# Rapid Likelihood Free Inference of Compact Binary Coalescences using Accelerated Hardware

**D. Chatterjee,**[1,2][†] **E. Marx,**[1,2] **W. Benoit,**[3] **R. Kumar,**[4]
**M. Desai,**[1,2] **E. Govorkova,**[1,2] **A. Gunny,**[1,2] **E. Moreno,**[1,2]
**R. Omer,**[3] **R. Raikman,**[1,2] **M. Saleem,**[3] **S. Aggarwal,**[3]
**M. W. Coughlin,**[3] **P. Harris,**[1] **E. 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, U. Minnesota, Minneapolis, MN 55455, USA
[4]Department of Aerospace Engineering, IIT Bombay, Powai, Mumbai, 400076, India

E-mail: [†]`deep1018@mit.edu`

**Abstract.** We report a gravitational-wave parameter estimation algorithm, `AMPLFI`, based on likelihood-free inference using normalizing flows. The focus of `AMPLFI` is to perform real-time parameter estimation for candidates detected by machine-learning based compact binary coalescence search, `Aframe`. We present details of our algorithm and optimizations done related to data-loading and pre-processing on accelerated hardware. We train our model using binary black-hole (BBH) simulations on real LIGO-Virgo detector noise. Our model has $\sim 6$ million trainable parameters with training times $\lesssim 24$ hours. Based on online deployment on a mock data stream of LIGO-Virgo data, `Aframe` + `AMPLFI` is able to pick up BBH candidates and infer parameters for real-time alerts from data acquisition with a net latency of $\sim 6$s.

## 1. Introduction

It has been almost a decade since the discovery of gravitational waves (GWs) from compact binary mergers [1], with the last few years seeing a steady increase in the number of discovered GW events. While the first observing run of the Laser Interferometer Gravitational-wave Observatory (LIGO) reported only three events [1, 2], ‡ the number count stood at 90 within a span of five years [3].

‡ GW151012, initially labeled as low-significance, was later confirmed as a third event in O1.

Furthermore, the current ongoing fourth observing run (O4) of ground-based observatories LIGO/Virgo/KAGRA (LVK) has already reported more than one hundred events discovered online. § The trend is expected to continue with the instrument getting closer to design sensitivity in fifth observing run.‖

In parallel, the scope of multi-messenger astronomy (MMA) with GWs has seen a steady increase in terms of effort and infrastructure being invested for the joint follow-up of GW signals with EM and other high-energy astrophysics counterparts. The online alert infrastructure of the LVK currently reports GW discoveries along with follow-up data products in $\sim$ 30s after merger time [4]. Early-warning searches [5] that can potentially pick up low-mass BNS systems up to $\sim$ 1 minute before merger have been deployed online [6]. The alert distribution mechanisms, like NASA GCN,¶ have seen upgrades [7], and new alert brokers like SCiMMA[+] have become available for the community to use. Publicly available services like TreasureMap [8] have seen a steady adoption from observatories to share observed and scheduled fields to orchestrate observations. Tools like SkyPortal [9] and TOM-Toolkit [10] have been developed to aid target-of-opportunity followup.

All this development comes at a time when the number of GW discoveries have significantly increased corresponding to the improvement in sensitivity of Advanced LIGO [11, 12], Advanced Virgo [13], and KAGRA [14] instruments, and the sensitivity of time-domain telescope facilities allow for unprecedented discovery rates. However, identifying gravitational-wave counterparts jointly have been extremely challenging. The discovery of GWs and multi-wavelength EM emission from the merger of the binary neutron star (BNS), GW170817, [15, 16] remains the first and only success story, albeit a rarity, with most subsequent candidates likely to be at much farther distances [17].

One primary step toward improving follow-up campaigns is the availability of fast, real-time Bayesian parameter estimation (PE) of compact binary coalescence (CBCs) to provide accurate data products for GW follow-up. The computationally expensive part of stochastic sampling techniques, like nested sampling currently in use, involve the repeated computation of the likelihood. Techniques like reduced-order-quadrature (ROQ) [18, 19] and focused-ROQ [20] have been developed in view of making real-time PE as fast as possible. This is currently used in LVK to deliver update alerts from Bayesian parameter estimation on the timescale of several minutes to hours. Other techniques like the use of accelerated hardware for stochastic sampling has been reported in [21, 22], and mesh-free approximation for sky-localization, reported in [23, 24].

More recently, likelihood-free inference (LFI), using variational methods, have emerged as a different paradigm with flexible neural network approximators being used to learn the posterior or the likelihood. Their use has been demonstrated on GW data [25, 26], in particular with posterior estimation using normalizing flows [27], such

§ See https://gracedb.ligo.org/superevents/public/O4/ for the most updated list.
‖ See https://observing.docs.ligo.org/plan/ for observing plans.
¶ https://gcn.nasa.gov/
[+] https://scimma.org/

as in the DINGO algorithm. However, in order to relay discovery alerts for prompt followup, the combination of *search and inference* needs to be considered together.* Also considering a live, real-time system design, several overhead costs like data transfer, file input/output operations, communicating data to a remote server, and so on are often overlooked in isolated analyses, but show up in the overall time-to-alert. It is also worth highlighting that traditionally in GW data analysis, the search and PE components have been treated separately – match-filtering searches pick up the candidates from the data stream using suitable detection statistic, but also provide important context like the best matching template and the signal-to-noise ratio time series, which is then used to compute sky-localization maps [29] and EM-bright source properties [30] sent out in the sub-minute alerts by the LVK [4]. The results are then updated based on Bayesian PE results, in few hours timescale.

Although machine-learning techniques like LFI bring promise, large model size and/or long training times can be a barrier for operations. Also, given the slowly changing background over the course of days to weeks, the algorithm should be re-trainable from a previous model state, preferably without investing on expensive and dedicated online hardware for this purpose. This is currently lacking for online models like DINGO, which report 10-day training time on a NVIDIA A100 GPU [27].

In this work, we try to address the points highlighted above in the context of fast online search and parameter estimation for MMA with GW. We report `AMPLFI`,♯ a PE algorithm based on LFI using normalizing flows. The primary focus of `AMPLFI` is to run alongside neural-network based CBC search `Aframe` [31], and compute GW alert data products like skymaps and other use source-properties to be sent out with LVK discovery alerts. Though the core principles of LFI and its application are similar to efforts mentioned above, the technical implementation is independent and focused toward online inference. In particular, there are several elements of GW data analysis that are re-implemented as a part of `ml4gw` (codebase: `https://github.com/ML4GW/ml4gw`), designed for running on accelerated hardware like GPUs for fast and efficient training and inference. Some common set of tools from `ml4gw` are used by both `Aframe` (codebase: `https://github.com/ML4GW/aframev2/`) and `AMPLFI` (codebase: `https://github.com/ML4GW/amplfi`), the latter being the focus of this work.

We outline the rest of the paper as follows. In section 2, we motivate our design principles toward running search and PE together. In section 3, we mention optimizations related to data pre-processing and implementing simulations on accelerated hardware which ensure that most of the computation is occurring on the GPU. In section 4, we present the details of a data embedding network which is used to summarize the data. This embedding is pre-trained using a self-supervised method to create data summary marginalizing parameters that are not of interest (at least from GW alerts standpoint). In section 5, we give the details of our normalizing flow implementation. We present results and benchmarks in section 6, before concluding in

---

* This is done in case of stochastic signals offline, for example, see [28].
♯ **A**ccelerated **M**ultimessenger **P**arameter estimation using **LFI**; pronounced "amp-li-fy"

section 7.

## 2. `Aframe` + `AMPLFI`

In order to build as fast of a system as we can, we have made a number of design choices when building the `Aframe` + `AMPLFI` framework that we highlight in the following:

- A modular design to perform search and PE. The search for GW signals in this case is done by `Aframe`. Candidates from `Aframe` provide an estimate of the time of arrival and the significance via a false-alarm-rate (FAR). This is unlike traditional match-filtering searches that provide, in addition, the best matching template, and the corresponding signal-to-noise time series that is used by other annotation algorithms to provide skymaps [29] and source properties [30] of binary systems. In the proposed framework, once a segment of data is found to be of high-significance i.e., containing a GW signal, the parameters are inferred using `AMPLFI`. Hence, Bayesian parameter estimation results are available along with the discovery of the candidate.

- Data is held in GPU memory to minimize overheads in communication between different components in the low-latency alert infrastructure [4]. For example, `Aframe` runs as a service, maintaining a buffer †† of the data in GPU memory. Once a trigger occurs, the relevant segment is passed to `AMPLFI` for inference. Based on model size of `Aframe` and `AMPLFI`, both are able to be served on a single GPU like NVIDIA A30. This reduces any inference overheads as the data is kept on the same device. However, the models may be served as separate micro-services in case the model size or running on a single GPU turns out to be a barrier.

- The accuracy of the results are suited for "online" purposes i.e. the aim is to provide data products like skymaps, and source properties for online LVK alerts. Therefore, we restrict ourselves to GW waveform models that capture the inspiral-merger-ringdown phases, but do not focus on physics of higher-modes, spin precession etc., and prioritize fast inference for data products required for follow-up.

For `AMPLFI` we use a normalizing flow to learn the posterior distribution directly using simulations of binary black hole signals (BBHs). We also make some optimizations compared to previous efforts in light of an online inference algorithm:

- We use real detector data from the LIGO GW instruments during training. Most previous efforts use simulated, colored Gaussian noise.

- We use efficient data loading and whitening tools to minimize the data transfers back and forth between CPU and GPUs (or other accelerators). We elaborate this below in section 3.

- We re-implement CBC waveform generation on the GPU memory to directly generate waveforms on-the-fly during training.

