

STROOPWAFEL: Simulating rare outcomes from astrophysical populations, with application to gravitational-wave sources*

Floor S. Broekgaarden,^{1–5†} Stephen Justham,^{7,8,1,2,5} Selma E. de Mink,^{6,1,2}
Jonathan Gair,^{9,10,5} Ilya Mandel,^{3,4,11,5} Simon Stevenson,^{4,12}
Jim W. Barrett,¹³ Alejandro Vigna-Gómez,^{11,3,4,5} Coenraad J. Neijssel^{11,14}

¹*Anton Pannekoek Institute for Astronomy, University of Amsterdam, Postbus 94249, 1090 GE Amsterdam, The Netherlands*

²*GRAPPA, University of Amsterdam, Science Park 904, 1098 XH Amsterdam, The Netherlands*

³*Monash Centre for Astrophysics, School of Physics and Astronomy, Monash University, Clayton, Victoria 3800, Australia*

⁴*The ARC Center of Excellence for Gravitational Wave Discovery – OzGrav, Hawthorn VIC 3122, Australia*

⁵*Dark Cosmology Centre, Niels Bohr Institute, University of Copenhagen, Juliane Maries Vej 30, DK-2100, København ø, Denmark*

⁶*Harvard-Smithsonian Center for Astrophysics, 60 Garden Street, Cambridge, MA 02138*

⁷*School of Astronomy & Space Science, University of the Chinese Academy of Sciences, Beijing 100012, China*

⁸*National Astronomical Observatories, Chinese Academy of Sciences, Beijing 100012, China*

⁹*School of Mathematics, University of Edinburgh, The King's Buildings, Peter Guthrie Tait Road, Edinburgh, EH9 3FD, UK*

¹⁰*Max Planck Institute for Gravitational Physics (Albert Einstein Institute), Am Mühlenberg 1, Potsdam-Golm 14476, Germany*

¹¹*Birmingham Institute for Gravitational Wave Astronomy and School of Physics and Astronomy, University of Birmingham Birmingham, B15 2TT, United Kingdom*


¹²*Center for Astrophysics and Supercomputing, Swinburne University of Technology, Hawthorn VIC 3122, Australia*

¹³*Klarna Bank AB (publ). Sveavägen 46, 111 34 Stockholm*

¹⁴*Albert-Einstein-Institut, Max-Planck-Institut für Gravitationsphysik, D-30167 Hannover, Germany*

15 October 2019

ABSTRACT

Gravitational-wave observations of double compact object (DCO) mergers are providing new insights into the physics of massive stars and the evolution of binary systems. Making the most of expected near-future observations for understanding stellar physics will rely on comparisons with binary population synthesis models. However, the vast majority of simulated binaries never produce DCOs, which makes calculating such populations computationally inefficient. We present an importance sampling algorithm, STROOPWAFEL, that improves the computational efficiency of population studies of rare events, by focusing the simulation around regions of the initial parameter space found to produce outputs of interest. We implement the algorithm in the binary population synthesis code COMPAS, and compare the efficiency of our implementation to the standard method of Monte Carlo sampling from the birth probability distributions. STROOPWAFEL finds ~ 25 – 200 times more DCO mergers than the standard sampling method with the same simulation size, and so speeds up simulations by up to two orders of magnitude. Finding more DCO mergers automatically maps the parameter space with far higher resolution than when using the traditional sampling. This increase in efficiency also leads to a decrease of a factor ~ 3 – 10 in statistical sampling uncertainty for the predictions from the simulations. This is particularly notable for the distribution functions of observable quantities such as the black hole and neutron star chirp mass distribution, including in the tails of the distribution functions where predictions using standard sampling can be dominated by sampling noise. 

Key words: gravitational waves – stars: evolution – binaries: general – methods: numerical – methods: statistical

1 INTRODUCTION

The direct detection of gravitational waves originating from merging binary black holes (BHs) and neutron stars (NSs) has opened up a new window on the Universe, and marked the birth of gravitational-wave astrophysics as a new field of research (Abbott et al. 2016,

*STROOPWAFEL: Simulating The Rare Outcomes Of Populations
With AIS For Efficient Learning.

†E-mail: fbroekgaarden@g.harvard.edu

2017). At the time of writing the first two observing runs of Advanced LIGO and Virgo have been completed (The LIGO Scientific Collaboration et al. 2018). A few dozen detections are expected during the third observing run, and we can anticipate hundreds of detections per year when the next generation of detectors with higher sensitivities come online (Abbott et al. 2018).

The detections are starting to reveal the properties of the population of merging binary BHs and NSs. The distributions of the inferred masses and spins contain valuable information about their origin. Distinguishing different theories for their formation and learning about the complex physical processes that govern the lives of their possible massive-star progenitors requires comparing observed populations with theoretical predictions.

Theoretical simulations of the population of merging DCOs are challenging because gravitational-wave events represent an extremely rare outcome of binary evolution. From a thousand massive binary systems typically only of order one, or less, yields a double compact object. A meaningful comparison with population observations requires simulating a statistically significant sample of events. When sampling from the birth distributions, which is a form of sampling commonly used in binary population synthesis, this often means we need to sample at least many millions of binary systems. For example, Kruckow et al. (2018) find that their DCO merger rates converge only when simulating $N \geq 3 \times 10^8$ binaries (and for BH–BH mergers their statistical noise remains at the 2 percent level even with $N = 10^9$ samples).

To make it feasible to simulate such large numbers of systems, all present-day simulations pay a high price. In many studies computational speed is ensured by using highly approximate algorithms that treat the physical processes in a simplified way. Another way to keep the computational costs reasonable is to limit the total number of simulations, restricting the exploration of the impact of the uncertain physical input assumptions beyond a few variations.

The recent detections bring to light a further challenge as we start to ask questions about rare subsets of the already rare gravitational-wave events. One example is the subset of heavy binary black hole mergers with total system masses in excess of $50 M_{\odot}$. Such systems produce loud GW signals and can thus be observed over large volumes (Fishbach & Holz 2017). The majority of currently observed BH–BH mergers have total masses above $50 M_{\odot}$ (The LIGO Scientific Collaboration et al. 2018). However, they are very rare in simulations of binaries sampled from the expected distribution of initial conditions. Most current theoretical predictions for the extreme tails of the mass distribution are heavily affected by Poisson noise, resulting from under sampling. A second example of an astrophysically important rare subset of the rare gravitational-wave events are the NS–NS systems that merge within about 50 Myrs from the moment the NS–NS is formed. Early NS–NS mergers are important as they are candidate sources for the observed early r-process enriched ultra-faint dwarf galaxies such as Reticulum II and Tucana III (e.g. Safarzadeh et al. 2019). A third example are the subset of BH–NS mergers with sufficiently similar components masses that there is significant tidal ejection from the merger to produce electromagnetic counterparts (Foucart et al. 2018). Obtaining statistically accurate predictions for the extreme tails of the distribution functions for rare but astrophysically important subpopulations is currently a challenge for most simulations.

Earlier studies have proposed improvements of the efficiency of population synthesis studies. Kolb (1993) and later Politano (1996) implemented a transformation function using Jacobian matrices to map known birth rates of cataclysmic variables directly into present day populations. Kalogera (1996) adopted this method, developed

an analytical model for the kick prescription and showed that it is possible to obtain similar expressions for several observable distribution functions (Kalogera & Webbink 1998; Kalogera 2000). More recently, Andrews et al. (2018) implemented Markov Chain Monte Carlo (MCMC) methods to efficiently simulate populations of binaries matching specified evolutionary endpoints, whereas Barrett et al. (2017) and Taylor & Gerosa (2018) use Gaussian process emulators to predict outputs of the binary population synthesis model for parametrised choices of physical assumptions that have not been simulated. However, current binary population synthesis models have complex output functions containing natural bifurcations (e.g., small changes in the initial mass of a star can lead to drastic changes in the final mass in the simulations). Moreover, binary population synthesis simulation output spaces often contain stochastic behaviour (e.g., due to the randomly drawn neutron star natal kick). Such discontinuities pose a challenge for MCMC methods and Gaussian process regression emulators, as they rely on a certain smoothness in order to converge and produce independent samples.

In this paper we present a new algorithm, STROOPWAFEL^{*,1}. We have designed STROOPWAFEL to improve the efficiency of simulations of rare astrophysical events, and so to enable accurate simulations of populations of extremely rare events at reasonable computational cost. The algorithm first explores the initial parameter space until it finds a preliminary population of systems of interest. This exploration is done by stochastically sampling from the birth distributions. STROOPWAFEL then concentrates the later sampling towards regions in the initial parameter space that are in the vicinity of the initial parameters of the interesting binaries found during the exploration phase. This is an example of “Adaptive Importance Sampling” (AIS), described further in the next section.

We focus here on the application of the study of DCO mergers as gravitational-wave sources, but our algorithm is much more broadly applicable. The user can specify any target population of interest. An advantage of the algorithm is that it can handle the bifurcations and stochastic behaviour that naturally occur in the physical prescriptions in binary population synthesis simulations, and which lead to discontinuous output surfaces. Finally, with our algorithm we can easily derive the uncertainties on the estimated parameters which can be a challenge for sampling methods such as MCMC and Gaussian Process regression emulators.

This paper is organized as follows. In Section 2 we describe the algorithm and provide expressions for how to calculate statistical estimates and the uncertainties from the simulations. We further derive the optimal relative duration of the exploratory phase and the total number of simulations, given the rareness of the target population. In Section 3 we provide a demonstration of our algorithm. We apply it to population synthesis simulations of double compact object mergers. In Section 4 we outline caveats and future directions for further improvement and refinement of the algorithm. We conclude and summarise in Section 5.

2 METHOD

Our algorithm is conceptually simple. It uses a strategy that may be familiar from playing the classic game *Battleship*. The aim of

¹ All data that has been used for this study, accompanied with a Jupyter notebook that reproduces the main results, can be found here [\[link\]](#). The code for the STROOPWAFEL algorithm will be made publicly available after acceptance. Early inquiries can be addressed to the lead author.

this game is to guess the coordinates of the ships of the other player, which are placed on a regular discrete grid. Most players will probably start with an exploration phase randomly trying different coordinates. After one or more successful “hits” most players will change their search strategy and instead try to refine their search by trying coordinates that are close to the successful hits until they uncover the full location and orientation of the ship. It has been shown that this is a more successful strategy compared to searching completely randomly throughout the entire game (e.g. Jones 1977).

Our algorithm, STROOPWAFEL, follows a conceptually similar strategy, but instead of aiming to win a game of Battleship the algorithm is designed to improve the efficiency for simulating populations of rare events (that is, rare outcomes from the space of initial conditions). Successful hits in this analogy are finding systems of interest that are part of a certain target population. These may be systems that result in DCO mergers or anything that the user specifies. We improve the efficiency by focusing on areas of the initial parameter space near to those which produced outcomes of interest during a prior, exploratory, sampling phase. Instead of Monte Carlo sampling from the birth probability distributions, STROOPWAFEL uses information from that exploratory sampling phase to create an alternative distribution function, from which it then samples.

This class of Monte Carlo methods is generically called “Importance Sampling” (Kahn & Harris 1951; Kahn & Marshall 1953). Since we do not know in advance which areas of the initial parameter space should receive extra attention we use *adaptive* importance sampling (AIS), for which see, e.g., Torrie & Valleau (1977); Hesterberg (1995); Ortiz & Pack Kaelbling (2013); Pennanen & Koivu (2006); Cappé et al. (2004); Cornuet et al. (2012). The nature of the AIS algorithm makes it straightforward to tune the focus of the simulation on a specific target population or function of interest (Cappé et al. 2008). Such AIS algorithms also allow for straightforward calculations of the sampling uncertainties. The STROOPWAFEL implementation of AIS is similar to that in Cornuet et al. (2012), but includes a new method to guide the fraction of the computational effort that should be spent on the exploratory phase (see Sect 2.2.4).

Whilst the concept of our algorithm is not complicated, there are some mathematical details involved in making the implementation efficient and robust. For example, some of the subtlety involved is in making sure not to concentrate *too* closely on locations which previously led to success. If we only look exactly where we have looked before, then we don’t learn anything new.

Section 2.1 introduces binary population synthesis as a mapping between input and output parameter spaces, along with some notation which is useful for the description of our algorithm. Section 2.2 presents the key details of STROOPWAFEL. We explain how we shape the adaptive sampling distribution from the information found in an initial exploratory phase, and how to optimally combine the samples from both the exploratory and adapted phases to estimate the population quantities of interest. We also describe how STROOPWAFEL self-consistently determines how long the exploratory phase should last as a fraction of the simulation time, based on continually updated estimates of the rareness of the target population. Section 2.3 illustrates the practical characteristics of our AIS algorithm in an idealised way, providing an explanatory summary of the behaviour of STROOPWAFEL for users who do not wish to learn all the mathematical details.

2.1 Definition of concepts and symbols

Binary population synthesis models the population observables for a particular class of event, under a set of assumptions about the

physics. Predicting such an output population typically involves simulating many individual systems from their initial conditions. Only a small fraction of those simulated systems may produce outcomes which are of interest for that study.

Selecting which specific points in the input space (i.e., the initial conditions) to simulate into the output space (i.e., the observables) is a key part of population synthesis. This process is called sampling, and must appropriately take into account the relative frequency of different initial conditions. Ideally, it should also efficiently explore the initial parameter space. Examples of these initial parameters are the initial masses of the two stars, $m_{1,i}$ and $m_{2,i}$, and the initial separation a_i between the two stars. For a given initial composition, these three dimensions are often regarded as adequate initial conditions. However, more generally, these input conditions may be distributed over many dimensions.

Each initial binary system $\mathbf{x}_i \in \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ in a binary population synthesis can thus be written as $\mathbf{x}_i = (m_{1,i}, m_{2,i}, a_i, \dots)$, which has a combined birth distribution

$$\pi(\mathbf{x}_i) = \pi(m_{1,i}, m_{2,i}, a_i, \dots). \quad (1)$$

This distribution of initial conditions is often taken as the Monte Carlo sampling distribution². In practice, simulations of binary-star populations which aim to study outcomes such as mergers between BHs and NSs do not sample from the full range of initial conditions of stellar binaries. Such simulations ignore stars whose mass is too low to produce a BH or NS, which is a simple form of importance sampling. The normalisation of π actually used for the sampling is then corrected to take this into account when predicting event rates.

For each initial binary, \mathbf{x}_i , the final state of the binary \mathbf{y}_f is determined using the binary population synthesis model u ,

$$\mathbf{y}_f = u(\mathbf{x}_i). \quad (2)$$

In many cases a simulation is run to study binaries that evolve to a certain target subtype T , e.g., maybe T is the population of binary black hole mergers. The following indicator function describes whether a binary \mathbf{y}_f is of interest:

$$\mathbb{1}_T(\mathbf{y}_f) := \begin{cases} 1 & \text{if } \mathbf{y}_f \in T \text{ (a hit)} \\ 0 & \text{if } \mathbf{y}_f \notin T \text{ (a miss),} \end{cases} \quad (3)$$

which equals 1 if \mathbf{x}_i simulates to the target binary system T (a hit) and zero if not (a miss). Combining equations (2) & (3) gives the function

$$\phi(\mathbf{x}_i) := \mathbb{1}_T(u(\mathbf{x}_i)), \quad (4)$$

which is a shorthand notation to describe whether an initially drawn binary evolved into a binary of the target population.

The samples from the initial parameter space that produced a binary of the target population (i.e., $\phi(\mathbf{x}_i) = 1$) can then be given by the set

$$\mathbf{x}_T := (m_{1,T}, m_{2,T}, a_T, \dots). \quad (5)$$

At the end of a simulation, the properties of the model population and the statistical uncertainties on those predicted properties

² Binary population synthesis simulations often sample in the initial mass ratio, $m_{2,i}/m_{1,i}$, rather than the mass of the initially least-massive star. This includes the simulations we present later; see also Appendix C. Whether sampling in $m_{2,i}$ or the mass ratio it has also been common to assume that the initial parameters are independent of each other, e.g.: $\pi(\mathbf{x}_i) = \pi(m_{1,i}) \cdot \pi(m_{2,i}/m_{1,i}) \cdot \pi(a_i) \cdot \dots$

This assumption of separability may not be valid, as found by Abt et al. (1990); Moe & Di Stefano (2017); Klencki et al. (2018).

