


Associating host galaxy candidates to massive black hole binaries resolved by pulsar timing arrays

Janna M. Goldstein ¹★, Alberto Sesana,¹ A. Miguel Holgado ² and John Veitch ^{1,3}

¹*School of Physics and Astronomy and Institute for Gravitational Wave Astronomy, University of Birmingham, Edgbaston, Birmingham B15 2TT, UK*

²*Department of Astronomy and National Center for Supercomputing Applications, University of Illinois at Urbana-Champaign, Urbana, IL 61801, USA*

³*School of Physics and Astronomy, University of Glasgow, Glasgow G12 8QQ, UK*

Accepted 2019 February 7. Received 2019 February 5; in original form 2018 December 6

ABSTRACT

We propose a novel methodology to select host galaxy candidates of future pulsar timing array (PTA) detections of resolved gravitational waves (GWs) from massive black hole binaries (MBHBs). The method exploits the physical dependence of the GW amplitude on the MBHB chirp mass and distance to the observer, together with empirical MBH mass–host galaxy correlations, to rank potential host galaxies in the mass–redshift plane. This is coupled to a null-stream-based likelihood evaluation of the GW amplitude and sky position in a Bayesian framework that assigns to each galaxy a probability of hosting the MBHB generating the GW signal. We test our algorithm on a set of realistic simulations coupling the likely properties of the first PTA resolved GW signal to synthetic all-sky galaxy maps. For a foreseeable PTA sky-localization precision of 100 deg^2 , we find that the GW source is hosted with 50 per cent (90 per cent) probability within a restricted number of $\lesssim 50$ ($\lesssim 500$) potential hosts. These figures are orders of magnitude smaller than the total number of galaxies within the PTA sky error-box, enabling extensive electromagnetic follow-up campaigns on a limited number of targets.

Key words: black hole physics – gravitational waves – pulsars: general.

1 INTRODUCTION

Multimessenger astronomy with gravitational waves (GWs) has long been anticipated as one of the ‘Holy Grails’ for the understanding of the Universe. After a long wait, the first spectacular confirmation of its potential came with the detection of GWs from GW170817 (Abbott et al. 2017a), a binary neutron star coalescence (BNS) at about 40 Mpc distance, accompanied by a bright electromagnetic (EM) signal observed at all wavelengths (Abbott et al. 2017c). The wealth of fresh information brought by this event has been the key to confirming several theoretical speculations, from the short gamma-ray burst–BNS merger connection (Abbott et al. 2017d) to the synthesis through r-processes of the heavy elements permeating the Universe (Chornock et al. 2017), and opened a new way to do cosmology with standard sirens (Schutz 1986; Abbott et al. 2017b; Fishbach et al. 2018). All of this has been achieved thanks to the excellent sky-localization and distance information provided by LIGO–Virgo, an intense follow-up campaign, and the presence of a bright, distinct EM counterpart that could be easily singled out from other possible candidates. The small size of the sky-localization error-box was crucial, since it allowed

systematic scanning of a relatively low number of possible galaxy hosts.

Realizing the full potential of multimessenger astronomy might prove more difficult in the low-frequency band relevant to space-based interferometers such as LISA (Amaro-Seoane et al. 2017) and pulsar timing arrays (PTAs Verbiest et al. 2016), where the expected loudest sources involve inspiral and merger of massive black hole binaries (MBHBs) at cosmological distances (Sesana, Vecchio & Colacino 2008; Klein et al. 2016). Merging MBHBs are not per se expected to produce EM signals, so multimessenger efforts need to rely on some distinctive signature in the emission of the gas that might be accreted by the system during the inspiral and final coalescence (Tang, MacFadyen & Haiman 2017). Even so, it is not clear what that signature would be and a range of possibilities have been proposed, from periodicity (e.g. Sesana et al. 2012) to peculiar spectral features (e.g. Tanaka, Menou & Haiman 2012), and EM chirps (e.g. Haiman 2017).

The situation is particularly challenging for PTAs. Besides detecting a stochastic gravitational wave background (GWB) produced by the superposition of many MBHB systems (e.g. Phinney 2001; Sesana et al. 2008; Ravi et al. 2012), PTAs also have the capability to detect and localize in the sky particularly loud MBHBs (Sesana, Vecchio & Volonteri 2009; Ravi et al. 2015; Rosado, Sesana & Gair 2015; Kelley et al. 2018). Mingarelli et al. (2017) predict that

* E-mail: jgoldstein@star.sr.bham.ac.uk

in 10 yr, the first resolved binary could be detected. Strategies for optimizing PTA for single source detection (by allocating observing time and targeted searched for new pulsars) have been proposed by e.g. Burt, Lommen & Finn (2011), Simon et al. (2014, by identifying ‘hot spots’ from nearby galaxy clusters), and Lam (2018). However both the prediction of Kelley et al. (2018) and the optimization for a detection of Lam (2018) are complicated by the difficult to model red noise of the pulsars.

When looking for an EM counterpart to the first PTA resolved binary detections, one faces three main problems. First, the sky-localization is expected to be relatively poor (of the order of hundreds of deg²; Sesana & Vecchio 2010; Goldstein et al. 2018). Secondly, the detected GW signal is likely to be monochromatic. The absence of observable frequency evolution (chirp) of the waveform prevents one from separating the source mass from the distance, since only the overall amplitude A and frequency f are measured. Last, the signal evolves slowly in time, with a periodicity of the order of years. Associated counterparts might be identified through peculiar features in the source luminosity or through potential peculiarities of the galaxy host (Sesana et al. 2012; Tanaka et al. 2012; Burke-Spolaor 2013). In any case there is no clear smoking-gun event such as a transient counterpart, as is the case for a BNS merger.

It is therefore crucial to find a way to identify the most promising host galaxy candidates among the millions of objects falling within the source sky location error-box. In this paper, we develop a Bayesian framework to identify the most likely hosts by matching the information contained in a hypothetical PTA detection to candidate galaxy properties. The key point around which our analysis is built is that individually resolvable sources in the PTA band necessarily have a large strain amplitude A (Rosado et al. 2015; Kelley et al. 2018), which can result only from particularly massive and/or nearby MBHBs. We show that this allows one to exclude at high confidence the vast majority of the galaxies in the error-box, significantly reducing the number of candidates.

To demonstrate this, we consider a synthetic PTA and inject GW signals with properties compatible to the first single sources to be detected by future PTAs, drawn by following the procedure described in Rosado et al. (2015). We then use the null-stream analysis developed in Goldstein et al. (2018) to construct the 3D likelihood function of the signal amplitude A and sky-localization θ, ϕ . We extract a mock catalogue of galaxies from the synthetic all sky maps obtained by Henriques et al. (2012) from the Millennium Simulation (Springel et al. 2005) and use Bayesian inference to rank host candidates.

The paper is organized as follows. In Section 2 we lay out the mathematical basis of our experiment, including the construction of a likelihood from null-streams and the Bayesian framework for the computation of a host galaxy probability. This framework is then applied in Section 3 to a number of representative simulations with results laid out in Section 4 and the main conclusions and outlook presented in Section 5.

2 MATHEMATICAL FRAMEWORK

2.1 Signal model and null-stream sky-localization

PTAs are capable of reconstructing the incoming direction of a deterministic GW source via triangulation (Babak & Sesana 2012; Boyle & Pen 2012), providing that three or more millisecond pulsars (MSPs) contribute to the detection. We consider, for simplicity, a circular, monochromatic MBHB. The emitted GW can be written

in the form (Jaranowski, Królak & Schutz 1998)

$$h^+(t) = A^+ \cos(2\psi) - A^\times \sin(2\psi) \quad (1)$$

$$h^\times(t) = A^+ \sin(2\psi) + A^\times \cos(2\psi), \quad (2)$$

where ψ is the GW polarization angle and

$$A^+ = A \frac{1}{2}(1 + \cos \iota) \cos(2\pi f t + \phi_0) \quad (3)$$

$$A^\times = A(\cos \iota) \sin(2\pi f t + \phi_0). \quad (4)$$

The two polarization amplitudes A^+, A^\times are modulated with the *observed* GW frequency f and are related to the intrinsic amplitude¹

$$A = 4 \frac{(G \mathcal{M}_z)^{5/3} (\pi f)^{2/3}}{D_1} \quad (5)$$

via the inclination angle to the line of sight ι . The amplitude A is a function of the source redshifted chirp mass

$$\mathcal{M}_z = (1+z)\mathcal{M} = (1+z) \frac{(M_1 M_2)^{3/5}}{(M_1 + M_2)^{1/5}}, \quad (6)$$

and of its luminosity distance

$$D_1 = (1+z)D_H \int_0^z \frac{dz'}{E(z')}. \quad (7)$$

In the above equations, M_1 and M_2 are the masses of the two black holes forming the binary, z is the source redshift, $D_H = c/H_0$, and $E(z) = \sqrt{\Omega_M(1+z)^3 + \Omega_\Lambda}$, with Ω_M and Ω_Λ being the fractional mass and cosmological constant energy content, H_0 the Hubble constant and assuming a standard flat Λ CDM universe (Planck Collaboration XIII 2016).

The GW induces into the pulse time of arrival a redshift of the form

$$z(t, \hat{\Omega}) = F^+(\hat{\Omega})h^+(t) + F^\times(\hat{\Omega})h^\times(t), \quad (8)$$

where the ‘antenna beam patterns’ F^+ and F^\times depend on the angle between the incoming GW direction $\hat{\Omega}$ and the known position of the MSP (see e.g. Anholm et al. 2009). In practice, PTAs are sensitive to the two wave polarizations h^+, h^\times that depend on the vector of parameters $(A, \iota, f, \psi, \phi_0, \theta, \phi)$, where we decomposed the incoming wave direction $\hat{\Omega}$ on to its (θ, ϕ) coordinates in the sky.

