# A machine-learning pipeline for real-time detection of gravitational waves from compact binary coalescences

**Ethan Marx**

`emarx@mit.edu`

Massachusetts Institute of Technology

**William Benoit**

University of Minnesota

**Alec Gunny**

Massachusetts Institute of Technology

**Rafia Omer**

University of Minnesota

**Deep Chatterjee**

Massachusetts Institute of Technology

**Ricco Venterea**

Cornell University

**Lauren Wills**

University of Minnesota

**Muhammed Saleem**

University of Minnesota

**Eric Moreno**

Massachusetts Institute of Technology

**Ryan Raikman**

Carnegie Mellon University

**Ekaterina Govorkova**

Massachusetts Institute of Technology

**Dylan Rankin**

University of Pennsylvania

**Michael Coughlin**

University of Minnesota   https://orcid.org/0000-0002-8262-2924

**Philip Harris**

Massachusetts Institute of Technology

**Erik Katsavounidis**

**Additional Declarations:** There is **NO** Competing Interest.

# A machine-learning pipeline for real-time detection of gravitational waves from compact binary coalescences

Ethan Marx[1,2]*, William Benoit[3]*, Alec Gunny[1,2]*, Rafia Omer[3], Deep Chatterjee[1,2], Ricco C. Venterea[3,5],

Lauren Wills[3], Muhammed Saleem[3], Eric Moreno[1,2], Ryan Raikman[1,6], Ekaterina Govorkova[1,2],

Dylan Rankin[4], Michael W. Coughlin[3], Philip Harris[1], Erik Katsavounidis[1,2]

[1] *Department of Physics, MIT, Cambridge, MA 02139, USA*

[2] *LIGO Laboratory, 185 Albany St, MIT, Cambridge, MA 02139, USA*

[3] *School of Physics and Astronomy, University of Minnesota, Minneapolis, Minnesota 55455, USA*

[4] *Department of Physics and Astronomy, University of Pennsylvania, Philadelphia, PA, 19104, USA*

[5] *Department of Astronomy, Cornell University, Ithaca, NY, 14853, USA*

[6] *Department of Physics, Carnegie Mellon University, Pittsburgh, PA, 15213*

**The promise of multi-messenger astronomy relies on the rapid detection of gravitational waves at very low latencies ($\mathcal{O}(1\,\mathrm{s})$) in order to maximize the amount of time available for follow-up observations. In recent years, neural-networks have demonstrated robust non-linear modeling capabilities and millisecond-scale inference at a comparatively small computational footprint, making them an attractive family of algorithms in this context. However, integration of these algorithms into the gravitational-wave astrophysics research ecosystem has proven non-trivial. Here, we present the first fully machine learning-based pipeline for the detection of gravitational waves from compact binary coalescences (CBCs) running in low-latency. We demonstrate this pipeline to have a fraction of the latency of traditional matched filtering search pipelines while achieving state-of-the-art sensitivity to higher-mass stellar binary black holes.**

Gravitational-wave astronomy has developed rapidly since the first direct detection of gravitational waves from a binary black hole merger in 2015[1], with new detections now a common occurrence[2]. With the fourth observing run (O4) of the International Gravitational-wave Network (IGWN), consisting of LIGO[3], Virgo[4], and KAGRA[5] already underway, and with future ground and space based detectors planned for various points in the next decade[6–8], ever more frequent discoveries of gravitational waves will enable follow-up observation of events across other cosmic messengers

---

*These authors contributed equally to this work

such as electromagnetic radiation and astrophysical neutrinos[9–14]. The insights we gain in this era of multi-messenger astrophysics will directly correlate with the volume and diversity of data we are able to collect.

While machine learning (ML) is ubiquitous in some areas of physics[15], it has only recently approached a stage of maturity in the gravitational-wave community. To date, there have been a number of machine learning models proposed for the detection of compact binary coalescences (CBCs); e.g.,[16–20]; but there are none currently running in O4[21] (though, ML-based unmodeled gravitational-wave searches have seen production usage[22]). This is both a product of well-known infrastructure hurdles separating the development and deployment of machine learning models[23], as well as a lack of standardized, astrophysically meaningful probes of the sensitivity of these models in the face of non-stationary and transient background noise.

The most well-modeled and frequently observed gravitational-wave events to date are the mergers of binary black hole (BBH) systems[2,24,25] Their comparatively high number of confirmed detections has given us reasonable models of their population statistics, allowing for astrophysically meaningful measures of search sensitivity. BBH mergers also benefit from a highly localized-in-time signal-to-noise ratio (SNR) profile relative to binary neutron star (BNS) mergers, which are in the sensitive band of the detectors much longer. Studying the ability of neural-networks to detect BBH mergers, and in particular what real time use in the IGWN detectors looks like in this context, represents an important first step towards developing a more thorough understanding of how, and whether, these algorithms can be applied to more challenging signals such as BNSs, and what tools and infrastructure would be required to do so.

Here, we present `Aframe`, a flexible pipeline for detection of BBH mergers using deep learning. The implementation presented here uses a 1D convolutional neural-network. Convolutional neural-networks have previously been shown to have potential for gravitational wave detection[26], and we use this architecture, along with aggressive data augmentation techniques, to achieve a sensitivity competitive with matched filtering CBC search pipelines while requiring a significantly lower latency. More broadly, `Aframe` encompasses a suite of tools for quickly implementing, testing, and deploying new ideas at scale in order to more confidently realize the potential of machine learning in service to gravitational wave astronomy.

Our neural-network architecture modifies a standard ResNet54[27], which maps fixed length time-series of gravitational wave strain from two interferometers (here, the Hanford and Livingston LIGO interferometers) to a scalar detection statistic indicating whether a signal is present in the input. Critically, we replace 2D with 1D convolutions to accommodate time-series input. In addition, we replace standard Batch Normalization layers (BN)[28], with Group Normalization (GN) layers[29]. While BN layers fit parameters to statistics calculated along the batch dimension, GN layers are fit to statistics calculated from groups of channels. This choice was motivated by differences in the statistical properties of batches during training and inference. During training, there are significantly more signals in each batch than during inference, where most of the batch consists of noise. Thus, during training, BN layers will learn spurious statistical properties that are not present at inference time. GN layers mitigate this problem by learning statistical properties of individual channels. We found that using GN layers improves the agreement between validation and test time metrics, as well as overall testing performance. Good agreement between validation and test metrics is essential for ensuring the best neural-network is being selected for deployment. The neural-network is trained by minimizing a binary cross entropy loss function with an Adam[30] optimizer. We use a one cycle learning rate scheduler with cosine annealing[31].

Analyzing data with `Aframe` involves loading and preprocessing timeseries data, breaking it up into short time segments, then passing these segments through the neural-network. The throughput associated with each of these steps can vary drastically, as can the hardware and software necessary to accelerate them. In order to optimize the total throughput of this system, we adopt an inference-as-a-service (IaaS) computing model in which neural-network inference is handled by a dedicated service, to which client applications can send inference requests remotely. Each step in our pipeline is then implemented and scaled independently to most efficiently leverage a fixed pool of heterogeneous computing resources. This model has been shown to be effective in optimizing ML inference in GW astronomy[32], provided that "snapshotting"[33] is used to cache overlapping input data on the server side to avoid redundant data transfer. We adopt this paradigm using an off-the-shelf IaaS implementation, Triton Inference Server[34], and use the ML inference framework TensorRT to accelerate the neural-network inference step. The ability to scale and distribute a workload is an important part of any search pipeline, and the authors are aware of only one other ML-based CBC detection algorithm that has focused on scalability to arbitrary resources[35]. In the sections below, we compare both our sensitivity and our throughput to this work.

