively. At a FAR of 1 per month the CNN-Coinc submission does not detect any signals, whereas the MFCNN still retains 60% of the sensitivity of the leading machine learning contribution.

On the first three datasets one can observe a steep gradient of the sensitivity curves at varying FARs for the MFCNN and CNN-Coinc submissions. At even higher FARs the curves level off again and return to a similar

Dataset	Group	Sensitivity [Mpc] at FAR = $x$ per month			Runtime [s]		
		$x = 100$	$x = 10$	$x = 1$	foreground	background	average
1	A: MFCNN	1586.90	852.18	779.21	42842	43820	43331
	B: PyCBC	2686.55	2550.57	2491.53	5406*	5092*	5249*
	C: CNN-Coinc	1372.30	1045.34	1001.55	14003	12996	13500
	D: TPI FSU Jena	2634.80	2472.31	2362.51	3758	3530	3644
	E: Virgo-AUTH	2511.95	2317.53	2116.38	5490	5520	5505
	F: cWB	N/A	N/A	N/A	N/A	N/A	N/A
2	A: MFCNN	1531.55	581.93	448.59	43431	40634	42033
	B: PyCBC	1719.98	1599.79	1543.79	157865**	161662**	159763**
	C: CNN-Coinc	554.32	443.58	373.64	14731	14976	14853
	D: TPI FSU Jena	1712.13	1544.28	1455.09	3920	3805	3862
	E: Virgo-AUTH	1608.97	1409.95	1242.37	5596	5748	5672
	F: cWB	N/A	N/A	N/A	N/A	N/A	N/A
3	A: MFCNN	885.46	472.96	355.83	37822	41251	39536
	B: PyCBC	1734.43	1630.35	1577.18	149025**	146683**	147854**
	C: CNN-Coinc	619.94	467.81	401.58	13628	14345	13986
	D: TPI FSU Jena	1727.73	1577.77	1487.67	3862	3621	3742
	E: Virgo-AUTH	1646.53	1494.98	1240.68	5450	5453	5451
	F: cWB	1461.56	1359.78	1252.09	5247***	N/A	N/A
4	A: MFCNN	1269.03	999.29	649.81	41942	46702	44322
	B: PyCBC	1722.43	1609.62	1544.33	162699**	163504**	163102**
	C: CNN-Coinc	960.70	620.85	0.00	12489	12431	12460
	D: TPI FSU Jena	257.87	0.00	0.00	3540	3487	3514
	E: Virgo-AUTH	1608.71	1400.30	1091.77	5462	5571	5516
	F: cWB	1406.88	1331.90	1229.14	4996***	N/A	N/A

TABLE IV. A summary of the analysis results for all submissions and all datasets. The columns labeled “Sensitivity” give the values for the sensitive distance at the three FARs  $10^2$  per month,  $10^1$  per month, and  $10^0$  per month rounded to the second decimal place. The values lie on the lines in Figure 2. The columns labeled “Runtime” list the time for evaluation of the foreground and background set in seconds, respectively. The runtime column labeled “average” lists the mean time obtained from evaluating the foreground and background data. Entries labeled “N/A” are not available, because they were not measured. The PyCBC times labeled with \* are only approximations. The analysis did not run on the challenge hardware but made use of a compute cluster. Shown times are the result of scaling the computational costs to 16 CPU cores. The PyCBC times labeled with \*\* are approximations obtained in the same manner as the approximations labeled with \*, but make use of a larger filter bank. The times of the cWB group marked with \*\*\* are approximations derived from dividing the CPU core-seconds reported by the search by 16 to normalize it to the challenge hardware.

slope observed at low FARs. The sudden increase leads to the MFCNN submission being more sensitive than the modeled PyCBC search by up to 15% on dataset 3 for FARs  $> 200$  per month. This behavior is not present in any of the other submissions and we were not able to find a clear explanation. However, we observe that both algorithms have different trigger rates on the foreground and background set. If the background is estimated from the foreground data only, the sensitivity of both algorithms drops sharply. All other algorithms are robust to this change. We show these sensitivity curves in Figure 3 of

the appendix.

For all datasets we compare the leading machine learning submission to the submission from PyCBC presented in section III B. We also compare it to the submission from cWB presented in section III F for datasets 3 and 4. These two are traditional, state-of-the-art search algorithms that have already been used successfully in past observation runs [1, 13, 124].

For dataset 1 we find that the machine learning search is able to achieve between 94% and 99% of the sensitivity obtained with PyCBC. These results are remark-

ably close and improve significantly on the findings from [74], which targeted a very similar dataset. However, the gap between the machine learning detection algorithm and the PyCBC search widens for lower FARs. Therefore, we expect that the PyCBC contribution will be able to attribute a substantially higher significance to many events. This is amplified by the ability of PyCBC to trivially increase the amount of data that can be used for background estimation by introducing time-slides between detectors [74, 90].

For dataset 2 the leading machine learning contribution gets even closer to the traditional algorithm from the PyCBC group. At low FARs  $\leq 20$  per month it retains  $\geq 93.5\%$  of the sensitivity achieved by the PyCBC submission. For high FARs  $\geq 200$  per month it even manages to outperform the PyCBC submission and is up to 1.5% more sensitive.

From dataset 2 to dataset 3 all submissions experience a slight increase of the measured sensitive distance. This may be surprising at first but can be explained by the distribution of the effective spin. For dataset 3 the spin orientations are distributed isotropically, which causes the average effective spin to be smaller than in dataset 2. This leads to few systems with large effective spin. The PyCBC search gains up to 3% in sensitivity at low FARs, although it loses about 1% in sensitivity at high FARs. A similar change can be observed in the submission from TPI FSU Jena. Since both the leading machine learning contribution and the PyCBC search gain similar amounts of sensitivity from dataset 2 to dataset 3 the comparison between the two does not change substantially. The submission from the TPI FSU Jena group is now up to 2.5% more sensitive at high FARs and still about 6% less sensitive at low FARs. The Virgo-AUTH, the MFCNN, and the CNN-Coinc submissions increase their sensitive distance by a larger fraction, suggesting that they benefit more from the signal population being closer to the distribution of signals in their training set. Dataset 3 is also the first dataset for which results from the cWB search are available. We find that cWB retains  $\geq 80\%$  of the sensitive distance obtained by PyCBC over all tested FARs. Subsequently the leading machine learning submission achieves a sensitive distance greater by 15% to 23% over the range of tested FARs.

For dataset 4 the leading machine learning contribution now comes from the Virgo-AUTH group. Compared to PyCBC their algorithm retains  $\geq 87\%$  of the sensitivity down to a FAR of 10 per month. For smaller FARs the sensitivity gap widens quickly. At a FAR of 1 per month the machine learning search achieves 70% of the sensitivity of PyCBC. The cWB submission evolves similarly to PyCBC and retains  $\geq 79\%$  of the sensitive distance. At high FARs the leading machine learning search manages a sensitive distance up to 27% larger than that of cWB. For low FARs the sensitive distance falls off quicker than that of cWB. At a FAR of 1 per month the cWB search is 12.5% more sensitive than the Virgo-AUTH submission. For lower FARs we expect this difference to be-

come larger, as the production level search algorithms are tuned for lower FARs than tested in this work. In comparison to the sensitivity difference on dataset 3 the machine learning submission from Virgo-AUTH does not retain as much sensitivity on real noise as the PyCBC or cWB submissions.

The results on dataset 1 demonstrate that machine learning detection algorithms are already capable of rivaling traditional search algorithms for simulated data at FARs  $\geq 1$  per month. A previous study [74] had identified the capability of machine learning searches to build an internal representation of the signal morphology as the main problem to achieve comparable sensitivities to traditional algorithms. Such a signal representation would allow the algorithms to compare detections in multiple detectors and require them to be consistent. The two leading machine learning algorithms in this challenge seem to have overcome this limitation, at least for high FAR detections.

