# Costs of Bayesian Parameter Estimation in Third-Generation Gravitational Wave Detectors: a Review of Acceleration Methods

Qian Hu[1, *] and John Veitch[1, †]

[1]*Institute for Gravitational Research, School of Physics and Astronomy,*
*University of Glasgow, Glasgow, G12 8QQ, United Kingdom*
(Dated: March 7, 2025)

Bayesian inference with stochastic sampling has been widely used to obtain the properties of gravitational wave (GW) sources. Although computationally intensive, its cost remains manageable for current second-generation GW detectors because of the relatively low event rate and signal-to-noise ratio (SNR). The third-generation (3G) GW detectors are expected to detect hundreds of thousands of compact binary coalescence events every year with substantially higher SNR and longer signal duration, presenting significant computational challenges. In this study, we systematically evaluate the computational costs of source parameter estimation (PE) in the 3G era by modeling the PE time cost as a function of SNR and signal duration. We examine the standard PE method alongside acceleration methods including relative binning, multibanding, and reduced order quadrature. We predict that PE for a one-month-observation catalog with 3G detectors could require billions to quadrillions of CPU core hours with the standard PE method, whereas acceleration techniques can reduce this demand to millions of core hours. These findings highlight the necessity for more efficient PE methods to enable cost-effective and environmentally sustainable data analysis for 3G detectors. In addition, we assess the accuracy of accelerated PE methods, emphasizing the need for careful treatment in high-SNR scenarios.

## I. INTRODUCTION

Since the first direct detection of gravitational waves (GWs) by the advanced LIGO [1] in 2015 [2], nearly 100 GW events have been published by the LIGO-Virgo-KAGRA (LVK) collaboration [3–7], with hundreds more expected from the current observing run. All of these GW events are identified as compact binary coalescences (CBCs) - mergers of compact objects such as black holes, neutron stars, and potentially exotic stars. Determining the properties of these sources is essential for various aspects of GW astronomy, including cosmology [8–10], astrophysical population [11–13], the equation of state of dense matter [14], stochastic background [15–17], among others.

CBC signals are characterized by the physical and geometric parameters of their sources, and the properties of the sources are inferred by estimating these parameters. A typical CBC system is described by at least 15 parameters, including two component masses, six component spin parameters, and seven extrinsic parameters describing the 3D sky location, inclination angle and polarization angle, and coalescence phase and time of the source. Additional parameters are required to account for specific phenomena, such as tidal deformability parameters for neutron stars [18], orbital eccentricity and eccentric anomaly for non-quasi-circular orbits [19, 20], and modification terms for alternative theories of gravity [21, 22].

Due to the transient nature of CBC signals, parameter estimation (PE) is solved under the Bayesian framework [23], in which the probability density distribution of the source parameters is given by the Bayes theorem, but not directly analyzable due to the high dimensionality of the parameter space. Stochastic sampling algorithms, Markov Chain Monte Carlo (MCMC) or Nested Sampling, are used to draw samples from the probability distribution, which represent the estimated source parameters. These methods, however, are computationally intensive. Currently, a typical PE run for a CBC source would cost hours to days, and the cost increases with the signal-to-noise ratio (SNR) and signal duration. A number of acceleration methods for GW source PE are proposed to mitigate the computational burden in stochastic sampling, including Relative Binning (RB, also referred to as Heterodyning) [24–26], Multibanding (MB)[27, 28], and Reduced Order Quadrature (ROQ) [29, 30]. Nevertheless, little work has been done to systematically investigate and compare the performance of these algorithms. The LVK data release still extensively uses the standard PE method without acceleration [3–6], as the cost is affordable given current detection rate, signal SNR and durations.

However, there is a growing concern regarding the computational burden for the proposed third-generation (3G) GW detectors [31, 32]. The 3G detectors, including Einstein Telescope (ET) [33, 34] and Cosmic Explorer (CE) [35, 36], are designed to improve the sensitivity by at least an order of magnitude over Advanced LIGO, and are expected to detect hundreds of thousands of CBC signals every year with significantly higher SNR and longer signal duration [34, 37]. While these detectors promise exciting scientific discoveries, the computational cost of analyzing all the detected signals using current methods is anticipated to be enormous, though it is not yet fully quantified.

In this paper, we present a detailed estimate of the computational costs of PE in the 3G era, considering both the standard PE method and accelerated methods. By conducting 1200 PE experiments and fitting the PE time cost as a function of SNR and duration, we estimate the total costs of PE for 3G detectors based on the simulated catalog in ET Mock Data Challenge 1 (MDC-1). We find that billions to quadrillions of CPU core hours is required to analyze one month of observation with the 3G detectors using the standard PE method, with most of the computational load coming from long binary neutron star (BNS) and neutron star-black hole (NSBH) signals. With acceleration methods, the cost could be reduced to millions of CPU core hours. We also compare the performance of different acceleration methods in terms of speed and accuracy. Our results identify ROQ as the fastest method overall, though it requires additional training before the PE. We also highlight accuracy issues in high-SNR scenarios, with RB showing the worst accuracy due to its reliance on fiducial parameters used in the heterodyning process.

This paper is organized as follows. Section II provides a brief introduction to Bayesian inference and stochastic sampling in Sec. II, including the standard methods in Sec. II A and acceleration methods in Sec. II B. The experiments are presented in Sec. III, with details on the experiment design in Sec. III A, fitting the PE cost in Sec. III B, and accuracy assessment in Sec. III C. Section IV applies the time cost fitting to the ET MDC-1 and estimates the total PE cost in the 3G era. The concluding remarks are provided in Sec. V.

## II. OVERVIEW OF BAYESIAN PARAMETER ESTIMATION

### A. Bayesian parameter estimation

#### 1. Bayes theorem

Bayesian inference gives us a means to quantify uncertainty about the source parameters using the concept of Bayesian updating. Starting with generic prior expectations about the source (which may be informed by astrophysical or geometrical considerations), we observe the transient GW signal, and then update our knowledge based on the new data. Mathematically, given a hypothesis $H$ (e.g. the waveform model) and a prior distribution of the parameters $p(\theta|H)$ (which can be non-informative), we can update the probability distribution of $\theta$ with observation data $d$ using Bayes' theorem:

$$p(\theta|d, H) = \frac{p(\theta|H)p(d|\theta, H)}{p(d|H)}, \qquad (1)$$

where

$$p(d|\theta, H) \propto \exp{-\frac{1}{2}(d - h(\theta)|d - h(\theta))} \qquad (2)$$

is the likelihood function, assuming the noise is stationary and Gaussian [38], $h(\theta)$ is the GW signal given parameter $\theta$, and $(\cdot|\cdot)$ denotes the inner product between two data sets in the frequency domain,

$$(a|b) = 4\,\mathrm{Re} \int_0^\infty \frac{a^*(f)b(f)}{S_n(f)} df, \qquad (3)$$

where Re denotes real part and $a^*$ denotes the complex conjugate of $a$. $S_n(f)$ is the power spectral density (PSD) of the detector noise. $p(d|H)$ is the evidence, a constant that is irrelevant for the shape of the posterior distribution, but is useful in model comparison [39]. The posterior probability distribution $p(\theta|d, H)$ is the updated knowledge and provides the estimation of the parameter $\theta$.

#### 2. Stochastic sampling

Although the posterior probability density function is computable in the entire parameter space (as long as the waveform model is valid) via Eq. 1, it is impossible to compute it all over the space due to the high dimensionality. Instead, stochastic sampling algorithms are used to obtain a set of samples that follows the distribution given by Eq. 1.

MCMC and nested sampling are two families of sampling algorithms. MCMC has a set of walkers randomly walking in the parameter space, with a higher probability of walking towards regions of higher target probability. Over a sufficiently long period of iterations, the trace of the walkers converges to the target probability distribution, which is designed to be the posterior probability. Nested sampling works with a set of live points drawn from the prior distribution. The live point with the lowest likelihood value will be discarded and a new point will be drawn from the space expanded by the remaining live points. Repeating this process, the live points will eventually form the posterior space.

A number of variants of MCMC and nested sampling algorithms have been used in GW astronomy. `dynesty` is the most widely used in recent LVK public data [5, 6], although many novel samplers are being developed to achieve faster speeds. For example, `nessai` [40, 41] is a nested sampling algorithm enhanced with machine learning and importance sampling, and it reduces the number of likelihood evaluations by an order of magnitude in GW source PE compared with `dynesty`. In this work, we use `nessai`, as we find it offers significantly improved speed and convergence for 3G PE tasks compared to `dynesty`.

### B. Acceleration methods

Stochastic sampling for GW PE is computationally intensive. Sampling algorithms may struggle to identify the correct region in the parameter space, requiring long

convergence times and numerous likelihood evaluations (waveform generation and vector manipulation). An increase in signal duration or SNR results in a narrower posterior, making it more challenging for the sampler to locate and converge on the posterior. Additionally, longer signal durations slow down waveform generation and vector manipulation involved in the likelihood evaluation. Both of these factors increase the computational cost of performing PE.