In Goldstein et al. (2018) we developed a null-stream-based analysis (see also Zhu et al. 2015 and Hazboun & Larson 2016) that, among other things, can be used to infer the amplitude and incoming direction of the GW source. Since for an individual GW source there are only two polarizations, but an array of N allows measurement of N independent time (or frequency) series, it is possible to apply a matrix transformation that ‘collapses’ the signal into two of these time series. This nulls the signal contribution in all the others, hence constructing $N - 2$ null-streams. Formally, the transformation takes the form (see Goldstein et al. 2018, for

¹This definition of A is equivalent to the definition with a prefactor of 2 instead of 4 – which is also seen in the literature, e.g. in Babak et al. (2016) – as that definition is accompanied by an additional factor of 2 in equations (3) and (4).

details):

$$\mathbf{M} \mathbf{d} = \begin{pmatrix} h^+ \\ h^\times \\ \eta_1 \\ \vdots \\ \eta_{N-2} \end{pmatrix} + \mathbf{M} \mathbf{n} \equiv \mathbf{h} + \mathbf{M} \mathbf{n}, \quad (9)$$

where \mathbf{d} represents the original N time series of the N pulsars (including signal and noise \mathbf{n}), \mathbf{M} is the matrix transformation, $\eta_i = 0$ are the null streams, and \mathbf{h} is the combined vector of GW polarizations and null streams. In practice this amounts to the construction of N linear combinations of the timing residuals so that the GW signal is present only in two of them and null in all the others.

The null streams can then be used to construct the likelihood function

$$l = -\frac{1}{2} \left((\mathbf{M} \mathbf{d} - \mathbf{h})^\top ((\mathbf{M}^{-1})^\top \Gamma \mathbf{M}^{-1}) (\mathbf{M} \mathbf{d} - \mathbf{h}) \right) + \text{norm}, \quad (10)$$

where Γ is the inverse of the covariance matrix appropriate for the expected noise of the detector. For a signal detected at frequency f , marginalization of the likelihood over the parameters ι , ψ , ϕ , yields the 3D likelihood function $\mathcal{L}(A, \theta, \phi)$. For an example $\mathcal{L}(A, \theta, \phi)$, see Fig. 1.

In this work, we use the Goldstein et al. (2018) null-stream pipeline to obtain $\mathcal{L}(A, \theta, \phi)$. However in principle any method could be used to localize the source, as long as it can provide a joint likelihood on the sky location and amplitude of the signal. The framework for candidate host galaxy selection, which is introduced in the following section, is written in term of a generic input $\mathcal{L}(A, \theta, \phi)$.

2.2 Bayesian inference for galaxy host

Our goal is to combine the likelihood information $\mathcal{L}(A, \theta, \phi)$ with individual galaxy properties to assess the probability of each given galaxy to be the host of the detected GW source. The question we want to answer in practice is: given the detection of a signal with 3D likelihood described by $\mathcal{L}(A, \theta, \phi)$, what is the probability that a galaxy G_i described by a set of observed parameters λ – known with prior probability $p(\lambda|G_i)$ – is the host of the signal source? To answer this question we need a theoretical model that connects the strength and location of a putative GW signal to observable galaxy parameters.

Since MBHBs reside in the centre of galaxies, the sky coordinates of each specific galaxy (θ_G, ϕ_G) coincide with the sky coordinates of the putative GW source. We therefore have $\theta_G = \theta$ and $\phi_G = \phi$. Furthermore, we see from equations (5) and (6) that the GW amplitude A depends on the source chirp mass \mathcal{M} and luminosity distance D_1 . This latter can be easily measured from the galaxy spectroscopic redshift by assuming a fiducial cosmology. Whereas \mathcal{M} can be written in terms of the total binary mass M and mass ratio $q = M_2/M_1$ (with $M_2 \leq M_1$) as: $\mathcal{M} = M q^{3/5}/(1+q)^{6/5}$, we can assume the total mass to be related to the bulge mass via an $M - M_b$ relation of the form

$$\log_{10} \left(\frac{M}{M_\odot} \right) = \alpha + \beta \log_{10} \left(\frac{M_b}{10^{11} M_\odot} \right) \quad (11)$$

which connects the total binary mass to the observable galaxy bulge stellar mass M_b . If we group the $M - M_b$ constants α and β with

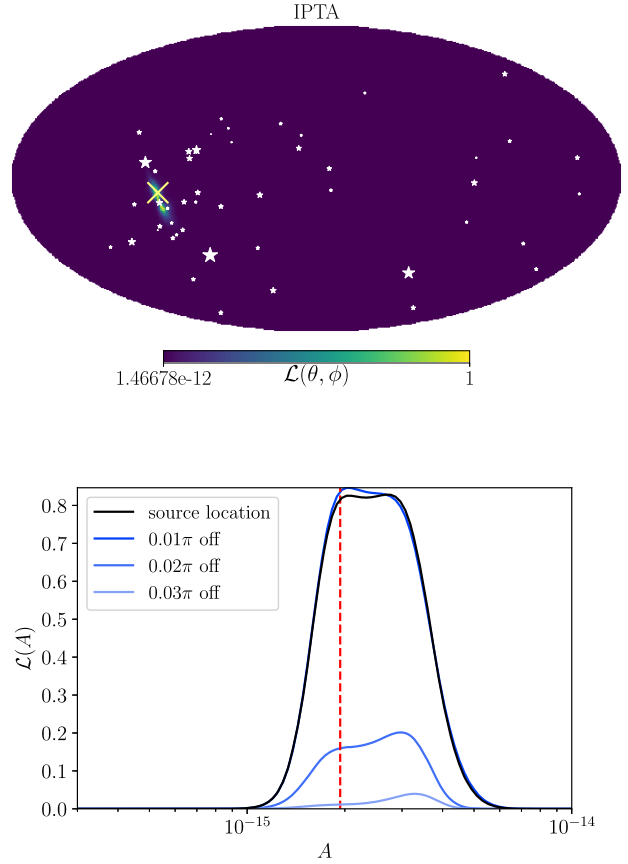


Figure 1. Example of $\mathcal{L}(A, \theta, \phi)$ as output by the null-stream pipeline. The injected signal is for source A, with $S/N = 12$ (see Section 3). Top: Likelihood marginalized over A (i.e. $\mathcal{L}(\theta, \phi)$ with an arbitrary normalization). The IPTA pulsars are marked with stars, where the size of the star corresponds to the noise level of the pulsar (with bigger stars for lower noise). The yellow cross indicates the position of the injected source. Bottom: $\mathcal{L}(A|\theta_s, \phi_s)$ at the source position (θ_s, ϕ_s) (in black) and at some offset positions $(\theta_s, \phi_s + \Delta)$ (in blue). The likelihoods are normalized only with respect to each other. The red dashed line is placed at the injected amplitude value.

the galaxy parameters, the vector of seven parameters

$$\lambda = (M_b, D_1, \theta, \phi, q, \alpha, \beta) \quad (12)$$

is sufficient to connect a specific galaxy to the GW strain. All of them but q , α , and β can be directly extracted from observations.

Formally the full calculation can be cast in term of Bayes' theorem. Let $P(G_i|d)$ be the probability of galaxy G_i being the host galaxy, given some data d , then:

$$P(G_i|d) = \frac{P(G_i)}{P(d)} P(d|G_i) = \frac{P(G_i)}{P(d)} \int p(d|\lambda) p(\lambda|G_i) d\lambda, \quad (13)$$

where $P(d) = \sum_i P(d|G_i)$ is the likelihood of the data marginalized over all galaxies (or evidence). $P(G_i)$ is the prior probability of G_i being the host, which we take to be a constant, having no reason a priori to prefer any particular galaxy. Of interest is the shape of the distribution of $P(G_i|d)$, so disregarding the constant prefactor $P(G_i)/P(d)$, we are left with the likelihoods $P(d|G_i)$.

The likelihood of a specific galaxy G_i to be the host of the GW source is given by the integral in equation (13) and is composed of the probability of the data given the source parameters $p(d|\lambda)$, times the prior distribution on these parameters $p(\lambda|G_i)$, integrated over all the relevant variables given in equation (12).

What is needed is an operational form for $p(d|\lambda)$. First, the amplitude A is independent on θ, ϕ and so

$$\begin{aligned} p(d|\lambda) &= p(d|M_b, D_1, \theta, \phi, q, \alpha, \beta) \\ &= p(d|A, \theta, \phi) p(A|M_b, D_1, q, \alpha, \beta). \end{aligned} \quad (14)$$

Secondly, A is a direct function of the chirp mass \mathcal{M} and distance only, we can therefore write

$$p(A|M_b, D_1, q, \alpha, \beta) = p(A|\mathcal{M}, D_1) p(\mathcal{M}|M_b, q, \alpha, \beta). \quad (15)$$

Last, \mathcal{M} is a function of q and M , and the latter is related to M_b by the $M - M_b$ relation. We therefore have

$$p(\mathcal{M}|M_b, q, \alpha, \beta) = p(\mathcal{M}|M, q) p(M|M_b, \alpha, \beta). \quad (16)$$

Putting the chain together we get

$$\begin{aligned} p(d|G_i) &= \int p(d|A, \theta, \phi) p(A|\mathcal{M}, D_1) p(\mathcal{M}|M, q) \\ &\quad p(M|M_b, \alpha, \beta) p(M_b, D_1, \theta, \phi, q, \alpha, \beta|G_i) \\ &\quad dM_b dD_1 d\theta d\phi dq d\alpha d\beta dM. \end{aligned} \quad (17)$$

We can now specify the individual elements of equation (17) for practical computational purposes.