Inference is performed at a rate of 4 Hz (not to be confused with the neural-network throughput, see the discussion of computational requirements below). In other words, we pass windows of data to our neural-network for inference such that each window is shifted by 0.25 s. This inference sampling rate reduces the overall compute load without sacrificing search sensitivity (see Sec. 4 of Methods). These neural-network predictions are then integrated over time using a 1 s top hat filter (see Fig. 1). Because the neural-network is trained to encode time translation invariance (see Sec. 2 of Methods), we expect to see a consistently high neural-network responses when analyzing astrophysical signals. Thus, integration provides a mechanism to promote consistently high outputs while rejecting short transients that may correspond to non-astrophysical sources. Finally, the integrated time-series of neural-network predictions is clustered to avoid yielding multiple triggers for the same event. The maximum integrated value over an 8 s window is taken as the detection statistic corresponding to a candidate event.
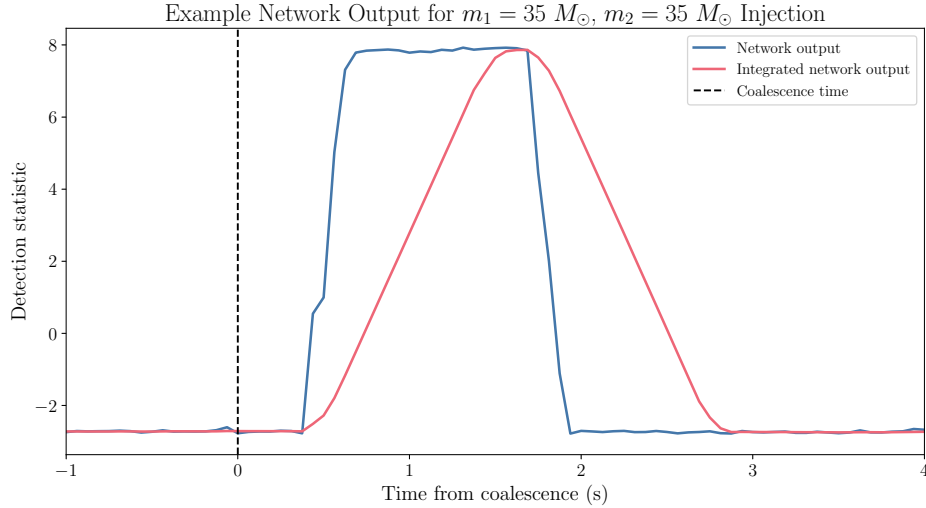


Figure 1: Example neural-network prediction and integrated neural-network prediction for a $m_1 = 35 M_\odot, m_2 = 35 M_\odot$ signal injection. The coalescence time is plotted as the vertical dashed black line. The brief gap between coalescence time and beginning of neural-network activation is due to the fact that we do not inject the coalescence time in the first or last 0.25 s of the window during training.

To demonstrate `Aframe`'s readiness for real-time deployment, we compare its sensitivity to search pipelines used in production by IGWN. For our pipeline, estimating sensitive volume requires analyzing simulated GW events "injected" into strain data, and analyzing background livetime produced by "timeslides." Performing timeslides is a standard way of empirically estimating

4

the background (i.e. the distribution of noise events) for a search pipeline which analyzes a network of detectors. In brief, the strain from one detector is shifted in time by an amount greater than the gravitational wave travel time between the detectors ($\sim 10\,\text{ms}$ for the two LIGO detectors). Therefore, any reported triggers could not have been caused by an astrophysical event. In this analysis, the Hanford strain data is held fixed and the Livingston data is shifted in 1 s increments until the required background livetime is accumulated. Then, a false alarm rate (FAR) can be assigned to injected events by dividing the number of background events with detection statistic greater than the event of interest, with livetime analyzed. All GW detections reported in the third Gravitational-Wave Transient Catalog (GWTC-3)[2] were excised from the background.

**Comparison with existing searches.** A useful metric to measure the sensitivity of search algorithms is the *sensitive volume*. Sensitive volume measures the volume over which some astrophysical population of sources distributed uniformly in co-moving volume is detectable at a given false alarm rate (FAR). Sensitive volume was used to measure the sensitivity of search pipelines in GWTC-3. This provides an astrophysically meaningful benchmark to compare the performance of `Aframe` to the performance of traditional searches. More details on the sensitive volume calculation can be found in Sec. 3 of Methods. Fig. 2 compares `Aframe`'s sensitive volume as a function of FAR with the sensitivity of the MBTA[37], PyCBC[38], GstLAL[39,40] and cWB[41] searches as reported in GWTC-3[2]. We note that the template banks used by MBTA, GstLAL, and PyCBC-Broad in the GWTC-3 analysis contain waveforms outside of the 5–100 $M_\odot$ range searched by `Aframe`. In principle, these searches could increase their sensitivities in the 5–100 $M_\odot$ range by removing these templates. This is evident when comparing the performance of PyCBC-BBH and PyCBC-Broad in Fig. 2. For the future, we encourage production level LVK CBC pipelines to publish BBH-specific sensitivities against which developing ML pipelines can benchmark.

In the 35-35 $M_\odot$ mass distribution, `Aframe` has a larger sensitive volume than the GWTC-3 configurations of all searches, and is comparable in the 35-20 $M_\odot$ mass bin, for the FARs considered in this analysis. As source masses decrease further, so does `Aframe`'s performance relative to existing pipelines. This is in part due to our neural-network architectures inability to model the lower frequency features of these low mass signals. While the architecture implements global pooling layers, the convolution layers use a kernel length of 3 samples. Improvements to neural-network architecture design, such as utilizing dilated convolutions that can better model these lower frequency

5

features will help to improve performance at these mass ranges.

Previous studies of ML-based gravitational wave detection algorithms tend not to use sensitive volume as a metric, preferring instead to use traditional ML metrics such as receiver operating characteristic (ROC) curves (an exception is[42], which uses a non-astrophysical prior and a Euclidean volume distribution). This makes direct comparison difficult, as these metrics depend on the parameter distributions of tested events. For the sake of completeness, in Fig. 3 we present our own ROC curve and find that, compared to previous works[35,43], we achieve nearly three orders of magnitude of improvement in true positive rate at a false positive rate of $\sim 10^{-6}$ for an SNR threshold of 6.23, where most astrophysical events are. However, we encourage future studies to use sensitive volume to astrophysically motivated distributions as the measure of performance.

**Detecting Astrophysical Candidates in GWTC-3.** The testing period we use contains 9 astrophysical candidate events reported as significant detections in GWTC-3. While we evaluated our algorithm's performance using "timeslides" of this data (see Sec. 4 of Methods), we also analyzed the unshifted (or "zero-lag") data to determine if our algorithm detects these known candidates. The results of this analysis are shown in Table 1. We detect all 9 candidates, with 8 of the 9 candidates detected at a false alarm rate of less than 1 per year, the minimum possible value for this analysis. For the final event, our reported false alarm rate, 14 per year, is of a similar magnitude to the false alarm rate reported by the GWTC-3 pipelines at 2.8 per year. Additionally, during this period, we do not report any non-catalog candidates with a false alarm rate less than 5 per month.

**Latency and Computational Requirements.** Training the neural-network with a single NVIDIA 16 GB Tesla V100 GPU takes approximately 43 hours, and once trained, the neural-network can continue to be used for months without retraining; see the discussion of algorithm longevity in Sec. 2 for details. For inference, we utilize a Triton inference server[34] that is hosted on a NVIDIA DGX server containing eight 16 GB Tesla V100 GPUs (See Sec. 4 for details on inference configuration). Altogether, analyzing the one year of background data and one year of injections used in this analysis to create Fig. 2 takes approximately 4 hours, corresponding to a throughput of about 500 seconds of data from a two detector network analyzed per second per GPU. This corresponds to an order of magnitude improvement in throughput compared with previous work by Huerta et al[35] and a factor of $\sim 2.5$ compared with Chatruvedi et al[43]. This improvement is
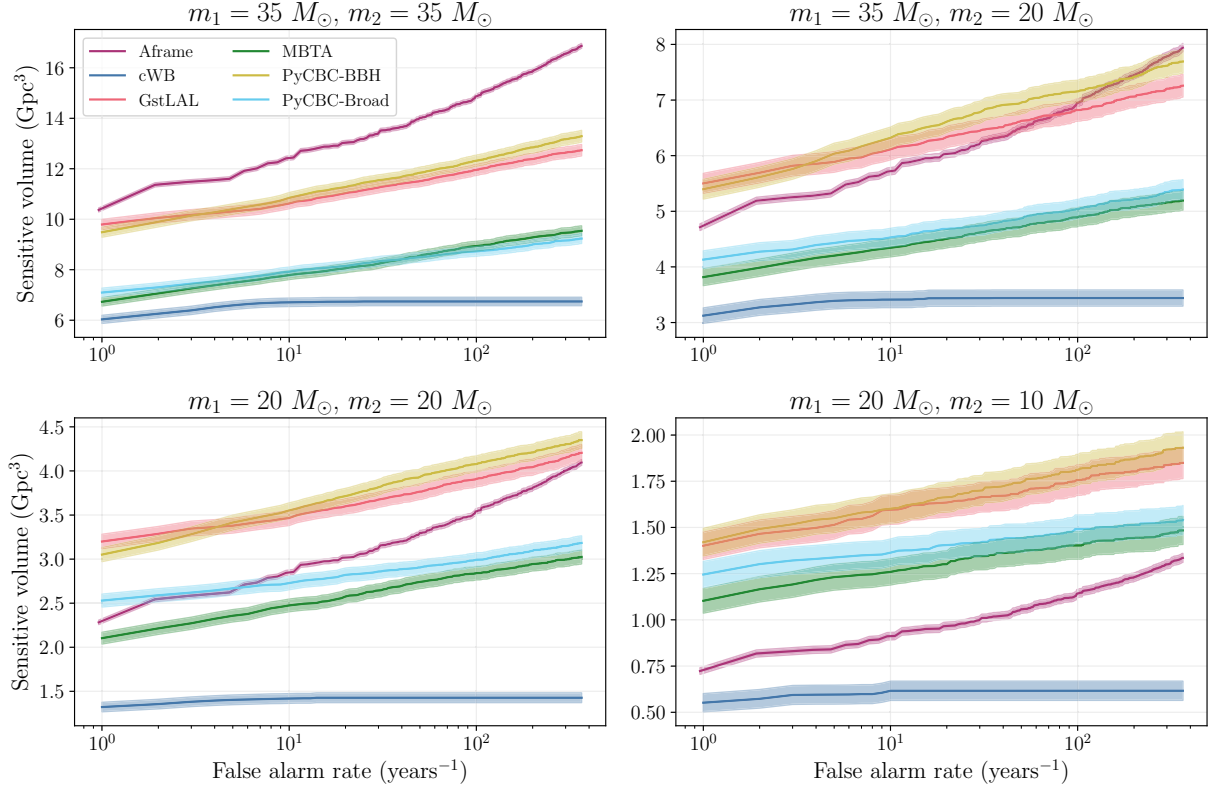
6