Several acceleration methods have been proposed to speed up the PE. In addition to the general improvement of the sampler like `nessai`, there are methods specifically targeting the CBC PE problem, aiming to reduce the data size and simplify likelihood evaluations. In this section, we introduce three methods that we will examine in this work: Relative Binning (RB), Multibanding (MB), and Reduced Order Quadrature (ROQ). All of them have been implemented in the GW PE Python package `bilby` [42]. We will also mention other fast PE methods in Sec. II B 4 but they are not involved in our assessment due to either having completely different mechanisms, being a mixture of many methods, only obtaining a subset of parameters, or having an incomplete software implementation.

### 1. Relative Binning

Relative Binning (RB) [24, 25], also known as heterodyning [26], is based on the idea that many likelihood evaluations occur near the maximum likelihood point during PE, and that the GW waveform changes smoothly as a function of the source parameters. Therefore, if a fiducial source parameter that is close to the true value can be assigned before PE (which is possible given information from matched filtering detection), the PE can focus on in the nearby region where likelihood evaluations can be simplified by exploiting the smooth variations in the waveform as the fiducial parameter is varied. Since the smooth variation with respect to the fiducial waveform requires a lower bandwidth to describe, a coarser frequency resolution is sufficient to reconstruct waveforms near the fiducial waveform This reduces the data size and therefore the computational cost of computing each likelihood. The new frequency bins with coarse frequency resolution can be chosen such that the GW signal only contains a few cycles within each bin, which is computable using the post-Newtonian (PN) phase expansion [43].

Considering frequency domain waveforms, we approximate the ratio between the fiducial waveform $h_0(f)$ and an arbitrary waveform $h(f)$ in a frequency bin b to the linear order

$$r(f) = \frac{h(f)}{h_0(f)} = r_0(h,\mathrm{b}) + r_1(h,\mathrm{b})\,(f - f_\mathrm{m}(\mathrm{b})) + \cdots, \quad (4)$$

where is $f_\mathrm{m}$ is the central frequency of the bin. $r_0$ and $r_1$ are bin-dependent coefficients that measure how $h$ devi-

ates from $h_0$, and can be efficiently derived from the values of $r(f)$ at the edges of bins. Since $(d|d)$ is a constant, the likelihood (Eq. 2) only requires the computation of $(d|h)$ and $(h|h)$, which can be approximated using $r(f)$. Define

$$\begin{aligned} A_0(\mathrm{b}) &= (d|h_0)_\mathrm{b}, \quad A_1(\mathrm{b}) = (d|(f - f_\mathrm{m})h_0)_\mathrm{b} \\ B_0(\mathrm{b}) &= (h_0|h_0)_\mathrm{b}, \quad B_1(\mathrm{b}) = (h_0|(f - f_\mathrm{m})h_0)_\mathrm{b}, \end{aligned} \quad (5)$$

where $(\ldots|\ldots)_\mathrm{b}$ is the inner product inside the frequency bin b, i.e., the integral limits are the minimum and the maximum frequencies of the bin. $(d|h)$ and $(h|h)$ can be calculated as

$$\begin{aligned} (d|h)^{\mathrm{RB}} &= \sum_\mathrm{b} (d|h)_\mathrm{b} = \sum_\mathrm{b} (d|rh_0)_\mathrm{b} \\ &\approx \sum_\mathrm{b} \left( A_0(\mathrm{b})r_0^*(h,\mathrm{b}) + A_1(\mathrm{b})r_1^*(h,\mathrm{b}) \right), \end{aligned} \quad (6)$$

and similarly

$$(h|h)^{\mathrm{RB}} \approx \sum_\mathrm{b} \left( B_0(\mathrm{b})\,|r_0(h,\mathrm{b})|^2 + 2B_1(\mathrm{b})\mathrm{Re}\left[r_0(h,\mathrm{b})r_1^*(h,\mathrm{b})\right] \right). \quad (7)$$

For any waveform $h$, it is only necessary to evaluate the waveform at the edges of the bins to obtain $r_0$ and $r_1$, which speeds up the likelihood evaluations. In addition, focusing on regions around the fiducial parameters significantly reduces the number of likelihood evaluations, although this may bring accuracy issues.

Examples of applications of RB can be found in Refs [44–48]. In particular, Dai *et al.* [46] showed that PE for GW170817 with aligned spin waveform model `IMRPhenomD_NRTidal` [49] and `TaylorF2` [43] and RB can be done with 150 CPU core hours. Wong *et al.* [48], also using aligned spin waveform model `IMRPhenomD` (`_NRTidal`), demonstrated a combination of RB and a gradient-based sampler, and achieved 2 hours of sampling time for GW150914 and 1 day for GW170817 with 400 CPU cores.

### 2. Multibanding

The frequencies of chirp signals from CBC sources increase during the binary inspiral. GW signals are present in the data from the low-frequency cutoff $f_{\min}$ of the GW detector ($\sim 20\,\mathrm{Hz}$ for current detectors and $\sim 5\,\mathrm{Hz}$ for 3G) up to the maximum frequency $f_{\max}$ of hundreds of Hz up to $\sim 2\,\mathrm{kHz}$, depending on the source masses. The full-bandwidth GW data (i.e. strain data with sampling rate of $4096\,\mathrm{Hz}$ or $16384\,\mathrm{Hz}$) is excessive for the low-frequency early inspiral stage, therefore the data can be adaptively down-sampled as long as sampling rate remains above the Nyquist frequency. This approach is known as multibanding or adaptive frequency resolution.

Several different MB schemes have been proposed; see, e.g., Refs [27, 28]. In this work, we adopt the scheme used in Morisaki [28], in which the time-domain data is

divided into bands with geometrically decreasing lengths $(T, T/2, T/4, \ldots, 4s)$. The starting and ending frequencies of each band are chosen such that (i) the data in the time-domain band can cover the signal within the frequency band, with time-frequency relation given by the 0PN equation for the lowest mass in the prior and (ii) the waveform smoothly vanishes at the left edge of each time-domain band with the window function applied. The frequency bands and relevant coefficients are determined once the prior is set, requiring no additional calculations before or during PE.

MB leads to coarser frequency grids on which the GW waveform needs to be evaluated and thus reduces the data size, accelerating waveform generation and likelihood evaluation. This requires waveform models that can be calculated at any frequency, typically with closed-form expressions in the frequency domain, such as the `IMRPhenom` waveforms [50, 51]. However, MB does not reduce the number of likelihood evaluations. Morisaki [28] shows that MB can speed up PE by a factor of 20-50 for signals starting from 20 Hz, and by a factor of 100–500 for signals starting from 5 Hz, in line with the asymptotic scaling $5f_{\max}/3f_{\min}$ found in [27]. MB can be combined with ROQ [28, 52] and RB [47, 53] to further enhance the performance.

In addition to PE, the similar idea has been explored in detection [54] and source localization [55].

### 3. Reduced Order Quadrature

ROQ [29, 30] utilizes the reduced order modeling (ROM) [56–58] of GW waveforms, which aims to accelerate waveform evaluation. In the frequency domain, the ROM of a waveform model can be written as the linear addition of $N$ bases, where $N$ is much smaller than the original length of the waveform:

$$h(f;\theta) \approx h^{\mathrm{ROM}}(f;\theta) = \sum_{J=1}^{N} h\left(\mathcal{F}_J;\theta\right) B_J(f), \quad (8)$$

where $B_J(f)$ is a set of bases that span the signal space obtained via a greedy algorithm [56] or singular value decomposition [57]. $\{\mathcal{F}_J\}$ is a set of the most representative frequencies that can be used to reconstruct the entire waveform determined by empirical interpolation [56]. The $h^{\mathrm{ROM}}(f;\theta)$, as an approximation of the original waveform model $h(f;\theta)$, only requires waveform evaluation at $\{\mathcal{F}_J\}$ and a set of predetermined bases $\{B_J(f)\}$, and is therefore faster to calculate. For low-dimensional parameter spaces, it is possible to build a surrogate model $h_J(\theta)$ to approximate the waveform value at a specific frequency node $h\left(\mathcal{F}_J;\theta\right)$. $h_J(\theta)$ can takes the form of polynomials [56], splines [57], or even neural networks [59]. However, PE for precessing systems has a high-dimensional parameter space in which surrogate models cannot be easily built, hence it is not suitable for our task. We therefore still need a waveform model that can be evaluated at any frequencies, like we do in MB.

The ROQ likelihood can be expressed using the ROM. The non-constant terms in the likelihood are:

$$(d|h)^{\mathrm{ROQ}} = \mathrm{Re} \sum_{J=1}^{N} h\left(F_J;\theta\right) \omega_J\left(t_c\right), \quad (9)$$

and

$$(h|h)^{\mathrm{ROQ}} = \sum_{I=1}^{N} \sum_{J=1}^{N} h^*\left(F_I;\theta\right) h\left(F_J;\theta\right) \psi_{IJ}, \quad (10)$$

where $\omega_J\left(t_c\right)$ and $\psi_{IJ}$ are integration weights that depend on the basis functions, data, and PSD, and their definitions can be found in Ref. [60]. To perform ROQ, basis functions $\{B_J(f)\}$ need to be calculated under a prior distribution of the intrinsic parameters (which can be broader than the prior used in PE) and the integration weights need to be tailored to the specific data. This precalculation stage can take anywhere from minutes to hours, depending on factors such as signal length and the prior range. However, the most computationally expensive step, basis construction, only needs to be done once for a wide range of sources, as long as the source falls within the prior used for basis construction. ROQ significantly accelerates the likelihood calculation by reducing it to linear and quadratic summations with waveforms evaluated at a reduced frequency grid.