(1) $p(\lambda|G_i) = p(M_b, D_1, \theta, \phi, q, \alpha, \beta|G_i)$ describes the prior knowledge of each galaxy property and the underlying $M - M_b$ constants. We assume that all five galaxy parameters – so excluding α and β – are independent so that the prior can be factorized as $p(\lambda|G_i) = \prod_{j=1}^5 p(\lambda_j|G_i)$. In particular,

(i) M_b in real surveys is generally obtained from the galaxy luminosity via bulge–disc decomposition. M_b is then computed from the bulge luminosity by assuming a stellar mass function. Typical uncertainties in this procedure can be up to a factor of two (Longhetti & Saracco 2009). Nonetheless, as a first approximation, we take M_b to be known exactly, reducing the prior $p(M_b)$ to a delta function (so the integral over M_b drops out).

(ii) D_1 is computed from the spectroscopic redshift of the galaxy z via equation (7). Uncertainties on the cosmological parameters $H_0, \Omega_M, \Omega_\Lambda$ are of the order of a few per cent (Planck Collaboration XIII 2016) and weak lensing is subdominant for the $z < 1$ galaxies relevant here (Shapiro et al. 2010). We therefore also assume D_1 to be known exactly, reducing the prior $p(D_1)$ to a delta function, dropping the integration over D_1 from the likelihood marginalization.

(iii) θ, ϕ are generally determined with arcsecond precision, which for any practical purposes can be treated as delta functions as well.

(iv) q , the binary mass ratio, is essentially undetermined. We therefore use a broad log flat prior between $-2 \leq \log_{10}(q) \leq 0$ (i.e. $0.01 \leq q \leq 1$).

The impact of changing the adopted priors in the calculation are discussed in Section 4.3.

(2) $p(d|A, \theta, \phi)$ is directly proportional to the likelihood in the 3D amplitude–sky location space $\mathcal{L}(A, \theta, \phi)$ returned as a numerical function with finite resolution by our null-stream-based parameter estimation pipeline. Given the values of A, θ , and ϕ from the priors, we select the numerical value from the sky pixel at (θ, ϕ) and the closest sampled amplitude to A . The sampling range (10^{-17} – 10^{-14}) is big enough to cover the area of interest, so for values of A outside this range, the likelihood is set to zero.

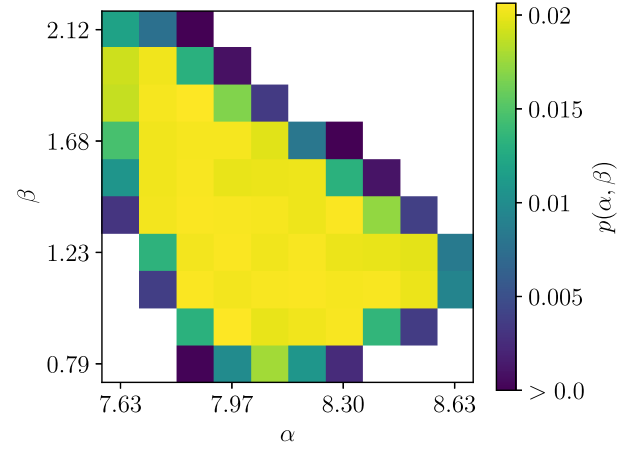


Figure 2. Prior on the $M - M_b$ constants (α, β) constructed from the compilation of $M - M_b$ relations in Middleton et al. (2018), see the text in Section 2.2. The prior is binned in a 10×10 regular grid with $\alpha \in [7.63, 8.63]$ and $\beta \in [0.79, 2.14]$. The pixels are normalized such that their sum is one. Some combinations (α, β) have zero prior weight and are masked in white.

(3) $p(A|\mathcal{M}, D_1)$ is determined by the GW quadrupole formula. Given the system chirp mass and distance, the amplitude is univocally determined by equation (5). We can thus write

$$p(A|\mathcal{M}, D_1) = \delta \left(A - 4 \frac{(G\mathcal{M}_z)^{5/3} (\pi f)^{2/3}}{D_1} \right). \quad (18)$$

(4) $p(\mathcal{M}|M, q)$ is similarly computed from the mathematical definition of \mathcal{M} in terms of M, q as

$$p(\mathcal{M}|M, q) = \delta \left(\mathcal{M} - \frac{Mq^{3/5}}{(1+q)^{6/5}} \right). \quad (19)$$

(5) $p(M|M_b, \alpha, \beta)$ is a core ingredient of the calculation. The possibility of ranking galaxy hosts stems from the simple fact that extremely massive black holes are hosted in extremely massive galaxies, a relation that has to be handled with care. Once a specific $M - M_b$ relation of the form given by equation (11) with intrinsic dispersion ϵ is given, the MBH total mass probability is described by a lognormal prior

$$p(M|M_b) = \frac{1}{\sqrt{2\pi\epsilon^2}} \exp \left\{ - \frac{\left[\log \frac{M}{M_\odot} - \left(\alpha + \beta \log \frac{M_b}{10^{11} M_\odot} \right) \right]^2}{2\epsilon^2} \right\} \quad (20)$$

that we integrate from -3ϵ to $+3\epsilon$ around the minimum and maximum expectation values of M [the range of M values being due to the spread in (α, β)].

The $M - M_b$ relation is quite uncertain, as demonstrated by its many different flavours found in the literature. Using the compilation of $M - M_b$ relations of Middleton et al. (2018), we construct an observationally motivated prior distribution in (α, β) by the following procedure. We make many random draws of the pair (α, β) uniformly from the ranges $\alpha \in [7.63, 8.63]$ and $\beta \in [0.79, 2.14]$, and consider the pair valid if the resulting $M - M_b$ line falls within the region enclosed by the compiled sample of relations in the range $10^6 M_\odot < M < 10^{10} M_\odot$. The resulting probability distribution $p(\alpha, \beta)$ is shown in Fig. 2. We then marginalize over the parameters (α, β) in the computation of $p(d|G_i)$ in equation (17). We

assume $\epsilon = 0.3$ throughout, which is the typical relation dispersion value reported in the literature.

Using these assumptions, equation (17) reduces to a 4D integral over M , q , α , and β . In the following, we show results where $M - M_b$ is always marginalized over (α , β) and we discuss the impact of assuming a specific scaling relation in Section 4.3. In practice, we transform variables to $^{10}\log(M)$ and $^{10}\log(q)$ and perform the numerical integration in log space for these parameters.

3 PRACTICAL IMPLEMENTATION

3.1 Source selection

To test our method, we simulate plausible future detections of single sources in PTA. We turn to work by Rosado et al. (2015) who have studied large-scale simulations of MBHB populations and the resulting GW signals that could be detected by PTAs. They construct 20 000 models (with different observed MBH mass functions, pair fractions, and MBH–galaxy relations) and drew several Monte Carlo realizations of each model to build realistic MBHB populations. They then considered the sensitivity of several PTAs as a function of time and used simple detection statistics to declare detection of either individual MBHBs or the overall stochastic background. Although they find that it is more likely that the background is detected first, eventually, individual sources can also be confidently identified. For each of the simulations, they record the properties of the first MBHB to be individually resolved by the PTA under consideration. Therefore, their procedure informs the likely parameters of the first resolvable MBHBs. We use it here to get the parameters for our test injections, as follows.

The signal-to-noise ratio (S/N) of a circular MBHB in an array of M pulsars can be written as

$$S/N = \left[\sum_{i=1}^M (S/N_i)^2 \right]^{1/2}, \quad (21)$$

where the S/N in the i -th pulsar is

$$(S/N_i)^2 = \frac{A^2}{4\pi^2 f^2 S_i} \mathcal{R}(\vec{\delta}). \quad (22)$$

Here, A is the GW amplitude given by equation (5) and f is the observed GW frequency. $\mathcal{R}(\vec{\delta})$ is a factor of the order of unity that depends on the geometry of the system – including source sky location and inclination, wave polarization angle, and pulsar sky location – and on the duration of the PTA observation T (see Rosado et al. 2015 for the full expression). S_i is the noise in the i -th pulsar that we consider to be of the form

$$S_i = 2\Delta t \sigma_i^2 + S_{h,\text{rest}}, \quad (23)$$

where the first term on the rhs is the rms noise level of the timing residuals and the second term is the level of confusion noise given by all other sources contributing to the overall GW signal.

To select suitable individual sources, we construct a mock version of the IPTA using the 49 pulsars of IPTA DR1 (Verbiest et al. 2016). We consider the actual sky location and rms noise σ_i of each pulsar and assume bi-weekly observations ($\Delta t = 2$ weeks) for a timespan of $T = 10$ yr. Next, we generate 50 realizations of a realistic population of circular, GW-driven MBHBs, based on one of the models presented in Sesana (2013). The number of realizations is chosen to produce a sample of individually resolvable sources that is large enough to give us freedom to pick sources in desired regions in the sky (see below). In particular we use a fairly optimistic

Table 1. Properties of the three test sources selected for this study.

Source	$\mathcal{M} [M_\odot]$	z	f [nHz]	A
A	3.18×10^9	0.62	7.44	0.96×10^{-15}
B	5.36×10^9	0.57	5.94	2.05×10^{-15}
C	3.69×10^9	0.18	5.18	2.40×10^{-15}

model resulting in a characteristic GW strain $h_c = A(f/\text{yr}^{-1})$ with $A \approx 1.3 \times 10^{15}$, which is just at the edge of the most recent PTA limits (Lentati et al. 2015; Shannon et al. 2015; Verbiest et al. 2016; Arzoumanian et al. 2018).

In each model realization, we select the loudest GW sources one-by-one and use all remaining MBHBs to consistently compute $S_{h,\text{rest}}$. All potentially resolvable GW sources had $S/N < 2$ in the adopted set-up. This is a good sanity check for our simulation; in fact it is expected that no observable sources result from this procedure, given that no single MBHB has been detected to date either. To increase the S/N , we suppress the noise by multiplying each rms residual σ_i by a fudge factor $\eta < 1$. After decreasing η to 0.2, we observe ≈ 30 sources (in 50 GW signal realizations) at $S/N \gtrsim 5$. We select three of those sources, which we name A, B, and C.

