Correcting misspecification of calibration uncertainties in gravitational-wave data analysis with efficient reweighting

Tomasz Baka^{1,2}, Mick Wright^{1,2}, Isobel Romero-Shaw³, Christopher P L Berry⁴, Carl-Johan Haster^{5,6}, Charlie Hoy⁷, Colm Talbot⁸, Peter T. H. Pang^{2,1}, Vivien Raymond⁷, Michael J. Williams⁹, Gregory Ashton¹⁰, Vladimi Bossilkov¹¹, Louis Dartez¹¹, Noah Manning¹², John Veitch⁴, Aaron Zimmerman¹³, Ling Sun¹⁴, and Will Farr^{15,16}

¹Institute for Gravitational and Subatomic Physics (GRASP), Utrecht University, 3584 CC Utrecht, The Netherlands ²Nikhef, 1098 XG Amsterdam, The Netherlands ³H.H. Wills Physics Laboratory, Tyndall Avenue, Bristol BS8 1TL, United Kingdom ⁴Institute for Gravitational Research, University of Glasgow, University Avenue, Glasgow, G12 8QQ, United Kingdom ⁵Department of Physics and Astronomy, University of Nevada, Las Vegas, 4505 South Maryland Parkway, Las Vegas, NV 89154, USA ⁶ Nevada Center for Astrophysics, University of Nevada, Las Vegas, NV 89154, USA ⁷ Gravity Exploration Institute, Cardiff University, Cardiff CF24 3AA, United Kingdom ⁸Kavli Institute for Cosmological Physics, The University of Chicago, Chicago, IL 60637, USA ⁹Institute of Cosmology and Gravitation, University of Portsmouth, Portsmouth, PO1 3FX, UK ¹⁰Department of Physics, Royal Holloway, University of London ¹¹LIGO Livingston Observatory, Livingston, Louisiana ¹²Center for Computational Relativity and Gravitation, Rochester Institute of Technology, Rochester, New York 14623, USA ¹³ Weinberg Institute, University of Texas at Austin, Austin, TX 78712, USA ¹⁴OzGrav-ANU, Centre for Gravitational Astrophysics, Research School of Physics and Research School of Astronomy & Astrophysics, The Australian National University, ACT 2601, Australia ¹⁵Department of Physics and Astronomy, Stony Brook University, Stony Brook NY 11794 and ¹⁶Center for Computational Astrophysics, Flatiron Institute, 162 5th Ave, New York NY 10010 (Dated: November 4, 2025)

Ground-based gravitational-wave detectors require calibration to map the measured output of the interferometer to reconstructed strain. This calibration has uncertainties that vary across the frequency band of the detectors; these uncertainties are included in analyses that infer the source properties of gravitational-wave signals. During the fourth observing run of the LIGO-Virgo-KAGRA network, it was discovered that there was an error in how uncertainties for LIGO calibration were handled in these analysis, with a mismatch between the output of the calibration calculation and the assumed input for the source-property inference. For data from the first three observing runs, this should have a small impact on results, as the calibration uncertainties had near-zero means. We present an efficient method to reweight inference results to correct the calibration uncertainties. We verify the accuracy of reweighting through comparison with a selected subset of rerun analyses. Using the reweighting algorithm, we are able to confirm that the error in calibration conventions does not significantly impact any conclusions of previous analyses, with the most significant difference being for the inferred right ascension for GW150914.

I. INTRODUCTION

Since the first gravitational-wave (GW) observation in 2015 [1], the LIGO–Virgo—KAGRA (LVK) Collaboration has observed 218 probable compact binary coalescence (CBC) candidates as of the fourth Gravitational-Wave Transient Catalog (GWTC-4) [2–4], with over 200 candidates detected during the ongoing fourth observing run (O4) [5]. The new information obtained from analyzing these signals has enabled us to better understand the astrophysical population of black holes [6], to place constraints on nuclear equation of state [7], to measure the expansion of the Universe in a new way [8], and to test general relativity (GR) [9], among other advances.

The GW data are measured by the two Advanced LIGO [10] detectors, the Advanced Virgo [11] detector,

or the KAGRA detector [12]. These are Michelson interferometers with Fabry–Pérot resonant cavities. They measure relative differential arm displacement called the strain

$$d(t) = \frac{\Delta L_x - \Delta L_y}{L} \,, \tag{1}$$

where L is the arm length of the interferometer. The resulting discrete time series d(t) is the primary data product for all GW analysis.

What the detectors actually record are the laser power fluctuations at the photodetector. The photodetector data is transformed to a calibrated strain using the measured response of the GW detectors [13–16]. This procedure is subject to both statistical and systemic errors. If we call our frequency domain calibrated data d(f) and the hypothetical perfectly calibrated data $d_{\star}(f)$, those

are related by

$$d_{\star}(f) = \eta_R(f)d(f), \qquad (2)$$

where $\eta_R(f)$ is a frequency-dependent calibration correction factor evaluated at a given time.¹ The calibration methods employed by LVK produce uncertainties on this factor in the form of a distribution.

For the parameter-estimation (PE) process, we need to define the likelihood function $\mathcal{L}(d|\vec{\theta})$ where $\vec{\theta}$ are the parameters of our model. From Eq. (2)

$$d = \frac{1}{\eta_B} (h(\vec{\theta}) + n_\star) = \frac{h(\vec{\theta})}{\eta_B} + n, \qquad (3)$$

where $h(\vec{\theta})$ is the data-analysis template and $n = n_{\star}/\eta_R$ is the measured detector noise. We need to account for potential miscalibration error by modifying each waveform template we calculate during the PE process, even though the effects of calibration on PE are small [17, 18].

During O4, an issue was discovered in the application of this recalibration process in all the LVK's analysis codes. For the LIGO and KAGRA [19] detectors, the uncertainty is reported for η_R , while the Virgo [20–22] detector reports that for $1/\eta_R$. It was incorrectly assumed that LIGO followed the same convention as Virgo and gave the uncertainty on $1/\eta_R$. As a result, instead of accounting for the calibration uncertainty, this introduces an additional error into the analysis.

In this work, we present a method of correcting for this calibration-uncertainty error for the previously produced posteriors from the first three observing runs (O1–O3) by transforming the posterior samples and then reweighting them to the correct likelihood. In Sec. II, we describe different possible approaches to reweighting and argue which is the best. In Sec. III, we verify that we have chosen the best reweighting method. By applying it to GWTC-3 PE results [23] and to the test of modified dispersion relation (MDR) [9, 24], we show that the error leads only to small changes in the posteriors, even when observations from multiple signals are combined.