ROQ can bring speed-up factors of more than ten times depending on the signal duration, and is widely-used in GW PE, see e.g. [29, 52, 60, 61]. Notably, Ref. [60] illustrates ROQ's application to long signals from BNS in 3G detectors: using a restricted ROQ prior and excluding the Earth's rotation effects, PE for a 90-minute-long BNS signal of SNR of 2400 can be done with 20 CPU hours of precalculation and 1600 CPU hours of sampling. Bayesian inference for such signals would be prohibitively slow ($> 10^7$ CPU hours) using the standard PE method.

### 4. Other methods

There are a number of other approaches for fast PE. One common technique is marginalization, which reduces computational cost by integrating over certain parameters. The luminosity distance, coalescence time, and coalescence phase (for the dominant $\ell = m = 2$ mode of the signal) can be analytically marginalized in the posterior [23, 39], allowing them to be excluded from the sampling process and reconstructed after the other parameters are sampled. Taking the matched filtering SNR timeseries as input, numerical marginalization of more parameters is demonstrated in `cogwheel` [45], achieving minutes to a few hours of PE time on a single CPU. `RIFT` [62, 63] also makes use of likelihood marginalization, evaluating the marginal likelihood on a grid of intrinsic parameters, and then constructing a fast interpolator to explore the posterior. Marginalization techniques

are also used for quickly obtaining subsets of parameters such as for source localization [64–66].

Machine learning methods based on generative models, such as `DINGO` [47, 67–69] based on normalizing flows [70, 71] and `Vitamin` [72] based on variational autoencoders [73], have shown great potential in fast PE. Generative models learn the data distribution and can quickly draw samples from the posterior in seconds once trained. The samples can be further improved with importance sampling at an additional time cost of a few minutes [68]. `DINGO` has been applied to the inference of LVK events and presents fast speed and good accuracy, including the inference with slow waveforms [69] that could be extremely expensive for traditional Bayesian inference. There are also 3G-specific applications in the literature, such as for overlapping signals [74] and long BNS signals in 3G detectors [53]. The speed advantages of machine learning approaches make them particularly well-suited for catalog-level tasks, a critical need in the 3G era with the huge increase in signal numbers.

## III. EXPERIMENTS

To estimate the cost of Bayesian inference in the 3G era, we need to understand the relationship between signal properties (duration, SNR, etc.) and the associated sampling cost. To do this, we perform 1200 PE runs with varying settings and fit the resulting relation. In this section, we describe the experiment setup, outline the process for obtaining the time cost fitting, and discuss the accuracy of different acceleration methods.

### A. Setup

We consider a detector network contains one triangular 10km-ET operating at the ET-D design sensitivity [75] at the Virgo site, and one 40km CE detector with the CE2 design sensitivity [35] at the LIGO Hanford site. Detector noise is assumed to be stationary and Gaussian, and different events have different noise realizations. The starting frequency is 5 Hz and data is sampled at 2048 Hz.

We simulate three types of binary black hole (BBH) sources: GW151226-like ($14+8M_\odot$, duration=256s), GW150914-like ($36+30M_\odot$, duration=64s), and GW190521-like ($85+66\ M_\odot$, duration=16s). They represent typical low, medium, and high mass BBHs in the population and produce long, medium, and short duration signals, respectively. We fix their intrinsic parameters (based on GWTC-3 [6] and GWTC-2.1 [5] results) and randomly simulate extrinsic parameters except for coalescence time and luminosity distance. The coalescence time is fixed at GPS time $t_c = 0\,s$ as this choice should not influence the result. The luminosity distance is scaled to adjust the network SNR to values in the set $(12, 30, 100, 500, 1000)$, ranging from near-threshold detections to the high-SNR golden events. For each SNR

and source type, 20 events with different extrinsic parameters are simulated. Each event is analyzed with four PE methods: the "standard" method (no acceleration method applied), RB, MB, and ROQ, resulting in a total of 1200 PE runs.

The signals are simulated and analyzed with the `IMRPhenomPv2` [50, 51] waveform model, which leads to 15 parameters to estimate. Although `IMRPhenomPv2` is not the most accurate waveform model to date, it is very fast to evaluate and can be efficiently parallelized for calculation. We may hope that by the time of 3G observing, advances in hardware and algorithms will reduce the cost of waveform evaluation, so this can be considered a conservative choice.

- Chirp mass $\mathcal{M}$: Uniform in $[8.5, 9.5]M_\odot$, $[26, 34]M_\odot$, and $[50, 100]M_\odot$, for GW151226-like, GW150914-like, and GW190521-like sources, respectively.

- Mass ratio $q$: Uniform in $[0.25, 1]$.

- Spin dimensionless magnitudes $a_1, a_2$: Uniform in $[0, 0.8]$.

- Spin tilt $\theta_1, \theta_2$: Uniform in cosine (isotropic spin).

- Spin angles $\varphi_{12}, \varphi_{JL}$: Uniform in $[0, 2\pi]$ rad.

- Inclination angle $\theta_{JN}$: Uniform in cosine (isotropic orientation).

- Right ascension $\alpha$: Uniform in $[0, 2\pi]$ rad.

- Declination angle $\delta$: Uniform in sine (isotropic on the sky).

- Polarization angle $\psi$: Uniform in $[0, \pi]$ rad.

- Coalescence phase $\phi_c$: Uniform in $[0, 2\pi]$ rad.

- Coalescence time $t_c$: Uniform in $[-0.1, 0.1]s$.

- Luminosity distance $d_L$: Uniform in the volume between $[d_L/2, 2d_L]$, where $d_L$ is the injected value.

We employ the sampler `nessai` for faster convergence, though future improvements to the sampler may further enhance performance. We set number of live points `nlive`=2000, which is sufficient for all signals in our simulation. Larger `nlive` may be required for higher SNRs or narrower posteriors such as for BNS signals, which would increase the sampling time. Experiments are run with `ncpu`=32 CPU cores with the AMD EPYC 7443 CPU model. The likelihood is evaluated in parallel but the training and sampling parts of `nessai` are not parallelized and only use one core. As a result, the total sampling CPU core time is calculated as total sampling wall time + (`ncpu`-1)×likelihood evaluation wall time. [76].

The ROQ bases are calculated using `PyROQ` package [77], with the same prior used in PE. Using reconstruction error tolerance of $10^{-8}$ for ROM bases and

$10^{-10}$ for quadratic bases, we obtain 392, 183, 136 ROM bases, and 184, 96, 112 quadratic bases for GW151226-like, GW150914-like, GW190521-like sources, respectively. The data length is significantly compressed compared to the original lengths of 64k, 32k, and 16k data points. With a single CPU core, ROQ training takes less than one hour for GW150914-like and GW190521-like sources, but takes five days for the GW151226-like source. Although it is possible to construct ROQs with a wider prior for more general applications, this would increase the training time and reduce the compression rate. Therefore, we focus the ROQ construction on the simulated signals. The precalculation time cost for building the ROQs is not included in the ROQ sampling time.

The fiducial parameter for RB is chosen to be close to the injection parameter. The error in $\mathcal{M}$ is randomly sampled within $\pm 0.1\%$, while errors for $q$ and $d_L$ are within $\pm 5\%$. The errors for $\alpha$, $\delta$, and $\theta_{JN}$ are sampled within $\pm 10\%$, and the coalescence time has no error (0%). Parameters such as spins and angles, which are difficult to estimate accurately before PE, are set to zero in the fiducial parameters. The $< 0.1\%$ error for $\mathcal{M}$ and $< 5\%$ error for $q$ may not be highly realistic, however, we observe convergence issues in PE with RB if the fiducial parameter is too far from the true value, particularly in high-SNR events. This is because RB is effective only within a region near the fiducial parameters, and it may fail to reach the narrow posterior region. In our experiments, five RB runs failed to converge and were excluded from the time cost estimate.

The code for the experiments can be found in the GitHub repository `3gpemethods` [78].

### B. Bilinear fitting for time costs

We calculate the CPU core time for stochastic sampling in our experiments, and the results are shown in Fig. 1. Overall, ROQ exhibits the fastest speed among all PE methods, completing sampling within 2 CPU days for most events in our simulation. MB and RB typically finish within 5 CPU days. In contrast, the standard method can take tens to hundreds of CPU days. Acceleration methods provide a speed-up factor of up to several hundred times, especially in high-SNR and long-signal scenarios.

There is a clear trend that sampling time increases with both SNR and signal duration, with the scaling being approximately linear in the logarithmic scale. Therefore, we use the following bilinear function to fit the relation between sampling CPU days $D$, signal duration $T$ (in second) and network SNR:

$$\log_{10} D = a \log_{10} T + b \log_{10} \text{SNR} + c, \qquad (11)$$

where the coefficients $a, b, c$ for different PE methods are determined by linear regression and are listed in Table I. We also calculate the $R^2$ score, the coefficient of determination, a higher value of which indicates a better fit.