Figure 2: Sensitive volume vs FAR for four different mass distributions. Masses are specified in the source frame. Each mass is drawn from a log-normal distribution with a mean of the value given above each plot and a width of 0.1. `Aframe` demonstrates state-of-the-art sensitivity at higher masses, but loses performance relative to traditional search pipelines at lower masses. The sensitive volume of the other pipelines was calculated using data from a GWTC-3 data release[44].

due to the use of a more efficient neural-network architecture, as well as the IaaS model described above.

With trained neural-network weights in hand, the requirements for online deployment are much smaller. A single NVIDIA 24 GB A30 GPU is sufficient for real-time inference at an inference sampling rate of 2048 Hz, which provides sufficient resolution for coalescence time estimation. The total memory required to hold both the neural-network and data is 4.6 GB. The computational latency of the neural-network is less than 10 milliseconds. In practice, the latency of our algorithm is dominated by pre- and post-processing steps that bring the total latency to
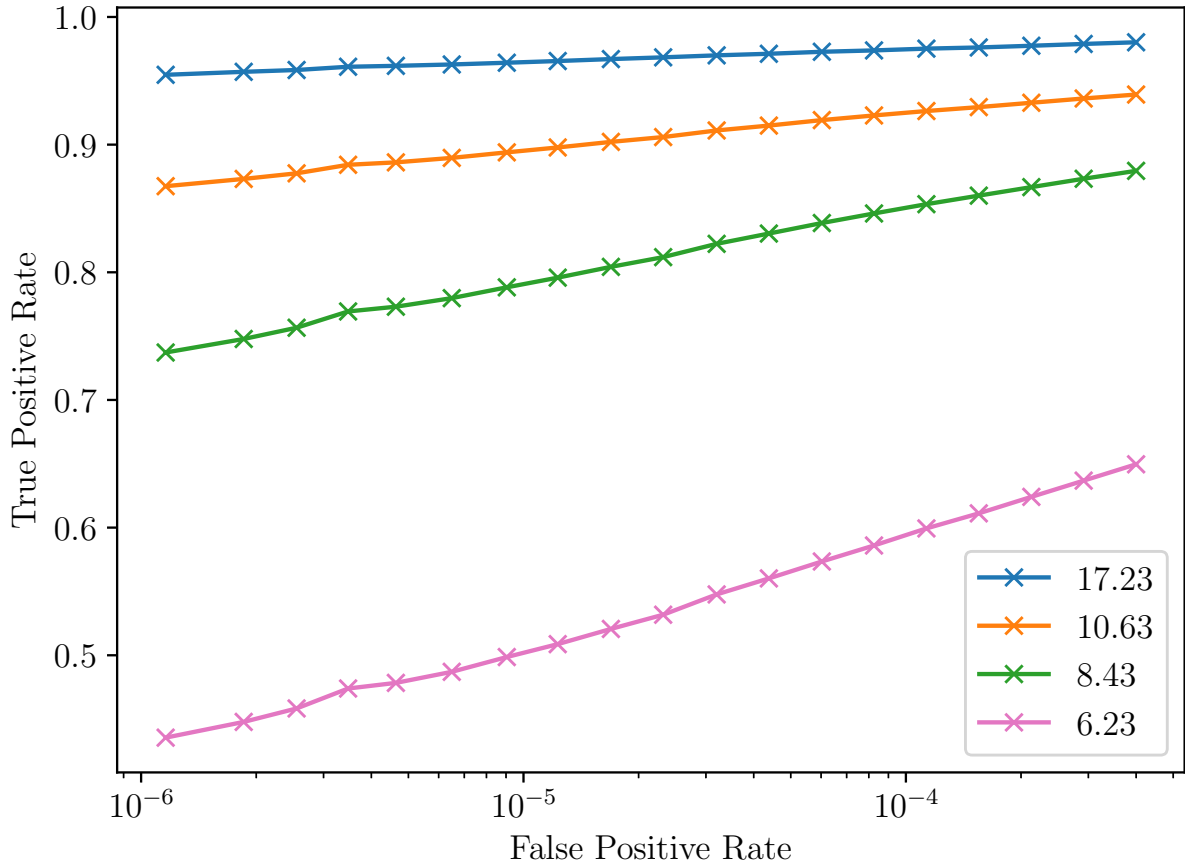
7

Figure 3: ROC curves for waveforms in different SNR bins in our testing dataset, described in Sec. 1. Each bin contains waveforms with SNRs at or above the given value.

approximately 3.1 s. For a detailed accounting of sources of latency within `Aframe`, see Sec. 4. In production, additional latency is incurred uploading events to the GRAvitational-wave Candidate Event DataBase (GraceDB) [a]. This latency is not included in this 3.1 s estimate. In addition, a recent study[46] used a real-time mock data challenge replay of O3 data to benchmark pipeline latencies, including GraceDB processing. Analyzing this data stream, we find a median (90%) event reporting latency of 3.9 s (4.3 s), in good agreement with our latency budget. Matched filtering pipelines report a median (90%) latency of 12.3 s (41.4 s).

**Discussion.** We have implemented a machine-learning based CBC search pipeline that is

---

[a] https://gracedb.ligo.org/

capable of low-latency use in a production setting. Through robust data augmentation techniques and extensive work in developing software infrastructure (Sec. 5), our algorithm achieves a sensitivity that is competitive with established search pipelines for higher mass BBHs. Work remains to improve the algorithm's performance on lower mass BBH systems. We leave these investigations to future work.

There are a number of extensions we plan to investigate in future work. Our algorithm is currently limited to the use of data from exactly two interferometers, and this limits our flexibility. In this work we trained our neural-network on data from the two LIGO interferometers, but we could benefit from the ability to include Virgo and KAGRA data. This could take the form of a four-detector model, or could be a suite of pairwise models that work in unison. Additionally, allowing for single-detector analysis would be beneficial for instances where only one detector is online. Further, low-latency alerts are less important for BBHs than binary neutron star (BNS) and neutron star-black hole (NSBH) mergers, where electromagnetic counterparts are more likely. The detection of these mergers with neural-networks is more challenging due to the greater length of time these signals spend in the sensitive band of the detector. Still, preliminary explorations indicate that our framework can adapt to address this problem.

| Event | $m_1(M_\odot)$ | $m_2(M_\odot)$ | Aframe | cWB | GstLAL | MBTA | PyCBC-BBH | PyCBC-Broad |
|---|---|---|---|---|---|---|---|---|
| GW190512_180714 | $23.2^{+5.6}_{-5.6}$ | $12.5^{+3.5}_{-2.6}$ | $< 0.97$ | 0.88 | $< 1.0 \times 10^{-5}$ | 0.038 | $< 1.1 \times 10^{-4}$ | $1.1 \times 10^{-4}$ |
| GW190513_205428 | $36.0^{+10.6}_{-9.7}$ | $18.3^{+7.4}_{-4.7}$ | $< 0.97$ | – | $1.3 \times 10^{-5}$ | 0.11 | 0.044 | 19 |
| GW190514_065416 | $40.9^{+17.3}_{-9.3}$ | $28.4^{+10.0}_{-10.1}$ | 14 | – | 450 | – | 2.8 | – |
| GW190517_055101 | $39.2^{+13.9}_{-9.2}$ | $24.0^{+7.4}_{-7.9}$ | $< 0.97$ | 0.0065 | 0.0045 | 0.11 | $3.5 \times 10^{-4}$ | 0.0095 |
| GW190519_153544 | $65.1^{+10.8}_{-11.0}$ | $40.8^{+11.5}_{-12.7}$ | $< 0.97$ | $3.1 \times 10^{-4}$ | $< 1.0 \times 10^{-5}$ | $7.0 \times 10^{-5}$ | $< 1.1 \times 10^{-4}$ | $< 1.0 \times 10^{-4}$ |
| GW190521 | $98.4^{+33.6}_{-21.7}$ | $57.2^{+27.1}_{-30.1}$ | $< 0.97$ | $2.0 \times 10^{-4}$ | 0.20 | 0.042 | 0.0013 | 0.44 |
| GW190521_074359 | $43.4^{+5.8}_{-5.5}$ | $33.4^{+5.2}_{-6.8}$ | $< 0.97$ | $1.0 \times 10^{-4}$ | $< 1.0 \times 10^{-5}$ | $1.0 \times 10^{-5}$ | $< 2.3 \times 10^{-5}$ | $< 1.8 \times 10^{-5}$ |
| GW190527_092055 | $35.6^{+18.7}_{-8.0}$ | $22.2^{+9.0}_{-8.7}$ | $< 0.97$ | – | 0.23 | – | 19 | – |
| GW190602_175927 | $71.8^{+18.1}_{-14.6}$ | $44.8^{+15.5}_{-19.6}$ | $< 0.97$ | 0.015 | $< 1.0 \times 10^{-5}$ | $3.0 \times 10^{-4}$ | 0.013 | 0.29 |