Relevant parameters of the selected sources are listed in Table 1 and their location in the sky, relative to the IPTA pulsars, can be seen in Fig. 4. We have intentionally picked three sources in areas of different IPTA pulsar density. Because the response functions depend on the angular distance between the pulsar and the GW propagation direction (equation 8), the localization behaviour is different for sources that are close to (good) pulsars than for those in relatively empty regions of the sky (see also Section 4.1). Parameters listed in Table 1 are consistent with distributions shown in Fig. 6 of Rosado et al. (2015). The first resolvable sources are likely to be at relatively low frequencies (few nHz) and can come from MBHBs at moderate redshifts (up to $z \approx 1$).

3.2 Source injection and likelihood evaluation

Each source is injected into a synthetic PTA based on IPTA data release 1 (Verbiest et al. 2016). The sky location and relative white noise level for each pulsar are kept the same as in IPTA DR1 (see their table 4 under Residual rms). Practical limitations on the method of Goldstein et al. (2018) mean the cadence and observation time of each pulsar has to be the same, so these are averaged over. We adjust the total observation time and/or reduce the noise in each pulsar by a constant factor to set the S/N of an injected source at the values 7, 10, 12, and 15 (see Table 2). We choose 7 as the smallest S/N value because it ensures a confident detection according to the \mathcal{F} statistic adopted by Rosado et al. (2015) (and used in this work). Assuming a typical PTA and a false alarm probability of 0.001, a source with $S/N = 7$ has a detection probability of ≈ 0.9 . For each set-up, a likelihood $\mathcal{L}(A, \theta, \phi)$ is obtained using three different realizations of random white noise in the null stream pipeline. Summarizing, we run a total of 36 simulations featuring:

- (i) *three* different sources: A, B, C;
- (ii) *four* values of detection $S/N = 7, 10, 12,$ and 15 ;
- (iii) *three* independent white noise realizations.

The likelihood is evaluated on a 3D grid in amplitude (A) and sky location (θ, ϕ). A is evenly sampled in log space, assuming a log flat prior between 10^{-17} and 10^{-14} . The location parameters θ (polar coordinate from $0-\pi$) and ϕ (azimuthal coordinate from

Table 2. Adjustments made to the simulated IPTA-like array in order to fix S/N of the three injected sources A, B, and C. The pulsar locations are kept the same as in IPTA DR1 (Verbiest et al. 2016), as are the relative white noise levels of each pulsar. For $S/N \geq 10$, the noise is decreased by a constant factor in all pulsars. The cadence ΔT and observation time T are averaged over for all pulsars. Then T is adjusted to set the S/N at specific values, keeping ΔT as close to the IPTA DR1 value as the Goldstein et al. (2018) method allows.

Source	Per cent rms		A		B		C	
	S/N	IPTA	T (yr)	ΔT (s) $\times 10^5$	T (yr)	ΔT (s) $\times 10^5$	T (yr)	ΔT (s) $\times 10^5$
7	100		12.8	2.12	10.7	2.03	12.2	2.76
10	80		21.3	2.71	16.0	2.33	12.2	2.11
12	80		29.8	2.64	26.7	2.69	18.4	2.20
15	70		34.0	2.52	32.0	2.70	24.5	2.54

$0-2\pi$) are sampled over using a grid of equal area pixels. This grid is constructed with the HEALpix algorithm (Górski et al. 2005) via healpy.² HEALpix allows the user to define a grid refinement parameter n , which results in a number of pixels $N_{\text{pix}} = 12n^2$. We choose $n = 32$, giving $N_{\text{pix}} = 12288$ pixels of approximately equal area of 3.36 deg^2 . For the likelihood calculation we use θ and ϕ at the middle point of each pixel.

The sky error-box Ω_{90} is determined as the (smallest) area in the sky containing 90 per cent of the total likelihood. For its practical computation, the likelihood is first marginalized over A , which gives $\mathcal{L}(\theta, \phi)$ at each sky location. Pixels are then ranked in an array $j = 1, \dots, N_{\text{pix}}$ in order of decreasing likelihood and their cumulative likelihood is calculated. Ω_{90} is then composed by the first K pixels (i.e. $j = 1, \dots, K$) enclosing 90 per cent of the total likelihood. For the sky area containing Ω_{90} , we implement the next level of HEALpix grid refinement ($n = 64$) that results in a smoother likelihood, evaluated on smaller pixels of 0.84 deg^2 .

3.3 Mock galaxy catalogue for host selection

Having determined $\mathcal{L}(A, \theta, \phi)$, we need to draw a set of properties of potential hosts from a realistic galaxy population. To this purpose, we use a mock realization of the observed sky extracted from the Millennium Run (Springel et al. 2005). The simulation evolves dark matter particles over a volume $(500/h \text{ Mpc})^3$, reconstructing the clustering of dark matter haloes. Semi-analytic galaxy formation models are then used to populate haloes with galaxies, tracking their star formation, accretion, and merger history.

Although not 'state of the art', the large volume of the Millennium Run (683.7 Mpc side; Springel et al. 2005), compared to more recent large-scale, fully hydrodynamical, simulations such as Illustris (105.6 Mpc side; Vogelsberger et al. 2014) and EAGLE (100 Mpc side; McAlpine et al. 2016), is relevant for our work. It ensures more statistical variation in the resulting galaxies, and in particular, a better sampling of the high mass tail of the distribution, which is where the best candidate galaxies reside. We use the simulated sky maps constructed by Henriques et al. (2012) that employ the semi-analytic model of Guo et al. (2011), which has been shown to reproduce a number of observed properties of galaxies, including luminosity function, morphology, and clustering.

The sky maps are flux-limited to $i < 21.0$ (see Henriques et al. 2012, for full details). This results in galaxy catalogues that are complete down to stellar masses of $\approx 10^{11} M_{\odot}$ at $z = 0.5$ and $\approx 4 \times 10^{11} M_{\odot}$ at $z = 1$. We will show in Section 4 that all credible hosts are above these completeness limits. We downloaded

all galaxies with stellar masses of $5 \times 10^{10} M_{\odot}$ and higher at $z \leq 1$, which resulted in about 50 million objects. For each galaxy we store the bulge mass M_b , coordinates in the sky (θ, ϕ) and apparent redshift z . The latter is then converted to D_1 by assuming our fiducial cosmology (flat Λ CDM with $h_0 = 0.73$, $\Omega_M = 0.25$). This information, together with a prior on the MBHB mass ratio q and the aforementioned assumptions for the $M - M_b$ relation, is all we need to perform the calculation outlined in Section 2.2.

To limit data size, only galaxies that fall within Ω_{90} are considered, which contain most of the relevant information. The simplifying assumption is made that one of the galaxies in Ω_{90} is the true source of the PTA signal, but there is a 10 per cent probability it falls outside the error-box. For each galaxy, the likelihood of being the GW source host is finally computed via equation (17), where A is determined by the injected sources and all relevant galaxy parameters are given by the mock catalogues and have prior distributions as described in Section 2.2.

4 RESULTS AND DISCUSSION

For each experimental set-up (injected source and S/N with three random noise realizations as in Section 3), we use the null stream pipeline to obtain $\mathcal{L}(A, \theta, \phi)$ and determine Ω_{90} , the results of which we discuss here first in Section 4.1. Then, we perform the calculation as described in Section 2.2 for each galaxy in Ω_{90} . This produces a population of $p(d|G_i)$ from which we can obtain a cumulative likelihood distribution. These results are shown in Section 4.2.

4.1 Sky localization

First we look at the behaviour of Ω_{90} with increasing S/N for the three different sources, which is shown in Fig. 3. The expected trend $\Omega_{90} \propto S/N^{-2}$ is roughly followed by all sources, albeit not perfectly, due to the small numbers of performed simulations for each case. An exception is source A at $S/N < 10$, which shows a much steeper slope. Although this is consistent with the 'transition zone' identified in Goldstein et al. (2018) – signalling the S/N at which the data start to be informative – sources B and C do not behave the same way.

We conjecture that this is related to the specific position of the sources, relative to the pulsars (see Fig. 4). When the source is close to the location of the best pulsars (like A), the combined S/N from all pulsars at the marginal detection level ($S/N \approx 7$) is mostly due to the contribution of these few, good pulsars (or possibly only one good pulsar). The other pulsars have a very low individual S/N . Therefore, the source is effectively triangulated by very few pulsars

²healpy.readthedocs.io

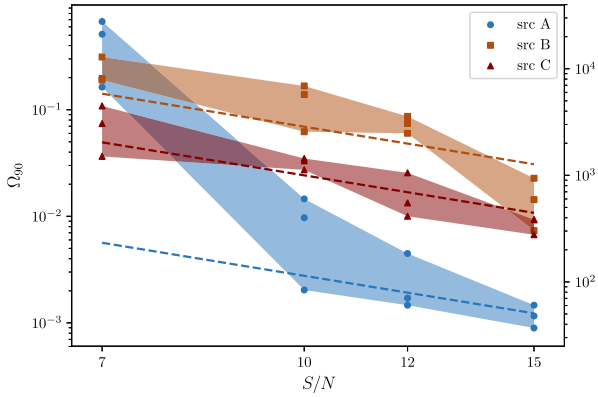


Figure 3. Sky-localization accuracy for the chosen sources A, B, and C at signal-to-noise ratios 7, 10, 12, and 15. At each S/N , a marker indicates Ω_{90} for each of three runs with different noise realizations. The dashed lines give the best fit of $\Omega_{90} \propto (S/N)^{-2}$ for the points at $S/N \geq 10$.