The $R^2$ score exceeds 0.5 for all PE methods except RB, which may be caused by its convergence issues since the choice of fiducial parameter has impact on the convergence speed. The regression results are shown as dashed lines in Fig. 1, and the accuracy is sufficient for a rough, order-of-magnitude level estimation. As a crosscheck of the extrapolation, Ref. [60] reported a 90-minute long BNS signal with SNR of 2400 can be analyzed with 1600 CPU core hours using the ROQ technique, while our fitting predicts $\sim 1400$ CPU hours. Given that we are underestimating the time cost (see discussions in Sec. IV B), we consider this a good agreement.

| Method | $a$ | $b$ | $c$ | $R^2$ |
|---|---|---|---|---|
| Standard | 1.282 | 0.243 | -1.341 | 0.941 |
| RB | 0.283 | 0.099 | -0.608 | 0.397 |
| MB | 0.239 | 0.224 | -0.708 | 0.596 |
| ROQ | 0.085 | 0.216 | -0.654 | 0.512 |

TABLE I. Coefficients $a, b, c$ from bilinear regression in Eq. 11 for sampling time. The last column gives the coefficient of determination $R^2$ score ranging from 0 to 1, with higher value indicating a better match between prediction and data.

Looking into the coefficients, we notice the standard method has the largest scaling factors ($a$ and $b$) with respect to the signal duration and SNR, indicating that its time cost grows most rapidly, which aligns with our expectations. ROQ has the lowest scaling factor (the coefficient $a$) with the signal duration, highlighting its good efficiency in speeding up likelihood evaluations. RB shows the lowest scaling factor (the coefficient $b$) with the signal SNR. For the runs that finish quickly (for example, low SNR and short-duration events), we note that the total cost is not dominated by the likelihood but by other operations of the `nessai` sampler.

It is worth mentioning that the effect of lower frequency cutoff $f_{\text{low}}$ can be estimated roughly using $a$. If $f_{\text{low}}$ changes to 2Hz or 3Hz, the SNR would not change significantly, but the duration of the signal would scale as $f_{\text{low}}^{-8/3}$. Using Eq. 11, the sampling time $D$ will be *multiplied* by a factor:

$$D_{f_{\text{low}}} = (\frac{5}{f_{\text{low}}})^{\frac{8a}{3}} D_{5\text{Hz}}, \qquad (12)$$

where $D_{5\text{Hz}}$ is the result given by Eq. 11 and $a$ is given in Table I. Substituting the values, the PE time cost of the standard method would increase roughly 22 and 6 times for $f_{\text{low}}$ of 2Hz, 3Hz, respectively. For the accelerated methods, the time increase is between 10% and 100%. Note that the increase is underestimated, since the changes in SNR are not accounted for.

### C. Accuracy

We assess the accuracy of the acceleration methods by comparing the posterior distributions obtained us-
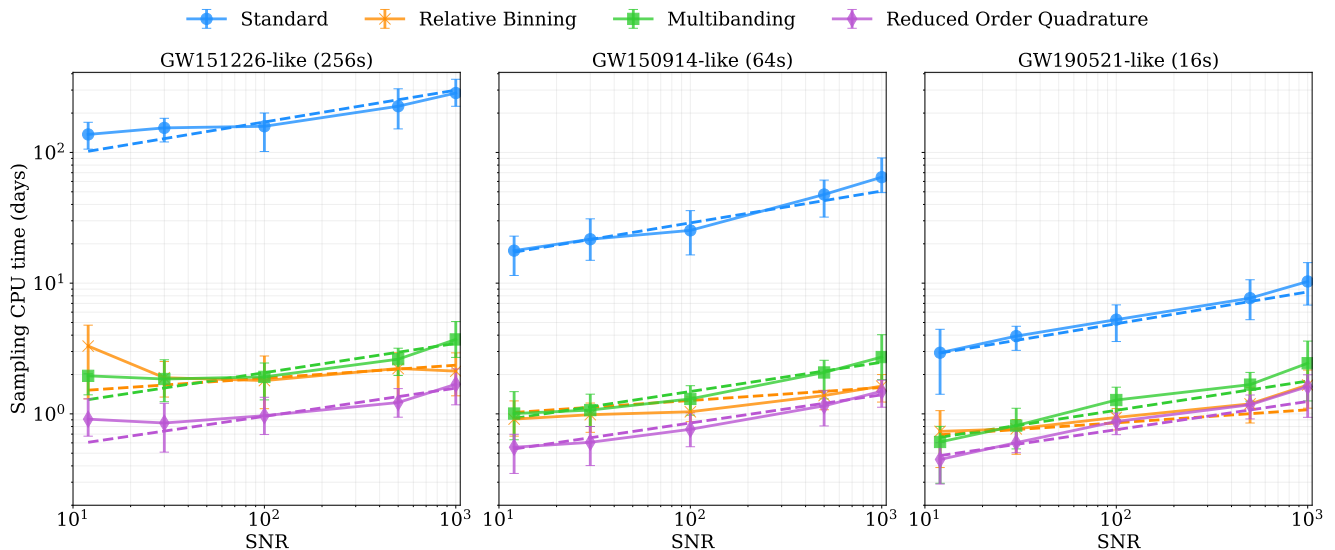
FIG. 1. The CPU core time of stochastic sampling in our experiments, along with the bilinear fitting. The error bars represent the 16% to 84% percentiles of the sampling time across 20 events for each configuration, with the mean indicated in the center. The dashed lines correspond to the bilinear fits of the sampling time cost.

ing each method to those obtained with the standard method. For each parameter in each event, we compute the Jensen-Shannon Divergence (JSD) between the posterior distribution from an acceleration method and the standard method. The JSD is calculated with a base of 2, so it ranges from 0 to 1, with the unit of *bits*. A JSD value of 0 indicates that the two posterior distributions are identical, while a value of 1 means that they are completely different, with no overlap in the result.

The Jensen-Shannon Divergence (JSD) is shown in Fig. 2. For ROQ and MB, most parameters are estimated with good accuracy, with typical JSD ranging from $10^{-3}$ to $10^{-2}$ bits. We observe that the error increases with SNR, suggesting the likelihood needs to be approximated with a higher accuracy in high SNR cases. For instance, increasing the number of bases in ROQ may help address this issue. RB shows the poorest accuracy among all PE methods, despite having relatively accurate fiducial parameters. Its accuracy deteriorates when the posterior is narrow, which happens in high SNR events and long signals (where the chirp mass posterior is narrow). Providing more accurate fiducial parameters may alleviate this issue, but extra work is needed to determine the required accuracy of the fiducial parameter before performing a full PE. A comprehensive investigation of PE accuracy in the 3G era is beyond the scope of this paper and should be addressed in future work.

## IV. COSTS ESTIMATES FOR 3G DETECTORS

With the estimated time cost given by Eq. 11, we now predict the total cost of Bayesian inference for 3G detectors using a realistic source population. We will first
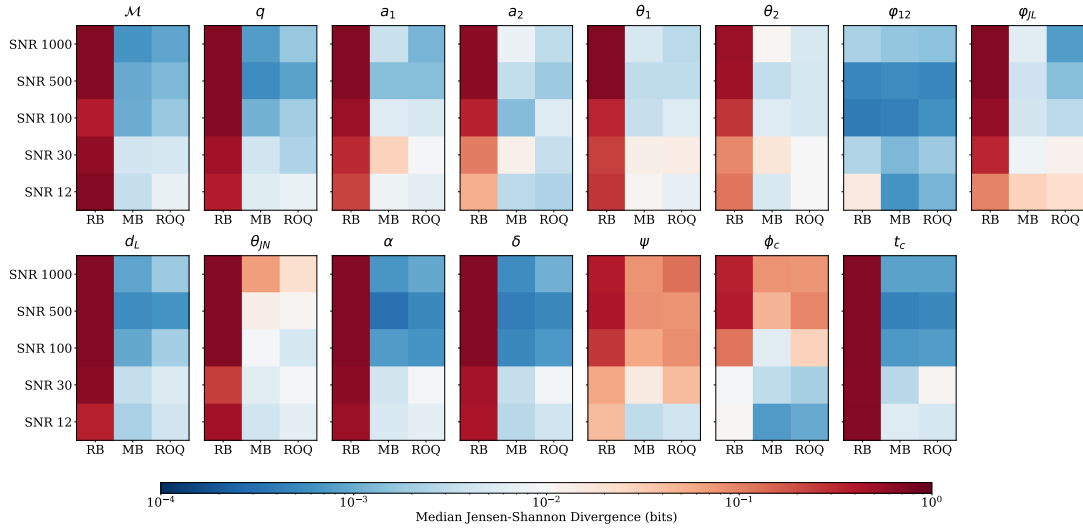
introduce the mock data catalog, ET MDC-1, and then demonstrate how the total cost is estimated.
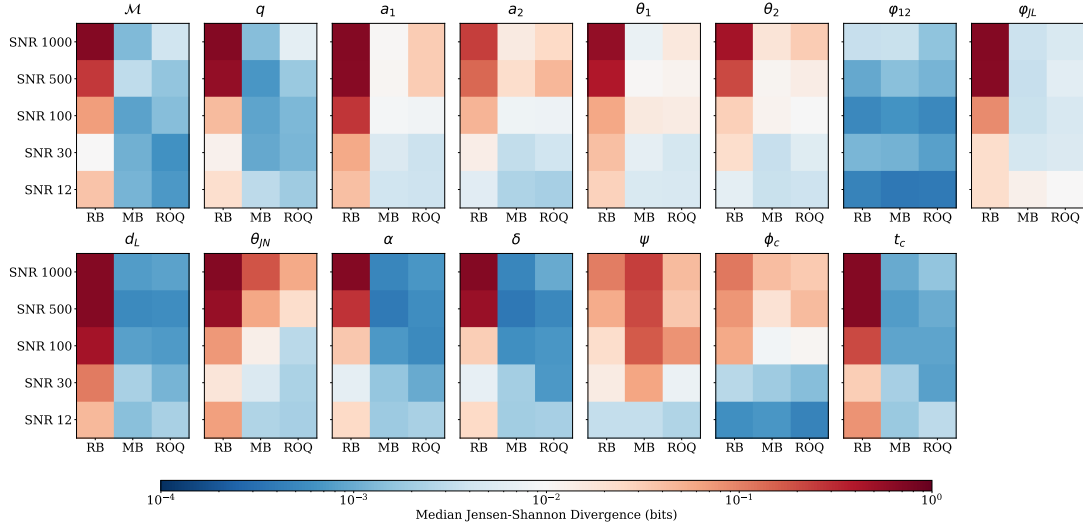
### A. Overview of ET MDC-1