For dataset 2 we expected machine learning searches to decline in sensitivity more strongly than traditional searches. This expectation was provoked by the short duration of data that is processed by most machine learning searches at each step. As the signals injected into dataset 2 are of longer duration than those used in dataset 1, the machine learning algorithms inherently lose some amount of sensitivity due to considering only small parts of the signal. We estimate this loss to account for at most a 1% difference in sensitivity. However, we observe the opposite effect for the two leading machine learning algorithms, which get even closer in sensitivity to the PyCBC submission compared to dataset 1. This may be caused by the distribution of signals in the training data used for the machine learning algorithms. Since both algorithms optimized for dataset 4, most signals in the training data will have non-zero spin. Therefore, the challenge set for dataset 2 is closer in nature to the training data, which may have introduced a bias that leads to higher sensitivities for spinning systems or in other words a slightly reduced sensitivity to non-spinning systems.

Dataset 3 was intended to test if machine learning searches are capable of outperforming traditional algorithms for precessing systems and signals carrying higher order mode information. We do not find substantial evidence in support of this hypothesis from the sensitivity curves. However, the challenge set 3 contains only very few signals with strong evidence for precession and higher order modes, as most signals are still relatively short. The impact on the overall sensitivity from these signals is, therefore, minor. Surprisingly, the leading machine learning search is still on par with PyCBC and manages to be significantly more sensitive even at the lowest tested FARs than the unmodeled cWB search. It must be noted that the cWB submission was not optimized for the parameter space used in this challenge. We, thus, expect this gap to narrow if more effort were to be used to tune the cWB pipeline.

The change in the relative difference in sensitivity be-

tween the PyCBC submission and the leading machine learning contribution, as well as the change in difference to the cWB submission, from dataset 3 to dataset 4 suggests that many machine learning algorithms currently used by the community are not yet capable of treating real noise as well as sophisticated traditional algorithms. We suspect that one major factor may be non-Gaussian noise artifacts that are misclassified as signals by machine learning algorithms, while the traditional searches excise them from the data or reject them on other bases. Another reason may be the non-stationary character of the noise that may lead to different sensitivities at different times. However, this would have also been a factor in dataset 3, where the PSDs used to simulate the noise change over the duration of the challenge set. However, since the leading machine learning search does retain sensitivity at all FARs it must have learned to reject most non-Gaussian noise artifacts, which is in line with expectations from studies carried out at higher FARs [62, 64, 66, 69].

## B. Found and missed injections

We generate found-missed plots for all submissions and show a few selected ones. The ones not included in this paper can be found in the associated data release [84]. These plots highlight specific areas in parameter space where the machine learning searches are already competitive and those where more work is required. Specifically, we provide plots for chirp-mass  $M_c$  versus decisive effective chirp-distance  $D_{c,\text{eff}}$ ,  $\tau_0$  versus mass-ratio  $q$ , and the effective precession spin  $\chi_p$  [125] versus inclination with respect to the line of sight  $\theta_{jn}$ . To first order  $\tau_0$  is the time to merger from the lower frequency cutoff of the waveform [126, 127]. The decisive effective chirp-distance is a measure for how strong the signal can be observed in the detector that has the worse sensitivity due to source location and orientation. The effective chirp-distance is the chirp-distance at which a source with the same intrinsic parameters and sky location but an optimal orientation would have been observed from at the same amplitude as the injected signal. The decisive effective chirp-distance is then the larger of the two effective chirp-distances from the two detectors. Therefore, the  $M_c/D_{c,\text{eff}}$  plot informs about the ability to detect signals as a function of the SNR in the detector that is less sensitive to the signal. We also include information on the ranking statistic like quantity returned for each detected event, to highlight how strongly it is correlated to the SNR. The  $\tau_0$  versus  $q$  plot highlights how well long and short duration signals are recovered. It also gives information on the mass ratio, which is an important parameter for the strength of precession effects. The main plot used to determine the impact precession effects and higher order modes have on the detectability of signals is the  $\chi_p$  versus  $\theta_{jn}$  plot.

In Figure 4 we show the found injections from dataset

1 in the  $M_c$ - $D_{c,\text{eff}}$ -plane for the PyCBC and TPI FSU Jena submissions, respectively. Both plots clearly show that closer injections are generally attributed a higher confidence to be a real signal. This indicates that the ranking statistic like quantities for both algorithms are actually correlated with the signal strength. Similar correlations can be observed for all submissions. From Figure 4 we find that the chirp-mass distribution from the TPI FSU Jena submission favors chirp-masses in the region  $M_c \in [20, 35]$ , which is not true for the PyCBC submission. We attribute this bias to the training set, which contained signals drawn from the distributions used for dataset 3 and 4. The probability distribution of the chirp-mass for these sets is shaped such that about 51% of signals are being drawn from the mass range  $20 M_\odot$  to  $35 M_\odot$ . A similar bias is not so evident for the other machine learning submissions but may be masked by other effects. The PyCBC submission uses a uniform prior on the chirp mass and thus avoids this bias.

In Figure 5 we compare the found injections from dataset 2 in the  $\tau_0$ - $q$ -plane for the PyCBC and TPI FSU Jena submissions. The plots show that the two searches are competitive in the comparable mass region and identify similar signals. The main difference between the two searches can be observed in the  $\tau_0$  distribution of two signals. Most of the signals with large values for  $\tau_0$ , i.e. long duration signals, are missed by the TPI FSU Jena submission. These crucially include many signals that the PyCBC submission identifies with relatively high confidence. Therefore, the short duration of the input windows used by the TPI FSU Jena submission still seem to be a limiting factor for the sensitivity. This limitation will likely be more severe if longer duration signals from sources like BNS or NSBH systems were considered.

In Figure 6 we compare the  $\theta_{jn}$  and  $\chi_p$  values of the injections from dataset 3 that are found by one algorithm but missed by the other. The two algorithms come from the PyCBC group and the TPI FSU Jena group. If either algorithm adapted better to signals with strong precession or higher order modes content, we would expect to see a clustering from that search in the scatter plot. However, we do not observe this clustering, which backs up our observation from the sensitivity curves that the machine learning algorithm from the TPI FSU Jena group has not learned a better representation of precessing systems or signals with higher order mode content than the modeled PyCBC search, which only includes non-precessing signals in its template bank. However, the amount of impact precession or higher order modes have on the detectability of short duration signals used in this study are small. A real test of this hypothesis would require the analysis of long duration signals.

## C. Runtimes

All runtimes in this section are given in terms of wall-clock times obtained on equivalent hardware, which is



listed in Table II. The runtimes are largely independent of the dataset for all submissions. We, therefore, discuss them only in summary. An overview of the times can be found in Table IV. They were measured by applying each algorithm to the foreground and background of each challenge set. We report the time between the algorithm call and it returning. To avoid bottlenecks, all files were transferred to the local storage of the individual compute nodes before calling the algorithm. The output was also written to said local storage and transferred back only after the algorithm returned. It should be noted that the runtimes are heavily dependent on the amount of optimization of the algorithms. The main objective for this challenge was the sensitivity and not the runtime.

The PyCBC and cWB submissions are exceptions as their runtimes were not measured on the same hardware. Instead they were run on compute clusters making heavy use of parallelized work over multiple CPUs. The times reported here are approximations by normalizing the compute time to 16 CPU cores available in the compute nodes used for this challenge. Furthermore, for the evaluation of dataset 1 PyCBC used a different template bank than those for dataset 2 to 4 was used. This bank was substantially smaller, resulting in faster evaluation. cWB times were reported to us only on the foreground data in CPU core seconds.

We find that of the machine learning algorithms the submission from the TPI FSU Jena group is the fastest, evaluating an entire month of archival data in about 1 h. It utilizes a single GPU when evaluating the network. The second fastest algorithm is the submission from the Virgo-AUTH group. It evaluates a month of data in 1.5 h on a single GPU and is thus about 50% slower than the fastest algorithm. Notably, the two fastest algorithms are also the two most sensitive machine learning searches. The algorithm from the CNN-Coinc group requires almost 4 h on a single GPU to evaluate the same amount of data but is significantly less sensitive. However, none of these algorithms are limited by the GPU performance. The differences in execution time can be mainly attributed to the difference in optimization of the pre-processing steps. The submission from the MFCNN group on the other hand does not apply any pre-processing directly. They instead use a neural network to carry out this computation. They operate on all 8 available GPUs and manage to evaluate the month of data in  $\approx 11.5$  h.