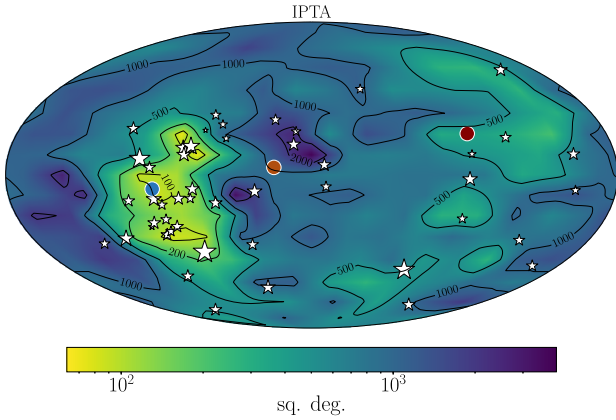


Figure 4. Localization capability of an IPTA-like array of pulsars for a source at fixed $S/N = 12$. This map is interpolated from 192 localization values obtained by injecting a source at 192 locations forming a grid of equal sky area pixels (a HEALPix grid with $n = 4$; Górski et al. 2005). The IPTA pulsars are marked with stars, where the size of the star corresponds to the noise level of the pulsar (with bigger stars for lower noise). The circles indicate the positions of sources A (blue, left), B (orange, middle), and C (red, right).

making localization poor. At higher total $S/N \approx 10$, more pulsars contribute to the triangulation as their individual S/N increases. As such, there is a steep improvement in sky-localization, steeper than the canonical $(S/N)^{-2}$ slope.

Conversely, when the source is far away from the majority of the best pulsars (like B and C), a detection with $S/N \approx 7$ already requires contribution from several different pulsars, making triangulation more effective. After this transition (the shaded area crossing in Fig. 3), the standard S/N scaling continues for source A as well.

Apart from the trend, the localization accuracy of the three sources vary by a factor of ~ 20 between them. This is due to both the inhomogeneous distribution of pulsars in the sky (Sesana & Vecchio 2010) as well as the different quality of the pulsars in the arrays (Babak et al. 2016), which is expected to cause a difference in localization. The best localization, at high S/N , is achieved for source A, sitting in the ‘sweet spot’ of the array (where most of the pulsars, including the best ones, are). However, there is not simply a monotonic increase of Ω_{90} for sources further away, since the

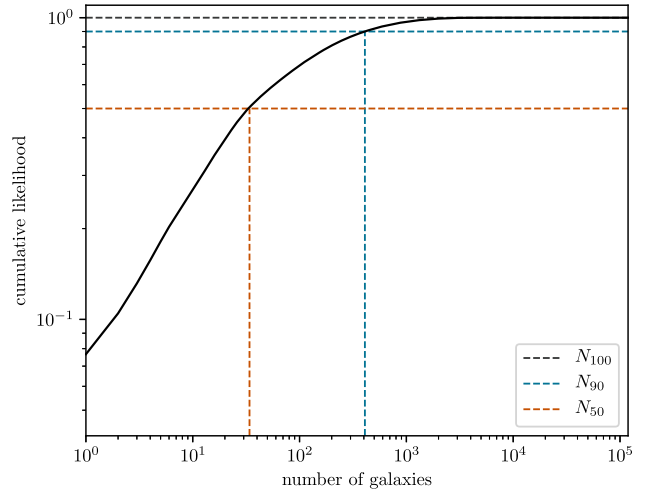


Figure 5. Cumulative likelihood of $p(d|G_i)$. The likelihood data d is for the IPTA set-up (as described in the text) with source A at $S/N = 15$ (one of the random noise realizations). Vertical dashed lines identify the number of galaxies making up 50 per cent (orange) and 90 per cent (blue) of the total likelihood.

furthest source C has a better localization than source B. This is also expected since, due to the shape of the PTA response function, sources that are antipodal to the sky region that is best covered by the array are better localized than sources that are orthogonal to that region (see e.g. fig. 10 in Sesana & Vecchio 2010).

A further investigation of this is visualized in Fig. 4. Here we inject a source with the same parameters as A at 192 different locations in the sky into white noise, using a synthetic IPTA-like array. The S/N is set to 12 everywhere, by scaling the amplitude A of the GW signal. The map shows the resulting localization Ω_{90} at each point. A dipolar structure of Ω_{90} is notable, where sources near the ‘sweet spot’ of clustered pulsars – which includes most of the best pulsars – and to a lesser extent, sources near the antipodal point are localized better than sources in between. This is related to the quadrupolar nature of GWs, which results in a pulsar response function that has this antipodal symmetry, as was also shown by Sesana & Vecchio (2010).

In any case, the huge scatter in Ω_{90} warns of a potential risk of an anisotropic sky coverage of the pulsars in the array. Should the loudest resolvable GW sources be positioned at unfavourable locations, their detection, even at moderate $S/N \approx 12$, would allow sky-localization accuracies of about 2000 deg^2 only (an area containing ~ 2 million galaxies in our catalogue before any selection), jeopardizing any effort to identify a possible EM counterpart.

4.2 Host candidate population

4.2.1 Number of credible host candidates

Our main results consist of a set of $p(d|G_i)$ for the galaxies $\{G_i\}$ within Ω_{90} for each experimental set-up (Section 3). First, we compute the cumulative likelihood distribution from these $p(d|G_i)$. We then define N_x to be the minimum number of galaxies needed to sum to x per cent of the total likelihood $\sum_i p(d|G_i)$. Specifically, we look at N_{50} and N_{90} as proxies for the expected number of candidate host galaxies.

An example can be seen in Fig. 5 for source A at $S/N = 15$ (the first random noise realization). Within $\Omega_{90} \approx 60 \text{ deg}^2$, there

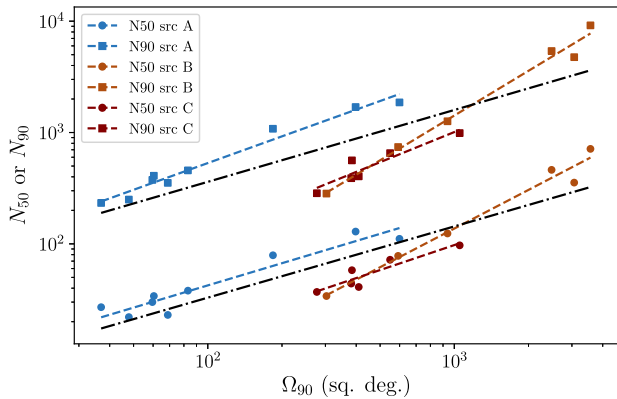


Figure 6. Number of candidate galaxies adding up to 50 per cent (N_{50} , circle markers) or 90 per cent (N_{90} , square markers) of the total likelihood to host the detected source, versus the sky-localization accuracy Ω_{90} for that detection. Results are shown for source A (blue) at $S/N = 10, 12,$ and $15,$ and sources B (orange) and C (red) at $S/N = 12$ and $15.$ For each S/N three noise realizations give a cluster of points at similar Ω_{90} values. The dashed lines show fitted power laws per source (see Table 3 for the best-fitting powers). The dot-dashed lines are fits to N_{50} for all sources, with power 0.64, and to N_{90} , with power 0.65.

Table 3. Best-fitting powers for the power-law fits to N_{50} and N_{90} as in Fig. 6. These are obtained by minimizing the sum of squared errors on the log N_x values.

Source	N_{50} power	N_{90} power
A	0.66	0.80
B	1.15	1.34
C	0.74	0.90
All	0.64	0.65

are $\sim 1.2 \times 10^5$ galaxies $\{G_i\}$ in our mock catalogue that would make detailed follow-ups for host identification impractical. The potential benefit of our technique is apparent from the fact that of those galaxies, only $N_{90} = 409$ make up 90 per cent of $p(d|G_i)$, and $N_{50} = 34$ make up 50 per cent of $p(d|G_i)$.

The collection of N_{50} and N_{90} of all experimental cases for which we obtained results can be found in Fig. 6. We can fit a power law as $N_x = c(\Omega_{90}/\Omega_{90}^*)^p$, with parameters c , p (the power) and Ω_{90}^* , and x being either 50 or 90. By minimizing the sum of squared differences between the predicted log values and the log of the data points, we obtain best-fitting powers 0.64 and 0.65, for N_{50} and N_{90} , respectively. Although naively one would expect a linear proportionality between Ω_{90} and the number of potential hosts, there is a significant scatter on the relation.

Tighter fits are obtained by treating the points for different injected sources separately, with best-fitting powers as in Table 3. These numbers show that fits to individual source data points are generally steeper and closer to the expected linear dependence. One of the causes of the shallower global fit appears to be the larger N_{50} and N_{90} for source A with respect to sources B and C at sky-localizations of $\approx 300 \text{ deg}^2$, as shown in Fig. 6. (Source A has $S/N = 10$ around this localization accuracy, while source B and C have higher $S/N = 15$. Consequently, N_{90} and N_{50} for source A includes galaxies with a lower bulge mass than for B and C, resulting in a larger N_{90} and N_{50}).