ET MDC-1 contains one month of simulated observation with the 3G detector ET [79]. The population model used in ET MDC-1 includes 6812 BBH, 61217 BNS, and 2029 NSBH sources over the month, and varying numbers of detections possible depending on the detector network. We consider three types of detector networks: ET, ET-CE, and ET-CE-CEL (ET-2CE). Here, 'CE' is assumed to be located at the LIGO Hanford site and 'CEL' at the LIGO Livingston site, with the same PSD used in the PE experiments. We calculate the network SNR for each network configuration and present the cumulative distribution of the SNR in Fig. 3. Setting SNR=8 as the detection threshold, a single ET could detect most BBH sources but would miss a certain fraction of BNS and NSBH sources. Adding CEs to the network significantly enlarges the detection range of BNS and NSBH sources.

We also show the cumulative distribution of the signal duration in Fig. 3. The signal duration is defined as the time it takes for the binary system to evolve from 5 Hz to the merger, computed using the 0PN equation. We add an additional 4s to the duration of each signal to account for the FFT corruption in real data analysis. Most BBH signals have durations less than 100s, while NSBH and BNS signals are mostly longer than 100s. Some BNS signals can last for over an hour.
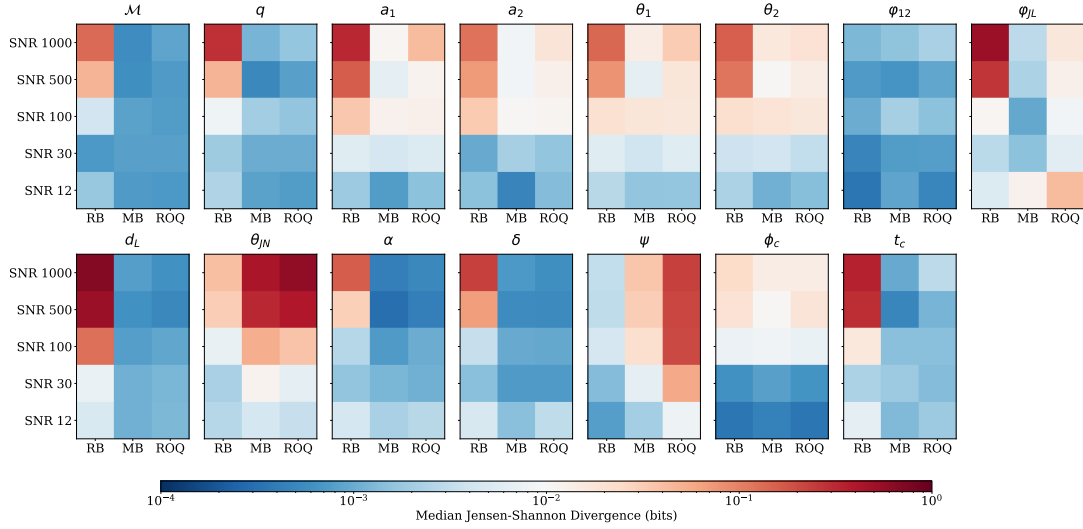
Note that ET MDC-1 only provides the GW strain data for ET. However, in this study, we focus solely on the signal duration and the (optimal) SNR, which are

(a)GW151226-like



(b)GW150914-like



(c)GW190521-like

FIG. 2. The Jensen-Shannon Divergence (JSD) between posteriors obtained using acceleration methods and the standard method. Each subfigure corresponds to a different source type, and each patch represents a different parameter. The color of each grid represent the value of JSD in *bits*. A JSD of 0 indicates that the two posterior distributions are identical, while a JSD of 1 indicates they are completely different.
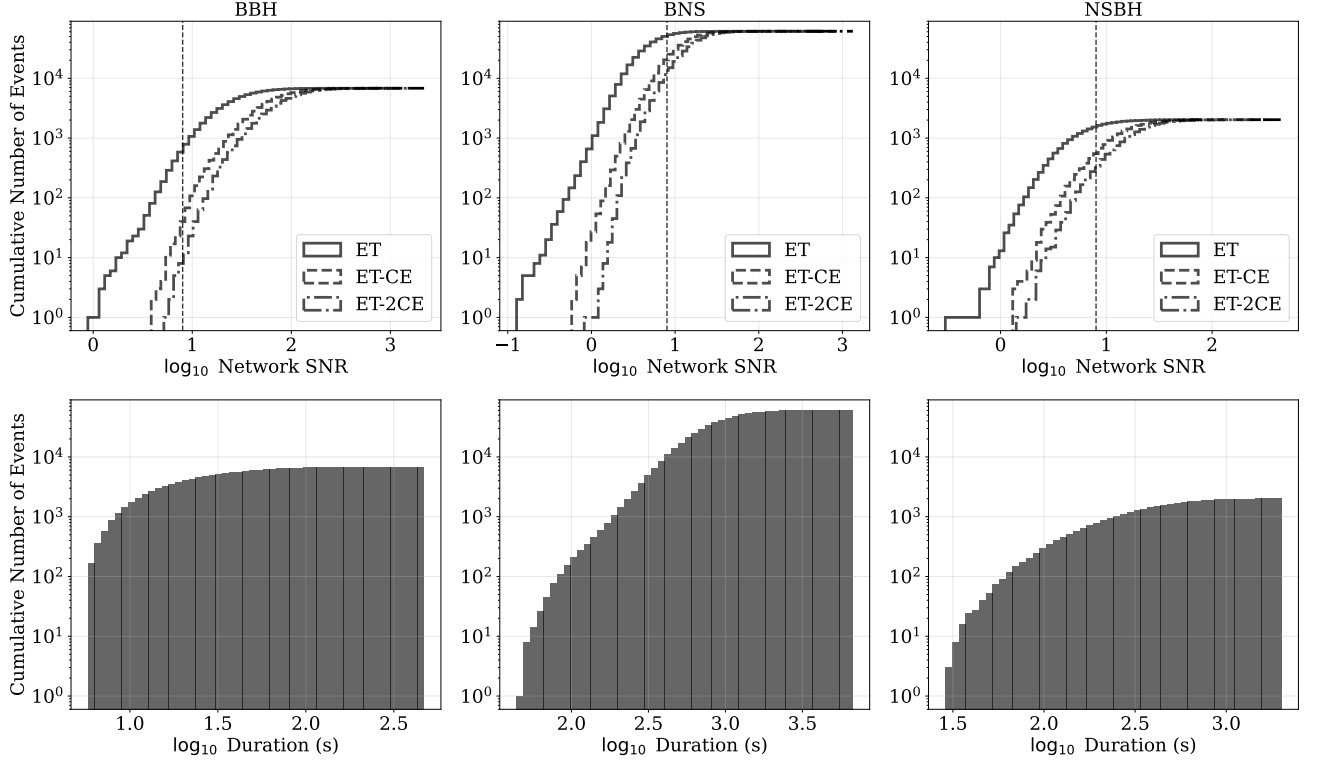
FIG. 3. **Top:** Cumulative distributions of network SNR of different sources (in different columns) and different detector networks (in different linestyles). The black dashed lines mark the SNR=8 threshold. **Bottom:** Cumulative distributions of signal durations, measured from 5 Hz, for different source types.

independent of the actual noise realization.

### B. Total PE cost estimates

Instead of directly using GW strain data, we simply use the signal durations and SNRs based on the MDC catalog to predict the sampling time with Eq. 11. However, we have overlooked some factors that also affect the sampling time. We discuss their expected contributions here:

- Parameter space dimensionality: Our experiments are based on 15-D PE for BBH systems. For NSBH and BNS systems, there are one and two more parameters for the tidal deformability, respectively. Moreover, future analyses may incorporate eccentricity, adding two more dimensions. Consequently, we are *underestimating* the total time cost for 3G. The scaling of sampling time with dimensionality is not entirely certain. Although it is shown that the number of likelihood evaluations scales exponentially with the dimensionality for `nessai` [41], the exact scaling depends on the new information (relevant to posterior and prior widths) brought by the new dimensions. Some samplers can achieve polynomial scaling [80], which may mitigate the in-

crease in time cost due to the increased dimensionality.