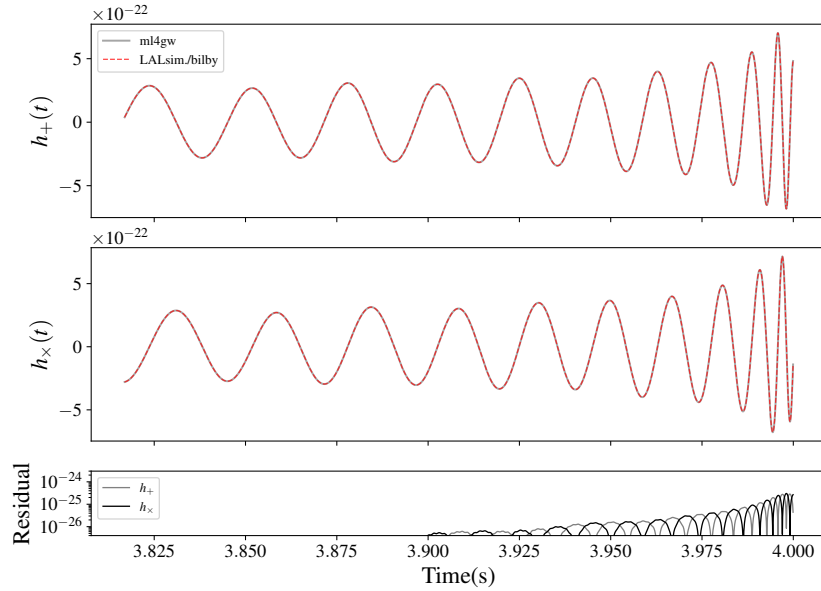†† A snapshotter that only sends new data segments into GPU memory

Figure 1: A comparison of time-domain `IMRPhenomD` between that implemented in this study, as a part of the `ml4gw` library. The parameters of the waveform is $\mathcal{M} = 26\ M_\odot, q = 1.0, D_L = 1000$ Mpc, $\chi_{1,2} = 0.0$. We find that while there are differences in the waveform strain, the residuals are below three orders of magnitude compared to the signal for most of the evolution, except the final few cycles where it is two orders of magnitude lower.

## 3. Simulations on Accelerated Hardware

### 3.1. Waveform model

We use the `IMRPhenomD` phenomenological waveform model for our simulations [32]. This waveform model contains the full inspiral-merger-ringdown physics, starting with the inspiral phase up to 3.5 post-newtonian order in GW phase (known as TaylorF2; see [33] for a review), and using an ansatz for the merger and ringdown, fitting them to numerical relativity results. One limitation of this waveform model is the restriction to aligned spins i.e. BH spin components perpendicular to the orbital plane and therefore no precession. Current online PE using stochastic sampling techniques use the `IMRPhenomPv2` waveform approximant, which contains precessing effects. Furthermore, high mass BBH systems use the `IMRPhenomXPHM` approximant, which also includes higher modes of radiation. However, we note that inference like sky-localization is insensitive to such effects. Also EM-brightness of a binary depends primarily on the aligned spin components aside from the mass ratio. Hence, the use of aligned-spin is justified for online purposes. In the future, however, we plan to implement and integrate the `IMRPhenomPv2` waveform model with our workflow.

In Figure 1, we show the time-domain strain of a representative BBH system based on our `IMRPhenomD` implementation and compare it with that implemented in
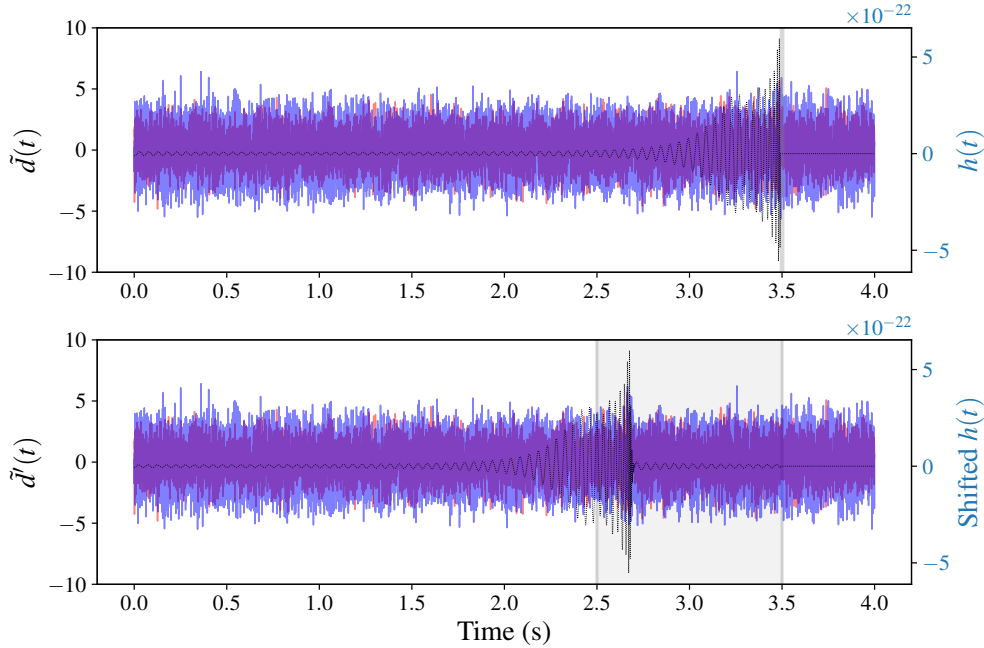
Figure 2: The figure shows the whitened time-domain background strain from Hanford and Livingston (HL) in two different colors. This is a stretch of data from May 2019 (early O3). A simulated BBH signal is injected; the waveform is overlayed. In the bottom panel, we see the same background strain with a time-shifted signal injected, the parameters of which are otherwise the same as the top panel. The time shift is random up to 1s compared to the top panel, indicated by the shading. We summarize our data views like $d$ and $d'$, and embedding them jointly. Subsequently, for LFI, we only use $d'$-like views.

`lalsimulation`. The latter is a part of the LIGO Algorithm Library [34], and provides the core components of GW data analysis with LVK data. We find consistency between our implementation and that in `lalsimulation`, with the residual errors below the signal by a couple of orders of magnitude throughout the evolution. Such residuals are unlikely to impact the results since the statistical errors of posterior is greater than systematic error due to such differences. We present some comparison results in section 6. Though we show comparison with a single representative system in Figure 1, several other combination of parameters are tested for consistency with `lalsimulation` as a part of unit-tests of the `ml4gw` codebase (`https://github.com/ML4GW/ml4gw`).

### 3.2. Data generation on the GPU

Generally, neural-network models are trained using batches (also called mini-batches) of data that is pre-processed on the CPU and then transferred to the GPU (or other co-processor) to carry out the forward/backward passes, and updating the model weights. However, this may leave the GPU under utilized if the pre-processing and data transfer between the CPU/GPU takes greater time compared to the operations to train the model. This is especially important in the context of LFI since efficient training relies on

providing unique combinations of parameters and data to approximate the distribution. We therefore take a different approach by performing the data generation on the GPU which allows generation of batches of data on the GPU, which is faster and at the same time avoids the data transfer overheads. Additionally, all data pre-processing, like fourier transforms, power spectra estimation, data whitening, are carried out on the GPU. This ensures consistent GPU utilization. Also, we can take advantage of the fact that GPU architecture today provide large memory. Our workflow involves:

- Transferring a chunk of two detector (Hanford and Livingston, subsequently HL) time-domain strain data, typically quarter of a day, sampled at 2048 Hz to the GPU before commencing training. The power spectra is fit to this background chunk. During training, $N$ small background chunks, each of 4s duration, are lazily loaded from the total training chunk, where $N$ is the training batch size.

- We generate $N$ points from our parameter prior and generate the `IMRPhenomD` waveforms directly on the GPU, as mentioned in section 3.1. This step is fast, for example, generating $N = 1000$ waveforms $\sim 0.15$s on a NVIDIA A40 GPU.

- We then inject the signals into the background chunks, obtaining the data batch. We whiten the batch using the estimated power spectra and pass it along with the parameters for training/validation/testing. Examples of the whitened data with the injected waveform overlayed is shown in the panels of Figure 2.

We call this implementation `InMemoryDataset` in `ml4gw`. We also note that though we have used GPUs as the co-processors in this work, our software framework can also be ported to other accelerators, like TPUs or HPUs, supported by the `pytorch-lightning` [35] framework that we use.

## 4. Embedding Network

Our input to the neural network model is a 2-channel Hanford-Livingston (HL) 4s whitened time-domain strain. This is projected into a lower-dimensional representation using a embedding network before performing LFI. The coherent analysis of both channels of data is important for some aspects of GW parameter estimation, like sky-localization since it depends on the difference in time of arrival in the different instruments. Our embedding network resembles the ResNet architecture used in `Aframe`. The implementation closely resembles that of the `torchvision` library, with some differences. Firstly, 1-D convolutions are used for time-series, instead of 2-D variants used for images. We use group normalization [36] instead of batch-normalization. Also, in `Aframe`, the architecture closely resembles a 34-layer residual network [37]. We, however, avoid the final 512-channel stack of convolution blocks (see Figure. 3 in [37]) since we do not find performance improvement after including the same. Thus our layer stacks contain blocks of 64, 128, and 256-channel residual convolution blocks. The number of convolution layers in each block is determined by hyper-parameter optimization (HPO) using Variance-Invariance-Covariance Regularization (VICReg) [38]

loss, detailed below in section 4.1. Details about the HPO are presented in Appendix A. The best configuration resembles an analogous 24-layer ResNet. A final fully-connected layer projects to an representation dimension $D_\gamma = 8$, which is also determined as a part of the HPO.

*4.1. Self-supervised learning of nuisance parameters*

We pre-train the embedding network to marginalize over uncertainties in arrival time up to 1 second. This is done since the *peak* of the detection statistic reported by `Aframe` may differ from the true arrival time up to tens of milliseconds. We choose 1-second as a conservative upper bound for the same. The pre-training is done via self-supervised learning (SSL) by identifying two "views" of the data as being the similar, and training the embedding network to minimize VICReg. Examples of two different views, $d$ and $d'$, as shown in Figure 2, where the upper panel shows a signal that is injected at a fixed reference time, while the lower panel shows a time-shifted signal i.e. all signal parameters being the same except the time of arrival, which is chosen randomly up to 1s in this case. Two batches of views are then forward-modeled through the embedding network and projected down to the resulting space. This projection, $\mathbf{\Gamma}$, is performed in two different steps via the ResNet, $f$ mentioned above, which projects the inputs in to an 8 dimensional space $\gamma \in \mathbb{R}^8$, then expanding this projected space using another fully-connected network, $h$, which takes $\gamma$ to a 24-dimensional space $x \in \mathbb{R}^{24}$. The resulting composition is given by $\mathbf{\Gamma} \equiv h \circ f$,

$$\gamma = f(\mathbf{d}); \ \gamma' = f(\mathbf{d}'); \ x = h(\gamma); \ x' = h(\gamma'). \tag{1}$$

We follow the prescription mentioned in [38] and compute the $\mathcal{L}_{\text{VICReg}}$ loss in the expanded dimension as,

$$\mathcal{L}_{\text{VICReg}}(x, x') = \lambda_1 \, \text{MSE}(x, x') + \lambda_2 \left[ \sqrt{\text{Var}(x) + \epsilon} + \sqrt{\text{Var}(x') + \epsilon} \right] + \\ \lambda_3 \left[ C(x) + C(x') \right]. \tag{2}$$

Here, MSE is the mean-squared error between the two projected views. The second term involves the variances of the individual batches, regularized by a tolerance to prevent collapsing to zero. Finally, the third term is the quadrature sum of the off-diagonal entries in the individual covariance matrices of the views. The $\lambda_{1,2,3}$ are relative weights of each term, which we also select after hyper-parameter tuning.