II. METHODS

A. Implementation of calibration uncertainty

In LVK GW data analysis, the calibration error is folded into the analysis by its inverse $\alpha = 1/\eta_R$. The standard likelihood then takes form [1]

$$\ln \mathcal{L} = -\frac{1}{2} \langle n | n \rangle + C = -\frac{1}{2} \langle d - \alpha h | d - \alpha h \rangle + C, \quad (4)$$

following Eq. (3), where $\langle \cdot | \cdot \rangle$ indicates a noise-weighted inner product [25, 26]. The calibration factor α is then split between its magnitude and phase as

$$\alpha(f) = (1 + \delta A(f)) e^{i\delta\psi(f)}, \qquad (5)$$

where $\delta A = \delta \psi = 0$ corresponds to no calibration correction (the strain is perfectly calibrated). The O1–O3 PE analyses was performed on the data recalibrated from the initial calibration model [23, 27], and therefore, calibration uncertainties were small and centered near zero.

In PE analysis the calibration-error factors $\delta A(f)$ and $\delta \psi(f)$ are modeled as cubic splines [28],

$$\delta A(f) = \text{Spline}(f; \delta A_i),$$
 (6)

$$\delta \psi(f) = \text{Spline}(f; \delta \psi_i),$$
 (7)

where δA_i , $\delta \psi_i$ are the values at fixed frequency points these are the calibration parameters in the PE analysis. For each of these parameters, we set up Gaussian priors,

$$\pi(\delta A_i) = \mathcal{N}(\mu_i, \sigma_i), \qquad (8)$$

$$\pi(\delta\psi_i) = \mathcal{N}(\mu_i', \sigma_i'), \qquad (9)$$

chosen to best fit the error estimates on the calibration correction factor [14–16]. These priors have historically been assumed to be for α in the PE analysis, but were set for η_R in LIGO calibration. Confusing the two for small calibration uncertainties gives the mean of the prior the wrong sign, and does not affect the width (see Sec. II F). In general, misspecification of the priors will lead to the analysis not properly accounting for calibration uncertainty.

B. General description of reweighting

To update results with incorrect calibration priors, analyses could be rerun with updated conventions. However, this is computationally expensive. An alternative is to reweight existing posteriors.

Let us begin by describing the general procedure of reweighting posteriors. Starting with samples $\{\theta\}$, drawn from a posterior p, our aim is to obtain samples, $\{\vartheta\}$, drawn from our target posterior p':

$$\vartheta \sim p'(\vartheta|d)d\vartheta$$
 (10)

$$\theta \sim p(\theta|d)d\theta$$
, (11)

where the posteriors are conditioned on the observed data d. In principle, the two parameter sets may be drawn from differing distributions. However, we may create a map between the samples:

$$\vartheta = \vartheta(\theta) = F(\theta), \tag{12}$$

where we require F to be a single-valued function whose image is the whole domain of $p'(\vartheta|d)$.

With this map and the requirement, we do not have to sample from p' to obtain ϑ —we may simply take our

We will drop the explicit frequency dependence and assume that the data is always in the frequency domain, unless specified otherwise.

original samples and use the mapping to transform them and assign each sample a weighting factor:

$$w(\vartheta) = \frac{p'(\vartheta|d)d\vartheta}{p(\theta|d)d\theta}.$$
 (13)

We can then use rejection sampling to obtain explicitly samples drawn from the new posterior, or apply the weighting when calculating any derived statistical quantities.

Bayes' Theorem states:

$$p(\theta|d) = \frac{\mathcal{L}(d|\theta)\pi(\theta)}{\mathcal{Z}}, \qquad (14)$$

where \mathcal{L} is the likelihood, π is the prior and \mathcal{Z} is the evidence—a normalization factor. Applying this, we may expand our definition of the weights:

$$w(\vartheta) = \frac{\mathcal{L}'(d|\vartheta)\pi'(\vartheta)d\vartheta}{\mathcal{L}(d|\theta)\pi(\theta)d\theta},\tag{15}$$

noting that we have ignored the evidences as these contribute merely a multiplicative factor that does not affect the weighting process.

This procedure is valid so long as the original posterior encloses the target posterior and we are not restricted in our transformation choice between ϑ and θ . In practical applications, however, there is a finite number of samples and it is possible for the reweighting procedure to be sufficiently inefficient as to reject the vast majority of samples.

In the rest of this section, we examine in detail three reweighting methods, arranged by complexity and discuss their applicability to the task of reweighting to account for the calibration-envelope error.

C. Prior reweighting

The simplest approach to reweighting the posterior is reweighting the prior—requiring only the analysis samples and the priors used and avoiding the need of recomputing the likelihood as is the case for the more complex methods.

In this approach, the likelihood is treated as correct—the calibration envelope should be applied to the template and thus the incorrect interpretation has resulted in incorrect priors. As such, $\mathcal{L}' = \mathcal{L}$, $\vartheta = F(\theta) = \theta$, and Eq. (15) simplifies to:

$$w(\theta) = \frac{\pi'(\theta)}{\pi(\theta)}. (16)$$

The form of the calibration correction given in Eq. (5) shows that the incorrect interpretation was:

$$\frac{1}{\alpha(f)} \approx (1 - \delta A(f)) e^{-i\delta\psi(f)}, \tag{17}$$

for small amplitude corrections $\delta A_i \ll 1$. The original priors on the calibration-envelope parameters are given by Gaussian distributions—see Eq. (8)—and the change in the prior distribution is a change in the sign of the mean of the distribution.

The reweighting factor for any single calibration parameter (single amplitude or phase node), referred here as x, is then:

$$w(x) = \frac{\mathcal{N}(-\mu, \sigma)}{\mathcal{N}(\mu, \sigma)}$$

$$= \exp\left(\frac{(x+\mu)^2}{2\sigma^2} - \frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$= \exp\left(\frac{2x^2 + 2\mu x + \mu^2 - 2x^2 + 2\mu x - \mu^2}{2\sigma^2}\right)$$

$$= \exp\left(2\frac{\mu x}{\sigma^2}\right). \tag{18}$$

We can quantify the performance of the reweighting in terms of the effective sample size, which is the effective number of independent samples after the reweighting process and is given by:

$$n_{\text{eff}} = \frac{\left(\sum_{i=1}^{N} w_i\right)^2}{\sum_{i=1}^{N} w_i^2},$$
 (19)

where N is the total number of samples and w_i is the importance weight of the i^{th} sample. From the effective sample size, we then define

$$\epsilon = \frac{n_{\text{eff}}}{N} = \frac{\langle w \rangle^2}{\langle w^2 \rangle}, \tag{20}$$

as the reweighting efficiency, where $\langle \dots \rangle$ denotes average over all the samples. This quantity may vary between 0—in which case all samples are lost in the reweighting process—and 1—all samples are kept.