For dataset 1 the PyCBC submission has a runtime comparable to that of the submission from the Virgo-AUTH group. On all other datasets it requires roughly 43 h to evaluate the month of data. The large difference in runtime between the datasets is caused by the smaller template bank that is used only for dataset 1. Contrary to the machine learning algorithms, the PyCBC submission did not utilize GPUs and ran on CPUs only. However, PyCBC is a production level search pipeline and as such has been optimized to run on CPUs. It is not limited by the pre-processing but rather by the matched fil-

ter operation. It should be noted that PyCBC is still the most sensitive search presented here and gains in computational efficiency could be obtained by reducing the number of templates. This would effectively trade off search sensitivity for lower computational cost.

The PyCBC submission is implemented on the CPU as a GPU implementation is inherently more difficult to optimize. GPUs, on the other hand, are usually far more efficient from a cost to performance and energy to performance standpoint [128]. One advantage of machine learning algorithms is that they make use of well optimized libraries such as PyTorch [95] or TensorFlow [129] that utilize GPUs for their computations. This makes the implementation of search algorithms on GPUs relatively straightforward and allows researchers to focus on optimizing the sensitivity of their algorithm rather than having to spend time on optimizing the algorithmic implementation.

The runtimes in this challenge are measured under the assumption that the lowest required FAR is 1 per month. In a real search lower FARs are beneficial especially for rare signals. Therefore, most traditional searches are tuned to be most sensitive at FARs well below the level tested in this challenge. PyCBC for instance can extend its background by introducing time-slides [90], thereby potentially lowering the FARs of detected events. This process is a trivial operation that requires a fraction of the computational cost of the actual filtering stage. If machine learning algorithms are not specifically designed to allow for a similar approach, lowering the FARs of detections requires multiple complete re-evaluations of the time-shifted data. This would in turn lead to a linear increase in the computational cost, i.e. lowering the potential FAR of an event by an order of magnitude would lead to an order of magnitude increase in the computational cost.

## VI. CONCLUSIONS

In this paper we have presented the results of the first Machine Learning Gravitational Wave Search Mock Data Challenge (MLGWSC-1). The study compiled curves showing the sensitive distances from 4 different machine learning submissions and compared them to 2 state-of-the-art traditional search algorithms; the modeled PyCBC [90] pipeline and the unmodeled coherent wave burst search [22, 115]. We established a common dataset and means for evaluation. We hope that other researchers will continue to make use of the resources presented in this work to allow for quantitative comparisons between different machine learning approaches and to traditional filtering techniques. As research continues and machine learning search algorithms become more sensitive, we want to motivate other groups to host new challenges, focusing on other parts of parameter space or targeting different observing strategies.

The key observations of this challenge are:

1. Machine learning search algorithms are competitive in sensitivity compared to state-of-the-art searches on simulated data and the limited parameter space explored in this challenge.
2. Most of the tested machine learning algorithms struggle to effectively handle real noise, which is contaminated with non-Gaussian noise artifacts.
3. Traditional search algorithms are capable of detecting signals at lower FARs, thus making detections more confident.
4. The tested machine learning searches struggle to identify long duration signals.

Therefore, the main challenges for current machine learning searches are the operation on real noise, the confidence in detections due to comparatively high FARs, and the detection of long duration signals. The last of those three is a major hurdle to confidently detect signals from BNS and NSBH systems. Improvements in any of these fields would be beneficial. Specifically, we identify the following key research areas:

1. Improve the ability to compare signal parameters, or representations thereof, between detectors to check for consistency and reject noise artifacts.
2. Improve the ability to calculate large amounts of background, for instance by designing algorithms that can trivially evaluate time-slides of the input data.
3. Increase the duration of data that is processed by machine learning algorithms to enable the detection of long duration signals.

This challenge shows the potential of machine learning algorithms to act as GW detection pipelines. We have shown that these algorithms are competitive in a realistic scenario to state-of-the-art searches today. They operate at low computational cost and allow for a trivial implementation of the algorithms on highly efficient GPUs, rather than relying on CPUs. We believe that this work justifies more research on this topic, especially in areas where machine learning may have a tangible impact on the rapid identification of GWs.

However, we do acknowledge that the research carried out here operates on a limited parameter space. Moreover, the targeted parameter space is not the computationally expensive part of the search space of traditional searches. About 1% of the total size of the template bank used in [14] is dedicated to the area this study searches. To have the greatest impact on real searches machine learning algorithms need to be extended to target either the low mass region, where signals are long and the computational cost of matched filtering rises rapidly, or the high mass region where signals and noise artifacts are difficult to distinguish.

We also want to mention that we did not receive a submission utilizing one of the most promising neural network architectures for GW detection of the recent past. A WaveNet based architecture, that uses dilated convolutions, has been reported to do well for this kind of task [64, 67, 130]. We also did not receive submission based on many other neural network architectures that have been used in the past, such as autoencoders [72, 79, 80, 131], inception networks [46, 68], or two dimensional convolutions that analyze time-frequency decompositions [69]. We hope that some of these approaches will be adopted to the requirements of this challenge and evaluated on the datasets presented here, to allow for a quantitative comparison.

Future mock data challenges could target longer duration signals, concentrating on BNS and NSBH systems. These are potentially EM-bright and would, therefore, be of particular interest. Furthermore, these signals stem from regions of parameter space where traditional searches are computationally expensive to run. For even longer signals, sub-solar mass black holes could be targeted. Existing modeled searches in those regions make use of several million templates and are computationally limited [132]. Another avenue may be very massive systems, which can be difficult to distinguish from noise artifacts. Finally, we recommend that future mock data challenges drop the notion of a foreground and background set and only provide data files containing injections. This would eliminate further sources of error and be more true to a realistic application, where no true GW-free background exists.

## VII. ACKNOWLEDGEMENTS

We want to thank Narenraju Nagarajan and Pascal Müller for their valuable scientific input and contributions to the code of this challenge.

We acknowledge the Max Planck Gesellschaft and the Atlas cluster computing team at Albert-Einstein Institut (AEI) Hannover for support.

O.Z. thanks the Carl Zeiss Foundation for the financial support within the scope of the program line "Breakthroughs".

The MFCNN team members would like to acknowledge that the submission was supported by the Peng Cheng Laboratory Cloud Brain (No. PCL2021A13).

The Virgo-AUTH team members would like to acknowledge the support provided by the IT Center of the Aristotle University of Thessaloniki (AUTH) throughout the progress of this work, as results presented in this work have been produced, in part, using the AUTH High Performance Computing Infrastructure and Resources, and thank the COST network CA17137 "G2Net" for support. P.I. acknowledges support by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme under grant agreement No. 759253.

The cWB team gratefully acknowledges the computational resources provided by LIGO-Virgo. This material is based upon work supported by NSF’s LIGO Laboratory, which is a major facility fully funded by the National Science Foundation. This research has made use of data, software and/or web tools obtained from the Gravitational Wave Open Science Center, a service of LIGO Laboratory, the LIGO Scientific Collaboration and the Virgo Collaboration. The work by S. K. was supported by NSF Grants No. PHY 1806165 and PHY 2110060.

This publication is based upon work from COST Action CA17137, supported by COST (European Cooperation in Science and Technology).

EAH is supported by Laboratory Directed Research and Development (LDRD) funding from Argonne National Laboratory, provided by the Director, Office of Science, of the U.S. Department of Energy under Contract No. DE-AC02-06CH11357, and the U.S. National Science Foundation Awards OAC-2209892 and OAC-1931561.

CM is supported by the Science and Technology Research Council (grant no. ST/V005634/1) and the European Cooperation in Science and Technology (COST) action CA17137.

F.O. was supported by the Max Planck Society’s Independent Research Group Programme.

