Cosmological inference using gravitational wave standard sirens: A mock data analysis

Rachel Gray®, ^{1,*} Ignacio Magaña Hernandez®, ^{2,†} Hong Qi®, ^{3,‡} Ankan Sur®, ^{4,5,§} Patrick R. Brady, ² Hsin-Yu Chen®, ⁶ Will M. Farr®, ^{7,8} Maya Fishbach®, ⁹ Jonathan R. Gair, ^{10,11} Archisman Ghosh®, ^{4,12,13,14} Daniel E. Holz®, ⁹ Simone Mastrogiovanni, ¹⁵ Christopher Messenger®, ¹ Danièle A. Steer®, ¹⁵ and John Veitch® ¹ ¹SUPA, University of Glasgow, Glasgow G12 800, United Kingdom ²University of Wisconsin-Milwaukee, Milwaukee, Wisconsin 53201, USA ³Cardiff University, Cardiff CF24 3AA, United Kingdom ⁴Nikhef, Science Park 105, 1098 XG Amsterdam, Netherlands ⁵Nicolaus Copernicus Astronomical Center, Polish Academy of Sciences, 00-716 Warsaw, Poland ⁶Black Hole Initiative, Harvard University, Cambridge, Massachusetts 02138, USA ⁷Department of Physics and Astronomy, Stony Brook University, Stony Brook, New York 11794, USA ³Center for Computational Astronomy, Flatiron Institute, New York, New York 10010, USA ⁹University of Chicago, Chicago, Illinois 60637, USA 10 School of Mathematics, University of Edinburgh, Edinburgh EH9 3FD, United Kingdom ¹¹Max Planck Institute for Gravitational Physics (Albert Einstein Institute), Potsdam-Golm 14476, Germany ¹²Delta Institute for Theoretical Physics, Science Park 904, 1090 GL Amsterdam, Netherlands ³Lorentz Institute, Leiden University, PO Box 9506, Leiden 2300 RA, Netherlands ¹⁴GRAPPA, University of Amsterdam, Science Park 904, 1098 XH Amsterdam, Netherlands ¹⁵Laboratoire Astroparticule et Cosmologie, CNRS, Université de Paris, 75013 Paris, France

(Received 4 October 2019; accepted 5 May 2020; published 8 June 2020)

The observation of binary neutron star merger GW170817, along with its optical counterpart, provided the first constraint on the Hubble constant H_0 using gravitational wave standard sirens. When no counterpart is identified, a galaxy catalog can be used to provide the necessary redshift information. However, the true host might not be contained in a catalog which is not complete out to the limit of gravitational-wave detectability. These electromagnetic and gravitational-wave selection effects must be accounted for. We describe and implement a method to estimate H_0 using both the counterpart and the galaxy catalog standard siren methods. We perform a series of mock data analyses using binary neutron star mergers to confirm our ability to recover an unbiased estimate of H_0 . Our simulations used a simplified universe with no redshift uncertainties or galaxy clustering, but with different magnitude-limited catalogs and assumed host galaxy properties, to test our treatment of both selection effects. We explore how the incompleteness of catalogs affects the final measurement of H_0 , as well as the effect of weighting each galaxy's likelihood of being a host by its luminosity. In our most realistic simulation, where the simulated catalog is about three times denser than the density of galaxies in the local universe, we find that a 4.4% measurement precision can be reached using galaxy catalogs with 50% completeness and ~250 binary neutron star detections with sensitivity similar to that of Advanced LIGO's second observing run.

DOI: 10.1103/PhysRevD.101.122001

I. INTRODUCTION

The idea that gravitational waves (GW) detections can be used for the inference of cosmological parameters, such as the Hubble constant (H_0), was first proposed over three decades ago by Bernard Schutz [1]. The key to this process

rachel.gray@ligo.org ignacio.magana@ligo.org hong.qi@ligo.org ankan.sur@ligo.org is that GW signals from compact binary coalescences (CBCs) act as standard sirens, in the sense that they provide a self-calibrated luminosity distance to the source. This can be obtained directly from the GW signal, and is therefore entirely independent of the cosmic distance ladder [2–10]. With the addition of redshift information for each source we then have the required input for cosmological inference.

At the time of writing, the current percent level state-of-the-art electromagnetic (EM) measurements of H_0 are in tension with each other. The Planck experiment

uses measurements of cosmic microwave background (CMB) anisotropies and provides a value of $H_0=67.4\pm0.5~{\rm km\,s^{-1}\,Mpc^{-1}}$ [11]. The supernovae, H_0 , for the equation of state of dark energy (SH0ES) experiment measures distances to Type Ia supernovae standard candles making use of the cosmic distance ladder, and gives $H_0=74.03\pm1.42~{\rm km\,s^{-1}\,Mpc^{-1}}$ [12]. These two independent measurements of H_0 are in tension at the level of $\sim 4.4-\sigma$ [12]. While the early-universe Planck measurements are also favored by measurements using supernovae calibrated with baryon acoustic oscillations [13], and the SH0ES results agree with local gravitational lensing measurements by the H0LiCOW Collaboration [14], calibration of supernovae using the Tip of the Red Giant Branch yields H_0 midway between the two [15].

This indicates the possibility that at least one of these measurements is subject to unknown systematics, or it could be an indication of new physics causing the discrepancy between the local measurements and the nonlocal (early universe) CMB based measurement. This makes a GW standard siren measurement of H_0 particularly interesting, as this will provide an alternative local constraint on H_0 . In this manner, the use of GWs as standard sirens may allow us to arbitrate the current situation, indicating either a bias in the current measurements, or pointing toward new physics.

The detection of the binary neutron star (BNS) event GW170817 [16], together with its optical counterpart [17,18] led to the first standard siren measurement of H_0 [19]. The counterpart associated with GW170817 allowed for the identification of its host galaxy, NGC4993, and hence a direct measurement of its redshift, which in turn resulted in the inferred value $H_0 = 70^{+12}_{-8} \text{ km s}^{-1} \text{ Mpc}^{-1}$. Future counterpart standard siren measurements are expected to constrain H_0 to the percent level [3–7].

Central to the aims of this paper is the case where an EM counterpart is not observed, and how H_0 inference can still be performed. In particular, the method proposed by Schutz in 1986 [1,20] allows the use of galaxy catalogs to provide redshift information for potential host galaxies within the event's GW sky-localization. The idea is that, by marginalizing over the possible discrete values of redshift for each GW detection we account for uncertainty as to which galaxy is the true host. By combining the information from many GW events, the contributions from the true host galaxies will grow since they will all share the same true H_0 . Contributions from the others will statistically average out, leading to a constraint on H_0 and possibly other cosmological parameters.

Over the course of the first observing run (O1) and the second observing run (O2) a total of 11 GW events were detected by the advanced LIGO and Virgo detectors: 10 are binary black hole (BBH) events and one is the abovementioned BNS event GW170817 [21]. The "galaxy catalog" method has been independently applied to both

the BNS event GW170817 (without assuming NGC4993 is the host) [22], and the BBH event GW170814 [23] resulting in posterior probability distributions on H_0 where the posterior from GW170814 was broader than (but consistent with) that obtained from GW170817. The difference in the widths of the H_0 constraints is an expected result due to the larger localization volume associated with GW170814, and the high number of galaxies it contained. Using the detections from O1 and O2, multiple GW events have been combined to give the latest standard siren measurement of H_0 [24] using the methodology presented in this paper.

Predictions suggest that it will be possible to constrain H_0 to less than 2% within 5 years of the start of the third observing run (O3) and to 1% within a decade, though this is dependent on the number of events observed with EM counterparts [6], and this may change as our understanding of astrophysical rates improves, and would require the detector amplitude calibration error to be measured to better than this precision. Simulations in [6] and [22], which assume complete catalogs based on realistic large-scale structure simulations, find that for BNSs without counterparts, the convergence is $40\%/\sqrt{N}$. The convergence found there for BBHs is much slower, as BBHs are typically detected at greater distances with larger localization volumes.

The prospects of identifying a transient EM counterpart will certainly increase, and correspondingly, the number of candidate host galaxies in a catalog will decrease, with improved event sky-localizations as future GW observatories join the detector network [25]. With the Japanese detector KAGRA [25] having joined O3 in early 2020, and LIGO-India approved for construction [26], the next decade of standard siren cosmology is set to be very exciting.

O3 began on April 1, 2019 and consists of 11 months' worth of data. The sensitivities of the LIGO and Virgo detectors have improved since O2, leading to an increased detection rate of GW candidates¹ [27]. This is the first observing run for which there will be 3 detectors operating for the entirety of the run. Having more detectors improves the duty-cycle of the network, i.e., the fraction of run time for which one or more detectors in the network is online, and also increases the rate of three-detector detections, which will likely be better localized on the sky than the two-detector ones. This is important, both in terms of performing EM follow-up for EM counterparts practically [28], and for reducing the number of possible host galaxies for events in the case where a counterpart is not observed.

This paper presents the Bayesian framework behind the GWCOSMO code, a product of the LIGO and Virgo

¹In the first half of O3 the detectors averaged the detection of one GW candidate per week. If all of these candidates are ultimately identified as real GW events, then O3 within its first two months will have exceeded the total number of detections of O1 and O2.

Collaborations (LVC) which was used to measure H_0 using detected GW events from O1 and O2 [24]. The method detailed in this paper is also expected to be implemented in future LIGO/Virgo/KAGRA standard siren measurements. We present results from a series of mock data analyses (MDAs) which were designed specifically to test this method's robustness against some of the most common pitfalls, in particular, GW selection effects which affect all H_0 measurements, and EM selection effects, which are relevant in the context of galaxy catalogs. This method builds upon the Bayesian framework first presented in [20] which has subsequently been extended, modified and independently derived by multiple authors [5,6,22,23,29]. The framework here is broadly equivalent to that in [6,22], however the mathematics and implementation differ, most notably in the treatment of EM selection effects. With specific care regarding selection effects we outline methods for constraining H_0 using both the "galaxy catalog" and "EM counterpart" approaches.

This paper is the first to explicitly test the robustness of a coded implementation of this methodology through use of galaxy catalogs which are incomplete and do not contain all of the GW host galaxies. Additionally, the GW data used in these MDAs were produced using an end-to-end simulation, including searching for "injected" signals in real detector data followed by a full parameter estimation to obtain the GW posterior samples [30,31], making this the most realistic set of simulated GW data to be used to explore GW cosmology to date. The analyses start with the most simplistic scenario, and increase in complexity with each iteration in order to ensure that the GWCOSMO code is able to pass each level satisfactorily before moving onto the next.

This paper is structured as follows. Section II presents the Bayesian framework used to estimate the posterior on H_0 . Section III discusses the design and preparation of the MDAs. In Sec. IV we present our results. We conclude in Sec. V giving a detailed discussion of results and providing guidance for future work. Some of the details of the Bayesian method have been set aside to be discussed in an Appendix.

II. METHODOLOGY

The late-time cosmological expansion in a Friedmann-Lemaître-Robertson-Walker universe is characterized by the Hubble-Lemaître parameter as a function of the redshift *z*,

$$H(z) = H_0 \sqrt{\Omega_{\rm m} (1+z)^3 + \Omega_{\rm k} (1+z)^2 + \Omega_{\Lambda}},$$
 (1)

where H_0 is the Hubble constant, the rate of expansion in the current epoch, and $\Omega_{\rm m}$ and Ω_{Λ} are the fractional matter density (including baryonic and cold dark matter) and fractional dark energy density (assumed to be due to a

cosmological constant) respectively; Ω_k is the fractional curvature energy density which is identically zero for a "flat" universe consistent with observations. Additionally, we have the constraint $\Omega_m + \Omega_k + \Omega_\Lambda = 1$ for all the components contributing to the energy density of universe at the present epoch.

The expansion history of the universe maps to a "red-shift-distance relation" associating the redshift z of observable sources to their luminosity distance $d_L(z)$ (see e.g., [32]) as,

$$d_L(z) = \frac{c(1+z)}{H_0} \int_0^z \frac{H_0}{H(z')} dz', \tag{2}$$

for a flat universe. From the relation between observed z and d_L to sources (EM sources such as variable stars or supernovae, or GW sources), one can measure the cosmological parameters appearing in H(z). With knowledge of the other cosmological parameters $\{\Omega_{\rm m},\Omega_{\rm k},\Omega_{\Lambda}\}$ coming from independent observations, the redshift-distance relation can be used to measure H_0 . We would like to note that with prior knowledge on the other cosmological parameters coming from EM observations, the measurement made with GW detections are not strictly independent measurements.