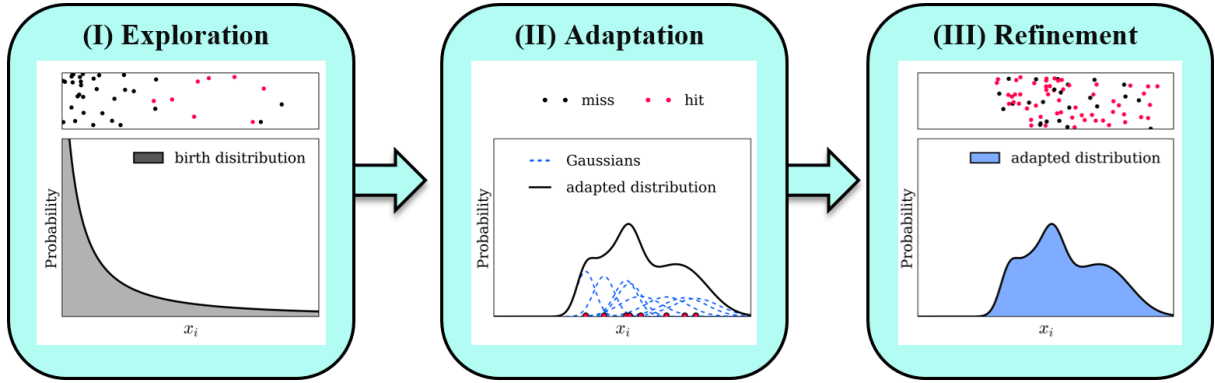


Figure 1. Illustration of the STROOPWAFEL algorithm. In the algorithm, (I) we first draw random binaries from the birth distribution π until a small population of events of interest (hits) is found. (II) We construct Gaussian distributions around each of the previously found successful events. We create an adapted instrumental distribution $q(\mathbf{x})$ from the mixture of Gaussian distributions. We scale the width of each Gaussian with the local sampling density. (III) We draw the remaining samples from this adapted distribution which focuses around the target population. The top panels show a random draw of samples from the corresponding distribution in the lower panel. The samples are assigned a random scatter in the y -direction for the visualization.

can both be determined using the standard Monte Carlo estimator (Fermi & Richtmyer 1948; Metropolis & Ulam 1949). For example, the relative formation rate of the target population, \mathcal{R}_T , is estimated with

$$\mathbb{E}_\pi[\mathcal{R}_T] \approx \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}_i), \quad (6)$$

where N is the total number of samples used, \mathbb{E} is the notation for the estimated mean and the subscript π in \mathbb{E}_π denotes that the samples \mathbf{x}_i are distributed following the birth distribution π (cf. Eq. 1). We shall refer to this relative formation rate, \mathcal{R}_T , throughout this section. Mathematically it is a fractional volume from the initial binary parameter space, weighted by the probability of forming a binary system at each part of that initial parameter space. This quantity is not a physical rate, but gives a formation rate for the population of interest as a fraction of the total number of initial binary systems formed. So it only differs from a true formation rate by a physical normalisation. We consider it appropriately intuitive to keep referring to this as a fractional, or relative, rate.

2.2 Adaptive sampling algorithm to increase efficiency of simulation

Our algorithm consists of three main steps, as illustrated in Fig. 1:

(I) Exploration

We first explore the parameter space by sampling directly from the birth distribution π until eventually a sufficient population of events of interest is found.

(II) Adaptation

We construct multivariate Gaussian distributions in the initial parameter space around each of the events of interest found during the exploration phase. We scale the widths of each of the Gaussians with the local sampling density. We create the adapted sampling distribution q , from here on referred to as the *instrumental* distribution, by combining the Gaussians into a mixture distribution.

(III) Refinement

We draw the samples for the remaining simulations from this instrumental distribution. Each sample is assigned a weight so that the predicted population appropriately reflects the birth distribution π .

The rest of this subsection explains these steps in more detail.

2.2.1 The instrumental distribution

When the exploratory phase has ended, the set of binaries of the target population \mathbf{x}_T (i.e., hits) contains $N_{T,\text{expl}}$ binaries that were found using a total of N_{expl} samples. In the STROOPWAFEL algorithm these samples are then used to create an adapted instrumental distribution $q(\mathbf{x})$, which is focused around the areas in the initial parameter space that produced the binaries of interest during the exploratory phase. The remaining binaries are thereafter sampled from this instrumental distribution. To obtain unbiased estimates of the target population, weights w_i are incorporated for each sample as is standard in importance sampling

$$w_i = \frac{\pi(\mathbf{x}_i)}{q(\mathbf{x}_i)}, \quad (7)$$

where π is the distribution of initial conditions, as given in Eq. (1).

The instrumental distribution $q(\mathbf{x})$ can be chosen to be any probability distribution function, but a robust instrumental distribution is characterized by the following criteria:

- The weights w_i are always finite and well defined. That is, $q(\mathbf{x}_i) = 0$ implies $\pi(\mathbf{x}_i)\phi(\mathbf{x}_i) = 0$ for all i .
- The instrumental distribution is efficient if $q(\mathbf{x})$ is close to the (unknown) target distribution of the binary population synthesis study, i.e., when the instrumental distribution $q(\mathbf{x})$ is proportional to $|\phi(\mathbf{x})|\pi(\mathbf{x})$, as shown by Kahn & Marshall (1953).
- It should be computationally inexpensive to generate random samples from $q(\mathbf{x})$ as well as to calculate the probability $q(\mathbf{x}_i)$ for each sample \mathbf{x}_i .

In order to achieve these properties the instrumental distribution, $q(\mathbf{x})$, in STROOPWAFEL is chosen to be a mixture³ of $N_{T,\text{expl}}$ Gaussian distributions q_k given by

$$q(\mathbf{x}) = \frac{1}{N_{T,\text{expl}}} \sum_{k=1}^{N_{T,\text{expl}}} q_k(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (8)$$

³ In this context, “mixture” has a standard mathematical meaning. A sample drawn from q is drawn from each Gaussian q_k with a probability $N_{T,\text{expl}}$ instead of taking the sum of normally distributed samples. The sum of two jointly normally distributed random variables will still have a normal distribution (even if the means are not the same) whereas a *mixture* of two normally distributed variables will have two peaks (assuming the means are far enough apart).

where each q_k contributes $1/N_{\text{T,expl}}$ to the mixture distribution.

However, when drawing from $q(\mathbf{x})$, some samples will fall outside the physical range of the parameter space Ω (e.g., when drawing a binary with a negative stellar mass). Such samples can immediately be rejected and redrawn. By doing so, we in practice sample from the normalized physical mixture distribution

$$\widetilde{q}(\mathbf{x}) = \frac{1}{(1 - F_{\text{rej}})} q(\mathbf{x}) \mathbb{1}_{\Omega}(\mathbf{x}) \quad (9)$$

where $\mathbb{1}_{\Omega}(\mathbf{x})$ is the indicator function that equals 1 when the sample \mathbf{x} lies in the physical range of the parameter space and 0 if not. The factor F_{rej} is the fraction of samples from $q(\mathbf{x})$ that are drawn outside of the physical parameter space. The factor $1/(1 - F_{\text{rej}})$ thus corrects for the normalization of $\widetilde{q}(\mathbf{x})$. It is computationally inexpensive to draw samples from $\widetilde{q}(\mathbf{x})$ since F_{rej} can be estimated once with a Monte Carlo simulation, and one can draw randomly from each Gaussian $q_k(\mathbf{x})$ separately.

The Gaussian distributions $q_k(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ are parametrized by their means $\boldsymbol{\mu}_k$ and covariance matrices $\boldsymbol{\Sigma}_k$. The covariance matrix, $\boldsymbol{\Sigma}_k$, determines the width of the Gaussian distributions. We adopt a diagonal covariance matrix

$$\boldsymbol{\Sigma}_k = \begin{bmatrix} \sigma_{1,k}^2 & 0 & \dots \\ 0 & \ddots & \\ \vdots & & \sigma_{d,k}^2 \end{bmatrix}, \quad (10)$$

where d is the number of dimensions of the initial parameter space.

We scale the width of each Gaussian, given by the covariance matrix $\boldsymbol{\Sigma}_k$, with the average distance to the next sampled binary \mathbf{x}_i in the initial parameter space, estimated via the local density of the prior distribution π . This allows the algorithm to construct broader Gaussian distributions (i.e., with larger σ_k) around previously found hits that lie in the regions of the parameter space that are less densely explored. If the one dimensional marginalised prior of the j -th parameter is π_j then the standard deviation $\sigma_{j,k}$ is given by:

$$\sigma_{j,k} = \kappa \frac{1}{\pi_j(x_k) N_{\text{expl}}^{1/d}}, \quad (11)$$

where N_{expl} represents the number of samples used for the exploration phase, the power $1/d$ scales this number to the effective number of samples per dimension and the factor $\pi_j(x_k)$ scales the width to the density of samples in the exploration phase around the previously found hit x_k . We also introduce a free parameter, κ , that scales the widths of the Gaussian distributions. This enables us to regulate how tightly the mixture of Gaussian distributions covers the parameter space near the successful binaries \mathbf{x}_T . In this paper we adopt $\kappa = 2$, which we chose following tests with a toy model (for which, see Section 2.2.5 and Appendix B).

2.2.2 Combining samples from the exploratory phase and refinement phase

It is desirable to make use of the samples from both the exploration and refinement sampling phases. The optimal way to achieve this is somewhat subtle. In principle this could be done by merging the samples and weights into a combined estimate (see Chapter 14 Robert & Casella 2013). However, Veach & Guibas (1995), Hesterberg (1995), and later Owen & Zhou (2000) showed that using *deterministic multiple mixture weights* is an efficient and robust way of combining the samples. This approach uses the fact that the combined samples from the exploratory phase and refined sampling

symbol	description
u	binary population synthesis model
\mathbf{x}_i	initial state of a binary system
\mathbf{y}_f	final state of a binary system
\mathbf{T}	target subpopulation of binaries of interest
\mathbf{x}_T	set of hits from the exploratory phase
$\pi(\mathbf{x})$	distribution of the initial conditions
Ω	physical parameter range of the simulation
$q(\mathbf{x})$	instrumental distribution
$\widetilde{q}(\mathbf{x})$	normalized instrumental distribution
N	total number of samples in the simulation
N_{expl}	number of binaries in the exploratory phase
N_{ref}	number of binaries in the refinement phase
N_T	number of systems of the target population
$N_{\text{T,expl}}$	number of hits in the exploratory phase
f_{expl}	fraction of samples in the exploration phase
κ	scale factor for the widths of the Gaussians
w_i	statistical weight of sample \mathbf{x}_i
\widetilde{w}_i	recomputed statistical weight

Table 1. Summary of the parameters that are used throughout this paper. Hits refer to binaries of the chosen target population.

phase can be represented by a mixture sampling distribution $Q(\mathbf{x})$ from both phases

$$Q(\mathbf{x}) = f_{\text{expl}} \pi(\mathbf{x}) + (1 - f_{\text{expl}}) \widetilde{q}(\mathbf{x}), \quad (12)$$

where $f_{\text{expl}} = N_{\text{expl}}/N$ is the fraction of samples spent on the exploratory phase. By analogy with Eq. (7), the weights of all the N samples can be recalculated with

$$\widetilde{w}_i = \frac{\pi(\mathbf{x}_i)}{Q(\mathbf{x}_i)}. \quad (13)$$

This use of deterministic multiple mixture weights is not fundamental to STROOPWAFEL. Our motivation for using deterministic multiple mixture weights is conservative, to increase the stability against potential sampling artefacts. One of the samples drawn from the importance function $q(\mathbf{x})$ may occasionally be extremely large. Such extreme weights could remain so large as to be problematic when merging the samples with the original weights w_i – no matter how efficient the sampling is in the exploratory phase (Cornuet et al. 2012). Use of deterministic multiple mixture weights suppresses this potential rare difficulty.

The multiple mixture weights approach ignores the distribution from which a given draw was sampled. This does not affect the estimators for the predicted values, although it does introduce a very small bias to the uncertainty estimators, which we confirmed to be negligible in our toy model tests. Recalculating the weights in this way yields comparable or better estimates than those which are obtained when merging the samples or using inverse-variance weighting for our adaptive importance sampling algorithm. Indeed, He & Owen (2014) derived a bound for the variance of the *balance heuristic* for such estimators that combine samples from different distributions and found that this is an efficient way of combining samples. See also sect. 3 in Veach & Guibas (1995) and sect. 2 in Cornuet et al. (2012) for a more detailed discussion.

2.2.3 Calculating statistical estimates using the adaptive distribution

At the end of each run the properties of the target population, such as the rate of formation \mathcal{R}_T of members of the target population, and distribution functions of population observables, can be determined

by standard Monte Carlo estimates. Because we have drawn the samples from a different distribution than the birth distribution we have to incorporate weights to make sure that the estimators for these quantities reflect the correct formation probabilities. For example, the relative formation rate \mathcal{R}_T of the target population within the simulation is estimated by the mean

$$\mathcal{R}_T \approx \mathbb{E}_Q[\mathcal{R}_T] = \frac{1}{N} \sum_{i=1}^N \phi(\mathbf{x}_i) \widetilde{w}_i. \quad (14)$$

The uncertainty in this rate is represented by the variance about the mean by

$$\mathbb{V}_Q[\mathcal{R}_T] = \frac{s_Q^2[\phi(\mathbf{x})\widetilde{w}(\mathbf{x})]}{N} \approx \frac{\sum_i^N \phi(\mathbf{x}_i)^2 \widetilde{w}_i^2}{N} - \mathcal{R}_T^2. \quad (15)$$

where s_Q is the (sample) standard deviation for samples drawn from the mixture sampling distribution $Q(\mathbf{x})$. This equation is known as the asymptotic variance. The other moments or statistical estimates for the target population can be similarly calculated.

2.2.4 The duration of the exploratory phase

An important choice in adaptive importance sampling algorithms is deciding when to switch from the exploratory phase to the refinement phase. This choice can have a substantial impact on the performance of the algorithm. Leaving the exploratory phase too early can result in missing important regions of the initial parameter space which produce systems in the target population. On the other hand, switching to the refinement sampling phase too late will miss out on the advantages of the algorithm, as most time will be spent sampling from the birth distributions instead of the more efficient adapted distribution.

A method often used in adaptive importance sampling algorithms to determine the fraction of samples that should be spent on exploring the parameter space is by using the effective sample size: $(\sum_k w_k)^2 / \sum_k (w_k^2)$, (ESS, [Hesterberg 1995](#); [Liu 2008](#)). This is a measure of efficiency and corresponds to the equivalent number of independent samples drawn from the prior distribution. However, it can be difficult to know in advance what a good value for the ESS should be. Instead, since STROOPWAFEL is a two-step adaptive algorithm, we can directly derive a value for f_{expl} by using the estimated rareness of the population, which we self-consistently calculate during the exploration phase.

Here we estimate the optimal fraction f_{expl} of the total number of samples N that we should spend on the exploratory phase. The challenge is that we do not know in advance how good our adaptive sampling distribution will be. Here, as a simplified proxy for the imperfect sampling distribution, we assume that the adaptive sampling distribution is determined sufficiently well that it perfectly matches the target distribution over most of the parameter space, but that a small region of the target parameter space could be missing samples due to a limited number of samples drawn during the exploratory phase, and thus have an adaptive sampling probability of zero (see [Appendix A](#) for details). In other words, we divide the volume of the input parameter space which successfully produces systems of interest into two parts, one which we assume we have accurately found and one which remains missing. We then find f_{expl} such that the event rate uncertainty is minimised; specifically, we require that the contribution to the event rate from potentially undiscovered islands is smaller than, or no worse than similar to, the sampling uncertainty in the rate contributed by the islands which are successfully found.