- Data size: The PE experiments use the ET-CE network that contains four data streams (three from ET and one from CE), but with a different network the number of data streams changes. We introduce the *effective signal duration* $T_{\mathrm{eff}} = \frac{N_{\mathrm{stream}}}{4} T$ to account for this change, where $N_{\mathrm{stream}} = 3, 4, 5$ for ET, ET-CE, ET-2CE, respectively. In other words, we consider the signal duration solely as a factor of data size, and equate adding/removing detectors with stretching/cutting the data. We expect this to roughly balance the change in the number of data streams.

- Sampling rate: We use a sampling rate of 2048 Hz in all experiments, but for 3G analysis this number could increase, especially for signals with higher modes and low mass systems. We are therefore *underestimating* the time cost.

- Sampling algorithm and configuration: Our investigation a variant of the nested sampling algorithm, which is commonly used in current LVK analyses. Our sampler configuration is fixed in all experiments. However, for longer and louder signals, a larger `nlive` may be required, which slows

down the PE. Since our experiments only extend to T=256 s, we expect this could happen in signals that last for hours. In this sense, we are *underestimating* the total time cost. We also expect that results would vary with different choices of sampling method, for example MCMC does not require the code to sample the entire prior before converging on the posterior and so may run faster if initial parameter estimates are available [81].

- Overlapping signals: In realistic 3G detector scenarios, signals from different sources could overlap in the time domain [82–87], requiring a more sophisticated analysis such as joint PE. We do not consider this effect in the experiments, leading to an *underestimate* of time cost.

- Long signal effects: For signals longer than $\sim 10$ minutes, several additional effects need to be taken into consideration, such as the Earth's rotation and variations in PSD. These may bring more calculations during PE and thus we are *underestimating* of time cost.

- *Ad hoc* ROQ: We restrict the prior of building ROQ to a small range, which reduces the required number of ROQ bases and lead to *underestimation* of time cost. However, on the other hand, ROQ could be improved with MB [52] and so compress the likelihood better, and this is not included in our experiments.

- Duty cycle and evolving sensitivity: GW detectors may not operate continuously throughout the entire month, and their sensitivity may not reach the design level at the start of the observation period. As a result, we may *overestimate* the total time cost. However, the duty cycle can be easily included by applying a constant factor, and the sensitivity does not influence the validity of Eq. 11. Therefore, our methods and results can be readily extended to more realistic observational scenarios.

Although it appears we may have been underestimating the time cost, it is difficult to predict future advancements in sampling techniques that could further accelerate PE. Given these caveats, we can say that with current technology, we provide an *optimistic* estimate of the future PE cost.

For each source in the MDC that is considered to be detected (network SNR>8), we use Eq. 11 and the effective duration to estimate its PE time cost. The cumulative histogram of the PE time cost is shown in the top row of Fig. 4. Most events, including BNS and loud signals, can be analyzed within 1000 CPU core hours using acceleration methods. However, when using the standard method, the time cost becomes prohibitively large, reaching up to $10^{10}$ CPU core hours. The total time cost to analyze one month of observation is shown in the bottom row of Fig. 4. Without acceleration methods, the

total CPU core hours required can reach $10^{15}$, which is clearly unrealistic. Acceleration methods can reduce this number to millions: we should expect them to play an important role in the future data analysis.

We do not fully explore different detector configurations in this study (e.g., L-shaped ET and CE-only networks). However, since the estimation of the PE time cost is purely a function of signal SNR and duration, we do not expect the total PE time cost to change significantly for different detector configurations. For example, for a single CE detector, the total time cost should lie between the estimations of a single ET and ET-CE network, and a 2CE network should lie between the ET-CE and ET-2CE network. The lower frequency cutoff affects the total PE cost approximately by the factor given in Eq. 12.

The LVK Computing Infrastructure [88] currently operates with fewer than 50 thousand CPU cores. Therefore, the millions of CPU hours required for the one-month observation would take several days to process on the entire computing cluster, consuming approximately 1 GWh of electricity and costing hundreds of thousands of USD in electricity charges. It is important to note that this is the cost of performing PE only once; in LVK, each event may be analyzed multiple times before the final results are obtained for publication. Considering potential future upgrades to computing clusters (e.g., the Worldwide LHC Computing Grid, which operates with more than 500,000 CPU cores [89]), we conclude that while full Bayesian PE is technically feasible, it is neither budget-friendly nor environmentally sustainable. Therefore, a more efficient PE method - one that surpasses the methods discussed in this paper - will be crucial for GW astronomy in the 3G era.

## V. CONCLUSIONS AND DISCUSSIONS

The computational challenge for the next-generation GW detectors has drawn significant attention in recent years, yet the precise computational cost has been largely an unexplored domain. In this paper, we investigated the computational cost associated with one of the most computationally intensive aspects of GW astronomy: Bayesian PE with stochastic sampling in the 3G era. We started from simulated signals that can be analyzed with a manageable cost and obtained a relationship between the sampling time and signal duration and SNR. With this relationship, we predicted the total PE cost for the event catalog from the ET MDC-1.

Our experiments include the standard PE method and three accelerated methods: relative binning, multi-banding, and reduced order quadrature. The standard method has been shown to be impractically slow for the 3G catalog: using this approach, some events would require $10^{10}$ CPU hours to analyze, and analyzing the entire one-month catalog could demand between $10^{13}$ and $10^{15}$ CPU hours. In contrast, the acceleration methods
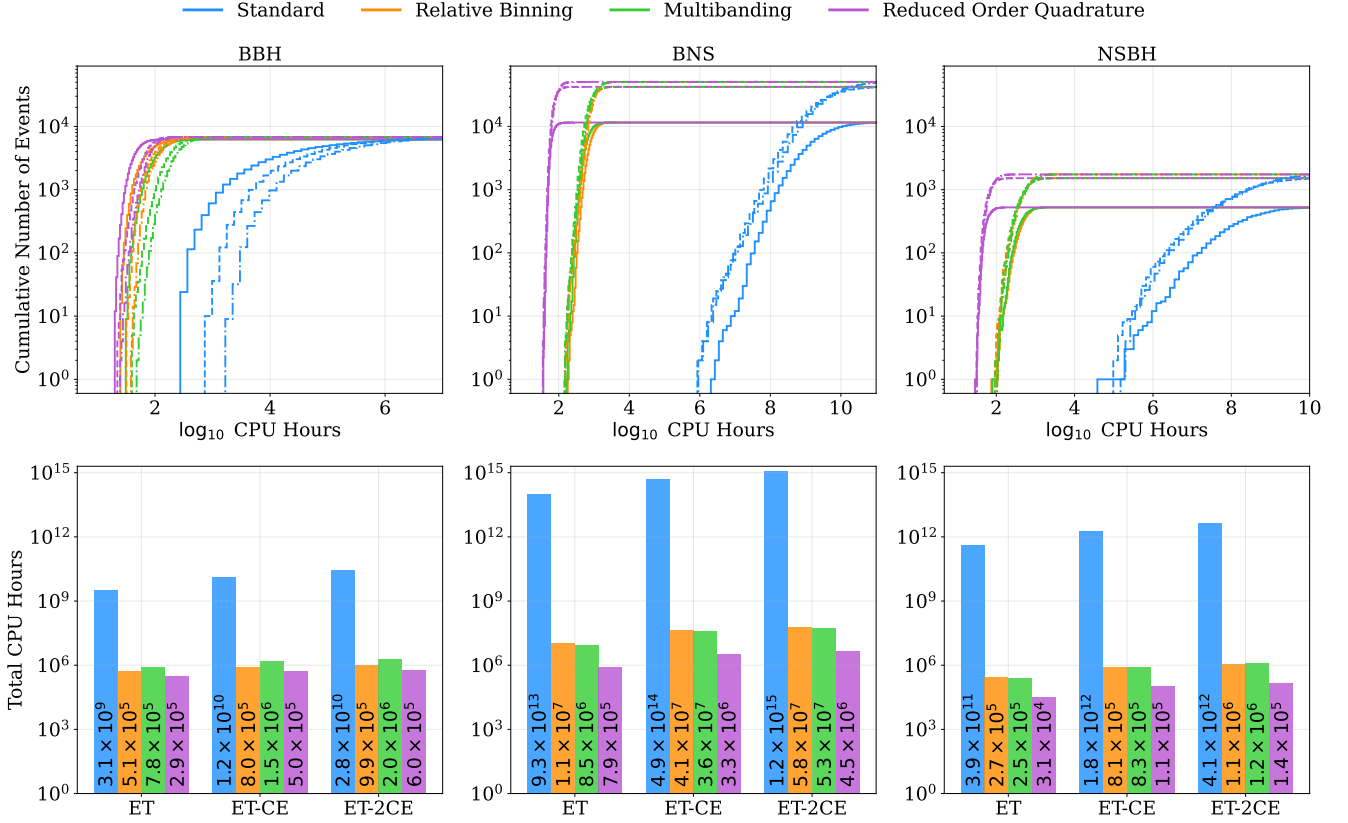
FIG. 4. **Top:** Cumulative distributions of estimated PE cost for detected events (SNR>8, in one-month observation) in CPU hours with different detector networks. The solid line represents ET-only network, dashed line represents the ET-CE network, and solid-dashed line represents the ET-2CE network. Different source types are in separate columns, with different colors representing various PE methods. **Bottom:** Total CPU hours required to perform Bayesian PE for the one-month observation with different detector networks. The numbers indicate the corresponding CPU core hours for each bar.

significantly reduce the time cost. Most events can be analyzed within 1000 CPU hours, and the total cost for the entire catalog is reduced to millions of CPU hours. However, this is still far from ideal, as millions of CPU hours per month remain a substantial burden on computing infrastructure, electricity costs, and environmental impact. To address this, more efficient methods must be developed, such as advanced samplers, improved data compression, faster waveform evaluations, parameter marginalization, and machine learning techniques.