If we assume that the posteriors of calibration parameters are equal to their Gaussian priors², then we can calculate the reweighting efficiency exactly. The reweighting factor for any single calibration parameter x is given by

 $^{^2}$ In GWTC-3 all posteriors of calibration parameters were dominated by the prior [29, 30].

$$\langle w \rangle = \int dx \, w(x) \mathcal{N}(\mu, \sigma)$$

$$= \int dx \, \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-(x-\mu)^2 + 4x\mu}{2\sigma^2}\right)$$

$$= \exp\left(\frac{4\mu^2}{\sigma^2}\right) \qquad (21)$$

$$\langle w^2 \rangle = \int dx \, w(x)^2 \mathcal{N}(\mu, \sigma)$$

$$= \int dx \, \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-(x-\mu)^2 + 8x\mu}{2\sigma^2}\right)$$

$$= \exp\left(\frac{12\mu^2}{\sigma^2}\right), \qquad (22)$$

and so the efficiency per parameter is

$$\epsilon = \frac{n_{\text{eff}}}{N} = \frac{\langle w \rangle^2}{\langle w^2 \rangle} = \exp\left(\frac{-4\mu^2}{\sigma^2}\right).$$
(23)

For the calibration priors $\mu \approx \sigma/2$ on average, and there are up to 40 independent calibration priors (10 frequency nodes and 2 real values per node per detector that observes a GW event; Virgo detector does not count as its calibration priors are symmetric) with similar weights. This gives a total efficiency $\epsilon \sim \exp(-40)$. No samples survive the reweighting process.

D. Likelihood reweighting

A significant drawback of the prior reweighting approach is that it is sensitive to the calibration prior even if the joint posterior of the non-calibration parameters is not sensitive to the calibration at all. We may mitigate this problem by interpreting the original prior as correct and the likelihood as being the source of the calibration error.

In this approach, the corrected likelihood has the form:

$$\mathcal{L}'(d|\theta) = \mathcal{L}(d|g(\theta)), \tag{24}$$

where

$$g(\vartheta) = \begin{cases} (1+\vartheta)^{-1} - 1 \approx -\vartheta, & \vartheta = \delta A_i, \\ -\vartheta, & \vartheta = \delta \psi_i, \end{cases}$$
 (25)

follows from the form of the calibration correction Eq. (5). Similarly to the prior reweighting approach, the samples remain untransformed, but this time, the priors remain the same. This simplifies Eq. (15) to

$$w(\theta) = \frac{\mathcal{L}(d|g(\theta))}{\mathcal{L}(d|\theta)} \approx \frac{\mathcal{L}(d|-\theta)}{\mathcal{L}(d|\theta)}.$$
 (26)

However, this method is also insufficient to handle the reweighting. To illustrate this, consider a toy model

where the calibration correction is simply a constant amplitude correction across all of the frequencies, i.e., $\alpha(f) = 1 + \theta = 1 + \theta$. For small calibrations, the correct calibration likelihood (4) takes the form:

$$\ln \mathcal{L}' = -\frac{1}{2} \langle d - (1+\theta)h|d - (1+\theta)h \rangle + C$$

$$\approx -\frac{1}{2} \langle d - h|d - h \rangle + C + \langle d - h|h \rangle \theta, \qquad (27)$$

and the original likelihood takes the form

$$\ln \mathcal{L} \approx -\frac{1}{2} \langle d - h | d - h \rangle + C - \langle d - h | h \rangle \theta.$$
 (28)

The weights in this case may then be expressed as:

$$\ln w(\theta) = \ln \mathcal{L}' - \ln \mathcal{L} = 2 \langle d - h | h \rangle \theta, \tag{29}$$

where θ is the amplitude sample and the other factor is the product of the waveform template with the residual. Assuming both factors are independently normally distributed,

$$\langle d - h | h \rangle \sim \mathcal{N}(0, \rho) \,,$$
 (30)

$$\theta \sim \mathcal{N}(\mu, \sigma)$$
. (31)

The distribution of their product is given by

$$\ln w \sim \mathcal{N}(0, 2\mu\rho) + p(\ln w), \qquad (32)$$

where

$$p(\ln w) = \frac{K_0 \left(\frac{|\ln w|}{2\sigma\rho}\right)}{2\pi\sigma\rho}, \qquad (33)$$

and

$$K_0(x) = \int_0^\infty dt \frac{\cos(xt)}{\sqrt{t^2 + 1}}$$
 (34)

is a modified Bessel function of the second kind. The spread in weights is then given by

$$\Delta \ln w(\theta) = \sqrt{4\mu^2 \rho^2 + \text{Var}\left[p(\ln w)\right]} = 2\rho \sqrt{\mu^2 + \sigma^2}$$
(35)

In an example case of the GW200129_065458, an O3 observation that has relatively high signal-to-noise ratio [23], $\rho \approx 26.7$, $\sigma \approx 0.013$ and $\mu \approx 0.005$ meaning that $\Delta \ln w(\theta) \approx 0.7$ and the resampling efficiency is $\epsilon \approx 10\%^3$. While much more efficient than prior reweighting, this approach can still reject most of the posterior samples.

 $^{^3}$ This was computed numerically as the distribution of w is poorly approximated by log-normal distribution.

E. Sample transformation and likelihood reweighting

The failure of the previous method is that when the calibration posterior differs significantly from the calibration prior—i.e., in the scenario where the calibration has a measurable effect on the likelihood—the original calibration samples do not sufficiently cover the correct region of parameter space, leading to a significant loss of samples.

In the previous methods, we have left the samples untransformed (i.e., $\theta=\vartheta$). However, we are permitted to choose an endomorphic transformation of the calibration parameters to move the samples from the location of the original posterior towards the expected location of the corrected posterior. To choose a suitable such transformation function $F(\theta)$, we force it to follow the symmetry of the prior

$$\pi(\theta)d\theta = \pi(\theta)d\theta. \tag{36}$$

The calibration parameter priors are Gaussian, Eq. (8), and as such, we may immediately see that the desired transformation function is a reflection along the mean:

$$F(\theta) = 2\mu - \theta \,, \tag{37}$$

which will yield

$$\pi(\theta)d\theta = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mu - 2\mu + \vartheta)^2}{2\sigma^2}\right) |d(2\mu - \vartheta)|$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\mu - \vartheta)^2}{2\sigma^2}\right) d\vartheta$$

$$= \pi(\vartheta)d\vartheta. \tag{38}$$

Employing this transformation function simplifies Eq. (15) to:

$$w(\vartheta) = \frac{\mathcal{L}'(d|\vartheta)}{\mathcal{L}(d|\theta)} = \frac{\mathcal{L}'(d|F(\theta))}{\mathcal{L}(d|\theta)}$$
$$= \frac{\mathcal{L}(d|g(F(\theta)))}{\mathcal{L}(d|\theta)} \approx \frac{\mathcal{L}(d|\theta - 2\mu)}{\mathcal{L}(d|\theta)}, \quad (39)$$

where $g(\theta)$ is the same function as defined in Eq. (25).