We assume that we sample from the mixture distribution $Q(\mathbf{x})$,

and aim to estimate the rate \mathcal{R}_T with N total samples. After simulating N samples we assume we have identified a target binary forming region with total weight z_1 , whereas a region with weight z_2 is yet undiscovered such that the estimated rate of the target population is

$$\mathcal{R}_T = \underbrace{z_1}_{\text{identified}} + \underbrace{z_2}_{\text{unidentified}} \approx \mathbb{E}_Q[\mathcal{R}_T]. \quad (16)$$

The uncertainty on this rate estimate is described by the variance, which we can approximate with (see [Appendix A](#) for more details)

$$\mathbb{V}_Q[\mathcal{R}_T] \approx \frac{1}{N} \left[\frac{z_1}{\frac{(1-f_{\text{expl}})}{z_1} + f_{\text{expl}}} + \frac{z_2}{f_{\text{expl}}} - (z_1 + z_2)^2 \right]. \quad (17)$$

The smallest uncertainty on the rate \mathcal{R}_T is obtained for the value of f_{expl} where the variance $\mathbb{V}_Q(\mathcal{R}_T)$ is the lowest. By taking the derivative of \mathbb{V}_Q with respect to f_{expl} , and then finding the roots of the derivative, we find that this minimum occurs at

$$f_{\text{expl}} = 1 - \frac{z_1(\sqrt{1-z_1} - \sqrt{z_2})}{\sqrt{1-z_1}(\sqrt{z_2(1-z_1)} + z_1)}. \quad (18)$$

To make practical use of this during our simulations, we need ongoing live estimates for z_1 and z_2 . For the region which has been identified, we adopt $z_1 \approx \mathbb{E}_\pi[\mathcal{R}_T]$ during the simulation using [Eq. \(14\)](#). We approximate the target population region that is yet undiscovered, z_2 , by $z_2 \approx \frac{1}{(f_{\text{expl}}N)}$. This represents the weight of stochastic sampling noise in the exploratory phase when using a total of $f_{\text{expl}}N$ samples. Moreover, this choice of f_{expl} ensures that the estimated uncertainty on the rate estimate is always comparable or larger than the uncertainty of missing a region, $1/f_{\text{expl}}N$.

While the running estimate of z_2 on the right hand side of [Eq. 18](#) is a function of f_{expl} , rather than explicitly solving for f_{expl} , we choose to iteratively approach the optimal solution over the course of the exploratory phase. The resulting f_{expl} is similar to the f_{expl} that is obtained if we had used the dependency of z_2 on f_{expl} when solving for the minimum in [Eq. 17](#).

We note that we have implicitly assumed that the adaptive sampling phase is perfectly efficient, i.e., that all the draws in the adaptive sampling phase find a member of the target output population. Here that is a conservative assumption, since a less-than-perfect efficiency will increase the sampling uncertainty with respect to the known islands (i.e., z_1). Therefore, imperfect efficiency in the adaptive sampling phase decreases the chance that the uncertainty from undiscovered islands is significant.

STROOPWAFEL internally uses [Eq. 18](#) to estimate f_{expl} . For clarity here, we can additionally assume that z_1 and z_2 are much smaller than 1, i.e., that the target population is a rare outcome of the initial conditions. Then we obtain the simplified equation

$$f_{\text{expl}} \approx 1 - \frac{z_1}{z_1 + \sqrt{z_2}}. \quad (19)$$

From this simplified equation it becomes clear that we recover the intuitively correct limit for extremely rare events, i.e. if $z_1 = \mathbb{E}_Q[\mathcal{R}_T] \rightarrow 0$, we find $f_{\text{expl}} \rightarrow 1$, which suggests that we should spend all our simulation time on exploration, as expected. On the other hand once we find at least 1 hit in the exploratory phase we find $f_{\text{expl}} \neq 1$, and so the variance of our rate estimate is expected to decrease when drawing some of the samples from the adapted distribution, compared to taking all samples from the birth distribution.

Lastly, from [Eq. 17](#) it also becomes clear that $f_{\text{expl}} \approx 0.5$ once we have found $N_{T,\text{expl}} \sim \sqrt{N_{\text{expl}}}$ target binaries. In other words,

$f_{\text{expl}} \sim 1$ if $N \leq 1/\mathcal{R}_T^2$ and therefore the total number of samples N should generally be similar to or larger than $1/\mathcal{R}_T^2$.

2.2.5 Determining the free parameter κ from tests with a toy model

We present results from the application of our method to astrophysical simulations in Section 3. Here we explore the methodology with a toy model to test the performance of the algorithm and determine the value of the free parameter κ . In principle $\kappa = 1$ could be adopted. Smaller values of κ will increase the efficiency of the STROOPWAFEL algorithm, but increase the chance of missing an important region of the output surface because the Gaussian distributions q_k are too narrow and do not cover the output surface well. Excessively large values of κ , meanwhile, will decrease the efficiency of finding samples of interest in the refinement phase and lower the gain of STROOPWAFEL. After performing tests with a toy model, as described in Appendix B, we adopt the value $\kappa = 2$. However, when applying STROOPWAFEL in higher dimensions the optimal value for κ may well change.

2.2.6 Summary of STROOPWAFEL algorithm

The algorithm for STROOPWAFEL, combining the methods discussed in this section, is summarized in Algorithm 1.

2.3 Characteristic behaviour

Here we use the analytic derivations from Section 2.2 to illustrate the characteristics of our algorithm in idealised cases. This is intended to help users of STROOPWAFEL understand the expected behaviour without needing to master the details presented above.

Key variables for this illustration are:

- \mathcal{R}_T , the formation rate of the population under study. This is expressed as the fraction of binaries, when drawn from initial conditions following the birth probability distribution, that yield target systems.
- N , the total number of binary systems (i.e. samples) used in a given population simulation, which is chosen by the user.
- f_{expl} , the fraction of the total number of samples that should optimally be spent on the exploration phase. This is chosen automatically by STROOPWAFEL during the exploration phase (see Sect. 2.2.4), when the algorithm estimates the formation rate \mathcal{R}_T .

Figure 2 presents derived quantities as a function of the fractional rate of the target population \mathcal{R}_T . This Figure, and Figure 3, include points representing simulated astrophysical populations, specifically for different subsets of DCO mergers. These simulations are described further in Section 3. The values from those simulations are included here to give context to the analytic expectations for the performance of STROOPWAFEL.

The top panel of Figure 2 shows the optimal f_{expl} , for simulations with a total number of samples of $N = 10^4, 10^5, 10^6$ and 10^7 . For a rarer target population, a larger fraction of the total number of samples should be spent on the exploratory phase. This is because it takes longer to determine a good sampling distribution when \mathcal{R}_T is low. A more common target population can be optimally simulated with a relatively small exploratory phase, since we expect this will be enough to build up a good adaptive distribution. STROOPWAFEL estimates \mathcal{R}_T during the exploration phase, when it samples from the birth probability distribution and

Algorithm 1: STROOPWAFEL algorithm

```

1  $i = 0$ ;
2
3 (I) Exploration:
4  $f_{\text{expl}} = 1$ ;
5 while  $i/N \leq f_{\text{expl}}$  do
6    $i += 1$ ;
7   draw new sample  $x_i \sim \pi(x)$ ;
8   evaluate sample  $y_f = u(x_i)$ ;
9   if  $y_f \in T$  then
10    counthits  $+= 1$ ;
11     $\mathbf{x}_T \leftarrow \mathbf{x}_i$  (add hit to the found collection of hits);
12    update estimate  $f_{\text{expl}}$  iteratively using Eq. 18
         $z_1 = \mathbb{E}_\pi[\mathcal{R}_T]$  and  $z_2 = \frac{1}{(f_{\text{expl}} N)}$ ;
13   end
14 end
15
16 (II) Adaptation:
17 set  $\boldsymbol{\mu} = \mathbf{x}_T$ ;
18 Calculate  $\boldsymbol{\Sigma}$  by determining  $\sigma_{j,k}(\mathbf{x}_k)$  for all
     $k = 1, \dots, N_T$ ;
19 This gives  $q(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ ;
20
21 (III) Refinement:
22 while  $N_{\text{expl}} \leq i \leq N$  do
23    $i += 1$ ;
24   draw new sample  $x_i \sim q(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathbb{1}_\Omega(\mathbf{x})$ ;
25   evaluate sample  $y_f = u(x_i)$ ;
26 end
27
28 Post processing:
29 calculate  $F_{\text{rej}}$  and mixture weights  $\tilde{w}_i = \frac{\pi(\mathbf{x}_i)}{Q(\mathbf{x}_i)}$ ;
30 calculate desired population quantities such as the rate
     $\mathcal{R}_T$ ;

```

self-consistently calculates f_{expl} based on the estimated \mathcal{R}_T and user-chosen N .

The bottom panel of Figure 2 shows the expected statistical uncertainty in the event rate predicted by the population simulation, for simulations with a total number of samples of $N = 10^4$ and 10^6 . By statistical uncertainty we mean the uncertainty that arises from using a finite number of samples. In standard Monte Carlo simulation this is also sometimes referred to as Poisson error, given for traditional sampling by $1/\sqrt{N_T}$. We estimate the minimum STROOPWAFEL uncertainty in Fig. 2 as $1/\sqrt{N_T}$, where N_T is now the total number of objects of interest found jointly in the exploration and refinement phases. The uncertainty as computed through Eq. 15 will be slightly greater, since the distribution of weights means that the effective sample size of the target population is lower than N_T . These analytic estimates appear to be consistent with our numerical calculations (see also Section 3.1).

This is not the physical uncertainty since the model used for the simulation might still be wrong. In practice the efficiency in the refinement phase will not be perfect, i.e., not all samples drawn in that phase will find an outcome from the target population. So the expected uncertainty from STROOPWAFEL will lie in the shaded region shown in Fig. 2. STROOPWAFEL efficiency gains will be greatest for rare events and large N , allowing a greater fraction of time to be spent in the efficient refinement phase.

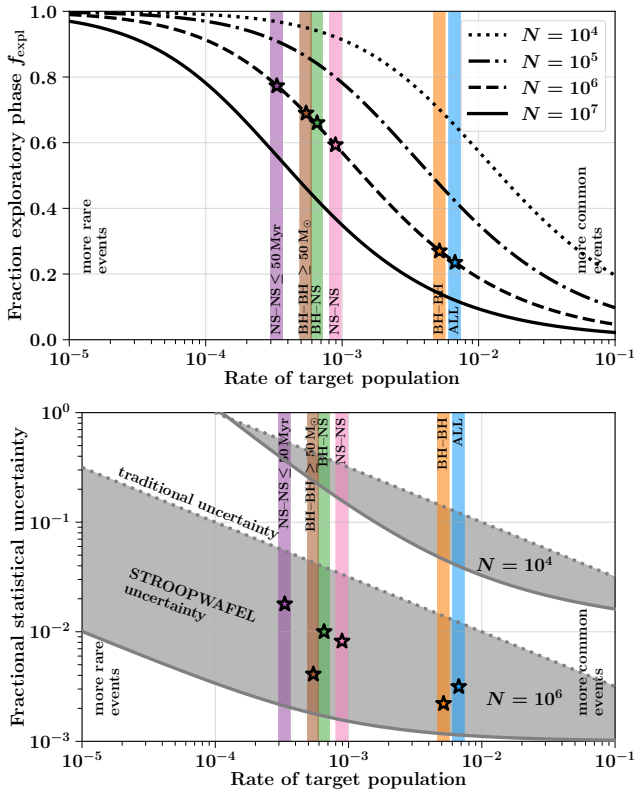


Figure 2. **Top panel:** the fraction of total samples that should be spent on the exploratory phase versus the fractional rate, \mathcal{R}_T , of the target population in a simulation. Different curves show the effect of varying the total number of samples N in the simulation. **Bottom panel:** expected sampling uncertainty on the predicted event rate versus the fractional rate, for two different choices of N . The dashed lines show the expected uncertainty from “traditional” Monte Carlo sampling (i.e., $\sqrt{\nabla \pi[\mathcal{R}_T]/\mathcal{R}_T}$). The solid lines show the minimum possible statistical uncertainty from STROOPWAFEL sampling, i.e., if the efficiency in the refinement phase is 1. In practice the statistical rate uncertainty provided when using STROOPWAFEL will lie in the area shaded in grey. **All panels:** coloured vertical bars indicate the fractional rate of the target populations, and star symbols show the corresponding values of these parameters, for the six simulations with $N = 10^6$ described in Section 3.

Comparisons to observational data will typically be made using distribution functions of predicted quantities (e.g., component masses), not just event rates. We later demonstrate the improvements provided by our algorithm for predictions of distribution functions. Nonetheless this overall decrease in statistical rate uncertainty for fixed sample number in a simulation is indicative of the improvements enabled by applying STROOPWAFEL to a target population.

Figure 3 shows the increase in the number of simulated binaries of interest versus traditional Monte Carlo sampling from a birth distribution for a simulation with fixed $N = 10^6$. The efficiency in the refinement phase is not known in advance. We show predictions for a refinement phase efficiency of 1 and 0.1; as long as the total number of successful samples is dominated by those drawn during the refinement phase, the maximum possible gain is roughly proportional to the refinement phase efficiency. The value for the efficiency of the refinement phase varies between $\sim 3.4 \cdot 10^{-2}$ and $3.7 \cdot 10^{-1}$ in our example astrophysical simulations (see Section 3 and Table 2).

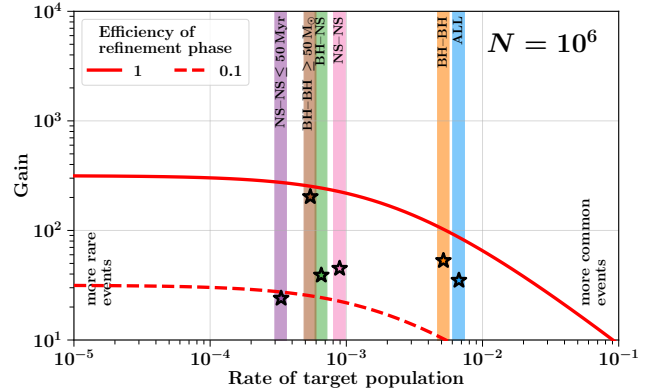


Figure 3. The ratio of the number of target binaries found with STROOPWAFEL to the number found when Monte Carlo sampling from the birth distribution – i.e., the multiplicative gain achieved by STROOPWAFEL – as a function of rareness of the target population. The red curves give this gain for two different efficiencies in the refinement phase, as labelled. The gain is shown for simulations with a total of $N = 10^6$ samples. Coloured vertical bars and star symbols present the values for the six simulations described in Section 3, as given in the last column of Table 2.

3 RESULTS

In this section we demonstrate the power and advantages of STROOPWAFEL. Our algorithm could be applied to many sampling routines, but the illustration here uses the binary population synthesis code COMPAS (Stevenson et al. 2017; Barrett et al. 2018; Vigna-Gómez et al. 2018). The physical assumptions and parameter settings we adopt are briefly summarised in Appendix C.

We combine STROOPWAFEL and COMPAS to model six different target populations. Four are simulations of subtypes of DCOs that merge in a Hubble time: (1) all DCO mergers (i.e., BH–NS, NS–NS and BH–BH), (2) BH–BH mergers, (3) NS–NS mergers and (4) BH–NS mergers. Additionally we model two simulations of extremely rare events by focusing on a subset of the above, namely (5) BH–BH mergers with total system masses in excess of $m_{\text{tot}} \geq 50 M_\odot$ and (6) NS–NS mergers that merge within $t_c \leq 50$ Myrs from the moment of the DCO formation (where t_c is the coalescence time). A summary of the results for these simulations can be found in Table 2.

In this section we present detailed results from simulations 1–4, as these target populations are those most commonly discussed in the literature. We just present the key findings for simulations 5 and 6. For each target population we compare a simulation using our sampling algorithm to one which uses birth distribution Monte Carlo sampling, which for conciseness we will typically call traditional sampling. Both the STROOPWAFEL and the traditional simulations sample $N = 10^6$ initial binaries.