Table 1: Masses in units of $M_\odot$, and false alarm rates in units of inverse years from `Aframe`, cWB, GstLAL, MBTA, and PyCBC-BBH for the known events in our testing set. Masses come from Table VIII of GWTC-2.1[25], and FARs from Table XV of GWTC-3[2]. As our analysis examined only one year of background, our minimum FAR is one per year.

1. Abbott et al. Observation of gravitational waves from a binary black hole merger. *Phys. Rev. Lett.* **116**, 061102 (2016).

2. Abbott, R. *et al.* Gwtc-3: Compact binary coalescences observed by ligo and virgo during the second part of the third observing run. *Physical Review X* **13** (2023). URL `http://dx.doi.org/10.1103/PhysRevX.13.041039`.

3. The LIGO Scientific Collaboration *et al.* Advanced ligo. *Classical and Quantum Gravity* **32**, 074001 (2015). URL `https://dx.doi.org/10.1088/0264-9381/32/7/074001`.

4. Acernese, F. *et al.* Advanced virgo: a second-generation interferometric gravitational wave detector. *Classical and Quantum Gravity* **32**, 024001 (2014). URL `https://dx.doi.org/10.1088/0264-9381/32/2/024001`.

5. Akutsu, T. *et al.* Overview of KAGRA: Detector design and construction history. *Progress of Theoretical and Experimental Physics* **2021**, 05A101 (2020). URL `https://doi.org/10.1093/ptep/ptaa125`. `https://academic.oup.com/ptep/article-pdf/2021/5/05A101/37974994/ptaa125.pdf`.

6. Dwyer, S. *et al.* Gravitational wave detector with cosmological reach. *Phys. Rev. D* **91**, 082001 (2015). URL `https://link.aps.org/doi/10.1103/PhysRevD.91.082001`.

7. Punturo, M. *et al.* The einstein telescope: a third-generation gravitational wave observatory. *Classical and Quantum Gravity* **27**, 194002 (2010). URL `https://dx.doi.org/10.1088/0264-9381/27/19/194002`.

8. Edwards, T. e. a. LISA: Study of the Laser Interferometer Space Antenna: Final Technical Report (FTR), Dornier Satellitensysteme GmbH. Tech. Rep., DSS Report No. LI-RP-DS-009 (2000).

9. Tanvir, N. R. *et al.* A 'kilonova'associated with the short-duration -ray burst grb 130603b. *Nature* **500**, 547–549 (2013). URL `https://doi.org/10.1038/nature12505`.

10. Ascenzi, S. *et al.* A luminosity distribution for kilonovae based on short gamma-ray burst afterglows. *Monthly Notices of the Royal Astronomical Society* **486**, 672–690 (2019). URL `https://doi.org/10.1093/mnras/stz891`. `https://academic.oup.com/mnras/article-pdf/486/1/672/28320190/stz891.pdf`.

11. Jin, Z.-P. *et al.* A kilonova associated with grb 070809. *Nature Astronomy* **4**, 77–82 (2020). URL https://doi.org/10.1038/s41550-019-0892-y.

12. Rastinejad, J. C. *et al.* A kilonova following a long-duration gamma-ray burst at 350 mpc. *Nature* **612**, 223–227 (2022). URL https://doi.org/10.1038/s41586-022-05390-w.

13. Albert, A. *et al.* Search for neutrino counterparts to the gravitational wave sources from ligo/virgo o3 run with the antares detector. *Journal of Cosmology and Astroparticle Physics* **2023**, 004 (2023). URL https://dx.doi.org/10.1088/1475-7516/2023/04/004.

14. Abbasi, R. *et al.* Icecube search for neutrinos coincident with gravitational wave events from ligo/virgo run o3. *The Astrophysical Journal* **944**, 80 (2023). URL https://dx.doi.org/10.3847/1538-4357/aca5fc.

15. Harris, P. *et al.* Physics community needs, tools, and resources for machine learning (2022). 2203.16255.

16. Baltus, G. *et al.* Convolutional neural networks for the detection of the early inspiral of a gravitational-wave signal. *Phys. Rev. D* **103**, 102003 (2021). URL https://link.aps.org/doi/10.1103/PhysRevD.103.102003.

17. Verma, C., Reza, A., Gaur, G., Krishnaswamy, D. & Caudill, S. Can convolution neural networks be used for detection of gravitational waves from precessing black hole systems? *arXiv preprint arXiv:2206.12673* (2022).

18. Krastev, P. G. Real-time detection of gravitational waves from binary neutron stars using artificial neural networks. *Physics Letters B* **803**, 135330 (2020). URL https://www.sciencedirect.com/science/article/pii/S0370269320301349.

19. George, D. & Huerta, E. A. Deep neural networks to enable real-time multimessenger astrophysics. *Phys. Rev. D* **97**, 044039 (2018). URL https://link.aps.org/doi/10.1103/PhysRevD.97.044039.

20. Nousi, P. *et al.* Deep residual networks for gravitational wave detection. *Phys. Rev. D* **108**, 024022 (2023). URL https://link.aps.org/doi/10.1103/PhysRevD.108.024022.

21. Online pipelines - igwn — public alerts user guide 2023 (2023). URL `https://emfollow.docs.ligo.org/userguide/analysis/searches.html`.

22. Skliris, V., Norman, M. R. K. & Sutton, P. J. Real-time detection of unmodelled gravitational-wave transients using convolutional neural networks (2022). `2009.14611`.

23. Paleyes, A., Urma, R.-G. & Lawrence, N. D. Challenges in deploying machine learning: A survey of case studies. *ACM Comput. Surv.* **55** (2022). URL `https://doi.org/10.1145/3533378`.

24. Abbott, B. P. *et al.* Gwtc-1: A gravitational-wave transient catalog of compact binary mergers observed by ligo and virgo during the first and second observing runs. *Phys. Rev. X* **9**, 031040 (2019). URL `https://link.aps.org/doi/10.1103/PhysRevX.9.031040`.

25. The LIGO Scientific Collaboration and the Virgo Collaboration *et al.* Gwtc-2.1: Deep extended catalog of compact binary coalescences observed by ligo and virgo during the first half of the third observing run (2022). `2108.01045`.

26. Gebhard, T. D., Kilbertus, N., Harry, I. & Schölkopf, B. Convolutional neural networks: A magic bullet for gravitational-wave detection? *Phys. Rev. D* **100**, 063015 (2019). URL `https://link.aps.org/doi/10.1103/PhysRevD.100.063015`.

27. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778 (2016).

28. Ioffe, S. & Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift (2015). `1502.03167`.

29. Wu, Y. & He, K. Group normalization. In *Proceedings of the European Conference on Computer Vision (ECCV)* (2018).

30. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv e-prints* arXiv:1412.6980 (2014). `1412.6980`.

31. Smith, L. N. & Topin, N. Super-convergence: Very fast training of residual networks using large learning rates. *CoRR* **abs/1708.07120** (2017). URL `http://arxiv.org/abs/1708.07120`. `1708.07120`.

32. Gunny, A. *et al.* Hardware-accelerated inference for real-time gravitational-wave astronomy. *Nature Astronomy* **6**, 529–536 (2022). URL `https://doi.org/10.1038/s41550-022-01651-w`.

33. Gunny, A. *et al.* A software ecosystem for deploying deep learning in gravitational wave physics. In *Proceedings of the 12th Workshop on AI and Scientific Computing at Scale Using Flexible Computing Infrastructures*, FlexScience '22, 9–17 (Association for Computing Machinery, New York, NY, USA, 2022). URL `https://doi.org/10.1145/3526058.3535454`.

**Correspondence** Correspondence and requests for materials should be addressed to Ethan Marx (emarx@mit.edu) and Will Benoit (benoi090@umn.edu).