Previous work using LFI reported other techniques in particular, group equivariance, to tackle this [39]. We, however, take a different approach since parameters like time of coalescence are not as important for follow-up as masses and sky location of the signal, and marginalize them in our data summary. For more details on this technique, the reader is referred to [40]. This technique can be extended to other nuisance parameters in case of GWs, for example the coalescence phase. This is left to future work.

Table 1: Prior distributions of parameters. Note that the distance prior is a power-law with index 2 i.e. uniform in volume; cosmological effects are not included.

| Parameter | Prior |
|---|---|
| $\mathcal{M}$ | Uniform(10, 100)$M_\odot$ |
| $q$ | Uniform(0.125, 1) |
| $D_L$ | Uniform in Vol.(100, 3000) Mpc ($\sim D_L^2$) |
| $\theta_{\mathrm{JN}}$ | Sine(0, $\pi$) |
| $\alpha$ (RA) | Uniform(0, $2\pi$) |
| $\delta$ (Dec.) | Cosine($-\pi/2$, $\pi/2$) |
| $\phi_c$ (Coal. phase) | Uniform(0, $2\pi$) |
| $\psi$ (Pol. angle) | Uniform(0, $\pi$) |

## 5. Posterior Estimation

Posterior estimation in LFI involves learning the posterior, $p(\mathbf{\Theta}|\mathbf{d})$, using an approximator, $q_\varphi(\mathbf{\Theta}|\mathbf{d})$, using simulations $\{\mathbf{\Theta}_i, \mathbf{d}_i\}$. The parameters $\varphi$ are adjusted to maximize the likelihood of the simulations, which is mathematically equivalent to minimizing the Kullback-Leibler (KL) divergence between the true posterior and the approximator. The loss function used is,

$$-\ln \mathcal{L}(\varphi) = -\frac{1}{N_{\mathrm{sims.}}} \sum_{i \in \mathrm{sims.}} \ln q_\varphi(\mathbf{\Theta}_i|\mathbf{d}_i), \tag{3}$$

where the simulations are forward modeled to calculate their likelihood, which is then maximized during training. The density evaluations are done by learning a set of variable transforms that take the original variables $\mathbf{\Theta}$ to variables of a simpler base distribution, like a standard normal, which we use here. Several techniques are used to build flexible transforms and preserve the probability density at each stage. We refer the reader to a review article on normalizing flows and the different implementations [41]. In our case the parameter space is 8-dimensional,

$$\mathbf{\Theta} = \{\mathcal{M}, q, D_L, \theta_{\mathrm{JN}}, \alpha, \delta, \psi, \phi_c\}. \tag{4}$$

The prior distribution used for generating the simulations is mentioned in Table 1. Note that the time of coalescence is not a part of parameter set since it is marginalized over. Although the `IMRPhenomD` supports BH spins, we have ignored it for the current work. This will be relaxed in subsequent versions of the code and reported in a future work. The data, $\mathbf{d}$, consists of 4s sampled at 2048 Hz of whitened time-domain strain containing a BBH signal, as illustrated in Figure 2. When training the normalizing flow, we condition the parameters on the data representation, $\gamma$, as shown Eq. (1). This uses only the ResNet, $f$. The expander, $h$, is not required subsequently. Also, we leave the weights of $f$ to change further as a part of training the normalizing flow. Hence, our

normalizing flow maximizes,

$$-\ln \mathcal{L}(\varphi) = -\frac{1}{N_{\text{sims.}}} \sum_{i \in \text{sims.}} \ln q_\varphi(\boldsymbol{\Theta}_i | f(\mathbf{d}_i)), \qquad (5)$$

where the difference between Eq. (3) vs. Eq. (5) is the conditioning on the data summary, pre-trained to marginalize time of arrival. It should be mentioned that the pre-training step with VICReg loss is significantly cheaper compared to training the normalizing flow. In our experiments, we found the pre-training requiring few-tens of epochs of training, which took less than an hour to reach early-stopping condition on a NVIDIA A40 GPU.

## 5.1. Autoregressive Flows for LFI

Our normalizing flow implementation uses inverse auto-regressive transforms [42]. This kind of auto-regressive transforms can be sampled in one forward pass. However, evaluating the density requires $D$-forward passes, where $D$ is the dimensionality of the parameter space i.e. $D = 8$ in our case. Masked auto-regressive transforms [43] on the other hand use similar masked linear layers [44], but on the contrary the density evaluation takes a single forward pass and sampling is $D$-times as expensive. Although auto-regressive flows are *universal approximators* [41], our choice of using inverse auto-regressive flow (IAF) is because our inference requires fast sampling, which is achieved in a single pass with the IAF. In addition to IAF, coupling transforms were also considered. However, in our experiments, we found such transforms to perform less optimally when constrained to the same number of trainable parameters.

Our transforms are implemented using the open source library `pyro` [45]. The complete transform is composed of 60 individual affine-autoregressive transforms. Each transform has 6 masked linear layers; each layer having 100 hidden units. More complex transform functions like monotonic splines and neural-network with positive weight exist in the literature. However, we use the affine transforms for simplicity and lower number of trainable parameters. We train the network with a batch size of 800, with 200 batches per epoch using the AdamW optimizer [46] with initial learning rate of 1e−3 and weight decay of 2e−3. The learning rate is scheduled to reduce by a factor of 10 upon plateauing of the validation loss with a patience of 10 epochs. These configurations are decided after hyper-parameter tuning over a combination of several hundred parameter combinations detailed in Appendix B. Our trainer is scheduled to terminate training once the validation loss saturates with a patience of 50 epochs. In terms of training time, we find training $\gtrsim 200$ epochs with the above configuration takes $\sim 20 - 24$ hours depending on a single 40GB NVIDIA A40/40GB NVIDIA A100 GPU. We note that because our dataset is generated on-the-fly, distributed training does not provide any across multiple devices. We did not find significant differences in model performance by training across one vs. multiple devices. In terms of number of trainable parameters, the embedding network contains $\sim 2.6$ million parameters, while the auto-regressive transforms contain $\sim 3.2$ million parameters, totaling to $\sim 6$ million parameters. We

would like to note that with on-the-fly data generation on the GPU, a typical training run sees $\gtrsim 200$ epochs $\times$ 200 batches per epoch $\times$ 800 batchsize $\sim 32$ million unique simulations. This implies that unlike most "large" neural-networks in the literature today, our network is not over-specified in the sense that training data samples exceed the number of trainable parameters by several factors.

In terms of inference, average sampling time for drawing 20,000 posterior samples, conditioned on new data, takes $\sim 0.05$s on NVIDIA A40 GPU. The same on a Intel Core i7 with 16 cores, takes $\sim 1-2$s. However, it should be noted to create a sky-localization map in the HEALPix format [47], used in the GW data analysis, the density needs to be evaluated across all pixel coordinates. We find the average time to drawn 20K samples, and then evaluate the density across all pixels on the sky for a HEALPix resolution of NSIDE $= 32$ is $\sim 0.6$s. This is due to the choice of the inverse auto-regressive flow, sampling is possible via one forward pass, but evaluating the density is done sequentially across each component. Doubling the resolution, i.e. using NSIDE $= 64$ takes $\sim 1.2$s and NSIDE $= 128$ takes $\sim 2.4$s.

## 6. Results and Performance

We show example posteriors in Figures 3 and 4 from representative higher and lower mass BBH signals. In both cases, the signal injection is performed in a background different from that used during training. The same signal is injected 20 times at different background segments and samples are drawn. The posteriors are shown in the blue colormap. We also perform inference on the same system via nested sampling using `Bilby` [48, 49] and `Dynesty` [50] with 1500 live points and identical priors as `AMPLFI`, repeating 5-times on different background segments. This is shown in the red colormap. We see that there is good consistency among the `AMPLFI` posteriors i.e., distributions in blue colormap fall on top of each other. For higher masses, like the $\mathcal{M} = 45\ M_\odot$, shown in Figure 3, we find the recovery accuracy of the chirp mass to be comparable to nested sampling. Also, considering all nested sampling runs, the mass-ratio accuracy from `AMPLFI` is comparable. However, extrinsic parameters like the inclination, or the sky location are not recovered with the same accuracy as nested sampling, although broad features like the "ring" pattern in the skymap are evident. The same is also true for the inclination posterior, where the inference is degenerate between the true inclination angle, and its supplementary angle. In case of the nested sampling runs, though this degeneracy is broken, the inference may not always select the right "peak" for the inclination. Thus, overall, we find consistency of the `AMPLFI` results with the true parameters of the injection and with nested sampling results, though not as accurate for some parameters.