This research has made use of data or software ob-

tained from the Gravitational Wave Open Science Center (gw-openscience.org), a service of LIGO Laboratory, the LIGO Scientific Collaboration, the Virgo Collaboration, and KAGRA. LIGO Laboratory and Advanced LIGO are funded by the United States National Science Foundation (NSF) as well as the Science and Technology Facilities Council (STFC) of the United Kingdom, the Max-Planck-Society (MPS), and the State of Niedersachsen/Germany for support of the construction of Advanced LIGO and construction and operation of the GEO600 detector. Additional support for Advanced LIGO was provided by the Australian Research Council. Virgo is funded, through the European Gravitational Observatory (EGO), by the French Centre National de Recherche Scientifique (CNRS), the Italian Istituto Nazionale di Fisica Nucleare (INFN) and the Dutch Nikhef, with contributions by institutions from Belgium, Germany, Greece, Hungary, Ireland, Japan, Monaco, Poland, Portugal, Spain. The construction and operation of KAGRA are funded by Ministry of Education, Culture, Sports, Science and Technology (MEXT), and Japan Society for the Promotion of Science (JSPS), National Research Foundation (NRF) and Ministry of Science and ICT (MSIT) in Korea, Academia Sinica (AS) and the Ministry of Science and Technology (MoST) in Taiwan.

- 
- [1] B. P. Abbott *et al.* (LIGO Scientific, Virgo), *Phys. Rev. Lett.* **116**, 061102 (2016), [arXiv:1602.03837 \[gr-qc\]](#).
- [2] J. Aasi *et al.* (LIGO Scientific), *Class. Quant. Grav.* **32**, 074001 (2015), [arXiv:1411.4547 \[gr-qc\]](#).
- [3] F. Acernese *et al.* (VIRGO), *Class. Quantum Grav.* **32**, 024001 (2015), [arXiv:1408.3978 \[gr-qc\]](#).
- [4] B. P. Abbott *et al.* (LIGO Scientific, Virgo), *Phys. Rev. Lett.* **119**, 161101 (2017), [arXiv:1710.05832 \[gr-qc\]](#).
- [5] B. P. Abbott *et al.* (LIGO Scientific, Virgo), *Phys. Rev. X* **9**, 031040 (2019), [arXiv:1811.12907 \[astro-ph.HE\]](#).
- [6] A. H. Nitz, C. Capano, A. B. Nielsen, S. Reyes, R. White, D. A. Brown, and B. Krishnan, *Astrophys. J.* **872**, 195 (2019), [arXiv:1811.01921 \[gr-qc\]](#).
- [7] A. H. Nitz, T. Dent, G. S. Davies, S. Kumar, C. D. Capano, I. Harry, S. Mozzon, L. Nuttall, A. Lundgren, and M. Tápai, *Astrophys. J.* **891**, 123 (2020), [arXiv:1910.05331 \[astro-ph.HE\]](#).
- [8] T. Venumadhav, B. Zackay, J. Roulet, L. Dai, and M. Zaldarriaga, *Phys. Rev. D* **100**, 023011 (2019), [arXiv:1902.10341 \[astro-ph.IM\]](#).
- [9] T. Venumadhav, B. Zackay, J. Roulet, L. Dai, and M. Zaldarriaga, *Phys. Rev. D* **101**, 083030 (2020), [arXiv:1904.07214 \[astro-ph.HE\]](#).
- [10] R. Abbott *et al.* (LIGO Scientific, Virgo), *Phys. Rev. X* **11**, 021053 (2021), [arXiv:2010.14527 \[gr-qc\]](#).
- [11] R. Abbott *et al.* (LIGO Scientific, VIRGO), (2021), [arXiv:2108.01045 \[gr-qc\]](#).
- [12] A. H. Nitz, C. D. Capano, S. Kumar, Y.-F. Wang, S. Kastha, M. Schäfer, R. Dhurkunde, and M. Cabero, *Astrophys. J.* **922**, 76 (2021), [arXiv:2105.09151 \[astro-ph.HE\]](#).
- [13] R. Abbott *et al.* (LIGO Scientific, VIRGO, KAGRA), (2021), [arXiv:2111.03606 \[gr-qc\]](#).
- [14] A. H. Nitz, S. Kumar, Y.-F. Wang, S. Kastha, S. Wu, M. Schäfer, R. Dhurkunde, and C. D. Capano, (2021), [arXiv:2112.06878 \[astro-ph.HE\]](#).
- [15] B. P. Abbott *et al.* (KAGRA, LIGO Scientific, Virgo), *Living Rev. Rel.* **23**, 3 (2020).
- [16] C. Cahillane and G. Mansell, “Review of the advanced ligo gravitational wave observatories leading to observing run four,” (2022), [arXiv:2202.00847 \[gr-qc\]](#).
- [17] C. Messick *et al.*, *Phys. Rev. D* **95**, 042001 (2017), [arXiv:1604.04324 \[astro-ph.IM\]](#).
- [18] T. Dal Canton, A. H. Nitz, B. Gadre, G. S. Davies, V. Villa-Ortega, T. Dent, I. Harry, and L. Xiao, (2020), [arXiv:2008.07494 \[astro-ph.HE\]](#).
- [19] T. Adams, D. Buskulic, V. Germain, G. M. Guidi, F. Marion, M. Montani, B. Mours, F. Piergiovanni, and G. Wang, *Class. Quant. Grav.* **33**, 175012 (2016), [arXiv:1512.02864 \[gr-qc\]](#).
- [20] B. Allen, W. G. Anderson, P. R. Brady, D. A. Brown, and J. D. E. Creighton, *Phys. Rev. D* **85**, 122006 (2012), [arXiv:gr-qc/0509116](#).
- [21] S. Klimentenko, S. Mohanty, M. Rakhmanov, and G. Mitselmakher, *Phys. Rev. D* **72**, 122002 (2005), [arXiv:gr-qc/0508068](#).
- [22] S. Klimentenko *et al.*, *Phys. Rev. D* **93**, 042004 (2016), [arXiv:1511.05999 \[gr-qc\]](#).
- [23] T. Akutsu *et al.* (KAGRA), *Nature Astron.* **3**, 35 (2019), [arXiv:1811.08079 \[gr-qc\]](#).