# Methods

## 1 Data

**Strain.** We train and validate our neural-network using open data from the Gravitational Wave Open Science Center (GWOSC)[47] between times 2019-04-29T13:29:25 and 2019-05-09T13:29:25, corresponding to a ten calendar day period at the beginning of the O3 observing run. The strain data is resampled to 2048 Hz for better computational efficiency. For each interferometer, we query the openly available science mode flag to remove segments with poor data quality. We then select segments for which the science mode flag is active for both the Hanford and Livingston LIGO interferometers. This amounts to approximately 4.7 days of coincident livetime. We reserve the segments that total a minimum of 15,000 seconds at the end of this period for validating the neural-network throughout the training process.

For evaluating the performance reported in Fig. 2, we select data satisfying the above criteria between times 2019-05-09T13:29:25 and 2019-06-08T13:29:25, corresponding to a 30 day period immediately after the training period. This amounts to approximately 18 days of coincident livetime. During evaluation, timeslides of this data are created such that the total desired background time is achieved. We emphasize that no data used for evaluating the performance of the neural-network was used during training or validation. In addition, we train the neural-network only with data from before the testing period. This mimics the data availability scenario for real-time application.

**Waveforms.** We use `bilby`[48] to simulate 100,000 eight second long BBH waveforms at 2048 Hz with the IMRPhenomPv2 approximant[49]. Out of these, 75,000 waveforms are used to train the neural-network, and the remaining 25,000 are reserved for validation. To simulate a waveform, a probability distribution is specified on each of the parameters that define a compact binary merger, and random samples are drawn from each. The distribution set used in this work is based on one used for GWTC-3[50] during O3 to assess the sensitivity of CBC search pipelines, and is described in Table 2. The sampled parameters are used to compute the time-domain strain for each polarization, $h_+$ and $h_\times$. The sampled component mass values are defined in the source frame, so conversion to detector frame quantities is performed before generation. The interferometer responses of the intrinsic polarizations are calculated during the training process to allow for real-time data augmentations, as described below in Sec. 2.

The same distributions are used to simulate signals for the testing dataset. Enough waveforms

15

are generated to fill the background timeslides with the waveform coalescence points spaced $24\,\mathrm{s}$ apart. As the signals are only $8\,\mathrm{s}$ long, they do not overlap. During the signal generation process, we perform rejection sampling and keep only signals that have an SNR greater than 4. This ensures that computation is not wasted on signals we do not expect to detect[44]. Rejection sampling reduces the uncertainty of a sensitive volume estimate for a fixed amount of analyzed injections (see Sec. 3). In total, we generate $\sim 45{,}000{,}000$ waveforms. Of these, $\sim 3\%$ percent are used for testing and $\sim 97\%$ are rejected.

## 2  Training

We apply several data augmentation techniques during the training process with the goal of providing robust, high entropy data that encodes physics-based knowledge for discriminating signals from noise. Below, we will describe how a training batch is composed, as well as the hyper-parameters that control the composition of the batches.

**Noise sampling.**  Sampled at $2048\,\mathrm{Hz}$, the entire training dataset is unable to fit onto a single $16\,\mathrm{GB}$ V100 GPU at once. Thus, efficient out-of-memory data-loading is required to fully utilize the extent of our strain dataset. To do this, we sample strain windows directly from disk during the training procedure. The length of each noise window sampled from disk is $10.5\,\mathrm{s}$. The first $8\,\mathrm{s}$ is used to estimate the power spectral density (PSD) used for whitening. We use Welch's method to estimate the PSD. The remaining $2.5\,\mathrm{s}$ of the window is whitened in the frequency-domain, and transformed back to time-domain. Due to whitening filter settle-in, $0.5\,\mathrm{s}$ of data is corrupted on both ends of the window and removed. Thus, only $1.5\,\mathrm{s}$ of data is actually analyzed by the neural-network. The PSD estimation, filter construction, and whitening are all done with `PyTorch`[51] modules to enable GPU-accelerated computation[20]. We use a training batch size of 384, which was chosen such that we fully utilize the GPU memory available. Our out-of-memory data-loading is sufficiently fast to support these batch sizes without bottle-necking the pre-processing or neural-network modules.

Noise instances are sampled independently in time from each interferometer. Thus, a noise instance from one interferometer can be paired with many different instances from the other interferometer. This combinatorially increases the amount of unique two-detector noise instances available for optimizing the network. Next, each noise instance has probability $p_{\mathrm{invert}}$ to be inverted $(h(t) \to -h(t))$ and, independently, probability $p_{\mathrm{reverse}}$ to be reversed $(h(t) \to h(-t))$[53]. Again, the

| Parameter | Description | Prior | Limits | Units |
|-----------|-------------|-------|--------|-------|
| $m_1$ | Mass of primary | $m_1^{-2.35}$ | $(5, 100)$ | $M_\odot$ |
| $m_2$ | Mass of secondary | $m_2$ | $(5, m_1)$ | $M_\odot$ |
| $z$ | Redshift | Comoving | $(0, 2)$ | - |
| $\psi$ | Polarization angle | Uniform | $(0, \pi)$ | rad. |
| $a_{1,2}$ | Dimensionless spin magnitude | Uniform | $(0, 0.998)$ | - |
| $\theta_{1,2}$ | Spin tilt | Sine | $(0, \pi)$ | rad. |
| $\phi_{12}$ | Relative spin azimuthal angle | Uniform | $(0, 2\pi)$ | rad. |
| $\phi_{JL}$ | Spin phase angle | Uniform | $(0, 2\pi)$ | rad. |
| $\phi$ | Orbital phase | Uniform | $(0, 2\pi)$ | rad. |
| RA | Right ascension | | $(0, 2\pi)$ | rad. |
| Dec | Declination | Cosine | $(-\pi/2, \pi/2)$ | rad. |
| $\theta_{JN}$ | Inclination angle | Sine | $(0, \pi)$ | rad. |

Table 2: Priors on parameters used to generate waveforms for both the training and testing sets. The prior is derived from that used in GWTC-3 to assess search pipelines. The component mass distributions are defined in the source frame. 'Comoving' refers to uniform in comoving volume.

| Parameter | Description | Prior | Limits | Best Value |
|-----------|-------------|-------|--------|------------|
| $lr_{\mathrm{max}}$ | Maximum learning rate | Log Uniform | $(10^{-4.5}, 10^{-2})$ | $5.8 \times 10^{-4}$ |
| $N_{\mathrm{ramp}}$ | Number of epochs over which learning rate increases | Uniform | $(2, 50)$ | 23 |
| $p_{\mathrm{signal}}$ | Probability of batch element containing a signal | Uniform | $(0.2, 0.6)$ | 0.277 |
| $p_{\mathrm{swap}}$ | Probability of swap augmentation | Uniform | $(0, 0.15)$ | 0.014 |
| $p_{\mathrm{mute}}$ | Probability of mute augmentation | Uniform | $(0, 0.3)$ | 0.055 |
| SNR steps | Number of batches over which SNR scheduler decays | Uniform | $(1, 2500)$ | 989 |

Table 3: Priors and descriptions of hyperparameters searched over. The best value corresponds to the neural-network from the hyperparameter search that produced the highest validation score across all epochs. A neural-network trained with these hyperparameters was used to evaluate results reported in Fig. 2. Details on hyperparameters can be found in Sec. 2

inversion and reversal augmentations increase the amount of unique noise instances in our training data. For transient noise, these augmentations increase the variety of morphologies provided during training, allowing for better generalization to unseen testing data. We fix $p_{\text{invert}}$ and $p_{\text{reverse}}$ to 0.5.

**Signal Injection.** Once a batch of noise instances is generated, simulated BBH signals are added into each 2.5 s unwhitened window with probability $p_{\text{signal}} = 0.277$ and labeled as signals; this signal probability is one of six hyperparameters that we search over (see Table 3 and the discussion of hyperparameters below). The procedure for injecting signals is as follows: first, intrinsic polarization time-series are randomly sampled from the training waveform bank. Next, random extrinsic parameters (right ascension, declination, polarization angle, and SNR) are sampled. The first three of these are sampled from the priors described in Table 2; We will discuss the method of SNR sampling in the following paragraph. Intrinsic polarization time-series are then projected onto the interferometers and re-scaled to the sampled SNR. Randomly sampling extrinsic parameters at training time allows each intrinsic time-series to be injected from a variety of sky localizations and distances throughout the training procedure. We found that standard CPU implementations of projecting intrinsic polarizations onto interferometers created bottlenecks that severely limited utilization of GPU resources. We eliminated this bottleneck by developing a `PyTorch`[51] implementation so that projection can be accelerated using GPUs by a factor of $\sim 200$. Finally, the interferometer responses are added into the noise instances. The coalescence time of the merger is randomly placed so that it falls at least 0.25 s from either edge of the 1.5 s whitened noise instance. We enforce this padding because we found that having the coalescence point too close to the left edge of the window makes it more difficult for the neural-network to learn, since much of the signal SNR would lie outside the window. The random placement of the coalescence time encodes time translational invariance so that the neural-network can identify signals with the coalescence time at different locations throughout the window.