To illustrate the improved efficiency of this method, we adopt the same toy model and follow a similar approach as the pure likelihood reweighting case above. In this case, the likelihoods are given by

$$\ln \mathcal{L}' \approx -\frac{1}{2} \langle d - h | d - h \rangle + C + \langle d - h | h \rangle \theta \qquad (40)$$

$$\ln \mathcal{L} \approx \ln \mathcal{L}' - 2 \langle d - h | h \rangle \mu, \tag{41}$$

the expression for the weights takes form

$$\ln w(\theta) = \ln \mathcal{L}' - \ln \mathcal{L} = 2 \langle d - h | h \rangle \mu, \tag{42}$$

and again assuming Gaussian distribution for $\langle d-h|h\rangle$ from Eq. (30), the distribution of weights is given by the log-normal distribution

$$ln w(\theta) \sim \mathcal{N}(0, 2\mu\rho).$$
(43)

The reweighting efficiency is then

$$\langle w \rangle = \exp\left(2\mu^2 \rho^2\right) \,, \tag{44}$$

$$\langle w^2 \rangle = \exp\left(8\mu^2 \rho^2\right) \,, \tag{45}$$

$$\epsilon = \frac{\langle w \rangle^2}{\langle w^2 \rangle} = \exp\left(-4\mu^2 \rho^2\right) \,.$$
 (46)

Again, using the case of GW200129_065458 as an example, $\rho \approx 26.7$, $\sigma \approx 0.013$, and $\mu \approx 0.005$, so the reweighting efficiency becomes $\epsilon \approx 90\%$, which is high enough to retain most samples.

Having selected this approach, we may now summarise the entire reweighting process that may be applied to correct GW analyses for the calibration error. The procedure is as follows for BILBY [31, 32] PE results:

- Retrieve the original posterior from the Result object.
- Retrieve the original likelihood. If the run was performed with marginalisation of any of the parameters, then this must be recomputed, as by default BILBY Result objects store only the marginalised likelihoods and not the true likelihoods.
- 3. Perform transformations of samples:
 - (a) Transform each calibration sample, θ , using Eq. (37).
 - (b) For each sample of the calibration amplitude, θ , replace it with $(1 + \theta)^{-1} 1$.
 - (c) For each sample of the calibration phase, θ , replace it with $-\theta$.
- 4. Compute the likelihood for each sample point.
- 5. Perform steps 3(b) and 3(c) again to undo them.
- 6. Using the recomputed old likelihood \mathcal{L} and newly computed likelihood \mathcal{L}' , assign the weights $w = \exp(\ln \mathcal{L}' \ln \mathcal{L})$.
- 7. These weighted samples now describe the posterior under the correct calibration model.
- 8. Adjust the evidence to account for the reweighting $\ln Z' = \ln Z + \ln \sum_i^N w_i \ln N$

F. Expected impact on calibration parameter posteriors

Having demonstrated the efficiency of the sample transformation method, we now demonstrate that the method transforms the calibration parameter posteriors closer to the correct posteriors, in the case of small calibration corrections.

Let $\pi(\theta) = \mathcal{N}(\mu, \sigma^2)$ be the prior on a calibration parameter, and let $\mathcal{L}_1(d|\theta)$ and $\mathcal{L}_2(d|\theta)$ be the likelihoods

for the original and the correct calibration models respectively. For small calibrations, $\mathcal{L}_2(d|\theta) = \mathcal{L}_1(d|-\theta)$. The respective posteriors are thus given by:

$$p_1(\theta) = \mathcal{L}_1(d|\theta)\pi(\theta),$$

$$p_2(\theta) = \mathcal{L}_2(d|\theta)\pi(\theta).$$
(47)

For simplicity during this demonstration, we assume that the marginalised calibration likelihoods are also normally distributed similarly to the prior:

$$\mathcal{L}_1(d|\theta) = \mathcal{N}(\mu', \sigma'^2)$$

$$\mathcal{L}_2(d|\theta) = \mathcal{L}_1(d|-\theta) = \mathcal{N}(-\mu', \sigma'^2), \qquad (48)$$

but with standard deviations much wider than those of the priors ($\sigma' \gg \sigma$). This latter statement derives from previous analyses demonstrating that the calibration posteriors are dominated by information from the prior [29, 30].

The product of two normal distributions is another normal distribution:

$$\mathcal{N}(\mu_A, \sigma_A^2) \times \mathcal{N}(\mu_B, \sigma_B^2)$$

$$\propto \mathcal{N}\left(\frac{\mu_A \sigma_B^2 + \mu_B \sigma_A^2}{\sigma_A^2 + \sigma_B^2}, (\sigma_A^{-2} + \sigma_B^{-2})^{-1}\right)$$
(49)

The means (μ_1, μ_2) and standard deviations (σ_1, σ_2) of the posteriors are therefore

$$\mu_{1} = \frac{\mu \sigma'^{2} + \mu' \sigma^{2}}{\sigma^{2} + \sigma'^{2}} \approx \mu + (\mu' - \mu) \frac{\sigma^{2}}{\sigma'^{2}}$$

$$\sigma_{1}^{2} = (\sigma^{-2} + \sigma'^{-2})^{-1} \approx \sigma^{2} \left(1 - \frac{\sigma^{2}}{\sigma'^{2}}\right), \quad (50)$$

for the original posterior and

$$\mu_2 \approx \mu + (-\mu' - \mu) \frac{\sigma^2}{\sigma'^2} = 2\mu - \mu_1 - 2\mu \frac{\sigma^2}{\sigma'^2}$$

$$\approx 2\mu - \mu_1$$

$$\sigma_2^2 \approx \sigma^2 \left(1 - \frac{\sigma^2}{\sigma'^2} \right) = \sigma_1^2, \tag{51}$$

for the corrected posterior. From this expression, we see the expectation that the corrected calibration posteriors will be mirrored around the mean of the prior compared with the original posteriors, justifying our choice of transformation.

III. RESULTS

To ascertain how much PE results are affected by the calibration issue, we investigate the effects on two datasets: parameter estimation results from GWTC-3 [23] and the test for a MDR in GWTC-3 [9].

GWTC-3 consists of 90 CBC signals with a probability of astrophysical origin greater than 0.5 [23].