- [24] I. Harry, S. Privitera, A. Bohé, and A. Buonanno, *Phys. Rev. D* **94**, 024012 (2016), arXiv:1603.02444 [gr-qc].
- [25] I. Harry, J. Calderón Bustillo, and A. Nitz, *Phys. Rev. D* **97**, 023004 (2018), arXiv:1709.09181 [gr-qc].
- [26] R. Dhurkunde and A. H. Nitz, (2022), arXiv:2207.14645 [astro-ph.IM].
- [27] A. M. Deiana *et al.*, (2021), arXiv:2110.13041 [cs.LG].
- [28] M. AlQuraishi, *Current Opinion in Chemical Biology* **65**, 1 (2021), mechanistic Biology \* Machine Learning in Chemical Biology.
- [29] K. T. Butler, D. W. Davies, H. Cartwright, O. Isayev, and A. Walsh, *Nature* **559**, 547 (2018).
- [30] L. Gray, T. Klijnsma, and S. Ghosh, (2020), arXiv:2003.08013 [cs.CV].
- [31] M. Cranmer, A. Sanchez-Gonzalez, P. Battaglia, R. Xu, K. Cranmer, D. Spergel, and S. Ho, (2020), arXiv:2006.11287 [cs.LG].
- [32] T. E. Smidt, M. Geiger, and B. K. Miller, *Phys. Rev. Res.* **3**, L012002 (2021), arXiv:2007.02005 [cs.LG].
- [33] M. P. Oxley, M. Ziatdinov, O. Dyck, A. R. Lupini, R. Vasudevan, and S. V. Kalinin, *npj Computational Materials* **7**, 65 (2021).
- [34] M. Dax, S. R. Green, J. Gair, M. Deistler, B. Schölkopf, and J. H. Macke, (2021), arXiv:2111.13139 [cs.LG].
- [35] M. Zevin *et al.*, *Class. Quant. Grav.* **34**, 064003 (2017), arXiv:1611.04596 [gr-qc].
- [36] D. George, H. Shen, and E. A. Huerta, (2017), arXiv:1706.07446 [gr-qc].
- [37] S. Bahaadini, V. Noroozi, N. Rohani, S. Coughlin, M. Zevin, J. R. Smith, V. Kalogera, and A. Katsaggeles, *Info. Sci.* **444**, 172 (2018).
- [38] A. J. K. Chua and M. Vallisneri, *Phys. Rev. Lett.* **124**, 041102 (2020), arXiv:1909.05966 [gr-qc].
- [39] H. Gabbard, C. Messenger, I. S. Heng, F. Tonolini, and R. Murray-Smith, *Nature Physics* (2021), 10.1038/s41567-021-01425-7.
- [40] D. Chatterjee, S. Ghosh, P. R. Brady, S. J. Kapadia, A. L. Miller, S. Nissanke, and F. Pannarale, *Astrophys. J.* **896**, 54 (2020), arXiv:1911.00116 [astro-ph.IM].
- [41] M. Dax, S. R. Green, J. Gair, J. H. Macke, A. Buonanno, and B. Schölkopf, *Phys. Rev. Lett.* **127**, 241103 (2021), arXiv:2106.12594 [gr-qc].
- [42] A. McLeod, D. Jacobs, C. Chatterjee, L. Wen, and F. Panther, (2022), arXiv:2201.11126 [astro-ph.IM].
- [43] C. Dreissigacker, R. Sharma, C. Messenger, R. Zhao, and R. Prix, *Phys. Rev. D* **100**, 044009 (2019), arXiv:1904.13291 [gr-qc].
- [44] F. Morawski, M. Bejger, and P. Cieciel\ag, *Mach. Learn. Sci. Tech.* **1**, 025016 (2020), arXiv:1907.06917 [astro-ph.IM].
- [45] A. L. Miller *et al.*, *Phys. Rev. D* **100**, 062005 (2019), arXiv:1909.02262 [astro-ph.IM].
- [46] C. Dreissigacker and R. Prix, *Phys. Rev. D* **102**, 022005 (2020), arXiv:2005.04140 [gr-qc].
- [47] B. Beheshtipour and M. A. Papa, *Phys. Rev. D* **101**, 064009 (2020), arXiv:2001.03116 [gr-qc].
- [48] B. Beheshtipour and M. A. Papa, *Phys. Rev. D* **103**, 064027 (2021), arXiv:2012.04381 [gr-qc].
- [49] T. S. Yamamoto, A. L. Miller, M. Sieniawska, and T. Tanaka, *Phys. Rev. D* **106**, 024025 (2022), arXiv:2206.00882 [gr-qc].
- [50] C. Chatterjee, L. Wen, K. Vinsen, M. Kovalam, and A. Datta, *Phys. Rev. D* **100**, 103025 (2019), arXiv:1909.06367 [astro-ph.IM].
- [51] M. Cabero, A. Mahabal, and J. McIver, *Astrophys. J. Lett.* **904**, L9 (2020), arXiv:2010.11829 [gr-qc].
- [52] S. Jadhav, N. Mukund, B. Gadre, S. Mitra, and S. Abraham, *Phys. Rev. D* **104**, 064051 (2021), arXiv:2010.08584 [gr-qc].
- [53] M. J. Williams, J. Veitch, and C. Messenger, *Phys. Rev. D* **103**, 103006 (2021), arXiv:2102.11056 [gr-qc].
- [54] T. Mishra, B. O'Brien, V. Gayathri, M. Szczepanczyk, S. Bhaumik, I. Bartos, and S. Klimenko, *Phys. Rev. D* **104**, 023014 (2021), arXiv:2105.04739 [gr-qc].
- [55] D. Lopez, V. Gayathri, A. Pai, I. S. Heng, C. Messenger, and S. K. Gupta, (2021), arXiv:2112.06608 [gr-qc].
- [56] T. Mishra, B. O'Brien, M. Szczepańczyk, G. Vedovato, S. Bhaumik, V. Gayathri, G. Prodi, F. Salemi, E. Milotti, I. Bartos, and S. Klimenko, *Phys. Rev. D* **105**, 083018 (2022).
- [57] C. McIsaac and I. Harry, (2022), arXiv:2203.03449 [gr-qc].
- [58] S. Khan and R. Green, *Phys. Rev. D* **103**, 064015 (2021), arXiv:2008.12932 [gr-qc].
- [59] P. Nousi, S.-C. Fragkouli, N. Passalis, P. Iosif, T. Apostolatos, G. Pappas, N. Stergioulas, and A. Tefas, (2021), arXiv:2107.04312 [cs.LG].
- [60] S.-C. Fragkouli, P. Nousi, N. Passalis, P. Iosif, N. Stergioulas, and A. Tefas, (2022), arXiv:2203.08434 [astro-ph.IM].
- [61] D. George and E. A. Huerta, *Phys. Rev. D* **97**, 044039 (2018), arXiv:1701.00008 [astro-ph.IM].
- [62] D. George and E. A. Huerta, *Phys. Lett. B* **778**, 64 (2018), arXiv:1711.03121 [gr-qc].
- [63] H. Gabbard, M. Williams, F. Hayes, and C. Messenger, *Phys. Rev. Lett.* **120**, 141103 (2018), arXiv:1712.06041 [astro-ph.IM].
- [64] T. D. Gebhard, N. Kilbertus, I. Harry, and B. Schölkopf, *Phys. Rev. D* **100**, 063015 (2019), arXiv:1904.08693 [astro-ph.IM].
- [65] P. G. Krastev, *Phys. Lett. B* **803**, 135330 (2020), arXiv:1908.03151 [astro-ph.IM].
- [66] P. G. Krastev, K. Gill, V. A. Villar, and E. Berger, *Phys. Lett. B* **815**, 136161 (2021), arXiv:2012.13101 [astro-ph.IM].
- [67] W. Wei, A. Khan, E. A. Huerta, X. Huang, and M. Tian, *Phys. Lett. B* **812**, 136029 (2021), arXiv:2010.15845 [gr-qc].
- [68] M. B. Schäfer, F. Ohme, and A. H. Nitz, *Phys. Rev. D* **102**, 063015 (2020), arXiv:2006.01509 [astro-ph.HE].
- [69] W. Wei and E. A. Huerta, *Phys. Lett. B* **816**, 136185 (2021), arXiv:2010.09751 [gr-qc].
- [70] W. Wei, E. A. Huerta, M. Yun, N. Loutrel, M. A. Shaikh, P. Kumar, R. Haas, and V. Kindratenko, *Astrophys. J.* **919**, 82 (2021), arXiv:2012.03963 [gr-qc].
- [71] E. A. Huerta *et al.*, *Nature Astron.* **5**, 1062 (2021), arXiv:2012.08545 [gr-qc].
- [72] F. Morawski, M. Bejger, E. Cuoco, and L. Petre, *Mach. Learn. Sci. Tech.* **2**, 045014 (2021), arXiv:2103.07688 [astro-ph.IM].
- [73] M. B. Schäfer, O. Zelenka, A. H. Nitz, F. Ohme, and B. Brüggemann, (2021), arXiv:2106.03741 [astro-ph.IM].
- [74] M. B. Schäfer and A. H. Nitz, (2021), arXiv:2108.10715 [astro-ph.IM].
- [75] A. Khan and E. A. Huerta, (2021), arXiv:2112.07669 [astro-ph.IM].
- [76] W.-H. Ruan, H. Wang, C. Liu, and Z.-K. Guo, (2021), arXiv:2111.14546 [astro-ph.IM].