So while there is clearly a relation between the size of the sky error-box and the number of candidate host galaxies, scatter

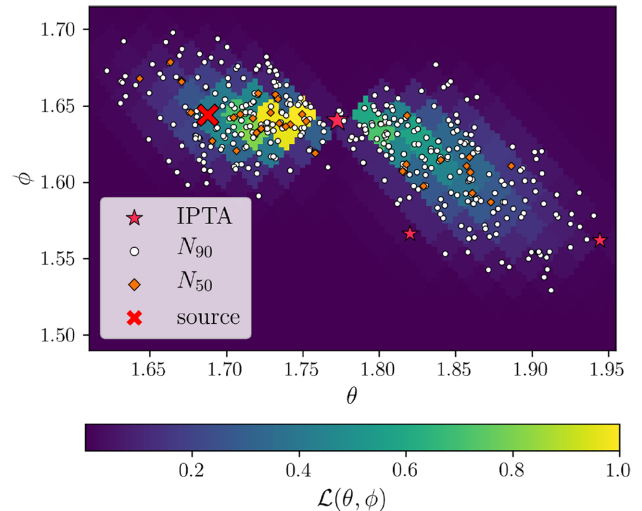


Figure 7. Locations of the best candidate host galaxies on top of the sky location likelihood for the injected source A (located at the red cross). The PTA has pulsar locations (pink stars) and relative noise levels of the IPTA DR1, but is adjusted such that the total $S/N = 15.0$ (see the text). The 34 best candidates sum to 50 per cent of the likelihood to be the host galaxy (N_{50} in the orange diamonds) and an additional 375 sum to 90 per cent (N_{90} in the white circles). For this example, the $M_{\text{BH}} - M_{\text{bulge}}$ relation is marginalized over priors obtained from the literature (see the text).

is caused by factors related to the detailed source properties. Nonetheless, as a rule of thumb, we expect that for a resolvable PTA signal located in the sky with a precision of $\approx 100 \text{ deg}^2$, we can identify few hundreds (few tens) galaxies in which the source sits with 90 per cent (50 per cent) confidence. Compared to all galaxies with stellar mass $> 5 \times 10^{11} M_{\odot}$ at $z < 1$ falling in the error-box, these numbers restrict the pool of realistic hosts by nearly three (four) orders of magnitude, making realistic detailed follow-up campaigns feasible.

We also calculated $p(d|G_i)$ and N_x for source A at $S/N = 7$, which has a very poor sky-localization of about $2.8 \times 10^4 \text{ deg}^2$ (67 per cent of the sky). The expected number of candidate hosts becomes very large, and also disobeys the trend discussed above. We conjecture that this is due to the localization likelihood distribution not having a single peak for the low S/N case, so potential hosts are allowed to be anywhere in the localization error-box, which is most of the sky.

4.2.2 Host candidate sky distribution and clustering

Apart from the number of galaxies that make up a significant fraction of the likelihood $\sum_i p(d|G_i)$, we can also look at the properties of these galaxies. The parameters from the mock galaxy catalogue are M_b , D_1 , θ , and ϕ . First, the sky locations of galaxies within N_{50} or N_{90} for the example case (source A at $S/N = 15$) are shown in Fig. 7. They are plotted on top of the localization likelihood $\mathcal{L}(\theta, \phi)$ of the injected source. The galaxies follow the shape of the localization area because we only used galaxies within Ω_{90} . Moreover, it can be seen that there is a relatively higher concentration of N_{50} galaxies in the highest likelihood pixels. Hence, $\mathcal{L}(\theta, \phi)$ must contribute more to the selection of candidates than simply what we get from selecting the ones in Ω_{90} .

We further investigate this statement using the clustering of good candidate galaxies – the N_{50} galaxies – for all of the experimental cases. Fig. 8 simultaneously shows a measure of the concentration of the localization likelihood $\mathcal{L}(\theta, \phi)$ and of the concentration of

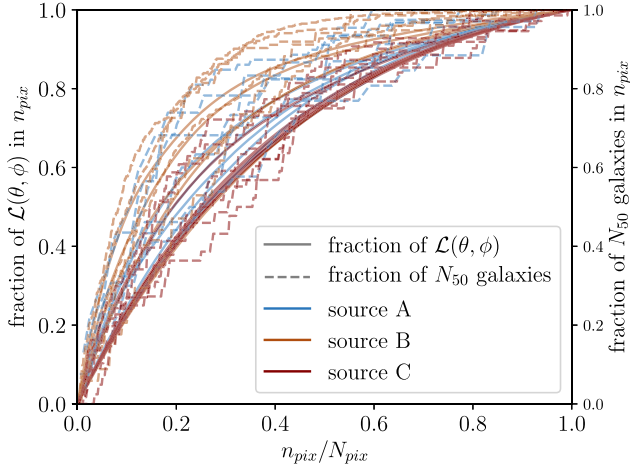


Figure 8. Comparison between the concentration of the sky-localization likelihood and of the locations of good candidate host galaxies. With the fractional area of Ω_{90} on the x -axis, the fractional localization likelihood in this area on the first y -axis (left, solid lines), and the fractional number of N_{50} galaxies on the second y -axis (right, dashed lines) (see the text for details.) The quantities are normalized between 0 and 1 so that all experimental cases fit on the same scale. This plot includes all three injected sources (A in blue, B in orange, and C in red), for $S/N = 12, 15$.

N_{50} galaxies. The sky error-box Ω_{90} consists of a number of pixels N_{pix} that are sorted in descending $\mathcal{L}(\text{pix}_i)$ order. Starting with the best pixel, we iteratively increase this number by adding the next best pixel. The size of the included area is recorded as the fraction of the number of pixels over the total in Ω_{90} , i.e. $n_{\text{pix}}/N_{\text{pix}}$. The concentration of the localization likelihood then is the likelihood in n_{pix} as a fraction of the total, i.e. $\sum_i^n \mathcal{L}(\text{pix}_i) / \sum_i^N \mathcal{L}(\text{pix}_i)$.

We compare this with the concentration of good candidate hosts, as the fractional number of N_{50} galaxies in the selected pixels. The distributions are spread out, but there is no significant difference between the sky likelihood and candidate host concentration, i.e. the host probability follows the sky-localization distribution. We therefore conclude that it is valuable to include detailed sky-localization information when selecting candidate host galaxies, rather than only making a selection based on the total sky-localization area.

4.2.3 Host candidate mass and redshift

Secondly, we consider the other two parameters from the catalogue, the bulge mass M_b , and luminosity distance D_1 . Fig. 9 shows their distribution among candidate hosts for the example case, where D_1 has been converted into redshift. This figure best visualizes the key idea behind our method. Since $A \propto M^{5/3}/D_1$ – and there is a proportionality $M \propto M_b^\beta$ and an almost linear proportionality between D_1 and z at $z < 1$ – there is only a stripe in the mass–redshift plane defining the region of possible galaxy hosts. Moreover, since the first detection of a resolved PTA source will necessarily involve a very strong signal from a very massive binary system, this region lies at the highest masses. Due to the steep decay of the high mass end of the galaxy mass function, only few credible host candidates can be identified.

In the example shown, galaxies belonging to N_{50} or N_{90} are bound by a line of slope $3/(5\beta)$ in the $\log D_1(z) - \log M_b$ plane (where β is the $M - M_b$ constant marginalized over our prior), as expected by the GW amplitude scaling. There is however a large mixing of galaxies with different likelihoods in this plane due to their specific

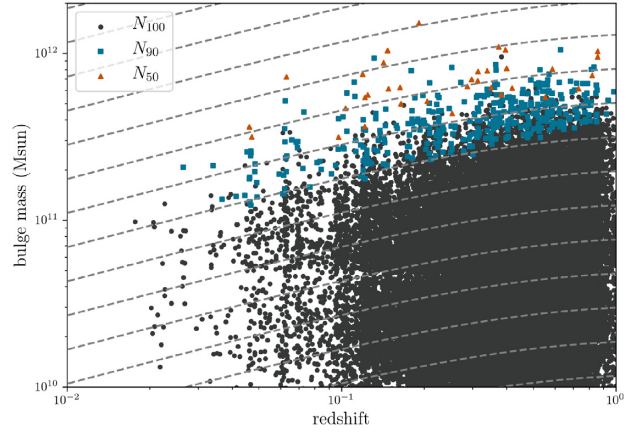


Figure 9. Distribution of bulge masses and redshifts of the candidate host galaxies of the example case source A with $S/N = 15$. The blue squares mark galaxies that make up N_{90} and the orange triangles mark the candidates which make up N_{50} . All other galaxies that fall within the sky-localization error-box Ω_{90} , but form the lowest 10 per cent of the total likelihood $\sum_i p(d|G_i)$, are marked with (dark) grey circles. The dashed grey lines are lines of constant GW amplitude (as in equation 5).

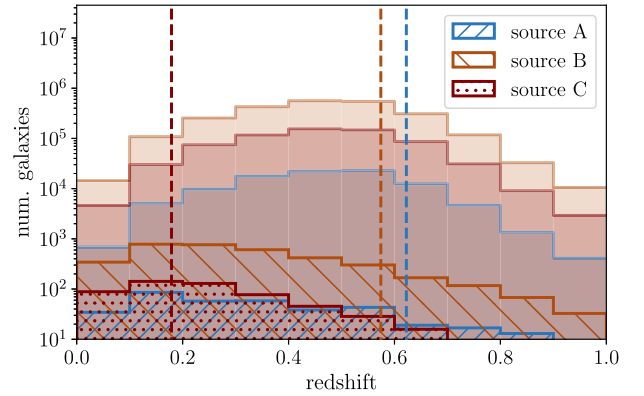


Figure 10. Logarithmic histogram of the redshifts of candidate host galaxies per source. The counts from the six experimental cases with $S/N = 12, 15$, and 3 noise realizations are averaged over. The foreground (hatched) histograms are N_{90} candidates, and the background (filled) histograms are all (i.e. N_{100}) candidates from the selected sky error-box. Injected redshift values for each source are indicated by a dashed line (see also Table 1).

sky location. For example, there are a few very massive galaxies that fall into the lowest 10 per cent of the likelihood, which is due to an unfavoured sky position. Note that there are N_{50} candidates across the whole range of redshifts in our sample.