We also find that the inference for lower mass, $\mathcal{M} = 15\ M_\odot$ BBH system, shown in Figure 4, is worse compared to higher mass. For example, in Figure 4, we see that the $\mathcal{M}$ posteriors recovered by `AMPLFI` is broader compared to that recovered from nested sampling. The extrinsic parameter recovery follows similar pattern as the high-mass
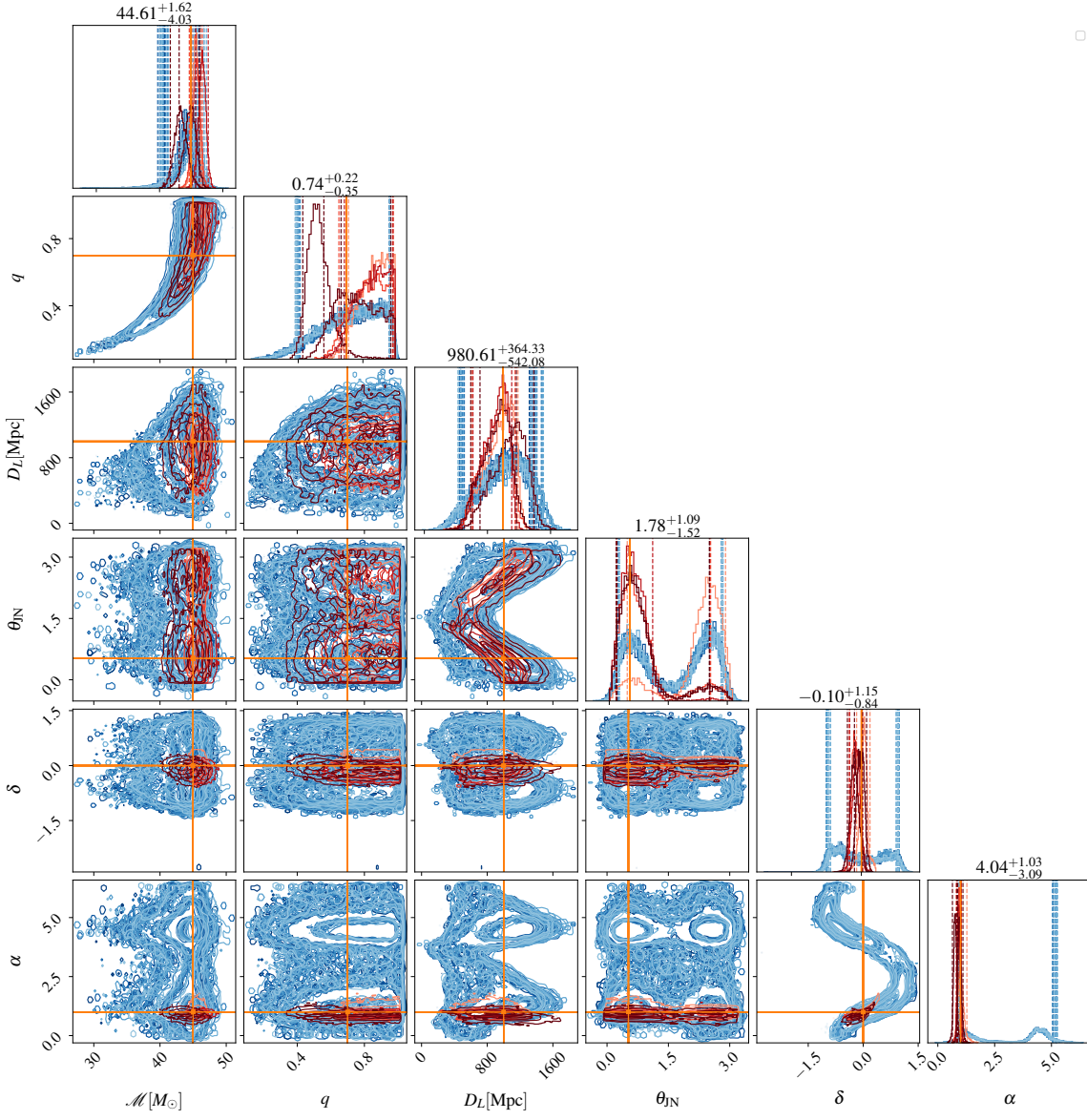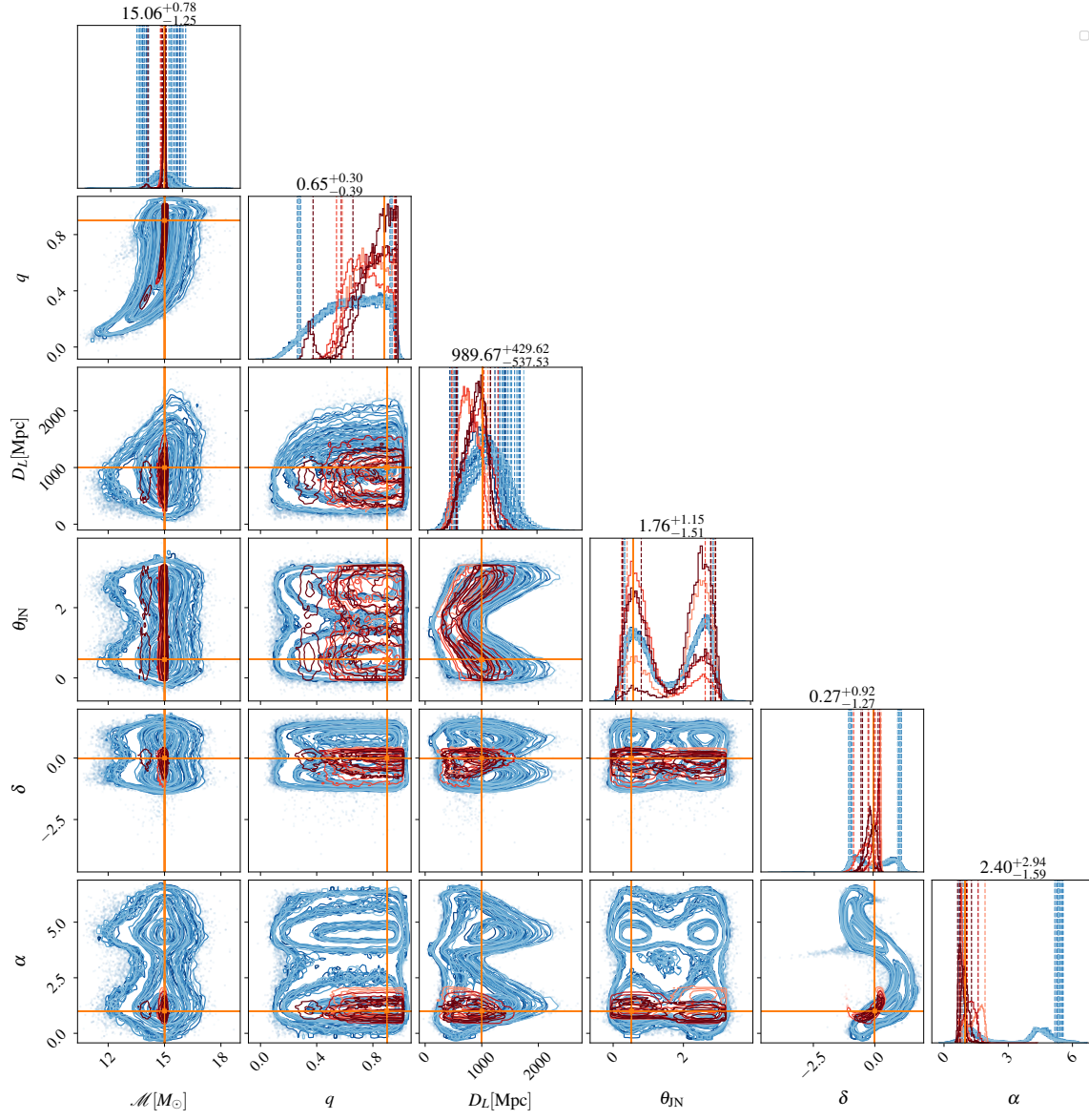
Figure 3: Example posterior for a signal with parameters $\{\mathcal{M} = 45 M_\odot, q = 0.7, D_L = 1000 \text{ Mpc}\}$ injected in 20 different background instances, sampled using `AMPLFI` is shown in sky blue . All posteriors are consistent with one another. Posteriors from the same signal injected in 5 different background instances (same background stretch as the `AMPLFI` injections) and analyzed via nested sampling with Bilby, is overlayed in varying Orange-red colors. We find that parameters like $\mathcal{M}$ and $q$ are consistent in terms of detection uncertainties across different runs. Extrinsic parameters, especially the sky-location, though consistent with the true parameters, shows larger uncertainty with `AMPLFI`.
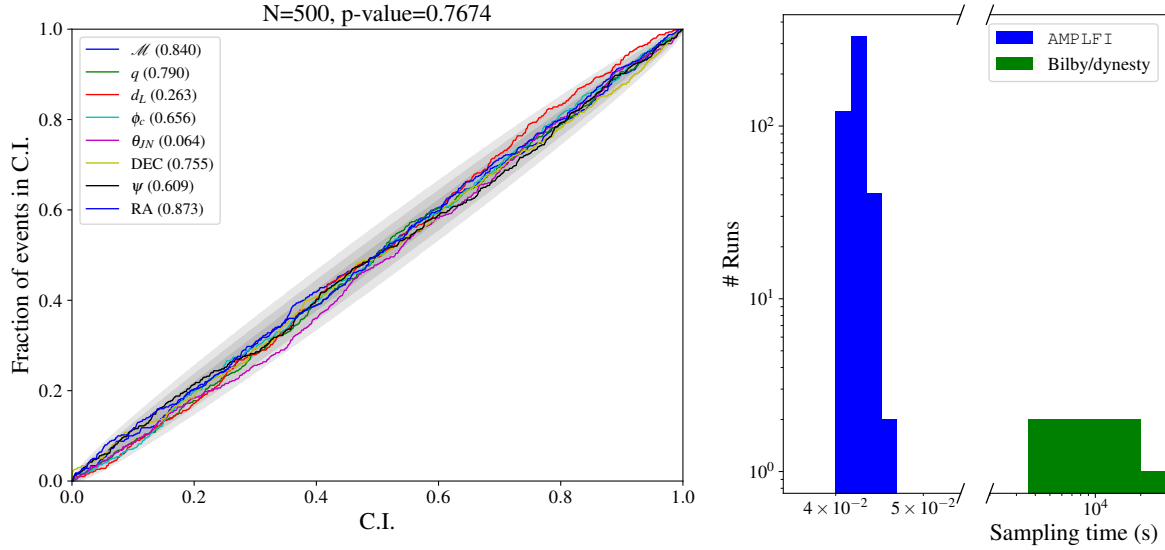
Figure 4: Example posterior for a signal with parameters $\{\mathcal{M} = 15M_\odot, q = 0.9, D_L = 1000 \text{ Mpc}\}$ injected in 20 different background instances, sampled using `AMPLFI` is shown in sky blue. All posteriors are consistent with one another. Posteriors from the same signal injected in 5 different background instances (same background stretch as the `AMPLFI` injections) and analyzed via nested sampling with Bilby, is overlayed in varied Orange-red colors. We find that parameters like $\mathcal{M}$ and $q$ are consistent in terms of detection uncertainties across different runs. Extrinsic parameters, especially the sky-location, though consistent with the true parameters, shows larger uncertainty with `AMPLFI`.

Figure 5: **Left**: Percentile-percentile (PP) plot showing recovery accuracy for 500 BBH injections performed in a testing background, different from training background. The different lines track the cumulative fraction of events within a corresponding confidence interval for the parameters mentioned in Eq. 4. **Right**: Sampling times for `AMPLFI` vs. nested sampling runs done on injections using Bilby, with identical waveform model and prior settings. The nested sampling runs were done with a CPU pool size of 24, and correspond to the runs using in Figure 4. The standard GW likelihood model is used. The `AMPLFI` sampling times correspond to the 500 injections used for the P-P plot on the left.

example. Though broad features of the sky location like the "ring" pattern for two detector is evident, the recovery is not at the level of nested sampling results. The greater consistency of the intrinsic parameters compared to the extrinsic parameters suggests one potential avenue of improvement being to further augment our dataloader in the extrinsic parameters $(\alpha, \delta, \psi)$ which is used to project the signal onto the GW antennae. Signal projection as implemented in `ml4gw` performs this operation on-the-fly on the GPU, hence oversampling the corresponding priors and creating a larger batch is feasible without compromising data generation time, and will be considered in a future implementation.