- [77] P. Chaturvedi, A. Khan, M. Tian, E. A. Huerta, and H. Zheng, (2022), [arXiv:2201.11133 \[gr-qc\]](#).
- [78] G. Baltus, J. Janquart, M. Lopez, H. Narola, and J.-R. Cudell, (2022), [arXiv:2205.04750 \[gr-qc\]](#).
- [79] E. A. Moreno, B. Borzyszkowski, M. Pierini, J.-R. Vliant, and M. Spiropulu, *Mach. Learn. Sci. Tech.* **3**, 025001 (2022), [arXiv:2107.12698 \[gr-qc\]](#).
- [80] P. Bacon, A. Trovato, and M. Bejger, (2022), [arXiv:2205.13513 \[gr-qc\]](#).
- [81] E. Cuoco *et al.*, *Mach. Learn. Sci. Tech.* **2**, 011002 (2021), [arXiv:2005.03745 \[astro-ph.HE\]](#).
- [82] E. A. Huerta and Z. Zhao, (2021), [10.1007/978-981-15-4702-7\\_47-1](#), [arXiv:2105.06479 \[astro-ph.IM\]](#).
- [83] E. G. O. EGO, “G2net gravitational wave detection kaggle challenge,” (2021).
- [84] M. Schäfer and O. Zelenka, “Mlgwsc-1 github repository,” (2021).
- [85] The HDF Group, “Hierarchical Data Format, version 5,” (1997-2022), <https://www.hdfgroup.org/HDF5/>.
- [86] G. Pratten *et al.*, *Phys. Rev. D* **103**, 104056 (2021), [arXiv:2004.06503 \[gr-qc\]](#).
- [87] L. S. Collaboration, “LIGO Algorithm Library - LAL-Suite,” free software (GPL) (2018).
- [88] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning* (MIT Press, 2016) <http://www.deeplearningbook.org>.
- [89] R. Abbott *et al.* (LIGO Scientific, Virgo), *SoftwareX* **13**, 100658 (2021), [arXiv:1912.11716 \[gr-qc\]](#).
- [90] S. A. Usman *et al.*, *Class. Quant. Grav.* **33**, 215004 (2016), [arXiv:1508.02357 \[gr-qc\]](#).
- [91] B. Abbott *et al.* (LIGO Scientific), *Phys. Rev. D* **77**, 062002 (2008), [arXiv:0704.3368 \[gr-qc\]](#).
- [92] H. Wang, S. Wu, Z. Cao, X. Liu, and J.-Y. Zhu, *Phys. Rev. D* **101**, 104003 (2020), [arXiv:1909.13442 \[astro-ph.IM\]](#).
- [93] H. Wang, “Mfcnn docker image,” (2022).
- [94] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang, “Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems,” (2015).
- [95] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, in *Advances in Neural Information Processing Systems 32*, edited by H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett (Curran Associates, Inc., 2019) pp. 8024–8035.
- [96] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan, Q. V. Le, and H. Adam, “Searching for mobilenetv3,” (2019).
- [97] A. H. Nitz, A. Lenon, and D. A. Brown, *Astrophys. J.* **890**, 1 (2019), [arXiv:1912.05464 \[astro-ph.HE\]](#).
- [98] A. H. Nitz, T. Dal Canton, D. Davis, and S. Reyes, *Phys. Rev. D* **98**, 024050 (2018), [arXiv:1805.11174 \[gr-qc\]](#).
- [99] L. Nuttall *et al.*, *Class. Quant. Grav.* **32**, 245005 (2015), [arXiv:1508.07316 \[gr-qc\]](#).
- [100] M. Cabero *et al.*, *Class. Quant. Grav.* **36**, 15 (2019), [arXiv:1901.05093 \[physics.ins-det\]](#).
- [101] D. Davis and M. Walker, *Galaxies* **10**, 12 (2022).
- [102] B. Allen, *Phys. Rev. D* **71**, 062001 (2005), [arXiv:gr-qc/0405045](#).
- [103] A. H. Nitz, *Class. Quant. Grav.* **35**, 035016 (2018), [arXiv:1709.08974 \[gr-qc\]](#).
- [104] A. H. Nitz, T. Dent, T. Dal Canton, S. Fairhurst, and D. A. Brown, *Astrophys. J.* **849**, 118 (2017), [arXiv:1705.01513 \[gr-qc\]](#).
- [105] G. S. Davies, T. Dent, M. Tápai, I. Harry, C. McIsaac, and A. H. Nitz, *Phys. Rev. D* **102**, 022004 (2020), [arXiv:2002.08291 \[astro-ph.HE\]](#).
- [106] S. Husa, S. Khan, M. Hannam, M. Pürrer, F. Ohme, X. Jiménez Forteza, and A. Bohé, *Phys. Rev. D* **93**, 044006 (2016), [arXiv:1508.07250 \[gr-qc\]](#).
- [107] S. Ioffe and C. Szegedy, in *International conference on machine learning* (PMLR, 2015) pp. 448–456, [arXiv:1502.03167 \[cs.LG\]](#).
- [108] M. Schäfer, “Cnn-coinc github repository,” (2022).
- [109] P. Welch, *IEEE Transactions on Audio and Electroacoustics* **15**, 70 (1967).
- [110] N. Passalis, A. Tefas, J. Kannianen, M. Gabbouj, and A. Iosifidis, *IEEE Transactions on Neural Networks and Learning Systems* **31**, 3760 (2019).
- [111] K. He, X. Zhang, S. Ren, and J. Sun, in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016) pp. 770–778.
- [112] M. Drago, S. Klimentko, C. Lazzaro, E. Milotti, G. Mitselmakher, V. Necula, B. O’Brian, G. A. Prodi, F. Salemi, M. Szczepanczyk, S. Tiwari, V. Tiwari, G. V. G. Vedovato, and I. Yakushin, *SoftwareX* **14**, 100678 (2021).
- [113] S. Klimentko, G. Vedovato, V. Necula, F. Salemi, M. Drago, R. Poulton, E. Chassande-Mottin, V. Tiwari, C. Lazzaro, B. O’Brian, M. Szczepanczyk, S. Tiwari, and V. Gayathri, “cwb pipeline library: 6.4.1,” (2021).
- [114] cWB Development Team, “coherent WaveBurst Homepage,” (2021).
- [115] S. Klimentko and G. Mitselmakher, *Class. Quant. Grav.* **21**, S1819 (2004).
- [116] T. Chen and C. Guestrin, in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’16 (ACM, New York, NY, USA, 2016) pp. 785–794.
- [117] C. Talbot and E. Thrane, *The Astrophysical Journal* **856**, 173 (2018).
- [118] A. Bohé *et al.*, *Physical Review D* **95** (2017), [10.1103/physrevd.95.044028](#).
- [119] K. Chandra, J. Calderón Bustillo, A. Pai, and I. Harry, (2022), [arXiv:2207.01654 \[gr-qc\]](#).
- [120] J. Aasi *et al.* (LIGO Scientific, VIRGO), *Class. Quant. Grav.* **32**, 115012 (2015), [arXiv:1410.7764 \[gr-qc\]](#).
- [121] B. P. Abbott *et al.* (LIGO Scientific, Virgo), *Class. Quant. Grav.* **33**, 134001 (2016), [arXiv:1602.03844 \[gr-qc\]](#).
- [122] B. P. Abbott *et al.* (LIGO Scientific, Virgo), *Class. Quant. Grav.* **37**, 055002 (2020), [arXiv:1908.11170 \[gr-qc\]](#).
- [123] D. Davis, M. Trevor, S. Mozzon, and L. K. Nuttall, (2022), [arXiv:2204.03091 \[gr-qc\]](#).
- [124] B. P. Abbott *et al.* (LIGO Scientific, Virgo), *Phys. Rev. D* **93**, 122004 (2016), [Addendum: *Phys. Rev. D* **94**, 069903 (2016)], [arXiv:1602.03843 \[gr-qc\]](#).
- [125] P. Schmidt, F. Ohme, and M. Hannam, *Phys. Rev. D* **91**, 024043 (2015), [arXiv:1408.1810 \[gr-qc\]](#).
- [126] T. Cokelaer, *Phys. Rev. D* **76**, 102004 (2007), [arXiv:0706.4437 \[gr-qc\]](#).
- [127] M. Maggiore, *Gravitational Waves. Vol. 1: Theory and*

- Experiments*, Oxford Master Series in Physics (Oxford University Press, 2007).
- [128] R. Dhurkunde, H. Fehrmann, and A. H. Nitz, *Phys. Rev. D* **105**, 103001 (2022), [arXiv:2110.13115 \[astro-ph.IM\]](#).
- [129] M. Abadi and et al., “TensorFlow: Large-scale machine learning on heterogeneous systems,” (2015), software available from [tensorflow.org](#).
- [130] A. Schmitt, K. Fu, S. Fan, and Y. Luo, in *Proceedings of the 2nd International Conference on Computer Science and Software Engineering*, CSSE 2019 (Association for Computing Machinery, New York, NY, USA, 2019) p. 73–78.
- [131] H. Shen, D. George, E. A. Huerta, and Z. Zhao, (2019), [10.1109/ICASSP.2019.8683061](#), [arXiv:1903.03105 \[astro-ph.CO\]](#).
- [132] A. H. Nitz and Y.-F. Wang, (2022), [arXiv:2202.11024 \[astro-ph.HE\]](#).