At low redshifts $z \ll 1$, the redshift-distance relation can be approximately described by the linear Hubble relation,

$$d_L(z) \approx cz/H_0,\tag{3}$$

which contains H_0 but is independent of the other cosmological parameters. With this approximate linear relation at low redshifts, any measurement of H_0 with GWs is independent of the values of the other cosmological parameters.

A. Standard sirens

The amplitude of the observed strain is inversely proportional to the luminosity distance to the GW source. For compact binary sources in quasicircular orbits, the two polarizations of the gravitational wave signal can be written to leading order as a function of frequency f as [33]

$$\tilde{h}_{+}(f) \propto \frac{\mathcal{M}_{z}^{5/6}}{2d_{L}} (1 + \cos^{2}(\iota)) f^{-7/6} \exp(i\phi(\mathcal{M}_{z}, f))$$
 (4)

$$\tilde{h}_{\times}(f) \propto \frac{\mathcal{M}_{z}^{5/6}}{d_{L}} \cos(i) f^{-7/6} \exp\left(i\phi(\mathcal{M}_{z}, f) + i\pi/2\right)$$
 (5)

where $\phi(\mathcal{M}_z,t)$ is the phase of the signal. The redshifting of the signal is accounted for by using the parameter $\mathcal{M}_z \equiv \mathcal{M}(1+z)$, the "redshifted chirp mass," to describe the signal as observed in the detector. Since \mathcal{M}_z appears in both the phase and the amplitude, and in practice is more strongly constrained by $\phi(\mathcal{M}_z,f)$, the dominant uncertainty on the signal amplitude results from the uncertainties

on luminosity distance, d_L and inclination angle ι . Each detector sees a linear combination of the two polarizations, $\tilde{h}(t) = F_+ \tilde{h}_+ + F_{\times} \tilde{h}_{\times}$, where $F_{+,\times}$ are the antenna response functions of the detector, which vary over the sky position and polarization angle of the source. Given multiple detectors at distant sites it is possible to simultaneously infer the parameters of the source, and therefore find a direct estimate of its luminosity distance [34]. This makes compact binaries self-calibrated luminosity distance indicators or "standard sirens" unlike EM distance indicators which need to undergo calibration via multiple rungs of the cosmic distance ladder. The redshift of the GW source, also required for cosmological inference, remains degenerate with the source's mass, contained within \mathcal{M}_{z} , and needs to be estimated in alternate ways. The precision of the d_L estimate is limited because of correlations with other parameters, particularly the inclination angle ι [5]. In this work we simulate these effects as part of our end-to-end analysis, described in Sec. III.

B. Galaxy information

There are multiple ways in which EM observations can provide complementary redshift² information. A BNS event may be detected in coincidence with an EM counterpart, which can be associated with the host galaxy to provide a direct measurement of the redshift of the source. More generically, a GW event may not have a detected EM counterpart, in which case one needs to fall back on the method outlined by Schutz [1] and use potential host galaxies within the event's sky localization region for the redshift information for the source. Two possibilities come up: (i) to use available galaxy catalogs, or (ii) to conduct dedicated EM follow-up on the event's sky region, mapping the galaxies within that area to as great a depth as possible to maximize the redshift information available.

When using galaxy catalogs to provide the prior redshift information, the possibility that the host galaxy lies beyond the reach of the catalog must be taken into account. EM telescopes are flux limited, which means that galaxy catalogs are inherently biased toward containing objects which are brighter and/or nearer-by (although there may be other selection effects due to galaxy color or size, depending on the catalog). These EM selection effects must be accounted for. Carrying out dedicated EM follow-up will, to some degree, mitigate this issue, as it will allow for far deeper coverage over a small section of the sky. For nearby events, the possibility that the host galaxy lies above the telescope's upper threshold may be negligible. However, the time and resources required for dedicated EM follow-up

means that the default approach for GW events observed without counterparts will be to use preexisting catalogs.

In either case, the uncertainty associated with each galaxy's redshift must be taken into account, including the redshift error due to the galaxy's peculiar velocity, v_n , and, in cases where the redshift is estimated photometrically, a much larger uncertainty due to the photometric algorithm. Peculiar velocities are significant for nearby galaxies. The effect of the peculiar velocity on the measurement of H_0 may be small if there are a large number of potential host galaxies in the GW event's sky-localization, but for a small number of galaxies, and for the counterpart case, this effect is particularly noticeable. For GW170817 at a nearby distance of about 40 Mpc, the peculiar velocity contribution was large as 10% of the total observed redshift [19], and different procedures of reconstructing the peculiar velocity field led to residual uncertainties on the redshift of between 2% and 8% [19,37–40]. The impact on H_0 measurement of peculiar velocities and their reconstruction is of topical interest, and has been the subject of several recent studies including [41-43]. Photometric redshifts on the other hand are important slightly farther away due to lack of spectroscopic data in galaxy catalogs. The "photoz" are estimated using fitting and machine learning algorithms [44,45], which often have large $\mathcal{O}(1)$ fractional uncertainties associated with them. While various caveats and subtleties for a realistic measurement have been outlined in [24], the impact of photo-z uncertainties on H_0 measurement is not precisely quantified in literature yet. Our present mock data analyses ignore these crucial redshift uncertainties altogether, and the impact of their magnitudes, profiles, and other systematic artefacts are left aside for possible future study.

C. Bayesian framework

This section presents an overview of the Bayesian framework of the GWCOSMO methodology. Parameters which appear explicitly in this overview are defined in Table I, while Table III in Appendix A 2 provides an extended list of parameter definitions, alongside a network diagram which demonstrates the conditional dependence of these parameters (see Fig. 9).

The posterior probability on H_0 from N_{det} GW events is computed as follows:

$$p(H_0|\{x_{\text{GW}}\}, \{D_{\text{GW}}\})$$

$$\propto p(H_0)p(N_{\text{det}}|H_0) \prod_{i}^{N_{\text{det}}} p(x_{\text{GW}i}|D_{\text{GW}i}, H_0) \quad (6)$$

where $\{x_{\rm GW}\}$ is the set of GW data, $D_{\rm GW}$ indicates that the event was detected as a GW and $p(H_0)$ is the prior on H_0 . For a given H_0 , the term $p(N_{\rm det}|H_0)$ is the probability of detecting $N_{\rm det}$ events. It depends on the intrinsic astrophysical rate of events in the source frame, $R = \frac{\partial N}{\partial V \partial T}$.

²There are ways of obtaining the redshift independent of EM observations, by using known population properties such as the mass distribution [10,35], or the neutron star equation-of-state [36].

TABLE I. A summary of the parameters present in the methodology.

Parameter	Definition
$\overline{H_0}$	The Hubble constant.
$N_{\rm det}$	The number of events detected during the observation period.
x_{GW}	The GW data associated with some GW source, s.
$D_{ m GW}$	Denotes that a GW signal was detected, i.e., that x_{GW} passed some detection statistic threshold ρ_{th} .
g	Denotes that a galaxy is (G) , or is not (\bar{G}) , contained within the galaxy catalog.
$x_{\rm EM}$	The EM data associated with some EM counterpart.
D_{EM}	Denotes that an EM counterpart was detected, i.e., that x_{EM} passed some threshold.

The total number of expected events is given by $N_{\rm det} = R\langle VT \rangle$, where $\langle VT \rangle$ is the average of the surveyed comoving volume multiplied by the observation time. By choosing a scalefree prior on rate, $p(R) \propto 1/R$, the dependence on H_0 drops out [46]. For simplicity this approximation is made throughout the analysis and therefore $p(N_{\rm det}|H_0)$ is absent from further expressions.

The remaining term factorizes into likelihoods for each detected event. Using Bayes' theorem we can write it as,

$$p(x_{\text{GW}}|D_{\text{GW}}, H_0) = \frac{p(D_{\text{GW}}|x_{\text{GW}}, H_0)p(x_{\text{GW}}|H_0)}{p(D_{\text{GW}}|H_0)}$$
$$= \frac{p(x_{\text{GW}}|H_0)}{p(D_{\text{GW}}|H_0)}, \tag{7}$$

where we set $p(D_{\rm GW}|x_{\rm GW},H_0)=1$, since the analysis is only carried out when the signal-to-noise ratio (SNR), ρ , associated with $x_{\rm GW}$ passes some detection statistic threshold $\rho_{\rm th}$ —it is a prerequisite that the event has been detected. Calculating $p(D_{\rm GW}|H_0)$ requires integrating over all possible realizations of GW events, with a lower integration limit of $\rho_{\rm th}$:

$$p(D_{\rm GW}|H_0) = \int_{\rho > \rho_{\rm th}}^{\infty} p(x_{\rm GW}|H_0) dx_{\rm GW}.$$
 (8)

For explicit details on the calculation of $p(D_{\rm GW}|H_0)$ see Appendix A 5. The term $p(D_{\rm GW}|H_0)$ depends on properties of the GW source population (e.g., the mass distribution), but in this work, for simplicity, it is assumed that the population properties are known exactly.

1. The galaxy catalog method

In the galaxy catalog case, the EM information enters the analysis as a prior, made up of a series of possibly smoothened delta functions³ at the redshift, right ascension (RA) and declination (dec) of the possible source locations. As we are in the regime where (especially for BBHs) galaxy catalogs cannot be considered complete out to the distances to which GW events are detectable, we have to consider the possibility that the host galaxy is not contained within the galaxy catalog due to being dimmer than the apparent magnitude threshold. In order to do so, we marginalize the likelihood over the case where the host galaxy is, and is not, in the catalog (denoted by G and \bar{G} respectively):

$$\begin{split} p(x_{\rm GW}|D_{\rm GW},H_0) &= \sum_{g=G,\bar{G}} p(x_{\rm GW}|g,D_{\rm GW},H_0) p(g|D_{\rm GW},H_0), \\ &= p(x_{\rm GW}|G,D_{\rm GW},H_0) p(G|D_{\rm GW},H_0) \\ &+ p(x_{\rm GW}|\bar{G},D_{\rm GW},H_0) p(\bar{G}|D_{\rm GW},H_0). \end{split}$$

While theoretically equivalent to and consistent with the methodology presented in [6,22], the mathematics and implementation here differ, most notably in the treatment of EM selection effects, and our focus on whether the host galaxy is contained within the galaxy catalog or not, rather than calculating a "completeness fraction" in order to weight the in-catalog and out-of-catalog likelihood contributions. This, alongside the modeling of EM selection effects using an apparent magnitude threshold, which has not been done before, accounts for the main differences between this derivation and those presented in earlier works. The methodology presented here aligns directly with the implementation of the GWCOSMO code. We leave the details of this derivation to Appendix A 2.

2. The counterpart method

The method outlined above is for the galaxy catalog case, in which no EM counterpart is observed, or expected. We also consider the case where we observe an EM counterpart. The main difference is the inclusion of a likelihood term for the EM counterpart data, mirroring that of the GW data.

The likelihood in this case, which is the term within the product in Eq. (6), is given by:

$$p(x_{\text{GW}}, x_{\text{EM}} | D_{\text{GW}}, D_{\text{EM}}, H_0)$$

$$= \frac{p(x_{\text{GW}}, x_{\text{EM}} | H_0) p(D_{\text{GW}}, D_{\text{EM}} | x_{\text{GW}}, x_{\text{EM}}, H_0)}{p(D_{\text{GW}}, D_{\text{EM}} | H_0)},$$

$$= \frac{p(x_{\text{GW}} | H_0) p(x_{\text{EM}} | H_0)}{p(D_{\text{EM}} | D_{\text{GW}}, H_0) p(D_{\text{GW}} | H_0)},$$
(10)

³While uncertainties on the galaxy sky-coordinates can be safely ignored, the error on the redshift can be modeled with a Gaussian or a more complicated distribution.

where $x_{\rm EM}$ refers to the EM counterpart data and $D_{\rm EM}$ denotes that the counterpart was detected. In the numerator we have assumed that the GW and EM data are independent of each other and so the joint GW-EM likelihood factors out. $p(D_{\rm GW}, D_{\rm EM}|x_{\rm GW}, x_{\rm EM}, H_0)$ is further factorized as $p(D_{\rm EM}|D_{\rm GW}, x_{\rm GW}, x_{\rm EM}, H_0)p(D_{\rm GW}|x_{\rm GW}, x_{\rm EM}, H_0)$. The first term is equal to 1, as this method is only used when we have observed an EM counterpart, meaning that by definition $x_{\rm EM}$ has passed some threshold for detectability set by EM telescopes. The second term also goes to 1, due to the same threshold argument as in Sec. II C.