The overall gain that is obtained when using STROOPWAFEL depends on the simulation and the initial efficiency of the ‘traditional’ method that the algorithm is compared with. For example, the choices for the initial parameter space can change how much the sampling can be improved when using STROOPWAFEL. We use settings that are commonly used in population synthesis studies of DCO mergers.

The remainder of the section is structured as follows. Section 3.1 demonstrates the increased efficiency of STROOPWAFEL at finding binaries from the target population. Section 3.2 discusses how that increased efficiency can be used to speed up simulations.

Section 3.3 shows how our algorithm produces better resolution of the target population. Section 3.4 describes how our sampling method leads to smaller statistical uncertainties in predicted population distribution functions. Section 3.5 shows how STROOPWAFEL becomes even more important for recovering tails of distribution functions and when considering observational bias. Section 3.6 discusses how STROOPWAFEL handles well the bifurcations and discontinuities in the binary population synthesis parameter space.

3.1 On the gain of generating binaries of the target distribution

We find that the number of binaries of these target populations increases by factors of about 25 – 200 when using our STROOPWAFEL sampling algorithm compared to simulations with traditional sampling. The panels in Fig. 4 showcase this by presenting the number of systems formed from the target population as a function of total number of sampled binaries, for both our sampling method and traditional birth distribution Monte Carlo sampling. For these four target populations the gains are between ~ 35 –55. For the two additional extremely rare target populations we find gains of 24 and 203. The gains are also shown in the last column of Table 2, and Figure 3.

At the beginning of each simulation, during the STROOPWAFEL exploratory phase, the two sampling methods produce similar number of binaries of interest (i.e., hits), only different by random chance. The duration of that exploratory phase is determined by f_{expl} , as derived in Section 2.2.4. For our simulated target populations, using $N = 10^6$ samples, f_{expl} ranges between ≈ 0.2 and 0.8 (see f_{expl} in Table 2). The algorithm then switches to the more focused refinement phase, using the information from the hits found during the exploratory phase. During this refinement phase our sampling algorithm is 45 – 650 times more efficient at finding hits (see the sixth column in Table 2).

The difference in efficiency gains between the populations originates mostly from two effects. First, the different rarenesses of the target populations (e.g., for these assumptions BH–NS mergers are more rarely produced than BH–BH mergers) influences how much the efficiency increases during the refinement phase and also the duration of the exploratory phase. Both are important factors for determining the overall gain in efficiency. Second, the structure of the output surfaces influences how well the Gaussian mixture distribution covers the regions of interest in the output space. A more stochastic or discontinuous output space (e.g., many small islands of hits) will lead to smaller efficiencies in the refinement phase of STROOPWAFEL. This effect is most noticeable in the different gains between the BH–BH systems with total masses over $50 M_{\odot}$ and the NS–NS systems that merge within 50 Myrs. Supernova remnants receive a random natal kick in our simulations. These kicks induce stochastic and discontinuous behaviour into the output surfaces, leading to lower refinement-phase efficiency and gain. The kicks are typically larger for NSs than BHs, and so affect NS–NS simulations the most. Conversely, the gain is greatest for the BH–BH merger populations which are least affected by these stochastic kicks (as shown in Figs. 2 and 3).

The largest overall gain shown in Table 2, and Fig. 3 is for modelling a BH–BH population with total mass over $50 M_{\odot}$. For this population it would be possible to increase the efficiency of Monte Carlo sampling from the birth distributions by making thoughtful changes to the boundaries of the initial parameter space, so the factor of ≈ 200 improvement we show for STROOPWAFEL is higher than

would arise when comparing to more carefully-targeted use of standard Monte Carlo sampling in this case. However, well-informed choices in the initial parameter space would also benefit STROOPWAFEL by increasing the efficiency of the initial exploration phase. Moreover, one-dimensional cuts to the parameter space would become increasingly inefficient when sampling in higher dimensions. STROOPWAFEL automatically finds the regions of interest, and avoids the risk of incorrect choices in restricting the initial parameter space.

The increase in the number of events decreases the sampling uncertainty in the predicted event rates. Although the standard uncertainty from Poisson noise decreases with the square root of the number of target systems found, i.e., as $1/\sqrt{N_T}$, in our weighted sampling case it also depends on the variance in the weights (see Eq. 15). We find that our sampling algorithm results in $\approx 3 - 10.5$ times smaller sampling uncertainties compared to traditional sampling for the same total of samples $N = 10^6$. This is presented in Figure 5, which shows the fractional statistical uncertainty estimate on the rate estimate, i.e., $\sqrt{\mathbb{V}[\mathcal{R}_T]}/\mathbb{E}[\mathcal{R}_T]$ from each simulation.

3.2 Speeding up simulations

Instead of using STROOPWAFEL to obtain more information from a simulation with the same number of samples, one could alternatively aim for a certain precision in the predicted event rates. In that case STROOPWAFEL can be used to speed up the simulation, since this precision will be reached using a fraction of the number of samples required when using traditional sampling. Traditional sampling would require 25 – 200 more simulations than STROOPWAFEL to achieve the same number of target binaries, and a factor of around 10 – 100 times more simulations to achieve the same rate estimate precision (these speed-up factors differ because the statistical uncertainty depends on the variance in the weights as well as the number of target samples).

The speed-up factor further depends on the computational cost of simulating samples from the chosen distribution. It might be that the binaries of interest require more or less computational time than other binaries. Therefore the speed-up when using the adaptive distribution Q depends on the science case of interest. In the simulations performed for this study the average computational cost (in CPU time) of simulating typical individual binaries sampled from the adaptive distribution Q was up to a factor of 2 smaller than for individual binaries sampled from the birth distribution. Therefore, the total speed-up was up to another factor of 2 larger in our simulations than from more efficient sampling alone.

More generally, we note that the gain or relative speed-up from using STROOPWAFEL will depend on the target population and the traditional method with which it is compared. First of all, the speed up from STROOPWAFEL will generally be greater (smaller) if the target population is more (less) rare. This is shown in Fig. 3. Equivalently, if one chooses a larger initial parameter space (e.g. sampling $m_{1,i}$ from $[1, 150] M_{\odot}$ instead of $m_{1,i}$ from $[5, 150] M_{\odot}$ used here), the gain would have been larger as the event of interest becomes rarer (assuming no binaries in the extended range form a binary of the target population). Secondly, in some binary population synthesis studies the primary mass is sampled uniformly in $\log m_{1,i}$ space. This is a form of importance sampling. The gain of using STROOPWAFEL (with uniform sampling in $\log m_{1,i}$ during the exploratory phase) could be lower than gains without this importance sampling if importance sampling makes the traditional Monte Carlo more efficient. Nevertheless, we would still expect a significant gain from STROOPWAFEL - especially if using that

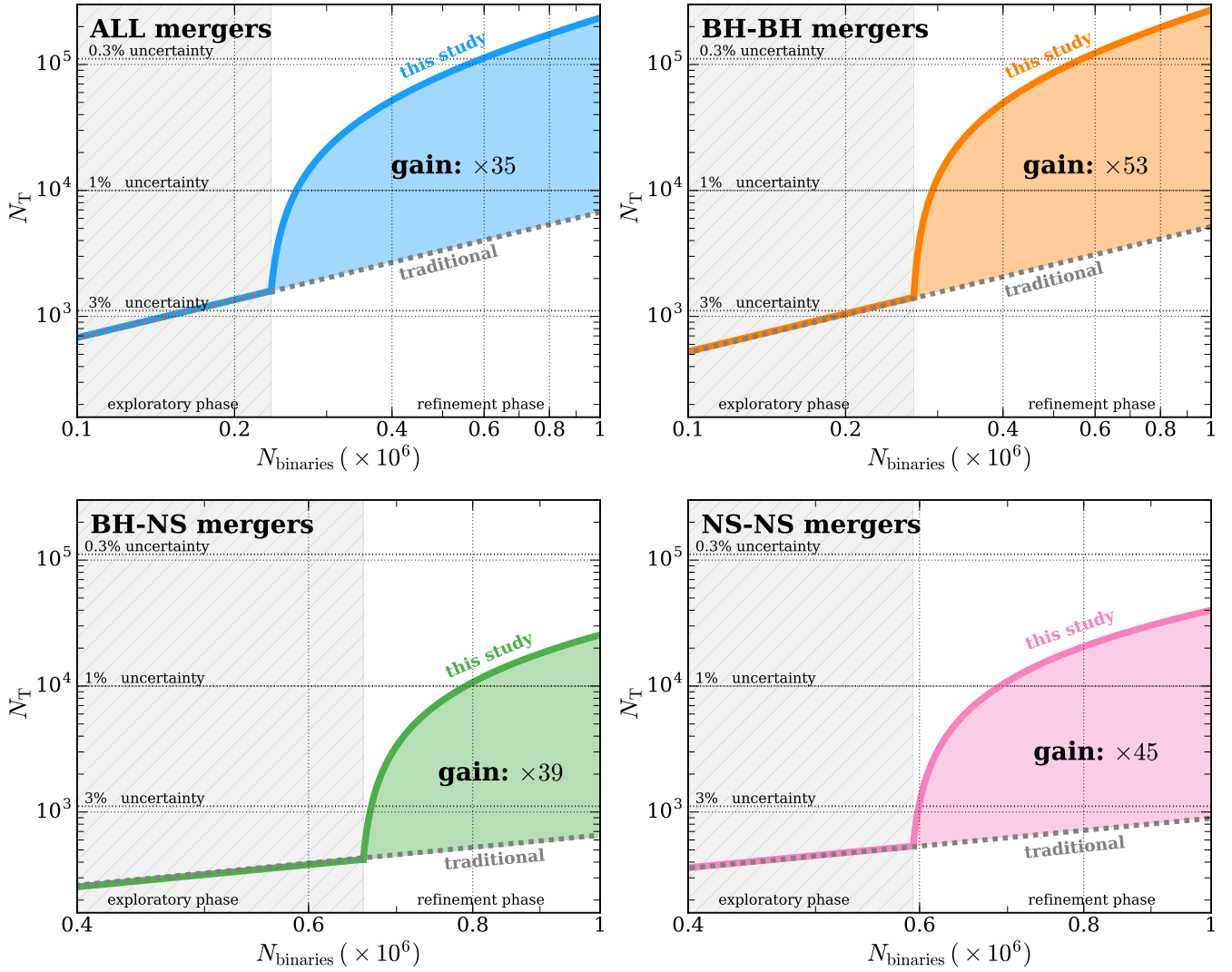


Figure 4. The number of simulated binaries N_T falling into the target population as a function of the total number of binaries N_{binaries} sampled for the traditional sampling method (gray dashed line) and the sampling method presented in this study (solid coloured line). The four panels show the simulations for each of the four target sub-populations. In each panel the duration of the exploratory phase is shown with a hashed gray area. In the background the standard Poisson fractional uncertainties of 0.3, 1 and 3% are shown with dashed lines.

nr	Target subpopulation	f_{expl}	efficiency exploratory	efficiency refinement	gain refinement	N_T traditional	N_T STROOPWAFEL	gain overall
1	All DCO mergers in a Hubble time	0.23	$6.78 \cdot 10^{-3}$	$3.05 \cdot 10^{-1}$	45×	$6.71 \cdot 10^3$	$2.35 \cdot 10^5$	35×
2	BH–BH mergers in a Hubble time	0.27	$5.25 \cdot 10^{-3}$	$3.69 \cdot 10^{-1}$	70×	$5.16 \cdot 10^3$	$2.71 \cdot 10^5$	53×
3	BH–NS mergers in a Hubble time	0.66	$6.36 \cdot 10^{-4}$	$7.38 \cdot 10^{-2}$	116×	$6.55 \cdot 10^2$	$2.55 \cdot 10^4$	39×
4	NS–NS mergers in a Hubble time	0.59	$9.03 \cdot 10^{-4}$	$9.71 \cdot 10^{-2}$	108×	$8.93 \cdot 10^2$	$4.00 \cdot 10^4$	45×
5	BH–BH mergers $m_{\text{tot}} \geq 50 M_{\odot}$	0.69	$5.45 \cdot 10^{-4}$	$3.55 \cdot 10^{-1}$	651×	$5.44 \cdot 10^2$	$1.10 \cdot 10^5$	203×
6	NS–NS mergers with $t_c \leq 50$ Myr	0.77	$3.43 \cdot 10^{-4}$	$3.38 \cdot 10^{-2}$	99×	$3.32 \cdot 10^2$	$7.95 \cdot 10^3$	24×

Table 2. Summary of the results from six target populations that are modelled in this paper to demonstrate our STROOPWAFEL algorithm. We list the fraction of samples spent in the exploratory phase, f_{expl} , and the efficiency of finding ‘hits’ in the exploratory and refinement phases. The gain in refinement is the ratio between the efficiency of finding samples of the target population during the refinement phase of STROOPWAFEL and traditional sampling (where the efficiency of traditional sampling is equal to the efficiency of the STROOPWAFEL exploratory phase). $N_{T,\text{traditional}}$ and $N_{T,\text{STROOPWAFEL}}$ represent the total number of systems of interest that are found by the end of the simulation (using a total of 10^6 samples). The last column is the overall gain that we found when using STROOPWAFEL compared to traditional Monte Carlo sampling from the birth distributions, which is defined by the ratio $N_{T,\text{STROOPWAFEL}} / N_{T,\text{traditional}}$.

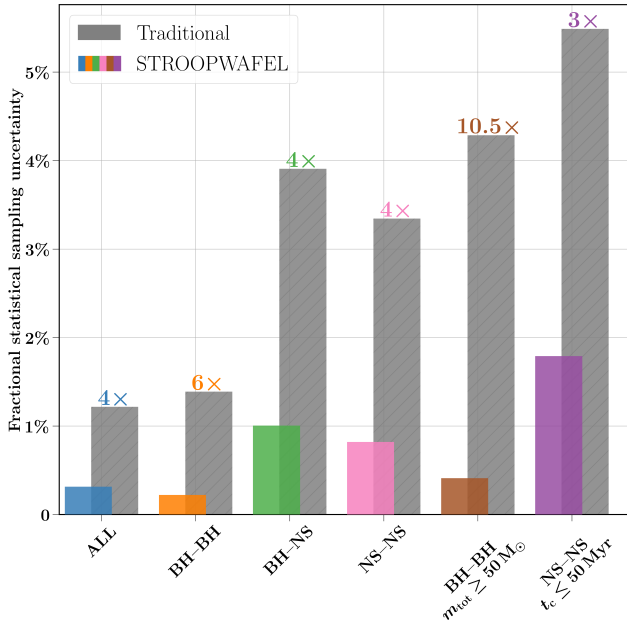


Figure 5. Sampling uncertainty estimate from each simulation of the target population. Gray bars show the uncertainty from traditional sampling methods whereas the coloured bars show the uncertainty from our sampling method STROOPWAFEL. All simulations use a total of 10^6 samples. The number shown on top of the traditional bar shows the factor in decrease in uncertainty from STROOPWAFEL compared to traditional sampling for that simulation.

form of importance sampling in the exploratory phase significantly decreases the duration of this phase.

3.3 Mapping the parameter space with higher resolution

The increase in computational efficiency from STROOPWAFEL leads to finding substantially more events of the target population which naturally enables a much higher resolution mapping of both the input and output parameter spaces. Figures 6 and 7 show examples of how the parameter space is explored in far greater detail using our sampling method compared to traditional birth distribution Monte Carlo sampling.

Figure 6 shows the location of the target population in the initial parameter space of primary mass $m_{1,i}$ and separation $\log a_i$ at birth. With our sampling method we obtain more detailed contours and more contour levels that map the initial parameter space with higher resolution. This leads to better knowledge of the initial conditions of a binary system that yield a binary of the target population. Physically the structures seen in the input parameter space correspond to the assumed physics of the different formation channels leading to compact-object mergers. More details are discussed in Stevenson et al. (2017); Vigna-Gómez et al. (2018).