**Curriculum Learning.** Curriculum learning is a technique for training machine learning models in which initially, easy to learn samples are provided as training data, and progressively harder samples are introduced over time. One way to apply this in the context of GW detection is to initially provide high SNR signals and gradually introduce lower SNR signals[20]. This allows the neural-network to quickly arrive at a minima of its parameter space before trying to optimize for the more realistic task. We begin with an SNR distribution that follows a power law, $p(\text{SNR}) \sim (\text{SNR})^{-3}$, with a minimum of $\text{SNR}_{\text{min}} = 12$ and a maximum of $\text{SNR}_{\text{max}} = 100$. The form of this distribution was chosen to roughly match the SNR distribution of of our astrophysically motivated

18

prior. Each time a new training batch is constructed, the minimum SNR bound of the distribution is decreased until we reach the ultimate lower bound of 4. This decrease happens uniformly over 989 batches, a value that was reached through a hyperparameter search.

**Glitch Mitigation.** Non-Gaussian noise transients, known as "glitches," can often mimic BBH signals and lead to high-significance false alarms. We implement two types of augmentations we call waveform *muting* and *swapping* to mitigate the impact of transient glitches. These augmentations respectively encode the concepts of coincidence and coherence that true astrophysical signals are expected to exhibit. The values of the parameters controlling these augmentations were determined by hyperparameter search; see below for more details.

*Muting*: For a fraction $p_{\mathrm{mute}} = 0.055$ of the training batch, we inject a BBH signal into only one of the interferometers and label these samples as noise. This teaches the neural-network that it is not enough for a BBH-like signal to be present in just one interferometer: coincidence between interferometers is a requirement for true astrophysical signals.

*Swapping*: For an independent fraction of the training batch, $p_{\mathrm{swap}} = 0.014$, we swap one of the interferometer responses with an interferometer response from different signal, and label these samples as noise. Thus, these windows will contain BBH waveforms with different intrinsic parameters in each interferometer. This motivates the neural-network to learn the concept of coherence: the time-frequency evolution of the signal must be identical in both interferometers.

**Algorithm Longevity.** Noise in gravitational wave interferometers is non-stationary. Therefore, the timescale over which a single trained neural-network will maintain its originally measured performance needs to be evaluated. Determining this timescale helps inform the cadence at which retraining is needed, if at all. To test the longevity of our algorithm, we construct several testing datasets at various intervals across O3. For each interval, we analyze the testing dataset with a neural-network trained using the first 10 days of O3 data. This is the same neural-network used to produce the results in Fig. 2. To separate the sensitivity of the neural-network from the sensitivity of the detectors, we do not measure sensitive volume, but instead look at the fraction of events with SNR $> 8$ that are detected at different FARs. This metric takes into account the variation in noise level across different time periods, though it does not account for all aspects of detector performance, such as the rate or morphology of glitches. At a FAR of 1 event per 2 months, a threshold comparable to the 1 event per 5 months used for releasing significant public alerts by

19

the IGWN[b], we see in Fig 4 that the fractional detection rate of the original neural-network does not decay with time. We note that the most significant background event across all weeks is found during week 2, corresponding to the sharp drop in detection fraction at a FAR of 1 per 2 months. Though there is some fluctuation from week to week, a single neural-network trained on a week's worth of data at the beginning of the observing run maintains sensitivity over the duration of the run.
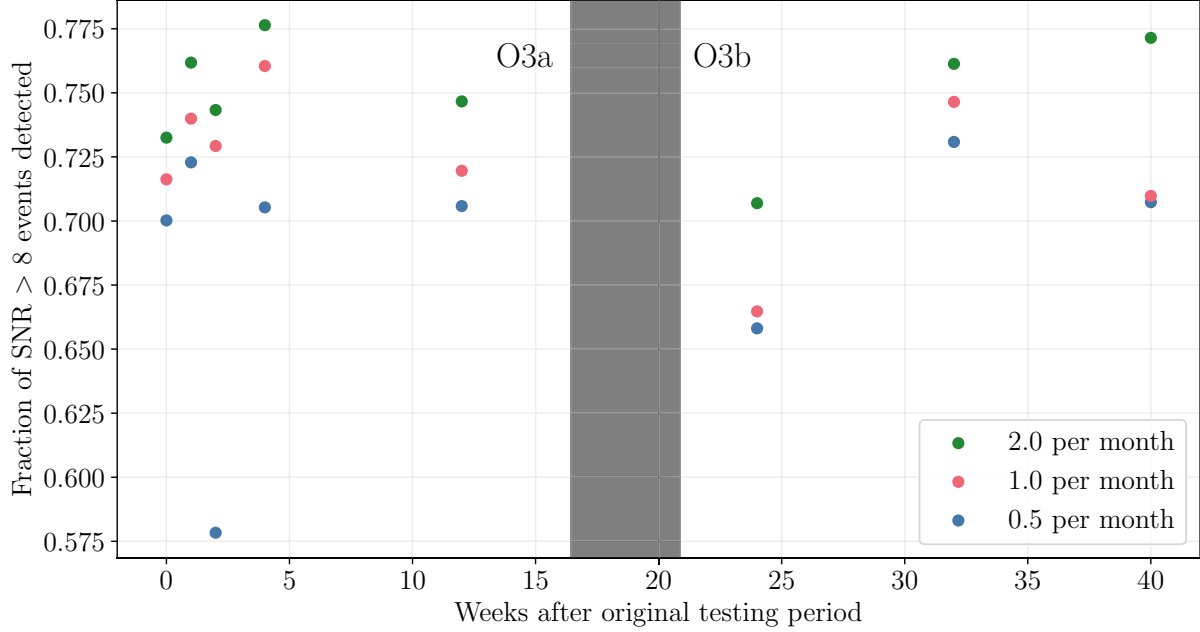


Figure 4: The fraction of SNR $> 8$ events detected at different false alarm rates during various weeks across a period of time during O3, beginning May 9th, 2019 and ending March 21st, 2020. Errors on detection fraction estimates are smaller than the plotted points.

**Validation.** We construct our validation procedure with the goal of establishing a strong correlation between validation and test metrics. This allows us to confidently pick the best performing neural-network during a hyperparameter search, as well as during individual training runs. To accomplish this, our validation procedure is designed to mimic the testing procedure as closely as possible. We reserve 15,000 seconds of strain data from immediately after the training period and 25,000 waveforms exclusively for neural-network validation during training. This data is not used at all for training the neural-network. This temporal choice of training and validation split mimics

---

[b]https://emfollow.docs.ligo.org/userguide/analysis/

the real-time production setting, where a deployed neural-network is only trained on past data.

To construct our validation set, we first create timeslides of the background data until at least 16 hours of livetime is accumulated. Similarly to training, this data is batched into 10.5 s windows, with the first 8 s used for whitening the final 2.5 s of each window. As with the training data, 0.5 s of data is cropped from each edge of the window after whitening. Next, we create a dataset of injections by adding waveforms from the validation waveform dataset into the background windows. We set a minimum detector-network SNR threshold of 4 for validation signals. Signals that are quieter are re-scaled to the SNR 4 threshold. The SNR is computed with respect to the PSD calculated from the first 8 s of the window. This rescaling procedure mimics the SNR-based rejection sampling performed for the testing dataset. We create 5 unique injection sets that have the coalescence point of each waveform at 0.25, 0.5, 0.75, 1.0, and 1.25 s within each whitened window. This ensures the validation metric covers a wider variety of scenarios.

The neural-network outputs a prediction for each window in the background and injection datasets. We use these predictions to calculate the area under the ROC curve (AUROC) up to a false positive rate (FPR) of $10^{-3}$, which is the final validation metric. We make this cut on the AUROC so that we are optimizing performance in the regime of low FARs. After the neural-network training has converged, the weights corresponding to the epoch with the highest validation score are used for testing.

**Hyperparameter Search.** The hyperparameters of our algorithm are optimized via a random search[54]. It is infeasible to search over all possible hyperparameters, so we selected those that we a-priori expect to have the greatest impact on the neural-network optimization process. These were the neural-network's maximum learning rate ($lr_{\mathrm{max}}$), the number of epochs over which the learning rate "ramps up" ($N_{\mathrm{ramp}}$) to $lr_{\mathrm{max}}$, $p_{\mathrm{signal}}$, $p_{\mathrm{mute}}$, $p_{\mathrm{swap}}$, and the number of steps over which SNR curriculum learning was performed. The priors on each of these parameters can be found in Table 3. 30 combinations of these parameters were randomly sampled and used to train a neural-network. Of these, the neural-network that reported the highest validation score was selected as the neural-network used for testing. The hyperparameters used to train this neural-network are reported in Table 3.