In this catalog, a number of PE analyses were employed. These analyses were performed using a variety of waveform models, predominantly: IMRPHENOMXPHM [33], IMRPHENOMNSBH [34], SEOBNRV4_ROM_NRTIDALV2_NSBH [35], and IMRPHENOMP_NRTIDAL [36, 37] and were performed either with the BILBY [31, 32, 38] Bayesian inference package using the DYNESTY nested sampler [39], or the LALINFERENCE package using its in-built Markov-chain Monte Carlo sampler [40–43]. Additionally, analyses using the more computationally expensive SEOBNRV4PHM [44] analyses were performed with the RIFT inference software [45–47].

RIFT analyses do not model the calibration uncertainty during the sampling process. Instead, RIFT introduces it in a post-processing step by reweighting the likelihood [17, 27]. As such, the error may be accounted for by rerunning the post-processing step with corrected calibration envelopes. Both Bilby and LALINFERENCE, however, model the calibration uncertainty by sampling over draws from the prior, and thus require one of the methods outlined in Sec. II to correct the results. Whilst this process could in principle be done for both pipelines, in the case of LALINFERENCE the saved results and metadata proved insufficient to accurately reconstruct the likelihood. As such, this work will focus only on reweighting BILBY results in which the likelihood object is directly saved and may thus be easily and accurately reconstructed. Of the GWTC-3 signals, only 9 were analysed with LAL-Inference instead of Bilby—those being GW170608. GW170187, GW190707_093326, GW190720_000836, GW190725_174728, GW190728_064510, GW190814. GW190924_021846 and GW190917_114630 [27, 48].

The second dataset used to measure the impact of the calibration-envelope error are the subset of 43 events selected to use for investigation of potential signatures of MDR based on GWTC-3 data performed in Baka $et\ al.$ [49]. This analysis was performed using the BILBY inference framework, the DYNESTY nested sampler, and the IMRPHENOMXPHM waveform model. Each event was analysed for 10 values of the α parameter defined below yielding a total of 430 PE posteriors for investigation.

To briefly summarise the MDR investigation, the dispersion relation is modified to [50, 51]:

$$E^2 = (pc)^2 + A_\alpha (pc)^\alpha , \qquad (52)$$

and we estimate the posterior of the amplitude parameter A_{α} . As this amplitude is a property of space-time and therefore common between all the signals, the investigation creates a combined posterior from the multiplication of individual event posteriors on the amplitude. As such, the MDR dataset provides a significant test of the impact of the calibration-envelope error by allowing us to probe even smaller systematic biases of the posterior that may not be indicated from examination of individual event posteriors.

A. Comparison of reweighting methodologies

We first demonstrate the performance of each of the reweighting methods outlined in Sec. II. From our theoretical discussion in that section, we would expect the performance of each method to improve in the order that they were discussed.

We quantify the performance of the method in terms of the reweighting efficiency defined in Eq. (20). Our aim is to select the methodology that maximises the value of ϵ .

In Fig. 1, we show the reweighting efficiency for each of the three methods on the 430 posteriors in the MDR dataset. We can see that the reweighting efficiency is following the theoretical predictions outlined in our previous discussion. For the prior reweighting approach, $\epsilon \in [4 \times 10^{-5}; 6 \times 10^{-3}]$ which, given the number of samples is of order 10³–10⁴, means that the resulting posterior consists of only a handful of samples. As such, this method fails completely. In the likelihood reweighting approach, the maximal efficiency is 0.897 with 90% of values above 0.123—below this value posteriors begin to suffer from low resolution. Of the 10% of cases below this value, a number do fail entirely with the worst-case scenario efficiency being 5×10^{-4} . Turning to the sample transformation and likelihood reweighting scenario, the situation markedly improves with the efficiency now ranging between 0.903 and 0.997, i.e., significant numbers of samples are retained which allows accurate representation of the new posterior.

This allows us to draw the conclusion that of the three techniques, only the sample transformation and likelihood reweighting approach is suitable and we may discard the other two. As such, in the following sections any reference to reweighted samples, posteriors, etc., refers specifically to reweighting done using this approach.

B. Effect on GWTC-3 PE posteriors

With our reweighting technique decided, we turn now to assess the impact of the calibration envelope error on the posteriors by quantification of the change in posteriors after the reweighting of the GWTC-3 PE dataset. This quantification is done via the Jesnsen–Shannon divergence (JSD) [52] between the two distributions. The JSD is defined between distributions p and q as:

$$JSD(p,q) = \frac{1}{2}D(p||m) + \frac{1}{2}D(q||m), \qquad (53)$$

where m = (p+q)/2 is the mixture distribution and

$$D(p||q) = \int dx \, p(x) \ln \left(\frac{p(x)}{q(x)}\right) \tag{54}$$

is the relative entropy [53]. For sampling performed with BILBY using the DYNESTY nested sampler, the JSD

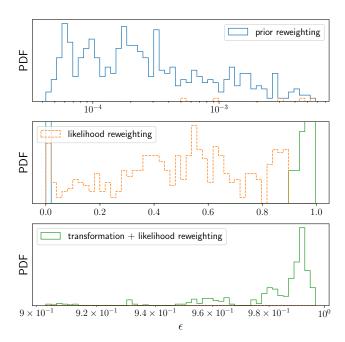


FIG. 1. Performance of different reweighting approaches quantified by their reweighting efficiency ϵ . Middle: All the methods plotted together. Top: Close-up of the low-efficiency region. Bottom: Close-up of the high-efficiency region.

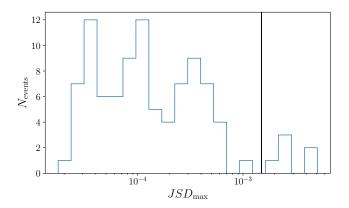


FIG. 2. Histogram of the maximum JSD among the one-dimensional posteriors for GWTC-3 candidates [23,27]. Only 6 analyses have JSD above $0.0015\,\mathrm{nat}$, denoted by the vertical line.

between one-dimensional marginalised posteriors is expected to be up to 0.0015 nat due to statistical fluctuations of the sampler [32]. As such, this is the criterion we adopt to determine how often the calibration envelope error leads to statistically distinct posteriors. For each of the events, we compute for all of the standard CBC parameters—masses, spins, distances, etc—the JSDs between the reweighted and original posteriors. The maximum of these is then compared with the criterion threshold value and if it is in excess of it, the posteriors are determined to be significantly different.