For simplicity, in this paper we make the assumption that the detection of an EM counterpart is flux-limited and, as in [19], that the detectability of EM counterparts extends well beyond the distance to which BNSs are detectable with O2-like LIGO and Virgo sensitivity. Following this, we make the assumption that the term $p(D_{\rm EM}|D_{\rm GW},H_0)\approx 1$, and leave a more rigorous analysis of the H_0 -dependence of this term for a future study.

In an ideal scenario, the observation of an EM counterpart will allow for the identification of one of the galaxies in the neighboring region as the host of the GW event. In the case where the EM counterpart cannot be unambiguously linked to a host galaxy, this uncertainty can also be taken into account. See Appendix A 4 for more details.

III. THE MOCK DATA ANALYSES

In this section we describe a series of mock data analyses (MDAs) that we use to test our implementation of the Bayesian formalism described in Sec. II and its ability to infer the posterior on H_0 under different conditions. For each case, the MDA consists of (i) simulated GW data, and (ii) a corresponding mock galaxy catalog. In all cases, we make several idealized assumptions regarding both the GW and galaxy data. On the GW side, the detection efficiency and the source population properties are assumed to be known exactly. On the galaxy side, the luminosity function and magnitude limit are also assumed to be known exactly in each case, so that the incompleteness correction can be calculated exactly. Further, we neglect the effects of large-scale structure and redshift uncertainties in the mock catalogs.

For each of the MDAs we use an identical set of simulated BNS events from the First Two Years of electromagnetic follow-up with Advanced LIGO and Virgo dataset [30,31].⁴ The set of BNS events comes from an end-to-end simulation of approximately 50,000 "injected"

events in detector noise corresponding to a sensitivity similar to what was achieved during O2. Only a subset (approximately 500 events) were "detected" by a network of two or three detectors with the GstLAL matched filter based detection pipeline [47]. From the above detections, 249 events were randomly selected (in a way that no selection bias was introduced), and these events underwent full Bayesian parameter estimation using the LALInference software library [34] to obtain gravitational wave posterior samples and skymaps. Consistency with the First Two Years parameter estimation results in terms of sky localization areas and 3D volumes was demonstrated in [48]. It is these 249 events of the First Two Years dataset and the associated GW data which we use for our analysis.

The galaxy catalogs for each iteration of the MDA described below are designed to test a new part of the GWCOSMO methodology in a cumulative fashion, starting with GW selection effects, adding in EM selection effects, and finally testing the ability to utilize the information available in the observed brightness of host galaxies, by weighting the galaxies with a function of their intrinsic luminosities.

The starting point for the galaxy catalogs is to take all 50,000 injected events from the First Two Years dataset and simulate a mock universe, which contain a galaxy corresponding to each injected event's sky location and luminosity distance, where the latter is converted to a redshift using a fiducial "simulated" H_0 value of 70 km s⁻¹ Mpc⁻¹. The First Two Years data was originally simulated in a universe where GW events followed a d_L^2 distribution, and there was no distinction between the source frame and the (redshifted) detector frame masses. Though not ideal, this data reasonably mimics a low redshift universe $(z \ll 1)$ in which the linear Hubble relation of Eq. (3) holds, and galaxies follow a z^2 distribution. We use the same linear relation for the generation of the MDA universe (i.e., a set of simulated galaxy catalog parameters) for each of the MDAs. It should be emphasized that the Bayesian method for estimation of H_0 outlined in Sec. II above is general, and can be extended to realistic scenarios with a nonlinear cosmology with $\{\Omega_{\rm m},\Omega_{\rm k},\Omega_{\Lambda}\}$ held fixed. So, in particular, the method is applicable for events which are detected at higher distances, where the low redshift approximation breaks down. The restriction to a linear cosmology in this paper comes only due to the use of the MDA dataset. We would like to note that by using a linear cosmology, we are not testing possible effects introduced by the presence of other cosmological parameters. The analysis at large redshifts may, for example, be sensitive to the values (or the assumed prior ranges) of the parameters like Ω_m and Ω_{Λ} .

The first four columns of Table II summarize the characteristics of each of the galaxy catalogs created and how they correspond to each MDA. We give a brief description for each of the cases below.

⁴The set of simulations in [31] are more realistic with the same injections in (recolored) detector data as opposed to Gaussian noise used in [30]. Correspondingly, the detection criterion is in terms of a false alarm rate (FAR) rather than a threshold on the SNR. This is an important distinction, particularly affecting events marginally close to the detection threshold. We use the simplified set of simulations in [30] noting potential caveats.

TABLE II. A summary of the main results. We quote the peak value and the 68.3% highest density error region for the posterior probability on H_0 for each of the MDAs combining all 249 events. The fractional uncertainty is defined as the half-width of the 68.3% highest density probability interval divided by the simulated value of $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$.

MDA	Host galaxy preference	Completeness ^a	$m_{ m th}$	Analysis assumption	$H_0 ({\rm km s^{-1} Mpc^{-1}})$	Fractional uncertainty
0	Known host			Direct counterpart	$69.08^{+0.79}_{-0.80}$	1.13%
1	Equal weights	100%		Unweighted catalog	$68.91^{+1.36}_{-1.22}$	1.84%
2a	Equal weights	75%	19.5	Unweighted catalog	$69.97^{+1.59}_{-1.50}$	2.21%
2b	Equal weights	50%	18	Unweighted catalog	$70.14^{+1.80}_{-1.67}$	2.48%
2c	Equal weights	25%	16	Unweighted catalog	$70.14^{+2.29}_{-2.18}$	3.20%
3a	Luminosity weighted	50%	14	Weighted catalog	$70.83^{+3.55}_{-2.72}$	4.48%
3b	Luminosity weighted	50%	14	Unweighted catalog	$69.50_{-3.24}^{+4.20}$	5.31%

^aThe completeness is calculated as a number completeness using Eq. (12) for MDAs 1 and 2, and as a luminosity completeness using Eq. (15) for MDA 3, out to a fiducial distance of 115 Mpc, such that it is indicative of the fraction of host galaxies which are inside the galaxy catalog in both cases.

A. MDA0: Known associated host galaxies

MDA0 is the simplest version of the MDAs, in which we identify with certainty the host galaxy for each GW event, and is equivalent to the direct counterpart case. As the galaxies are generated with no redshift uncertainties or peculiar velocities, the results will be (very) optimistic. This MDA provides the "best possible" constraint on H_0 using the 249 events, which then allows for comparison with the other MDAs.

B. MDA1: Complete galaxy catalog

The MDA1 universe consists of the full set of 50,000 galaxies out to $z \approx 0.1$ ($d_L \approx 428$ Mpc) in the original First Two Years dataset. This gives a galaxy number density of \sim 1 per 7000 Mpc³, which is \sim 35 times sparse compared to the actual density of galaxies in the local universe [49]. Additional galaxies are generated beyond the edge of the dataset universe, uniformly across the sky and uniformly in comoving volume, thereby extending the universe out to a radius of 2000 Mpc (z = 0.467 for $H_0 = 70$ km s⁻¹ Mpc⁻¹). This means that, even allowing H_0 to be as large as 200 km s⁻¹ Mpc⁻¹, the edge of the MDA universe is more than twice the highest redshift associated with the farthest detection (which is at \sim 270 Mpc). Each of the 249 detected BNS have a unique associated host galaxy contained within the MDA1 catalog. This catalog is thus *complete* in the sense that it contains every galaxy in the simulated universe. We refer to the MDA universe as MDA1 throughout the rest of the paper, and similarly for the subsequent MDAs.

MDA1 is designed to test our treatment of GW selection effects, by ensuring that given a set of sources and access to a complete catalog, our methodology and analysis produces a result consistent with the simulated value of H_0 .

C. MDA2: Incomplete galaxy catalog

MDA2 is designed to test our treatment of EM selection effects, by applying an apparent magnitude threshold to the MDA universe, such that a certain fraction of the host galaxies is not contained in it. This is a necessary consideration, given that we are in the regime where GW signals are being detected beyond the distance to which the current galaxy catalogs can be considered to be complete. This has been true for BBHs detections since O1, and is true of BNSs as well in O3.

In order to create the catalog for MDA2, we start with the initial MDA1 universe and assign luminosities to each of the galaxies within it. We assume that the luminosity distribution of the galaxy catalog is known to the observer throughout and follows a Schechter function of the form [50]

$$\phi(L)dL = n^* \left(\frac{L}{L^*}\right)^{\alpha} e^{-L/L^*} \frac{dL}{L^*},\tag{11}$$

where L denotes a given galaxy luminosity and $\phi(L)dL$ is the number of galaxies within the luminosity interval [L, L + dL]. The characteristic galaxy luminosity is given by $L^* = 1.2 \times 10^{10} h^{-2} L_{\odot}$ with solar luminosity $L_{\odot} = 3.828 \times 10^{26}$ W, and $h \equiv H_0/(100 \text{ km s}^{-1} \text{ Mpc}^{-1})$, $\alpha = -1.07$ characterizes the exponential drop off of the luminosity function, and n^* denotes the number density of objects in the MDA universe (in practice, this only acts as a normalization constant). The integral of the Schechter function diverges at $L \to 0$, requiring a lower luminosity cutoff for the dimmest galaxies in the universe which we set to $L_{\text{lower}} = 0.001L^*$. This choice is arbitrary for our purpose here, but small enough to include almost all objects classified as galaxies in real catalogs like GLADE [49].

⁵For MDA1 and for all subsequent MDAs, it has been tested that the artificial "edge of the universe" has no bearing on the results.

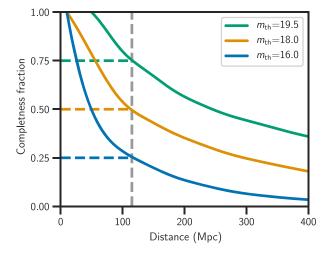
 $^{^6}$ We note that the parameter L^* of the Schechter luminosity function itself depends on H_0 , which we allow to vary and hence take into account within our formalism.

These luminosities are then converted to apparent magnitudes using $m \equiv 25-2.5\log_{10}(L/L^*)+5\log_{10}(d_L/\mathrm{Mpc})$, and an apparent magnitude threshold m_{th} is applied as a crude characterization of the selection function of an optical telescope observing only objects with $m < m_{\mathrm{th}}$. MDA2 is broken into three sub-MDAs, in order to test our ability to handle different levels of galaxy catalog completeness dictated by different telescope sensitivity thresholds. In each case, the catalog completeness is defined as the ratio of the number of galaxies inside the Catalog relative to the number of galaxies inside the MDA universe, out to a reference fiducial distance d_L ,

$$f_{\text{completeness}}(d_L) = \frac{\sum_{j}^{\text{MDA2}} \Theta(d_L - d_{L_j})}{\sum_{k}^{\text{MDA1}} \Theta(d_L - d_{L_k})}, \quad (12)$$

where the numerator is a sum over the galaxies contained within the MDA2 catalog out to some reference distance d_L , and the denominator is a sum over the galaxies in the MDA1 catalog.

Apparent magnitude thresholds of $m_{\rm th}=19.5$, 18, and 16 are chosen for the three sub-MDAs, which correspond to cumulative number completeness fractions of 75%, 50% and 25% respectively, evaluated at a distance of $d_L=115$ Mpc, chosen such that given the luminosity distance distribution of detected BNSs, the completeness fraction for the sub-MDA to this distance is roughly indicative of the percentage of host galaxies which remain inside the galaxy catalog. The left panel of Fig. 1 shows how the completeness of each of the MDA2 catalogs drop off as a function of distance.



D. MDA3: Luminosity weighting

MDA3 is designed to test the effect of weighting the likelihood of any galaxy being host to a GW event as a function of their luminosity. It is probable that the more luminous galaxies are also more likely hosts for compact binary mergers; the luminosity in blue (B-band) is indicative of a galaxy's star formation rate, for example, while the luminosity in high infrared (K-band) is a tracer of the stellar mass [51–53]. The bulk of the host probability is expected to be contained within a smaller number of brighter galaxies, effectively reducing the number of galaxies which need to be considered. Additional information from luminosity is thus expected to improve the constraint on H_0 by narrowing its posterior probability density.