#### Appendix A: Sensitivity plot foreground

#### Appendix B: Found-Missed plots

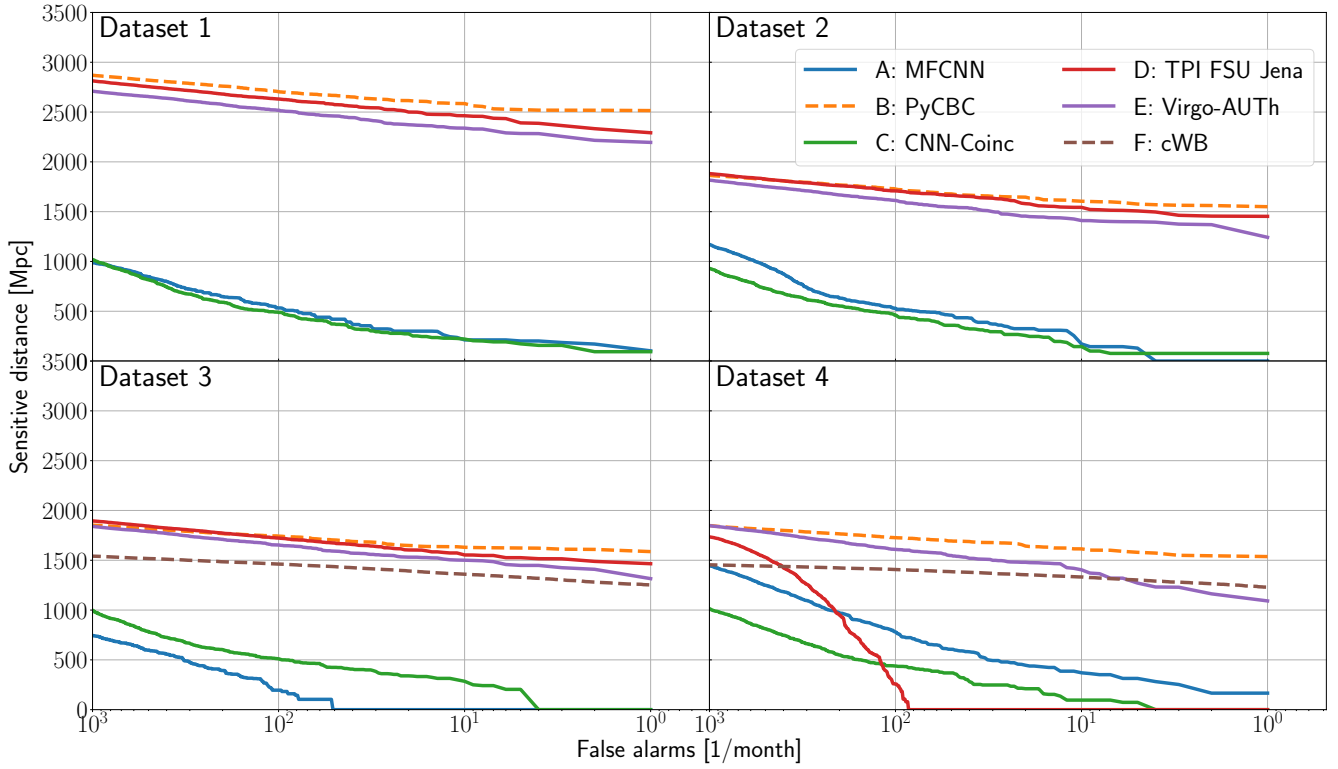
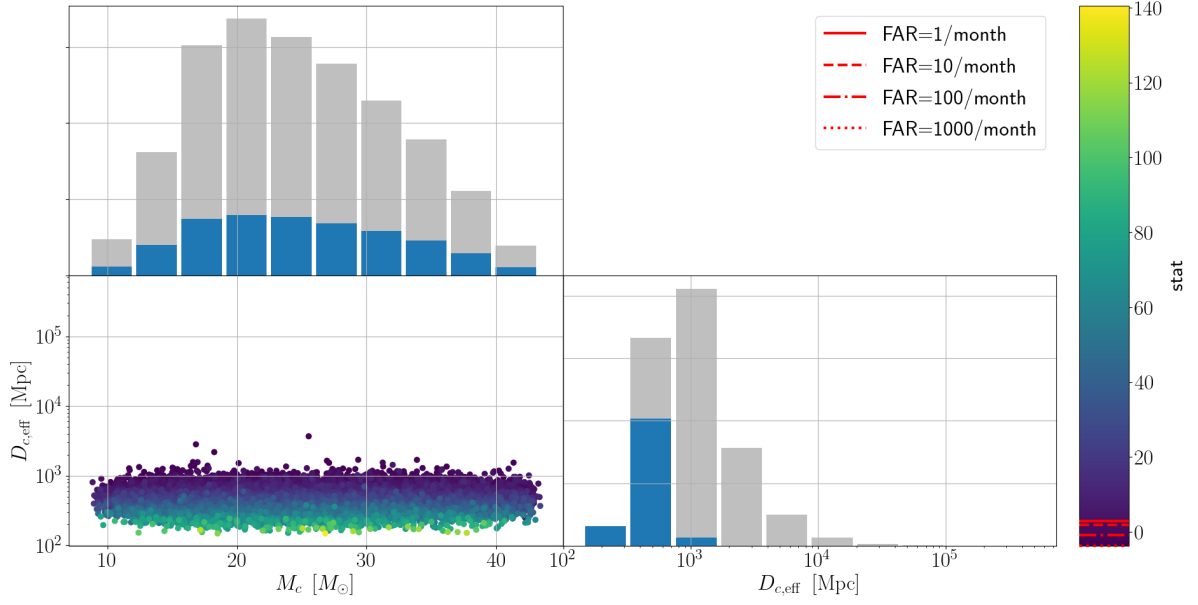


FIG. 3. The sensitive distances of all submissions and all four datasets as functions of the FAR. The sensitive distances are calculated using only the data from the foreground file. The FAR is determined from the false positives on that data. Submissions that made use of a machine learning algorithm at their core are shown with solid lines, others with dashed lines. This figure differs from Figure 2 as the algorithms from MFCNN and CNN-Coinc behave differently on the foreground and the background.