In Figure 5, we show parameter recovery consistency with true values via a percentile-percentile (PP) plot for 500 simulated BBHs. The parameters of these are sampled from the same prior as that used during training. As for the previous examples above, for these injections, the background is different from the training background. The diagonal trend of the plot shows that there is no bias in the inference with `AMPLFI` across the parameter space. In Figure 6, we show the searched area distributions for the two representative higher mass BBH and lower mass BBH placed at a fiducial distance of $D_L = 1000$ Mpc. We see the expected trend of getting a better search area statistic for higher mass system compared to a lower mass system due to their larger amplitude. However, the searched area is several $\mathcal{O}(1000)$ degrees for the two detector model because of the larger error-bars on the sky coordinates mentioned above.
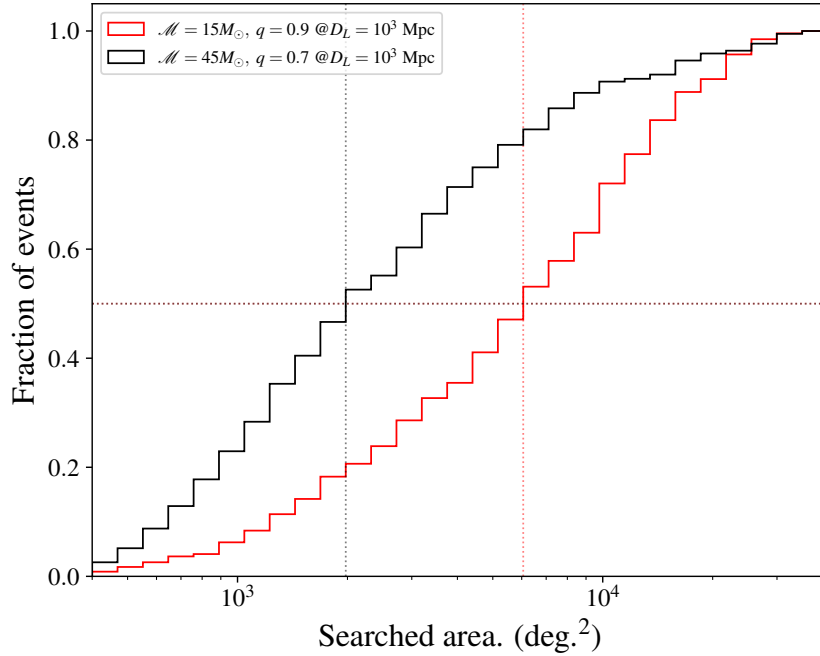
Figure 6: Searched area for a representative system with $\mathcal{M} = 45M_\odot$ and $\mathcal{M} = 15M_\odot$ at $D_L = 1$ Gpc done across the sky. We note that, as expected, larger chirp mass give better searched area as a result of being louder signals. However, the median searched area is $\mathcal{O}(1000)\text{deg}^2$.

## 6.1. Comparison with BAYESTAR skymaps

In this section, we run our model on several O3 events using the data from GWOSC and compare the sky-localization with the corresponding rapid localization method, BAYESTAR [29]. Since our model is trained on only 2-detector (HL) data, we re-run BAYESTAR on the events using only HL SNR timeseries for this comparison. We show the skymap, along with the 90% localization area, for several O3 events that were detected online in Figure 7. The left (right) panel shows the reconstruction using `AMPLFI` (BAYESTAR). The choice of the selected events is that their event parameters, as published in GWTC-3, lie within the prior range used by us during training, and the events are spread over the most of O3 from May 2019 to March 2020 (the event identifiers carry the date of discovery). While our training background is a half-day chunk at the start of O3, we test on events which had been detected over the duration of the entire run. This is done to obtain a qualitative assessment of the impact of changing background i.e., if the model performance greatly degrades when supplied with data several months away from the training background. We find that there is broad similarity in the skymaps between `AMPLFI` and BAYESTAR. There is no trend in terms of the skymaps being more/less constrained. However, given that the events used for Figure 7 range from May 2019 to March 2020, (recall that training background is

limited to half a day in May 2019), we conclude that the model validity is maintained for different background up to several months. This does not imply that the model once trained, may be optimal for the entire run. In fact, one of the requirements for `AMPLFI` is the ability for periodic re-training; however, we anticipate that such re-training may converge quickly given the preliminary observations mentioned here. Furthermore, the suitable re-training cadence is yet to be determined and will be reported in the future.

## 7. Conclusion and Future Work

In this work, we presented a GW parameter estimation algorithm, `AMPLFI`, using likelihood-free inference. This work is one of the efforts to integrate AI algorithms in GW data analysis using tools build as a part of `ml4gw` – like data cleaning using `DeepClean` [51], search of CBCs using `Aframe` [31], anomaly detection algorithm `GWAK` [52]. The use case of `AMPLFI` is to run alongside recently reported neural-network based CBC search, `Aframe`. The intended design is for both `Aframe` + `AMPLFI` to run online, preferably on the same hardware to minimize any communication overhead and reduce the time-to-alert. While to date we have focused on this integration of neural-networks in order to achieve very low-latency outputs, there is reason in principle that `AMPLFI` could not also be paired with any other search pipeline that provides an estimated time of arrival. An important future product of this research will be a standalone version of the software that can be run simply with the provision of trigger time, enabling easy adoption for searches that want to use it. Current real-time GW alert data products include sky-localization and EM-bright source properties apart from the significance, and a derived data-product, P-astro, based on rate of foreground and background triggers. We have not discussed the EM Bright source properties in this work since the analysis, so far, is limited to BBH signal which are expected to be EM-dark. However, the availability of the posterior samples in real-time allows for their straightforward computation by binning the posterior samples, or marginalizing them over several equations of state (see section C.2 of [4]).

As a part of routine end-to-end testing, the LVK has set up a streaming data channel with injections over a 40-days chunk of O3 background. This was used to profile latencies in several components in the alert infrastructure, reported in [4]. Preliminary work has been done toward the online deployment of `Aframe` + `AMPLFI` to analyze this data stream. Based on preliminary testing, we find the net latency of data acquisition by `Aframe`, evaluating significane, passing data to `AMPLFI`, followed by generating posteriors and skymaps takes $\sim 6$ seconds. At the time of writing, candidates are reported to a test instance of GraceDB, the candidate database used by the LVK, however, the view for the same is not public. As a follow-up to the methods reported here, the performance of `Aframe` + `AMPLFI` will be reported on this mock data challenge (MDC) in a future work. The dataset contains $\mathcal{O}(1000)$ BBHs, for which several match-filtering searches and annotation pipelines including BAYESTAR was run. This will contain a systematic comparison with BAYESTAR on the simulated BBHs in the MDC, along

with comparison with online PE results reported as a part of the same study.

Certain aspects of the model requires improvements, for example, the inference on the extrinsic parameters like sky location, and the extension to use spinning waveforms. This will be considered in a future work. Also, we note that the focus has been on BBHs due to their shorter signal duration. However, the main focus of MMA is binary neutron star (BNS) and neutron star black hole (NSBH) systems for which the signal duration can be $\mathcal{O}(\text{min})$ depending on the starting frequency. This makes the input arrays larger by an order of magnitude compared to the ones considered here, and therefore expensive in terms of memory and compute. However, low mass systems *inspiral* for most of that duration and the frequency evolution is described analytically, primarily at newtonian order. Thus, feature extraction from the time-series data, or alternative data representations like q-transforms can be a possible approach toward search and parameter estimation.

Finally, the framework for `AMPLFI` is not limited to CBC signals, and can be extended to burst signal morphologies like sine-Gaussians. This is relevant for running parameter estimation on candidates picked up by pipelines like `GWAK` that look for unmodeled events [52]. In this case, a sine-gaussian parameter estimation may lead to measurement of fundamental features like central frequency or duration.
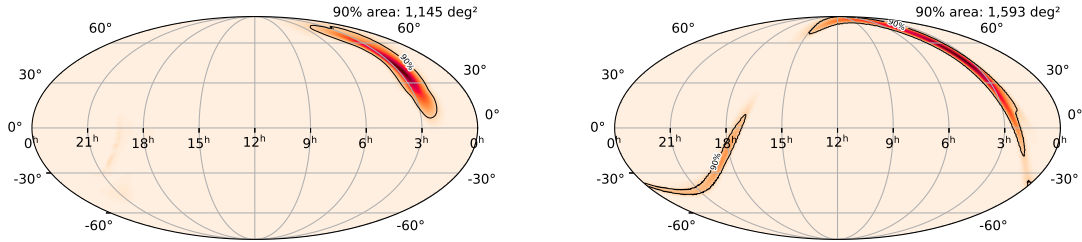
## 8. Acknowledgments

Academia Sinica (AS) and National Science and Technology Council (NSTC) in Taiwan.

S190512at

S190513bm

S190701ah

Figure 7: **Left**: `AMPLFI` skymaps **Right** Bayestar maps from events, using HL data, spanning different months of the third observing run, O3. We find that the skymaps are broadly consistent, although there is no clear trend of one being more constraining. However, it demonstrates model validity to differing backgrounds up to several months from that used in training.

S191215w

90% area: 2,706 deg²

90% area: 1,818 deg²

S200129m

90% area: 1,642 deg²

90% area: 363 deg²

S200311bg

90% area: 947 deg²

90% area: 772 deg²

Figure 7: Continued from above – **Left**: `AMPLFI` skymaps **Right** Bayestar maps.

Table A1: Top 10 hyper-parameter configurations for hyper-parameter optimization runs for the embedding network. The best trial configuration is shown in boldface.