For MDA3, the probability of a galaxy hosting a GW event is chosen to be proportional to the galaxy's luminosity. Because the GW events for these MDAs were generated in advance, and we are retroactively simulating the universe in which they exist, generating the MDA3 universe required some care: luminosities have to be assigned to the host galaxies *and* the nonhost galaxies in such a way that our choice of simulated luminosity weighting is correctly represented within the galaxy catalog.

As with MDA2, the luminosity distribution of the galaxies in the universe is assumed to follow the Schechter luminosity function as in Eq. (11) (referred to from now on as p(L)). However, the joint probability of a single galaxy having luminosity L and hosting a GW event (which emits a signal, s) is $p(L, s) \propto Lp(L)$, where we assume that the probability of a galaxy of luminosity L hosting a source is proportional to the luminosity itself. All host galaxies thus have luminosities sampled from Lp(L).

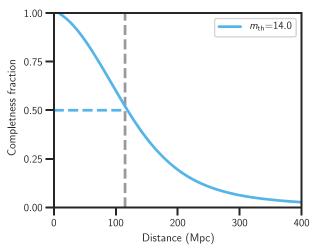


FIG. 1. Galaxy catalog completeness fractions for MDA2 and MDA3. Left panel: Galaxy number completeness fraction defined in Eq. (12) as a function of luminosity distance for the three MDA2 sub-catalogs. The lines in green, orange and blue correspond to the catalogs with $m_{\rm th} = 19.5$, 18, and 16 respectively; these correspond to completeness fractions of 75%, 50% and 25% out to a fiducial reference distance of 115 Mpc (shown as a vertical grey line). Right panel: The galaxy luminosity completeness fraction defined in Eq. (15) as a function of luminosity distance for the MDA3 catalog, with $m_{\rm th} = 14$. At the reference distance of 115 Mpc (vertical grey line), this is corresponds to a completeness fraction of ~50%.

In this context, we must consider all galaxies which hosted GW events, not just those from which a signal was detected. With this in mind, the overall luminosity distribution has the following form:

$$p(L) = \beta \frac{L}{\langle L \rangle} p(L) + (1 - \beta) x(L)$$
 (13)

where β is the fraction of host galaxies to total galaxies for the observed time period $(1 \ge \beta \ge 0)$, $L/\langle L \rangle$ is the normalized luminosity, and x(L) is the unknown luminosity distribution of galaxies which did not host GW events, which we can sample for a given value of β .

Rearranging to obtain the only unknown, x(L), gives

$$x(L) = \frac{p(L)}{1 - \beta} \left[1 - \beta \frac{L}{\langle L \rangle} \right],\tag{14}$$

and from this we see there is an additional constraint on β , because the term inside the brackets must be >0. The maximum value that β can take is given by $\beta_{\rm max} = \langle L \rangle / L_{\rm max}$, where $L_{\rm max}$ is the maximum luminosity from the Schechter function, and $\langle L \rangle$ is the mean. From the Schechter function parameters detailed in Sec. III C, $\beta_{\rm max} \approx 0.015$.

The original First Two Years data was generated by simulating ~50, 000 BNS events, of which ~500 were detected, of which 249 randomly selected detections underwent parameter estimation. The number of "hosting" and "nonhosting" galaxies have to be rescaled to represent this. Thus half of the original galaxies were denoted as hosts (including those associated with the 249 detected GW events). However, in order to satisfy the requirements for β , a greater density of nonhosting galaxies had to be added to the universe before luminosities could be assigned. Thus for MDA3, the density of galaxies is increased by a factor of 100, with the acknowledgement that this would lead to a broadening of the final posterior. MDA3 has a galaxy density of ~ 1 galaxy per 70 Mpc³, which is about 3 times denser than the actual density of galaxies in the local universe [49]. This also means that MDA3 is not directly comparable with the previous MDA versions, save MDA0. The galaxies which are hosts are assigned luminosities from Lp(L), and nonhosts from x(L) above.

In order to include EM selection effects, an apparent magnitude cut $m_{\rm th}$ of 14 is applied, such that the completeness of the galaxy catalog is ~50% out to the same fiducial distance of 115 Mpc as in MDA2. In this case, completeness is however defined in terms of the fractional luminosity contained in the catalog, rather than in terms of numbers of objects:

$$f_{\text{completeness}}(d_L) = \frac{\sum_{j}^{\text{MDA3}} L_j \Theta(d_L - d_{L_j})}{\sum_{k}^{\text{complete}} L_k \Theta(d_L - d_{L_k})}, \quad (15)$$

where the numerator is summed over the galaxies inside the MDA3 apparent magnitude-limited catalog, and the denominator is summed over the galaxies in the whole MDA3 universe. This is shown in the right panel of Fig. 1. As the host galaxies are luminosity weighted, the cumulative luminosity completeness is representative of the percentage of BNS event hosts inside the catalog.

IV. RESULTS

In this section we summarize the results for the mock data analyses described in Sec. III. We show the combined posteriors on H_0 for each MDA, discuss the convergence to the simulated value of $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$ and calculate the precision of the combined measurement under each set of conditions. In Table II we list the measured values of the Hubble constant for the combined 249 event posterior (maximum *a posteriori* and 68.3% highest density posterior intervals) all computed with a uniform prior on H_0 in the range of [20, 200] km s⁻¹ Mpc⁻¹, as well as the corresponding fractional uncertainties for each of the MDAs.

A. MDA0: Known associated host galaxies

We first consider the simple case where we identify the true host galaxy for every event and determine the resulting 249-event combined H_0 posterior. Figure 2 presents the results of this analysis. The likelihoods for each individual GW event are shown (normalized relative to each other but scaled with respect to the combined posterior for clarity) shaded by the event's optimal SNR in the detector network, as defined in [54]. In this case, each likelihood is informative, having a clearly-defined peak corresponding to finding the likely values of H_0 for the known galaxy redshift. Each curve traces the information in the corresponding d_L distribution, which is usually unimodal, but in some cases may have two or more peaks [30,31]. We see that the peaks of the individual likelihoods do not necessarily correspond to the true value $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$, but there is always support for it, leading to the combined posterior, which is overlaid in thick purple. This gives us a statistical estimate for the maximum a posteriori value and 68.3% maximum-density credible interval for H_0 as $69.08^{+0.79}_{-0.80}~{\rm km\,s^{-1}\,Mpc^{-1}}$. The final result combining all the 249 events have converged to a relatively symmetric "Gaussian" distribution [55].

The result of MDA0 provides us with the best possible H_0 estimate given the set of GW detections, since this case corresponds to perfect knowledge of the host galaxies. This gives us a benchmark against which other versions of the MDA can be compared. Since this is a best-case scenario, we have the least statistical uncertainty in the final result, making any systematic bias more apparent than for the subsequent MDAs. For the combined result with 249

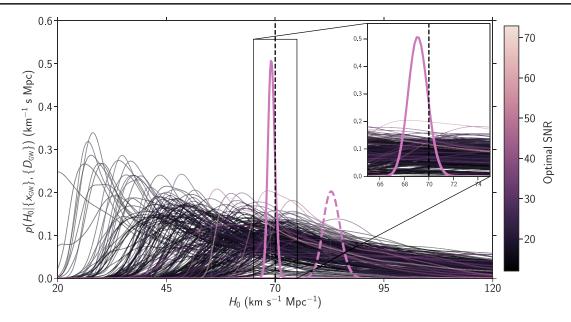


FIG. 2. Individual and combined results for MDA0 (known host galaxy or direct counterpart case). The solid thick purple line shows the combined posterior probability density on H_0 , while the dashed line shows the combined posterior when GW selection effects are neglected. Individual likelihoods (normalized and then scaled by an arbitrary value), for each of the 249 events, are shown as thin lines with shades corresponding to their optimal SNR. The simulated value of H_0 is shown as a vertical dashed line.

events, the simulated value is contained within the support of the posterior distribution of H_0 .

MDA0 demonstrates the importance of correctly accounting for GW selection effects. We are biased toward detecting sources which are nearer-by, and which are optimally orientated (closer to face-on). If an analysis is performed without taking into consideration the denominator $p(D_{GW}|H_0)$ of Eq. (7), which corrects for this, the posterior density on H_0 converges to a value different from its simulated value of $70 \text{ km s}^{-1} \text{ Mpc}^{-1}$. This can be seen in Fig. 2, where the dashed purple line shows the MDA0 combined posterior for all 249 events, neglecting GW selection effects entirely. We leave a detailed exploration of what level of accuracy in the GW selection function is required in order to move beyond 249 BNS-with-counterpart events, and simply note that in this case, it is sufficient enough that any biases which could affect the next stages of the MDA do not arise from the GW selection effects.

B. MDA1: Complete galaxy catalog

The next more complex case is MDA1, where we assume no counterpart was observed, and resort to using a galaxy catalog. MDA1 uses a *complete* galaxy catalog containing all potential hosts—an optimistic scenario, in which EM selection effects do not need to be considered. The results with MDA1 already show a wider posterior distribution on H_0 (68.91 $^{+1.36}_{-1.22}$ km s⁻¹ Mpc⁻¹) because of lack of certainty of the host galaxy (Fig. 3). The introduction of this uncertainty means that the contributions from each event will be smoothed out, depending on the size of the event's

sky localization and the number of galaxies within it. As can be seen in Fig. 4, there is a far higher proportion of events for which the likelihood is relatively broad and less informative, in comparison to Fig. 2. However, many events clearly have a small number of galaxies in their skyarea, and hence still show clear peaks.

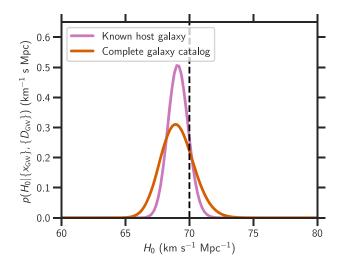


FIG. 3. Comparison of the galaxy catalog method with the known host galaxy case. Joint posterior probability density on H_0 using all 249 events for MDA0 (known host galaxy) and MDA1 (complete galaxy catalog) are shown respectively in purple and red. For this set of simulations, uncertainty with the galaxy catalog is only about 1.63 times larger than with known host galaxies.

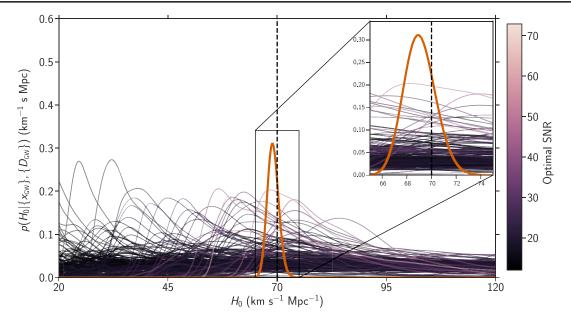


FIG. 4. Individual and combined results for MDA1 (complete galaxy catalog). The thick red line shows the combined posterior probability density on H_0 . Individual likelihoods (normalized and then scaled by an arbitrary value), for each of the 249 events, are shown as thin lines with shades corresponding to their optimal SNR. The simulated value of H_0 is shown as a vertical dashed line. Many of the individual likelihoods do not have sharp features, however the final result converges to the simulated value with redshift information present in the galaxy catalogs. This demonstrates the applicability of our methodology.

C. MDA2: Incomplete galaxy catalog

The next most complex scenario is the case where we have incomplete galaxy catalogs, limited by an apparent magnitude threshold. This gives us the first case where accounting for EM selection effects is important. To investigate this, we consider three galaxy catalogs, with apparent magnitude thresholds of $m_{\rm th}=19.5, 18$ and 16, with respective completeness fractions of 75%, 50% and 25% in addition to the complete catalog for MDA1 (see Sec. III C for details). The combined 249-event posterior distributions on H_0 are shown in Fig. 5.

As the catalogs become less complete, the combined H_0 posterior becomes wider. This is because the probability that the host galaxy is inside the catalog decreases. The contribution from the galaxies within the catalog is reduced, and the uninformative contribution from the out-of-catalog term in Eq. (9) increases. This is visible in the individual likelihoods shown in Fig. 6, where instead of decreasing toward zero at high values of H_0 , many of the individual likelihoods tend toward a constant. This is because, in the absence of EM data, and with the linear Hubble relation assumed in this work, the number of unobserved galaxies increases without limit as d_L^2 . This is seen mostly for events at high distances (where the host has a lower probability of being in the catalog), or for welllocalized events where there is no catalog support at the relevant redshifts within the event's sky area. However, enough events are detected at low distances, where the catalogs are more complete and so provide informative redshift information, to produce an upper constraint on H_0 . We estimate $H_0 = 69.97^{+1.59}_{-1.50}$, $70.14^{+1.80}_{-1.67}$, and $70.14^{+2.29}_{-2.18}$ km s⁻¹ Mpc⁻¹ respectively for galaxy catalogs of 75%, 50%, and 25% completeness. See Sec. IV E for a more in depth comparison of how galaxy catalog completeness affects posterior width.