B: PyCBC Dataset 1



D: TPI FSU Jena Dataset 1

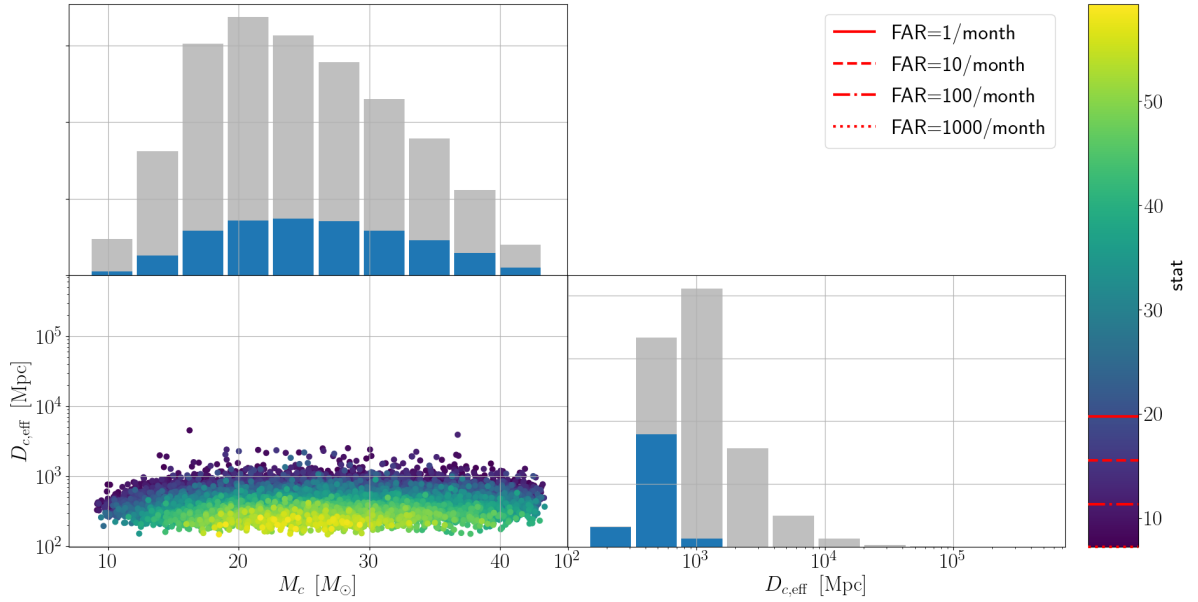


FIG. 4. The injections from dataset 1 identified by the PyCBC and TPI FSU Jena submissions with a FAR  $\leq 10^3$  per month in the chirp-mass  $M_c$  versus decisive effective chirp-distance  $D_{c,eff}$  plane. The blue bars in the histograms show the one dimensional marginal distributions of the found injections. The gray bars show the distribution of injected signals, including those missed by the search. The color shows the “stat” value attributed to the injection by the algorithm. The red lines in the colorbar highlight the thresholds on the “stat” to achieve different FARs.



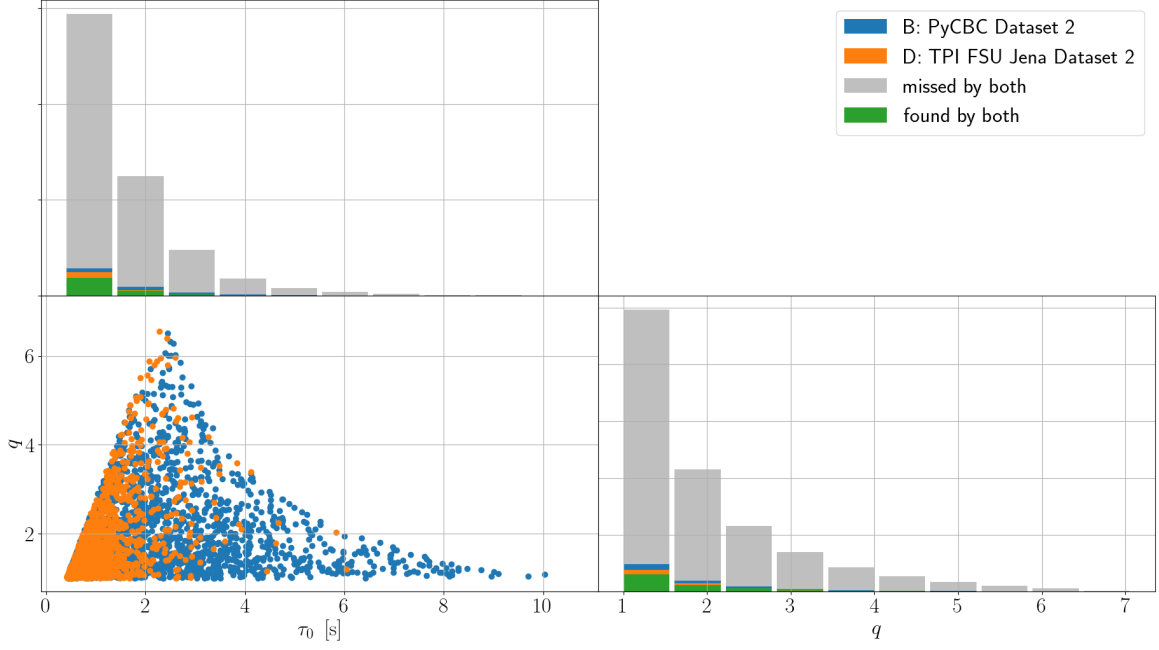


FIG. 5. The injections from dataset 2 identified by the PyCBC and TPI FSU Jena submissions with a FAR  $\leq 10^3$  per month in the signal duration  $\tau_0$  versus mass ratio  $q$  plane. The scatter plot shows injections that are found only by one of the two algorithms. Injections that are missed or found by both are only shown in the 1D marginal distributions.

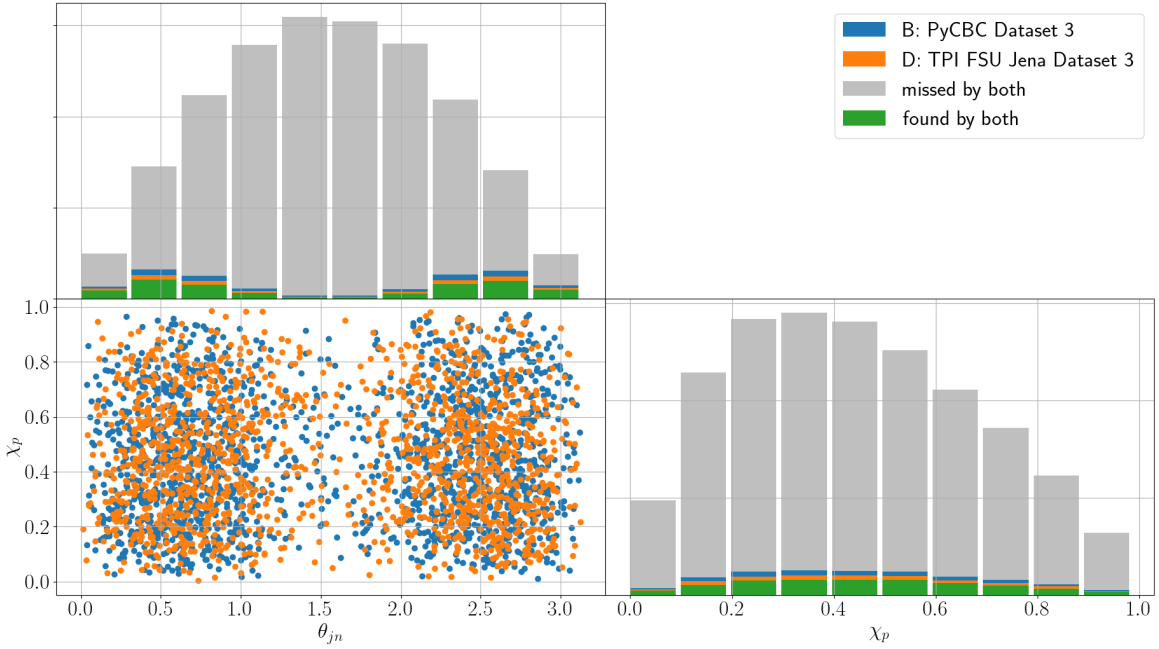


FIG. 6. The injections from dataset 3 identified by the PyCBC and TPI FSU Jena submissions with a FAR  $\leq 10^3$  per month in the inclination to spin axis  $\theta_{jn}$  to  $\chi_p$  plane. The scatter plot shows injections that are found only by one of the two algorithms. Injections that are missed or found by both are only shown in the 1D marginal distributions.