The redshift of the injected source A is 0.62 (Table 1), so it is not a surprise that candidate hosts for this source have redshifts across the whole range 0–1. To explore this further, we look at the redshifts of candidate host galaxies for all injected sources and S/N values. Fig. 10 shows a number of histograms of z on a logarithmic scale. For each source A, B, and C, the results of $S/N = 12$ and 15 (with three noise realizations per S/N) are combined. We make a comparison between the redshift distribution of the candidate galaxies pre-selected within the sky error-box (the background histograms), and the N_{90} candidates selected with our method (the foreground, hatched histograms).

Compared to the prior distribution, lower redshifts are preferred. However, for all injected sources, there are a significant number of

candidates at redshifts $z > 0.6$. Even though the injected redshift for source C ($z = 0.18$) is much lower than four sources A and B, the redshift distributions of candidate hosts differ only slightly, which reflects the fact that redshift is degenerate with mass in our method.

The turnover in the total number of galaxies in the error-box seen in Fig. 10 at $z > 0.5$ is due to the $i = 21$ flux limit of the adopted galaxy catalogues, which result in severe incompleteness of lower-mass galaxies. Fig. 9, however, shows that typical galaxies belonging to N_{90} have $M_b > 2 \times 10^{11}$ at $z = 0.5$ and $M_b > 4 \times 10^{11}$ at $z = 1$. The adopted catalogue is therefore complete in the mass–redshift range where potential GW galaxy hosts live.

Fig. 9 also shows that the distribution of credible galaxy hosts of resolvable PTA sources peaks at $z \lesssim 0.2$, whereas two out of three of our selected signals (A and B) lie around $z \approx 0.6$. GW sources were picked according to their sky location, therefore A, B, and C are not an unbiased sample and are not necessarily representative of the actual redshift distribution of the first GW resolved signals. However, there are several systems at $z > 0.5$ in the sample of 30 resolvable sources found in Section 3, and also Rosado et al. (2015) found that the peak of the first resolved PTA sources is at $z \approx 0.5$.

High z sources are more common despite there being less potential host galaxies at such redshifts. This indicates that the likelihood of a galaxy to be a host is not only connected to its sky location and its position in the M – z plane, considered in this work. The other key parameter is likely to be the absolute galaxy mass (regardless of redshift). There is evidence – both from observations and from cosmological simulations (see e.g. results compiled in fig. 1 of Sesana et al. 2016) – that the galaxy merger rate at low redshift is a strong function of the galaxy mass, with massive galaxies merging more often.

Since the MBHB population simulated in Section 3 consistently takes this fact into account, the resulting MBHB population is naturally skewed towards high masses. Conversely, our host selection method only picks galaxies based on the GW amplitude given the combination of redshift and bulge mass, and therefore chooses relatively more lighter galaxies. However, because of these candidates’ lower masses, they are less likely to have undergone a major merger (and hence host a GW source) compared to the few more massive ones picked at higher redshift. This suggests that combining our method with a (prior) probability of hosting an MBHB based on galaxy mass only (Rosado & Sesana 2014; Mingarelli et al. 2017) can somewhat break the intrinsic mass–redshift degeneracy, further reducing the numbers of credible galaxy hosts.

4.3 Assumptions and approximations

Although simulations performed in this work are realistic in many aspects, few assumptions and choices had to be made to make their total runtime manageable.

Several assumptions were made in the connection of the chirp mass of the GW source to the bulge mass of the host galaxy. First, we assumed a log-flat prior on $-2 \leq \log q \leq 0$, based on the broad q distribution of merging binaries found in cosmological simulations (Kelley, Blecha & Hernquist 2017). Although this is not necessarily representative of the q distribution of real MBHBs, we tested that different choices have only a minor impact on the results (see also Holgado et al. 2018; Inayoshi, Ichikawa & Haiman 2018; Sesana et al. 2018).

Secondly, we did not consider errors in the measurements of galaxy M_b and D_l . The latter does not matter; for any practical purposes, galaxy redshifts can be determined almost exactly, and

estimates of D_l are only affected by galaxy-peculiar velocities and uncertainties in the knowledge of the cosmological parameters, resulting in a negligible few per cent error. Conversely, the former can be significant, as bulge mass determination can be uncertain within a factor of two. This is likely to impact our results, spreading the host probability distribution thus returning more host candidates. Some tests on a limited number of set-ups found that including an uncertainty of a factor of two on the galaxy bulge mass results in roughly a factor of two more candidate hosts galaxies.

Last, we marginalized over the uncertainty in the $M - M_b$ relation. Assuming a specific $M - M_b$ relation instead can affect our results, especially if the relations predict relatively higher or lower black hole masses than the marginalized relation. As an example, we ran some test cases assuming the $M - M_b$ relation from Kormendy & Ho (2013), which associates relatively higher black hole masses given the galaxy bulge mass. The number of candidate host galaxies in these cases is increased by a factor ranging between ~ 3 and ~ 8 with respect to the marginalized $M - M_b$ case. Conversely, for a ‘pessimistic’ $M - M_b$ relation such as Shankar et al. (2016) – which predicts relatively lower black hole masses especially for high-mass galaxies – the number of candidates is a factor ~ 2 to ~ 4 lower.

Due to computational limitations, we ran a limited number of simulations. Although we checked robustness of the results against the specific noise realization, we only picked one sky location for each source. This may make cosmic variance a factor in the determination of the number of galaxy hosts. To test this, for a selected GW source, we performed some rigid rotations of the Millennium sky and counted N_{50} and N_{90} for each of them. Although numbers vary, the scattering is consistent with that observed in Fig. 6.

An important assumption of our method is that the true host of the detected GW signal is present in the galaxy catalogue. This is guaranteed only for complete catalogues. Real catalogues based on observations never are, and the simulated catalogue from Henriques et al. (2012) reflects this by selecting galaxies based on observational criteria. This results in a number of missing galaxies – more towards higher redshifts. However, for the most part these are the small galaxies (which are more difficult to observe) and are not relevant host candidates. Since at redshifts $z \lesssim 1$ only the most massive galaxies are selected in N_{90} (see Fig. 9); this is unlikely to affect the results for N_{90} and N_{50} , but it is a possible source of error. As there are good candidates up to $z = 1$, it is also possible that there are a small number of potential N_{90} galaxies at $z > 1$ that were not included.

Finally, it should be kept in mind that we selected the 90 per cent sky location credible region. By selecting N_{50} and N_{90} in this region, the actual probability to find the true host in these sets is $0.9 \times 0.5 = 0.45$ and $0.9 \times 0.9 = 0.81$, respectively.

5 CONCLUSION

In this paper, we proposed a novel methodology to select host galaxy candidates of the first individual GW sources observed by pulsar timing arrays. Since PTA source localization is expected to be of several deg^2 at best, up to several million galaxies might end up in the sky error-box. Classifying the most promising host candidates is therefore of paramount importance to increase the chances of true host identification via dedicated follow-ups. Our method exploits the GW strength dependence on chirp mass and distance, together with empirical MBH mass–host galaxy correlation, to rank galaxies in the mass–redshift plane. We frame this concept in the Bayesian language, together with the null-stream-based sky-

localization method developed in Goldstein et al. (2018), to assign each galaxy a probability of hosting the MBHB generating a specific GW signal.

To test our method, we performed realistic simulations by drawing GW sources from detailed MBHB population models based on observed merging galaxies, by employing the actual IPTA pulsar sky locations and rms values to build the array, and by selecting host candidates based on formation and evolution models. We considered different GW source sky positions and detection S/N and investigated the ensemble of credible host galaxy candidates. In particular, we defined N_{50} and N_{90} to be the smallest numbers of galaxies having a collective 50 and 90 per cent chance of being the true host of the GW source, respectively, assuming the true host is among the prior selection of candidates. Our key results can be summarized as follows:

- (i) N_{50} and N_{90} are, respectively, nearly four and three orders of magnitude smaller than the number of galaxies with stellar mass $M_* > 5 \times 10^{10} M_\odot$ at $z < 1$ found in the 90 per cent confidence sky location region Ω_{90} ;
- (ii) N_{50} and N_{90} should roughly be proportional to Ω_{90} . We find a sublinear proportionality, although with large scatter;
- (iii) despite the large scatter, a useful rule of thumb is that for $\Omega_{90} = 100 \text{ deg}^2$, $N_{50} \lesssim 50$ and $N_{90} \lesssim 500$;
- (iv) although the distribution of potential hosts peaks around $z < 0.2$, it has a long tail that extends up to $z \lesssim 1$.

Our methodology can therefore effectively select the most likely host galaxy candidates, which might have a major impact on future multimessenger observations of MBHBs. For typical PTA sky-localization precision of hundreds of deg^2 , instead of following up millions of galaxies, we can choose to accept the risk of missing the true host with 55 per cent (19 per cent) probability and monitor only the ≈ 100 (1000) most promising ones. There is significant uncertainty on these numbers, mainly due to the uncertainty in the $M - M_b$ relation (see Section 4.3).

The applicability of our method obviously relies on the availability of photometric and spectroscopic data from all-sky surveys necessary to identify potential galaxy hosts and to estimate their stellar (and bulge) masses. Since the most credible galaxy candidates are necessarily very massive (and/or particularly nearby), relatively shallow surveys are sufficient for this scope. Catalogues from SDSS (Alam et al. 2015; covering $\approx 1/4$ of the sky), Pan-STARRS (Kaiser et al. 2002; $\approx 3/4$ of the sky), LSST (LSST Science Collaboration 2009; $\approx 1/2$ of the sky), and *Gaia* (Gaia Collaboration 2016, all sky) will provide enough imaging, photometric, and (possibly) spectroscopic information for reliable mass estimates via, e.g. spectral energy distribution fitting (see e.g. Longhetti & Saracco 2009; Duncan et al. 2014).