Our exercise demonstrates that we need to know (or assess) the completeness of galaxy catalogs, and put in an

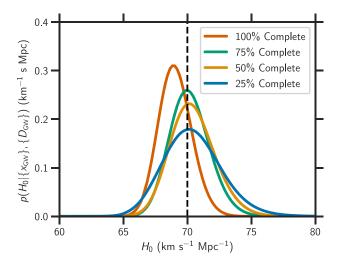


FIG. 5. Comparison of results with varying galaxy catalog completeness. In MDA2, the simulated apparent magnitude threshold is varied to obtain galaxy catalogs of 100%, 75%, 50%, and 25% completeness. The corresponding posterior probability densities on H_0 using all 249 events are shown in red, green, yellow, and blue respectively.

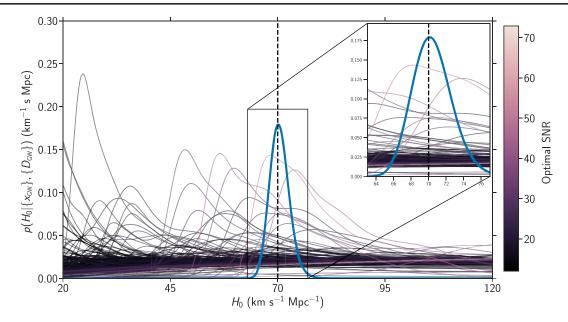


FIG. 6. Individual and combined results for MDA2 with a 25% complete galaxy catalog. The thick blue line shows the combined posterior probability density on H_0 . Individual likelihoods (normalized and then scaled by an arbitrary value), for each of the 249 events, are shown as thin lines with shades corresponding to their optimal SNR. The simulated value of H_0 is shown as a vertical dashed line. Compared to MDA0 (Fig. 2) and MDA1 (Fig. 4), fewer individual likelihoods are peaked here. Although the final H_0 estimate is less precise, the results converge to the simulated value, demonstrating the applicability of our methodology to threshold-limited galaxy catalogs of about 25% completeness.

appropriate out-of-catalog term in the analysis. For any of the MDA2 catalogs, if we assume that the galaxy catalog is complete, when in reality it is not, we get a posterior distribution on H_0 which is inconsistent with a value of $70 \text{ km s}^{-1} \text{ Mpc}^{-1}$. This is because the well-localized events for which the host is not inside the catalog do not have support for the correct value of H_0 . In real catalogs, galaxy clustering might ensure that there are nearby bright galaxies in the catalog, partially mitigating this bias.

D. MDA3: Luminosity weighting

Until now we have considered all galaxies in our catalog to be equally likely to host a gravitational-wave source. In MDA3 we analyze the case described in Sec. III D where this is no longer true by constructing a galaxy catalog such that the probability of any single galaxy hosting a GW source is directly proportional to its luminosity. MDA3 includes the same EM selection effects as MDA2, in the sense that the catalog is magnitude limited. The completeness of this catalog, defined in terms of luminosity rather than numbers of galaxies, as defined in Eq. (15), is 50% out to 115 Mpc. This is indicative that approximately 50% of the detected GW events have host galaxies inside the catalog.

To investigate the importance of luminosity weighting, MDA3 was analyzed twice under different assumptions, given in Eq. (A3). In the first, the analysis was matched to the known properties of the galaxy catalog, such that the probability of any galaxy hosting a GW event was

proportional to its luminosity. In the second, we feigned ignorance and ran the analysis with the assumption that each galaxy was equally likely to be host to a GW event (as was true in MDAs 1 and 2). This allows us to determine the effect of ignoring galaxy weighting with this dataset. The combined H_0 posteriors for both cases are shown in Fig. 7. The estimated values of the Hubble constant are $70.83^{+3.55}_{-2.72}$ km s⁻¹ Mpc⁻¹ (assuming hosts are luminosity weighted) and $69.50^{+4.20}_{-3.24}$ km s⁻¹ Mpc⁻¹ (assuming equal weights). By weighting the host galaxies with the correct function of their luminosities, which happens to be known in this case, the constraint on H_0 improves—the uncertainty narrows by a factor of 1.2, compared to the case in which equal weights are assumed. Both results are consistent with the fiducial H_0 value of 70 km s⁻¹ Mpc⁻¹. In the limit of a far greater number of events, one might expect to see a bias emerge in the case in which the assumptions in the analysis do not match those with which the catalog was simulated. The luminosity weighting of host galaxies, by its very nature, increases the probability that the host galaxy is inside the galaxy catalog; assuming equal weighting gives disproportionate weight to the contribution that comes from beyond the galaxy catalog. However, for the 249 BNS events considered here, the final posteriors are too broad to be able to detect any kind of bias.

E. Comparison between the MDAs

So far we have focused on individual event likelihoods and combined results for all 249 events. Our dataset also

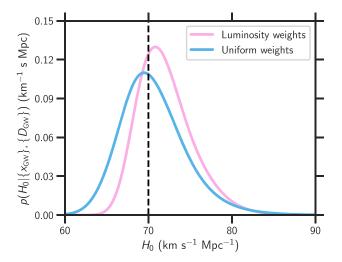


FIG. 7. Comparison of results with and without luminosity weighting. In MDA3, by construction, the probability of any galaxy hosting a GW event is proportional to its luminosity. The pink curve shows the posterior probability density on H_0 for the case where we take this into account in our analysis as a weighting by the galaxy's luminosity. The blue curve shows the posterior density for the case where we ignore this extra information, and treat every galaxy as equally likely to be hosts. Luminosity weighting improves the precision in the results by a factor of 1.2 for this set of simulations.

allows us to assess the convergence for the combined Hubble posterior as we add events. We calculate the intermediate combined posteriors as a function of the number of events, and show the resulting convergence in Fig. 8. We plot the fractional H_0 uncertainty (defined here as the half-width of the 68.3% credible interval divided by H_0 , $\Delta_{H0}^{68.3\%}/2H_0$), against the number of events we include in a randomly selected group. The scatter between realizations of the group is indicated by the error bars, which encompass 68.3% of their range. There is a considerable variation between different realizations, for the incomplete catalogs. For example, of the 100 realizations we used, for 25% completeness and 40 events, there are groups that give ~10% precision, but others that give ~70% precision.

With a sufficiently large number of events, we expect a $1/\sqrt{N}$ scaling of the uncertainty with the number of events [5,6]. To check whether this behavior is indeed true, we fit the results for each MDA to the expected scaling, obtaining the coefficient of $1/\sqrt{N}$ by maximizing its likelihood given the fractional uncertainties and their variances from the different realizations. The coefficient of the scaling is automatically dominated by the fractional uncertainties at large N where the variances are small. We show this scaling for each MDA as a set of dashed lines in Fig. 8.

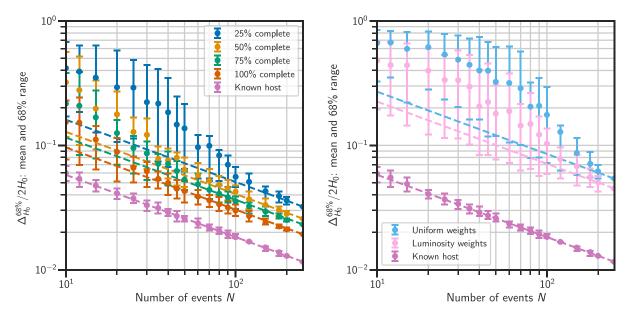


FIG. 8. Fractional uncertainty in H_0 as a function of the number N of the events for the combined H_0 posteriors. The fractional uncertainty in H_0 is defined as the half-width of the 68.3% highest probability density interval divided by 70 km s⁻¹ Mpc⁻¹, and is shown as the plotted dots for all cases. The error bars contain 68% of the scatter arising from different realizations of the events. (left) In purple, red, green, yellow and blue we show the associated host galaxy case (MDA0), complete galaxy catalog (MDA1) case, and the 75%, 50% and 25% completeness cases; we find a fractional H_0 uncertainty of 1.13%, 1.84%, 2.21%, 2.48% and 3.20% respectively for the combined H_0 posterior from 249 events. (right) convergence for MDA3 (event probability proportional to galaxy luminosity), analyzed with luminosity-weighted likelihood (pink) or equally-weighted likelihood (light blue). We find fractional H_0 uncertainties of 4.48% and 5.31% respectively. MDA0 (purple) is included for reference. We plot the expected $1/\sqrt{N}$ scaling behavior for large values of N for all cases with the dashed lines. This scaling behavior is met by all MDAs as the number of events reaches 249, but for the less informative, lower completeness MDAs the trend is slower to emerge. This is even more evident in MDA3, where the density of galaxies is 100 times greater, producing more potential hosts for each event. This is mitigated somewhat by the effect of luminosity-weighting the potential hosts (pink).

It can be seen that for each MDA, the results converge to the expected $1/\sqrt{N}$ scaling. The number of events required before this behavior is reached is dependent on the amount of EM information available on average for each event, in agreement with the results of [6]. The direct counterpart case is always on the trend after $\mathcal{O}(10)$ events, and shows a $\sim 18\%/\sqrt{N}$ convergence, comparable to and consistent with the results in [6,7]. With the most complete galaxy catalogs, if the host galaxy is not directly identified it will take tens of events before this behavior is reached. However, even the least complete catalog (25%) appears to have reached this behavior by the time all 249 events are combined. It should be noted that as the catalogs for MDAs 1 and 2 were not simulated realistically, their low density relative to the density of the universe means that these numbers should not be taken as predictions of how fast $1/\sqrt{N}$ may be reached (except, perhaps, in the counterpart case, although one should bear in mind that even for that case, peculiar velocities and redshift uncertainties have been neglected). Even with a galaxy catalog which is 25% complete, MDA2 gives a result which is only about a factor of 3 times worse than the counterpart case.

As the density of galaxies in MDA3 was increased by 2 orders of magnitude over MDAs 1 and 2, the final posteriors cannot be directly compared between MDAs. However, by plotting the equivalent convergence figure for MDA3 (including the "known host" case as a reference, see Fig. 8), the impact of increasing the density of galaxies in the universe on the rate at which the posterior converges on the $1/\sqrt{N}$ behavior becomes clear. When there are more host galaxies, the results are overall less precise, and take longer to reach the $1/\sqrt{N}$ trend. As expected, using luminosity-weighting of potential host galaxies as an assumption in the analysis concentrates the probability to a smaller number of galaxies, leading to a more precise result.

F. Limited robustness studies

Our results are expected to be sensitive to the luminosity distribution parameters—if one uses values of the Schechter function parameters α and L^* in the analysis which are different from the ones used to simulate the galaxy catalogs, one would expect to end up with a bias in the results. With variations of these parameters within their current measurement uncertainties, we have however demonstrated that the resulting variation in the final result is small compared to the statistical uncertainties reached with the current set of MDAs. Furthermore we have also demonstrated that our results are robust against a small $\mathcal{O}(1)$ variation in the value of the telescope sensitivity threshold m_{th} .

V. CONCLUSIONS AND OUTLOOK

The H_0 measurement using GW standard sirens has been demonstrated with recent events, including both the

counterpart method for GW170817 [19], and the galaxy catalog method [22,23]. These approaches are combined in the analysis of both BNS and BBH events from the first two observing runs of the advanced detector network [24], using the method described in this paper. Future measurements will rely on a combination of counterpart and catalog methods, as appropriate for each new detected event, with catalog incompleteness playing an important role for the more distant, yet more common, BBHs. This paper outlines a coherent approach that tackles both of these scenarios, including treatment of selection effects in both GWs (due to the limited sensitivity of GW detectors) and EM (due to the flux-limitations of EM observing channels). We performed a series of MDAs to validate our method using up to 249 observed events. For each of the MDAs analyzed, the final posterior on H_0 is found to be consistent with the value of $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$ used to simulate the MDA galaxy catalogs, demonstrating that our method can produce sufficiently unbiased results for treating these numbers of events, in our simulations.