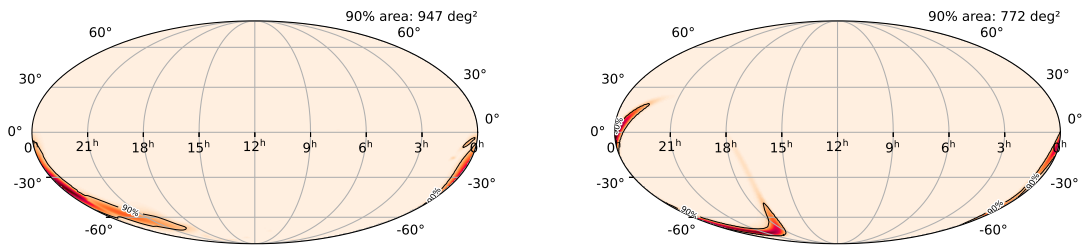| ResNet conf. | kernel size | $\lambda_1$ | $\lambda_2$ | $\lambda_3$ | $D_\gamma$ | LR | Momentum | wt. decay | $N$ | VICReg. |
|---|---|---|---|---|---|---|---|---|---|---|
| **( 5 , 3 , 3 )** | **5** | **1** | **1** | **5** | **8** | $\mathbf{7.16 \cdot 10^{-4}}$ | $\mathbf{8.07 \cdot 10^{-5}}$ | $\mathbf{4.42 \cdot 10^{-4}}$ | **3** | **0.48** |
| ( 4 , 3 , 3 ) | 7 | 1 | 1 | 1 | 11 | $8.97 \cdot 10^{-4}$ | $5.75 \cdot 10^{-5}$ | $9.01 \cdot 10^{-5}$ | 3 | 0.528 |
| ( 4 , 5 , 3 ) | 7 | 1 | 1 | 1 | 9 | $2.08 \cdot 10^{-4}$ | $2.60 \cdot 10^{-4}$ | $3.72 \cdot 10^{-4}$ | 3 | 0.539 |
| ( 4 , 3 , 3 ) | 7 | 1 | 1 | 1 | 10 | $2.62 \cdot 10^{-4}$ | $3.69 \cdot 10^{-5}$ | $1.16 \cdot 10^{-5}$ | 3 | 0.543 |
| ( 4 , 5 , 3 ) | 5 | 5 | 1 | 1 | 8 | $3.48 \cdot 10^{-4}$ | $6.10 \cdot 10^{-4}$ | $7.37 \cdot 10^{-4}$ | 5 | 0.553 |
| ( 5 , 5 , 4 ) | 3 | 5 | 1 | 1 | 8 | $1.31 \cdot 10^{-4}$ | $1.75 \cdot 10^{-5}$ | $1.67 \cdot 10^{-4}$ | 3 | 0.578 |
| ( 3 , 4 , 4 ) | 3 | 5 | 1 | 5 | 9 | $5.38 \cdot 10^{-4}$ | $1.42 \cdot 10^{-5}$ | $1.67 \cdot 10^{-5}$ | 5 | 0.610 |
| ( 4 , 5 , 3 ) | 3 | 1 | 1 | 5 | 10 | $9.76 \cdot 10^{-4}$ | $2.69 \cdot 10^{-4}$ | $4.20 \cdot 10^{-5}$ | 4 | 0.627 |
| ( 4 , 3 , 4 ) | 3 | 1 | 1 | 1 | 9 | $1.37 \cdot 10^{-4}$ | $5.46 \cdot 10^{-5}$ | $1.33 \cdot 10^{-5}$ | 3 | 0.677 |
| ( 4 , 5 , 3 ) | 5 | 5 | 5 | 1 | 8 | $1.64 \cdot 10^{-4}$ | $9.47 \cdot 10^{-4}$ | $1.61 \cdot 10^{-5}$ | 3 | 0.707 |

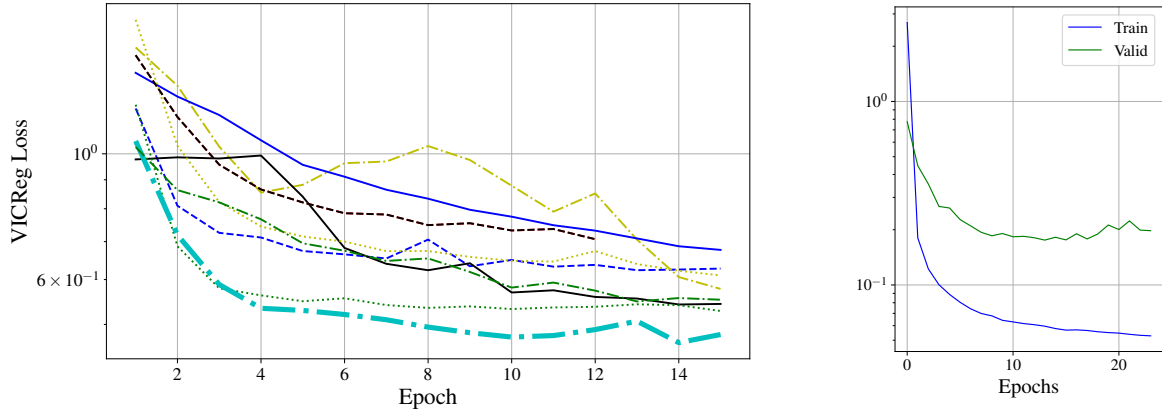## Appendix A. Hyperparameter Tuning of the Embedding Network



Figure A1: **Left**: Avg. VICReg validation loss as a function of training epoch from the top 10 HPO runs in Table A1. The best configuration is plotted in thick solid line; corresponding configuration is boldfaced entry in Table A1. **Right**: Training/Validation loss for best model configuration from the left panel until early-stopping.

To determine the configuration to be used for the embedding network mentioned in Sec. 4, we perform an extensive hyperparameter optimization over the search space of the layers of ResNet, the convolutional kernel size of the ResNet, the dimensionality of the representation, $D_\gamma$, the dimensionality of the expanded space where the VICReg loss is computed; this is tune by a factor, $N$, i.e. the dimensionality of the expanded representation is $N \times D_\gamma$. We use Stochastic gradient descent optimizer and also sample over the learning rate, the weight decay and the momentum terms of the optimizer. A total of $\sim 250$ training runs were performed using asynchronous hyperbanding with early-stopping [55] using the `ray.tune` library [56]. This technique stops poor performing trials allowing more favorable trials to continue. We use a grace period of

Table B1: Top 10 hyper-parameter configurations for run involving 30 epochs. Validation loss is noted at the end of the 30th epoch.

| # transforms | # blocks | # hidden feat. | LR | batch size | wt. decay | val. loss |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| **60** | **6** | **100** | **0.00129** | **800** | **0.00241** | **8.52** |
| 100 | 6 | 150 | 0.000636 | 1000 | 0.000263 | 8.63 |
| 60 | 6 | 120 | 0.00242 | 800 | 0.00471 | 8.65 |
| 60 | 8 | 150 | 0.000488 | 1000 | 0.000553 | 8.70 |
| 80 | 6 | 120 | 0.000442 | 1000 | 0.00161 | 8.71 |
| 60 | 7 | 150 | 0.00103 | 1000 | 0.000271 | 8.72 |
| 80 | 7 | 120 | 0.000318 | 800 | 0.000173 | 8.78 |
| 80 | 6 | 100 | 0.00093 | 1000 | 0.00332 | 8.83 |
| 80 | 8 | 100 | 0.000873 | 1200 | 0.00453 | 8.86 |
| 60 | 6 | 120 | 0.000297 | 1000 | 0.0509 | 8.87 |

3 epochs before half of ongoing trials are stopped. The experiment was carried over 8 workers over 4 A40 GPU taking $\sim 40$ hours. We show the top 10 trial configuration in Table A1 and show the epoch-average VICReg validation loss for the same in the left panel of Figure A1. The right panel of the figure shows the training/validation of the best model configuration trained until early-stopping condition is met.

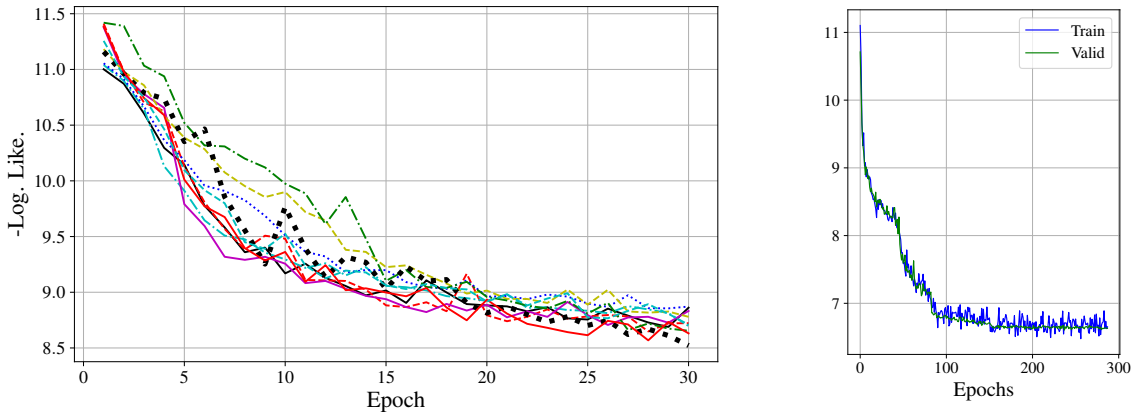## Appendix B. Hyperparameter Tuning of the Normalizing Flow



Figure B1: **Left**: Validation loss from the top 10 HPO runs in Table B1. **Right**: Training/Validation loss for best model from the left panel until early-stopping.

Hyper-parameter optimization is done for $\sim 100$ configurations involving the number of transforms, the configuration of each transform, learning rate, batch size, optimizer weight decay and momentum parameters shown in Table B1 using asynchronous hyper-banding with early stopping [55] using the `ray.tune` library [56]. This stops under-performing runs in favor of allowing better performing runs to continue. We carried out the HPO runs over 4 A40 GPUs which took $\sim 15$ hours. The average

validation loss over validation epoch for the top 5 runs are shown in the left panel of Figure B1. The right panel of the figure shows the training/validation loss of the best model configuration trained until early-stopping condition is met. While all the runs were allowed to run for 30 epochs, most of them are stopped early. The validation loss for the top ten performing runs are shown in the figure. We use the topmost configuration in Table B1 for the results of the paper. The training and validation for this configuration, trained until early-stopping is shown on right panel of Figure B1.

# References

[1] Abbott et al. Observation of gravitational waves from a binary black hole merger. *Phys. Rev. Lett.*, 116(6):061102, 2016.

[2] Abbott et al. GW151226: Observation of Gravitational Waves from a 22-Solar-Mass Binary Black Hole Coalescence. *Phys. Rev. Lett.*, 116(24):241103, 2016.