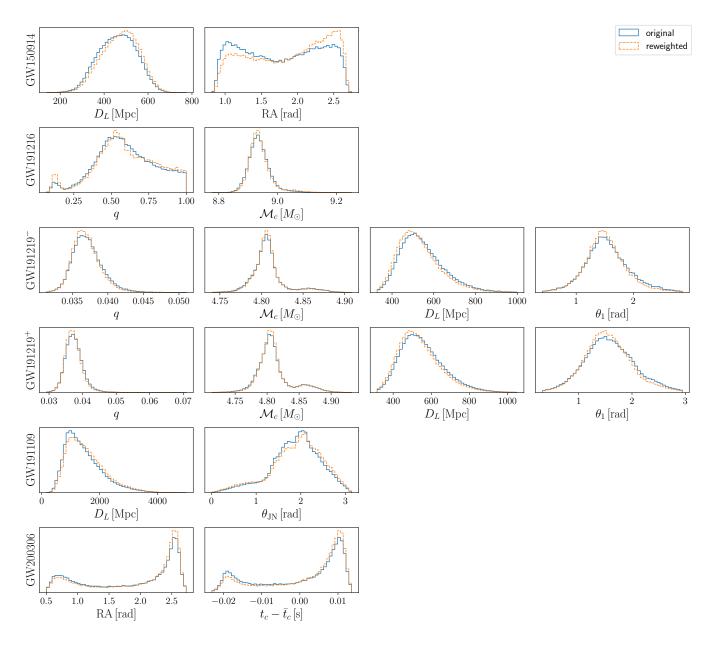


FIG. 3. Comparison of all the incorrect and the reweighted posteriors of GWTC-3 signals with JSD > 0.0015 nat between the posteriors. For GW191219_163120, both the analysis with low- ($^-$) and high-spin ($^+$) priors were affected by the calibration issue. The parameters shown are luminosity distance (D_L), right ascension (RA), mass ratio (q), chirp mass (\mathcal{M}), tilt angle of the primary spin (θ_1), angle between the line-of-sight and the total angular momentum ($\theta_{\rm JN}$) and time of coalescence ($t_{\rm c}$) relative to its mean.

The results of this process are shown Fig. 2, in which the threshold value is noted as a vertical line. The vast majority of results are below the threshold value indicating that their posteriors are unaffected by the calibration envelope error. There are 6 analyses, however, which do show significant differences in their posteriors: GW150914, GW191216_213338, GW191219_163120 (this is the case for both the low and high spin prior analyses—high and low spin prior analyses are carried out for events that may contain neutron stars), GW191109_010717, and GW200306_093714. For these, we show the one-

dimensional posteriors that have significant differences in Fig. 3. The differences are relatively minor—with the exception of the right-ascension posterior for GW150914, in which the dominant peak of the bimodal posterior distribution is reversed. Sky localization is particularly impacted by calibration uncertainty [17, 54], and the localization for GW150914 is understood to vary with the adopted calibration uncertainty [1, 55], so this is unsurprising.

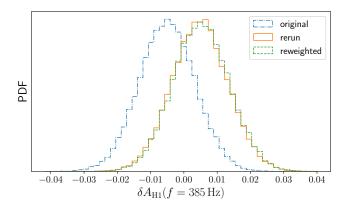


FIG. 4. Fractional calibration amplitude $(\delta A_{\rm H\,I})$ posterior for the LIGO Hanford detector at the spline node at 385 Hz. The reweighting procedure transforms the original posterior to closely match the posterior obtained with the correct calibration model.

1. GW150914

As GW159014 is the signal with the most noticeable difference between the original and reweighted posteriors, we select it to test the validity of the reweighting method. This is done by performing two PE investigations on the strain data: one with the original incorrect treatment of the calibration uncertainty and one with the correct treatment. We otherwise retain all configuration settings between the two runs. We then apply the reweighting algorithm to the former run to verify that the posteriors match the latter run.

Fig 4 shows the comparison between the original and rewighted posteriors from the run with incorrect treatment of the calibration uncertainty and the posterior from the run with the correct treatment for specifically the fractional calibration amplitude δA of the LIGO Hanford detector at the spline node located at 385 Hz. Examining the two directly sampled posteriors, we see that they have peaks away from $\delta A = 0$ with opposite sign. This neatly illustrates why the two other reweighting methods demonstrate such poor efficiency in this case the two posteriors do not overlap, so simple reweighting would not be able to transform one into the other. By transforming the samples first, we have moved them to approximately the correct position which can then be more effectively reweighted. Indeed, we see that the reweighted posterior appropriately matches with the posterior obtained from the analysis with correct calibrationuncertainty treatment.

Fig 5 shows similar posterior comparisons for the two non-calibration parameters: the luminosity distance $D_{\rm L}$ and the right ascension RA. We confirm the validity of our reweighting method here as the reweighted posterior is an accurate reflection of the posterior obtained when the correct calibration uncertainty treatment is applied. Whilst a similar pattern would hold for the other

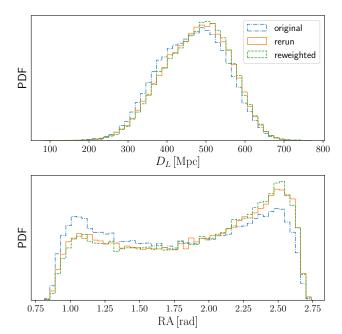


FIG. 5. Comparison of GW1509154 posteriors between an incorrect run, reweighted posteriors and the rerun with the correct calibration model. The reweighting procedure accurately transforms the shape of the posterior.

parameters, these are not shown as the differences between the parameters is insufficient to be visually noticeable. Similar investigations of few other events yield similar results—the reweighting process yields posteriors that match the corrected results.

Even though the posteriors are affected by the calibration error, the calibration envelopes of GW150914 are not significantly antisymmetric, as can be seen in Fig. 6. The envelopes become more antisymmetric at higher frequencies, but their means are always within one standard deviation of 0.

C. Effect on GWTC-3 MDR posteriors

Turning to the MDR analyses, as noted in this investigation the independent observations of A_{α} for each event are combined together to yield the final posterior. This can be written as:

$$p(A_{\alpha}|\boldsymbol{d}) \propto \pi(A_{\alpha})^{1-N} \prod_{i=1}^{N} p(A_{\alpha}|\boldsymbol{d_i}),$$
 (55)

where d is the combined data from each of the N individual observations, d_i , and $\pi(A_{\alpha})$ is the prior. As multiple posteriors are combined together, the error from the incorrect calibration uncertainty may compound, which could show differences in such a combined posterior even if they were not apparent in individual event posteriors.

Having verified in the previous subsection that the reweighting process produces posteriors representative of

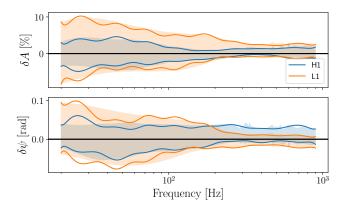


FIG. 6. Comparison of GW1509154 calibration uncertainty envelope on $\alpha=1/\eta_R$ used to construct the priors (shaded region) with the recovery of the recalibration factor after PE. Both areas are regions within one standard deviation of the mean. The calibration factor is split between the fractional amplitude deviation δA and the phase deviation $\delta \psi$.

those from analyses with the correct calibration treatment, Fig. 7 shows the original and reweighted combined A_{α} posteriors. The JSD between the two posteriors for each A_{α} exceeds the 0.0015 nat threshold for individual posteriors, with the largest distance being 0.023 nat for the $\alpha=0$ distribution. We also see a corresponding shift in the quantiles at which GR ($A_{\alpha}=0$) is recovered—with a maximal shift of 1.8% for the $\alpha=-1$ case.