Note that a positive host identification chance increase of less than a factor of two comes at the expense of following up a factor of 10 more galaxies. The follow-up strategy can therefore be optimized based on the future number of resolved PTA sources and on available observing facilities. Reducing the number of credible host is critical mostly because our knowledge of MBHB signatures is poor (see e.g. Dotti, Sesana & Decarli 2012). One therefore has to collect all possible hints to build up confidence that the true host have been found. This might require, for example, multiple photometric and spectroscopic optical and IR follow-up of the candidates to unveil any observational hint of an accreting MBHB, deep-field imaging to assess the presence of post merger features such as stellar tails and shells (e.g. Lotz et al. 2008), integral field spectroscopy to identify the presence of a ‘dry’ MBHB via kinematic signatures in the stellar

distribution (Meiron & Laor 2013), deep X-ray observations to unveil the presence of an obscured AGN and its possible high-energy signatures (Koss et al. 2018), and many more.

The upcoming ELT (Gilmozzi & Spyromilio 2007) and *JWST* (Gardner et al. 2006) will be particularly suited for the optical and near-infrared follow-ups mentioned above, whereas the X-ray satellite *Athena* (Nandra et al. 2013) can potentially survey the 100 most probable hosts within less than 1 d of observation time. Clearly, the fewer the candidates, the more extensive the follow-up campaign can be, thus enhancing the chances of a positive detection. Archival data can also be used to identify hints of, e.g. periodic variability matching the frequency of the GW source. This can be done in the optical and, possibly, in X-ray with LSST and eROSITA (Merloni et al. 2012) archival data, respectively.

Finally, the mismatch between the credible host redshift distribution identified with our method and the expected distribution of the first PTA sources predicted by Rosado et al. (2015) indicates that a more efficient galaxy host selection can be performed when the mass-dependent galaxy merger probability is folded into the calculation (see also Mingarelli et al. 2017). By doing so, the mass-redshift degeneracy intrinsic in our method might be alleviated, further decreasing the number of credible hosts. We plan to further pursue this line of investigation in future work.

ACKNOWLEDGEMENTS

We thank A. Vecchio for useful comments. AS is supported by the Royal Society. JV is supported by the Science and Technology Funding Council grant ST/K005014/1. AMH is supported by the Department Of Energy, National Nuclear Security Administration, Steward Science Graduate Fellowship under grant number DE-NA0003864. The methods for this work are implemented using the Python programming language,³ and make extensive use of the NumPy/SciPy library (Jones et al. 2001; van der Walt, Colbert & Varoquaux 2011).

REFERENCES

- Abbott B. P. et al., 2017a, *Phys. Rev. Lett.*, 119, 161101
- Abbott B. P. et al., 2017b, *Nature*, 551, 85
- Abbott B. P. et al., 2017c, *ApJ*, 848, L12
- Abbott B. P. et al., 2017d, *ApJ*, 848, L13
- Alam S. et al., 2015, *ApJS*, 219, 12
- Amaro-Seoane P. et al., 2017, preprint (arXiv:1702.00786)
- Anholm M., Ballmer S., Creighton J. D. E., Price L. R., Siemens X., 2009, *Phys. Rev. D*, 79, 084030
- Arzoumanian Z. et al., 2018, *ApJ*, 859, 47
- Babak S., Sesana A., 2012, *Phys. Rev. D*, 85, 044034
- Babak S. et al., 2016, *MNRAS*, 455, 1665
- Boyle L., Pen U.-L., 2012, *Phys. Rev. D*, 86, 124028
- Burke-Spolaor S., 2013, *Class. Quantum Gravity*, 30, 224013
- Burt B. J., Lommen A. N., Finn L. S., 2011, *ApJ*, 730, 17
- Chornock R. et al., 2017, *ApJ*, 848, L19
- Dotti M., Sesana A., Decarli R., 2012, *Adv. Astron.*, 2012, 940568
- Duncan K. et al., 2014, *MNRAS*, 444, 2960
- Fishbach M., Gray R., Hernandez I. M., Magaña I., Qi H., Sur A., 2018, LIGO Scientific Collaboration, Virgo Collaboration
- Gaia Collaboration, 2016, *A&A*, 595, A1
- Gardner J. P. et al., 2006, *Space Sci. Rev.*, 123, 485
- Gilmozzi R., Spyromilio J., 2007, *Messenger*, 127, 11
- Goldstein J. M., Veitch J., Sesana A., Vecchio A., 2018, *MNRAS*, 477, 5447

³www.python.org

- Górski K. M., Hivon E., Banday A. J., Wandelt B. D., Hansen F. K., Reinecke M., Bartelmann M., 2005, *ApJ*, 622, 759
- Guo Q. et al., 2011, *MNRAS*, 413, 101
- Haiman Z., 2017, *Phys. Rev. D*, 96, 23004
- Hazboun J. S., Larson S. L., 2016, preprint([arXiv:1607.03459](https://arxiv.org/abs/1607.03459))
- Henriques B. M. B., White S. D. M., Lemson G., Thomas P. A., Guo Q., Marleau G.-D., Overzier R. A., 2012, *MNRAS*, 421, 2904
- Holgado A. M., Sesana A., Sandrinelli A., Covino S., Treves A., Liu X., Ricker P., 2018, *MNRAS*, 481, L74
- Inayoshi K., Ichikawa K., Haiman Z., 2018, *ApJ*, 863, L36
- Jaranowski P., Królak A., Schutz B. F., 1998, *Phys. Rev. D*, 58, 063001
- Jones E., et al., 2001, SciPy: Open source scientific tools for Python. Available at: <http://www.scipy.org/>
- Kaiser N. et al., 2002, in Tyson J. A., Wolff S., eds, Proc. SPIE Conf. Ser. Vol. 4836, Survey and Other Telescope Technologies and Discoveries. SPIE, Bellingham. p. 154
- Kelley L. Z., Blecha L., Hernquist L., 2017, *MNRAS*, 464, 3131
- Kelley L. Z., Blecha L., Hernquist L., Sesana A., Taylor S. R., 2018, *MNRAS*, 477, 964
- Klein A. et al., 2016, *Phys. Rev. D*, 93, 024003
- Kormendy J., Ho L. C., 2013, *ARA&A*, 51, 511
- Koss M. J. et al., 2018, *Nature*, 563, 214
- Lam M. T., 2018, *ApJ*, 868, 33
- Lentati L. et al., 2015, *MNRAS*, 453, 2576
- Longhetti M., Saracco P., 2009, *MNRAS*, 394, 774
- Lotz J. M., Jonsson P., Cox T. J., Primack J. R., 2008, *MNRAS*, 391, 1137
- LSST Science Collaboration et al., 2009, LSST Science Book, Version 2.0, preprint ([arXiv:0912.0201](https://arxiv.org/abs/0912.0201))
- McAlpine S. et al., 2016, *Astron. Comput.*, 15, 72
- Meiron Y., Laor A., 2013, *MNRAS*, 433, 2502
- Merloni A. et al., 2012, the German eROSITA Consortium, preprint ([arXiv:1209.3114](https://arxiv.org/abs/1209.3114))
- Middleton H., Chen S., Del Pozzo W., Sesana A., Vecchio A., 2018, *Nat. Commun.*, 9, 573
- Mingarelli C. M. F. et al., 2017, *Nat. Astron.*, 1, 886
- Nandra K. et al., 2013, preprint ([arXiv:1306.2307](https://arxiv.org/abs/1306.2307))
- Phinney E. S., 2001, preprint ([arXiv:astro-ph/0108028](https://arxiv.org/abs/astro-ph/0108028))
- Planck Collaboration XIII, 2016, *A&A*, 594, A13
- Ravi V., Wyithe J. S. B., Hobbs G., Shannon R. M., Manchester R. N., Yardley D. R. B., Keith M. J., 2012, *ApJ*, 761, 84
- Ravi V., Wyithe J. S. B., Shannon R. M., Hobbs G., 2015, *MNRAS*, 447, 2772
- Rosado P. A., Sesana A., 2014, *MNRAS*, 439, 3986
- Rosado P. A., Sesana A., Gair J., 2015, *MNRAS*, 451, 2417
- Schutz B. F., 1986, *Nature*, 323, 310
- Sesana A., 2013, *MNRAS*, 433, L1
- Sesana A., Vecchio A., 2010, *Phys. Rev. D*, 81, 104008
- Sesana A., Vecchio A., Colacino C. N., 2008, *MNRAS*, 390, 192
- Sesana A., Vecchio A., Volonteri M., 2009, *MNRAS*, 394, 2255
- Sesana A., Roedig C., Reynolds M. T., Dotti M., 2012, *MNRAS*, 420, 860
- Sesana A., Shankar F., Bernardi M., Sheth R. K., 2016, *MNRAS*, 463, L6
- Sesana A., Haiman Z., Kocsis B., Kelley L. Z., 2018, *ApJ*, 856, 42
- Shankar F. et al., 2016, *MNRAS*, 460, 3119
- Shannon R. M. et al., 2015, *Science*, 349, 1522
- Shapiro C., Bacon D. J., Hendry M., Hoyle B., 2010, *MNRAS*, 404, 858
- Simon J., Polin A., Lommen A., Stappers B., Finn L. S., Jenet F. A., Christy B., 2014, *ApJ*, 784, 60
- Springel V. et al., 2005, *Nature*, 435, 629
- Tanaka T., Menou K., Haiman Z., 2012, *MNRAS*, 420, 705
- Tang Y., MacFadyen A., Haiman Z., 2017, *MNRAS*, 469, 4258
- van der Walt S., Colbert S. C., Varoquaux G., 2011, *Comput. Sci. Eng.*, 13, 22
- Verbiest J. P. W. et al., 2016, *MNRAS*, 458, 1267
- Vogelsberger M. et al., 2014, *Nature*, 509, 177
- Zhu X.-J. et al., 2015, *MNRAS*, 449, 1650

This paper has been typeset from a $\text{\TeX}/\text{\LaTeX}$ file prepared by the author.