Meanwhile, Fig. 7 shows the higher resolution mapping from STROOPWAFEL for the output space of the final masses of the compact objects $m_{1,f}$ and $m_{2,f}$ in each DCO. We plot on top the gravitational-wave events found from O1 and O2 data from The LIGO Scientific Collaboration et al. (2018) and Zackay et al.

(2019)⁴. The simulations with our STROOPWAFEL algorithm again yield higher resolutions and more systems of the target populations in the regions overlapping with the observations. This is important in order to compare observations and theory and test the physical assumptions in our models.

Figure 7 shows that even in the STROOPWAFEL simulation, there are relatively few samples consistent with the 90% credible regions for some of the highest mass gravitational-wave events observed in O1 and O2. The samples shown here are not weighted for the sensitivity of gravitational-wave interferometers. In addition, they correspond to a particular model chosen in our simulation, and in practice many variations of this model need to be explored for a meaningful comparison with observations (e.g., Barrett et al. 2018). For example, the metallicity in this model is fixed to $Z = 0.001$ and we expect to form more massive black holes at lower metallicities (e.g., Spura et al. 2015; Belczynski et al. 2016; Neijssel et al. 2019). In addition, the most massive BH–BH mergers may have formed through a different formation channel than the classical isolated binary evolution via the common-envelope phase which is simulated with COMPAS. One example of such a different formation channel is chemically homogeneous evolution, explored by Mandel & de Mink (2016); de Mink & Mandel (2016); Marchant et al. (2016). Another example is dynamical formation of BH–BH mergers in a dense stellar environment (see Mandel & Farmer 2018; Mapelli 2018 for reviews). Another possible explanation is that some of these highest-mass events are instead from instrumental origin as they also have a relatively high false-alarm-rate (The LIGO Scientific Collaboration et al. 2018; Zackay et al. 2019).

The BH–NS and NS–NS panels in Fig. 7 show discontinuities in the NS remnant mass, with obvious gaps in the NS mass range. This is a consequence of discontinuities in the delayed Fryer et al. (2012) model describing the mapping from the carbon–oxygen core mass to the remnant (NS) mass, which is used in the Fiducial COMPAS model. Although such discontinuities may be physical, this mapping does not reproduce the mass distribution of Galactic NS–NS binaries (Vigna-Gómez et al. 2018). For the purposes of this paper, our main message is how such features in model populations can be clarified by improved sampling.

3.4 Smaller variances in distribution functions

The most important consequence of the increase in sampling efficiency enabled by STROOPWAFEL is that it leads to a significant decrease in the statistical uncertainty of the predictions for the output parameter spaces. That is, it improves the precision in the predicted population observables.

Figure 8 illustrates this improvement. The left panel in Fig. 8 shows the number of binaries of the target population N_T found within a certain chirp mass bin for the BH–NS simulation. The chirp mass $m_{\text{chirp}} = (m_{1,f} m_{2,f})^{3/5} / (m_{1,f} + m_{2,f})^{1/5}$ is a combination of the masses of the DCOs that is particularly accurately measured with gravitational-wave observations. For the same histogram bin widths, i.e., the same DCO chirp mass resolution, our improved sampling leads to more binaries of the target distribution per bin, and hence yields smaller fractional sampling uncertainties for each histogram bin. This is shown in the right panel of Fig. 8, which displays the normalised chirp mass distributions from traditional and STROOPWAFEL sampling. The error bars showing the

⁴ Publicly available data can be found at <https://www.gw-openscience.org/catalog/GWTC-1-confident/html/>.

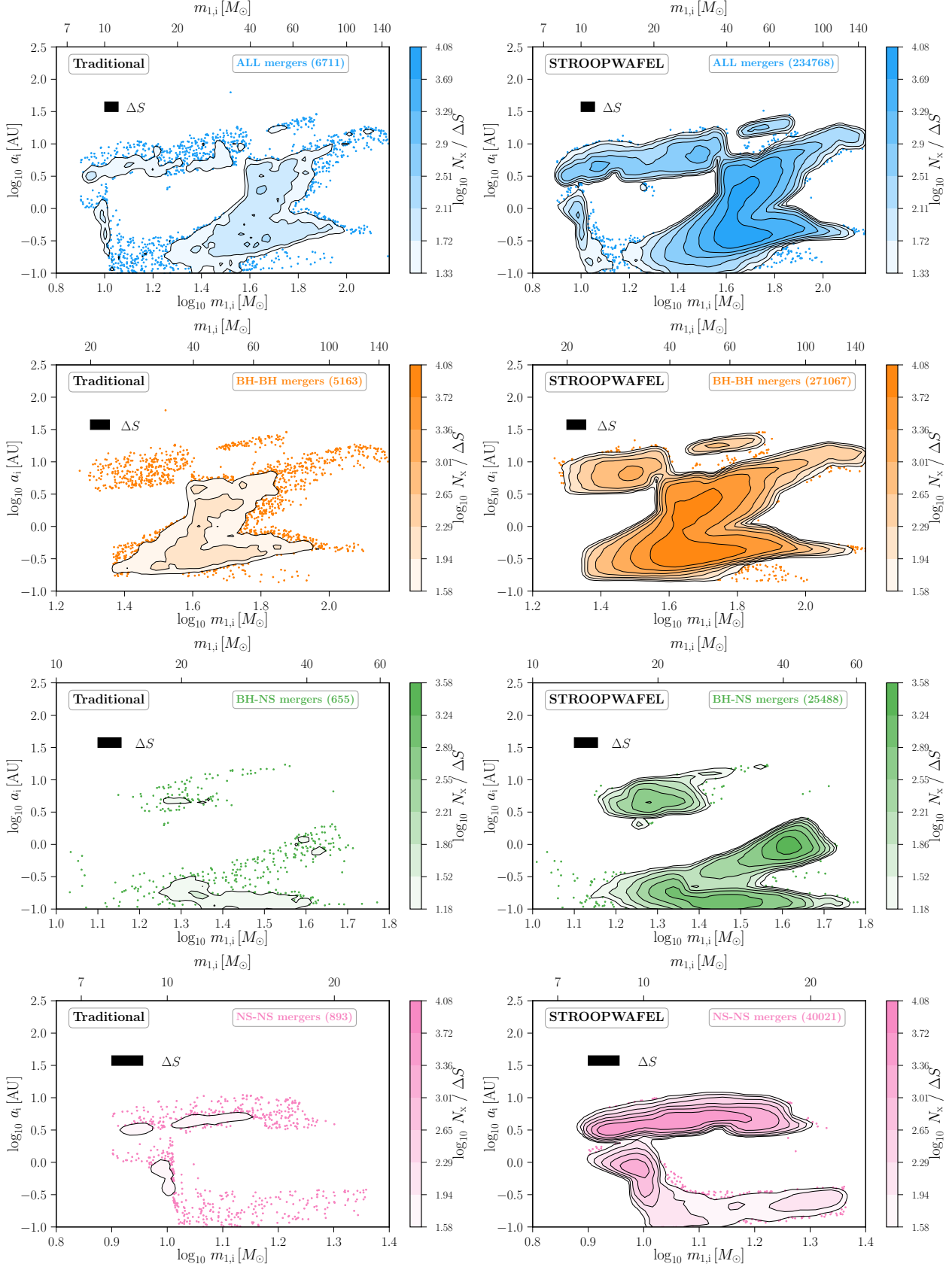


Figure 6. Contour plots of the locations in $\log m_{1,i}$ and $\log a_i$ space of the hits \mathbf{x}_T (i.e., binaries of the target population) found in each simulation when using traditional birth distribution Monte Carlo sampling (left panels) and the sampling method STROOPWAFEL developed in this study (right panels). Contours represent a constant density of binaries of the target population found per unit area in $\log m_{1,i} - \log a_i$ space. The colour gradient indicates the number of samples per area ΔS , the size of which is shown with a black rectangle. If the density is below the level of our lowest contour we plot the individual points. The four different panels from top to bottom represent the first four target populations shown in Table 2. The total number of hits N_T found in each simulation is quoted in parentheses. The metallicity assumed in all simulations is $Z = 0.001$.

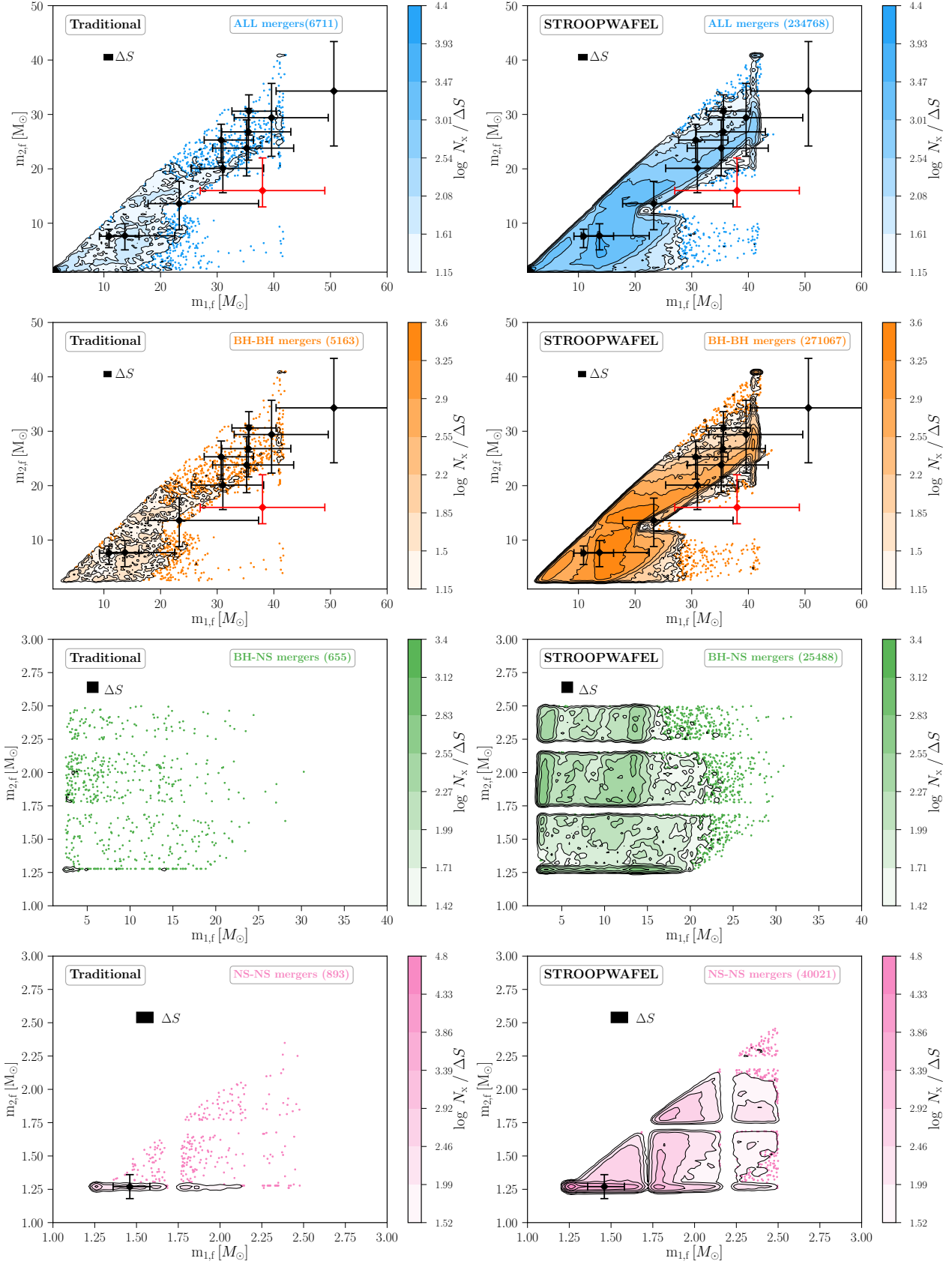


Figure 7. Similar to Figure 6 but now for the output parameters: the final compact object masses $m_{1,f}$ and $m_{2,f}$ of the DCO. We overplot the gravitational-wave observations from O1 and O2 from The LIGO Scientific Collaboration et al. (2018) in black and from Zackay et al. (2019) in red. Error bars indicate 90% credible regions around the median. The metallicity assumed in all simulations is $Z = 0.001$; selection effects of gravitational-wave detectors are not accounted for.

statistical sampling uncertainty on each bin are much smaller for our sampling algorithm, leading to better predictions for the distribution functions.

3.5 Recovering tails of distribution functions

The need for more efficient sampling algorithms in binary population synthesis simulations of gravitational-wave source progenitors becomes even more evident when we consider the observational biases of gravitational-wave detectors. These biases, which generally favour high-mass systems with greater gravitational-wave amplitudes, over-emphasise the rare and frequently under-sampled tails of the simulated distributions. An example is shown in Fig. 9, where we plot the predicted distribution of chirp masses for BH–NS systems estimated using traditional birth distribution Monte Carlo or by using the STROOPWAFEL algorithm. The shown distributions are weighted by the sensitivity of gravitational-wave interferometers, approximated as a bias dependent on the primary DCO mass $\propto m_{1,f}^{2.2}$ (Fishbach & Holz 2017). We also show 1- and 2- σ confidence intervals which are calculated by bootstrapping the samples 1000 times. Our algorithm produces much smoother distribution predictions with much smaller sampling uncertainties compared to traditional sampling methods for the same number of samples simulated. In particular, Figure 9 demonstrates that simulations using the traditional Monte Carlo sampling from the birth distributions under-sample the high-mass end of the population. This will be particularly significant when comparing population models to observations.

Figure 9 corresponds to a particular model choice; variations of the model have to be considered in order to compare with observations. The displayed distribution is from a simulation at a single metallicity of $Z = 0.001$, while a range of metallicities will contribute to the observed BH–NS merger population. An integration over the metallicity-dependent cosmic star formation history is therefore required. The properties of BH–NS mergers will be explored with COMPAS in future work.

3.6 Handling bifurcations and stochasticity

One of the most important results is that our sampling algorithm STROOPWAFEL handles well the bifurcations and stochasticity that naturally occur in the parameter spaces of binary population synthesis simulations. This discontinuous behaviour is visible in Figs. 6 and 7 by the disconnected contours and the offset of the location of some individual points from those contours. Such offset points physically relate to extremely rare formation channels or tails of distribution functions, while the ridges in the birth parameter space relate to bifurcations in the fate of the binary. Not only does STROOPWAFEL recover the irregularly shaped structures in the parameter space, our algorithm also finds these more scattered points.

4 DISCUSSION

We have demonstrated that the performance of STROOPWAFEL is substantially superior to traditional Monte Carlo sampling from the birth probability distributions. For the types of rare events simulated in Section 3, the gain is already so large that the current implementation of our algorithm can contribute to drastic speed-ups of binary population synthesis simulations. Hence we have postponed some natural improvements to STROOPWAFEL until later, but we discuss

them here. After those, we discuss additional potential applications for our algorithm.

4.1 The exploratory phase

During the exploratory phase in STROOPWAFEL the initial parameter space is sampled by drawing random binaries from the priors (as in traditional birth distribution Monte Carlo sampling) until N_{expl} events of interest are found. There are several features of the exploratory phase that could be optimised and improved.

- We now use sampling from the birth distribution π for drawing the random binaries in the exploratory phase. Future improvements of STROOPWAFEL could use more efficient sampling algorithms in the exploratory phase. Examples include (1) using importance sampling in the exploratory phase when there is an existing guess at a more efficient sampling distribution, or (2) implementing techniques such as Latin hypercube sampling (LHS, McKay et al. 1979; Eglajs & Audze 1977; Iman et al. 1980, 1981). LHS is a Monte Carlo method that generates near-random samples which are more equally distributed throughout the initial parameter space by placing only one sample in each row and column of the Latin square⁵. By doing so, it could improve the sampling in the exploratory phase as the probability of the randomly drawn samples being clustered decreases slightly.
- The duration of the exploratory phase is now determined with f_{expl} , which is optimised for the uncertainties on the rates of the target distribution. A future improvement would be to determine f_{expl} based on the uncertainty in the simulated output distribution function. See also Sect. 4.3.
- The exploratory phase duration is optimised under the simplifying assumption that the instrumental distribution will match the target distribution except for some missing regions in parameter space. The optimisation could also consider the level of fluctuation in the instrumental sampling distribution (i.e., the variance in the weights).
- If the structures in the parameter space have a known minimum volume, we could use this to derive a better informed estimate for the uncertainty contributing from the probability of missing such structures in the exploratory phase. This seems unlikely to apply to binary-star population synthesis, but might be relevant for other applications of STROOPWAFEL.