GW selection effects are inherent in every version of the MDA and were corrected for by the term $p(D_{\rm GW}|H_0)$ in the denominator of Eq. (7). EM selection effects are addressed in MDAs 2 and 3 by the out-of-catalog terms containing \bar{G} in Eq. (9). In both these MDAs, in spite of having an apparent magnitude-limited galaxy catalog, we are able to accurately infer H_0 without any bias. MDA2 further demonstrates our ability to account for missing host galaxies down to a level where only 25% of events have hosts inside the catalog. Even in this case, we converge to the correct H_0 value, to the level of precision which could be reached by 249 events.

MDA3 demonstrates a clear tightening of the posterior distribution when we can assume that GW events trace the galaxy luminosities, compared to the case in which we treat all galaxies as equally likely hosts. The "uniform weights" analysis of MDA3 remains consistent with the simulated H_0 value. Hence we are unable to conclude whether an incorrect assumption would lead to a biased result, as one might expect. We used only 249 events for our MDAs. With enough events of comparable nature the bias would be detected. Future work will expand these studies to include a larger numbers of simulated GW events, and will be able to discern smaller sources of systematic effects.

Although the galaxy-catalog standard-siren measurement of H_0 is less precise than the counterpart measurement, it is still able to constrain H_0 , but requires at least an order of magnitude more events in order to reach a comparable accuracy (in the most realistic case of MDA3). These MDAs have validated our method and implementation in simplified scenarios. However future work will be needed to improve on this in several directions, to test its applicability to BBHs (which are detectable out to much farther distances), realistic cosmology, and real galaxy catalogs [6,24].

In both the counterpart and galaxy catalog cases, the lack of redshift uncertainties and peculiar velocities implies that the contributions from individual galaxies are a lot more precise than they would be in reality. Moreover, the simulated catalogs in MDAs 1 and 2 have a low density of galaxies compared to the universe, making them more informative than real catalogs. MDA3, with a galaxy density of 1 galaxy per 70 Mpc³, comes closest to the actual density of galaxies in the local universe of ~ 1 galaxy per 200 Mpc³ [49]. In this scenario there is still a clear convergence toward the simulated H_0 value. In comparison to actual catalogs such as GLADE [49], the apparent magnitude threshold of 14 is very low, so we expect a real catalog-only analysis to fall somewhere between MDAs 2 and 3. We caution the reader that with tens of events, the precision of results can vary by almost an order of magnitude depending on the particular realization of the detected population, before eventually converging to the expected $1/\sqrt{N}$ behaviour [5,6]. Analyzing more realistic catalogs will also require a sky-varying EM selection function, as the magnitude threshold varies significantly on the sky according to the design of particular surveys.

The galaxy distribution in these simulated catalogs is uniform in comoving volume. Although it has not been studied here, clustering of galaxies is expected to improve the constraint on H_0 (see, e.g., [6,56]), since even when the host is not in the catalog, it is likely that there will be observed galaxies nearby.

The Advanced LIGO—Virgo second observing run [21] has confirmed that BBH systems are detected at higher rates than BNSs. Since their greater mass allows them to be observed at much greater distances, where galaxy catalogs are incomplete, the catalog method including EM selection effects is particularly important. With the catalog of GW events expected to expand at an increasing rate in future observing runs, our analysis will evolve to meet the challenges that come with it, and give us the fullest picture of cosmology as revealed by gravitational waves.

ACKNOWLEDGMENTS

We thank members of the LIGO-Virgo Collaboration for valuable discussions pertaining to the writing of this paper, and in particular Nicola Tamanini for a careful reading of the manuscript. A. G. additionally thanks P. Ajith, Walter Del Pozzo, Anuradha Samajdar, and Chris Van Den Broeck for discussion at various stages of the work. R. G., C. M.

and J. V. are supported by the Science and Technology Research Council (Grant No. ST/L000946/1). I. M. H. is supported by the NSF Graduate Research Fellowship Program under grant No. DGE-17247915. I. M. H. also acknowledges support from NSF Grant No. PHY-1607585. H. Q. is supported by Science and Technology Facilities Council (Grant No. ST/T000147/1). A. S. thanks Nikhef for its hospitality and support from the Amsterdam Excellence Scholarship (2016-2018). H. Y. C. was supported by the Black Hole Initiative at Harvard University, through a grant from the John Templeton Foundation. M. F. and D. E. H. were supported by NSF grant No. PHY-1708081. They were also supported by the Kavli Institute for Cosmological Physics at the University of Chicago through an endowment from the Kavli Foundation. A. G. is supported by the research programme of the Netherlands Organisation for Scientific Research (NWO). D. E. H. gratefully acknowledges support from the Marion and Stuart Rice Award. We are grateful for computational resources provided by the Leonard E Parker Center for Gravitation, Cosmology and Astrophysics at the University of Wisconsin-Milwaukee, and those provided by Cardiff University, and funded by an STFC grant supporting UK Involvement in the Operation of Advanced LIGO. This article has been assigned LIGO document number LIGO-P1900017.

APPENDIX: DETAILED METHODOLOGY

1. A note on luminosity weighting and redshift evolution

The probability for a galaxy to host a GW event is not uniform over all the galaxies present in the catalog. Indeed, brighter galaxies are supposed to have an higher star-formation rate and hence have an higher probability to host a GW event. Also galaxies at higher redshifts may be more likely to be hosts, as mergers are expected to be more frequent [57]. Our prior belief for a galaxy at redshift z, sky position Ω and absolute and relative magnitudes M, m, to host a GW source s can be expressed as

$$p(z, \Omega, M, m|s, H_0)$$

$$= p(m|z, \Omega, M, s, H_0)p(z, \Omega, M|s, H_0), \quad (A1)$$

where if we assume that z, Ω and M are conditionally independent given s, H_0 ,

$$\begin{split} p(z,\Omega,M,m|s,H_0) &= \delta(m-m(z,M,H_0)) p(z|s) p(\Omega) p(M|s,H_0), \\ &= \delta(m-m(z,M,H_0)) \frac{p(s|z) p(z)}{p(s)} p(\Omega) \frac{p(s|M,H_0) p(M|H_0)}{p(s|H_0)}. \end{split} \tag{A2}$$

In the last equation we used the explicit relation between apparent magnitude and z, M and H_0 . The probability p(z) is the prior distribution of galaxies in the universe, taken to be uniform in comoving volume-time, $p(\Omega)$ is the prior on galaxy sky location, assumed uniform over the celestial

sphere, and $p(M|H_0)$ is the distribution of absolute magnitudes represented by the Schechter function. In the sections below we will show that the terms p(s) and $p(s|H_0)$ cancel out with other terms, and so their exact form does not need to be considered. $p(s|M,H_0)$ can take the form

$$p(s|M,H_0) \propto \begin{cases} L(M(H_0)) & \text{if GW hosting probability is proportional to luminosity} \\ \text{constant} & \text{if GW hosting probability is independent of luminosity}. \end{cases}$$
 (A3)

We refer to the above equation as luminosity weighting. The term p(s|z) represents the probability for the merger rate to depend on the redshift,

$$p(s|z) \propto \begin{cases} \text{function}(z) & \text{if rate evolves with redshift} \\ \text{constant} & \text{if rate is does not evolve with redshift.} \end{cases}$$
 (A4)

For the MDAs in this paper with $z \ll 1$, we assume a constant rate model but a more generic model with $p(s|z) \propto (1+z)^{\lambda}$ can be used with detections at higher redshifts. This was the case of [24], for example, in which a $p(s|z) \propto (1+z)^3$ was assumed.

2. A detailed breakdown of the galaxy catalog case

This section presents a more detailed look into the galaxy catalog method presented in Sec. II C 1. The approach is summarized in Fig. 9, a network diagram which shows how each of the parameters of this extended derivation fit together and their dependencies on each other. The parameters which appear in this diagram, and in the following subsections, are defined in Table III.

The subsections below provide derivations of the individual components of Eq. (9). Note that in the cases where the integration boundaries are not specified, they can be assumed to cover the full parameter space.

a. Likelihood when host is in catalog: $p(x_{GW}|G,D_{GW},H_0)$

The likelihood in the case where the host galaxy is inside the galaxy catalog, $p(x_{\rm GW}|G,D_{\rm GW},H_0)$, can be obtained from the marginalization over redshift, sky location, absolute magnitude and apparent magnitude. If we assume that the GW data, $x_{\rm GW}$, is independent of the galaxy catalog G, m and M we can write

$$p(x_{\rm GW}|G, D_{\rm GW}, s, H_0) = \frac{1}{p(D_{\rm GW}|G, s, H_0)} \iiint p(x_{\rm GW}|z, \Omega, s, H_0) p(z, \Omega, M, m|G, s, H_0) dz d\Omega dM dm.$$
 (A5)

The probability density function $p(z, \Omega, M, m|G, s, H_0)$ is taken as a sum of delta functions with specific z, Ω and m corresponding to the location of each galaxy in the catalog. This can be further factorized as

$$p(z, \Omega, M, m|G, s, H_0) = \frac{p(s|z, \Omega, M, m, G, H_0)\delta(M - M(z, m, H_0))p(z, \Omega, m|G)}{p(s|G, H_0)},$$
(A6)

where we have assumed again a relation between the apparent magnitude, redshift, H_0 and absolute magnitude. This allows us to integrate over the absolute magnitude in Eq. (A5) and obtain

$$p(x_{\text{GW}}|G, D_{\text{GW}}, s, H_0) = \frac{1}{p(D_{\text{GW}}|G, s, H_0)p(s|G, H_0)} \iiint p(x_{\text{GW}}|z, \Omega, s, H_0)p(s|z, \Omega, M(z, m, H_0), m, G, H_0) \times p(z, \Omega, m|G)dzd\Omega dm. \tag{A7}$$

Remembering that $p(z, \Omega, m|G)$ represents the distribution of the galaxies in the catalog, we can replace the integral above with a sum over the galaxies.

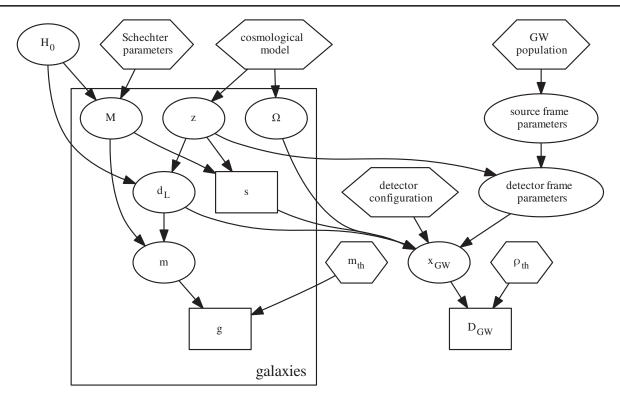


FIG. 9. A network diagram showing how the main parameters of the methodology interlink. Circular nodes denote ordinary parameters. Hexagonal nodes denote assumed knowns. Rectangular nodes denote binary flags. The arrows indicate the dependence of each parameter on the parameters which feed into them. The parameters grouped in the "galaxies" cluster are those which can be evaluated for every galaxy in the universe.

$$p(x_{\text{GW}}|G, D_{\text{GW}}, s, H_0) = \frac{1}{p(D_{\text{GW}}|G, s, H_0)p(s|G, H_0)} \frac{1}{N} \sum_{i=1}^{N} p(x_{\text{GW}}|z_i, \Omega_i, s, H_0)p(s|z_i)p(s|M(z_i, m_i, H_0)), \quad (A8)$$

where we have factorized $p(z_i|s)$ and $p(M(z_i, m_i, H_0)|s)$, together with the term $p(s|z, \Omega, M(z, m, H_0), m, G, H_0)$. Finally expanding the denominator $p(D_{GW}|G, s, H_0)$ in the same way, we can recover the likelihood for the "in catalog" part of the galaxy catalog method.

$$p(x_{\text{GW}}|G, D_{\text{GW}}, s, H_0) = \frac{\sum_{i=1}^{N} p(x_{\text{GW}}|z_i, \Omega_i, s, H_0) p(s|z_i) p(s|M(z_i, m_i, H_0))}{\sum_{i=1}^{N} p(D_{\text{GW}}|z_i, \Omega_i, s, H_0) p(s|z_i) p(s|M(z_i, m_i, H_0))}.$$
(A9)

Notably, in the case the galaxies in the catalogs are provided along with their redshift uncertainties $p(z_i)$, these can be implemented in the above equations as:

$$p(x_{\text{GW}}|G, D_{\text{GW}}, s, H_0) = \frac{\sum_{i=1}^{N_{\text{gal}}} \int p(x_{\text{GW}}|z_i, \Omega_i, s, H_0) p(s|z_i) p(s|M(z_i, m_i, H_0)) p(z_i) dz_i}{\sum_{i=1}^{N_{\text{gal}}} \int p(D_{\text{GW}}|z_i, \Omega_i, s, H_0) p(s|z_i) p(s|M(z_i, m_i, H_0)) p(z_i) dz_i}.$$
(A10)

b. Probability the host galaxy is in the galaxy catalog: $p(G|D_{GW},H_0)$ and $p(\bar{G}|D_{GW},H_0)$

The probability that the host galaxy is inside the galaxy catalog, given that a GW signal was detected, can be expressed as

$$p(G|D_{GW}, s, H_0) = \iiint p(G|z, \Omega, M, m, D_{GW}, s, H_0) p(z, \Omega, M, m|D_{GW}, s, H_0) dz d\Omega dM dm,$$

$$= \iiint \Theta[m_{th} - m] \frac{p(D_{GW}|z, \Omega, M, m, s, H_0) p(z, \Omega, M, m|s, H_0)}{p(D_{GW}|s, H_0)} dz d\Omega dM dm,$$

$$= \frac{1}{p(D_{GW}|s, H_0)} \iiint \Theta[m_{th} - m] p(D_{GW}|z, \Omega, s, H_0) p(z, \Omega, M, m|s, H_0) dz d\Omega dM dm. \tag{A11}$$

TABLE III. A summary of the parameters present in the network diagram, Fig. 9.