Among the acceleration methods we tested, we found that ROQ provides the best speed while maintaining reasonable accuracy. RB and MB have similar speeds, approximately twice as slow as ROQ. However, we observed accuracy issues with RB, potentially due to its dependency on fiducial parameters. We also emphasize that systematic errors in high SNR events need to be carefully controlled. For example, more ROQ bases may be required in high-SNR events for better accuracy. Further investigations into the accuracy of PE in the 3G era are necessary, and we plan to address this in future work.

As discussed in Sec. IV B, the numbers presented in this work represent an optimistic estimate of extending current analysis techniques to future data. We hope that the methods and results in this work can serve as a reference and baseline benchmark for the research and development of 3G detectors. When novel data analysis methods are developed in the future, whether for more realistic data (e.g. including overlapped signals and noise variations) or a faster speed, their computational cost can be assessed and compared with the results presented here. For analyses at the catalog level, the computing efficiency must be carefully considered, as the cost for processing a single event will be multiplied across the entire 3G catalog.

[1] J. Aasi *et al.* (LIGO Scientific), Class. Quant. Grav. **32**, 074001 (2015), arXiv:1411.4547 [gr-qc].

[2] B. P. Abbott *et al.* (LIGO Scientific, Virgo), Phys. Rev. Lett. **116**, 061102 (2016), arXiv:1602.03837 [gr-qc].

[3] B. P. Abbott *et al.* (LIGO Scientific, Virgo), Phys. Rev. X **9**, 031040 (2019), arXiv:1811.12907 [astro-ph.HE].

[4] R. Abbott *et al.* (LIGO Scientific, Virgo), Phys. Rev. X **11**, 021053 (2021), arXiv:2010.14527 [gr-qc].

[5] R. Abbott *et al.* (LIGO Scientific, VIRGO), Phys. Rev. D **109**, 022001 (2024), arXiv:2108.01045 [gr-qc].

[6] R. Abbott *et al.* (KAGRA, VIRGO, LIGO Scientific), Phys. Rev. X **13**, 041039 (2023), arXiv:2111.03606 [gr-qc].

[7] R. Abbott *et al.* (KAGRA, VIRGO, LIGO Scientific), Astrophys. J. Suppl. **267**, 29 (2023), arXiv:2302.03676 [gr-qc].

[8] B. P. Abbott *et al.* (LIGO Scientific, Virgo, 1M2H, Dark Energy Camera GW-E, DES, DLT40, Las Cumbres Observatory, VINROUGE, MASTER), Nature **551**, 85 (2017), arXiv:1710.05835 [astro-ph.CO].

[9] M. Soares-Santos *et al.* (DES, LIGO Scientific, Virgo), Astrophys. J. Lett. **876**, L7 (2019), arXiv:1901.01540 [astro-ph.CO].

[10] B. P. Abbott *et al.* (LIGO Scientific, Virgo, VIRGO), Astrophys. J. **909**, 218 (2021), arXiv:1908.06060 [astro-ph.CO].

[11] B. P. Abbott *et al.* (LIGO Scientific, Virgo), Astrophys. J. Lett. **882**, L24 (2019), arXiv:1811.12940 [astro-ph.HE].

[12] R. Abbott *et al.* (LIGO Scientific, Virgo), Astrophys. J. Lett. **913**, L7 (2021), arXiv:2010.14533 [astro-ph.HE].

[13] R. Abbott *et al.* (KAGRA, VIRGO, LIGO Scientific), Phys. Rev. X **13**, 011048 (2023), arXiv:2111.03634 [astro-ph.HE].

[14] B. P. Abbott *et al.* (LIGO Scientific, Virgo), Phys. Rev. Lett. **121**, 161101 (2018), arXiv:1805.11581 [gr-qc].

[15] B. P. Abbott *et al.* (LIGO Scientific, Virgo), Phys. Rev. Lett. **118**, 121101 (2017), [Erratum: Phys.Rev.Lett. 119, 029901 (2017)], arXiv:1612.02029 [gr-qc].

[16] B. P. Abbott *et al.* (LIGO Scientific, Virgo), Phys. Rev. Lett. **118**, 121102 (2017), arXiv:1612.02030 [gr-qc].

[17] B. P. Abbott *et al.* (LIGO Scientific, Virgo), Phys. Rev. Lett. **120**, 091101 (2018), arXiv:1710.05837 [gr-qc].

[18] L. Wade, J. D. E. Creighton, E. Ochsner, B. D. Lackey, B. F. Farr, T. B. Littenberg, and V. Raymond, Phys. Rev. D **89**, 103012 (2014).

[19] M. A. Shaikh, V. Varma, H. P. Pfeiffer, A. Ramos-Buades, and M. van de Meent, Phys. Rev. D **108**, 104007 (2023), arXiv:2302.11257 [gr-qc].

[20] A. Ramos-Buades, A. Buonanno, and J. Gair, Phys. Rev. D **108**, 124063 (2023), arXiv:2309.15528 [gr-qc].

[21] C. K. Mishra, K. G. Arun, B. R. Iyer, and B. S. Sathyaprakash, Phys. Rev. D **82**, 064010 (2010), arXiv:1005.0304 [gr-qc].

[22] R. Abbott *et al.* (LIGO Scientific, VIRGO, KAGRA), (2021), 10.48550/arXiv.2112.06861, arXiv:2112.06861 [gr-qc].

[23] J. Veitch *et al.*, Phys. Rev. D **91**, 042003 (2015), arXiv:1409.7215 [gr-qc].

[24] B. Zackay, L. Dai, and T. Venumadhav, (2018), 10.48550/arXiv.1806.08792, arXiv:1806.08792 [astro-ph.IM].

[25] N. Leslie, L. Dai, and G. Pratten, Phys. Rev. D **104**, 123030 (2021), arXiv:2109.09872 [astro-ph.IM].

[26] N. J. Cornish, "Fast Fisher Matrices and Lazy Likelihoods," (2010), arXiv:1007.4820 [gr-qc].

[27] S. Vinciguerra, J. Veitch, and I. Mandel, Class. Quant. Grav. **34**, 115006 (2017), arXiv:1703.02062 [gr-qc].

[28] S. Morisaki, Phys. Rev. D **104**, 044062 (2021), arXiv:2104.07813 [gr-qc].

[29] P. Canizares, S. E. Field, J. Gair, V. Raymond, R. Smith, and M. Tiglio, Phys. Rev. Lett. **114**, 071104 (2015), arXiv:1404.6284 [gr-qc].

[30] R. Smith, S. E. Field, K. Blackburn, C.-J. Haster, M. Pürrer, V. Raymond, and P. Schmidt, Phys. Rev. D **94**, 044031 (2016), arXiv:1604.08253 [gr-qc].

[31] P. Couvares *et al.*, "Gravitational Wave Data Analysis: Computing Challenges in the 3G Era," (2021), arXiv:2111.06987 [gr-qc].

[32] S. Bagnasco *et al.*, EPJ Web Conf. **295**, 04015 (2024), arXiv:2312.11103 [gr-qc].

[33] M. Punturo *et al.*, Class. Quant. Grav. **27**, 194002 (2010).

[34] M. Branchesi *et al.*, JCAP **07**, 068 (2023), arXiv:2303.15923 [gr-qc].

[35] D. Reitze and et al, (2019), 10.48550/arXiv.1907.04833, arXiv:1907.04833 [astro-ph, physics:gr-qc].

[36] M. Evans *et al.*, (2023), 10.48550/arXiv.2306.13745, arXiv:2306.13745 [astro-ph.IM].

[37] S. Borhanian and B. S. Sathyaprakash, Phys. Rev. D **110**, 083040 (2024), arXiv:2202.11048 [gr-qc].

[38] L. S. Finn, Phys. Rev. D **46**, 5236 (1992), arXiv:gr-qc/9209010.

[39] E. Thrane and C. Talbot, Publications of the Astronomical Society of Australia **36**, e010 (2019), arXiv:1809.02293 [astro-ph.IM].

[40] M. J. Williams, J. Veitch, and C. Messenger, Mach. Learn. Sci. Tech. **4**, 035011 (2023), arXiv:2302.08526 [astro-ph.IM].