4.2 The refined sampling phase

The Gaussians which are used to form the instrumental distribution are currently constructed using diagonal covariances (Σ_k). These then remain unchanged throughout the refinement phase of adaptive importance sampling – even though much more information becomes available about the distribution of hits in the initial parameter space. A potential future improvement is to update the instrumental sampling distribution during the refinement phase. In principle this might be done locally, with only the samples drawn from each of the individual Gaussians used to update the corresponding element of the instrumental distribution. Doing so would avoid a potentially expensive nearest-neighbour search, as the tree is automatically built for free by the sampling already being performed. See also for example the AMIS algorithm described in Cornuet et al. (2012).

⁵ A Latin square of order n is an arrangement of n different variables in a $n \times n$ array such that each variable occurs exactly once in every row or column (Euler 1782).

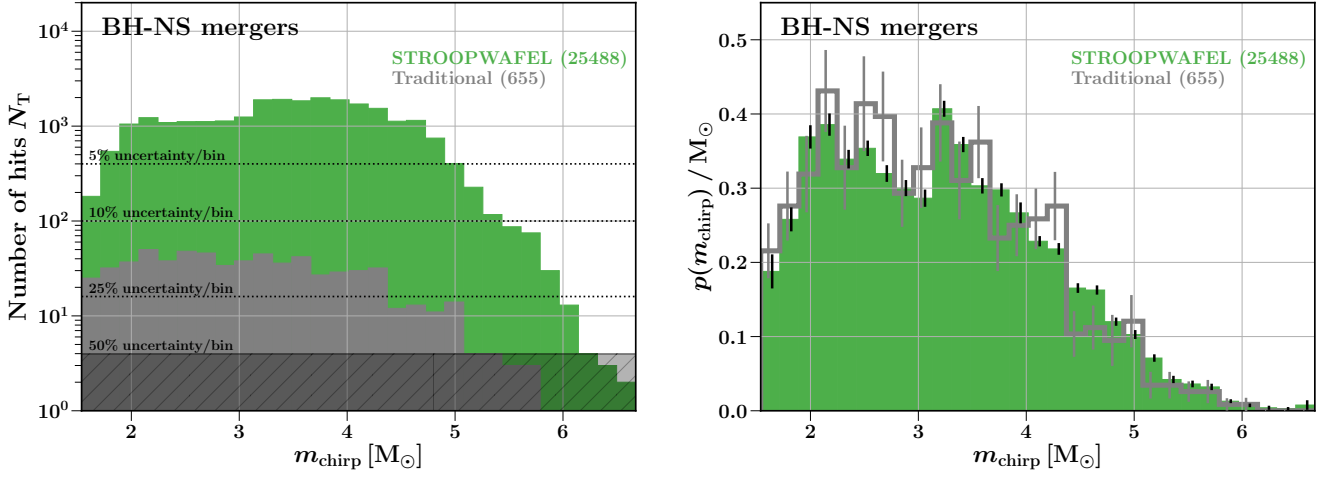



Figure 8. **Left panel:** histograms of the number of target BH–NS binaries N_T found per chirp mass bin m_{chirp} for traditional Monte Carlo sampling (grey) and the STROOPWAFEL sampling method presented in this study (green). The standard Monte Carlo fractional uncertainties (i.e. Poisson noise) are shown with dashed lines in the background. We mark everything below 4 events (i.e. 50% uncertainty) as statistically insignificant as it is consistent with no hits within 2 standard deviations. **Right panel:** the BH–NS chirp mass probability distribution; STROOPWAFEL results have been re-sampled with weights from Eq. (13). The metallicity assumed in the simulations is $Z = 0.001$. The bin width is approximately $0.2M_\odot$ and is constant between the traditional and STROOPWAFEL algorithm. 

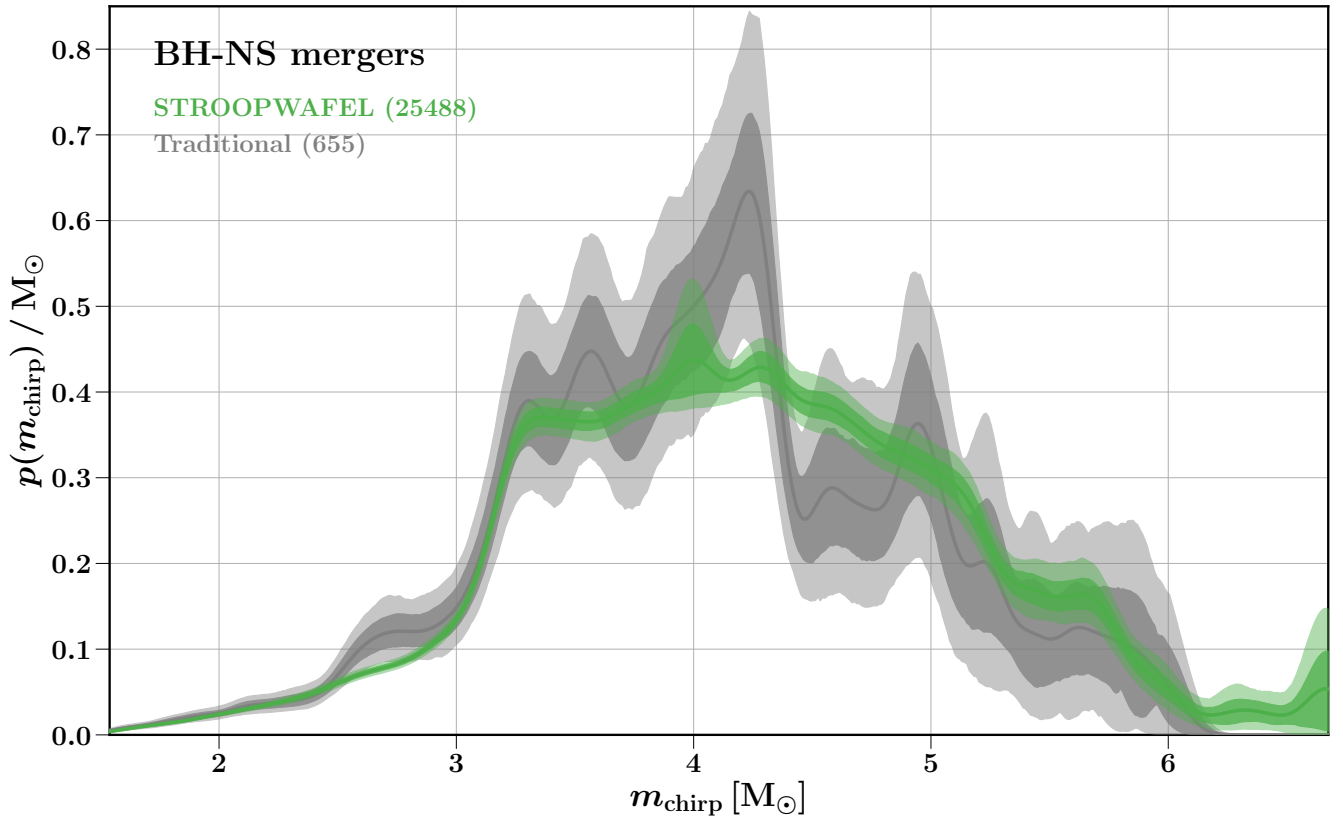


Figure 9. Predicted distribution of the chirp mass of the merging BH–NS population using STROOPWAFEL (green) and traditional (grey) sampling. In both cases the simulation uses $N = 10^6$ samples and the distributions are weighted by the sensitivity of gravitational-wave interferometers using Fishbach & Holz (2017). Shaded regions show the 1- and 2- σ confidence intervals which are calculated by bootstrapping the samples 1000 times. This distribution is for a particular set of model assumptions, including a single metallicity $Z = 0.001$, and an integration over a metallicity-dependent cosmic star formation history is required for comparisons with observations. The same `scipy` kernel density estimator smoothing with a dimensionless kernel density estimator factor of about 0.1 is used for traditional and STROOPWAFEL distributions (see also Appendix D).

Adaptive distribution choices beyond a mixture of Gaussians could also be explored.

4.3 Adapting to uncertainty in distribution functions

Observational selection effects must be applied to model predictions in order to statistically compare models against observations. These selection effects are generally applied after population synthesis models are generated, and may place significant weight on rarely-formed systems in the tails of output distribution functions (e.g., higher-mass DCO binaries). Even though STROOPWAFEL can greatly improve the overall number of systems produced from a simulation, there may still be relatively few systems in these tails.

In principle, we could tune STROOPWAFEL to produce a model population weighted towards any observational population distribution, i.e., optimising for observational selection effects and spending less time on systems which do not contribute to the observed sample. This can be achieved by incorporating selection effects directly in the instrumental distribution rather than applying them to STROOPWAFEL outputs. This can be practically implemented by changing the instrumental distribution weights. At the moment all Gaussians contribute equally to the mixture distribution with a weight $1/N_{T,expl}$ but instead the contribution of each Gaussian can be weighted with the probability of observing the system to focus the simulation on systems that are more likely to contribute to the observable population. An extreme example of this approach would be re-defining the target population to be an even rarer subset of the initial target population, e.g., tails of a distribution function. STROOPWAFEL could also be used to sample from regions of the initial parameter space giving rise to properties consistent with specific observed systems (see also [Andrews et al. 2018](#)).

The current implementation of STROOPWAFEL might be thought of as something like adaptive mesh refinement, familiar from hydrodynamics, applied to the phase space of binary population synthesis. This potential future development of STROOPWAFEL would be refining in the space of predicted observables. This development of STROOPWAFEL could naturally be applied to modelling any population for comparison to observations, not only intrinsically rare populations.

4.4 STROOPWAFEL in higher dimensions

The demonstrations in this paper have all used STROOPWAFEL to sample in a three-dimensional birth parameter space of the two component masses and the initial orbital separation. STROOPWAFEL can be readily applied to sample in more dimensions. The scaling parameter κ may well have a different optimal value for higher dimensions, and should be investigated before being applied to higher dimensions, although we anticipate moving away from using diagonal covariances (see subsection 4.2). Additional potential dimensions to add to the space of initial conditions include, e.g., initial compositions, the initial eccentricity of the system, or the spins of the stars. Moreover, in COMPAS systems are labelled from the start with vectors representing normalised versions of the supernova kicks that will be applied during compact-object formation (see also, e.g., [Andrews et al. 2018](#)); each kick adds three dimensions to the parameter space. Potentially, applying STROOPWAFEL to the kick vectors can be promising since the kick magnitudes and directions can significantly affect the fates of the systems. In our simulations this would be especially important to increase the gain in simulating NS–NS mergers, as the kicks contribute most to the current stochastic output surfaces for this target population.

4.5 Combining STROOPWAFEL with MCMC or Gaussian process regression emulators on continuous spaces

STROOPWAFEL could be applied directly to the combined parameter space of initial parameters of an individual system (e.g., the initial masses and separations) and hyper-parameters describing the model assumptions (e.g., wind-driven mass loss rates, common-envelope physics).

Alternatively, STROOPWAFEL could be combined with other methods for exploring the parameter space, such as emulators based on Gaussian process regression (see, e.g., [Barrett et al. 2017](#); [Taylor & Gerosa 2018](#)). Some of the parameters map continuously to the output space, while others exhibit discontinuities. STROOPWAFEL distinguishes itself in handling bifurcations and stochastic output surfaces. On the other hand, emulators can be more efficient in sampling parameters that are smooth. Thus an intended future development is to combine STROOPWAFEL with such methods to obtain the best overall efficiency.

STROOPWAFEL output samples could also be converted into probability distributions using Gaussian mixture models based on Dirichlet processes ([Del Pozzo et al. 2018](#)).

5 SUMMARY AND CONCLUSIONS

We have presented a new sampling algorithm that aims to improve the efficiency of simulating rare events in astrophysical populations, and demonstrated its utility for binary population synthesis of gravitational-wave merger populations. Our algorithm STROOPWAFEL adaptively improves the sampling distribution to focus more computational time on the target population. Some key findings of our investigation are:

- (i) Using STROOPWAFEL we find a factor of about $25\text{--}200\times$ more systems of interest in simulations of a certain length, as compared to Monte Carlo sampling from the birth distributions. To simulate the same number of events of interest with such commonly-used Monte Carlo sampling would require up to two orders of magnitude more computational time. This gain will improve binary population synthesis simulations by making it computationally feasible both to include more details of the relevant massive-star physics and to explore a greater number of variations of the physical assumptions of the model.
- (ii) The increase in efficiency of STROOPWAFEL leads to higher-resolution mapping of both the input and output parameter space. This reduces the sampling uncertainty by factors of ≈ 3 to 10 for our simulations with 10^6 total samples.
- (iii) STROOPWAFEL improvements are particularly significant when simulating extremely rare events or tails of distribution functions, such as the most massive BH–BH mergers or early NS–NS mergers.
- (iv) One of the core strengths of STROOPWAFEL is that it can handle well the bifurcations and discontinuities that naturally occur in the parameter spaces of binary population synthesis simulations. Such stochasticity often poses a challenge for applying sampling and emulation methods such as Markov Chain Monte Carlo and Gaussian process regression emulators that rely on smoothness to converge and produce independent samples.

Future improvements to the STROOPWAFEL algorithm (dis-

cussed in Section 4) should be able to further improve its performance. This could make it more realistic for next-generation binary population synthesis simulations to include detailed stellar evolution models whilst also exploring more variations in the model physics and assumptions. Such improvements will help in comparing population models to population data, and so help to constrain the physics of evolutionary processes occurring on timescales too long to directly observe.

ACKNOWLEDGEMENTS

We thank especially T. Fragos, D. Szécsi and M. Renzo for constructive discussions and comments. We also thank C. Berry, R. Willcox and M. Wassink for their helpful suggestions on this paper. We thank D. Stops for technical support. We thank the anonymous referee for their helpful comments and careful feedback on this manuscript. FSB, AVG, IM, SJ and JG thank the Kavli Foundation, Niels Bohr institute and DARK Cosmology Centre in Copenhagen for their hospitality and for organizing the Kavli summer school in gravitational-wave astrophysics 2017 where part of this work has been performed. The work has been performed under the Project HPC-EUROPA3 (INFRAIA-2016-1-730897), with the support of the EC Research Innovation Action under the H2020 Programme; in particular, FSB gratefully acknowledges the support of the School of Mathematics, University of Edinburgh and Birmingham Institute for Gravitational Wave Astronomy, University of Birmingham and the computer resources and technical support provided by EPCC. FSB also acknowledges support through the McKinsey excellence grant and Kapteyn grant. SJ and IM thank the Aspen Center for Physics for hospitality whilst part of this work was performed, supported by National Science Foundation grant PHY-1607611. SJ is further grateful for partial support of this work via a grant from the Simons Foundation. SdM and FSB acknowledge funding by the European Union's Horizon 2020 research and innovation programme from the European Research Council (ERC) (Grant agreement No. 715063), and by the Netherlands Organisation for Scientific Research (NWO) as part of the Vidi research program BinWaves with project number 639.042.728. SS is supported by the Australian Research Council Centre of Excellence for Gravitational Wave Discovery (OzGrav), through project number CE170100004. AVG acknowledges support from Consejo Nacional de Ciencia y Tecnología (CONACYT).