## 3 Sensitive Volume

A key metric in understanding a search algorithm's performance is the *sensitive volume*, which is a measure of the region of space in which a pipeline is expected to detect merging binaries. The sensitive volume as a function of the FAR is defined by

$$V(\mathcal{F}) = \int d\mathbf{x} \, d\theta \, \epsilon(\mathcal{F}; \mathbf{x}, \theta)\phi(\mathbf{x}, \theta) \tag{1}$$

where $\phi$ is the distribution of events over spatial coordinates $\mathbf{x}$ and binary system parameters $\theta$, and $\epsilon$ is the detection efficiency of the pipeline at a false alarm rate $\mathcal{F}$[55]. Generally, this quantity is estimated using Monte-Carlo integration by drawing waveforms from a population model, injecting them into a background, and counting how many produce triggers below a given false alarm rate threshold. If the samples are drawn from within the redshifted volume[56] $V_0$, with

$$V_0 = \int_{z_{\min}}^{z_{\max}} dz \frac{dV_c}{dz} \frac{1}{1+z} \tag{2}$$

where $dV_c/dz$ is the differential comoving volume, then the sensitive volume is approximately

$$V(\mathcal{F}) \approx V_0 \frac{N(\mathcal{F})}{N_{\mathrm{draw}}} \tag{3}$$

where $N(\mathcal{F})$ is the number of signals detected at a FAR less than $\mathcal{F}$ and $N_{\mathrm{draw}}$ is the number of injected events.

It is often desired to quantify the sensitivity of an algorithm to different populations. For example, an algorithm's sensitivity may vary with different source masses. Through the technique of importance sampling, it is possible to use one injection set from a broad population to calculate the sensitive volume for several populations. Each injection is weighted by the ratio of the probability of having been drawn from the injected distribution to that of the population distribution of interest[57]:

$$V_{\mathrm{pop}}(\mathcal{F}) \approx \frac{V_0}{N_{\mathrm{inj}}} \sum_{i=1}^{N(\mathcal{F})} \frac{p_{\mathrm{pop}}(\theta_i)}{p_{\mathrm{inj}}(\theta_i)} \tag{4}$$

The Monte-Carlo uncertainty on this estimation is[58]

$$(\delta V_{\mathrm{pop}})^2 = \frac{V_0^2}{N_{\mathrm{inj}}^2} \sum_{i=1}^{N(\mathcal{F})} \left( \frac{p_{\mathrm{pop}}(\theta_i)}{p_{\mathrm{inj}}(\theta_i)} \right)^2 - \frac{V_{\mathrm{pop}}^2}{N_{\mathrm{inj}}} \tag{5}$$

The SNR-based rejection performed during the generation of test set waveforms is done to improve this uncertainty. Waveforms that are sampled but have an SNR less than 4 are not injected; however, they still count towards $N_{\text{draw}}$. The cut is placed such that any waveforms below the SNR cutoff are not expected to be recovered at any reasonable FAR, and so would not contribute to the sensitive volume: whether injected or not, their weight would be zero. This procedure allows us to effectively draw many times more samples than are actually injected, greatly reducing the uncertainty on the sensitive volume. For this analysis, we re-weight to the same population distributions used in the sensitive volume analysis conducted in GWTC-3[2], log-normal distributions about central masses of interest with widths of 0.1. In addition, we enforce time difference of no more than 0.25 s between the recovered and injected coalescence times. This time difference corresponds to the resolution available at an inference sampling rate of 4 Hz. This time resolution can be reduced by increasing the inference sampling rate.

## 4 Inference

Our inference pipeline is an ensemble of three models: a snapshotter[33], a whitener, and the neural-network itself. Clients send streaming updates of strain data to a snapshotter. The snapshotter sends the latest state to the whitening module. Finally, batches of whitened data are constructed and analyzed by `Aframe`, producing predictions. The length of the state maintained by the snapshotter is determined by the length of the timeseries used to estimate the PSD, the batch size, and the inference sampling rate. For our analysis, the whitening module uses the first 64 seconds of the snapshotter state to estimate the PSD and build a whitening filter. The remaining data is whitened, and half a second is cropped from both edges to remove the effects of filter settle-in. The whitened data is then unfolded into a batch of overlapping windows. We use a batch size of 128 windows, and, as an inference sampling rate of 4 Hz was used, each 1.5 s window overlaps its neighbors by 1.25 s. This batch of windows is passed to the neural-network for prediction. Lastly, neural-network predictions are aggregated client-side and post-processed via the integration and clustering described above.

For an online analysis, the pre- and post-processing steps incur a total latency of approximately 3.1 s, see Table 4 for a summary. The most significant source of latency in the online analysis comes from waiting for data to exist such that we can crop the edges after resampling and whitening. All other computational steps (data reading/writing, data transfer to/from GPU, whitening, event identification, etc.) take less than 0.4 seconds combined, while the inference step itself takes less

23

than 10 milliseconds. An additional source of latency is the means by which live data is made available during an observation run. New data is written out in 1 s segments. Thus, depending on where the coalescence point of an event falls within one of these segments, it may be necessary to wait for a full additional second for a file to be written before event identification can occur. This factor is not included in Table 4, nor is the time it takes for data to become available, or the time taken to upload a candidate event to GraceDB, as none of these processes is within our control.

A critical parameter is the inference sampling rate. The inference sampling rate controls the stride between consecutive windows seen by the neural-network. Too small of an inference sampling rate, and astrophysical events may be skipped over. Too large, and computing resources are wasted on redundant inferences. We examined the impact of the inference sampling rate on our sensitivity by repeating trials of our inference procedure at inference sampling rates of 1, 2, 4, 8, 16 and 64 Hz. For this analysis, we accumulated two months worth of timeslide data for each trial. Fig. 5 shows a subset of the results of this analysis. Algorithms mostly perform within their statistical error. However, at low FARs the 1 Hz analysis has a small performance dip in the 35-35 mass bin. Because analyses performed at 4 Hz require 16 times fewer inference requests than 64 Hz without sacrificing performance, we use an inference sampling rate of 4 Hz for the analyses in this paper.

## 5   Data and Software Availability

All code used to produce results in this work is publicly available. The `Aframe` project repository can be found at `https://github.com/ML4GW/aframe`.

In addition, two open source libraries, `ml4gw`[c] and `hermes`[d] were developed to support this work. The `ml4gw` library contains PyTorch utilities for efficient on-GPU data-loading, whitening, PSD estimation and other data processing techniques common to GW analysis. The `hermes` library contains utilities for deploying models in the IaaS paradigm via Triton Inference Servers.

---

[c]`https://github.com/ML4GW/ml4gw`
[d]`https://github.com/ML4GW/hermes`

| Latency Source | Latency (s) |
|---|---|
| Coalescence point exiting training kernel padding | 0.25 |
| Cropping corruption from whitening filter | 0.50 |
| Cropping corruption from resampling to 2048 Hz | 1.0 |
| Integrating neural-network output | 1.0 |
| Reading data and transferring to GPU | $1.03^{+0.06}_{-0.05} \times 10^{-2}$ |
| Estimating PSD and whitening | $8.77^{+1.35}_{-0.31} \times 10^{-4}$ |
| Performing inference on whitened data | $9.63^{+0.38}_{-0.32} \times 10^{-3}$ |
| Integrating and aggregating neural-network output | $3.42^{+0.02}_{-0.01} \times 10^{-1}$ |
| Identifying candidate events in integrated output | $1.40^{+0.62}_{-0.43} \times 10^{-4}$ |
| Total | $3.114^{0.006}_{0.001}$ |