Parameter	Definition			
$\overline{H_0}$	The Hubble constant			
x_{GW}	The GW data associated with some GW source, s.			
$D_{ m GW}$	Denotes that a GW signal was detected, i.e., that $x_{\rm GW}$ passed some detection statistic threshold $\rho_{\rm th}$.			
g	Denotes that a galaxy is (G) , or is not (\bar{G}) , contained within the galaxy catalog.			
S	Denotes that a GW signal was emitted.			
M	Absolute magnitude.			
z	Redshift.			
Ω	Sky location (right ascension and declination).			
d_L	Luminosity distance.			
m	Apparent magnitude.			
$m_{ m th}$	Apparent magnitude threshold of the galaxy catalog.			
$ ho_{ m th}$	SNR threshold of the detector network.			
Cosmological model	The cosmological model assumed for the analysis. Typically a Friedmann-Lemaître-Robertson-Walker universe.			
Schechter parameters	The parameters which characterize the assumed absolute magnitude distribution of galaxies in the universe.			
GW population	The assumed underlying population of GW sources.			
Source frame parameters	Source frame parameters of a GW source, e.g., component masses, spins, inclination and polarization.			
Detector frame	As above, but redshifted into the detector frame.			
parameters				
Detector configuration	The network set up, including which detectors are included in the search and their noise floors.			

If we assume that the galaxy catalog is apparent magnitude-limited, such that only galaxies which are observed above some detection threshold are contained within it, we can approximate $p(G|z, \Omega, M, m, D_{GW}, s, H_0)$ as a Heaviside step around the detection threshold $m = m_{\text{th}}$.

$$p(G|D_{\mathrm{GW}}, s, H_0) = \frac{1}{p(D_{\mathrm{GW}}|s, H_0)} \iiint \Theta[m_{\mathrm{th}} - m] p(D_{\mathrm{GW}}|z, \Omega, s, H_0) p(z, \Omega, M, m|s, H_0) dz d\Omega dM dm. \tag{A12}$$

We now expand $p(z, \Omega, M, m|s, H_0)$ as in Eq (A2) and we obtain

$$p(G|D_{\text{GW}}, s, H_0) = \frac{1}{p(s)p(s|H_0)} \frac{1}{p(D_{\text{GW}}|s, H_0)} \int_0^{z(M, m_{\text{th}}, H_0)} dz \int d\Omega \int dM p(D_{\text{GW}}|z, \Omega, s, H_0) p(s|z) p(z) p(\Omega)$$

$$\times p(s|M, H_0) p(M|H_0). \tag{A13}$$

The term $p(D_{GW}|s, H_0)$ can be expanded in a similar way and finally gives the probability for the host galaxy to be in the catalog.

$$p(G|D_{GW}, s, H_0) = \frac{\int_0^{z(M, m_{th}, H_0)} dz \int d\Omega \int dM p(D_{GW}|z, \Omega, s, H_0) p(s|z) p(z) p(\Omega) p(s|M, H_0) p(M|H_0)}{\iiint p(D_{GW}|z, \Omega, s, H_0) p(s|z) p(z) p(\Omega) p(s|M, H_0) p(M|H_0) dz d\Omega dM}.$$
(A14)

As the probabilities of being in the catalog and not in the catalog must be complementary, we have,

$$p(\bar{G}|D_{GW}, s, H_0) = 1 - p(G|D_{GW}, s, H_0). \tag{A15}$$

c. Likelihood when host is not in catalog: $p(x_{GW}|\bar{G},D_{GW},H_0)$

We follow an approach similar to the one presented in Appendix A 2 a. We expand

$$p(x_{\rm GW}|\bar{G},D_{\rm GW},s,H_0) = \frac{1}{p(D_{\rm GW}|\bar{G},s,H_0)} \iiint p(x_{\rm GW}|z,\Omega,s,H_0) \frac{p(\bar{G}|z,\Omega,M,m,s,H_0)p(z,\Omega,M,m|s,H_0)}{p(\bar{G}|s,H_0)} dz d\Omega dM dm, \tag{A16}$$

The prior term, $p(z, \Omega, M, m|s, H_0)$ can now be expanded as it was in Eq (A2). Substituting this in, and utilizing a Heaviside step function to represent the galaxy catalog's apparent magnitude threshold for $p(\bar{G}|z, \Omega, M, m, s, H_0)$,

$$p(x_{\text{GW}}|\bar{G}, s, H_0) = \frac{1}{p(s)p(s|H_0)} \frac{1}{p(\bar{G}|s, H_0)} \int_{z(H_0, m_{\text{th}}, M)}^{\infty} dz \int d\Omega \int dM p(x_{\text{GW}}|z, \Omega, s, H_0) p(s|z)$$

$$\times p(z)p(\Omega)p(s|M, H_0)p(M|H_0).$$
(A17)

Expanding the denominator, $p(D_{GW}|\bar{G}, s, H_0)$, in the same way gives an equivalent term,

$$p(D_{\text{GW}}|\bar{G}, s, H_0) = \frac{1}{p(s)p(s|H_0)} \frac{1}{p(\bar{G}|s, H_0)} \int_{z(H_0, m_{\text{th}}, M)}^{\infty} dz \int d\Omega \int dM p(D_{\text{GW}}|z, \Omega, s, H_0)$$

$$\times p(s|z)p(z)p(\Omega)p(s|M, H_0)p(M|H_0). \tag{A18}$$

And substituting this back into Eq (A16) finally gives,

$$p(x_{\text{GW}}|\bar{G}, D_{\text{GW}}, s, H_0) = \frac{\int_{z(M, m_{\text{th}}, H_0)}^{\infty} dz \int d\Omega \int dM p(x_{\text{GW}}|z, \Omega, s, H_0) p(s|z) p(z) p(\Omega) p(s|M, H_0) p(M|H_0)}{\int_{z(M, m_{\text{th}}, H_0)}^{\infty} dz \int d\Omega \int dM p(D_{\text{GW}}|z, \Omega, s, H_0) p(s|z) p(z) p(\Omega) p(s|M, H_0) p(M|H_0)}.$$
(A19)

3. The catalog patch case

While in general the galaxy catalog method derived in A 2 was for use with a galaxy catalog which covers the entire sky, a small modification allows the use of catalogs which only cover a patch of sky, as long as the patch can be specified using limits in right ascension and declination. If we represent the sky area covered by the catalog as $\Omega_{\rm cat}$, and the area outside the catalog as $\Omega_{\rm rest}$, such that $\Omega_{\rm cat} + \Omega_{\rm rest}$ covers the whole sky, this can be written as follows:

$$\begin{split} p(x_{\rm GW}|D_{\rm GW},H_0) &= \int p(x_{\rm GW}|\Omega,D_{\rm GW},H_0)p(\Omega)d\Omega, \\ &= \int^{\Omega_{\rm cat}} p(x_{\rm GW}|\Omega,D_{\rm GW},H_0)p(\Omega)d\Omega \\ &+ \int^{\Omega_{\rm rest}} p(x_{\rm GW}|\Omega,D_{\rm GW},H_0)p(\Omega)d\Omega. \end{split} \tag{A20}$$

The first term is equivalent to the regular galaxy catalog case, but with limits on the integral over Ω , while the second term has no G and \bar{G} terms, and covers the rest of the sky from redshift 0 to ∞ .

4. Direct and pencil beam counterpart cases

The "direct" method assumes that the counterpart has been unambiguously linked to the host galaxy of the GW event, such that the redshift and sky location of that galaxy can be taken to be that of the GW event with certainty, see Eq. (10). Instead the numerator is calculated by evaluating the GW likelihood at the delta-function location of the counterpart in z and Ω , and the term in the denominator is evaluated as:

$$p(D_{\text{GW}}|H_0) = \iiint p(D_{\text{GW}}|z, \Omega, H_0)p(z)$$
$$\times p(\Omega)p(M|H_0)dzd\Omega dM, \qquad (A21)$$

for priors p(z) and $p(\Omega)$ (note that this is independent of galaxy catalog data).

The "pencil-beam" method makes the assumption that while the sky location of the galaxy associated with the counterpart is that of the GW event, we may not make a direct association to a known galaxy but to a set of potential candidate hosts. We can use the EM constrained sky localization and therefore return to the question of whether the host is within or beyond the galaxy catalog. In this case, the likelihood takes the same form as in the galaxy catalog

case, but evaluated along the line of sight of the candidate counterparts.

5. GW selection effects

Equation (8) in Sec. II C can be written as:

$$p(D_{\rm GW}|H_0) = \int p(D_{\rm GW}|x_{\rm GW}, H_0) p(x_{\rm GW}|H_0) dx_{\rm GW}. \tag{A22}$$

where $p(D_{\rm GW}|x_{\rm GW},H_0)$ is a binary quantity which is 1 if the SNR of $x_{\rm GW}$ passes ρ_{th} , and 0 otherwise.

Looking at the individual components of Eq. (9) in their expanded forms [e.g., Eq. (A10), (A14) and (A19)], $p(D_{\rm GW}|H_0)$ only appears in an expanded form, where it is additionally conditioned on z and Ω . Calculating $p(D_{\rm GW}|z,\Omega,H_0)$ requires integrating over all realizations of GW events (detected and not), for a range of z, Ω and H_0 values, and applying a detection threshold (ρ_{th}) which all events must pass in order to be deemed detected.

Practically, Monte-Carlo integration can be used:

$$p(D_{\text{GW}}|z, \Omega, H_0) = \frac{1}{N_{\text{samples}}} \sum_{i=1}^{N_{\text{samples}}} p(D_{\text{GW}i}|x_{\text{GW}i}, z, \Omega, H_0).$$
(A23)

where $x_{\text{GW}i}$ corresponds to an event, the parameters of which have been randomly drawn from the prior distributions of parameters which affect an event's detectability (mass, inclination, polarization, and sky location) and the event's ρ_i is calculated for specific values of z and H_0 .

$$p(D_{\text{GW}i}|x_{\text{GW}i}, z, \Omega, H_0) = \begin{cases} 1, & \text{if } \rho > \rho_{th} \\ 0, & \text{otherwise.} \end{cases}$$
 (A24)

which gives a smooth function for $p(D_{\rm GW}|z,\Omega,H_0)$, which drops from 1 to 0 over a range of z,Ω and H_0 values.

For the MDA, we use a ρ_{th} of 8 in each detector (assuming every event was detected by two detectors) and the 2016 PSD from [30] [Fig. 1], and evaluate $p(D_{\text{GW}i}|x_{\text{GW}i},z,\Omega,H_0)$ for 5000 samples, such that the integral converges. For this analysis, we assume that the probability of detection is averaged over the course of the entire simulated observation period, such that the dependence of D_{GW} on Ω is smeared out over the course of many days. We approximate this to mean that $p(D_{\text{GW}}|z,H_0)$ is uniform over the sky (ignoring the mild declination dependence which would remain after the rotation of the Earth is taken into account). Figure 10 shows how the probability of detection behaves as a function of z for different values of H_0 .