However, the calibration uncertainty error is not the only uncertainty that accumulates when combining individual posteriors together. Errors in kernel density estimation (KDE) will also accumulate in the combined posterior. Notably, the decreases in effective sample size affects the bandwidth used for KDE according to Scott's rule:

$$b = \sigma n_{\text{eff}}^{-1/5},\tag{56}$$

where b is the bandwidth and σ is the standard deviation of the samples. The bandwidth strongly affects the estimates obtained from the KDE and its effect on the combined MDR posteriors was remarked on in Baka $et\ al.$ [49]. Considering this, we may conclude that the results from both the original and reweighted combination are, overall, consistent despite the increase in JSD.

ACKNOWLEDGEMENTS

T.B. is supported by the research program of the Netherlands Organisation for Scientific Research (NWO). C.H. thanks the UKRI Future Leaders Fellowship for support through the grant MR/T01881X/1. C. T. is supported by the Eric and Wendy Schmidt AI in Science Postdoctoral Fellowship, a Schmidt Sciences program. M.J.W. acknowledges support from STFC grants ST/X002225/1, ST/Y004876/1 and the University of

Portsmouth. J. V. acknowledges support from STFC grant ST/V005634/1. A.Z. was supported by NSF Grant PHY-2308833. L. S. is supported by the Australian Research Council Centre of Excellence for Gravitational Wave Discovery (OzGrav), Project Number CE230100016, and the Australian Research Council Discovery Early Career Researcher Award, Project Number DE240100206. This material is based upon work supported by the US National Science Foundation's (NSF's) LIGO Laboratory which is a major facility fully funded by the NSF, as well as the Science and Technology Facilities Council (STFC) of the United Kingdom, the Max-Planck-Society (MPS), and the State of Niedersachsen/Germany for support of the construction of Advanced LIGO and construction and operation of the GEO 600 detector. Additional support for Advanced LIGO was provided by the Australian Research Council. Virgo is funded, through the European Gravitational Observatory (EGO), by the French Centre National de Recherche Scientifique (CNRS), the Italian Istituto Nazionale di Fisica Nucleare (INFN) and the Dutch Nikhef, with contributions by institutions from Belgium, Germany, Greece, Hungary, Ireland, Japan, Monaco, Poland, Portugal, Spain. KAGRA is supported by Ministry of Education, Culture, Sports, Science and Technology (MEXT), Japan Society for the Promotion of Science (JSPS) in Japan; National Research Foundation (NRF) and Ministry of Science and ICT (MSIT) in Korea; Academia Sinica (AS) and National Science and Technology Council (NSTC) in Taiwan. The authors are grateful for computational resources provided by the LIGO Laboratory and supported by NSF Grants PHY-0757058 and PHY-0823459.

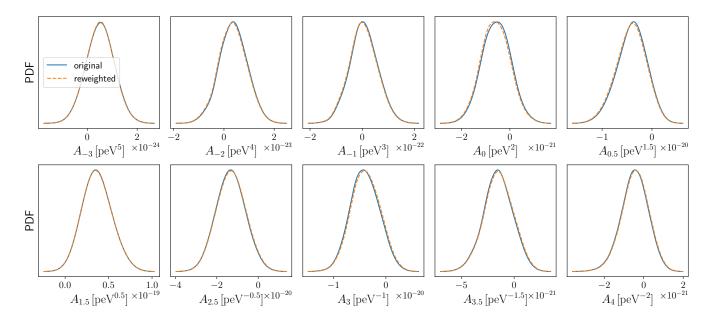


FIG. 7. Combined posteriors on the dispersion amplitude A_{α} for 43 signals used in MDR test in GWTC-3. The calibration issue causes only a negligible difference between the original and the reweighted posteriors.

- B. P. Abbott, R. Abbott, T. D. Abbott et al., Phys. Rev. Lett. 116, 241102 (2016), arXiv:1602.03840 [gr-qc].
- [2] The LIGO Scientific Collaboration, the Virgo Collaboration, the KAGRA Collaboration et al., arXiv e-prints, arXiv:2508.18080 (2025), arXiv:2508.18080 [gr-qc].
- [3] The LIGO Scientific Collaboration, the Virgo Collaboration, the KAGRA Collaboration et al., arXiv e-prints , arXiv:2508.18081 [gr-qc].
- [4] The LIGO Scientific Collaboration, the Virgo Collaboration, the KAGRA Collaboration et al., arXiv e-prints, arXiv:2508.18082 (2025), arXiv:2508.18082 [gr-qc].
- [5] T. L. S. Collaboration, the Virgo Collaboration and KA-GRA, Ligo/virgo/kagra public alerts, GraceDB (2025).
- [6] The LIGO Scientific Collaboration, the Virgo Collaboration and the KAGRA Collaboration, arXiv e-prints, arXiv:2508.18083 (2025), arXiv:2508.18083 [astro-ph.HE].
- [7] B. P. Abbott, R. Abbott, T. D. Abbott et al., Phys. Rev. Lett. 121, 161101 (2018), arXiv:1805.11581 [gr-qc].
- [8] The LIGO Scientific Collaboration, the Virgo Collaboration and the KAGRA Collaboration, arXiv e-prints, arXiv:2509.04348 (2025), arXiv:2509.04348 [astro-ph.CO].
- [9] R. Abbott, H. Abe, F. Acernese et al., arXiv e-prints, arXiv:2112.06861 (2021), arXiv:2112.06861 [gr-qc].
- [10] J. Aasi, B. P. Abbott, R. Abbott et al., Classical and Quantum Gravity 32, 074001 (2015), arXiv:1411.4547 [gr-qc].
- [11] F. Acernese, M. Agathos, K. Agatsuma et al., Classical and Quantum Gravity 32, 024001 (2015), arXiv:1408.3978 [gr-qc].