Software:

COMPAS (Stevenson et al. 2017; Barrett et al. 2018; Vigna-Gómez et al. 2018)

Python available from python.org

matplotlib (Hunter 2007)

NumPy (van der Walt et al. 2011)

ipython/jupyter (Perez & Granger 2007; Kluyver et al. 2016)

This research has made use of data, software and/or web tools obtained from the Gravitational Wave Open Science Center (<https://www.gw-openscience.org>), a service of LIGO Laboratory, the LIGO Scientific Collaboration and the Virgo Collaboration. LIGO is funded by the U.S. National Science Foundation. Virgo is funded by the French Centre National de Recherche Scientifique (CNRS), the Italian Istituto Nazionale della Fisica Nucleare (INFN) and the Dutch Nikhef, with contributions by Polish and Hungarian institutes.

REFERENCES

Abbott B. P., et al., 2016, *Phys. Rev. Lett.*, 116, 061102

Abbott B. P., et al., 2017, *Phys. Rev. Lett.*, 119, 161101

Abbott B. P., et al., 2018, *Living Reviews in Relativity*, 21, 3

Abt H. A., 1983, *ARA&A*, 21, 343

Abt H. A., Gomez A. E., Levy S. G., 1990, *ApJS*, 74, 551

Andrews J. J., Zezas A., Fragos T., 2018, *The Astrophysical Journal Supplement Series*, 237, 1

Barrett J. W., Mandel I., Neijssel C. J., Stevenson S., Vigna-Gómez A., 2017, in Brescia M., Djorgovski S. G., Feigelson E. D., Longo G., Caviuoti S., eds, IAU Symposium Vol. 325, *Astroinformatics*. pp 46–50 ([arXiv:1704.03781](https://arxiv.org/abs/1704.03781)), doi:10.1017/S1743921317000059

Barrett J. W., Gaebel S. M., Neijssel C. J., Vigna-Gómez A., Stevenson S., Berry C. P. L., Farr W. M., Mandel I., 2018, *MNRAS*, 477, 4685

Belczynski K., Kalogera V., Bulik T., 2002, *ApJ*, 572, 407

Belczynski K., Holz D. E., Bulik T., O'Shaughnessy R., 2016, *Nature*, 534, 512

Cappé O., Guillin A., Marin J. M., Robert C. P., 2004, *Journal of Computational and Graphical Statistics*, 13, 907

Cappé O., Douc R., Guillin A., Marin J.-M., Robert C. P., 2008, *Statistics and Computing*, 18, 447

Cornuet J.-M., Marin J.-M., Mira A., Robert C. P., 2012, *Scandinavian Journal of Statistics*, 39, 798

de Mink S. E., Mandel I., 2016, *MNRAS*, 460, 3545

Del Pozzo W., Berry C. P. L., Ghosh A., Haines T. S. F., Singer L. P., Vecchio A., 2018, *MNRAS*, 479, 601

Dominik M., Belczynski K., Fryer C., Holz D., Berti E., Bulik T., Mandel I., O'Shaughnessy R., 2012, *Astrophys. J.*, 759, 52

Eglajs V., Audze P., 1977, *Problems of Dynamics and Strengths*, 35, 104

Euler L., 1782, *Euler Archive - All Works*, 530, 9, 85

Fermi E., Richtmyer R. D., 1948, *J. 10.2172/4423221*

Fishbach M., Holz D. E., 2017, *ApJ*, 851, L25

Foucart F., Hinderer T., Nisanke S., 2018, *Phys. Rev. D*, 98, 081501

Fryer C. L., Belczynski K., Wiktorowicz G., Dominik M., Kalogera V., Holz D. E., 2012, *ApJ*, 749, 91

Goldberg D., Mazeh T., 1994, *A&A*, 282, 801

He H. Y., Owen A. B., 2014, *arXiv e-prints*, p. [arXiv:1411.3954](https://arxiv.org/abs/1411.3954)

Hesterberg T., 1995, *Technometrics*, 37, 185

Hunter J. D., 2007, *Computing in Science and Engineering*, 9, 90

Hurley J. R., Pols O. R., Tout C. A., 2000, *MNRAS*, 315, 543

Hurley J. R., Tout C. A., Pols O. R., 2002, *MNRAS*, 329, 897

Iman R. L., Davenport J. M., Zeigler D. K., 1980, Technical report, Latin hypercube sampling (program user's guide).[LHC, in FORTRAN]. Sandia Labs., Albuquerque, NM (USA)

Iman R. L., Helton J. C., Campbell J. E., 1981, *Journal of quality technology*, 13, 174

Jones N. D., 1977, *DAIMI Report Series*, 6

Kahn H., Harris T. E., 1951, *National Bureau of Standards applied mathematics series*, 12, 27

Kahn H., Marshall A. W., 1953, *Journal of the Operations Research Society of America*, 1, 263

Kalogera V., 1996, *ApJ*, 471, 352

Kalogera V., 2000, *ApJ*, 541, 319

Kalogera V., Webbink R. F., 1998, *ApJ*, 493, 351

Klencki J., Moe M., Gladysz W., Chruslinska M., Holz D. E., Belczynski K., 2018, *A&A*, 619, A77

Kluyver T., et al., 2016, in *ELPUB*. pp 87–90

Kolb U., 1993, *A&A*, 271, 149

Kroupa P., 2001, *MNRAS*, 322, 231

Kruckow M. U., Tauris T. M., Langer N., Kramer M., Izzard R. G., 2018, *MNRAS*,

Liu J. S., 2008, *Monte Carlo strategies in scientific computing*. Springer Science & Business Media

Mandel I., Farmer A., 2018, *arXiv e-prints*, p. [arXiv:1806.05820](https://arxiv.org/abs/1806.05820)

Mandel I., de Mink S. E., 2016, *MNRAS*, 458, 2634

Mapelli M., 2018, *arXiv e-prints*, p. [arXiv:1809.09130](https://arxiv.org/abs/1809.09130)

Marchant P., Langer N., Podsiadlowski P., Tauris T. M., Moriya T. J., 2016, *A&A*, 588, A50

Mazeh T., Goldberg D., Duquennoy A., Mayor M., 1992, *ApJ*, 401, 265

McKay M. D., Beckman R. J., Conover W. J., 1979, *Technometrics*, 21, 239

- Metropolis N., Ulam S., 1949, *Journal of the American Statistical Association*, 44, 335
- Moe M., Di Stefano R., 2017, *The Astrophysical Journal Supplement Series*, 230, 15
- Neijssel C. J., et al., 2019, arXiv e-prints, p. arXiv:1906.08136
- Öpik E., 1924, Publications of the Tartu Astrofizika Observatory, 25, 1
- Ortiz L. E., Pack Kaelbling L., 2013, arXiv e-prints, p. arXiv:1301.3882
- Owen A., Zhou Y., 2000, *Journal of the American Statistical Association*, 95, 135
- Pennanen T., Koivu M., 2006, in , Monte Carlo and Quasi-Monte Carlo Methods 2004. Springer, pp 443–455
- Perez F., Granger B. E., 2007, *Computing in Science and Engineering*, 9, 21
- Politano M., 1996, *ApJ*, 465, 338
- Pols O., Hurley J., Tout C., 1998, in IAU Symposium. p. 607
- Robert C., Casella G., 2013, Monte Carlo statistical methods. Springer Science & Business Media
- Safarzadeh M., Ramirez-Ruiz E., Andrews J. J., Macias P., Fragos T., Scanapieco E., 2019, *ApJ*, 872, 105
- Sana H., et al., 2012, *Science*, 337, 444
- Smarr L., D. Blandford R., 1976, *The Astrophysical Journal*, 207, 574
- Spera M., Mapelli M., Bressan A., 2015, *MNRAS*, 451, 4086
- Stevenson S., Vigna-Gómez A., Mandel I., Barrett J. W., Neijssel C. J., Perkins D., de Mink S. E., 2017, *Nature Communications*, 8, 14906
- Stevenson S., Sampson M., Powell J., Vigna-Gómez A., Neijssel C. J., Szécsi D., Mandel I., 2019, arXiv e-prints, p. arXiv:1904.02821
- Taylor S. R., Gerosa D., 2018, *Phys. Rev. D*, 98, 083017
- The LIGO Scientific Collaboration et al., 2018, arXiv e-prints, p. arXiv:1811.12907
- Torrie G. M., Valleau J. P., 1977, *Journal of Computational Physics*, 23, 187
- Tout C. A., 1991, *MNRAS*, 250, 701
- van der Walt S., Colbert S. C., Varoquaux G., 2011, *Computing in Science and Engineering*, 13, 22
- Veatch E., Guibas L. J., 1995, in Proceedings of the 22Nd Annual Conference on Computer Graphics and Interactive Techniques. SIGGRAPH '95. ACM, New York, NY, USA, pp 419–428, doi:10.1145/218380.218498, <http://doi.acm.org/10.1145/218380.218498>
- Vigna-Gómez A., et al., 2018, *MNRAS*, 481, 4009
- Woosley S. E., 2017, *The Astrophysical Journal*, 836, 244
- Zackay B., Venumadhav T., Dai L., Roulet J., Zaldarriaga M., 2019, arXiv e-prints, p. arXiv:1902.10331

APPENDIX A: DERIVATION OF THE VARIANCE USED FOR OPTIMISING THE LENGTH OF THE STROOPWAFEL EXPLORATION PHASE

In this section we derive the expression for the variance (Eq. 17) on the rate estimate used to optimize f_{expl} . We can estimate the optimal fraction f_{expl} of samples that we should spend in the exploratory phase by taking into account the probability of not identifying a target population forming region in the exploratory phase. We assume that we sample from the mixture distribution $Q(x)$ which is given by Eq. 12:

$$Q(x) = f_{\text{expl}}\pi(x) + (1 - f_{\text{expl}})\tilde{q}(x), \quad (\text{A1})$$

where π is the prior (used for the exploratory phase sampling) and \tilde{q} is the mixture of Gaussians. We also assume we aim to estimate the rate \mathcal{R}_T with N total samples. After simulating these N samples we may have identified a target binary-forming region with total weight z_1 whereas a region with weight z_2 is yet unidentified such that the estimated rate of the target population is (Eq. 16)

$$\mathcal{R}_T = \underbrace{z_1}_{\text{identified}} + \underbrace{z_2}_{\text{unidentified}} \approx \mathbb{E}_Q[\mathcal{R}_T]. \quad (\text{A2})$$

The variance on \mathcal{R}_T , $\mathbb{V}_Q[\mathcal{R}_T]$, is a measure for the uncertainty

of the estimated rate $\mathbb{E}_Q[\mathcal{R}_T]$. Therefore, the optimal value of f_{expl} is one that minimizes the variance on \mathcal{R}_T . To determine this we first derive an estimate for the variance $\mathbb{V}_Q[\mathcal{R}_T]$.

Using the continuous definition for the variance we have

$$\mathbb{V}_Q[\mathcal{R}_T] = \frac{1}{N} \left[\int \phi(x)^2 \widetilde{w(x)}^2 Q(x) dx - \mathbb{E}_Q[\mathcal{R}_T]^2 \right]. \quad (\text{A3})$$

where the $1/N$ factor comes from taking the variance of the mean (average).

Since

$$\widetilde{w(x)} = \pi(x)/Q(x), \quad (\text{A4})$$

and $\mathbb{E}_Q[\mathcal{R}_T] = z_1 + z_2$, we find that

$$\mathbb{V}_Q[\mathcal{R}_T] = \frac{1}{N} \left[\int \phi(x)^2 \widetilde{w(x)}\pi(x) dx - (z_1 + z_2)^2 \right]. \quad (\text{A5})$$

Since we assumed the target binary forming region is equal to $z_1 + z_2$, by definition no binaries of the target population are found outside this region and thus $\phi(x) = 0$ outside z_1 and z_2 . Using this we can rewrite the integral in Equation A5 as

$$\int \phi(x)^2 \widetilde{w(x)}\pi(x) dx = \int_{z_1} \widetilde{w(x)}\pi(x) dx + \int_{z_2} \widetilde{w(x)}\pi(x) dx. \quad (\text{A6})$$

We now assume that we have found enough binaries of our target population in our exploratory phase and that we don't have a bias for an output function such that the Gaussian mixture $\tilde{q}(x)$ is effectively flat over z_1 . In other words, we assume that on the target binary forming region z_1

$$\tilde{q}(x) \approx \frac{\pi(x)}{z_1}. \quad (\text{A7})$$

where $\pi(x)$ is the birth distribution.

Using this in Eq. A1, we approximate that $Q(x) \approx (1 - f_{\text{expl}})\pi(x)/z_1 + f_{\text{expl}}\pi(x)$ in z_1 such that

$$\widetilde{w(x)} \approx \frac{1}{\left(\frac{1-f_{\text{expl}}}{z_1}\right) + f_{\text{expl}}} \quad \text{for } x \text{ in } z_1. \quad (\text{A8})$$

In addition, we also assume that the our Gaussian mixture \tilde{q} is negligible outside of the target binary forming regions, i.e. $\tilde{q}(x) = 0$ outside of z_1 . In other words we assume that z_2 is far enough from z_1 that the probability is zero to sample it with \tilde{q} , and that we have completely missed it during the exploratory sampling. By doing so we obtain that on z_2 we have

$$\widetilde{w(x)} \approx \frac{1}{f_{\text{expl}}} \quad \text{for } x \text{ on } z_2 \quad (\text{A9})$$

Substituting Eqs. A8 and A9 into the integral expression of Eq. A6 then yields:

$$\int \phi(x)^2 \widetilde{w(x)}\pi(x) dx \approx \frac{z_1}{\left(\frac{1-f_{\text{expl}}}{z_1}\right) + f_{\text{expl}}} + \frac{z_2}{f_{\text{expl}}} \quad (\text{A10})$$

where we used $\int_{z_1} \pi(x) dx = z_1$ and $\int_{z_2} \pi(x) dx = z_2$.

We can now write the variance as

$$\mathbb{V}_Q[\mathcal{R}_T] \approx \frac{1}{N} \left[\frac{z_1}{\left(\frac{1-f_{\text{expl}}}{z_1}\right) + f_{\text{expl}}} + \frac{z_2}{f_{\text{expl}}} - (z_1 + z_2)^2 \right], \quad (\text{A11})$$

i.e., Eq. 17.

APPENDIX B: TOY MODEL

We construct a toy model that can be run without too much computational burden and is inspired by binary population synthesis simulations to study the performance of STROOPWAFEL. The advantage of a toy model is that the moments are analytically known and the toy model can be repeatedly evaluated at minimal computational cost. We use this toy model to investigate a suitable choice for the scale parameter κ in the width of Gaussian sampling distributions (see Eq. 11).

We build the toy model in the 3-dimensional parameter space defined by the initial parameters x_1, x_2 and x_3 . The distribution functions (and ranges) of x_1, x_2 and x_3 are chosen to be similar to the initial parameters $m_{1,i}, a_i$ and q_i , used for our binary population synthesis model (see Appendix C). Similarly to the birth distribution of a_i in the binary population synthesis code, we sample in $\log x_2$. The output of the toy model is constructed from a union of disconnected volumes D and an output function $\phi(\mathbf{x}_i)$ given by

$$\mathbb{1}_{D, \text{toy model}}(\mathbf{x}_i) = \begin{cases} 1 & \text{if } \mathbf{x}_i \in D \\ 0 & \text{otherwise,} \end{cases} \quad (\text{B1})$$

where $D = D_0 \cup D_1 \cup D_2$ is the union of three cuboids with the following vectors for the location of the center and the half-length of each cuboid in the x_1, x_2 and x_3 direction

$$D_0 = [20, 34, 0.3] \pm [1.9, 8.0, 0.1],$$

$$D_1 = [40, 1.0, 0.3] \pm [1.7, 0.6, 0.2],$$

$$D_2 = [34, 7.0, 0.8] \pm [1.8, 0.6, 0.1].$$

The two-dimensional projected distribution of the union D in x_3 and $\log x_2$ is shown in bright green in Figure B1. One might notice that the disconnected target regions look similar to the ‘boats’ in the game Battleship - which is often played with a strategy that is conceptually similar to STROOPWAFEL. The fractional rates produced by these regions are the integrals over their volumes of the local density of the birth probability distributions (the prior distribution). The fractional rate of D in the initial parameter space equals $V_D = 0.0013127$, where a prior (birth distribution) of $x_1 \propto x_1^{-2.3}$ is assumed on x_1 and flat priors are assumed on $\log x_2$, and x_3 . The value of V_D is chosen to be similar to the average yield of double compact object mergers in the publicly available simulations⁶ of Vigna-Gómez et al. (2018). In addition, the fractional rate produced from D_0 is similar to that from D_1 , whereas the fractional rate from D_2 is relatively 10 times smaller. (For these parameter scalings, the absolute volume of D_0 is larger than the volume of D_1 , but the different prior distributions weight the volume of D_1 more highly.)