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

[4] Sushant Sharma Chaudhary, Andrew Toivonen, Gaurav Waratkar, Geoffrey Mo, Deep Chatterjee, Sarah Antier, Patrick Brockill, Michael W. Coughlin, Reed Essick, Shaon Ghosh, Soichiro Morisaki, Pratyusava Baral, Amanda Baylor, Naresh Adhikari, Patrick Brady, Gareth Cabourn Davies, Tito Dal Canton, Marco Cavaglia, Jolien Creighton, Sunil Choudhary, Yu-Kuang Chu, Patrick Clearwater, Luke Davis, Thomas Dent, Marco Drago, Becca Ewing, Patrick Godwin, Weichangfeng Guo, Chad Hanna, Rachael Huxford, Ian Harry, Erik Katsavounidis, Manoj Kovalam, Alvin K.Y. Li, Ryan Magee, Ethan Marx, Duncan Meacher, Cody Messick, Xan Morice-Atkinson, Alexander Pace, Roberto De Pietri, Brandon Piotrzkowski, Soumen Roy, Surabhi Sachdev, Leo P. Singer, Divya Singh, Marek Szczepanczyk, Daniel Tang, Max Trevor, Leo Tsukada, Verónica Villa-Ortega, Linqing Wen, and Daniel Wysocki. Low-latency gravitational wave alert products and their performance at the time of the fourth ligo-virgo-kagra observing run. *Proceedings of the National Academy of Sciences*, 121(18), April 2024.

[5] Kipp Cannon, Romain Cariou, Adrian Chapman, Mireia Crispin-Ortuzar, Nickolas Fotopoulos, Melissa Frei, Chad Hanna, Erin Kara, Drew Keppel, Laura Liao, Stephen Privitera, Antony Searle, Leo Singer, and Alan Weinstein. Toward early-warning detection of gravitational waves from compact binary coalescence. *The Astrophysical Journal*, 748(2):136, March 2012.

[6] Ryan Magee, Deep Chatterjee, Leo P. Singer, Surabhi Sachdev, Manoj Kovalam, Geoffrey Mo, Stuart Anderson, Patrick Brady, Patrick Brockill, Kipp Cannon, Tito Dal Canton, Qi Chu, Patrick Clearwater, Alex Codoreanu, Marco Drago, Patrick Godwin, Shaon Ghosh, Giuseppe Greco, Chad Hanna, Shasvath J. Kapadia, Erik Katsavounidis, Victor Oloworaran, Alexander E. Pace, Fiona Panther, Anwarul Patwary, Roberto De Pietri, Brandon Piotrzkowski, Tanner Prestegard, Luca Rei, Anala K. Sreekumar, Marek J. Szczepańczyk, Vinaya Valsan, Aaron Viets, Madeline Wade, Linqing Wen, and John Zweizig. First demonstration of early warning gravitational-wave alerts. *The Astrophysical Journal Letters*, 910(2):L21, April 2021.

[7] Scott Barthelmy, Eric Burns, Dakota Dutko, Meredith Gibb, Victor Gonzalez-Leon, Tess Jaffe, Ryan Lorek, Israel Martinez, Tom McGlynn, Judy Racusin, David Simpson, Leo Singer, Teresa Sheets, Alan Smale, Dongguen Tak, Gcn, Tach, and Heasarc. Introducing new GCN Kafka broker and web site for transient alerts, https://gcn.nasa.gov. *GRB Coordinates Network*, 32419:1, July 2022.

[8] Samuel Wyatt, Aaron Tohuvavohu, Iair Arcavi, David Sand, Michael Lundquist, D. Andrew Howell, Curtis McCully, Austin Riba, Jamison Burke, and Gravitational Wave Treasure Map Team. Announcing the GW Treasure Map. *GRB Coordinates Network*, 26244:1, November 2019.

[9] Stéfan J. van der Walt, Arien Crellin-Quick, and Joshua S. Bloom. SkyPortal: An astronomical data platform. *Journal of Open Source Software*, 4(37), may 2019.

[10] R. A. Street, M. Bowman, E. S. Saunders, and T. Boroson. General-purpose software for managing astronomical observing programs in the LSST era. In Juan C. Guzman and Jorge Ibsen, editors, *Software and Cyberinfrastructure for Astronomy V*, volume 10707 of *Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series*, page 1070711, July 2018.

[11] A. Buikema, C. Cahillane, G. L. Mansell, C. D. Blair, R. Abbott, C. Adams, R. X. Adhikari, A. Ananyeva, S. Appert, K. Arai, J. S. Areeda, Y. Asali, S. M. Aston, C. Austin, A. M. Baer, M. Ball, S. W. Ballmer, S. Banagiri, D. Barker, L. Barsotti, J. Bartlett, B. K. Berger, J. Betzwieser, D. Bhattacharjee, G. Billingsley, S. Biscans, R. M. Blair, N. Bode, P. Booker, R. Bork, A. Bramley, A. F. Brooks, D. D. Brown, K. C. Cannon, X. Chen, A. A. Ciobanu, F. Clara, S. J. Cooper, K. R. Corley, S. T. Countryman, P. B. Covas, D. C. Coyne, L. E. H. Datrier, D. Davis, C. Di Fronzo, K. L. Dooley, J. C. Driggers, P. Dupej, S. E. Dwyer, A. Effler, T. Etzel, M. Evans, T. M. Evans, J. Feicht, A. Fernandez-Galiana, P. Fritschel, V. V. Frolov, P. Fulda, M. Fyffe, J. A. Giaime, K. D. Giardina, P. Godwin, E. Goetz, S. Gras, C. Gray, R. Gray, A. C. Green, E. K. Gustafson, R. Gustafson, J. Hanks, J. Hanson, T. Hardwick, R. K. Hasskew, M. C. Heintze, A. F. Helmling-Cornell, N. A. Holland, J. D. Jones, S. Kandhasamy, S. Karki, M. Kasprzack, K. Kawabe, N. Kijbunchoo, P. J. King, J. S. Kissel, Rahul Kumar, M. Landry, B. B. Lane, B. Lantz, M. Laxen, Y. K. Lecoeuche, J. Leviton, J. Liu, M. Lormand, A. P. Lundgren, R. Macas, M. MacInnis, D. M. Macleod, S. Márka, Z. Márka, D. V. Martynov, K. Mason, T. J. Massinger, F. Matichard, N. Mavalvala, R. McCarthy, D. E. McClelland, S. McCormick, L. McCuller, J. McIver, T. McRae, G. Mendell, K. Merfeld, E. L. Merilh, F. Meylahn, T. Mistry, R. Mittleman, G. Moreno, C. M. Mow-Lowry, S. Mozzon, A. Mullavey, T. J. N. Nelson, P. Nguyen, L. K. Nuttall, J. Oberling, Richard J. Oram, B. O'Reilly, C. Osthelder, D. J. Ottaway, H. Overmier, J. R. Palamos, W. Parker, E. Payne, A. Pele, R. Penhorwood, C. J. Perez, M. Pirello, H. Radkins, K. E. Ramirez, J. W. Richardson, K. Riles, N. A. Robertson, J. G. Rollins, C. L. Romel, J. H. Romie, M. P. Ross, K. Ryan, T. Sadecki, E. J. Sanchez, L. E. Sanchez, T. R. Saravanan, R. L. Savage, D. Schaetzl, R. Schnabel, R. M. S. Schofield, E. Schwartz, D. Sellers, T. Shaffer, D. Sigg, B. J. J. Slagmolen, J. R. Smith, S. Soni, B. Sorazu, A. P. Spencer, K. A. Strain, L. Sun, M. J. Szczepańczyk, M. Thomas, P. Thomas, K. A. Thorne, K. Toland, C. I. Torrie, G. Traylor, M. Tse, A. L. Urban, G. Vajente, G. Valdes, D. C. Vander-Hyde, P. J. Veitch, K. Venkateswara, G. Venugopalan, A. D. Viets, T. Vo, C. Vorvick, M. Wade, R. L. Ward, J. Warner, B. Weaver, R. Weiss, C. Whittle, B. Willke, C. C. Wipf, L. Xiao, H. Yamamoto, Hang Yu, Haocun Yu, L. Zhang, M. E. Zucker, and J. Zweizig. Sensitivity and performance of the advanced ligo detectors in the third observing run. *Phys. Rev. D*, 102:062003, Sep 2020.

[12] Aasi et al. Advanced ligo. *Classical and Quantum Gravity*, 32(7):074001, 2015.

[13] F. Acernese et al. Advanced Virgo: a second-generation interferometric gravitational wave detector. *Class. Quant. Grav.*, 32(2):024001, 2015.

[14] T Akutsu, M Ando, K Arai, Y Arai, S Araki, A Araya, N Aritomi, Y Aso, S Bae, Y Bae, L Baiotti, R Bajpai, M A Barton, K Cannon, E Capocasa, M Chan, C Chen, K Chen, Y Chen, H Chu, Y K Chu, S Eguchi, Y Enomoto, R Flaminio, Y Fujii, M Fukunaga, M Fukushima, G Ge, A Hagiwara, S Haino, K Hasegawa, H Hayakawa, K Hayama, Y Himemoto, Y Hiranuma, N Hirata, E Hirose, Z Hong, B H Hsieh, C Z Huang, P Huang, Y Huang, B Ikenoue, S Imam, K Inayoshi, Y Inoue, K Ioka, Y Itoh, K Izumi, K Jung, P Jung, T Kajita, M Kamiizumi, N Kanda, G Kang, K Kawaguchi, N Kawai, T Kawasaki, C Kim, J C Kim, W S Kim, Y M Kim, N Kimura, N Kita, H Kitazawa, Y Kojima, K Kokeyama, K Komori, A K H Kong, K Kotake, C Kozakai, R Kozu, R Kumar, J Kume, C Kuo, H S Kuo, S Kuroyanagi, K Kusayanagi, K Kwak, H K Lee, H W Lee, R Lee, M Leonardi, L C C Lin, C Y Lin, F L Lin, G C Liu, L W Luo, M Marchio, Y Michimura, N Mio, O Miyakawa, A Miyamoto, Y Miyazaki, K Miyo, S Miyoki, S Morisaki, Y Moriwaki, K Nagano, S Nagano, K Nakamura, H Nakano, M Nakano,

R Nakashima, T Narikawa, R Negishi, W T Ni, A Nishizawa, Y Obuchi, W Ogaki, J J Oh, S H Oh, M Ohashi, N Ohishi, M Ohkawa, K Okutomi, K Oohara, C P Ooi, S Oshino, K Pan, H Pang, J Park, F E Peña Arellano, I Pinto, N Sago, S Saito, Y Saito, K Sakai, Y Sakai, Y Sakuno, S Sato, T Sato, T Sawada, T Sekiguchi, Y Sekiguchi, S Shibagaki, R Shimizu, T Shimoda, K Shimode, H Shinkai, T Shishido, A Shoda, K Somiya, E J Son, H Sotani, R Sugimoto, T Suzuki, T Suzuki, H Tagoshi, H Takahashi, R Takahashi, A Takamori, S Takano, H Takeda, M Takeda, H Tanaka, K Tanaka, K Tanaka, T Tanaka, T Tanaka, S Tanioka, E N Tapia San Martin, S Telada, T Tomaru, Y Tomigami, T Tomura, F Travasso, L Trozzo, T Tsang, K Tsubono, S Tsuchida, T Tsuzuki, D Tuyenbayev, N Uchikata, T Uchiyama, A Ueda, T Uehara, K Ueno, G Ueshima, F Uraguchi, T Ushiba, M H P M van Putten, H Vocca, J Wang, C Wu, H Wu, S Wu, W-R Xu, T Yamada, K Yamamoto, K Yamamoto, T Yamamoto, K Yokogawa, J Yokoyama, T Yokozawa, T Yoshioka, H Yuzurihara, S Zeidler, Y Zhao, and Z H Zhu. Overview of KAGRA: Detector design and construction history. *Progress of Theoretical and Experimental Physics*, 2021(5):05A101, 08 2020.