[41] M. J. Williams, J. Veitch, and C. Messenger, Phys. Rev. D **103**, 103006 (2021), arXiv:2102.11056 [gr-qc].

[42] G. Ashton *et al.*, Astrophys. J. Suppl. **241**, 27 (2019), arXiv:1811.02042 [astro-ph.IM].

[43] A. Buonanno, B. Iyer, E. Ochsner, Y. Pan, and B. S. Sathyaprakash, Phys. Rev. D **80**, 084043 (2009), arXiv:0907.0700 [gr-qc].

[44] N. J. Cornish, Phys. Rev. D **103**, 104057 (2021), arXiv:2101.01188 [gr-qc].

[45] J. Roulet, J. Mushkin, D. Wadekar, T. Venumadhav, B. Zackay, and M. Zaldarriaga, Phys. Rev. D **110**, 044010 (2024), arXiv:2404.02435 [gr-qc].

[46] L. Dai, T. Venumadhav, and B. Zackay, "Parameter Estimation for GW170817 using Relative Binning," (2018), arXiv:1806.08793 [gr-qc].

[47] M. Dax, S. R. Green, J. Gair, N. Gupte, M. Pürrer, V. Raymond, J. Wildberger, J. H. Macke, A. Buonanno, and B. Schölkopf, "Real-time gravitational-wave inference for binary neutron stars using machine learning,"

(2024), arXiv:2407.09602 [gr-qc].

[48] K. W. K. Wong, M. Isi, and T. D. P. Edwards, Astrophys. J. **958**, 129 (2023), arXiv:2302.05333 [astro-ph.IM].

[49] T. Dietrich, A. Samajdar, S. Khan, N. K. Johnson-McDaniel, R. Dudi, and W. Tichy, Phys. Rev. D **100**, 044003 (2019), arXiv:1905.06011 [gr-qc].

[50] S. Husa, S. Khan, M. Hannam, M. Pürrer, F. Ohme, X. Jiménez Forteza, and A. Bohé, Phys. Rev. D **93**, 044006 (2016), arXiv:1508.07250 [gr-qc].

[51] S. Khan, S. Husa, M. Hannam, F. Ohme, M. Pürrer, X. Jiménez Forteza, and A. Bohé, Phys. Rev. D **93**, 044007 (2016), arXiv:1508.07253 [gr-qc].

[52] S. Morisaki, R. Smith, L. Tsukada, S. Sachdev, S. Stevenson, C. Talbot, and A. Zimmerman, Phys. Rev. D **108**, 123040 (2023), arXiv:2307.13380 [gr-qc].

[53] Q. Hu, J. Irwin, Q. Sun, C. Messenger, L. Suleiman, I. S. Heng, and J. Veitch, "Decoding Long-duration Gravitational Waves from Binary Neutron Stars with Machine Learning: Parameter Estimation and Equations of State," (2024), arXiv:2412.03454 [gr-qc].

[54] F. Aubin et al., Class. Quant. Grav. **38**, 095004 (2021), arXiv:2012.11512 [gr-qc].

[55] Q. Hu and J. Veitch, Astrophys. J. Lett. **958**, L43 (2023), arXiv:2309.00970 [gr-qc].

[56] S. E. Field, C. R. Galley, J. S. Hesthaven, J. Kaye, and M. Tiglio, Phys. Rev. X **4**, 031006 (2014), arXiv:1308.3565 [gr-qc].

[57] M. Pürrer, Class. Quant. Grav. **31**, 195010 (2014), arXiv:1402.4146 [gr-qc].

[58] J. Blackman, S. E. Field, C. R. Galley, B. Szilágyi, M. A. Scheel, M. Tiglio, and D. A. Hemberger, Phys. Rev. Lett. **115**, 121102 (2015), arXiv:1502.07758 [gr-qc].

[59] L. M. Thomas, G. Pratten, and P. Schmidt, Phys. Rev. D **106**, 104029 (2022), arXiv:2205.14066 [gr-qc].

[60] R. Smith et al., Phys. Rev. Lett. **127**, 081102 (2021), arXiv:2103.12274 [gr-qc].

[61] G. Morras, J. F. N. Siles, and J. Garcia-Bellido, Phys. Rev. D **108**, 123025 (2023), arXiv:2307.16610 [gr-qc].

[62] J. Lange, R. O'Shaughnessy, and M. Rizzo, "Rapid and accurate parameter inference for coalescing, precessing compact binaries," (2018), arXiv:1805.10457 [gr-qc].

[63] J. Wofford et al., Phys. Rev. D **107**, 024040 (2023).

[64] L. P. Singer and L. R. Price, Phys. Rev. D **93**, 024013 (2016), arXiv:1508.03634 [gr-qc].

[65] Q. Hu, C. Zhou, J.-H. Peng, L. Wen, Q. Chu, and M. Kovalam, Phys. Rev. D **104**, 104008 (2021), arXiv:2110.01874 [gr-qc].

[66] T. Tsutsui, K. Cannon, and L. Tsukada, Physical Review D **103** (2021), 10.1103/physrevd.103.043011.

[67] M. Dax, S. R. Green, J. Gair, J. H. Macke, A. Buonanno, and B. Schölkopf, Phys. Rev. Lett. **127**, 241103 (2021), arXiv:2106.12594 [gr-qc].

[68] M. Dax, S. R. Green, J. Gair, M. Pürrer, J. Wildberger, J. H. Macke, A. Buonanno, and B. Schölkopf, Phys. Rev. Lett. **130**, 171403 (2023), arXiv:2210.05686 [gr-qc].

[69] N. Gupte et al., (2024), 10.48550/arXiv.2404.14286, arXiv:2404.14286 [gr-qc].

[70] I. Kobyzev, S. J. Prince, and M. A. Brubaker, IEEE transactions on pattern analysis and machine intelligence

**43**, 3964 (2020).

[71] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan, Journal of Machine Learning Research **22**, 1 (2021).

[72] H. Gabbard, C. Messenger, I. S. Heng, F. Tonolini, and R. Murray-Smith, Nature Phys. **18**, 112 (2022), arXiv:1909.06296 [astro-ph.IM].

[73] A. Pagnoni, K. Liu, and S. Li, "Conditional variational autoencoder for neural machine translation," (2018), arXiv:1812.04405 [cs.CL].

[74] J. Langendorff, A. Kolmus, J. Janquart, and C. Van Den Broeck, Phys. Rev. Lett. **130**, 171402 (2023), arXiv:2211.15097 [gr-qc].

[75] S. Hild et al., Class. Quant. Grav. **28**, 094013 (2011), arXiv:1012.0908 [gr-qc].

[76] The sampler `nessai` provides total sampling wall time (the wall time to finish the whole analysis) and wall time spent evaluating likelihood. Since likelihood is evaluated in parallel while all other processes use a single CPU core, the total CPU time is calculated as (total sampling wall time-likelihood evaluation wall time)*1 + likelihood evaluation wall time*`ncpu`, which simplifies to the form in the main text.

[77] H. Qi and V. Raymond, Phys. Rev. D **104**, 063031 (2021), arXiv:2009.13812 [gr-qc].

[78] https://github.com/MarinerQ/3gpemethods.

[79] ET MDC1 data are publicly available from http://et-origin.cism.ucl.ac.be.

[80] W. J. Handley, M. P. Hobson, and A. N. Lasenby, Monthly Notices of the Royal Astronomical Society **453**, 4385–4399 (2015).

[81] T. Wouters, P. T. H. Pang, T. Dietrich, and C. Van Den Broeck, Phys. Rev. D **110**, 083033 (2024), arXiv:2404.11397 [astro-ph.IM].

[82] Y. Himemoto, A. Nishizawa, and A. Taruya, Physical Review D **104**, 044010 (2021), arXiv:2103.14816 [astro-ph, physics:gr-qc].

[83] E. Pizzati, S. Sachdev, A. Gupta, and B. Sathyaprakash, Physical Review D **105**, 104016 (2022), arxiv:2102.07692 [astro-ph, physics:gr-qc].

[84] P. Relton and V. Raymond, Physical Review D **104**, 084039 (2021), arXiv:2103.16225 [astro-ph, physics:gr-qc].

[85] P. Relton, A. Virtuoso, S. Bini, V. Raymond, I. Harry, M. Drago, C. Lazzaro, A. Miani, and S. Tiwari, Physical Review D **106**, 104045 (2022), arxiv:2208.00261 [astro-ph, physics:gr-qc].

[86] A. Samajdar, J. Janquart, C. V. D. Broeck, and T. Dietrich, Physical Review D **104**, 044003 (2021), arXiv:2102.07544 [astro-ph, physics:gr-qc].

[87] Q. Hu and J. Veitch, Astrophys. J. **945**, 103 (2023), arXiv:2210.04769 [gr-qc].

[88] S. Bagnasco, "The Ligo-Virgo-KAGRA Computing Infrastructure for Gravitational-wave Research," (2023), arXiv:2311.12559 [astro-ph.IM].

[89] J. Albrecht et al. (HEP Software Foundation), Comput. Softw. Big Sci. **3**, 7 (2019), arXiv:1712.06982 [physics.comp-ph].