We run repeated simulations varying the parameter κ in STROOPWAFEL. We fix the total number of samples to $N = 10^6$. We know the true value for the volume integral V_D and calculate for each simulation the deviation between the fractional rate estimate and this true weighted volume. The closer to zero this deviation is, the better the estimate. For each variation of κ we run 100 simulations.

The result of one such simulation for $\kappa = 2$ is shown in Fig. B1. In the STROOPWAFEL simulation the three islands are recovered with much better resolution than in traditional Monte Carlo sampling from the prior. The dark regions around the islands

D_0, D_1 and D_2 in the STROOPWAFEL simulation demonstrate that our method focuses more of the computational time around the regions of interest. In both simulations, the islands D_0 and D_1 contain more samples than D_2 , as expected.

Figure B2 shows, in blue, the 1σ deviation from the true value for STROOPWAFEL simulations as a function of κ . The result of repeated simulations with traditional birth distribution Monte Carlo sampling is shown as a reference on the right. Shades of green in the background show regions of 0.3%, 1% and 3% fractional sampling uncertainty on the rate estimate.

Excessively small values of κ lead to biases in the estimated rate of more than 10%. This is because overly narrow Gaussian distributions create ‘holes’ in the adapted sampling distribution, which then no longer completely cover or characterize the regions of interest. As a result the refinement phase misses part of D , leading to systematic underestimation of the true rate.

On the other hand, excessively large κ values decrease the efficiency of the refinement phase as many samples drawn from the mixture of Gaussians $q(\mathbf{x})$ fall outside the target regions. The scatter from the true rate estimate approaches the scatter obtained from the traditional Monte Carlo simulation as κ increases.

We find that $\kappa \gtrsim 1.5$ is a robust for all our test simulations and yields substantially better estimates on rates and distribution functions compared to traditional Monte Carlo simulations. We therefore adopt $\kappa = 2$ (the red dotted line in Figure B2) when sampling in three dimensions.

APPENDIX C: BINARY POPULATION SYNTHESIS MODEL SET-UP

We test the performance of the algorithm STROOPWAFEL by implementing our algorithm in the rapid binary population synthesis code COMPAS. COMPAS (Compact Object Mergers: Population Astrophysics and Statistics) is designed to study uncertainties in stellar binary evolution models and constrain them with observations, particularly those of gravitational-wave sources (Stevenson et al. 2017; Barrett et al. 2018; Vigna-Gómez et al. 2018). COMPAS interpolates between and extrapolates beyond stellar evolutionary tracks based on algorithms from the code SSE (Hurley et al. 2000), which rely on analytic fits of single star evolution from Pols et al. (1998). COMPAS relies on an approximate and parametrized treatment of the physical processes, including for binary interactions, and can typically compute a predicted final outcome for a single binary system within a second.

In this work the code is used to analyse DCO systems that form through isolated binary evolution, which often involves the common-envelope phase (e.g. Smarr & D. Blandford 1976). The approach we use to simulate a synthetic population of DCOs is similar to other binary population synthesis studies (including, e.g., Hurley et al. 2002; Belczynski et al. 2002; Dominik et al. 2012). We evolve a population of binary systems from their birth until they form a DCO system or otherwise either merge or disrupt. We then make a sub-selection of the DCOs that consist of two compact objects that merge within the age of the Universe through gravitational-wave emission and study the properties of this population.

In general, we follow the Fiducial model described in Vigna-Gómez et al. (2018). We mention the most important assumptions here. The birth distribution for the primary mass $m_{1,i}$ is chosen to be a power law distribution known as the initial mass function (IMF) where $p(m_{1,i}) \propto m_{1,i}^{-\alpha}$ with $\alpha = 2.3$ for massive stars (Kroupa 2001). For the simulations we draw $m_{1,i}$ in $[5, 150] M_\odot$.

⁶ Populations are available at <http://www.sr.bham.ac.uk/compas/data>.

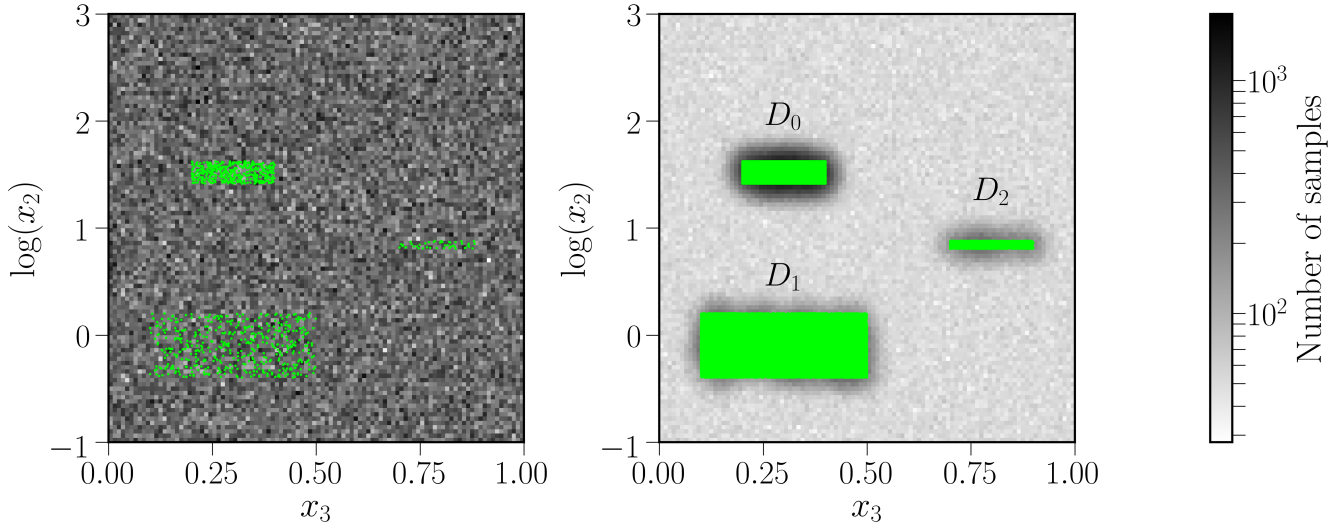


Figure B1. Toy model illustration: distribution of the 10^6 samples drawn in the birth distribution Monte Carlo (left-panel) and STROOPWAFEL (right-panel) simulation. The panels show two-dimensional projections of $\log x_2$ and x_3 from the three dimensional parameter space. In both figures the sampling density (gray) is shown through a two-dimensional histogram with 100×100 bins. Over plotted (green) are the samples that lie within the volume V_D and recover the rare outcome D . Dark regions surrounding the green areas of interest in the right plot indicate that our STROOPWAFEL algorithm focuses more of the computational time around the region of interest.

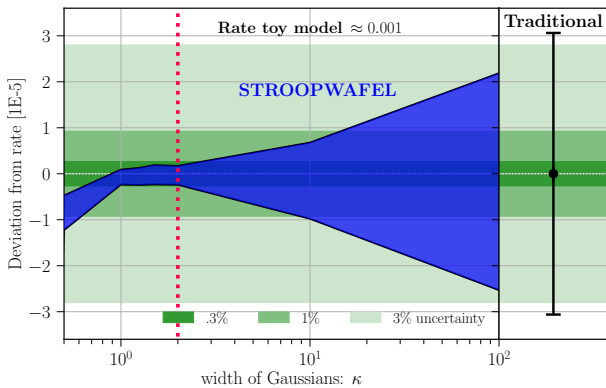


Figure B2. Toy model test result: 1σ deviations from the true volume integral V_D when estimating the fractional rate using STROOPWAFEL with different values for the scaling factor of the width of the Gaussians κ (left panel). The true fractional rate in the toy model is known and equals approximately 0.001. The deviations are calculated by running 100 repeated simulations with a total number of $N = 10^6$ samples per simulation. The green contours show a .3%, 1% and 3% fractional sampling uncertainty on the rate estimate of the target population. On the right the 1σ uncertainty of the traditional Monte Carlo sampled simulation is shown with an error bar; this matches the expected fractional uncertainty $1/\sqrt{N_T}$. In this work we adopt $\kappa = 2$ (red dotted line).

The initial mass ratio $q_i = m_{2,i}/m_{1,i}$ of binary systems is suggested from observations to follow a flat distribution (e.g. [Tout 1991](#); [Mazeh et al. 1992](#); [Goldberg & Mazeh 1994](#); [Sana et al. 2012](#)) given by $p(q_i) \propto 1$. We adopt $q_i \in [0, 1]$. The initial separation a_i is assumed to be flat in the log, also known as Öpik's law $p(a_i) \propto \frac{1}{a_i}$ ([Öpik 1924](#); [Abt 1983](#)). We choose $a_i \in [0.01, 1000]$ AU. We as-

sume that all our binaries have circular orbits at birth. These distributions and parameter ranges resemble commonly used settings for binary population synthesis simulations.

Our changes to the Fiducial model from [Vigna-Gómez et al. \(2018\)](#) are the following:

- We use a metallicity of $Z = 0.001$ for all our simulations.
- We use the DELAYED prescription for the core-collapse supernovae treatment from [Fryer et al. \(2012\)](#).
- We use a prescription for pair-instability supernovae and pulsational pair instability supernovae based on [Woosley \(2017\)](#). The implementation in COMPAS is described in [Stevenson et al. \(2019\)](#).

We fix the total number of binaries in each simulation to $N = 10^6$ both for when using birth distribution Monte Carlo and STROOPWAFEL sampled simulations. The total computational time for each of the birth distribution Monte Carlo simulations is approximately 180 CPU hours. The total computational time of the STROOPWAFEL simulations is up to a factor of 2 lower as a result of a decrease in average simulation time per sample in our sampling method compared to traditional sampling (see Section 3.2). This result rises from binaries that become a gravitational-wave source costing on average less computational time to simulate with COMPAS than other binaries.

APPENDIX D: BANDWIDTH VARIATIONS OF KERNEL DENSITY ESTIMATOR

Figure 8 and 9 use the same resolution (i.e., bandwidth or bin width) to estimate the chirp mass distribution function for the output of both the STROOPWAFEL and birth distribution Monte Carlo simulation. The bin width of a histogram or the bandwidth of a kernel density estimator can strongly influence the estimated distribution. Hence, in practice, the bandwidth should be adapted to the resolution available in the data.

We show in Figure D1 the predicted chirp mass distribution of BH–BH mergers using adapted resolutions for birth distribution Monte Carlo sampling and STROOPWAFEL sampling.

We estimate the adapted bandwidth using *Scott's Rule* which, in one dimension, scales as $\propto N_T^{-1/4}$. This is the default bandwidth choice in the `scipy` kernel density estimator function. For the STROOPWAFEL sampling we replace N_T with the effective sample size (ESS) given by $(\sum_i \hat{w}_i)^2 / \sum_i \hat{w}_i^2$ (which, in practice, is approximately equal to $N_{T, \text{STROOPWAFEL}}$). The top panel of Figure 9 shows the estimated chirp mass distribution from both sampling methods for a dimensionless kernel density estimator factor for the bandwidth of $\text{ESS}_{\text{STROOPWAFEL}}^{-1/4} \approx 0.044$. This bandwidth is too small for the $53\times$ smaller birth distribution Monte Carlo BH–BH population which therefore shows significant statistical noise fluctuations. The middle panel of Fig. 9 shows the estimated chirp mass distribution from both sampling methods for a bandwidth of $N_{T, \text{traditional}}^{-1/4} \approx 0.12$. This bandwidth causes smaller statistical fluctuations for the traditional sampling, but removes some of the more detailed features for the STROOPWAFEL sampled distribution. In the bottom panel the two plots are combined, showing the distributions with the relative bandwidths. The STROOPWAFEL obtains smaller uncertainties on the distribution as well as a higher resolution, which is a result from the higher number of BH–BH mergers found in this simulation.

This paper has been typeset from a \LaTeX file prepared by the author.

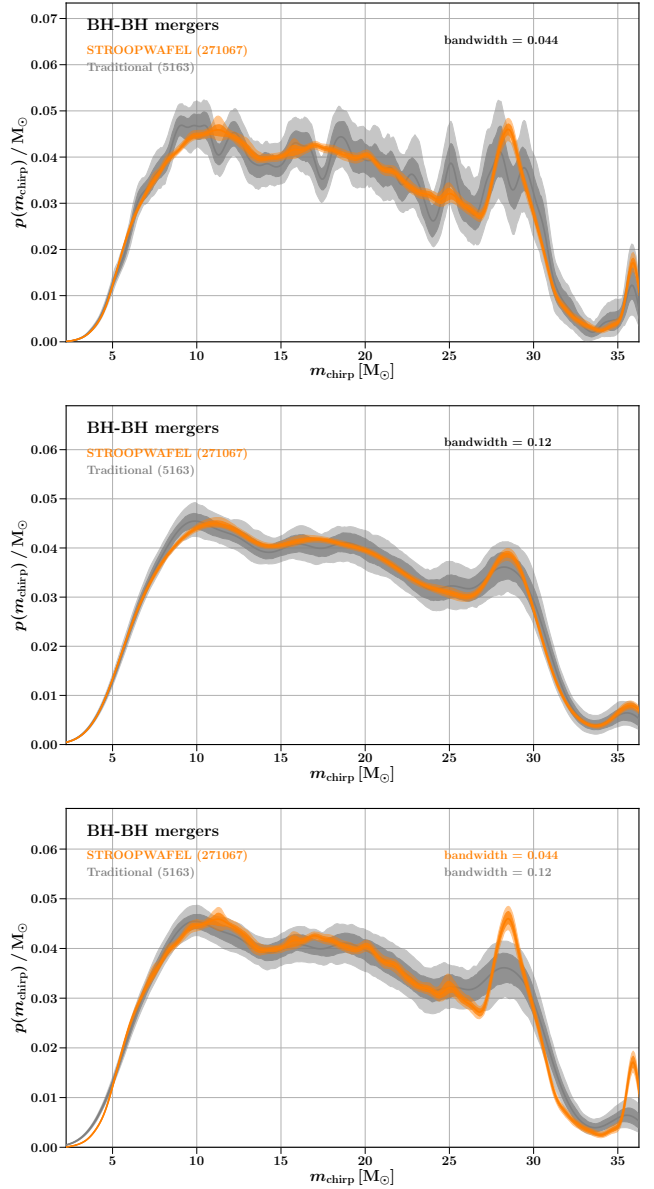


Figure D1. Predicted chirp mass distribution of the BH–BH merger population using STROOPWAFEL (orange) and traditional (grey) sampling. In all cases the simulation uses $N = 10^6$ samples and the distributions are weighted to the sensitivity of gravitational-wave interferometers using Fishbach & Holz (2017). Shaded regions show the 1- and 2- σ confidence intervals which are calculated by bootstrapping the samples 100 times. This distribution is for a particular set of model assumptions, including a single metallicity $Z = 0.001$, and an integration over a metallicity-dependent cosmic star formation history is required for comparisons with observations. The same `scipy` kernel density estimator smoothing with a kernel density estimator factor for the bandwidth of about 0.044 (top panel) and 0.12 (middle panel) is used for the traditional and STROOPWAFEL simulations. In the bottom panel a kernel bandwidth of about 0.044 is used for the STROOPWAFEL method whereas for the traditional method we use a bandwidth of about 0.12.