[15] B.P. Abbott, R. Abbott, T.D. Abbott, F. Acernese, K. Ackley, C. Adams, T. Adams, P. Addesso, R.X. Adhikari, V.B. Adya, and et al. GW170817: Observation of gravitational waves from a binary neutron star inspiral. *Physical Review Letters*, 119(16), Oct 2017.

[16] B. P. Abbott et al. Multi-messenger Observations of a Binary Neutron Star Merger. *Astrophys. J. Lett.*, 848(2):L12, 2017.

[17] Michael W. Coughlin. Lessons from counterpart searches in LIGO and Virgo's third observing campaign. *Nature Astronomy*, 4:550–552, June 2020.

[18] Priscilla Canizares, Scott E. Field, Jonathan Gair, Vivien Raymond, Rory Smith, and Manuel Tiglio. Accelerated gravitational-wave parameter estimation with reduced order modeling. *Phys. Rev. Lett.*, 114(7):071104, 2015.

[19] Soichiro Morisaki, Rory Smith, Leo Tsukada, Surabhi Sachdev, Simon Stevenson, Colm Talbot, and Aaron Zimmerman. Rapid localization and inference on compact binary coalescences with the advanced ligo-virgo-kagra gravitational-wave detector network, 2023.

[20] Soichiro Morisaki and Vivien Raymond. Rapid Parameter Estimation of Gravitational Waves from Binary Neutron Star Coalescence using Focused Reduced Order Quadrature. *Phys. Rev. D*, 102(10):104020, 2020.

[21] Colm Talbot, Rory Smith, Eric Thrane, and Gregory B. Poole. Parallelized Inference for Gravitational-Wave Astronomy. *Phys. Rev. D*, 100(4):043030, 2019.

[22] Kaze W. K. Wong, Maximiliano Isi, and Thomas D. P. Edwards. Fast gravitational wave parameter estimation without compromises, 2023.

[23] Lalit Pathak, Amit Reza, and Anand S. Sengupta. Fast likelihood evaluation using meshfree approximations for reconstructing compact binary sources. *Phys. Rev. D*, 108:064055, Sep 2023.

[24] Lalit Pathak, Sanket Munishwar, Amit Reza, and Anand S. Sengupta. Prompt sky localization of compact binary sources using a meshfree approximation. *Phys. Rev. D*, 109:024053, Jan 2024.

[25] Hunter Gabbard, Chris Messenger, Ik Siong Heng, Francesco Tonolini, and Roderick Murray-Smith. Bayesian parameter estimation using conditional variational autoencoders for gravitational-wave astronomy. *Nature Physics*, 18(1):112–117, December 2021.

[26] Uddipta Bhardwaj, James Alvey, Benjamin Kurt Miller, Samaya Nissanke, and Christoph Weniger. Sequential simulation-based inference for gravitational wave signals. *Phys. Rev. D*, 108:042004, Aug 2023.

[27] Maximilian Dax, Stephen R. Green, Jonathan Gair, Jakob H. Macke, Alessandra Buonanno, and Bernhard Schölkopf. Real-time gravitational wave science with neural posterior estimation. *Physical Review Letters*, 127(24), December 2021.

[28] Rory Smith and Eric Thrane. Optimal search for an astrophysical gravitational-wave background. *Physical Review X*, 8(2), April 2018.

[29] Leo P. Singer and Larry R. Price. Rapid bayesian position reconstruction for gravitational-wave transients. *Physical Review D*, 93(2), January 2016.

[30] Deep Chatterjee, Shaon Ghosh, Patrick R. Brady, Shasvath J. Kapadia, Andrew L. Miller, Samaya Nissanke, and Francesco Pannarale. A machine learning-based source property inference for compact binary mergers. *The Astrophysical Journal*, 896(1):54, June 2020.

[31] Ethan Marx, William Benoit, Alec Gunny, Rafia Omer, Deep Chatterjee, Ricco C. Venterea, Lauren Wills, Muhammed Saleem, Eric Moreno, Ryan Raikman, Ekaterina Govorkova, Dylan Rankin, Michael W. Coughlin, Philip Harris, and Erik Katsavounidis. A machine-learning pipeline for real-time detection of gravitational waves from compact binary coalescences. 2024.

[32] Sebastian Khan, Sascha Husa, Mark Hannam, Frank Ohme, Michael P¥"urrer, Xisco Jim¥'enez Forteza, and Alejandro Boh¥'e. Frequency-domain gravitational waves from nonprecessing black-hole binaries. II. A phenomenological model for the advanced detector era. *Phys. Rev. D*, 93(4):044007, 2016.

[33] Alessandra Buonanno, Bala R. Iyer, Evan Ochsner, Yi Pan, and B. S. Sathyaprakash. Comparison of post-newtonian templates for compact binary inspiral signals in gravitational-wave detectors. *Physical Review D*, 80(8), October 2009.

[34] LIGO Scientific Collaboration. LIGO Algorithm Library - LALSuite. free software (GPL), 2018.

[35] William Falcon and The PyTorch Lightning team. Pytorch lightning, March 2024.

[36] Yuxin Wu and Kaiming He. Group normalization, 2018.

[37] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015.

[38] Adrien Bardes, Jean Ponce, and Yann LeCun. Vicreg: Variance-invariance-covariance regularization for self-supervised learning, 2022.

[39] Maximilian Dax, Stephen R. Green, Jonathan Gair, Michael Deistler, Bernhard Schölkopf, and Jakob H. Macke. Group equivariant neural posterior estimation, 2023.

[40] Deep Chatterjee, Philip C. Harris, Maanas Goel, Malina Desai, Michael W. Coughlin, and Erik Katsavounidis. Optimizing likelihood-free inference using self-supervised neural symmetry embeddings. *Machine Learning and the Physical Sciences Workshop, NeurIPS*, 2023.

[41] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference, 2021.

[42] Diederik P. Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improving variational inference with inverse autoregressive flow, 2017.

[43] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation, 2018.

[44] Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. Made: Masked autoencoder for distribution estimation. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 881–889, Lille, France, 07–09 Jul 2015. PMLR.

[45] Eli Bingham, Jonathan P. Chen, Martin Jankowiak, Fritz Obermeyer, Neeraj Pradhan, Theofanis Karaletsos, Rohit Singh, Paul Szerlip, Paul Horsfall, and Noah D. Goodman. Pyro: Deep universal probabilistic programming. *Journal of Machine Learning Research*, 20(28):1–6, 2019.

[46] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019.

[47] K. M. Górski, E. Hivon, A. J. Banday, B. D. Wandelt, F. K. Hansen, M. Reinecke, and M. Bartelmann. HEALPix: A Framework for High-Resolution Discretization and Fast Analysis of Data Distributed on the Sphere. *The Astrophysical Journal*, 622:759–771, April 2005.

[48] Gregory Ashton et al. BILBY: A user-friendly Bayesian inference library for gravitational-wave astronomy. *Astrophys. J. Suppl.*, 241(2):27, 2019.

[49] I M Romero-Shaw, C Talbot, S Biscoveanu, V D'Emilio, G Ashton, C P L Berry, S Coughlin, S Galaudage, C Hoy, M Hübner, K S Phukon, M Pitkin, M Rizzo, N Sarin, R Smith, S Stevenson, A Vajpeyi, M Arène, K Athar, S Banagiri, N Bose, M Carney, K Chatziioannou, J A Clark, M Colleoni, R Cotesta, B Edelman, H Estellés, C García-Quirós, Abhirup Ghosh, R Green, C-J Haster, S Husa, D Keitel, A X Kim, F Hernandez-Vivanco, I Magaña Hernandez, C Karathanasis, P D Lasky, N De Lillo, M E Lower, D Macleod, M Mateu-Lucena, A Miller,

M Millhouse, S Morisaki, S H Oh, S Ossokine, E Payne, J Powell, G Pratten, M Pürrer, A Ramos-Buades, V Raymond, E Thrane, J Veitch, D Williams, M J Williams, and L Xiao. Bayesian inference for compact binary coalescences with ¡scp¿bilby¡/scp¿: validation and application to the first ligo–virgo gravitational-wave transient catalogue. *Monthly Notices of the Royal Astronomical Society*, 499(3):3295–3319, September 2020.

[50] Joshua S. Speagle. dynesty: a dynamic nested sampling package for estimating Bayesian posteriors and evidences. *Mon. Not. Roy. Astron. Soc.*, 493(3):3132–3158, 2020.

[51] Muhammed Saleem, Alec Gunny, Chia-Jui Chou, Li-Cheng Yang, Shu-Wei Yeh, Andy H. Y. Chen, Ryan Magee, William Benoit, Tri Nguyen, Pinchen Fan, Deep Chatterjee, Ethan Marx, Eric Moreno, Rafia Omer, Ryan Raikman, Dylan Rankin, Ritwik Sharma, Michael Coughlin, Philip Harris, and Erik Katsavounidis. Demonstration of machine learning-assisted real-time noise regression in gravitational wave detectors, 2023.

[52] Ryan Raikman et al. GWAK: gravitational-wave anomalous knowledge with recurrent autoencoders. *Mach. Learn. Sci. Tech.*, 5(2):025020, 2024.

[53] R. Abbott et al. Open Data from the Third Observing Run of LIGO, Virgo, KAGRA, and GEO. *Astrophys. J. Suppl.*, 267(2):29, 2023.

[54] Rich Abbott et al. Open data from the first and second observing runs of Advanced LIGO and Advanced Virgo. *SoftwareX*, 13:100658, 2021.

[55] Lisha Li, Kevin Jamieson, Afshin Rostamizadeh, Katya Gonina, Moritz Hardt, Benjamin Recht, and Ameet Talwalkar. Massively parallel hyperparameter tuning, 2018.

[56] Richard Liaw, Eric Liang, Robert Nishihara, Philipp Moritz, Joseph E Gonzalez, and Ion Stoica. Tune: A research platform for distributed model selection and training. *arXiv preprint arXiv:1807.05118*, 2018.