- [12] T. Akutsu, M. Ando, K. Arai et al., Nature Astronomy 3, 35 (2019), arXiv:1811.08079 [gr-qc].
- [13] B. P. Abbott, R. Abbott, T. D. Abbott et al., Phys. Rev. D 95, 062003 (2017), arXiv:1602.03845 [gr-qc].
- [14] C. Cahillane, J. Betzwieser, D. A. Brown et al., Phys. Rev. D 96, 102001 (2017), arXiv:1708.03023 [astro-ph.IM].
- [15] L. Sun, E. Goetz, J. S. Kissel et al., Classical and Quantum Gravity 37, 225008 (2020), arXiv:2005.02531 [astro-ph.IM].
- [16] L. Sun, E. Goetz, J. S. Kissel et al., arXiv e-prints , arXiv:2107.00129 (2021), arXiv:2107.00129 [astro-ph.IM].
- [17] E. Payne, C. Talbot, P. D. Lasky et al., Phys. Rev. D 102, 122004 (2020), arXiv:2009.10193 [astro-ph.IM].
- [18] S. Vitale, C.-J. Haster, L. Sun et al., Phys. Rev. D 103, 063016 (2021), arXiv:2009.10192 [gr-qc].
- [19] T. Akutsu, M. Ando, K. Arai et al., Progress of Theoretical and Experimental Physics 2021, 05A102 (2021), arXiv:2009.09305 [gr-qc].
- [20] T. Accadia, F. Acernese, M. Agathos et al., Classical and Quantum Gravity 31, 165013 (2014), arXiv:1401.6066 [gr-qc].
- [21] F. Acernese, T. Adams, K. Agatsuma et al., Classical and Quantum Gravity 35, 205004 (2018), arXiv:1807.03275 [gr-qc].
- [22] F. Acernese, M. Agathos, A. Ain et al., Classical and Quantum Gravity 39, 045006 (2022), arXiv:2107.03294 [gr-qc].
- [23] R. Abbott, T. D. Abbott, F. Acernese et al., Physical Review X 13, 041039 (2023), arXiv:2111.03606 [gr-qc].
- [24] B. P. Abbott, R. Abbott, T. D. Abbott et al., Phys. Rev. Lett. 118, 221101 (2017), arXiv:1706.01812

- [gr-qc].
- [25] L. S. Finn, Phys. Rev. D 46, 5236 (1992), arXiv:gr-qc/9209010 [gr-qc].
- [26] B. P. Abbott, R. Abbott, T. D. Abbott et al., Classical and Quantum Gravity 37, 055002 (2020), arXiv:1908.11170 [gr-qc].
- [27] R. Abbott, T. D. Abbott, F. Acernese et al., Phys. Rev. D 109, 022001 (2024), arXiv:2108.01045 [gr-qc].
- [28] W. Farr, B. Farr and T. Littenberg, Modelling calibration errors in CBC waveforms, Tech. Rep. DCC-T1400682 (LIGO, 2014).
- [29] L. S. Collaboration and V. Collaboration, 10.5281/zenodo.6513631 (2022).
- [30] L. S. Collaboration, V. Collaboration and K. Collaboration, 10.5281/zenodo.5546663 (2021).
- [31] G. Ashton, M. Hübner, P. D. Lasky et al., ApJS 241, 27 (2019), arXiv:1811.02042 [astro-ph.IM].
- [32] I. M. Romero-Shaw, C. Talbot, S. Biscoveanu et al., MN-RAS 499, 3295 (2020), arXiv:2006.00714 [astro-ph.IM].
- [33] G. Pratten, C. García-Quirós, M. Colleoni et al., Phys. Rev. D 103, 104056 (2021), arXiv:2004.06503 [gr-qc].
- [34] J. E. Thompson, E. Fauchon-Jones, S. Khan et al., Phys. Rev. D 101, 124059 (2020), arXiv:2002.08383 [gr-qc].
- [35] A. Matas, T. Dietrich, A. Buonanno et al., Phys. Rev. D 102, 043023 (2020), arXiv:2004.10001 [gr-qc].
- [36] T. Dietrich, S. Bernuzzi and W. Tichy, Phys. Rev. D 96, 121501 (2017), arXiv:1706.02969 [gr-qc].
- [37] T. Dietrich, S. Khan, R. Dudi et al., Phys. Rev. D 99, 024029 (2019), arXiv:1804.02235 [gr-qc].
- [38] R. J. E. Smith, G. Ashton, A. Vajpeyi et al., MNRAS 498, 4492 (2020), arXiv:1909.11873 [gr-qc].
- [39] J. S. Speagle, MNRAS 493, 3132 (2020), arXiv:1904.02180 [astro-ph.IM].

- [40] J. Veitch, V. Raymond, B. Farr et al., Phys. Rev. D 91, 042003 (2015), arXiv:1409.7215 [gr-qc].
- [41] C. Röver, R. Meyer and N. Christensen, Classical and Quantum Gravity 23, 4895 (2006), arXiv:gr-qc/0602067 [gr-qc].
- [42] M. van der Sluys, V. Raymond, I. Mandel et al., Classical and Quantum Gravity 25, 184011 (2008), arXiv:0805.1689 [gr-qc].
- [43] LIGO Scientific Collaboration, Virgo Collaboration and KAGRA Collaboration, LVK Algorithm Library - LAL-Suite, Free software (GPL) (2018).
- [44] S. Ossokine, A. Buonanno, S. Marsat et al., Phys. Rev. D 102, 044055 (2020), arXiv:2004.09442 [gr-qc].
- [45] C. Pankow, P. Brady, E. Ochsner et al., Phys. Rev. D 92, 023002 (2015), arXiv:1502.04370 [gr-qc].
- [46] J. Lange, R. O'Shaughnessy, M. Boyle et al., arXiv e-prints, arXiv:1705.09833 (2017), arXiv:1705.09833 [gr-qc].
- [47] D. Wysocki, R. O'Shaughnessy, J. Lange et al., Phys. Rev. D 99, 084026 (2019), arXiv:1902.04934 [astro-ph.IM].
- [48] B. P. Abbott, R. Abbott, T. D. Abbott et al., Physical Review X 9, 031040 (2019), arXiv:1811.12907 [astro-ph.HE].
- [49] T. Baka, B. Cirok, K. Haris et al., arXiv e-prints (2025), arXiv:2511.00497 [gr-qc].
- [50] S. Mirshekari, N. Yunes and C. M. Will, Phys. Rev. D 85, 024041 (2012), arXiv:1110.2720 [gr-qc].
- [51] N. Yunes, K. Yagi and F. Pretorius, Phys. Rev. D 94, 084002 (2016), arXiv:1603.08955 [gr-qc].
- [52] J. Lin, IEEE Transactions on Information Theory 37, 145 (1991).
- [53] S. Kullback and R. A. Leibler, The Annals of Mathematical Statistics 22, 79 (1951).
- [54] B. P. Abbott, R. Abbott, T. D. Abbott et al., Living Reviews in Relativity 23, 3 (2020).
- [55] B. P. Abbott, R. Abbott, T. D. Abbott et al., Physical Review X 6, 041015 (2016), arXiv:1606.04856 [gr-qc].