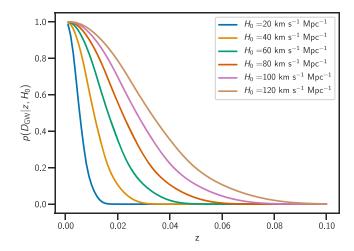


FIG. 10. Probability of detection, $p(D_{\text{GW}}|z, H_0)$, as a function of z for different values of H_0 . We assume a 2-detector network at O2-like sensitivity, for a population of binary neutron stars.

6. Prior mass distribution

An event's detectability is dependent on its observed (redshifted) detector-frame mass, M_z , but priors on the mass refer to their source-frame mass. When calculating $p(D|H_0)$ the masses are drawn from the priors on source mass, $p(M_1,M_2)$ and then converted to observed masses through the equation:

$$M_z = (1+z)M. \tag{A25}$$

However, when we use GW data in the form of posterior samples, the prior used to generate those is uniform on the redshifted mass, M_z [34]. Due to the way the MDA GW data was generated, with masses chosen on the detector-frame, rather than the source-frame, this was not something which had to be considered. With real GW data, as the redshift is linked directly to H_0 , it is necessary to take into account the redshifting of the masses explicitly.

In general, when calculating $p(D|H_0)$ for BBHs, the primary mass M_1 is drawn from a power-law with slope α , between limits $[a, b]M_{\odot}$. The secondary mass, M_2 is drawn from a uniform distribution between aM_{\odot} and M_1 [27], to give (for $\alpha \neq -1$):

$$p(M_1, M_2) = \frac{(\alpha + 1)M_1^{\alpha}}{bM_{\odot}^{(\alpha+1)} - aM_{\odot}^{(\alpha+1)}} \frac{1}{M_1 - aM_{\odot}}.$$
 (A26)

This is related to the redshifted mass by the Jacobian:

$$p(M_{1,z}, M_{2,z}) = p(M_1, M_2) \left| \frac{\partial(M_1, M_2)}{\partial(M_{1,z}, M_{2,z})} \right|,$$

$$= p(M_1, M_2) \left| \frac{1}{(1+z)^2} \right|. \tag{A27}$$

Substituting in our expression for $p(M_1, M_2)$:

$$\begin{split} p(M_{1,z},M_{1,z}) &= \frac{(\alpha+1)M_{1}^{\alpha}}{bM_{\odot}^{(\alpha+1)} - aM_{\odot}^{(\alpha+1)}} \frac{1}{M_{1} - aM_{\odot}} \frac{1}{(1+z)^{2}}, \\ &= \frac{(1+z)^{2}(\alpha+1)M_{1,z}^{\alpha}}{bM_{\odot,z}^{(\alpha+1)} - aM_{\odot,z}^{(\alpha+1)}} \frac{1}{M_{1,z} - aM_{\odot,z}} \frac{1}{(1+z)^{2}}, \\ &= \frac{(\alpha+1)M_{1,z}^{\alpha}}{bM_{\odot,z}^{(\alpha+1)} - aM_{\odot,z}^{(\alpha+1)}} \frac{1}{M_{1,z} - aM_{\odot,z}}. \quad (A28) \end{split}$$

The factor of $(1+z)^2$ cancels in the numerator and denominator. As all redshift (and hence H_0) dependence has been removed, no correction is required for the differing priors. For the case in which $\alpha=-1$, it can be shown that all redshift dependence falls out as well, meaning that as long as the prior mass distribution takes the form of a power law, no prior correction is required. This will not be the case for all mass distributions.

- [1] B. F. Schutz, Nature (London) **323**, 310 (1986).
- [2] D. E. Holz and S. A. Hughes, Astrophys. J. **629**, 15 (2005).
- [3] N. Dalal, D. E. Holz, S. A. Hughes, and B. Jain, Phys. Rev. D 74, 063006 (2006).
- [4] S. Nissanke, D. E. Holz, S. A. Hughes, N. Dalal, and J. L. Sievers, Astrophys. J. 725, 496 (2010).
- [5] S. Nissanke, D. E. Holz, N. Dalal, S. A. Hughes, J. L. Sievers, and C. M. Hirata, arXiv:1307.2638.
- [6] H.-Y. Chen, M. Fishbach, and D. E. Holz, Nature (London) 562, 545 (2018).
- [7] S. M. Feeney, H. V. Peiris, A. R. Williamson, S. M. Nissanke, D. J. Mortlock, J. Alsing, and D. Scolnic, Phys. Rev. Lett. 122, 061105 (2019).
- [8] E. Di Valentino, D. E. Holz, A. Melchiorri, and F. Renzi, Phys. Rev. D 98, 083523 (2018).
- [9] D. J. Mortlock, S. M. Feeney, H. V. Peiris, A. R. Williamson, and S. M. Nissanke, Phys. Rev. D 100, 103523 (2019).
- [10] W. M. Farr, M. Fishbach, J. Ye, and D. E. Holz, Astrophys. J. Lett. 883, L42 (2019).
- [11] N. Aghanim et al. (Planck Collaboration), arXiv:1807.06209.
- [12] A. G. Riess, S. Casertano, W. Yuan, L. M. Macri, and D. Scolnic, Astrophys. J. 876, 85 (2019).
- [13] E. Macaulay *et al.* (DES Collaboration), Mon. Not. R. Astron. Soc. **486**, 2184 (2019).
- [14] S. Birrer et al., Mon. Not. R. Astron. Soc. 484, 4726 (2019).
- [15] W. L. Freedman et al., Astrophys. J. 882, 34 (2019).
- [16] B. P. Abbott, R. Abbott, T. D. Abbott, F. Acernese, K. Ackley, C. Adams, T. Adams, P. Addesso *et al.* (LIGO Scientific and Virgo Collaborations), Phys. Rev. Lett. 119, 161101 (2017).
- [17] B. P. Abbott, R. Abbott, T. D. Abbott, F. Acernese, K. Ackley, C. Adams, T. Adams, P. Addesso *et al.*, Astrophys. J. **848**, L12 (2017).
- [18] M. Soares-Santos, D. E. Holz, J. Annis, R. Chornock, K. Herner, Berger et al., Astrophys. J. Lett. 848, L16 (2017).
- [19] B. P. Abbott, R. Abbott, T. D. Abbott, F. Acernese, K. Ackley, C. Adams, T. Adams, P. Addesso *et al.* (LIGO Scientific and Virgo Collaborations), Nature (London) 551, 85 (2017).
- [20] W. Del Pozzo, Phys. Rev. D 86, 043011 (2012).
- [21] B. P. Abbott *et al.* (LIGO Scientific and Virgo Collaborations), Phys. Rev. X **9**, 031040 (2019).
- [22] M. Fishbach *et al.* (LIGO Scientific and Virgo Collaborations), Astrophys. J. 871, L13 (2019).

- [23] M. Soares-Santos et al. (DES, LIGO Scientific, and Virgo Collaborations), Astrophys. J. 876, L7 (2019).
- [24] B. P. Abbott *et al.* (LIGO Scientific and Virgo Collaborations), arXiv:1908.06060.
- [25] B. P. Abbott, R. Abbott, T. D. Abbott, M. R. Abernathy, F. Acernese, K. Ackley, C. Adams, T. Adams, P. Addesso, and R. X. Adhikari (KAGRA, LIGO Scientific, and Virgo Collaborations), Living Rev. Relativity 21, 3 (2018).
- [26] T. Padma, Nature (London) https://doi.org/10.1038/d41586-019-00184-z (2019).
- [27] B. P. Abbott *et al.* (LIGO Scientific and Virgo Collaborations), Astrophys. J. Lett. 882, L24 (2019).
- [28] B. P. Abbott, R. Abbott, T. D. Abbott, S. Abraham, F. Acernese, K. Ackley, C. Adams, R. X. Adhikari, V. B. Adya, and C. Affeldt, Astrophys. J. 875, 161 (2019).
- [29] R. Nair, S. Bose, and T. D. Saini, Phys. Rev. D **98**, 023502 (2018).
- [30] L. P. Singer, L. R. Price, B. Farr, A. L. Urban, C. Pankow, S. Vitale, J. Veitch, W. M. Farr, C. Hanna, K. Cannon, T. Downes, P. Graff, C.-J. Haster, I. Mandel, T. Sidery, and A. Vecchio, Astrophys. J. 795, 105 (2014).
- [31] C. P. L. Berry et al., Astrophys. J. 804, 114 (2015).
- [32] D. W. Hogg, arXiv:astro-ph/9905116.
- [33] M. Maggiore, Gravitational Waves Volume 1: Theory and Experiments (Oxford University Press, Oxford, 2008).
- [34] J. Veitch et al., Phys. Rev. D 91, 042003 (2015).
- [35] S. R. Taylor, J. R. Gair, and I. Mandel, Phys. Rev. D 85, 023535 (2012).
- [36] C. Messenger and J. Read, Phys. Rev. Lett. 108, 091101 (2012).
- [37] C. M. Springob, C. Magoulas, M. Colless, J. Mould, P. Erdogdu, D. H. Jones, J. R. Lucey, L. Campbell, and C. J. Fluke, Mon. Not. R. Astron. Soc. **445**, 2677 (2014).
- [38] J. Carrick, S. J. Turnbull, G. Lavaux, and M. J. Hudson, Mon. Not. R. Astron. Soc. 450, 317 (2015).
- [39] J. Hjorth, A. J. Levan, N. R. Tanvir, J. D. Lyman, R. Wojtak, S. L. Schrøder, I. Mandel, C. Gall, and S. H. Bruun, Astrophys. J. 848, L31 (2017).
- [40] C. Guidorzi et al., Astrophys. J. 851, L36 (2017).
- [41] C. Howlett and T. M. Davis, Mon. Not. R. Astron. Soc. **492**, 3803 (2020).
- [42] S. Mukherjee, G. Lavaux, F. R. Bouchet, J. Jasche, B. D. Wandelt, S. M. Nissanke, F. Leclercq, and K. Hotokezaka, arXiv:1909.08627.

- [43] C. Nicolaou, O. Lahav, P. Lemos, W. Hartley, and J. Braden, arXiv:1909.09609.
- [44] J. De Vicente, E. Sánchez, and I. Sevilla-Noarbe, Mon. Not. R. Astron. Soc. 459, 3078 (2016).
- [45] I. Sadeh, F.B. Abdalla, and O. Lahav, Publ. Astron. Soc. Pac. 128, 104502 (2016).
- [46] M. Fishbach, D. E. Holz, and W. M. Farr, Astrophys. J. Lett. 863, L41 (2018).
- [47] C. Messick et al., Phys. Rev. D 95, 042001 (2017).
- [48] W. Del Pozzo, C. P. Berry, A. Ghosh, T. S. F. Haines, L. P. Singer, and A. Vecchio, Mon. Not. R. Astron. Soc. 479, 601 (2018).
- [49] G. Dálya, G. Galgóczi, L. Dobos, Z. Frei, I. S. Heng, R. Macas, C. Messenger, P. Raffai, and R. S. de Souza, Mon. Not. R. Astron. Soc. 479, 2374 (2018).
- [50] P. Schechter, Astrophys. J. 203, 297 (1976).
- [51] L. P. Singer, H.-Y. Chen, D. E. Holz, W. M. Farr, L. R. Price, V. Raymond, S. B. Cenko, N. Gehrels, J. Cannizzo,

- M. M. Kasliwal, S. Nissanke, M. Coughlin, B. Farr, A. L. Urban, S. Vitale, J. Veitch, P. Graff, C. P. L. Berry, S. Mohapatra, and I. Mandel, Astrophys. J. **829**, L15 (2016).
- [52] C. N. Leibler and E. Berger, Astrophys. J. **725**, 1202 (2010).
- [53] W. Fong, E. Berger, R. Chornock, R. Margutti, A. J. Levan, N. R. Tanvir, R. L. Tunnicliffe, I. Czekala, D. B. Fox, D. A. Perley, S. B. Cenko, B. A. Zauderer, T. Laskar, S. E. Persson, A. J. Monson, D. D. Kelson, C. Birk, D. Murphy, M. Servillat, and G. Anglada, Astrophys. J. 769, 56 (2013).
- [54] C. Cutler and E. E. Flanagan, Phys. Rev. D 49, 2658 (1994).
- [55] A. M. Walker, J. R. Stat. Soc. Ser. B 31, 80 (1969).
- [56] C. L. MacLeod and C. J. Hogan, Phys. Rev. D 77, 043512 (2008).
- [57] P. Madau and M. Dickinson, Annu. Rev. Astron. Astrophys. 52, 415 (2014).