Table 4: Sources of latency for an `Aframe` online analysis. For the items listed in the upper section this table, the latency does not come from performing the computation, but rather from needing to wait for the data to exist before the action can occur. Items in the lower section are computational steps, and we report the median timing of 9191 trials. The upper and lower error bars represent the 95th and 5th percentile, respectively. All measurements were taken on a dedicated NVIDIA A30 GPU.
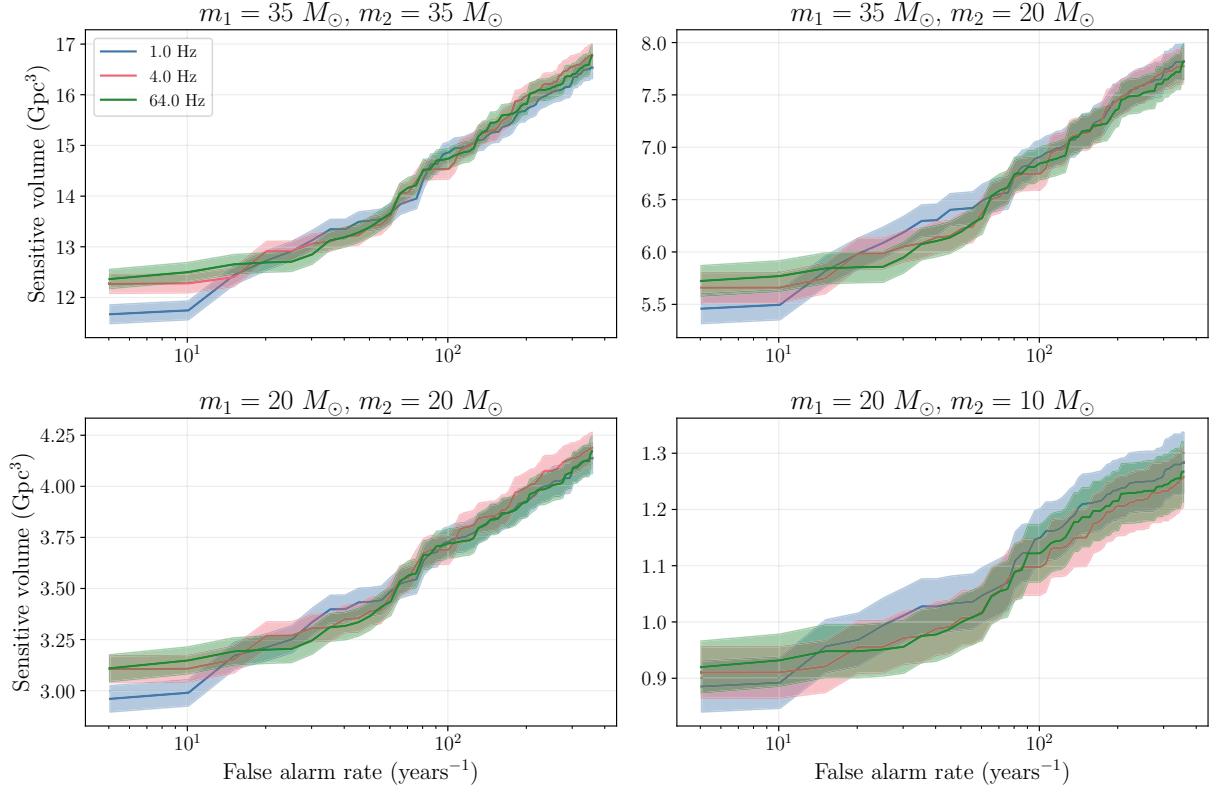
Figure 5: Sensitivity comparisons for the same neural-network run over the same data at different inference rates. For the purposes of clarity, only a subset of the tested rates are shown here. Except for the 1 Hz inference, all results are within error of each other for all mass combinations and FARs, including for rates not shown in this plot.

# References

34. NVIDIA Corporation. Triton Inference Server: An Optimized Cloud and Edge Inferencing Solution. (2023). URL `https://github.com/triton-inference-server/server`.

35. Huerta, E. A. *et al.* Accelerated, scalable and reproducible ai-driven gravitational wave detection. *Nature Astronomy* **5**, 1062–1068 (2021). URL `https://doi.org/10.1038/s41550-021-01405-0`.

36. Abbott, R. *et al.* Gwtc-3: Compact binary coalescences observed by ligo and virgo during the second part of the third observing run. *Physical Review X* **13** (2023). URL `http://dx.doi.org/10.1103/PhysRevX.13.041039`.

37. Aubin, F. *et al.* The MBTA pipeline for detecting compact binary coalescences in the third LIGO–virgo observing run. *Classical and Quantum Gravity* **38**, 095004 (2021). URL `https://doi.org/10.1088%2F1361-6382%2Fabe913`.

38. Canton, T. D. *et al.* Real-time search for compact binary mergers in advanced ligo and virgo's third observing run using pycbc live. *The Astrophysical Journal* **923**, 254 (2021). URL `https://dx.doi.org/10.3847/1538-4357/ac2f9a`.

39. Ewing, B. *et al.* Performance of the low-latency gstlal inspiral search towards ligo, virgo, and kagra's fourth observing run (2023). `2305.05625`.

40. Tsukada, L. *et al.* Improved ranking statistics of the GstLAL inspiral search for compact binary coalescences. *Physical Review D* **108** (2023). URL `https://doi.org/10.1103%2Fphysrevd.108.043004`.

41. Drago, M. *et al.* coherent waveburst, a pipeline for unmodeled gravitational-wave data analysis. *SoftwareX* **14**, 100678 (2021). URL `https://www.sciencedirect.com/science/article/pii/S2352711021000236`.

42. Schäfer, M. B. *et al.* First machine learning gravitational-wave search mock data challenge. *Phys. Rev. D* **107**, 023021 (2023). URL `https://link.aps.org/doi/10.1103/PhysRevD.107.023021`.

43. Chaturvedi, P., Khan, A., Tian, M., Huerta, E. A. & Zheng, H. Inference-optimized ai and high performance computing for gravitational wave detection at scale. *Frontiers in Artificial Intelligence* **5** (2022). URL http://dx.doi.org/10.3389/frai.2022.828672.

44. LIGO Scientific Collaboration, Virgo Collaboration & KAGRA Collaboration. GWTC-3: Compact Binary Coalescences Observed by LIGO and Virgo During the Second Part of the Third Observing Run — O3 search sensitivity estimates (2023). URL https://doi.org/10.5281/zenodo.7890437.

45. The LIGO Scientific Collaboration and the Virgo Collaboration *et al.* Gwtc-2.1: Deep extended catalog of compact binary coalescences observed by ligo and virgo during the first half of the third observing run (2022). 2108.01045.

46. Chaudhary, S. S. *et al.* Low-latency gravitational wave alert products and their performance in anticipation of the fourth ligo-virgo-kagra observing run (2023). 2308.04545.

47. The LIGO Scientific Collaboration and the Virgo Collaboration and the KAGRA Collaboration *et al.* Open data from the third observing run of ligo, virgo, kagra, and geo. *The Astrophysical Journal Supplement Series* **267**, 29 (2023). URL https://dx.doi.org/10.3847/1538-4365/acdc9f.

48. Ashton, G. *et al.* Bilby: A user-friendly bayesian inference library for gravitational-wave astronomy. *The Astrophysical Journal Supplement Series* **241**, 27 (2019). URL https://dx.doi.org/10.3847/1538-4365/ab06fc.

49. Hannam, M. *et al.* Simple model of complete precessing black-hole-binary gravitational waveforms. *Phys. Rev. Lett.* **113**, 151101 (2014). URL https://link.aps.org/doi/10.1103/PhysRevLett.113.151101.

50. Abbott, R. *et al.* Population of merging compact binaries inferred using gravitational waves through gwtc-3. *Phys. Rev. X* **13**, 011048 (2023). URL https://link.aps.org/doi/10.1103/PhysRevX.13.011048.

51. Paszke, A. *et al.* *PyTorch: An Imperative Style, High-Performance Deep Learning Library* (Curran Associates Inc., Red Hook, NY, USA, 2019).

52. Nousi, P. *et al.* Deep residual networks for gravitational wave detection. *Phys. Rev. D* **108**, 024022 (2023). URL https://link.aps.org/doi/10.1103/PhysRevD.108.024022.

53. Bini, S., Vedovato, G., Drago, M., Salemi, F. & Prodi, G. A. An autoencoder neural network integrated into gravitational-wave burst searches to improve the rejection of noise transients. *Classical and Quantum Gravity* **40**, 135008 (2023). URL https://doi.org/10.1088%2F1361-6382%2Facd981.

54. Bergstra, J. & Bengio, Y. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.* **13**, 281–305 (2012).

55. Schäfer, M. B., Ohme, F. & Nitz, A. H. Detection of gravitational-wave signals from binary neutron star mergers using machine learning. *Phys. Rev. D* **102**, 063015 (2020). URL https://link.aps.org/doi/10.1103/PhysRevD.102.063015.

56. Chen, H.-Y. *et al.* Distance measures in gravitational-wave astrophysics and cosmology. *Classical and Quantum Gravity* **38**, 055010 (2021). URL https://dx.doi.org/10.1088/1361-6382/abd594.

57. Tiwari, V. Estimation of the sensitive volume for gravitational-wave source populations using weighted monte carlo integration. *Classical and Quantum Gravity* **35**, 145009 (2018). URL https://dx.doi.org/10.1088/1361-6382/aac89d.

58. Farr, W. M. Accuracy requirements for empirically measured selection functions. *Research Notes of the AAS* **3**, 66 (2019). URL https://dx.doi.org/10.3847/2515-5172/ab1d5f.

59. Gunny, A. *et al.* A software ecosystem for deploying deep learning in gravitational wave physics. In *Proceedings of the 12th Workshop on AI and Scientific Computing at Scale Using Flexible Computing Infrastructures*, FlexScience '22, 9–17 (Association for Computing Machinery, New York, NY, USA, 2022). URL https://doi.org/10.1145/3526058.3535454.